

UCLA

UCLA Electronic Theses and Dissertations

Title

On the Contextual Unfairness of Modern Machine Learning: Graph Neural Networks to Large Language Models

Permalink

<https://escholarship.org/uc/item/1zt783t9>

Author

Subramonian, Arjun

Publication Date

2025

Peer reviewed|Thesis/dissertation

UNIVERSITY OF CALIFORNIA

Los Angeles

On the Contextual Unfairness of Modern Machine Learning:
Graph Neural Networks to Large Language Models

A dissertation submitted in partial satisfaction
of the requirements for the degree
Doctor of Philosophy in Computer Science

by

Arjun Subramonian

2025

© Copyright by
Arjun Subramonian
2025

ABSTRACT OF THE DISSERTATION

On the Contextual Unfairness of Modern Machine Learning:
Graph Neural Networks to Large Language Models

by

Arjun Subramonian

Doctor of Philosophy in Computer Science

University of California, Los Angeles, 2025

Professor Yizhou Sun, Co-Chair

Professor Kai-Wei Chang, Co-Chair

Graph neural networks (GNNs) and large language models (LLMs) have emerged as popular machine learning (ML) models for powering applications such as social recommendation on social media platforms and chat-based assistants, respectively. GNNs and LLMs are both *contextual* models: GNNs operate on social context in social networks, while LLMs process syntactic and semantic context in language. In conjunction with the proliferation of GNNs and LLMs, there is decreasing trust in the fairness of ML [MFP25]. Unfair ML models cause real-world harm, such as the reinforcement of stereotypes [BCZ16, BKD23] and discrimination in hiring. The unfairness of GNNs is exacerbated by social context (e.g., graph structure, message passing). However, this aspect is not explored in research on the fairness of traditional ML models and requires a deeper principled understanding. Moreover, the open-ended nature of LLM generations can make automatic evaluations of syntactic and semantic context-dependent unfairness difficult.

This dissertation tackles technical challenges in addressing the unfairness of GNNs and

LLMs. In the first part, we theoretically and empirically investigate different forms of GNN unfairness (i.e., imputation bias, preferential attachment bias, degree bias), and how they are affected by graph structure and the choice of graph filter. We further propose principled metrics and methods to alleviate GNN unfairness. In the second part of this dissertation, we assess the measurement validity of evaluations of LLM misgendering. In the final part, we return to the relatively simple setting of feedforward neural networks, and even in this setting, we identify and tackle major challenges in obtaining a precise analytical theory of how model design choices and data properties contribute to unfairness. Such a theory for GNNs and LLMs could aid in interpreting model outputs and designing stronger evaluation and mitigation methods for unfairness. Overall, this dissertation develops a principled understanding of and addresses the unfairness of modern ML models, towards preventing the further entrenchment of social inequalities and promoting justice.

The dissertation of Arjun Subramonian is approved.

Levent Sagun

Aditya Grover

Baharan Mirzasoleiman

Kai-Wei Chang, Committee Co-Chair

Yizhou Sun, Committee Co-Chair

University of California, Los Angeles

2025

To my family, younger self, and dark-eyed juncos. . .

“There are more of us than you know and together we will change the world.”

—Nithya Elsa Ramesh [[Ram21](#)]

TABLE OF CONTENTS

1	Introduction	1
1.1	Motivation	1
1.2	Research Questions, Objectives, and Challenges	5
1.3	Dissertation Structure	6
I	On the Unfairness of Graph Neural Networks	10
2	On the Discrimination Risk of Mean Aggregation Feature Imputation in Graphs	11
2.1	Introduction	11
2.2	Related work	14
2.3	Discrimination risk and model unfairness	16
2.4	Graph feature imputation	19
2.5	Discrimination risk of mean aggregation feature imputation	20
2.5.1	Analysis of Theorem 2	22
2.6	Fairer graph feature imputation	23
2.7	Experimental results and discussion	25
2.8	Conclusion	30
2.9	Broader Impacts	31
3	Networked Inequality: Preferential Attachment Bias in Graph Neural Network Link Prediction	32

3.1	Introduction	32
3.2	Related Work	34
3.3	Preliminaries	36
3.4	Theoretical Analysis	37
3.4.1	Symmetric Normalized Filter	38
3.4.2	Within-Group Fairness	40
3.4.3	Random Walk Normalized Filter	42
3.5	Fairness Regularizer	43
3.6	Experiments	44
3.6.1	Validating Theoretical Analysis	44
3.6.2	Within-Group Fairness	46
3.6.3	Fairness Regularizer	47
3.7	Conclusion	49
3.8	Broader Impacts	49
4	Theoretical and Empirical Insights into the Origins of Degree Bias in Graph Neural Networks	51
4.1	Introduction	51
4.2	Background and Related Work	53
4.3	Preliminaries	54
4.4	Test-Time Degree Bias	56
4.4.1	Random Walk Graph Filter	57
4.4.2	Symmetric Graph Filter	61
4.5	Training-Time Degree Bias	63

4.6	Principled Roadmap to Address Degree Bias	66
4.7	Conclusion	68
4.8	Broader Impacts	68
II	On the Unfairness of Large Language Models	71
5	Agree to Disagree? A Meta-Evaluation of LLM Misgendering	72
5.1	Introduction	72
5.2	Related Work	75
5.3	Evaluation Paradigms for LLM Misgendering	77
5.3.1	Pronoun Preliminaries	77
5.3.2	Probability-Based Evaluation	77
5.3.3	Generation-Based Evaluation	78
5.4	Experimental Setup	79
5.4.1	Models and Data	79
5.4.2	Converting Probability-Based to Generation-Based Evaluations	80
5.4.3	Converting Generation-Based to Probability-Based Evaluations	80
5.5	Agreement between Probability-Based and Generation-Based Evaluations	80
5.5.1	Metrics	81
5.5.2	Results	81
5.6	Human Evaluation	85
5.7	Recommendations	89
5.8	Conclusion and Future Work	89
5.9	Broader Impacts	90

III	Towards a Precise Theory of Machine Learning Unfairness	91
6	An Effective Theory of Bias Amplification	92
6.1	Introduction	92
6.1.1	Main Contributions	93
6.1.2	Related Work	94
6.2	Preliminaries	96
6.2.1	Data Distributions	96
6.2.2	Models and Metrics	97
6.3	Theoretical Analysis	99
6.3.1	Main Result: Ridge Regression with Random Projections	100
6.4	Bias Amplification	103
6.4.1	Isotropic Covariance	104
6.4.2	Regularization and Training Dynamics	106
6.5	Minority-Group Bias	110
6.6	Conclusion	112
6.7	Broader Impacts	113
7	Conclusion and Future Directions	114
7.1	Future Directions	115
7.2	Social Dimensions of Fairness	119
	Bibliography	120
A	Appendix for Chapter 2	168

A.1	Proofs	168
A.1.1	Proof of Lemma 1 (special case of total variation distance)	168
A.1.2	Proof of Theorem 1	169
A.1.3	Example mean aggregation imputation algorithms	171
A.1.4	Proof of Theorem 2	172
A.1.5	Extending Theorem 2	176
A.1.6	Proof of Theorem 3	176
A.1.7	Theorem 4	177
A.2	Additional experimental results	180
A.2.1	Datasets	180
A.2.2	Imputation algorithms	180
A.2.3	Models and training	180
A.2.4	Performance evaluation	181
A.2.5	Contraction coefficient	182
B	Appendix for Chapter 3	196
B.1	Proofs	196
B.1.1	Proof of Lemma 3.4.1	196
B.1.2	Proof of Lemma 3.4.2	197
B.1.3	Proof of Theorem 3.4.3	199
B.1.4	Lemma B.1.1 and Proof	200
B.1.5	Proof of Theorem 3.4.4	201
B.2	Approximation of $\Delta^{(b)}$	202
B.2.1	Approximation of $\Delta^{(b)}$ for Φ_s	202

B.2.2	Approximation of $\Delta^{(b)}$ for Φ_r	202
B.3	Datasets Used in §3.6.1	203
B.4	Datasets Used in §3.6.2	204
B.4.1	NBA Dataset	205
B.4.2	German Dataset	205
B.4.3	DBLP-Fairness Dataset	205
B.5	Models	207
B.6	Remaining Plots for §3.6.1	207
B.7	Additional Experiments	210
B.7.1	Additional Experiments for §3.6.1 (4-layer Encoders)	210
B.7.2	Additional Experiments for §3.6.1 (Hadamard Product and MLP LP Score Function)	212
B.7.3	Additional Experiments for §3.6.2	215
B.7.4	Additional Experiments for §3.6.3	215
B.8	Theory Pitfalls	216
B.9	Error Analysis of Φ_r Theoretic Scores	216
C	Appendix for Chapter 4	219
C.1	Overview of Theoretical Analyses of and Hypotheses for Degree Bias	219
C.1.1	Theoretical Analyses of Degree Bias	219
C.1.2	Hypotheses for Degree Bias	219
C.2	Proofs	221
C.2.1	Theorem 4.4.1	221
C.2.2	Theorem 4.4.2	221

C.2.3	Theorem 4.4.3	222
C.2.4	Lemma 4.5.1	223
C.2.5	Theorem 4.5.2	224
C.2.6	Theorem C.8.1	225
C.3	Datasets	226
C.4	Models	227
C.5	Additional Degree Bias Plots	228
C.6	Additional Visual Summaries of Theoretical Results	231
C.7	Additional Inverse Collision Probability Plots	239
C.8	Training-Time Degree Bias: Random Walk Graph Filter	241
C.9	Achieving Maximum Training Accuracy	241
C.10	Additional Training Loss Plots	243
C.11	Limitations and Future Directions	246
D	Appendix for Chapter 5	249
D.1	Limitations	249
D.2	Formal Details About Probability-Based and Generation-Based Evaluations	250
D.2.1	Generation-Based Evaluation	250
D.2.2	Probability-Based Evaluation	250
D.3	Experimental Details	251
D.3.1	MISGENDERED	251
D.3.2	RUFF	252
D.3.3	TANGO	252
D.3.4	Practical Challenges	253

D.3.5	Agreement Metrics	254
D.4	Theoretical Analysis of Divergence Between Evaluation Formats	256
D.4.1	Converting from Probability-Based to Generation-Based Evaluation	256
D.4.2	Converting from Generation-Based to Probability-Based Evaluation	257
D.5	Additional Experimental Results	258
D.5.1	MISGENDERED	258
D.5.2	TANGO	263
D.5.3	RUFF	266
D.6	Human Annotation Guidelines	271
D.6.1	Pronoun Annotation	271
D.6.2	Extraneous Gendered Language	271
D.6.3	Other Notes and Peculiarities	272
D.7	Qualitative Examples	273
D.7.1	Qualitative Examples of Human Disagreement with Generation-Based Evaluation Results	273
D.7.2	Qualitative Examples of Extraneous Gendered Language in Generations	274
D.8	Measuring Repetition	276
E	Appendix for Chapter 6	280
E.1	Technical Assumptions	280
E.2	Warm-Up: Classical Linear Model	281
E.2.1	Single Model Learned for Both Groups	281
E.2.2	Separate Model Learned Per Group	282
E.2.3	Phase Diagram	283

E.3	Proof of Theorem E.2.2	284
	E.3.1 Variance Term	285
	E.3.2 Bias Term	288
E.4	Proof of Theorem E.2.1	291
	E.4.1 Variance Terms	291
	E.4.2 Bias Terms	295
E.5	Proof of Theorem 6.3.1	302
	E.5.1 Computing $\mathbb{E}\bar{\text{tr}} r_j^{(1)}$	305
	E.5.2 Computing $\mathbb{E}\bar{\text{tr}} r_j^{(2)}$	306
	E.5.3 Computing $\mathbb{E}\bar{\text{tr}} r_j^{(3)}$	309
	E.5.4 Computing $\mathbb{E}\bar{\text{tr}} r_j^{(4)}$	310
E.6	Theorem 6.3.2	312
E.7	Solving Fixed-Point Equations for Theorem E.2.1	313
	E.7.1 Proportional Covariance Matrices	313
E.8	Corollary E.8.1	314
	E.8.1 Case 1: $\theta_0 = 0$	316
	E.8.2 Case 2: $\theta_0 > 0$	316
E.9	Experimental Details	319
	E.9.1 Synthetic Experiments	319
	E.9.2 Colored MNIST Experiments	321
E.10	Bias Amplification Plots	322
E.11	Power-Law Covariance	326
E.12	Proof of Corollary E.11.1	328

E.13 Bias Amplification During Training	329
E.14 Colored MNIST Plots	331
E.15 Minority-Group Bias Plots	334
E.16 Actionable Insights from Theory	337

LIST OF FIGURES

1.1	LLMs process syntactic and semantic context. We visualize the layer-9 attention of <code>bert-base-uncased</code> for two sentences that differ only in the pronoun (<code>they</code> or <code>she</code>) referring to <code>Mary</code> [Vig19]. The pronoun token attends more strongly to <code>Mary</code> when the pronoun is <code>she</code> , suggesting that LLMs capture unreliable gendered associations between pronouns and names [GSL24].	3
1.2	GNNs process social context. GNNs used by a company to filter job applicants may reject qualified women applicants while accepting less qualified men applicants due to their proximity to current company employees in social networks like LinkedIn [BLM14]. Even though network location can be a proxy for social category (e.g., gender), as of the writing of [BLM14], it is not illegal in the U.S. for companies to discriminate based on personal network.	4
2.1	A job applicant network. After applying Feature Propagation [RKG22], the nodes in Q will have feature values that are more distinct from the feature values of the nodes in R than the ground truth.	12
3.1	An academic collaboration network where nodes are Computer Science (CS) and Education (EDU) researchers, solid edges are current or past collaborations, and dashed edges are collaborations recommended by a GCN. Circular nodes are women and square nodes are men.	33
3.2	The plots display $\widehat{\Delta}^{(b)}$ vs. $\Delta^{(b)}$ for Φ_s for the NBA, German, and DBLP-Fairness datasets over all $b \in [B]$ and 10 random seeds. Each point corresponds to a different random seed, and the color of the point corresponds to the social group $S^{(b)}$. We compute $\widehat{\Delta}^{(b)}$ and $\Delta^{(b)}$ post-sigmoid using only the LP scores over the sampled (positive and negative) test edges. The plots display the NRMSE and PCC of $\widehat{\Delta}^{(b)}$ as a predictor of $\Delta^{(b)}$	45

3.3	The plots display the theoretic vs. GCN LP scores for the Cora, CS, and LastFMAsia datasets over 10 random seeds. (We include the plots for the remaining datasets in Appendix B.6.) The top row of plots corresponds to Φ_s , the bottom row to Φ_r . In the plots, each circle corresponds to a single pair of test nodes (between which we are predicting a link). The center of each circle represents the mean of the theoretic and GCN scores and its area captures the range of scores. The color of each circle indicates the social group to which the node pair belongs. The plots include: (1) the total number of test node pairs N ; (2) the number of social groups B ; (3) the dashed line of equality for easy comparison of the theoretic and GCN scores. For all the datasets, the tables display: (1) the mean/standard deviation of the GCN test AUC on LP; and (2) the mean/standard deviation of the range-normalized ¹ root-mean-square deviation (NRMSE) [Ott19] and Pearson correlation coefficient (PCC) of the theoretic LP scores as predictors of the GCN scores. The left table corresponds to Φ_s , the right to Φ_r	50
4.1	Test loss vs. degree of nodes in CiteSeer for RW, SYM, and ATT GNNs. High-degree nodes generally incur a lower test loss than low-degree nodes do. Error bars are reported over 10 random seeds; all error bars are 1-sigma and represent the standard deviation about the mean.	53
4.2	Inverse collision probability vs. degree of nodes in CiteSeer for RW, SYM, and ATT GNNs. Node degrees generally have a strong association with inverse collision probabilities.	61

4.3	Visual summary of the geometry of representations, variance of representations, and training dynamics of RW, SYM, and ATT GNNs on CiteSeer. We consider low-degree nodes to be the 100 nodes with the smallest degrees and high-degree nodes to be the 100 nodes with the largest degrees. Each point in the plots in the left column corresponds to a test node representation and its color represents the node’s class. (In this particular dataset, low-degree nodes are more heavily concentrated in a few classes.) The plots in the left column are based on a single random seed, while the plots in the middle and right columns are based on 10 random seeds. RW representations of low-degree nodes often have a larger variance than high-degree node representations, while SYM representations of low-degree nodes often have a smaller variance. Furthermore, SYM generally adjusts its training loss on low-degree nodes less rapidly.	70
5.1	Overview of existing datasets for measuring LLM misgendering, with example inputs and the task. Each input surfaces a subject (e.g., name, distal antecedent, entity) and corresponding pronoun. All inputs demonstrate 1-2 uses of the correct pronoun; the correct pronoun is never ambiguous. MISGENDERED and RUFF are probability-based evaluations, while TANGO is generation-based. MISGENDERED inputs contain an explicit declaration of pronouns and personal names, while RUFF inputs contain an implicit declaration and no personal names.	73
5.2	An example context from the generation-based evaluation dataset TANGO [OGD23] and a corresponding generation by Llama-3.2-1B with misgendering. The context surfaces a subject (Jaime) and base pronoun (they). The context and generation can be converted to a template to support probability-based evaluation.	75

5.3	An example template from the probability-based evaluation dataset MISGENDERED [HDS23]. The template surfaces a subject (Reise) and base pronoun (xe). The template can be converted to pre- and post-[MASK] contexts to support generation-based evaluation.	78
5.4	Variation and agreement for MISGENDERED. (a) Generation variation σ (Eq. 5.1) for each model and pronoun in the pre-[MASK] generation setting. As we sample 5 generations, $\sigma \in \{0, 0.4, 0.49\}$. The bar labels h , s , t , x correspond to he , she , they , xe . (b) Raw observed agreement v^{p_o} (Eq. 5.2) for each model and pronoun between the probability-based and pre-[MASK] generation-based evaluation results. Error bars represent the standard error of v^{p_o} (computed over dataset instances). The horizontal dashed line is the upper bound of v^{p_o}	83
5.5	Instance-level variation σ (Eq. 5.1) for each model and pronoun with TANGO. (a) Generation-Based variation. The bar labels h , s , t , x correspond to he , she , they , xe . (b) Probability-Based variation. As we exclude templates with no pronoun, we do not always have 5 templates per instance (see Figure D.3). Hence, we report the mean and standard deviation.	84
5.6	Agreement between human and automatic evaluation of misgendering in the pre-[MASK] generation setting. Many models fall short of human-human agreement (96%).	86
5.7	Human annotations of Llama-70B and OLMo-13B generations from the pre-[MASK] (left) and post-[MASK] (right) settings.	87
5.8	Proportion of generations with extraneous gendered words in the pre-[MASK] generation setting. MISGENDERED contains named subjects with pronoun declarations, which seem to elicit more extraneous gendered cues than RUFF, which contains occupations.	88

6.1	<i>ODD, EDD, and ADD phase diagrams for ridge regression with random projections.</i>	We plot the bias amplification phase diagrams with respect to ϕ (rate of features to samples) and ψ (rate of parameters to samples), as predicted by our theory for ridge regression with random projections (Theorems 6.3.1, 6.3.2). Red regions indicate theoretical predictions greater than 1 (i.e., bias amplification in the rightmost plot), while blue regions indicate theoretical predictions less than 1 (i.e., bias deamplification in the rightmost plot). Darkness indicates intensity. We consider isotropic covariance matrices: $\Sigma_1 = 2I_d, \Sigma_2 = I_d, \Theta = 2I_d, \Delta = I_d$. Additionally, $n = 1 \times 10^4, \sigma_1^2 = \sigma_2^2 = 1$. We further choose $\lambda = \lambda_1 = \lambda_2 = 1 \times 10^{-6}$ to approximate the minimum-norm interpolator. We show that bias amplification can occur even in the balanced data setting, i.e., when $p_1 = p_2 = 1/2$	104
6.2	<i>Our theory predicts that models can amplify bias even with balanced groups and without spurious correlations.</i>	We empirically validate our theory (Theorems 6.3.1 and 6.3.2) for <i>ODD, EDD</i> , and <i>ADD</i> under the setup described in Section 6.4.1, with $a_1 = 0.5, a_2 = 1, \sigma_1^2 = 1$, and $\sigma_2^2 = 1 \times 10^{-5}$. The solid lines capture empirical values while the corresponding lower-opacity dashed lines represent what our theory predicts. We plot <i>ODD</i> and <i>EDD</i> on the same scale for easy comparison, and include a black dashed line at $ADD = 1$ to contrast bias amplification vs. deamplification. We include all the plots with error bars in Appendix E.10.	107
6.3	<i>Our theory predicts that disparate label noise between groups deamplifies bias on Colored MNIST.</i>	We plot the <i>ODD</i> and <i>EDD</i> of a CNN over training time t for Colored MNIST. As t increases, the <i>ODD</i> is relatively low while the <i>EDD</i> is noticeably higher. The error bars capture the standard deviation computed over 10 random seeds.	109

6.4	Minority-group test risk can peak with different model sizes depending on the rate of features to samples. We empirically demonstrate that minority-group bias is affected by extraneous features. We validate our theory (Theorems 6.3.1 and 6.3.2) for together R_1, R_2 (i.e., single model learned for both groups) and separate R_1, R_2 (i.e., separate model learned per group) under the setup described in Section 6.4.2, with $a_1 = 2, b_2 = 0.2$, and $\pi = 0.5$. The solid lines capture empirical values while the corresponding lower-opacity dashed lines represent what our theory predicts. We include a black dashed line at $ADD = 1$ to contrast bias amplification vs. deamplification. All y-axes are on the same scale for easy comparison. All the plots with error bars are in Appendix E.15.	111
7.1	All the stages of the ML development lifecycle have an interdependent impact on model fairness.	116
7.2	In consecutive rounds of training, changes in the decision boundary of a GNN can cause certain test instances (indicated by arrows) to receive unstable predictions. Each point in the figure represents a test instance and its color corresponds to its class.	117
A.1	Heatmap of discrimination risks and maximum α (over all channels) of Feature Propagation for 36 combinations of unknown feature rates for each group in SBM. We use 0.5 relative group sizes and 0.5 inter- and intra-link rates.	183
A.2	Plots of discrimination risk and maximum α (over all channels) of Feature Propagation vs. relative size of group Q in SBM. We use 0.5 unknown feature rates for both groups and 0.5 inter- and intra-link rates.	184
A.3	Plots of discrimination risk and maximum α (over all channels) of Feature Propagation vs. ratio of inter-link rate to intra-link rate in SBM. We use 0.5 relative group sizes and 0.5 unknown feature rates for both groups.	184

A.4	Heatmap of discrimination risks and maximum α (over all channels) of Graph Regularization for 36 combinations of unknown feature rates for each group in SBM. We use 0.5 relative group sizes and 0.5 inter- and intra-link rates.	185
A.5	Plots of discrimination risk and maximum α (over all channels) of Graph Regularization vs. relative size of group Q in SBM. We use 0.5 unknown feature rates for both groups and 0.5 inter- and intra-link rates.	186
A.6	Plots of discrimination risk and maximum α (over all channels) of Graph Regularization vs. ratio of inter-link rate to intra-link rate in SBM. We use 0.5 relative group sizes and 0.5 unknown feature rates for both groups.	186
A.7	Heatmap of discrimination risks and maximum α (over all channels) of Neighbor Mean for 36 combinations of unknown feature rates for each group in SBM. We use 0.5 relative group sizes and 0.5 inter- and intra-link rates.	187
A.8	Plots of discrimination risk and maximum α (over all channels) of Neighbor Mean vs. relative size of group Q in SBM. We use 0.5 unknown feature rates for both groups and 0.5 inter- and intra-link rates.	188
A.9	Plots of discrimination risk and maximum α (over all channels) of Neighbor Mean vs. ratio of inter-link rate to intra-link rate in SBM. We use 0.5 relative group sizes and 0.5 unknown feature rates for both groups.	188
A.10	Heatmap of discrimination risks and maximum α (over all channels) of Global Mean for 36 combinations of unknown feature rates for each group in SBM. We use 0.5 relative group sizes. Note: Global Mean is not affected by graph structure.	189
A.11	Plots of discrimination risk and maximum α (over all channels) of Global Mean vs. relative size of group Q in SBM. We use 0.5 unknown feature rates for both groups. Note: Global Mean is not affected by graph structure.	190
B.1	Theoretic vs. GCN LP scores for collaboration network datasets.	207

B.2	Theoretic vs. GCN LP scores for citation network datasets.	208
B.3	Theoretic vs. GCN LP scores for online social network datasets.	209
B.4	Theoretic LP score vs. 4-layer Φ_s LP score for all network datasets.	211
B.5	Theoretic LP score vs. Φ_s LP score (with Hadamard product and MLP) for all network datasets.	214
B.6	The plots display $\widehat{\Delta}^{(b)}$ vs. $\Delta^{(b)}$ for 4-layer Φ_s for the NBA, German, and DBLP-Fairness datasets over all $b \in [B]$ and 10 random seeds.	215
B.7	Associations of absolute deviation with degree product and with feature similarity for CiteSeer.	216
B.8	Weak associations of max term with NRMSE and PCC of theoretic LP scores for Φ_r across all datasets described in §B.3.	217
B.9	Weak associations of mean Φ_r LP scores (over 10 random seeds) with degree of each incident node and product of degrees of both incident nodes. Colors correspond to different groups.	217
C.1	Test loss vs. degree of nodes in citation and collaboration network datasets for RW, SYM, and ATT GNNs. High-degree nodes generally incur a lower test loss than low-degree nodes do. Error bars are reported over 10 random seeds; all error bars are 1-sigma and represent the standard deviation about the mean.	229
C.2	Test loss vs. degree of nodes in online product and Wikipedia network datasets for RW, SYM, and ATT GNNs. High-degree nodes generally incur a lower test loss than low-degree nodes do. Error bars are reported over 10 random seeds; all error bars are 1-sigma and represent the standard deviation about the mean.	230
C.3	Visual summary of the geometry of representations, variance of representations, and training dynamics of RW, SYM, and ATT GNNs on Cora_ML.	232

C.4	Visual summary of the geometry of representations, variance of representations, and training dynamics of RW, SYM, and ATT GNNs on CS.	233
C.5	Visual summary of the geometry of representations, variance of representations, and training dynamics of RW, SYM, and ATT GNNs on Physics.	234
C.6	Visual summary of the geometry of representations, variance of representations, and training dynamics of RW, SYM, and ATT GNNs on Amazon Photo.	235
C.7	Visual summary of the geometry of representations, variance of representations, and training dynamics of RW, SYM, and ATT GNNs on Amazon Computers.	236
C.8	Visual summary of the geometry of representations, variance of representations, and training dynamics of RW, SYM, and ATT GNNs on chameleon.	237
C.9	Visual summary of the geometry of representations, variance of representations, and training dynamics of RW, SYM, and ATT GNNs on chameleon.	238
C.10	Inverse collision probability vs. degree of nodes in citation and collaboration network datasets for RW, SYM, and ATT GNNs. Node degrees generally have a strong association with inverse collision probabilities.	239
C.11	Inverse collision probability vs. degree of nodes in citation and collaboration network datasets for RW, SYM, and ATT GNNs. Node degrees generally have a strong association with inverse collision probabilities.	240
C.12	Mean absolute parameter gradient vs. training epoch for RW, SYM, and ATT GNNs on CiteSeer (over 10 random seeds). The training accuracy of SYM, RW, and ATT ultimately reach the accuracy of MAJ _{WL}	242
C.13	Mean absolute parameter gradient vs. training epoch for RW, SYM, and ATT GNNs on Photo and Computers. The training accuracy of SYM, RW, and ATT ultimately reach the accuracy of MAJ _{WL}	243

C.14	Mean absolute parameter gradient vs. training epoch for RW, SYM, and ATT GNNs on Cora_ML, CS, and Physics. The training accuracy of SYM, RW, and ATT ultimately reach the accuracy of MAJ _{WL}	244
C.15	Mean absolute parameter gradient vs. training epoch for RW, SYM, and ATT GNNs on chameleon and squirrel. The training accuracy of SYM, RW, and ATT ultimately reach the accuracy of MAJ _{WL}	245
D.1	(a) Generation variation σ (Eq. 5.1) for each model and pronoun in the post-[MASK] generation setting for MISGENDERED. Because we sample five generations per context, $\sigma \in \{0, 0.4, 0.49\}$. The bar labels h , s , t , x correspond to he , she , they , xe . (b) Raw observed agreement v^{p_o} (Eq. 5.2) for each model and pronoun between the probability-based and post-[MASK] generation-based evaluation results for MISGENDERED. The error bars represent the standard error of v^{p_o} (computed over dataset instances). The horizontal dashed line represents the upper bound of v^{p_o}	259
D.2	Disagreement (Eq. D.6) across all models and pronouns of the probability-based and pre and post-[MASK] generation-based evaluation results for MISGENDERED. Each point represents a latent beta distribution that models the probability of disagreement in results for a single model (marker color) and pronoun (marker shape). The dashed lines capture the critical values $\alpha = 1, \beta = 1, \alpha = \beta$	262
D.3	Mean rate (across the five generations per instance) at which TANGO generations lack pronouns (i.e., templates fail to be constructed for Prob-TANGO) for each model and pronoun. The error bars represent the standard error (computed over dataset instances). The horizontal dashed line represents the lower bound of the failure rate.	263

D.4	Raw observed agreement v^{p_o} (Eq. 5.3) for each model and pronoun between the probability- and generation-based evaluation results for TANGO. The error bars represent the standard error of v^{p_o} (computed over dataset instances). The horizontal dashed line represents the upper bound of v^{p_o}	264
D.5	Disagreement (Eq. D.6) across all models and pronouns of the probability- and generation-based evaluation results for TANGO. Each point represents a latent beta distribution that models the probability of disagreement in results for a single model (marker color) and pronoun (marker shape). The dashed lines capture the critical values $\alpha = 1, \beta = 1, \alpha = \beta$	265
D.6	Generation variation σ (Eq. 5.1) for each model and pronoun in the pre and post-[MASK] generation settings for RUFF. Because we sample five generations per context, $\sigma \in \{0, 0.4, 0.49\}$. The bar labels h , s , t , x correspond to he , she , they , xe	266
D.7	Raw observed agreement v^{p_o} (Eq. 5.2) for each model and pronoun between the probability-based and pre and post-[MASK] generation-based evaluation results for RUFF. The error bars represent the standard error of v^{p_o} (computed over dataset instances). The horizontal dashed line represents the upper bound of v^{p_o}	267
D.8	Disagreement (Eq. D.6) across all models and pronouns of the probability- and generation-based evaluation results for RUFF. Each point represents a latent beta distribution that models the probability of disagreement in results for a single model (marker color) and pronoun (marker shape). The dashed lines capture the critical values $\alpha = 1, \beta = 1, \alpha = \beta$	270
D.9	Agreement between human and automatic evaluation of misgendering in the post-[MASK] generation setting. Many models fall short of human-human agreement (96%).	273

E.1	<i>ODD</i>, <i>EDD</i>, and <i>ADD</i> phase diagrams for classical ridge regression.	We plot the bias amplification phase diagrams with respect to ϕ (rate of features to samples), as predicted by our theory for ridge regression without random projections (Theorems E.2.1, E.2.2). Dashed black lines indicate theoretical predictions. We consider isotropic covariance matrices: $\Sigma_1 = 2I_d, \Sigma_2 = I_d, \Theta = 2I_d, \Delta = I_d$. Additionally, $n = 1 \times 10^4, \sigma_1^2 = \sigma_2^2 = 1$. We further choose $\lambda = \lambda_1 = \lambda_2 = 1 \times 10^{-6}$ to approximate the minimum-norm interpolator. We observe that bias amplification can occur even in the balanced data setting, i.e., when $p_1 = p_2 = 1/2$, without spurious correlations.	284
E.2	We empirically demonstrate that bias amplification occurs and validate our theory (Theorems 6.3.1 and 6.3.2) for <i>ODD</i> , <i>EDD</i> , and <i>ADD</i> under the setup described in Section 6.4.1. The solid lines capture empirical values while the corresponding lower-opacity dashed lines represent what our theory predicts. We plot <i>ODD</i> and <i>EDD</i> on the same scale for easy comparison, and include a black dashed line at <i>ADD</i> = 1 to contrast bias amplification vs. deamplification. The error bars capture the range of the estimators over 25 random seeds.		323
E.3	We empirically demonstrate that bias amplification occurs and validate our theory (Theorems 6.3.1 and 6.3.2) for <i>ODD</i> , <i>EDD</i> , and <i>ADD</i> under the setup described in Section 6.4.1. The solid lines capture empirical values while the corresponding lower-opacity dashed lines represent what our theory predicts. We plot <i>ODD</i> and <i>EDD</i> on the same scale for easy comparison, and include a black dashed line at <i>ADD</i> = 1 to contrast bias amplification vs. deamplification. The error bars capture the range of the estimators over 25 random seeds.		324

E.4	We empirically demonstrate that bias amplification occurs and validate our theory (Theorems 6.3.1 and 6.3.2) for <i>ODD</i> , <i>EDD</i> , and <i>ADD</i> under the setup described in Section 6.4.1. The solid lines capture empirical values while the corresponding lower-opacity dashed lines represent what our theory predicts. We plot <i>ODD</i> and <i>EDD</i> on the same scale for easy comparison, and include a black dashed line at $ADD = 1$ to contrast bias amplification vs. deamplification. The error bars capture the range of the estimators over 25 random seeds.	325
E.5	Our theory predicts that bias amplification is larger for higher noise ratios than lower noise ratios. We observe that Corollary E.11.1 generally predicts the <i>ADD</i> profile with respect to the noise ratio c . The solid lines capture empirical values while the corresponding lower-opacity dashed lines represent what Theorem 6.3.1 predicts. We plot <i>ODD</i> and <i>EDD</i> on the same scale for easy comparison, and include a black dashed line at $ADD = 1$ to contrast bias amplification vs. deamplification. The error bars capture the range of the estimators over 25 random seeds. We consider the setup described in Appendix E.11 with $\psi = 0.5$, $\phi = 0.2$, and $\lambda = 1 \times 10^{-6}$	327
E.6	Our theory reveals that there may be an optimal regularization penalty to deamplify bias. We empirically demonstrate that bias amplification can be heavily affected by λ and validate our theory (Theorems 6.3.1 and 6.3.2) for <i>ODD</i> , <i>EDD</i> , and <i>ADD</i> under the setup described in Section 6.4.2. The solid lines capture empirical values while the corresponding lower-opacity dashed lines represent what our theory predicts. We include a black dashed line at $ADD = 1$ to contrast bias amplification vs. deamplification. The error bars capture the range of the estimators over 25 random seeds.	330

- E.7 **Our theory predicts that more disparate label noise between groups deamplifies bias on Colored MNIST.** We plot the *ODD* and *EDD* of a CNN for different label noise ratios $c = \sigma_2^2/\sigma_1^2$ for Colored MNIST. As c increases, the *EDD* generally increases while the *ODD* remains relatively low, which is predicted by our theory (see reasoning in Section 6.4.2). In our experiments, $\sigma_1^2 = 0.05$ stays fixed while σ_2^2 varies. For each value of c , the model is evaluated after $t = 80$ training steps and has a penultimate layer with dimension $m = 500$. The error bars capture the standard deviation computed over 10 random seeds. 332
- E.8 **Our theory predicts that a larger model size reduces bias on Colored MNIST in the single model setting.** We plot the *ODD* and *EDD* of a CNN for different model sizes m (where m is the dimension of the penultimate CNN layer) for Colored MNIST. As m increases, the *ODD* appears to decrease and plateau, which is in line with what our theory predicts in the regime where $\phi < 1$ (see analysis in Section 6.4.1). The *EDD* does not tend towards 0. In our experiments, $\sigma_1^2 = \sigma_2^2 = 0.05$. For each value of m , the model is evaluated after $t = 80$ training steps. The error bars capture the standard deviation computed over 10 random seeds. 333
- E.9 We empirically demonstrate that minority-group bias is affected by extraneous features. We validate our theory (Theorems 6.3.1 and 6.3.2) for together R_1, R_2 (i.e., single model learned for both groups) and separate R_1, R_2 (i.e., separate model learned per group) under the setup described in Section 6.4.2. The solid lines capture empirical values while the corresponding lower-opacity dashed lines represent what our theory predicts. We include a black dashed line at $ADD = 1$ to contrast bias amplification vs. deamplification. All y-axes are on the same scale for easy comparison. The error bars capture the range of the estimators over 25 random seeds. 334

E.10 We empirically demonstrate that minority-group bias is affected by extraneous features. We validate our theory (Theorems 6.3.1 and 6.3.2) for together R_1, R_2 (i.e., single model learned for both groups) and separate R_1, R_2 (i.e., separate model learned per group) under the setup described in Section 6.4.2. The solid lines capture empirical values while the corresponding lower-opacity dashed lines represent what our theory predicts. We include a black dashed line at $ADD = 1$ to contrast bias amplification vs. deamplification. All y-axes are on the same scale for easy comparison. The error bars capture the range of the estimators over 25 random seeds. 335

E.11 We empirically demonstrate that minority-group bias is affected by extraneous features. We validate our theory (Theorems 6.3.1 and 6.3.2) for together R_1, R_2 (i.e., single model learned for both groups) and separate R_1, R_2 (i.e., separate model learned per group) under the setup described in Section 6.4.2. The solid lines capture empirical values while the corresponding lower-opacity dashed lines represent what our theory predicts. We include a black dashed line at $ADD = 1$ to contrast bias amplification vs. deamplification. All y-axes are on the same scale for easy comparison. The error bars capture the range of the estimators over 25 random seeds. 336

LIST OF TABLES

2.1	Reconstruction error (RE), discrimination risk (DR), and test group membership identification accuracy (MI) of all models averaged over all 25 combinations of unknown feature rates of $\{0.1, 0.3, 0.5, 0.7, 0.9\}$ for each group in SBM. We use 0.5 relative group sizes and 0.5 inter- and intra-link rates.	29
2.2	Test accuracy (Acc) and statistical parity (SP) of all models averaged over all 25 combinations of unknown feature rates of $\{0.1, 0.3, 0.5, 0.7, 0.9\}$ for each group in German credit.	30
3.1	$\frac{1}{B} \sum_{b \in [B]} \Delta^{(b)}$ and the test AUC for the NBA, German, and DBLP-Fairness datasets with various settings of λ_{fair} . The left table corresponds to Φ_s , and the right to Φ_r	47
4.1	Five most popular hypotheses for the origins of degree bias proposed by papers. The remaining hypotheses can be found in Table C.2 in the appendix.	55
5.1	<i>MCC</i> agreement v^{MCC} (Eq. 5.2) between probability-based and pre-[MASK] generation-based evaluations, for each model and pronoun in MISGENDERED. We report the asymmetric 95% confidence interval, computed using SciPy [VGO20], except with xe and Mixtral-8x22B, as the model gets every instance correct in the probability-based setting.	82
5.2	<i>MCC</i> agreement v^{MCC} (Eq. 5.2) between probability- and generation-based evaluation for each model and pronoun in TANGO. We report the asymmetric 95% confidence interval, computed using SciPy [VGO20].	85

A.1	Reconstruction error (RE), discrimination risk (DR), and test group membership identification accuracy (MI) of all models averaged over relative sizes of group Q of $\{0.1, 0.3, 0.5, 0.7, 0.9\}$ in SBM . We use 0.5 unknown feature rates for both groups and 0.5 inter- and intra-link rates.	191
A.2	Reconstruction error (RE), discrimination risk (DR), and test group membership identification accuracy (MI) of all models averaged over all 25 combinations of inter- and intra-link rates of $\{0.1, 0.3, 0.5, 0.7, 0.9\}$ in SBM . We use 0.5 relative group sizes and 0.5 unknown feature rates for both groups.	192
A.3	equal opportunity (EO) averaged over all 25 combinations of unknown feature rates of $\{0.1, 0.3, 0.5, 0.7, 0.9\}$ for each group in German credit	193
A.4	Test accuracy (Acc) and statistical parity (SP) of all models averaged over all 25 combinations of unknown feature rates of $\{0.1, 0.3, 0.5, 0.7, 0.9\}$ for each group in Credit defaulter	194
A.5	equal opportunity (EO) of all models averaged over all 25 combinations of unknown feature rates of $\{0.1, 0.3, 0.5, 0.7, 0.9\}$ for each group in Credit defaulter	195
B.1	Summary of the datasets used in our experiments.	204
B.2	The test AUC of the 4-layer Φ_s encoders on the real-world network datasets, and the NRMSE and PCC of the theoretic LP scores as predictors of the Φ_s scores.	210
B.3	The test AUC of the Φ_s encoders with an f_{MLP} score function on the real-world network datasets, and the NRMSE and PCC of the theoretic LP scores as predictors of the Φ_s scores.	213
B.4	$\frac{1}{B} \sum_{b \in [B]} \Delta^{(b)}$ and the test AUC for the NBA, German, and DBLP-Fairness datasets with various settings of λ_{fair} . The left table corresponds to 4-layer Φ_s , and the right to 4-layer Φ_r	215

C.1	A taxonomy of GNN degree bias papers based on whether they theoretically analyze the origins of degree bias, explicitly linking a node’s degree to its test and training error.	219
C.2	Full taxonomy of the hypotheses for the origins of GNN degree bias proposed by papers.	220
C.3	Summary of the datasets used in our experiments.	227
D.1	<i>MCC</i> agreement v^{MCC} (Eq. 5.2) for each model and pronoun between the probability-based and post-[MASK] generation-based evaluation results for MISGENDERED. We report the asymmetric 95% confidence interval, computed using SciPy [VGO20], except for <code>xe</code> with Mixtral-8x22B, as the model gets every instance correct in the probability-based setting.	260
D.2	κ agreement v^κ (Eq. 5.2) for each model and pronoun between the probability-based and pre and post-[MASK] generation-based evaluation results for MISGENDERED. We report the 95% confidence interval, computed using <code>statsmodels</code> [SP10].	261
D.3	κ agreement v^κ (Eq. 5.3) for each model and pronoun between the probability- and generation-based evaluation results for TANGO. We report the 95% confidence interval, computed using <code>statsmodels</code> [SP10].	264
D.4	<i>MCC</i> agreement v^{MCC} (Eq. 5.2) for each model and pronoun between the probability-based and pre and post-[MASK] generation-based evaluation results for RUFF. We report the asymmetric 95% confidence interval, computed using <code>statsmodels</code> [SP10].	268
D.5	κ agreement v^κ (Eq. 5.2) for each model and pronoun between the probability-based and pre and post-[MASK] generation-based evaluation results for RUFF. The interval represents the 95% confidence interval, computed using <code>statsmodels</code> [SP10].	269

D.6	Examples of human disagreement with the results of automatic evaluation of pre-[MASK] generations for MISGENDERED and RUFF. The bolded text represents the context while the unbolded text represents the generation. The underlined text indicates the first pronoun in the generation (i.e., the pronoun considered in automatic misgendering evaluation), and the italicized text indicates instances of misgendering of the subject.	274
D.7	Examples of extraneous gendered language in generations for MISGENDERED and RUFF. All generations are pre-[MASK] unless otherwise specified. Bolded text represents the context while the unbolded text represents the generation. The underlined text indicates extraneous gendered terms, and the italicized text indicates misgendering of the subject.	275
D.8	Repetition rate (mean \pm standard deviation) of pre-[MASK] generations for Gen-MISGENDERED across different models and pronouns.	277
D.9	Repetition rate (mean \pm standard deviation) of post-[MASK] generations for Gen-MISGENDERED across different models and pronouns.	277
D.10	Repetition rate (mean \pm standard deviation) of generations for TANGO across different models and pronouns.	278
D.11	Repetition rate (mean \pm standard deviation) of pre-[MASK] generations for Gen-RUFF across different models and pronouns.	279
D.12	Repetition rate (mean \pm standard deviation) of post-[MASK] generations for Gen-RUFF across different models and pronouns.	279

ACKNOWLEDGMENTS

This dissertation would not have been possible without the contributions and support of innumerable incredible individuals. I am immensely grateful to Yizhou Sun for her caring and thoughtful mentorship and contagious passion for research and equity. She has provided me with the flexibility and guidance to learn more about and study topics that excite me. Moreover, by example, she has encouraged me to ask and think about important questions, and to prioritize people over projects. Furthermore, I am extremely grateful to Kai-Wei Chang for his generous mentorship and encouragement. His advice and holistic approach to science have greatly impacted my research. I have worked with both Yizhou Sun and Kai-Wei Chang since I was an undergraduate student at UCLA, and I owe so much to them for their unwavering support and kindness. Additionally, I would like to thank all my committee members: Aditya Grover, Baharan Mirzasoleiman, Kai-Wei Chang, Levent Sagun, and Yizhou Sun, for their helpful feedback and guidance throughout the PhD. Their research has inspired and continues to be incredibly influential on my interests. I would also like to thank all my amazing co-authors who made the research in this dissertation possible: Dietrich Klakow, Elvis Dohmatob, Jian Kang, Kai-Wei Chang, Levent Sagun, Preethi Seshadri, Samuel Bell, Vagrant Gautam, Yizhou Sun. Thanks to the Eugene V. Cota-Robles Fellowship, NSF NRT MENTOR Fellowship, Amazon Science Hub Fellowship, and Meta FAIR for supporting my research.

Thanks to everyone in the UCLA NLP fairness subgroup, which has been an invaluable community during my PhD: Ashima Suvarna, Christina Chance, Elaine Wan, Elia Ovalle, and Rebecca Pattichis. Their brilliant research and commitment to advancing justice are inspiring. Thanks to Shichang Zhang, Sunipa Dev, and Ziniu Hu for hand-holding me through my first research projects as an undergraduate student and fundamentally shaping my research skills. I would also like to thank my other UCLA collaborators and labmates: Amita Kamath, Ana Brendel, Da Yin, Daniel Israel, Derek Xu, Di Wu, Hritik Bansal, Jieyu Zhao, Junheng

Hao, Kareem Ahmed, Masoud Monajatipoor, Oliver Broadrick, Poorva Garg, Renato Geh, Tanmay Parekh, Yanqiao Zhu, Yunsheng Bai, Zongyu Lin, Zongyue Qin, and other people in Yizhou Sun, Kai-Wei Chang, Wei Wang, and Guy Van den Broeck’s labs. They are all so smart, inspiring, and kind. I am further grateful to my undergraduate and PhD mentees, from whom I have learned so much: Margaret Capetz, Naisha Agarwal, Pranav Subbaraman, and Steven Swee. Thanks to all the people in UCLA ACM AI, ACM Teach LA, ACM-W, and GWICS for being wonderful company and growing my interest in the intersection of AI and equity! Thanks to Helen Tran, Madelen Hem, and Joseph Brown for always patiently and promptly answering all my logistical questions.

Outside of UCLA, I have had the fortune to collaborate on research projects with numerous brilliant and caring individuals. I would like to thank Levent Sagun for their ever-supportive, hands-on, and person-first mentorship and advice, and for being such a positive driving force in responsible AI. I would also like to thank Jian Kang and Elvis Dohmatob for their hands-on and encouraging mentorship; owing to their detailed feedback and advice, I have grown significantly as a researcher. In addition, QAI has been a large and supportive community that has been invaluable during my PhD. So much love for all these people (you know who you are)!

I have been lucky to have a large number of informal mentors (and friends) during my PhD who went out of their way to give me advice, support my research and career, and bring me into various communities: Angelina Wang, Irene Solaiman, Jiahao Chen, Konstantina Palla, Leif Hancox-Li, Luca Soldaini, Maria Antoniak, Merve Gürel, Pranav Agrawal, Rida Qadri, Sara Beery, Sunipa Dev, Swabha Swayamdipta, William Agnew, and Zeerak Talat. I also owe so much to my peer mentors: Alissa Valentine, Ashwin Singh, Chantal Shaib, Ellin Zhao, Evani Radiya-Dixit, Hetvi Jethwani, Jennifer Chien, Kush Jain, Lucy Li, Maria Ryskina, Michelle Lin, Morris Alper, Niloofar Mireshghallah, Preethi Seshadri, Shaily Bhatt, Stephanie Milani, and Vagrant Gautam, who have spent generous amounts of time guiding me through tough situations and making me laugh. I am especially grateful to Vagrant Gautam

for always being incredibly kind, engaging in lovely research discussions, and above all, being a fantastic friend. I would additionally like to express gratitude to my unofficial CIFRE PhD labmates: Fabian Glöckle, Federico Baldassarre, Ismail Labiad, Joséphine Raugel, Juliette Decugis, Kruno Lehman, Pierre Fernandez, Pierre Orhan, Quentin Garrido, Tom Sander, Virginie Do, Wes Bouaziz, and many others, for making me feel so welcome in a new country.

I am grateful for all my friends outside academia with whom I have shared food, explored nature, and felt grounded: Esha Krishnamoorthy, Jaspreet Sandhu, Jenny Zhang, Karen Yi, Kausalya Kethu, Megha Ilango, Mounika Narayanan, Pragathi Venkatesh, Ruining Ding, Sav Bell, Sharvani Jha, Spurthi Rallapalli, and many others! Thanks a thousand times to my parents, sister, grandparents, aunts and uncles, and other family members for always being there for me during my ups and downs and providing me with the security to pursue my PhD.

VITA

Expected Ph.D. in Computer Science, University of California, Los Angeles

2021 B.S. in Computer Science, University of California, Los Angeles
Honors: Summa Cum Laude, Samueli Outstanding Bachelor of Science

SELECT PUBLICATIONS

Arjun Subramonian, Samuel J. Bell, Levent Sagun, Elvis Dohmatob. “An Effective Theory of Bias Amplification” (**ICLR 2025**)

Elvis Dohmatob, Yunzhen Feng, **Arjun Subramonian**, Julia Kempe. “Strong Model Collapse” (**ICLR 2025**)

Arjun Subramonian, Jian Kang, Yizhou Sun. “Theoretical and Empirical Insights into the Origins of Degree Bias in Graph Neural Networks” (**NeurIPS 2024**)

Arjun Subramonian*, Vagrant Gautam*, Dietrich Klakow, Zeerak Talat. “Understanding ‘Democratization’ in NLP Research” (**EMNLP 2024**, ***equal contribution**)

Shichang Zhang*, Ziniu Hu*, **Arjun Subramonian**, Yizhou Sun. “Motif-Driven Contrastive Learning of Graph Representations” (**TKDE**, ***equal contribution**)

Arjun Subramonian, Levent Sagun, Yizhou Sun. “Networked Inequality: Preferential Attachment Bias in Graph Neural Network Link Prediction” (**ICML 2024**)

Anaelia Ovalle, **Arjun Subramonian**, Vagrant Gautam, Gilbert Gee, Kai-Wei Chang. “Factoring the Matrix of Domination: A Critical Review and Reimagination of Intersectionality in AI Fairness” (**AIES 2023**)

Organizers Of QueerInAI, Nathan Dennler, Anaelia Ovalle, Ashwin Singh, Luca Soldaini, **Arjun Subramonian**, Huy Tu, William Agnew, Avijit Ghosh, Kyra Yee, Irene Font Peradejordi, Zeerak Talat, Mayra Russo, Jess de Jesus de Pinho Pinhal. “Bound by the Bounty: Collaboratively Shaping Evaluation Processes for Queer AI Harms” (**AIES 2023**)

Arjun Subramonian, Xingdi Yuan, Hal Daumé III, Su Lin Blodgett. “It Takes Two to Tango: Navigating Conceptualizations of NLP Tasks and Measurements of Performance” (**Findings of ACL 2023**)

Organizers of QueerInAI, Anaelia Ovalle, **Arjun Subramonian**, Ashwin Singh, Claas Voelcker, Danica J. Sutherland, Davide Locatelli, Eva Breznik, Filip Klubička, Hang Yuan, Hetvi J, Huan Zhang, Jaidev Shriram, Kruno Lehman, Luca Soldaini, Maarten Sap, Marc Peter Deisenroth, Maria Leonor Pacheco, Maria Ryskina, Martin Mundt, . . . “Queer In AI: A Case Study in Community-Led Participatory AI” (**FAccT 2023**)

Arjun Subramonian, Kai-Wei Chang, Yizhou Sun. “On the Discrimination Risk of Mean Aggregation Feature Imputation in Graphs” (**NeurIPS 2022**)

Sunipa Dev, Masoud Monajatipoor*, Anaelia Ovalle*, **Arjun Subramonian***, Jeff M Phillips, Kai-Wei Chang. “Harms of Gender Exclusivity and Challenges in Non-Binary Representation in Language Technologies” (**EMNLP 2021**, *equal contribution)

CHAPTER 1

Introduction

1.1 Motivation

Machine learning (ML) has proliferated rapidly over the last decade [MFP25], shaping the products that companies offer users. For example, the development of large language models has enabled companies to make powerful chat-based assistants accessible to the public [Ope, Ant], and the creation of neural filters for graphics editing software [Ado24] has allowed users to quickly transform their images without low-level edits. Simultaneously, ML is shaping how predictions about individuals are made at scale. ML powers targeted advertising [Eng22, Com23], automated hiring technologies [Aju19], and credit risk prediction systems [CTM19, Cap19].

In tandem with the proliferation of ML, there is decreasing trust in its fairness [MFP25]. Prior research has revealed innumerable fairness issues with ML, such as performance disparities across social groups and stereotyping. These issues can be attributed to: (1) marginalized social groups being underrepresented or misrepresented in the data on which models are trained, and (2) model design choices made by ML practitioners. Moreover, fairness issues cause real-world harm. For example, [BCZ16] discovered that word embeddings trained on Google News articles capture offensive gendered associations like men are to computer programmers as women are to homemakers. As such, resume screening systems used in hiring that are built atop word embeddings can discriminate against women candidates. [BG18b] found that commercial facial recognition systems misclassify darker-skinned women

in up to 34.7% of instances, compared to 0.8% for lighter-skinned men. This disparity in performance can make Black communities, and especially Black women, more susceptible to misclassification, interrogation, and harassment by law enforcement using facial recognition technology. More recently, [DMO21] and [OGD23] showed that large language models can misgender (i.e., use gender non-affirming names or pronouns) and reject transgender and non-binary individuals, psychologically harming and erasing these communities. In addition, [BKD23] and [WC24] uncovered that text-to-image models reinforce racial and gender stereotypes and power dynamics (e.g., portraying men as CEOs and women as their assistants). As ML models are increasingly deployed to make predictions about humans and their relationships at scale, it is critical to ensure that these models are fair and do not amplify social inequalities [OSG23].

In this dissertation, we study the unfairness of two modern ML models: graph neural networks (GNNs) and large language models (LLMs). Despite being trained in different manners, on distinct data modalities, these models operate on *social context and syntactic and semantic context*, respectively, which are important in today’s world. GNNs, when applied to social networks, harness the social context of individuals (e.g., neighbors, communities) to make predictions. They do so by extending traditional neural networks to leverage topological structure. Similarly, LLMs are often pretrained on massive language corpora to perform next-token prediction, by using the preceding syntactic and semantic context, e.g., identity markers of individuals (see Figure 1.1).

Graph Neural Networks. Many industrial systems and research projects alike leverage GNNs. Companies including Alibaba, Uber, and Pinterest have deployed GNNs [Ram24] to power product recommendation [ZZY19, AI19] and content discovery [YHC18] on their platforms. In addition, Twitter has applied graph learning to perform “personalized ads rankings, account follow-recommendation, offensive content detection, and search ranking” [EMP22]. Moreover, numerous GNN architectures have been proposed in the scientific

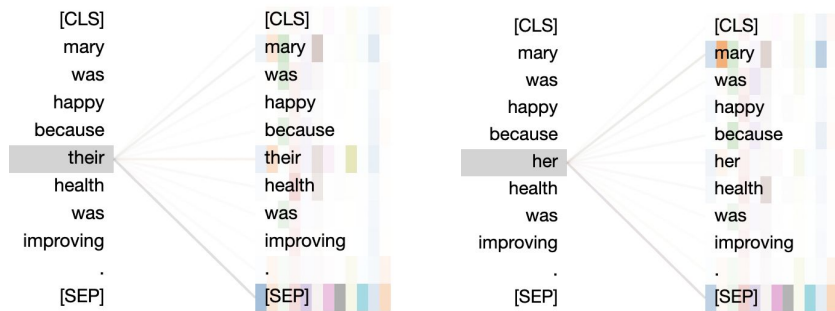


Figure 1.1: LLMs process syntactic and semantic context. We visualize the layer-9 attention of `bert-base-uncased` for two sentences that differ only in the pronoun (`they` or `she`) referring to Mary [Vig19]. The pronoun token attends more strongly to Mary when the pronoun is `she`, suggesting that LLMs capture unreliable gendered associations between pronouns and names [GSL24].

literature [SGT09, KW17, HYL17, VCC18], and “graph neural networks” has been a popular keyword at ML conferences [Ram24, Li23].

However, GNNs pose serious fairness concerns. For example, GNNs are often applied by social media companies to their user networks (where nodes are humans and edges capture social connections) to recommend new user connections. However, without concern for their diversity, such recommendations can create echo chambers [LWZ21] and unfairly stunt the connections of minoritized groups [SRC18]. In addition, companies that use GNNs to assist in hiring may unfairly favor job candidates who are already connected to current employees on social platforms like LinkedIn (see Figure 1.2) [BLM14]. Moreover, without attention to fairness, GNNs used by banks to assess creditworthiness based on prior financial connections can disproportionately reject loan applications from economically disenfranchised racial groups, amplifying inequity. Importantly, the unfairness of GNNs is exacerbated by *social context* (e.g., neighbors, communities) and their messaging passing. These issues are

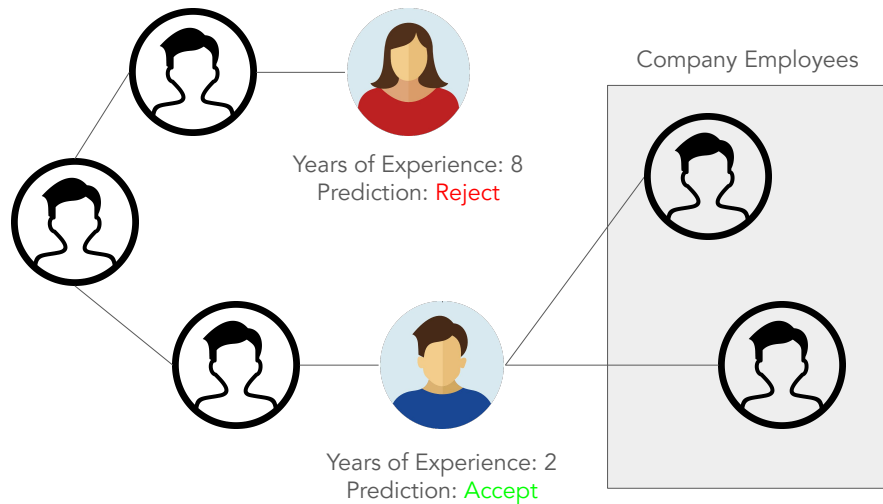


Figure 1.2: GNNs process social context. GNNs used by a company to filter job applicants may reject qualified women applicants while accepting less qualified men applicants due to their proximity to current company employees in social networks like LinkedIn [BLM14]. Even though network location can be a proxy for social category (e.g., gender), as of the writing of [BLM14], it is not illegal in the U.S. for companies to discriminate based on personal network.

not explored in research on the fairness of traditional ML and require a deeper principled understanding.

Large Language Models. LLM-powered assistants are becoming increasingly accessible to the public [Ai225, Dee]. However, LLMs pose numerous harms to marginalized communities, such as LGBTQIA+ people [Sub23]. LLMs are often trained on data that: (1) contain transphobic hate speech [Ung24]; (2) lack queer-affirmative language and representation of diverse genders and pronouns [DMO21]; and (3) are filtered of references to LGBTQIA+ identities [DSM21]. Hence, such models can regurgitate toxic language when presented with *semantic context* involving queer people, contributing to their alienation and erasure [OGD23].

1.2 Research Questions, Objectives, and Challenges

Questions. This dissertation aims to theoretically and empirically address the unfairness of GNNs and LLMs. We investigate three key research questions:

- **RQ1.** How can we develop a richer theoretical and empirical understanding of the context-based unfair behavior of GNNs and LLMs?
- **RQ2.** What technical challenges in tackling the unfairness of GNNs and LLMs are specific to these models and transcend traditional ML?
- **RQ3.** Can we develop a precise unifying theory that explains the unfairness of ML models more broadly?

Objectives and Challenges. In response to these questions, our research has the following objectives and domain-inherent challenges:

- **Graph Neural Networks:** Theoretically and empirically elucidate forms of GNN unfairness and *how* they arise from *social context* (e.g., graph structure, choice of graph filter). Propose principled metrics and methods to mitigate GNN unfairness.

In networks, nodes (unlike instances in traditional ML) are not assumed to have independently distributed data due to social phenomena like homophily (i.e., similar nodes are more likely to connect [HCT17]) and preferential attachment (i.e., highly-connected nodes are more likely to form new connections [AB02]). Additionally, unlike traditional neural networks, GNNs do not process nodes independently, instead making predictions by recursively transforming, aggregating, and combining the features of nodes and their neighbors. These differences can invalidate the applicability of previous ML fairness research, which often heavily relies on the assumption that data are independently distributed and processed [LLC24]. Furthermore, because of the interplay

of graph structure and node features in message passing, it can be difficult to dissect and measure their respective contributions to unfairness.

- **Large Language Models:** Assess the validity of evaluations of *syntactic and semantic context*-dependent LLM unfairness.

The unstructured and open-ended nature of LLM generations can make automated evaluations of LLM unfairness challenging. To overcome this challenge, template-, lexicon-, and auxiliary classifier-based evaluation methods have been proposed [GRB24], with varying levels of measurement validity. Validity issues are compounded by the general-purpose nature of LLMs, making evaluating downstream unfairness difficult.

- **Machine Learning More Broadly:** Illuminate and overcome technical challenges in developing a precise theory of *why* ML models are unfair.

Beyond GNNs and LLMs, there remain significant challenges in obtaining a precise analytical theory of how model design choices and data properties contribute to model unfairness, even in simple settings such as a multi-layer perceptron. Such challenges include deriving deterministic equivalents of complex rational expressions, without relying on bounds which can obscure unfair behavior in certain regimes. This requires powerful random matrix theory tools.

1.3 Dissertation Structure

This dissertation is organized into sections based on model type, with chapters making the following main contributions:

I. On the Unfairness of Graph Neural Networks: We study new forms of unfairness posed by GNNs that arise from graph structure and the choice of graph filter.

- **Chapter 2:** We show that applying feature imputation to graph data can amplify the

unfairness of GNNs trained on the data. To measure this phenomenon, we introduce the discrimination risk metric. We prove that a high discrimination risk can cause GNNs trained on imputed data to make disparate predictions for marginalized and dominant groups, yielding discrimination. We then theoretically and empirically show which graph properties result in imputation producing a high discrimination risk. We propose an efficient algorithm based on projected gradient descent to ensure that imputed features provably have a low discrimination risk with minimal reconstruction error. This chapter is based on [SCS22].

- **Chapter 3:** We prove that Graph Convolutional Networks (GCNs) have a preferential attachment bias, disproportionately predicting links with high-degree nodes. We further bridge GCN’s preferential attachment bias with a new form of unfairness based on disparities in link prediction scores (and thus social recommendation) between groups in social networks, and we propose a metric for this unfairness. Towards alleviating this unfairness, we develop a regularization-based training-time fairness algorithm. This chapter is based on [SSS24].
- **Chapter 4:** We prove *why* high-degree nodes tend to have a lower probability of misclassification by GNNs, showing that this bias arises from a variety of factors that are associated with a node’s degree (e.g., homophily of neighbors, diversity of neighbors). Furthermore, we show that during training, some GNNs may adjust their loss on low-degree nodes more slowly than on high-degree nodes. We connect our findings to previously-proposed hypotheses for the origins of degree bias, and we describe a principled roadmap to alleviate degree bias. This chapter is based on [SKS24].

II. On the Unfairness of Large Language Models: We investigate how the open-ended nature of LLM generations threatens the measurement validity of evaluations of unfairness.

- **Chapter 5:** We systematically assess whether probability- and generation-based

evaluations of LLM misgendering have *convergent* validity, that is, whether their results are in agreement. By automatically evaluating a suite of 6 models from 3 families, we find that these evaluation methods can disagree with each other at the instance, dataset, and model levels. Furthermore, with a human evaluation of 2400 LLM generations, we show that misgendering behavior is complex and goes far beyond pronouns, which automatic evaluations are not currently designed to capture, suggesting essential disagreement with human evaluations. We provide recommendations for future evaluations of LLM misgendering. This chapter is based on [SGS25].

III. Towards a Precise Theory of Machine Learning Unfairness: Despite the proliferation of GNNs and LLMs, we find that developing a precise theory that unifies different unfair model phenomena, even in simple settings, remains elusive. Such a theory can offer greater interpretability of unfair model predictions and aid in the design of better unfairness evaluation and mitigation methods. For example, the theory can explain common sources of seemingly disparate unfair model behaviors, towards the creation of mitigation methods that generalize across these behaviors. Moreover, the theory can interpolate between sparse empirical findings and help us reason about how bias emerges in under-explored settings.

- **Chapter 6:** We contribute a precise analytical theory of how model design choices and data distribution properties contribute to unfairness in the setting of simplified feedforward neural networks. Our theory offers a unified and rigorous explanation of ML bias, providing insights into phenomena such as bias amplification and minority-group bias in various feature and parameter regimes. Our theoretical predictions align with empirical observations reported in the literature on ML bias, as well as offer new insights. This chapter is based on [SBS25].

In **Chapter 7**, we discuss directions for future work centered around addressing ML fairness issues caused by scaling GNNs to large networks and LLM practitioner choices

during data curation, pretraining, alignment, and evaluation. We also discuss mechanistic connections between GNNs and LLMs that can support fairness research across the models. We additionally highlight the social and systemic dimensions of fairness. Overall, developing a principled understanding of and addressing the unfairness of modern ML models is paramount to ensure that ML does not further entrench social inequalities.

Part I

On the Unfairness of Graph Neural Networks

CHAPTER 2

On the Discrimination Risk of Mean Aggregation Feature Imputation in Graphs

2.1 Introduction

Many machine learning (ML) methods for graphs rely on fully-observed features for each node, which are not available for privacy reasons, as a consequence of exclusionary data collection practices, or due to the high expenses involved in feature annotation [Woo07, SWC09, LBC20]. As a result, the development of algorithms to leverage a graph’s structure and known node feature values to impute unknown or missing features has emerged as an important research area [SSU20, RKG22].

In human networks, nodes belonging to a marginalized group (e.g., on the basis of race, gender, disability, etc.) often have a disproportionate rate of unknown features compared to the dominant group because marginalized communities may be more reluctant to share their data, annotators erase their data, and they are sidelined in data collection [BS16, BHN19, FCD21, PS22]. Furthermore, node neighborhoods are often associated with group membership [PP20, DLJ22], especially in homophilic graphs where nodes belonging to the same group have a higher likelihood of being connected [LPB22]. Homophily can be due to social stratification [HCT17] or the limited collection of inter-links (i.e., edges between nodes belonging to different groups) [LPB22, LWZ21]. Moreover, known feature values can be tainted and proxies for group membership [BHN19, WZY19]. Hence, even if graph feature imputation algorithms do not have direct access to the group membership of a node, these

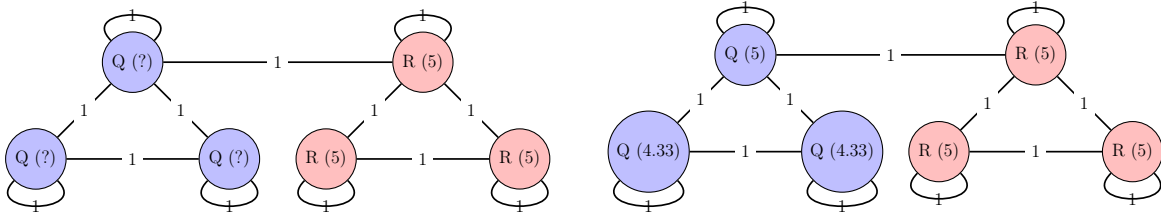


Figure 2.1: A job applicant network. After applying Feature Propagation [RKG22], the nodes in Q will have feature values that are more distinct from the feature values of the nodes in R than the ground truth.

algorithms are influenced by unknown feature rate disparities, graph structure, and known feature values. In fact, they can predict values for unknown features that cause the marginalized group’s feature values to be more distinct from the dominant group’s feature values than they are in reality.

To illustrate this phenomenon, let’s consider an automated candidate screening system based on the job applicant network shown in Figure 2.1, where nodes represent applicants and edges between applicants indicate that they have similar past work experiences. Each node has one feature: the number of years that the applicant has previously worked. Furthermore, each node belongs to one of two groups, the disabled community Q or the able-bodied community R ¹. All the nodes in R have a known feature value of 5 years, while all the nodes in Q have unknown feature values; however, in reality, all the nodes in Q also have a feature value of 5 years. Additionally, all the nodes in Q are connected to each other but have few links (if any) to nodes in R because of systemic barriers including hiring discrimination and a lack of accommodations [ASA15]. Consequently, after applying the graph feature imputation algorithm Feature Propagation [RKG22], because of the disparate rates of unknown features between the groups and structure of the graph, the nodes in Q will have, on average, feature values that are more distinct from the feature values of the nodes in R than the ground truth.

We call this distinction in imputed feature values between the marginalized and domi-

¹In reality, disability is a fluid and complex identity that should not be reduced to a binary [Bar10].

nant groups the **discrimination risk**. In this chapter, we present a theoretically-justified formulation of the discrimination risk of imputed features. We further prove that a higher discrimination risk can amplify the unfairness of a ML model applied to the imputed data, which can be especially dangerous, unethical, and illegal in high-stakes applications like automated candidate screening and loan approval [Aju19, ALZ21, WGJ21]. For instance, in the automated candidate screening example, a model applied to the imputed data may more easily learn to identify and reject disabled job applicants because they appear to have fewer years of work experience than able-bodied applicants, thereby reinforcing systemic discrimination against the disabled community [ASA15].

We also formalize a general graph feature imputation framework called **mean aggregation imputation** that encompasses common diffusion-based imputation algorithms in the literature. Subsequently, we theoretically and empirically characterize graphs in which applying mean aggregation feature imputation can yield a high discrimination risk, without making any assumptions about the underlying distributions of unknown features or graph structure. To the best of our knowledge, we are the first to study the effect of graph feature imputation on the fairness of models. This is challenging because we must consider biases stemming from graph structure.

Furthermore, we propose a simple algorithm to ensure mean aggregation-imputed features provably have a low discrimination risk, while minimally sacrificing reconstruction error (with respect to the imputation objective). We do so by viewing mean aggregation imputation through the lens of gradient descent and projecting imputed feature values onto the feasible space of feature values with low discrimination risk. We empirically evaluate the fairness and accuracy of our solution on synthetic and real-world credit networks, finding that it improves fairness without a significant loss in reconstruction error on the synthetic datasets but doesn't improve fairness on the real-world datasets. We close by discussing the limitations of our solution.

2.2 Related work

Feature imputation Feature imputation algorithms leverage known feature values to predict unknown feature values (and sometimes update known feature values). For example, unknown feature values may be filled as the mean of known values [RKG22]. However, more intricate feature imputation methods have been proposed in the ML, statistics, and epidemiology literature, with popular approaches including matrix completion [CR09, KW14, HML15, YJS18], nearest neighbors [LLP97, TCS01], multiple imputation via conditional models [RLH01, SB11, JHL23], and causal inference [SMP15, DL18]. Notably, while feature imputation may be applied to data with unknown feature values prior to the data being passed to a ML model, feature imputation is distinct from label prediction with missing data, wherein models work directly with unknown feature values [GJ93, WLX05, PDS05, XZC17, SST18, YMD20, TLM21, JZ21]; we do not consider the latter paradigm in this chapter. Works have extended feature imputation approaches for tabular data to incorporate graph structure [Hui14, KBB14]. Graph feature imputation using deep learning methods is also gaining traction [BKW17, HGL18, SSU20]. However, [RKG22] proposes a non-neural diffusion-based approach called Feature Propagation that imputes graph features by minimizing the graph’s Dirichlet energy.

Fairness of missing data and feature imputation Works have theoretically and empirically investigated the impact of missing data [WBS18, Fri20, FCD21, GAD21, ZL21a] and feature imputation on the fairness of ML models [Fri20, JWC22, ZL21b, WGJ21, LLC24]. These works consistently find that missing data can amplify biases, and some show that in practice, feature imputation can yield less unfair (relative to excluding missing data) but nevertheless discriminatory model outcomes. However, these works only study tabular data and do not consider biases that emerge from graph structure [PS22, DLJ22, ZZW25]. Furthermore, they often adapt models to directly work with missing data rather than mitigate the unfairness of feature imputation itself. Despite the prevalence of graph feature

imputation methods, to the best of our knowledge, there is no research on their influence on the fairness of models. Some works have explored fairness constraints in semi-supervised settings [ZZH20, CMT22], but they assume that node features are entirely available, which is not the case in feature imputation.

Fair graph machine learning We focus on group fairness, which ensures model predictions exhibit some form of parity between different groups [RSB19, PP20, BD20, ALZ21, MGW22, KZX22]. Works have studied (amongst other fairness formulations) statistical parity, wherein a model predicts the positive outcome at the same rate for different groups, and equal opportunity, in which the true positive rate of model predictions is equivalent across different groups [VR18]. In the automated candidate screening example, statistical parity would imply that all candidates have an equivalent likelihood to pass screening regardless of group membership, while equal opportunity would mean that, regardless of group membership (i.e., $\mathbb{P}(Z|S = Q) = \mathbb{P}(Z|S = R)$), candidates are *correctly* classified to pass screening at the same rate (i.e., $\mathbb{P}(Z|Y = 1, S = Q) = \mathbb{P}(Z|Y = 1, S = R)$). When ground-truth labels are tainted (e.g., the screening system is trained with sexist hiring data), we may prefer statistical parity to equal opportunity. While our theoretical study of discrimination risk is aligned with statistical parity, we empirically explore the effect of mean aggregation imputation with a lower discrimination risk on both the demographic parity and equal opportunity of models. Mechanisms for improving group fairness have been categorized into pre-processing [DLJ22], training-time [PP20, LPB22, CHG22], and post-processing [BH19, MWY20]. Our work is similar to pre-processing, as we seek to lower the discrimination risk of imputed training data towards improving model fairness. There exist many works on modifying graph structure to mitigate topology-induced biases [MWY20, SSH22, DLJ22, CHG22, WLL22, LWN22]. However, because graph semantics (especially for large graphs) are difficult to interpret, it is unclear if the solutions that these papers propose preserve the semantics of the original graph. On the other hand, our work does not modify graph structure and instead optimally

transforms imputed features to have a low discrimination risk. [DW21] investigates the fairness of graph neural networks in the presence of limited group membership information but does not consider imputation; in contrast, our work assumes group membership is fully available but features are not.

2.3 Discrimination risk and model unfairness

To understand how feature imputation could amplify the unfairness of a ML model, we now present a theoretically-justified metric called the discrimination risk, which applies beyond the setting of graphs and is agnostic to model architecture and labels. (We explore discrimination risk in the context of graphs in Section 2.5.) Suppose we have an arbitrary data distribution \mathcal{D} and two groups: a marginalized group Q and a dominant group R . For any data instance $(x, y, s) \sim \mathcal{D}$, let $x \in \mathcal{X}$ be the d -dimensional feature values of the instance (where x_i denotes the i -th entry of x), let $y \in \mathcal{Y}$ be its label, and let $s \in \{Q, R\}$ be its group membership. We assume that we can observe the group membership of any data instance and that no instance can belong to both the marginalized and dominant groups.

Definition 1 (Discrimination Risk) We define the discrimination risk of \mathcal{D} as:

$$\mathcal{R}_{\mathcal{D}} = \left\| \mathbb{E}_{(x,y,s) \sim \mathcal{D}}[x|s = Q] - \mathbb{E}_{(x,y,s) \sim \mathcal{D}}[x|s = R] \right\|_{\infty}, \quad (2.1)$$

We now explore the relevance of discrimination risk to model unfairness. Let \mathcal{D}' be the data distribution with ground-truth features and \mathcal{D} be the distribution with imputed features. We will show that if $\mathcal{R}_{\mathcal{D}} > \mathcal{R}_{\mathcal{D}'}$, imputation may amplify model unfairness.

To be concrete, let's again consider the example of the automated candidate screening system. Let $\mathbb{P}(S \in \{Q, R\})$ be the distribution over the group membership of job applicants. Furthermore, let $\mathbb{P}(X' \in \mathcal{X})$ be the distribution over the ground-truth feature values of applicants (e.g., number of years previously worked, highest degree, etc.) In contrast, let $\mathbb{P}(X \in \mathcal{X})$ be the distribution over imputed feature values of applicants. Let $\mathbb{P}(Y \in \mathcal{Y})$ be the

distribution over ground-truth applicant labels (e.g., whether an applicant should be hired, a screening score for an applicant, etc.) Now, let h' be a model trained on samples from \mathcal{D}' (without direct access to S), and h be another model trained on samples from \mathcal{D} (without direct access to S). Additionally, let $\mathbb{P}(Z' \in \mathcal{Y})$ be the distribution over the predictions of h' on instances sampled from $\mathbb{P}(X')$, and let $\mathbb{P}(Z \in \mathcal{Y})$ be the distribution over the predictions of h on instances sampled from $\mathbb{P}(X)$.

We have the dependencies $S \rightarrow X'$ and $S \rightarrow X$ because disability affects the number of years previously worked by a job applicant. Additionally, $S \rightarrow X$ because S can influence the feature imputation algorithm (as illustrated in Figure 2.1). Furthermore, we assume that h' and h have access to the feature values of an applicant, but not the applicant's group membership, and that the association of S with Z' through Y can be fully explained by X' . Thus, Z' is conditionally independent of S given X' . Similarly, we assume that the association of S with Z through Y can be fully explained by X , so Z is conditionally independent of S given X .

Because we are interested in how feature imputation may amplify model unfairness, we investigate when the statistical parity or total variation distance $d_{TV}(\mathbb{P}(Z|S = Q), \mathbb{P}(Z|S = R)) > d_{TV}(\mathbb{P}(Z'|S = Q), \mathbb{P}(Z'|S = R))$ [DHP12], where $\mathbb{P}(Z|S = Q)$ is the prediction distribution conditioned on group- Q membership. $d_{TV}(A, B)$ measures the distance between two probability distributions A and B as $\sup_{x \in \mathcal{F}} |A(x) - B(x)|$. Intuitively, $d_{TV}(\mathbb{P}(Z|S = Q), \mathbb{P}(Z|S = R))$ captures how much the prediction distribution of h absolutely differs between Q and R , and thus it quantifies the (statistical parity) unfairness of h . We now discuss the relationships of $d_{TV}(\mathbb{P}(Z'|S = Q), \mathbb{P}(Z'|S = R))$ to $d_{TV}(\mathbb{P}(X'|S = Q), \mathbb{P}(X'|S = R))$ and $d_{TV}(\mathbb{P}(Z|S = Q), \mathbb{P}(Z|S = R))$ to $d_{TV}(\mathbb{P}(X|S = Q), \mathbb{P}(X|S = R))$. We begin with the following lemma from [ZG19].

Lemma 1 (Corollary 17 from [ZG19]) By the Data Processing Inequality, $d_{TV}(\mathbb{P}(Z'|S = Q), \mathbb{P}(Z'|S = R)) \leq d_{TV}(\mathbb{P}(X'|S = Q), \mathbb{P}(X'|S = R))$ and $d_{TV}(\mathbb{P}(Z|S = Q), \mathbb{P}(Z|S = R)) \leq$

$d_{TV}(\mathbb{P}(X|S = Q), \mathbb{P}(X|S = R))$.

Please refer to Appendix A.1.1 for the proof of Lemma 1. At a high level, Lemma 1 states that the statistical parity unfairness of a model is upper-bounded by the statistical parity distance of the feature distributions between the groups. We now use Lemma 1 to prove the following theorem that connects $d_{TV}(\mathbb{P}(Z'|S = Q), \mathbb{P}(Z'|S = R))$ to $\mathcal{R}_{\mathcal{D}'}$ and $d_{TV}(\mathbb{P}(Z|S = Q), \mathbb{P}(Z|S = R))$ to $\mathcal{R}_{\mathcal{D}}$.

Theorem 1 Suppose $\mathbb{P}(X'|S = Q) = \mathcal{N}(\mu'_Q, \Sigma'_Q)$; $\mathbb{P}(X'|S = R) = \mathcal{N}(\mu'_R, \Sigma'_R)$; $\mathbb{P}(X|S = Q) = \mathcal{N}(\mu_Q, \Sigma_Q)$; and $\mathbb{P}(X|S = R) = \mathcal{N}(\mu_R, \Sigma_R)$. We then get the following bounds:

$$d_{TV}(\mathbb{P}(Z'|S = Q), \mathbb{P}(Z'|S = R)) \in \left[\frac{1}{\frac{4 \cdot \max\{\lambda_{max}(\Sigma'_Q), \lambda_{max}(\Sigma'_R)\}}{C' \cdot \mathcal{R}_{\mathcal{D}'}} + 1}, \sqrt{1 - \sqrt{\frac{\det \Sigma'_Q}{\det \Sigma'_R} \cdot e^{-\frac{C' \cdot \mathcal{R}_{\mathcal{D}'}}{\lambda_{min}(\Sigma'_R)} - tr(\Sigma_R^{-1} \Sigma'_Q) + d}}} \right];$$

$$d_{TV}(\mathbb{P}(Z|S = Q), \mathbb{P}(Z|S = R)) \in \left[\frac{1}{\frac{4 \cdot \max\{\lambda_{max}(\Sigma_Q), \lambda_{max}(\Sigma_R)\}}{C \cdot \mathcal{R}_{\mathcal{D}}} + 1}, \sqrt{1 - \sqrt{\frac{\det \Sigma_Q}{\det \Sigma_R} \cdot e^{-\frac{C \cdot \mathcal{R}_{\mathcal{D}}}{\lambda_{min}(\Sigma_R)} - tr(\Sigma_R^{-1} \Sigma_Q) + d}}} \right],$$

where $0 \leq C' \leq d$ and $0 \leq C \leq d$. Please refer to Appendix A.1.2 for the proof of Theorem 1.

This result suggests that minimizing $\mathcal{R}_{\mathcal{D}}$ can minimize the unfairness of the model h applied to the imputed data. Furthermore, it is possible that $\mathcal{R}_{\mathcal{D}} > \mathcal{R}_{\mathcal{D}'}$, in which case feature imputation may amplify the unfairness of a model. While we leverage strong generative assumptions in Theorem 1, it is plausible that $X'|S = Q$ and $X'|S = R$ are normally-distributed. Furthermore, mean aggregation imputation (Section 2.4) produces approximately normal $X|S = Q$ and $X|S = R$ (by the Central Limit Theorem). We additionally note that the lower bounds do not require the feature values to be normally distributed; the bounds only assume their distributions have finite covariance. Finally, for arbitrary distributions, matching even an infinite number of moments is not sufficient to bound their distance [LB00].

While $\mathcal{R}_{\mathcal{D}}$ applies beyond graphs and is agnostic to model architecture and labels, in practice, it is important to consider model complexity and task context. We also add that $\mathcal{R}_{\mathcal{D}}$ risk bears resemblance to the Average Treatment Effect studied in causality [Sek08]. Moreover, [JWC22] proposes a metric also called discrimination risk which quantifies how much the deviation of imputed feature values from the ground-truth feature values differs

across groups; this quantity is more aligned with the accuracy disparity (and thus equalized odds) of a model applied to the imputed data. In contrast, $\mathcal{R}_{\mathcal{D}}$ simply measures how much imputed features differ across groups, enabling its computation when ground-truth feature values are not available, and is aligned with the statistical parity of a model. We leave extending our definition of discrimination risk to other formulations of fairness as future work [VR18].

2.4 Graph feature imputation

We would like to investigate the discrimination risk of graph feature imputation. Prior to doing so, we present a general framework called mean aggregation imputation that encompasses common non-neural diffusion-based graph feature imputation algorithms in the literature.

Suppose we have an undirected weighted homogeneous graph $G = (V, E)$. Each node has d features, hence the node feature matrix $X \in \mathbb{R}^{N \times d}$, where $N = |V|$ (i.e., the cardinality of V). For simplicity, we assume that $d = 1$. The feature value is unknown for some nodes in G (denoted as the set U), and known for others (denoted as the set K). The feature value of each node is either known or unknown (i.e., $U \cup K = V$ and $U \cap K = \emptyset$). Let X_S refer to the feature values of the nodes in set S . Assume without loss of generality that $X = \begin{bmatrix} X_K \\ X_U \end{bmatrix}$. Furthermore, let $A \in \mathbb{R}^{N \times N}$ denote the weighted adjacency matrix of G , where A_{ij} is the nonnegative weight corresponding to the edge from node j to node i . Additionally, let $A_{S_1 S_2}$ denote the submatrix of A with rows belonging to the nodes in set S_1 and columns belonging to the nodes in set S_2 . Let $A := \begin{bmatrix} A_{KK} & A_{KU} \\ A_{UK} & A_{UU} \end{bmatrix}$. D is the diagonal degree matrix, i.e., $D_{ii} = \sum_{j=1}^N A_{ij}$ and $D := \begin{bmatrix} D_K & 0 \\ 0 & D_U \end{bmatrix}$.

Definition 2 (Mean Aggregation Feature Imputation) Denote the feature values at iteration t of mean aggregation feature imputation as $X^{(t)}$. Furthermore, let $X_S^{(t)}$ refer to the feature values of the nodes in set S at iteration t . Then, at each iteration t :

$$MX^{(t+1)} := \phi(MX^{(t)}) = \begin{bmatrix} \beta I_{|K|} & 0 \\ 0 & I_{|U|} \end{bmatrix} TMX^{(t)} + \begin{bmatrix} (1 - \beta)I_{|K|} & 0 \\ 0 & 0 \end{bmatrix} MX^{(0)}, \quad (2.2)$$

where $M : \mathbb{R}^N \rightarrow \mathbb{R}^N$ is a diagonal invertible map, $T \in \mathbb{R}^{N \times N}$ is a right-stochastic matrix, and $\beta \in [0, 1]$ is a regularization hyperparameter.

When $d > 1$, we apply mean aggregation feature imputation to each channel independently², and it only works with continuous (not discrete) features. Mean aggregation encompasses common graph feature imputation methods, including **Global Mean** (predicts unknown feature values as the uniform mean of known feature values), **Neighbor Mean** (predicts unknown feature values as the degree-weighted mean of the known feature values for neighboring nodes), **Feature Propagation** (predicts unknown feature values that minimize the Dirichlet energy of the graph while preserving known feature values), and **Graph Regularization** (predicts feature values via a smoothness constraint and a fitting constraint for the known features). For proofs, refer to Appendix A.1.3. Despite Neighbor Mean, Feature Propagation, and Graph Regularization being intended for homophilic graphs [RKG22], mean aggregation feature imputation encompasses algorithms that could perform well on heterophilic graphs as well with an appropriate choice of T [RKG22]. We also note that T cannot be A , as this might cause feature values to explode over multiple iterations.

2.5 Discrimination risk of mean aggregation feature imputation

To understand how mean aggregation feature imputation may amplify the unfairness of a ML model, we theoretically characterize graphs in which mean aggregation imputation increases

²Incorporating associations between features is a promising direction of research.

the discrimination risk, without making assumptions about the underlying distributions of unknown features or graph structure.

We begin by defining new notation. In particular, we first focus on the case of a single feature (i.e., $d = 1$), and extend our analysis to the case $d > 1$ in Appendix A.1.5. Let X_v be the feature of a node v . Let Q denote the set of nodes that belong to the marginalized group, and R denote the set of nodes that belong to the dominant group. We assume that $Q \cup R = V$ and $Q \cap R = \emptyset$. We define the discrimination risk after t iterations of mean aggregation imputation as:

$$\mathcal{R}^{(t)} := \left| \mathbb{E}_{q \sim Q}[X_q^{(t)}] - \mathbb{E}_{r \sim R}[X_r^{(t)}] \right|, \quad (2.3)$$

where the expectations are taken uniformly over the nodes in each set. Now, define $\tilde{X} := MX$. Then, a modified version of the discrimination risk after t iterations of mean aggregation imputation is $\tilde{\mathcal{R}}^{(t)} := \left| \mathbb{E}_{q \sim Q}[\tilde{X}_q^{(t)}] - \mathbb{E}_{r \sim R}[\tilde{X}_r^{(t)}] \right|$. In Theorem 2, we bound the discrimination risk $\tilde{\mathcal{R}}^{(t)}$ of the imputed features with respect to $\tilde{\mathcal{R}}^{(0)}$. We bound $\tilde{\mathcal{R}}^{(t)}$ rather than $\mathcal{R}^{(t)}$ for simplicity; however, we empirically validate that the bound properties also hold for $\mathcal{R}^{(t)}$ in Section 2.7. Define $\tilde{\mu}_Q^{(t)} := \mathbb{E}_{q \sim Q}[\tilde{X}_q^{(t)}]$ and $\tilde{\mu}_R^{(t)} := \mathbb{E}_{r \sim R}[\tilde{X}_r^{(t)}]$. Furthermore, let $\tilde{\sigma}^{(t)}$ denote the maximal deviation of the feature values at iteration t , i.e., $\forall q_1 \in Q, |\tilde{X}_{q_1}^{(t)} - \tilde{\mu}_Q^{(t)}| \leq \max_{q_2 \in Q} |\tilde{X}_{q_1}^{(t)} - \tilde{X}_{q_2}^{(t)}| \leq \tilde{\sigma}^{(t)}$ and $\forall r_1 \in R, |\tilde{X}_{r_1}^{(t)} - \tilde{\mu}_R^{(t)}| \leq \max_{r_2 \in R} |\tilde{X}_{r_1}^{(t)} - \tilde{X}_{r_2}^{(t)}| \leq \tilde{\sigma}^{(t)}$. Additionally, define $T_{S_1 \rightarrow S_2} := \sum_{b \in S_2} \sum_{a \in S_1} T_{ba}$.

Theorem 2 Let the contraction coefficient $\alpha := \left| 1 - \frac{T_{R \rightarrow Q \cap U} + \beta T_{R \rightarrow Q \cap K}}{|Q|} - \frac{T_{Q \rightarrow R \cap U} + \beta T_{Q \rightarrow R \cap K}}{|R|} \right|$.

Then, $\alpha \leq 1$, and:

$$\max \left\{ \alpha^t \tilde{\mathcal{R}}^{(0)} - 2 \left(\sum_{j=0}^{t-1} \alpha^j \right) \tilde{\sigma}^{(0)}, 0 \right\} \leq \tilde{\mathcal{R}}^{(t)} \leq \alpha^t \tilde{\mathcal{R}}^{(0)} + 2 \left(\sum_{j=0}^{t-1} \alpha^j \right) \tilde{\sigma}^{(0)}$$

$$\alpha < 1 \implies \lim_{t \rightarrow \infty} \tilde{\mathcal{R}}^{(t)} \leq \frac{2\tilde{\sigma}^{(0)}}{1 - \alpha}.$$

Please refer to Appendix A.1.4 for a proof of Theorem 2. Theorem 2 shows that the bounds on the discrimination risk contract more slowly (with more iterations of mean aggregation

feature imputation) with a larger α . Furthermore, the upper bound on the discrimination risk is larger with a larger α and may depend on the initial unknown feature values.

2.5.1 Analysis of Theorem 2

Theorem 2 allows for interesting interpretations of how graph properties like the rate of unknown features, group size, and graph structure affect the discrimination risk of mean aggregation imputation. Below, we successively vary each property (holding the other properties constant) and investigate its impact on α , and in turn the discrimination risk. We focus on Feature Propagation [RKG22], but our analysis may be easily extended to other mean aggregation imputation algorithms. We assume for simplicity that all edges have a weight of 1.

Unknown feature rates *A low unknown feature rate for both groups or disparate unknown feature rates across the groups can increase α , and thus the discrimination risk of mean aggregation-imputed features.* Suppose the intra-link rate $\mathbb{P}((u, v) \in E | u \in Q, v \in Q) = \mathbb{P}((u, v) \in E | u \in R, v \in R) = \frac{1}{2}$ and inter-link rate $\mathbb{P}((u, v) \in E | u \in Q, v \in R) = \mathbb{P}((u, v) \in E | u \in R, v \in Q) = \frac{1}{2}$. Furthermore, assume equal (relative) group sizes, i.e., $\frac{|Q|}{N} = \frac{|R|}{N} = \frac{1}{2}$. Then, $\frac{T_{R \rightarrow Q \cap U} + \beta T_{R \rightarrow Q \cap K}}{|Q|} = \frac{T_{R \rightarrow Q \cap U}}{N/2} = \frac{\sum_{q \in Q \cap U} \sum_{r \in R} T_{qr}}{N/2} = \frac{\sum_{q \in Q \cap U} \sum_{r \in R} D_{qq}^{-1} A_{qr}}{N/2}$. By decomposition, $D_{qq} = \sum_{u \in Q} A_{qu} + \sum_{v \in R} A_{qv} = \frac{1}{2}|Q| + \frac{1}{2}|R| = \frac{N}{2}$. Therefore, $\frac{T_{R \rightarrow Q \cap U} + \beta T_{R \rightarrow Q \cap K}}{|Q|} = \frac{\sum_{q \in Q \cap U} \sum_{r \in R} A_{qr}}{N^2/4} = \frac{\frac{1}{2}(|Q \cap U| \times |R|)}{N^2/4} = \frac{\frac{1}{2}(\mathbb{P}(v \in U | v \in Q) \cdot |Q| \times |R|)}{N^2/4} = \frac{1}{2} \mathbb{P}(v \in U | v \in Q)$, where $\mathbb{P}(v \in U | v \in Q)$ is the unknown feature rate for group Q . Similarly, $\frac{T_{Q \rightarrow R \cap U} + \beta T_{Q \rightarrow R \cap K}}{|R|} = \frac{1}{2} \mathbb{P}(v \in U | v \in R)$. Thus, $\alpha = |1 - \frac{1}{2} \mathbb{P}(v \in U | v \in Q) - \frac{1}{2} \mathbb{P}(v \in U | v \in R)|$. This aligns with [ZL21b]’s finding that imputation fairness can be influenced by the imbalance of feature missingness across groups, although [ZL21b] studies equalized odds rather than statistical parity fairness.

Group sizes *Group size alone may not affect α or the discrimination risk of mean aggregation-imputed features.* Suppose the intra-link rate $\mathbb{P}((u, v) \in E | u \in Q, v \in Q) =$

$\mathbb{P}((u, v) \in E | u \in R, v \in R) = \frac{1}{2}$ and inter-link rate $\mathbb{P}((u, v) \in E | u \in Q, v \in R) = \mathbb{P}((u, v) \in E | u \in R, v \in Q) = \frac{1}{2}$. Furthermore, assume equal unknown feature rates, i.e., $\mathbb{P}(v \in U | v \in Q) = \mathbb{P}(v \in U | v \in R) = \frac{1}{2}$. Then, $\frac{T_{R \rightarrow Q \cap U} + \beta T_{R \rightarrow Q \cap K}}{|Q|} = \frac{\sum_{q \in Q \cap U} \sum_{r \in R} T_{qr}}{|Q|} = \frac{\sum_{q \in Q \cap U} \sum_{r \in R} D_{qq}^{-1} A_{qr}}{|Q|}$. $D_{qq} = \frac{N}{2}$. Therefore, $\frac{T_{R \rightarrow Q \cap U} + \beta T_{R \rightarrow Q \cap K}}{|Q|} = \frac{\sum_{q \in Q \cap U} \sum_{r \in R} A_{qr}}{|Q| \cdot N/2} = \frac{\frac{1}{2}(|Q \cap U| \times |R|)}{|Q| \cdot N/2} = \frac{\frac{1}{2}(\frac{1}{2}|Q| \times |R|)}{|Q| \cdot N/2} = \frac{1}{2} \cdot \frac{|R|}{N}$. Similarly, $\frac{T_{Q \rightarrow R \cap U} + \beta T_{Q \rightarrow R \cap K}}{|R|} = \frac{1}{2} \cdot \frac{|Q|}{N}$. Thus, $\alpha = |1 - \frac{1}{2} \cdot \frac{|R|}{N} - \frac{1}{2} \cdot \frac{|Q|}{N}| = \frac{1}{2}$.

Graph structure *A low inter-link to intra-link ratio can increase α and the discrimination risk of mean aggregation-imputed features.* Suppose we have equal unknown feature rates, i.e.,

$\mathbb{P}(v \in U | v \in Q) = \mathbb{P}(v \in U | v \in R) = \frac{1}{2}$. Furthermore, assume equal (relative) group sizes, i.e., $\frac{|Q|}{N} = \frac{|R|}{N} = \frac{1}{2}$. Then, $\frac{T_{R \rightarrow Q \cap U} + \beta T_{R \rightarrow Q \cap K}}{|Q|} = \frac{\sum_{q \in Q \cap U} \sum_{r \in R} T_{qr}}{N/2} = \frac{\sum_{q \in Q \cap U} \sum_{r \in R} D_{qq}^{-1} A_{qr}}{N/2}$. $D_{qq} = \sum_{u \in Q} A_{qu} + \sum_{v \in R} A_{qv} = \mathbb{P}((u, v) \in E | u \in Q, v \in Q) |Q| + \mathbb{P}((u, v) \in E | u \in R, v \in Q) |R| = \frac{N}{2} [\mathbb{P}((u, v) \in E | u \in Q, v \in Q) + \mathbb{P}((u, v) \in E | u \in R, v \in Q)]$. Therefore, $\frac{T_{R \rightarrow Q \cap U} + \beta T_{R \rightarrow Q \cap K}}{|Q|} = \frac{\mathbb{P}((u, v) \in E | u \in R, v \in Q) (|Q \cap U| \times |R|)}{[\mathbb{P}((u, v) \in E | u \in Q, v \in Q) + \mathbb{P}((u, v) \in E | u \in R, v \in Q)] \cdot N^2/4} = \frac{\frac{1}{2}|Q| \times |R|}{[1 + \frac{\mathbb{P}((u, v) \in E | u \in R, v \in Q)}{\mathbb{P}((u, v) \in E | u \in Q, v \in Q)}] \cdot N^2/4} = \frac{1}{2} \cdot \frac{1}{1 + \frac{\mathbb{P}((u, v) \in E | u \in R, v \in Q)}{\mathbb{P}((u, v) \in E | u \in Q, v \in Q)}}$. Similarly, $\frac{T_{Q \rightarrow R \cap U} + \beta T_{Q \rightarrow R \cap K}}{|R|} = \frac{1}{2} \cdot \frac{1}{1 + \frac{\mathbb{P}((u, v) \in E | u \in Q, v \in R)}{\mathbb{P}((u, v) \in E | u \in R, v \in R)}}$.

Thus, $\alpha = \left| 1 - \frac{1}{2} \cdot \frac{1}{1 + \frac{\mathbb{P}((u, v) \in E | u \in R, v \in Q)}{\mathbb{P}((u, v) \in E | u \in Q, v \in Q)}} - \frac{1}{2} \cdot \frac{1}{1 + \frac{\mathbb{P}((u, v) \in E | u \in Q, v \in R)}{\mathbb{P}((u, v) \in E | u \in R, v \in R)}} \right|$. Because G is undirected, $\mathbb{P}((u, v) \in E | u \in R, v \in Q) = \mathbb{P}((u, v) \in E | u \in Q, v \in R)$ and $\mathbb{P}((u, v) \in E | u \in Q, v \in Q) = \mathbb{P}((u, v) \in E | u \in R, v \in R)$.

Ultimately, our theoretical (Section 2.5) and empirical (Appendix A.2.5) characterizations of α and the discrimination risk can be leveraged to audit real-world graph data for structural factors that contribute to the unfairness of mean aggregation feature imputation and ML models applied to the imputed data.

2.6 Fairer graph feature imputation

We propose a simple and effective solution to ensure mean aggregation feature imputation provably has a low discrimination risk, while minimally sacrificing reconstruction error (with respect to the imputation objective). At a high level, we want to constrain the discrimination

risk at every iteration t of imputation to be at most ϵ (i.e., $\forall t \in [0, \infty), \mathcal{R}^{(t)} \leq \epsilon$). We do this by viewing mean aggregation imputation through a gradient descent lens and projecting imputed feature values onto the feasible space of feature values with discrimination risk at most ϵ at each iteration [Sub22]. We focus on a single feature i , but our algorithms can be extended to more features by applying them to each feature separately.

We begin with the case where known feature values remain fixed (i.e., $\beta = 0$). Recall the iterative algorithm for mean aggregation feature imputation when $\beta = 0$:

$$MX^{(t+1)} := \phi(MX^{(t)}) = \begin{bmatrix} 0 & 0 \\ 0 & I_{|U|} \end{bmatrix} TMX^{(t)} + \begin{bmatrix} I_{|K|} & 0 \\ 0 & 0 \end{bmatrix} MX^{(0)}.$$

We see that $\forall t \in [0, \infty), X_K^{(t)} = X_K$. Furthermore, define $\Delta := I_N - M^{-1}TM$. In the case of Feature Propagation, $\Delta = I_N - D^{\frac{1}{2}}(D^{-1}A)D^{-\frac{1}{2}} = I_N - D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$ is the symmetric normalized Laplacian of G [RKG22]. We see that $X_U^{(t+1)} = (I_N - \Delta)_{UU}X_U^{(t)} - \Delta_{UK}X_K$. As discussed in [RKG22], we can view the update $X_U^{(t+1)} := X_U^{(t)} - \gamma(\Delta_{UU}X_U^{(t)} + \Delta_{UK}X_K)$ as an iteration of gradient descent (with step size $\gamma = 1$) on the objective function $\ell(x) = \frac{1}{2}x^T \Delta_{UU}x + X_K^T \Delta_{KU}x + \frac{1}{2}X_K^T \Delta_{KK}X_K$, where X_K is constant. For Feature Propagation, ℓ is the Dirichlet energy of G [RKG22]. We now present a theorem that shows how to perform mean aggregation feature imputation with a discrimination risk of at most ϵ when the known feature values remain fixed.

Theorem 3 (ϵ -Fair Imputation, $\beta = 0$) Vanilla mean aggregation feature imputation updates $X_U^{(t+1)} := X_U^{(t)} - \gamma(\Delta_{UU}X_U^{(t)} + \Delta_{UK}X_K) = Z_U^{(t)}$, where $\gamma = 1$. Let ϵ -fair mean aggregation feature imputation instead update $X_U^{(t+1)} := P_W Z_U^{(t)} + P_B$, where:

$$P_W = \begin{cases} I_{|U|}, & \mathcal{R}_K - \epsilon \leq c^T Z_U^{(t)} \leq \mathcal{R}_K + \epsilon \\ I_{|U|} - \frac{cc^T}{c^T c}, & \text{otherwise} \end{cases}, P_B = \frac{cc^T}{c^T c} \begin{cases} \mathcal{R}_K - \epsilon, & c^T Z_U^{(t)} < \mathcal{R}_K - \epsilon \\ \mathcal{R}_K + \epsilon, & c^T Z_U^{(t)} > \mathcal{R}_K + \epsilon \\ 0, & \text{otherwise} \end{cases}$$

$$\mathcal{R}_K = \frac{1}{|R|} \sum_{r \in R \cap K} X_r - \frac{1}{|Q|} \sum_{q \in Q \cap K} X_q$$

$$c \in \mathbb{R}^{|U|}, c^T Z_U^{(t)} = \frac{1}{|Q|} \sum_{q \in Q \cap U} Z_q^{(t)} - \frac{1}{|R|} \sum_{r \in R \cap U} Z_r^{(t)}$$

Then, assuming $0 \leq \lambda_{\min}(\Delta_{UU}) \leq \lambda_{\max}(\Delta_{UU}) < 1$ (where λ_{\min} and λ_{\max} are the minimum and maximum eigenvalues, respectively): 1) a unique optimal (with respect to ℓ) feasible solution X_U^* exists; 2) for fixed step size $\gamma = \frac{1}{\lambda_{\max}(\Delta_{UU})}$, ϵ -fair imputation converges as $\|X_U^{(t)} - X_U^*\|_2^2 \leq \left(1 - \frac{\lambda_{\min}(\Delta_{UU})}{\lambda_{\max}(\Delta_{UU})}\right)^t \|X_U^{(0)} - X_U^*\|_2^2$; 3) for fixed step size $\gamma \leq \frac{1}{\lambda_{\max}(\Delta_{UU})}$, ϵ -fair imputation converges to X_U^* .

Please refer to Appendix A.1.6 for a proof of Theorem 3. The convergence of our solution to the unique optimal (with respect to ℓ) feasible solution implies that our solution provably has a discrimination risk of at most ϵ while minimally sacrificing reconstruction error (with respect to the objective). Furthermore, because our solution simply interleaves projections into the mean aggregation imputation framework, it preserves the framework’s speed and scalability [RKG22]. We similarly have a solution when $\beta > 0$ (i.e., when the known node feature values do not remain fixed) which we present in Appendix A.1.7. Choosing ϵ , due to its uninterpretable nature, can be difficult in practice; we encourage work on making ϵ more intelligible.

2.7 Experimental results and discussion

We empirically evaluate the fairness and reconstruction error of our solution on various mean aggregation feature imputation algorithms and synthetic and real-world datasets. We find that while our solution yields improved fairness without a significant loss in reconstruction error on the synthetic datasets, there is not an improvement in fairness on the real-world datasets.

Datasets We construct undirected two-block synthetic networks (SBM) using `StochasticBlockModelDataset` from `PyTorch Geometric` [FL19] (where one block corresponds to the marginalized group Q and the other block to the dominant group R) with various (relative)

group sizes and inter- and intra-link rates (more information in Appendix A.2.1). SBM does not have a corresponding task, i.e., the nodes do not have labels. We also use the real-world `Credit defaulter` and `German credit` networks from [ALZ21] (there exist limited “natively” graph real-world datasets with sensitive attributes available). The `Credit defaulter` network consists of 30,000 nodes representing individuals, with edges between them indicating similar spending and payment patterns. The corresponding task is to predict whether an individual will default on their credit card payment or not, and the groups are those 25 years old or younger and those above the age of 25 (more information in Appendix A.2.1). The `German credit` network comprises 1,000 nodes representing clients in a German bank who are connected if they have similar credit accounts. The corresponding task is to predict whether a client has good or bad credit risk, and the groups are men and women (more information in Appendix A.2.1)³. We refrain from using the `Recidivism graph` from [ALZ21] so as not to support the development of carceral technology [Ham21].

Protocols and performance evaluation By default, none of the datasets contain unknown or missing features. Despite the real-world prevalence of missing features, most publicly-available graph datasets inherently do not have missing node features because current graph ML techniques predominantly rely on fully-observable features. Hence, similarly to [RKG22], to simulate diverse scenarios with unknown features, for each group, we independently at random mark node feature values as unknown with a different probability. Nodes, even within the same group, may have different unknown features. In this way, we simulate missing completely at random (MCAR) and missing at random (MAR) on our data. We experiment with all 25 combinations of unknown feature rates of $\{0.1, 0.3, 0.5, 0.7, 0.9\}$ for the groups. We choose this scheme to study the effect of disparate unknown feature rates across the groups, which reflects the real world [BS16, BHN19, FCD21, PS22]. Empirically characterizing the real-world distributions of node feature missingness requires further study.

³In reality, gender is neither binary nor static, and treating it as such can produce harms [DMO21]

We also encourage empirical work on other unknown feature schemes, including missingness based on node degree or marking all or none of the features for each node as unknown.

To impute unknown features, we use the vanilla mean aggregation imputation algorithms overviewed in Section 2.4: Global Mean (GM), Neighbor Mean (NM), Feature Propagation (FP), and Graph Regularization (GR) (with $\beta = 0.25$), and their ϵ -fair counterparts (for $\epsilon \in \{0.0, 0.025, 0.05\}$). For models, we use a linear classifier (*linear*), two-layer MLP (*mlp*), and two-layer Graph Convolutional Network (*gcn*) [KW17] (more information in Appendix A.2.3). We train all models on the imputed data, but validate and test on fully-observed data. Because there are no previous works (to the best of our knowledge) that directly address the unfairness of graph feature imputation, we do not have baselines against which to compare.

To evaluate imputed features for SBM, since we don't have labels, we employ relative reconstruction error **RE** [RKG22]. For the real-world datasets, we consider the test accuracy (**Acc**) of models applied to the imputed data. To evaluate group fairness, we compute the discrimination risk (**DR**) of the imputed data. For SBM, we also train models on the imputed data to predict group membership and calculate the test accuracy of the models on identifying group membership (which we refer to as **MI**) [ZLM18, PP20]. To evaluate group fairness for the real-world datasets, we use the test statistical parity (**SP**) of the models, defined as $|\mathbb{P}(Z = 1|S = Q) - \mathbb{P}(Z = 1|S = R)|$ [DHP12], and test equal opportunity (**EO**), defined as $|\mathbb{P}(Z = 1|S = Q, Y = 1) - \mathbb{P}(Z = 1|S = R, Y = 1)|$ [HPS16]. Please refer to Appendix A.2.4 for more details on the metrics. For all metrics, we report the mean and standard error over 5 runs using different random seeds. On each run, a new dataset (in the case of SBM) is generated, new splits are created, and a new model is trained.

Q1. Does the contraction coefficient α align with the discrimination risk of mean aggregation imputation across graphs with different properties? Figures A.1 to A.11 in the appendix show that for SBM, discrimination risk and α generally have a strong positive association over unknown feature rates, group sizes, and inter-link rate to intra-link

rate ratios, which substantiates Theorem 2. This association is weaker for Global Mean and Neighbor Mean. Refer to Appendix A.2.5 for more details.

Q2. Does ϵ -fair mean aggregation imputation (compared to regular mean aggregation imputation) improve the group fairness of a model applied to the imputed data? Table 2.1 shows that, for SBM, ϵ -fair mean aggregation imputation achieves comparable reconstruction error to its vanilla counterpart while greatly reducing the discrimination risk and test group membership identification accuracy for all models. As expected, we see that the discrimination risk of ϵ -fair imputation is at most ϵ , and discrimination risk and test group membership identification accuracy are positively associated, which substantiates Theorem 1. Furthermore, the reconstruction error for ϵ -fair FP and GR (which leverage graph structure and are thus more susceptible to graph structural bias) are much lower than that of the naïve ϵ -fair GM (which does not consider graph structure), but fair FP and GR reduce the discrimination risk and test group membership identification accuracy for all models to similar levels as fair GM. However, the test group membership identification accuracy generally decreases less as ϵ decreases for *mlp* and *gcn* than it does for *linear*, which suggests that minimizing the discrimination risk of imputed features is more effective at removing linearly-encoded group membership information than non-linearly encoded information. Furthermore, ϵ -fair feature imputation does not guard against group membership information that *gcn* learns via graph structure during training. We have similar findings when averaging over different relative group sizes (refer to Table A.1 in the appendix) and combinations of inter- and intra-link rates (refer to Table A.2 in the appendix).

In contrast, Tables 2.2 and A.4 (in the appendix) show that, for the real-world datasets, regular mean aggregation imputation and its ϵ -fair counterpart yield comparable test accuracy and statistical parity fairness for all models. In the case of `Credit defaulter`, our solution even appears to exacerbate the unfairness of *gcn*. We find similar results for equal opportunity, as shown in Tables A.3 and A.5 (in the appendix). We were unable to obtain similar unfairness

Table 2.1: Reconstruction error (**RE**), discrimination risk (**DR**), and test group membership identification accuracy (**MI**) of all models averaged over all 25 combinations of unknown feature rates of $\{0.1, 0.3, 0.5, 0.7, 0.9\}$ for each group in **SBM**. We use 0.5 relative group sizes and 0.5 inter- and intra-link rates.

Method	RE ↓	DR ↓	MI _{linear} ↓	MI _{mlp} ↓	MI _{gcn} ↓
0-Fair GM	1.21 ± 0.021	0 ± 0	0.602 ± 0.098	0.669 ± 0.019	0.504 ± 0.038
0.025-Fair GM	1.204 ± 0.021	0.021 ± 0.002	0.72 ± 0.087	0.683 ± 0.01	0.551 ± 0.069
0.05-Fair GM	1.196 ± 0.02	0.034 ± 0.005	0.736 ± 0.037	0.69 ± 0.014	0.535 ± 0.058
Regular GM	1 ± 0	0.051 ± 0.015	0.817 ± 0.02	0.817 ± 0.013	0.651 ± 0.089
0-Fair NM	1.19 ± 0.02	0 ± 0	0.599 ± 0.094	0.686 ± 0.021	0.507 ± 0.04
0.025-Fair NM	1.183 ± 0.02	0.02 ± 0.002	0.72 ± 0.084	0.706 ± 0.013	0.552 ± 0.059
0.05-Fair NM	1.175 ± 0.019	0.033 ± 0.004	0.734 ± 0.037	0.706 ± 0.02	0.539 ± 0.052
Regular NM	0.977 ± 0.002	0.048 ± 0.014	0.828 ± 0.021	0.818 ± 0.013	0.631 ± 0.094
0-Fair FP	1.184 ± 0.020	0 ± 0	0.6 ± 0.096	0.702 ± 0.02	0.505 ± 0.034
0.025-Fair FP	1.176 ± 0.020	0.018 ± 0.003	0.716 ± 0.076	0.724 ± 0.014	0.562 ± 0.058
0.05-Fair FP	1.169 ± 0.019	0.025 ± 0.006	0.72 ± 0.03	0.723 ± 0.018	0.531 ± 0.058
Regular FP	0.977 ± 0.003	0.028 ± 0.009	0.814 ± 0.023	0.817 ± 0.012	0.612 ± 0.079
0-Fair GR	1.006 ± 0.004	0 ± 0	0.588 ± 0.093	0.713 ± 0.022	0.511 ± 0.039
0.025-Fair GR	1.005 ± 0.004	0.023 ± 0.002	0.757 ± 0.055	0.741 ± 0.012	0.577 ± 0.078
0.05-Fair GR	1.003 ± 0.004	0.039 ± 0.006	0.772 ± 0.027	0.744 ± 0.016	0.538 ± 0.066
Regular GR	0.977 ± 0.003	0.021 ± 0.007	0.814 ± 0.024	0.821 ± 0.01	0.604 ± 0.08

results to those in [ALZ21], even when all features are known. More deeply understanding why our method does not work on the real-world datasets is an important and interesting future work; we would like to analyze the modularity of the real-world networks, as well as the distribution of node degrees, labels, and features across groups to diagnose sources of failure.

Table 2.2: Test accuracy (**Acc**) and statistical parity (**SP**) of all models averaged over all 25 combinations of unknown feature rates of $\{0.1, 0.3, 0.5, 0.7, 0.9\}$ for each group in German credit.

Method	$\text{Acc}_{\text{linear}} \uparrow$	$\text{Acc}_{\text{mlp}} \uparrow$	$\text{Acc}_{\text{gcn}} \uparrow$	$\text{SP}_{\text{linear}} \downarrow$	$\text{SP}_{\text{mlp}} \downarrow$	$\text{SP}_{\text{gcn}} \downarrow$
0.0-Fair GM	0.700 ± 0.006	0.707 ± 0.003	0.699 ± 0.002	0.051 ± 0.016	0.028 ± 0.004	0.011 ± 0.01
0.025-Fair GM	0.705 ± 0.003	0.707 ± 0.003	0.698 ± 0.002	0.044 ± 0.012	0.028 ± 0.007	0.02 ± 0.026
0.05-Fair GM	0.704 ± 0.007	0.708 ± 0.004	0.697 ± 0.002	0.034 ± 0.009	0.029 ± 0.003	0.013 ± 0.005
Regular GM	0.701 ± 0.002	0.708 ± 0.004	0.699 ± 0.001	0.043 ± 0.01	0.025 ± 0.005	0.006 ± 0.005
0.0-Fair NM	0.699 ± 0.005	0.706 ± 0.003	0.697 ± 0.003	0.053 ± 0.015	0.033 ± 0.007	0.007 ± 0.006
0.025-Fair NM	0.7 ± 0.006	0.706 ± 0.003	0.697 ± 0.002	0.041 ± 0.009	0.037 ± 0.007	0.015 ± 0.013
0.05-Fair NM	0.7 ± 0.006	0.706 ± 0.003	0.697 ± 0.002	0.046 ± 0.013	0.033 ± 0.005	0.01 ± 0.003
Regular NM	0.7 ± 0.003	0.708 ± 0.001	0.698 ± 0.001	0.044 ± 0.02	0.034 ± 0.007	0.016 ± 0.007
0.0-Fair FP	0.694 ± 0.021	0.713 ± 0.011	0.704 ± 0.007	0.012 ± 0.012	0.025 ± 0.025	0.019 ± 0.039
0.025-Fair FP	0.708 ± 0.012	0.706 ± 0.012	0.689 ± 0.027	0.024 ± 0.028	0.026 ± 0.026	0.026 ± 0.057
0.05-Fair FP	0.702 ± 0.03	0.706 ± 0.011	0.708 ± 0.046	0.078 ± 0.01	0.023 ± 0.026	0 ± 0
Regular FP	0.7 ± 0.024	0.708 ± 0.012	0.698 ± 0.01	0.063 ± 0.069	0.03 ± 0.02	0.007 ± 0.01
0.0-Fair GR	0.698 ± 0.005	0.703 ± 0.004	0.698 ± 0.001	0.04 ± 0.02	0.021 ± 0.003	0.006 ± 0.007
0.025-Fair GR	0.702 ± 0.004	0.702 ± 0.002	0.7 ± 0.001	0.041 ± 0.014	0.025 ± 0.002	0.008 ± 0.01
0.05-Fair GR	0.699 ± 0.004	0.703 ± 0.003	0.699 ± 0.003	0.034 ± 0.01	0.024 ± 0.004	0.005 ± 0.005
Regular GR	0.697 ± 0.003	0.703 ± 0.003	0.7 ± 0.001	0.038 ± 0.017	0.027 ± 0.005	0.01 ± 0.009

2.8 Conclusion

In this chapter, we prove that a higher discrimination risk can amplify the unfairness of a ML model applied to imputed data. We formalize a general graph feature imputation framework called mean aggregation imputation and theoretically and empirically characterize graphs in which applying the framework can yield a high discrimination risk. We propose a simple and effective solution to ensure mean aggregation-imputed features provably have a low discrimination risk, while minimally sacrificing reconstruction error (with respect to the imputation objective).

2.9 Broader Impacts

Our analysis and solution, like many fair ML algorithms, assume that groups are discrete and that group membership is known and static, which is not true in reality [Sub22, DMO21, Bar10]. Furthermore, we don't consider fairness at the intersections of different groups [KNR18, BG18b], or operationalizations of fairness beyond the parity of two non-overlapping groups [JW21]. Furthermore, while fairness is often framed as sufficient for the creation of ethical systems, this is often not the case. For instance, ϵ -fair mean aggregation imputation may be used to train a “fairer” model that diversifies news recommendations to social media users [LWZ21], but this model could recommend hostile or intolerant news sources to LGBTQIA+ users and cause psychological harm [Sub22].

CHAPTER 3

Networked Inequality: Preferential Attachment Bias in Graph Neural Network Link Prediction

3.1 Introduction

Link prediction (LP) using GNNs is increasingly leveraged to recommend friends in social networks [FML19, SLY21], as well as by scholarly tools to recommend academic literature in citation networks [XZH21]. In recent years, graph learning researchers have raised concerns about the unfairness of GNN LP [LWZ21, CHG22, LWN22]. This unfairness is often attributed to graph structure, including the stratification of social groups; for example, online networks are usually segregated by ethnicity [HCT17]. However, most fair GNN LP research has focused on dyadic fairness, i.e., satisfying some notion of parity between inter-group and intra-group link predictions. This formulation neglects: 1) LP dynamics within social groups [KA21]; and 2) the “rich get richer” effect, i.e., the prediction of links at a higher rate with high-degree nodes [BA99]. In the context of friend recommendation systems, the “rich get richer” effect can increase the number of links formed with high-degree individuals, which boosts their influence on other individuals in the network, and thus their power [BFS23].

In this chapter, we shed light on how degree bias in networks affects GCN LP [KW17]. We theoretically and empirically find that GCNs with a symmetric normalized graph filter have a within-group preferential attachment (PA) bias in LP. Specifically, GCNs often output LP scores that are approximately proportional to the geometric mean of the (within-group) degrees of the incident nodes when the nodes belong to the same social group. (We elaborate

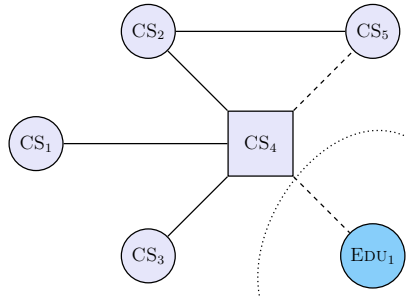


Figure 3.1: An academic collaboration network where nodes are Computer Science (CS) and Education (EDU) researchers, solid edges are current or past collaborations, and dashed edges are collaborations recommended by a GCN. Circular nodes are women and square nodes are men.

on PA and our motivation in the appendix of [SSS24].) We focus on GCNs with symmetric and random walk normalized graph filters because they are popular architectures for graph deep learning, and they provide us with a reasonable setting to develop a rigorous theory of PA bias in GNN LP while leveraging tools from spectral graph theory.

Our finding can have significant implications for the fairness of GCN LP. For example, consider links within the CS social group in the toy academic collaboration network in Figure 3.1. Because men in CS, on average, have a higher within-group degree ($\text{deg} = 3$) than women in CS ($\text{deg} = 1.25$), due to gender discrimination, a collaboration recommender system that uses a GCN can suggest men as collaborators at a higher rate. This has the detrimental effect of further concentrating research collaborations among men, thereby reducing the influence of women in CS and reinforcing their marginalization in the field [YF22]. Furthermore, considering this marginalization in the context of CS is important, as such marginalization may be less severe or different in EDU.

Our contributions are as follows:

1. We theoretically uncover that GCNs with a symmetric normalized graph filter have a within-group PA bias in LP (§3.4.1). We validate our theoretical analysis on diverse

real-world network datasets (e.g., citation, collaboration, online social networks) of varying size (§3.6.1). In doing so, we lay a foundation to study this previously-unexplored PA bias in the GNN setting.

2. We bridge GCN’s PA bias with unfairness in LP (§3.4.2, §3.6.2). We contribute a new within-group fairness metric for LP, which quantifies disparities in LP scores within social groups, towards combating the amplification of degree and power disparities. To our knowledge, we are the first to study the within-group fairness of GNNs.
3. We propose a training-time strategy to alleviate within-group unfairness (§3.5), and we assess its effectiveness on citation, online social, and credit networks (§3.6.3). Our experiments reveal that even for this new form of unfairness, simple regularization approaches can be successful.

3.2 Related Work

Degree Bias in GNNs Numerous papers have investigated how GNN performance is degraded for low-degree nodes on node representation learning and classification tasks [TYS20, LNF21, KZX22, XXH23, SJW23]. [LNF23] presents a generalized notion of degree bias that considers different multi-hop structures around nodes and proposes a framework to address it; in contrast to prior work, which focuses on *degree equal opportunity* (i.e., similar accuracy for nodes with the same degree), [LNF23] also studies *degree statistical parity* (i.e., similar prediction rates of each class for nodes with the same degree). Beyond node classification, [WD22] finds GNN LP performance disparities across nodes with different degrees: low-degree nodes often benefit from higher performance than high-degree nodes. In this chapter, we find that GCNs have a PA bias in LP, and present a new fairness metric which quantifies disparities in GNN LP scores within social groups. We focus on *group fairness* (i.e., parity between groups) rather than *individual fairness* (i.e., treating similar individuals similarly); this is because producing similar LP scores for similar-degree individuals does not

prevent high-degree individuals from unfairly amassing links, and thus power (see Figure 3.1). We further compare our work to prior degree bias works in the appendix of [SSS24].

Fair Link Prediction Prior work has investigated the unfairness of GNN LP [LWZ21, CHG22, LWN22], often attributing it to graph structure, (e.g., stratification of social groups). However, most of this research has focused on dyadic fairness, i.e., satisfying some notion of parity between inter-group and intra-group links. Like [WD22], we examine how degree bias impacts GNN LP; however, rather than focus on performance disparities across nodes with different degrees, we study GCN’s PA bias and LP score disparities across (sub)groups.

Within-Group Fairness Much previous work has studied within-group fairness, i.e., fairness over social subgroups (e.g., Black women, Indigenous men) defined over multiple axes (e.g., race, gender) [KNR18, FIK20, GGR21, WRR22]. The motivation of this chapter is that classifiers can be fair with respect to two social axes separately, but be unfair to subgroups defined over both these axes. While prior research has termed this phenomenon *intersectional* unfairness, we opt for *within-group* unfairness to distinguish it from the critical framework of Intersectionality [OSG23]. We study within-group fairness in the GNN setting. In particular, our theoretical and empirical findings reveal that GCN LP can further marginalize social subgroups; this relates to the “complexity” tenet of Intersectionality, which expresses that the marginalization faced by, e.g., Black women, is non-additive and distinct from the marginalization faced by Black men and white women [CB20].

Bias and Power in Networks A wealth of literature outside fair graph learning has examined how network structure enables discrimination and disparities in capital [FBB19, SHC20, ZHM21, BFS23]. [BLM14] describes how an individual’s position in a social network affects their access to jobs and public health information, as well as how they are surveilled. [SRC18] observes that high-degree accounts on Instagram overwhelmingly belong to men and recommendation algorithms further boost these accounts; complementarily, the authors

find that even a simple, random walk-based recommendation algorithm can amplify degree disparities between social groups in networks modeled by PA dynamics. Similarly, we investigate how GCN LP can amplify degree disparities in networks and further concentrate power among high-degree individuals.

3.3 Preliminaries

We have a simple, undirected n -node graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with doubly-weighted self-loops. The nodes have features $(\mathbf{x}_i)_{i \in \mathcal{V}}$, with each $\mathbf{x}_i \in \mathbb{R}^d$. We denote the adjacency matrix of \mathcal{G} as $\mathbf{A} \in \{0, 1\}^{n \times n}$ and the degree matrix as $\mathbf{D} = \text{diag} \left(\left(\sum_{j \in \mathcal{V}} \mathbf{A}_{ij} \right)_{i \in \mathcal{V}} \right)$, with $\mathbf{D} \in \mathbb{N}^{n \times n}$. We consider two L -layer GCN encoders: (1) $\Phi_s : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times d'}$ [KW17], which uses a symmetric normalized filter, and (2) $\Phi_r : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times d'}$, which uses a random walk normalized filter. Φ_s and Φ_r compute node representations as, $\forall i \in \mathcal{V}$:

$$\Phi_s \left((\mathbf{x}_j)_{j \in \mathcal{V}} \right)_i = \mathbf{s}_i^{(L)}, \quad \Phi_r \left((\mathbf{x}_j)_{j \in \mathcal{V}} \right)_i = \mathbf{r}_i^{(L)} \quad (3.1)$$

$$\forall l \in [L], \mathbf{s}_i^{(l)} = \sigma^{(l)} \left(\sum_{j \in \Gamma(i)} \frac{\mathbf{W}_s^{(l)} \mathbf{s}_j^{(l-1)}}{\sqrt{\mathbf{D}_{ii} \mathbf{D}_{jj}}} \right), \quad \mathbf{r}_i^{(l)} = \sigma^{(l)} \left(\sum_{j \in \Gamma(i)} \frac{\mathbf{W}_r^{(l)} \mathbf{r}_j^{(l-1)}}{\mathbf{D}_{ii}} \right), \quad (3.2)$$

where $(\mathbf{s}_i^{(0)})_{i \in \mathcal{V}} = (\mathbf{r}_i^{(0)})_{i \in \mathcal{V}} = (\mathbf{x}_i)_{i \in \mathcal{V}}$; $\Gamma(i)$ is the 1-hop neighborhood of i ; $\mathbf{W}_s^{(l)}$ and $\mathbf{W}_r^{(l)}$ are the weight matrices corresponding to layer l of Φ_s and Φ_r , respectively; for $l \in [L-1]$, $\sigma^{(l)}$ is a ReLU non-linearity; and $\sigma^{(L)}$ is the identity function. We now consider the first-order Taylor expansions of Φ_s and Φ_r around $(\mathbf{0})_{i \in \mathcal{V}}$:

$$\mathbf{s}_i^{(L)} = \sum_{j \in \mathcal{V}} \left[\frac{\partial \mathbf{s}_i^{(L)}}{\partial \mathbf{x}_j} \right] \mathbf{x}_j + \xi \left(\mathbf{s}_i^{(L)} \right), \quad \mathbf{r}_i^{(L)} = \sum_{j \in \mathcal{V}} \left[\frac{\partial \mathbf{r}_i^{(L)}}{\partial \mathbf{x}_j} \right] \mathbf{x}_j + \xi \left(\mathbf{r}_i^{(L)} \right), \quad (3.3)$$

where ξ is the error of the first-order approximations. This error is low when $(\mathbf{x}_i)_{i \in \mathcal{V}}$ are close to $\mathbf{0}$, which we validate empirically in §3.6.1. Furthermore, we consider an inner-product LP score function $f_{LP} : \mathbb{R}^{d'} \times \mathbb{R}^{d'} \rightarrow \mathbb{R}$:

$$f_{LP} \left(\mathbf{h}_i^{(L)}, \mathbf{h}_j^{(L)} \right) = \left(\mathbf{h}_i^{(L)} \right)^\top \mathbf{h}_j^{(L)}, \quad (3.4)$$

where $\mathbf{h}_i^{(L)}$ is the last-layer representation for node i . While it is common to use a vanilla GCN and inner-product score function for LP¹, researchers have proposed methods to improve the expressivity of node representations for LP by capturing subgraph information [ZC18, LWW20, CSR23]. Our theoretical findings remain relevant to methods that ultimately use a GCN to predict links (e.g., [ZC18, LWW20]), as we do not make assumptions about the features passed to the GCN (i.e., they could be distance encodings, SEAL node embeddings, etc.) Our results may also generalize to GNN architectures that use a degree-normalized graph filter, e.g., Graph Attention Networks [VCC18]. Studying the fairness of more expressive LP methods is an interesting direction for future research. Furthermore, although we only consider an inner-product LP score function in our theoretical analysis, we also run experiments with a Hadamard product and MLP score function (see Appendix B.7.2), and we find that our theoretical analysis is still relevant to and reasonably supports the experimental results.

3.4 Theoretical Analysis

We leverage spectral graph theory to study how degree bias affects GCN LP. Theoretically, we find that GCNs with a symmetric normalized graph filter have a within-group PA bias (§3.4.1), but GCNs with a random walk normalized filter may lack such a bias (§3.4.3). We further bridge GCN’s PA bias with unfairness in GCN LP, proposing a new LP within-group fairness metric (§3.4.2) and a simple training-time strategy to alleviate unfairness (§3.5). We empirically validate our theoretical results and fairness strategy in §3.6. We provide proofs for all theoretical results in Appendix B.1.

Our ultimate goal is to bound the expected LP scores $\mathbb{E} \left[f_{LP} \left(\mathbf{s}_i^{(L)}, \mathbf{s}_j^{(L)} \right) \right]$ and $\mathbb{E} \left[f_{LP} \left(\mathbf{r}_i^{(L)}, \mathbf{r}_j^{(L)} \right) \right]$ for nodes i, j in the same social group in terms of the degrees of i, j . We begin with Lemma 3.4.1, which expresses GCN representations (in expectation) as a linear combination of the initial node features. In doing so, we decouple the computation of GCN

¹https://github.com/pyg-team/pytorch_geometric/blob/master/examples/link_pred.py

representations from the non-linearities $\sigma^{(l)}$.

Lemma 3.4.1. *Similarly to [XLT18], assume that each path from node $i \rightarrow j$ in the computation graph of Φ_s is independently activated with probability $\rho_s(i)$, and similarly, $\rho_r(i)$ for Φ_r (see Appendix B.1.1 for a discussion of this assumption). Furthermore, suppose that $\mathbb{E} \left[\xi \left(\mathbf{s}_i^{(L)} \right) \right] = \mathbb{E} \left[\xi \left(\mathbf{r}_i^{(L)} \right) \right] = \mathbf{0}$, where the expectations are taken over the probability distributions of paths activating. We define $\alpha_j = \left(\prod_{l=L}^1 \mathbf{W}_s^{(l)} \right) \mathbf{x}_j$, and $\beta_j = \left(\prod_{l=L}^1 \mathbf{W}_r^{(l)} \right) \mathbf{x}_j$. Then, $\forall i \in \mathcal{V}$:*

$$\mathbb{E} \left[\mathbf{s}_i^{(L)} \right] = \sum_{j \in \mathcal{V}} \rho_s(i) \left(\mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} \right)_{ij}^L \alpha_j, \quad \mathbb{E} \left[\mathbf{r}_i^{(L)} \right] = \sum_{j \in \mathcal{V}} \rho_r(i) \left(\mathbf{D}^{-1} \mathbf{A} \right)_{ij}^L \beta_j. \quad (3.5)$$

Lemma 3.4.1 demonstrates that under certain assumptions (which we show to be reasonable in §3.6.1), the expected GCN representations can be expressed as a linear combination of the node features that depends on a normalized version of the adjacency matrix.

We now introduce social groups in \mathcal{G} into our analysis. Suppose that \mathcal{V} can be partitioned into B disjoint social groups $\{S^{(b)}\}_{b \in [B]}$, such that $\bigcup_{b \in [B]} S^{(b)} = \mathcal{V}$ and $\bigcap_{b \in [B]} S^{(b)} = \emptyset$. Furthermore, we define $\mathcal{G}^{(b)}$ as the induced connected subgraph of \mathcal{G} formed from $S^{(b)}$. (If a group comprises $C > 1$ connected components, it can be treated as C separate groups.) Let $\hat{\mathbf{A}}$ be a within-group adjacency matrix that contains links between nodes in the same group, i.e., $\hat{\mathbf{A}}$ contains the link (i, j) if and only if for some group $S^{(b)}$, $i, j \in S^{(b)}$. Without loss of generality, we reorder the rows and columns of $\hat{\mathbf{A}}$ and \mathbf{A} such that $\hat{\mathbf{A}}$ is a block matrix. Let $\hat{\mathbf{D}}$ be the degree matrix of $\hat{\mathbf{A}}$.

3.4.1 Symmetric Normalized Filter

We first focus on analyzing Φ_s . We introduce the notation $\mathbf{P} = \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}}$ for the symmetric normalized adjacency matrix. We further define $\hat{\mathbf{P}} = \hat{\mathbf{D}}^{-\frac{1}{2}} \hat{\mathbf{A}} \hat{\mathbf{D}}^{-\frac{1}{2}}$, which

has the form $\begin{bmatrix} \hat{\mathbf{P}}^{(1)} & & \mathbf{0} \\ & \ddots & \\ \mathbf{0} & & \hat{\mathbf{P}}^{(B)} \end{bmatrix}$. Each $\hat{\mathbf{P}}^{(b)}$ admits the orthonormal spectral decom-

position $\widehat{\mathbf{P}}^{(b)} = \sum_{k=1}^{|S^{(b)}|} \lambda_k^{(b)} \mathbf{v}_k^{(b)} \left(\mathbf{v}_k^{(b)}\right)^\top$. Let $\left(\lambda_k^{(b)}\right)_{1 \leq k \leq |S^{(b)}|}$ be the eigenvalues of $\widehat{\mathbf{P}}^{(b)}$ sorted in non-increasing order; the eigenvalues fall in the range $(-1, 1]$. By the spectral properties of $\widehat{\mathbf{P}}^{(b)}$, $\lambda_1^{(b)} = 1$. Following [Lov96], we denote the *spectral gap* of $\widehat{\mathbf{P}}^{(b)}$ as $\lambda^{(b)} = \max \left\{ \lambda_2^{(b)}, \left| \lambda_{|S^{(b)}|}^{(b)} \right| \right\} < 1$; $\lambda_2^{(b)}$ corresponds to the smallest non-zero eigenvalue of the symmetric normalized graph Laplacian. Let $\mathbf{P} = \widehat{\mathbf{P}} + \Xi^{(0)}$. If \mathcal{G} is highly modular or approximately disconnected, then $\Xi^{(0)} \approx \mathbf{0}$, albeit with positive and non-positive entries. Finally, we define the volume $\text{vol}(\mathcal{G}^{(b)}) = \sum_{k \in S^{(b)}} \widehat{\mathbf{D}}_{kk}$.

In Lemma 3.4.2, we present an inequality for the entries of \mathbf{P}^L in terms of the spectral properties of $\widehat{\mathbf{P}}$. We can then combine this inequality with Lemma 3.4.1 to bound $\mathbb{E} \left[\mathbf{s}_i^{(L)} \right]$, and subsequently $\mathbb{E} \left[f_{LP} \left(\mathbf{s}_i^{(L)}, \mathbf{s}_j^{(L)} \right) \right]$.

Lemma 3.4.2. *For $i, j \in S^{(b)}$:*

$$\left| \mathbf{P}_{ij}^L - \frac{\sqrt{\widehat{\mathbf{D}}_{ii} \widehat{\mathbf{D}}_{jj}}}{\text{vol}(\mathcal{G}^{(b)})} \right| \leq \zeta_s = (\lambda^{(b)})^L + \sum_{l=1}^L \binom{L}{l} \|\Xi^{(0)}\|_{op}^l \|\widehat{\mathbf{P}}\|_{op}^{L-l}, \quad (3.6)$$

where $\|\cdot\|_{op}$ is the operator norm. And for $i \in S^{(b)}, j \notin S^{(b)}$, $|\mathbf{P}_{ij}^L - 0| \leq \sum_{l=1}^L \binom{L}{l} \|\Xi^{(0)}\|_{op}^l \|\widehat{\mathbf{P}}\|_{op}^{L-l} \leq \zeta_s$.

The proof of Lemma 3.4.2 is similar to spectral proofs of random walk convergence. When L is small (e.g., 2 for many GCNs [KW17]) and $\|\Xi^{(0)}\|_{op} \approx 0$, $\sum_{l=1}^L \binom{L}{l} \|\Xi^{(0)}\|_{op}^l \|\widehat{\mathbf{P}}\|_{op}^{L-l} \approx 0$. Furthermore, with significant stratification between social groups [HCT17] and high expansion within groups [MM11, LLD09], $\lambda^{(b)} \ll 1$. In this case, $\zeta_s \approx 0$ and $\mathbf{P}_{ij}^L \approx \frac{\sqrt{\widehat{\mathbf{D}}_{ii} \widehat{\mathbf{D}}_{jj}}}{\text{vol}(\mathcal{G}^{(b)})}$ for $i, j \in S^{(b)}$. Combining Lemmas 3.4.1 and 3.4.2, Φ_s can oversmooth the expected representations to $\mathbb{E} \left[\mathbf{s}_i^{(L)} \right] \approx \rho_s(i) \sqrt{\widehat{\mathbf{D}}_{ii}} \cdot \sum_{j \in S^{(b)}} \frac{\sqrt{\widehat{\mathbf{D}}_{jj}}}{\text{vol}(\mathcal{G}^{(b)})} \alpha_j$ [Ker22, GRC23]. We use this knowledge to bound $\mathbb{E} \left[f_{LP} \left(\mathbf{s}_i^{(L)}, \mathbf{s}_j^{(L)} \right) \right]$ in terms of the degrees of i, j .

Theorem 3.4.3. *Following a relaxed assumption from [XLT18], for nodes $i, j \in S^{(b)}$, we*

assume that $\rho_s(i) = \rho_s(j) = \bar{\rho}_s(b)$. Then:

$$\left| \mathbb{E} \left[f_{LP} \left(\mathbf{s}_i^{(L)}, \mathbf{s}_j^{(L)} \right) \right] - C_0 \sqrt{\widehat{\mathbf{D}}_{ii} \widehat{\mathbf{D}}_{jj}} \right| \leq \zeta_s \bar{\rho}_s^2(b) \left(\sqrt{\widehat{\mathbf{D}}_{ii}} + \sqrt{\widehat{\mathbf{D}}_{jj}} \right) C_1 C_2 + \zeta_s^2 \bar{\rho}_s^2(b) C_2^2, \quad (3.7)$$

where: (3.8)

$$C_0 = \bar{\rho}_s^2(b) C_1^2, \quad C_1 = \left\| \sum_{k \in S^{(b)}} \frac{\sqrt{\widehat{\mathbf{D}}_{kk}}}{\text{vol}(\mathcal{G}^{(b)})} \alpha_k \right\|_2, \quad C_2 = \sum_{k \in \mathcal{V}} \|\alpha_k\|_2. \quad (3.9)$$

In simpler terms, Theorem 3.4.3 states that with social stratification and expansion, the expected LP score $\mathbb{E} \left[f_{LP} \left(\mathbf{s}_i^{(L)}, \mathbf{s}_j^{(L)} \right) \right] \propto \sqrt{\widehat{\mathbf{D}}_{ii} \widehat{\mathbf{D}}_{jj}}$ approximately when i, j belong to the same social group. This is because, as explained before Theorem 3.4.3, $\zeta_s \cong 0$, so the RHS of the bound is $\cong 0$. This demonstrates that in LP, GCNs with a symmetric normalized graph filter have a within-group PA bias. If Φ_s positively influences the formation of links over time, this PA bias can drive “rich get richer” dynamics within social groups [SRC18]. As shown in Figure 3.1 and §3.4.2, such “rich get richer” dynamics can engender group unfairness when nodes’ degrees are statistically associated with their group membership (§3.4.2). An association between node degree and group membership depends on group size and homophily; in particular, when a group has many nodes and intra-links (i.e., is homophilous), there may be more nodes with a high within-group degree. Beyond fairness, Theorem 3.4.3 reveals that GCNs do not align with theories that *social rank* influences link formation, i.e., the likelihood of a link forming between nodes is proportional to their degree *difference* [GSL18].

3.4.2 Within-Group Fairness

We further investigate the fairness implications of the PA bias of Φ_s in LP. We first introduce an additional set of social groups. Suppose that \mathcal{V} can also be partitioned into D disjoint social groups $\{T^{(d)}\}_{d \in [D]}$; then, we can consider intersections of $\{S^{(b)}\}_{b \in [B]}$ and $\{T^{(d)}\}_{d \in [D]}$. For example, revisiting Figure 3.1, S may correspond to academic discipline (e.g., CS, EDU) and T may correspond to gender (e.g., men, women). For simplicity, we let $D = 2$. We

measure the unfairness $\Delta^{(b)} : \mathbb{R}^{d'} \times \mathbb{R}^{d'} \rightarrow \mathbb{R}$ of LP for group b as:

$$\Delta^{(b)} \left(\mathbf{h}_i^{(L)}, \mathbf{h}_j^{(L)} \right) := \left| \mathbb{E}_{i,j \sim U((S^{(b)} \cap T^{(1)}) \times S^{(b)})} f_{LP} \left(\mathbf{h}_i^{(L)}, \mathbf{h}_j^{(L)} \right) - \mathbb{E}_{i,j \sim U((S^{(b)} \cap T^{(2)}) \times S^{(b)})} f_{LP} \left(\mathbf{h}_i^{(L)}, \mathbf{h}_j^{(L)} \right) \right|, \quad (3.10)$$

where $U(\cdot)$ is a discrete uniform distribution over the input set. $\Delta^{(b)}$ quantifies disparities in GCN LP scores within $S^{(b)}$ (with respect to $T^{(1)}$ and $T^{(2)}$). In other words, $\Delta^{(b)}$ measures differences in how GCNs allocate LP scores across subgroups, i.e., are links with nodes in one subgroup predicted at a higher rate than links with nodes in the other subgroup? Our metric is motivated by how GNN link predictions influence real-world link formation (e.g., GNN-based recommender systems use LP scores to rank suggested social connections), which has consequences for degree and power disparities. Based on Theorem 3.4.3 and Appendix B.2.1, when $\zeta_s \cong 0$, we can estimate $\Delta^{(b)} \left(\mathbf{s}_i^{(L)}, \mathbf{s}_j^{(L)} \right)$ as:

$$\widehat{\Delta}^{(b)} \left(\mathbf{s}_i^{(L)}, \mathbf{s}_j^{(L)} \right) = \frac{\bar{\rho}_s^2(b)}{|S^{(b)}|} \left\| \sum_{k \in S^{(b)}} \frac{\sqrt{\widehat{\mathbf{D}}_{kk}}}{\text{vol}(\mathcal{G}^{(b)})} \alpha_k \right\|_2^2 \quad (3.11)$$

$$\times \left| \sum_{j \in S^{(b)}} \sqrt{\widehat{\mathbf{D}}_{jj}} \underbrace{\left(\mathbb{E}_{i \sim U(S^{(b)} \cap T^{(1)})} \sqrt{\widehat{\mathbf{D}}_{ii}} - \mathbb{E}_{i \sim U(S^{(b)} \cap T^{(2)})} \sqrt{\widehat{\mathbf{D}}_{ii}} \right)}_{\text{degree disparity}} \right|. \quad (3.12)$$

This suggests that a large disparity in the degree of nodes in $S^{(b)} \cap T^{(1)}$ vs. $S^{(b)} \cap T^{(2)}$ can greatly increase the unfairness $\Delta^{(b)}$ of Φ_s LP. For example, in Figure 3.1, the large degree disparity within CS (between men and women) entails that a GCN collaboration recommender system applied to the network will have a large $\Delta^{(b)}$. We empirically validate these fairness implications on diverse network datasets in §3.6.2. While we consider pre-activation LP scores in Eqn. 3.10 (in line with prior work, e.g., [LWZ21]), we consider post-sigmoid scores $\sigma \left(f_{LP} \left(\mathbf{h}_i^{(L)}, \mathbf{h}_j^{(L)} \right) \right)$ (where σ is the sigmoid function) in §3.6.2 and §3.6.3, as this simulates how LP scores may be processed in practice.

Ultimately, within-group unfairness is characteristic of all GNN link prediction methods that: (1) predict scores for links with magnitudes that are positively associated with the degrees of their incident nodes, and (2) are applied to graphs where within-group membership

is associated with node degree.

3.4.3 Random Walk Normalized Filter

We now follow similar steps as with Φ_s to understand how degree bias affects LP scores for Φ_r . We redefine $\mathbf{P} = \mathbf{D}^{-1}\mathbf{A}$, $\widehat{\mathbf{P}} = \widehat{\mathbf{D}}^{-1}\widehat{\mathbf{A}}$, and the remaining notation from §3.4.1 accordingly for the random walk setting.

Theorem 3.4.4. *Let $\zeta_r = \max_{u,v \in \mathcal{V}} \sqrt{\frac{\widehat{\mathbf{D}}_{vv}}{\widehat{\mathbf{D}}_{uu}}} (\lambda^{(b)})^L + \sum_{l=1}^L \binom{L}{l} \|\Xi^{(0)}\|_{op}^l \|\widehat{\mathbf{P}}\|_{op}^{L-l}$. Furthermore, for nodes $i, j \in S^{(b)}$, assume that $\rho_r(i) = \rho_r(j) = \bar{\rho}_r(b)$. Combining Lemmas 3.4.1 and B.1.1 (in the appendix):*

$$\left| \mathbb{E} \left[f_{LP} \left(\mathbf{r}_i^{(L)}, \mathbf{r}_j^{(L)} \right) \right] - C_0 \right| \leq \zeta_r \bar{\rho}_r^2(b) C_1 C_2 + \zeta_r^2 \bar{\rho}_r^2(b) C_2^2, \quad (3.13)$$

$$\text{where:} \quad (3.14)$$

$$C_0 = \bar{\rho}_r^2(b) C_1^2, \quad C_1 = \left\| \sum_{k \in S^{(b)}} \frac{\widehat{\mathbf{D}}_{kk}}{\text{vol}(\mathcal{G}^{(b)})} \beta_k \right\|, \quad C_2 = \sum_{k \in \mathcal{V}} \|\beta_k\|_2. \quad (3.15)$$

In other words, if $\zeta_r \cong 0$, $\mathbb{E} \left[f_{LP} \left(\mathbf{r}_i^{(L)}, \mathbf{r}_j^{(L)} \right) \right]$ is approximately constant when i, j belong to the same social group. Based on Theorem 3.4.4 and Appendix B.2.2, we can estimate $\Delta^{(b)} \left(\mathbf{s}_i^{(L)}, \mathbf{s}_j^{(L)} \right)$ as $\widehat{\Delta}^{(b)} \left(\mathbf{s}_i^{(L)}, \mathbf{s}_j^{(L)} \right) = 0$. Theoretically, this would suggest that a large disparity in the degree of nodes in $S^{(b)} \cap T^{(1)}$ vs. $S^{(b)} \cap T^{(2)}$ does not increase the unfairness $\Delta^{(b)}$ of Φ_r LP. However, we find empirically that this is not the case (§3.6.1). Even so, we include theoretical results for the random walk filter to be more comprehensive with respect to filter choice, as well as be upfront about the limitations of our analysis in this case. We also seek to provide an example of how to apply our analysis to other filters, for researchers who would like to build on it in the future. For example, findings for the random walk filter could be relevant to the GAT filter [VCC18], which is also a row-stochastic matrix.

In summary, in §3.4, we build on prior analysis techniques for random walks and GNNs. At a high level, we: (1) simplify the GCN architecture to be a linear function by truncating

its Taylor expansion and considering node representations in expectation; (2) analyze the convergence of node representations via a spectral analysis of the convergence of short random walks within subgraphs (corresponding to social groups); and (3) use norm inequalities to estimate link prediction scores. Our analysis comprises numerous novel elements including:

1. Analyzing the convergence of random walks within subgraphs, which requires accounting for the rate at which probability mass escapes from the subgraphs. In contrast, random walk results in the literature usually concern the convergence of random walks over an entire graph.
2. Uncovering properties of short random walks on graphs, since most GNNs are shallow. In contrast, random walk results in the literature often concern the stationary distribution of random walks.
3. Concretely relating theoretical properties of random walks to the fairness of GCN link prediction.

3.5 Fairness Regularizer

We propose a simple training-time solution to alleviate within-group LP unfairness regardless of graph filter type and GNN architecture. In particular, we can add a fairness regularization term $\mathcal{L}_{\text{fair}}$ to our original GNN training loss [KAS11]:

$$\mathcal{L}_{\text{new}} = \mathcal{L}_{\text{orig}} + \lambda_{\text{fair}} \mathcal{L}_{\text{fair}} = \mathcal{L}_{\text{orig}} + \frac{\lambda_{\text{fair}}}{B} \sum_{b \in [B]} \Delta^{(b)}, \quad (3.16)$$

where λ_{fair} is a tunable hyperparameter that for higher values, pushes the GNN to learn fairer parameters. With our fairness strategy, we empirically observe a significant decrease in the average unfairness across groups $\frac{1}{B} \sum_{b \in [B]} \Delta^{(b)}$ without a severe drop in LP performance for GCN (§3.6.3).

3.6 Experiments

In this section, we empirically validate our theoretical analysis (§3.6.1) and the within-group fairness implications of GCN’s LP PA bias (§3.6.2) on diverse real-world network datasets of varying size. We further find that our simple training-time strategy to alleviate unfairness is effective on citation, online social, and credit networks (§3.6.3). We present experimental results with 4-layer GCN encoders and a Hadamard product with MLP LP score function in Appendix B.7, with similar conclusions.

3.6.1 Validating Theoretical Analysis

We validate our theoretical analysis on 10 real-world network datasets (e.g., citation, collaboration, online social networks), which we describe in Appendix B.3. Each dataset is natively intended for node classification; however, we adapt the datasets for LP, treating the connected components within the node classes as the social groups $S^{(b)}$. This design choice is reasonable, as in all the datasets, the classes naturally correspond to socially-relevant groupings of the nodes, or proxies thereof (e.g., in the LastFMAsia dataset, the classes are the home countries of users). Because we adopt the class labels for each dataset as the social group labels, the social groups are largely homophilic; this aligns with our assumptions when interpreting Theorems 3.4.3 and 3.4.4 that social groups are stratified in networks.

We train GCN encoders Φ_s and Φ_r for LP over 10 random seeds (see Appendix B.5 for more details). In Figure 3.3, we plot the theoretic² LP score that we derive in §3.4 against the GCN LP score *for pairs of test nodes belonging to the same social group* (including positive and negative links). In particular, for Φ_s , the theoretic LP score is $\bar{\rho}_s^2(b) \sqrt{\widehat{\mathbf{D}}_{ii} \widehat{\mathbf{D}}_{jj}} \left\| \sum_{k \in S^{(b)}} \frac{\sqrt{\widehat{\mathbf{D}}_{kk}}}{\text{vol}(\mathcal{G}^{(b)})} \alpha_k \right\|_2^2$ and the GCN LP score is $f_{LP}(\mathbf{s}_i^{(L)}, \mathbf{s}_j^{(L)})$ (see Theorem 3.4.3). In contrast, for Φ_r , the theoretic LP score is $\bar{\rho}_s^2(b) \left\| \sum_{k \in S^{(b)}} \frac{\widehat{\mathbf{D}}_{kk}}{\text{vol}(\mathcal{G}^{(b)})} \beta_k \right\|_2^2$ and the GCN LP score is $f_{LP}(\mathbf{r}_i^{(L)}, \mathbf{r}_j^{(L)})$ (see

²While our theoretic scores resulted from our theoretical analysis in §3.4, we reiterate that our results in §3.4 rely on the assumptions that we state and the theoretic score is not a ground-truth value.

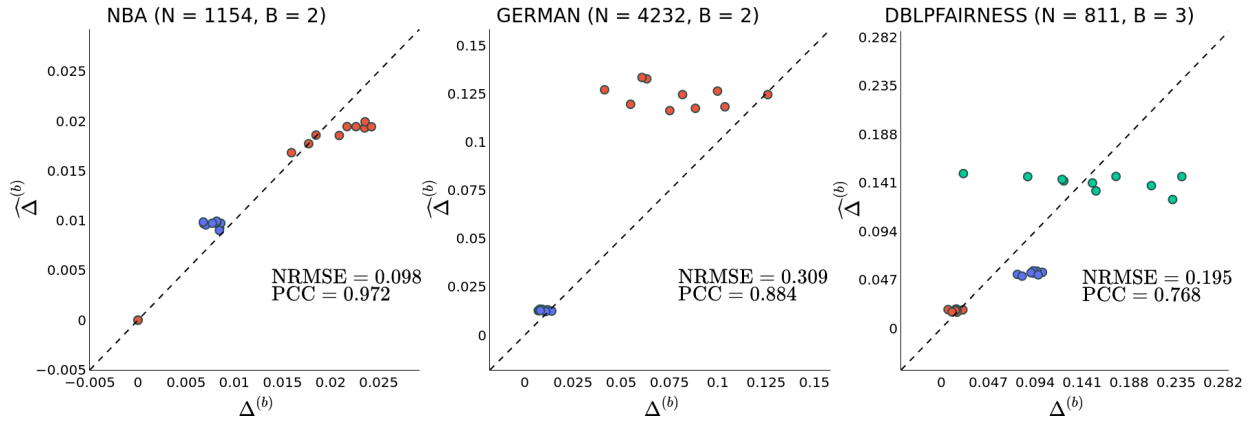


Figure 3.2: The plots display $\widehat{\Delta}^{(b)}$ vs. $\Delta^{(b)}$ for Φ_s for the NBA, German, and DBLP-Fairness datasets over all $b \in [B]$ and 10 random seeds. Each point corresponds to a different random seed, and the color of the point corresponds to the social group $S^{(b)}$. We compute $\widehat{\Delta}^{(b)}$ and $\Delta^{(b)}$ post-sigmoid using only the LP scores over the sampled (positive and negative) test edges. The plots display the NRMSE and PCC of $\widehat{\Delta}^{(b)}$ as a predictor of $\Delta^{(b)}$.

Theorem 3.4.4). For all the datasets, we estimate $\bar{\rho}_s^2(b)$ and $\bar{\rho}_r^2(b)$ separately for each social group $S^{(b)}$ as the slope of the least-squares regression line (through the data from $S^{(b)}$) that predicts the GCN score as a function of the theoretic score. Hence, we do not plot any pair of test nodes that is the only pair in $S^{(b)}$, as it is not possible to estimate $\bar{\rho}_s^2(b)$. Further, the test AUC is consistently high, indicating that the GCNs are well-trained. The large range of each color in the plots indicates a diversity of LP scores within each social group.

We visually observe that the theoretic LP scores are strong predictors of the Φ_s scores for each dataset, validating our theoretical analysis. This strength is further confirmed by the generally low NRMSE and high PCC (except for the EN dataset). However, we observe a few cases in which our theoretical analysis does not line up with our experiments:

1. Our theoretical analysis predicts that the LP score between two nodes i, j that belong to the same social group $S^{(b)}$ will always be non-negative; however, Φ_s can predict negative scores for pairs of nodes in the same social group. In this case, it appears that Φ_s relies more on the dissimilarity of (transformed) features than node degree.

2. For many network datasets (especially from the citation and online social domains), there exist node pairs (near the origin) for which the theoretic LP score underestimates the Φ_s score. Upon further analysis (see Appendix B.8), we find that the theoretic score is less predictive of the Φ_s score for nodes i, j when the product of their degrees (i.e., their PA score) or similarity of their features is relatively low.
3. It appears that the theoretic LP score tends to poorly estimate the Φ_s score when the Φ_s score is relatively high; this suggests that Φ_s may conservatively rely more on the (dis)similarity of node features than node degree when the degree is large.

We do not observe that the theoretic LP scores are strong predictors of the Φ_r scores, although there is still a moderate association between these variables. This could be because the error bound for the theoretic scores for Φ_r , unlike for Φ_s , has an extra dependence $\max_{u,v \in \mathcal{V}} \sqrt{\frac{\bar{D}_{uv}}{\bar{D}_{uu}}}$ on the degrees of the incident nodes (see ζ_r in Theorem 3.4.4). In contrast, the error bound for the theoretic scores for Φ_s (see ζ_s in Theorem 3.4.3) does not depend on this degree ratio. This ratio can be quite large in social networks (e.g., celebrities vs. new users in the Twitter follow network); we further confirm that this ratio is large for our datasets in Appendix B.9.

3.6.2 Within-Group Fairness

We now empirically validate the implications of GCN’s PA bias for within-group unfairness in LP. We run experiments on three network datasets: (1) the NBA social network [DW21], (2) the German credit network [ALZ21], and (3) a new DBLP-Fairness citation network that we construct. We describe these datasets in Appendix B.4, including $\{S^{(b)}\}_{b \in [B]}$ and $\{T^{(d)}\}_{d \in [D]}$.

We train 2-layer GCN encoders Φ_s for LP (see Appendix B.5). In Figure 3.2, for all the datasets, we plot $\widehat{\Delta}^{(b)}$ vs. $\Delta^{(b)}$ (see Eqns. 3.10, 3.12) for each $b \in [B]$. We qualitatively and quantitatively observe that $\widehat{\Delta}^{(b)}$ is moderately predictive of $\Delta^{(b)}$ for each dataset. This confirms our theoretical intuition (§3.4.2) that a large disparity in the degree of nodes in $S^{(b)} \cap T^{(1)}$ vs. $S^{(b)} \cap T^{(2)}$ can greatly increase the unfairness $\Delta^{(b)}$ of Φ_s LP; such unfairness

Table 3.1: $\frac{1}{B} \sum_{b \in [B]} \Delta^{(b)}$ and the test AUC for the NBA, German, and DBLP-Fairness datasets with various settings of λ_{fair} . The **left** table corresponds to Φ_s , and the **right** to Φ_r .

	λ_{fair}	$\frac{1}{B} \sum_{b \in [B]} \Delta^{(b)}$ (\downarrow)	Φ_s Test AUC (\uparrow)		λ_{fair}	$\frac{1}{B} \sum_{b \in [B]} \Delta^{(b)}$ (\downarrow)	Φ_r Test AUC (\uparrow)
NBA	4.0	0.000 \pm 0.001	0.753 \pm 0.002	NBA	4.0	0.000 \pm 0.000	0.585 \pm 0.030
NBA	2.0	0.004 \pm 0.003	0.752 \pm 0.003	NBA	2.0	0.000 \pm 0.000	0.584 \pm 0.032
NBA	1.0	0.007 \pm 0.004	0.752 \pm 0.003	NBA	1.0	0.000 \pm 0.000	0.581 \pm 0.034
NBA	0.0	0.013 \pm 0.005	0.752 \pm 0.003	NBA	0.0	0.000 \pm 0.000	0.583 \pm 0.034
DBLPFAIRNESS	4.0	0.072 \pm 0.018	0.741 \pm 0.008	DBLPFAIRNESS	4.0	0.053 \pm 0.015	0.715 \pm 0.010
DBLPFAIRNESS	2.0	0.095 \pm 0.025	0.756 \pm 0.007	DBLPFAIRNESS	2.0	0.060 \pm 0.016	0.731 \pm 0.009
DBLPFAIRNESS	1.0	0.110 \pm 0.033	0.770 \pm 0.010	DBLPFAIRNESS	1.0	0.065 \pm 0.022	0.746 \pm 0.009
DBLPFAIRNESS	0.0	0.145 \pm 0.020	0.778 \pm 0.007	DBLPFAIRNESS	0.0	0.090 \pm 0.028	0.758 \pm 0.011
GERMAN	4.0	0.012 \pm 0.006	0.876 \pm 0.017	GERMAN	4.0	0.029 \pm 0.011	0.830 \pm 0.024
GERMAN	2.0	0.028 \pm 0.017	0.889 \pm 0.017	GERMAN	2.0	0.031 \pm 0.019	0.843 \pm 0.027
GERMAN	1.0	0.038 \pm 0.016	0.897 \pm 0.014	GERMAN	1.0	0.019 \pm 0.012	0.864 \pm 0.020
GERMAN	0.0	0.045 \pm 0.013	0.912 \pm 0.009	GERMAN	0.0	0.015 \pm 0.005	0.883 \pm 0.009

can amplify degree disparities, worsening power imbalances in the network. Many points deviate from the line of equality; these deviations can be explained by the reasons in §3.6.1 and the compounding of errors.

3.6.3 Fairness Regularizer

We evaluate our solution to alleviate LP unfairness (§3.4.2). In particular, we add our fairness regularization term $\mathcal{L}_{\text{fair}}$ to the original training loss for the 2-layer Φ_s and Φ_r encoders. During each training epoch, we compute $\Delta^{(b)}$ post-sigmoid using only the LP scores over the sampled (positive and negative) training edges. In Table 3.1, we summarize the link prediction fairness ($\frac{1}{B} \sum_{b \in [B]} \Delta^{(b)}$) and performance (test AUC) for the NBA, German, and DBLP-Fairness datasets with various settings of λ_{fair} .

For both graph filter types, we generally observe a significant decrease in $\frac{1}{B} \sum_{b \in [B]} \Delta^{(b)}$ (without a severe drop in test AUC) for $\lambda_{\text{fair}} > 0.0$ over $\lambda_{\text{fair}} = 0.0$ (with the exception of Φ_r for German); however, the varying magnitudes by which $\frac{1}{B} \sum_{b \in [B]} \Delta^{(b)}$ decreases across the

datasets suggests that λ_{fair} may need to be tuned per dataset. As expected, we mostly observe a tradeoff between $\frac{1}{B} \sum_{b \in [B]} \Delta^{(b)}$ and the test AUC as λ_{fair} increases. Our experiments reveal that, regardless of graph filter type, even simple regularization approaches can alleviate this new form of unfairness. As this form of unfairness has not been previously explored, we have no baselines.

Our fairness regularizer can be easily integrated into model training, does not require significant additional computation, and directly optimizes for LP fairness. The time complexity of calculating the regularization term is $\mathcal{O}\left(\sum_{b=1}^B |S^{(b)} \cap T^{(1)}| \cdot |S^{(b)}| + |S^{(b)} \cap T^{(2)}| \cdot |S^{(b)}|\right)$, as we have already computed the LP scores for the cross-entropy loss term and simply need to sum them appropriately with respect to the groups and subgroups. Furthermore, the time complexity of computing gradients for the regularization term is on the same order as backpropagation for the cross-entropy loss term.

However, our fairness regularizer is not applicable in settings where model parameters cannot be retrained or finetuned. Hence, we encourage future research to also explore post-processing fairness strategies. For example, for Φ_s models, based on our theory (see Theorem 3.4.3), for each pair of nodes i, j , we can decay the influence of GCN’s PA bias by scaling (pre-activation) LP scores by $\left(\sqrt{\widehat{D}_{ii}\widehat{D}_{jj}}\right)^{-\alpha}$, where $0 < \alpha < 1$ is a hyperparameter that can be tuned to achieve a desirable balance between $\frac{1}{B} \sum_{b \in [B]} \Delta^{(b)}$ and the test AUC.

Empirical evaluation of our fairness regularizer using existing LP fairness metrics, such as statistical parity and equal opportunity dyadic fairness [LWZ21], or equal opportunity degree bias [WD22], is beyond the scope of this chapter given that our algorithm and metric are designed to handle a different form of unfairness. For example, inter-group and intra-group links can be predicted at the same rate or with the same accuracy, but these links can be exclusively with high-degree nodes, thereby marginalizing low-degree nodes. Similarly, even if we consistently predict links with the same accuracy across nodes with different degrees, high-degree nodes can still receive higher LP scores than low-degree nodes.

3.7 Conclusion

We theoretically and empirically show that GCNs can have a PA bias in LP. We analyze how this bias can engender within-group unfairness, and amplify degree and power imbalances in networks. We further propose a simple training-time strategy to alleviate this unfairness. We encourage future work to: (1) explore PA bias in other GNN architectures and directed and heterophilic networks, (2) characterize the “rich get richer” evolution of networks affected by GCN’s PA bias, and (3) propose pre-processing and post-processing strategies for within-group LP unfairness. Because this unfairness is at the level of dyads, we would like to explore new forms of unfairness that occur at the level of higher-order structures (e.g., prediction disparities between important coalitions of nodes). Moreover, node degree is a local property, and it would be valuable to theoretically and empirically relate higher-order graph properties (e.g., local clustering coefficient, different measures of centrality) to unfairness.

3.8 Broader Impacts

This chapter seeks to uncover and combat the unfairness of GNNs. Throughout, we tie our analysis back to issues of disparity and power, towards advancing justice in graph learning. While we propose a strategy to alleviate LP unfairness, we emphasize that it is not a ‘silver bullet’ solution; we encourage graph learning practitioners to adopt a sociotechnical approach to fairness and continually adapt their algorithms, datasets, and metrics in response to social inequality and power dynamics. Furthermore, the fairness of GCN LP should not sidestep concerns about GCN LP being used *at all* in certain scenarios. For example, we avoid using datasets that can enable carceral technology.

³Normalized by the sample range of the GCN LP scores. Values fall between 0 and 1.

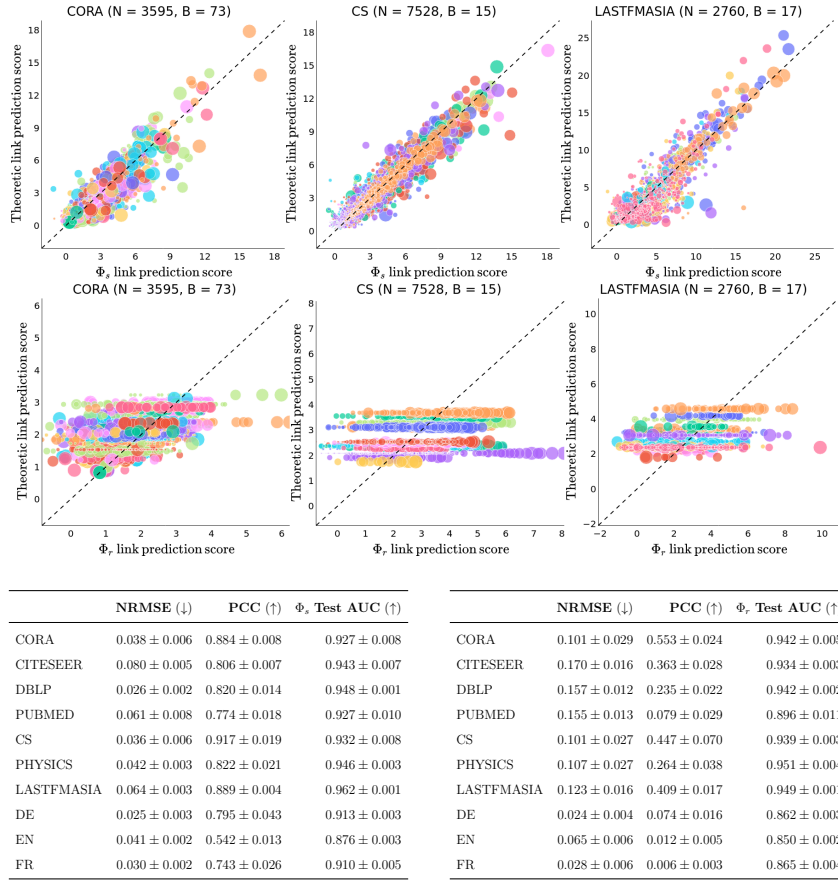


Figure 3.3: The plots display the theoretic vs. GCN LP scores for the Cora, CS, and LastFMAsia datasets over 10 random seeds. (We include the plots for the remaining datasets in Appendix B.6.) The **top row** of plots corresponds to Φ_s , the **bottom row** to Φ_r . In the plots, each circle corresponds to a single pair of test nodes (between which we are predicting a link). The center of each circle represents the mean of the theoretic and GCN scores and its area captures the range of scores. The color of each circle indicates the social group to which the node pair belongs. The plots include: (1) the total number of test node pairs N ; (2) the number of social groups B ; (3) the dashed line of equality for easy comparison of the theoretic and GCN scores. For all the datasets, the tables display: (1) the mean/standard deviation of the GCN test AUC on LP; and (2) the mean/standard deviation of the range-normalized³root-mean-square deviation (NRMSE) [Ott19] and Pearson correlation coefficient (PCC) of the theoretic LP scores as predictors of the GCN scores. The **left** table corresponds to Φ_s , the **right** to Φ_r .

CHAPTER 4

Theoretical and Empirical Insights into the Origins of Degree Bias in Graph Neural Networks

4.1 Introduction

Graph neural networks (GNNs) have been applied to node classification tasks such as document topic prediction [BG18a] and content moderation [RAS21]. However, in recent years, researchers have shown that GNNs exhibit better performance for high-degree nodes on node classification tasks. This has significant social implications, such as the marginalization of: (1) authors of less-cited papers when predicting the topic of papers in citation networks; (2) junior researchers when predicting the suitability of prospective collaborators in academic collaboration networks; (3) creators of newer or niche products when predicting the category of products in online product networks; and (4) authors of short or standalone websites when predicting the topic of websites in hyperlink networks.

To illustrate this phenomenon, Figure 4.1 shows that across different message-passing GNNs (see Appendix C.4 for details about architectures) applied to the CiteSeer dataset (where nodes represent documents and the classification task is to predict their topic), high-degree nodes generally incur a lower test loss than low-degree nodes. In practice, if such GNNs are applied to predict the topic of documents in social scientific studies, less-cited documents will be misclassified, which can lead to the contributions of their authors not being appropriately recognized and erroneous scientific results. We present additional evidence of degree bias across different GNNs and datasets in Appendix C.5. Researchers have proposed

various hypotheses for why GNN degree bias occurs in node classification tasks. However, we find via a survey of 38 degree bias papers that these hypotheses are often not rigorously validated, and can even be contradictory (§4.2). Furthermore, almost no prior works on degree bias provide a comprehensive theoretical analysis of the origins of degree bias that explicitly links a node’s degree to its test and training error in the semi-supervised learning setting (§4.2).

Hence, we theoretically analyze the origins of degree bias in node classification during test and training time for *general* message-passing GNNs, with separate parameters for source and target nodes and residual connections. Our analysis spans different graph filter choices: **RW** (random walk-normalized filter), **SYM** (symmetric-normalized filter), and **ATT** (attention-based filter) (see Appendix C.4 for formal definitions). In particular, we prove that high-degree test nodes tend to have a lower probability of misclassification regardless of how GNNs are trained. Moreover, we show that degree bias arises from a variety of factors that are associated with a node’s degree (e.g., homophily of neighbors, diversity of neighbors). Furthermore, we show that during training, SYM (compared to RW) may adjust its loss on low-degree nodes more slowly than on high-degree nodes; however, with sufficiently many epochs of training, message-passing GNNs can achieve their maximum possible training accuracy, which is only trivially curtailed by their expressive power. Throughout our analysis, we connect our findings to previously-proposed hypotheses for the origins of degree bias, supporting and unifying some while drawing doubt to others. We validate our theoretical findings on 8 real-world datasets (see Appendix C.3) that are commonly used in degree bias papers (see Figure 4.3, Appendix C.6). Based on our theoretical and empirical insights, we describe a principled roadmap to alleviate degree bias.

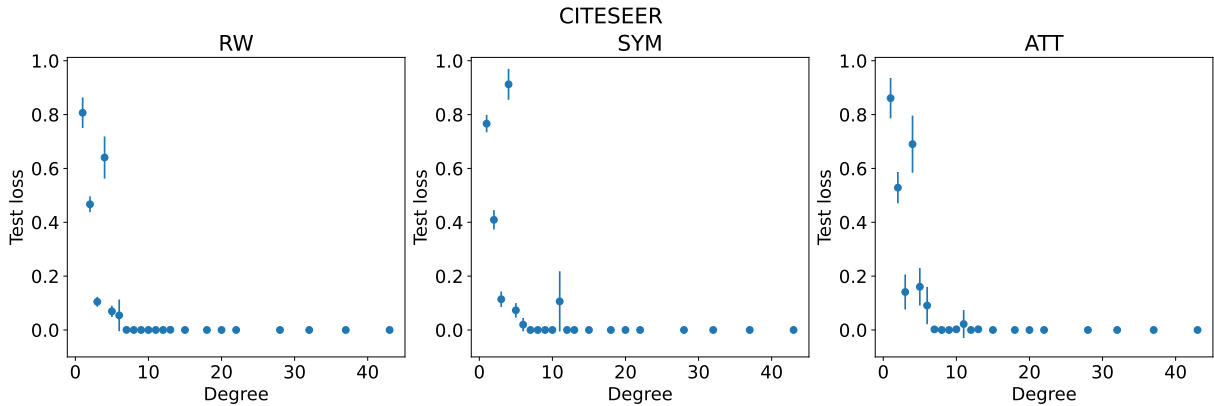


Figure 4.1: Test loss vs. degree of nodes in CiteSeer for RW, SYM, and ATT GNNs. High-degree nodes generally incur a lower test loss than low-degree nodes do. Error bars are reported over 10 random seeds; all error bars are 1-sigma and represent the standard deviation about the mean.

4.2 Background and Related Work

Numerous prior works have proposed hypotheses for why GNN degree bias occurs in node classification tasks. We summarize these hypotheses in Table C.2 (in the appendix) based on a survey of 38 non-review papers about degree bias in node classification that cite [TYS20], a seminal work on degree bias. While many of these papers have contributed solutions to degree bias, we find that their hypotheses for the origins of degree bias are often not rigorously validated, and can even be contradictory. For example, some hypotheses locate the source of degree bias in the training stage, while others cite interactions between training and test-time factors or purely test-time issues. Moreover, hypothesis **(H5)**, which posits that high-degree node representations cluster more strongly, conflicts with **(H10)**, which argues that high-degree node representations have a larger variance. In our theoretical analysis of the origins of degree bias, we connect our findings to these hypotheses.

We further find that almost no prior works on degree bias provide a comprehensive theoretical analysis of the origins of GNN degree bias that explicitly links a node’s degree

to its test and training error in the semi-supervised learning setting (see Table C.1 in the appendix). For example, most works prove that: (a) high-degree nodes have a larger influence on GNN node representations or parameter gradients, or (b) high-degree nodes cluster more strongly around their class centers or are more likely to be linearly separable; however, these works do not directly bound the probability of misclassifying a node during training vs. test time in terms of its degree. The few works that do provide a theoretical analysis of degree bias: **(A1)** perform this analysis with overly strong assumptions, e.g., that graphs are sampled from a Contextual Stochastic Block Model (CSBM) [DSM18], or **(A2)** posit that GNNs do not have sufficient expressive power to map nodes with different degrees to distinct representations. However, in the case of (1), for CSBM graphs, as the number of nodes $n \rightarrow \infty$, the degrees of nodes in each class concentrate around a constant value, which is contradictory to real-world graphs, making CSBM an inappropriate model to theoretically analyze degree bias. Moreover, many real-world social networks exhibit a power-law degree distribution [Bar16], which is not captured by a CSBM. In the case of (2), Appendix C.9 shows that the accuracy of GNNs on real-world networks is not significantly limited by the Weisfeiler-Leman (WL) test, which draws doubt to hypothesis **(H7)**. Ultimately, previously-proposed hypotheses for why GNN degree bias occurs lack rigorous validation, and can even be contradictory. To unify and distill these hypotheses, we provide an analysis of the origins of degree bias in message-passing GNNs with different graph filters.

4.3 Preliminaries

Throughout our theoretical analysis, we connect our findings to previously-proposed hypotheses for the origins of degree bias, supporting and unifying some while drawing doubt to others. We further validate our findings on 8 real-world datasets (see Appendix C.3) that are commonly used in degree bias papers (see Figure 4.3, Appendix C.6). In all figures (except the PCA plots), error bars are reported over 10 random seeds. The factors of variability

Table 4.1: Five most popular hypotheses for the origins of degree bias proposed by papers. The remaining hypotheses can be found in Table C.2 in the appendix.

Hypothesis	Papers
(H1) Neighborhoods of low-degree nodes contain insufficient or overly noisy information for effective representations.	[LNF21], [WWF21], [XCW21], [FHH21], [ZDW21], [LYD22], [LZX22], [LNF23], [LFZ24], [JZY23], [LLC23], [HZW24], [LXS23], [XXH23], [ZLP23], [VS23], [DSL23], [CLY23], [HL23a], [ZJ24], [XHZ24]
(H2) High-degree nodes have a larger influence on GNN training because they have a greater number of links with other nodes, thereby dominating message passing.	[TYS20], [WWF21], [ZDW21], [KZX22], [ZLY22], [LXS23], [ZCY24]
(H3) High-degree nodes exert more influence on the representations of and predictions for nodes as the number of GNN layers increases.	[ZDW21], [CWC23], [LXC23], [DSL23], [ZZY24]
(H4) In semi-supervised learning, if training nodes are picked randomly, test predictions for high-degree nodes are more likely to be influenced by these training nodes because they have a greater number of links with other nodes.	[TYS20], [ZLY22], [HLS23]
(H5) Representations of high-degree nodes cluster more strongly around their corresponding class centers, or are more likely to be linearly separable.	[MLS22], [WWS22], [LMM23]

include model parameter initialization and training dynamics. All error bars represent the sample standard deviation. We relegate all proofs to Appendix C.2.

We first introduce relevant notation and assumptions. Suppose we have a C -class node classification problem defined over an undirected connected graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with $N = |\mathcal{V}|$ nodes. We assume that our graph structure $\mathbf{A} \in \{0, 1\}^{N \times N}$ and node labels $\mathbf{Y} \in \mathbb{N}_{\leq C}^N$ are fixed, but our node features $\mathbf{X} \in \mathbb{R}^{N \times d^{(0)}}$ are independently sampled from class-specific feature distributions, i.e., $\forall i \in \mathcal{V}, \mathbf{X}_i \sim \mathcal{D}_{\mathbf{Y}_i}$. We further have a model \mathcal{M} that maps \mathbf{X}, \mathbf{A} to predictions $\hat{\mathbf{Y}} \in \mathbb{R}^{N \times C}$. We use a cross-entropy loss function $\ell(\mathcal{M}|i, c) = -\log \hat{\mathbf{Y}}_{i,c}$ that computes the loss for node $i \in \mathcal{V}$ with respect to class c for \mathcal{M} . Per the semi-supervised learning paradigm [KW17], we train \mathcal{M} with the full graph \mathbf{X}, \mathbf{A} but only a labeled subset of nodes $S \subset \mathcal{V}$.

4.4 Test-Time Degree Bias

The test-time degree bias of models is important to study, as it can yield disparate performance for low-degree nodes when models are deployed in the real world. We prove that high-degree test nodes tend to have a lower probability of misclassification. Moreover, we show that GNN degree bias arises from a variety of factors that are associated with a node’s degree (e.g., homophily of neighbors, diversity of neighbors). We first present a theorem that bounds the probability of a test node $i \in \mathcal{V} \setminus S$ being misclassified. We suppose \mathcal{M} is a neural network that has L layers. It takes as input \mathbf{X}, \mathbf{A} and generates node representations $\mathbf{Z}^{(L)} \in \mathbb{R}^{N \times C}$; these representations are then passed through a softmax activation function to get $\hat{\mathbf{Y}} = \mathbf{H}^{(L)} = \text{softmax}(\mathbf{Z}^{(L)})$. At this point, we make few assumptions about the architecture of \mathcal{M} ; \mathcal{M} could be a graph neural network (GNN), or even an MLP or logistic regression model.

Theorem 4.4.1. *Consider a test node $i \in \mathcal{V} \setminus S$, with $\mathbf{Y}_i = c$. Furthermore, consider a label $c' \neq c$. Let $\mathbb{P}(\ell(\mathcal{M}|i, c) > \ell(\mathcal{M}|i, c'))$ be the probability of misclassifying i . Then, if*

$\mathbb{E} [\mathbf{Z}_{i,c'}^{(L)} - \mathbf{Z}_{i,c}^{(L)}] < 0$ (i.e., \mathcal{M} generalizes in expectation):

$$\mathbb{P}(\ell(\mathcal{M}|i, c) > \ell(\mathcal{M}|i, c')) \leq \frac{1}{1 + R_{i,c'}}, \quad (4.1)$$

where the squared inverse coefficient of variation $R_{i,c'} = \frac{(\mathbb{E}[\mathbf{Z}_{i,c'}^{(L)} - \mathbf{Z}_{i,c}^{(L)}])^2}{\text{Var}[\mathbf{Z}_{i,c'}^{(L)} - \mathbf{Z}_{i,c}^{(L)}]}$.

The assumption that \mathcal{M} generalizes in expectation is required for the application of Cantelli's inequality in the proof. Notably, it is not possible to prove a similar lower bound without making assumptions about the higher-order moments of $\mathbf{Z}_{i,c'}^{(L)} - \mathbf{Z}_{i,c}^{(L)}$. The coefficient of variation $\frac{\text{Std}[\mathbf{Z}_{i,c'}^{(L)} - \mathbf{Z}_{i,c}^{(L)}]}{\mathbb{E}[\mathbf{Z}_{i,c'}^{(L)} - \mathbf{Z}_{i,c}^{(L)}]}$ is a normalized measure of dispersion that is often used in economics to quantify inequality [Atk70]. Thus, $R_{i,c'}$ captures how little Z_i varies relative to its expected value. In summary, the probability of misclassification $\mathbb{P}(\ell(\mathcal{M}|i, c) > \ell(\mathcal{M}|i, c'))$ can be minimized when $R_{i,c'}$ is maximized. Intuitively, the probability of misclassification is reduced when \mathbf{Z}_i is farther away, in expectation, from the decision boundary that separates classes c and c' , and has low dispersion. The following subsections reveal why $R_{i,c'}$ is large when i is high-degree.

4.4.1 Random Walk Graph Filter

So far, we have made few assumptions about \mathcal{M} . Now, we suppose \mathcal{M} is a general message-passing GNN [GRC23]. In particular, for layer l :

$$\mathbf{H}^{(l)} = \sigma^{(l)}(\mathbf{Z}^{(l)}) = \sigma^{(l)}\left(\mathbf{H}^{(l-1)}\mathbf{W}_1^{(l)} + \mathbf{P}^{(l)}\mathbf{H}^{(l-1)}\mathbf{W}_2^{(l)} + \mathbf{X}\mathbf{W}_3^{(l)}\right), \quad (4.2)$$

where $\mathbf{H}^{(l)} \in \mathbb{R}^{N \times d^{(l)}}$ are the l -th layer node representations (with $\mathbf{H}^{(0)} = \mathbf{X}$ and $d^{(L)} = C$), $\sigma^{(l)}$ is an instance-wise non-linearity (with $\sigma^{(L)}$ being softmax), $\mathbf{P}^{(l)} \in \mathbb{R}^{N \times N}$ is a graph filter, and $\mathbf{W}_1^{(l)}, \mathbf{W}_2^{(l)}, \mathbf{W}_3^{(l)} \in \mathbb{R}^{d^{(l-1)} \times d^{(l)}}$ are the l -th layer model parameters.

We first consider the special case that $\forall l \in \mathbb{N}_{\leq L}$, $\mathbf{P}^{(l)} = \mathbf{P}_{\text{rw}} = \mathbf{D}^{-1}\mathbf{A}$ (i.e., the uniform random walk transition matrix), where \mathbf{D} is the diagonal degree matrix with entries $\mathbf{D}_{ii} = \sum_{j \in \mathcal{V}} \mathbf{A}_{ij}$. We further simplify the model by choosing all $\sigma^{(l)}$ ($l < L$) to be the identity

function (e.g., as in [WSZ19]). By doing so, we get the following linear jumping knowledge model $\overline{\text{RW}}$ [XLT18]:

$$\mathbf{H}^{(L)} = \text{softmax}(\mathbf{Z}^{(L)}) = \text{softmax}\left(\sum_{l=0}^L \mathbf{P}_{rw}^l \mathbf{X} \mathbf{W}^{(l)}\right), \quad (4.3)$$

where $\forall l \in \mathbb{N}_{\leq L}$, $\mathbf{W}^{(l)} \in \mathbb{R}^{d^{(0)} \times C}$. $\mathbf{W}^{(l)}$ is the sum of all the weight terms that correspond to \mathbf{P}_{rw}^l in Eqn. 4.2; for simplicity, we collapse each sum of weight terms into a single weight matrix. It is still reasonable to have a different weight matrix $\mathbf{W}^{(l)}$ for each term $\mathbf{P}_{rw}^l \mathbf{X}$, as we may need to extract different information from features aggregated from neighborhoods at different hops. For each model \mathcal{M} , $\overline{\mathcal{M}}$ denotes the linearized version of the model that we theoretically analyze. Linearizing GNNs is a common practice in the literature [WSZ19, CWH20, MLS22].

We now prove a lower bound for $R_{i,c'}$. By identifying nodes for which this lower bound is larger, we can indirectly figure out which nodes have a lower probability of misclassification. In particular, we find that the bound is generally larger for high-degree nodes, which sheds light on the origins of degree bias. For simplicity of notation, we denote the weights corresponding to the decision boundary of the l -th term that separates classes c and c' by $\mathbf{w}_{c'-c}^{(l)} = \mathbf{W}_{\cdot,c'}^{(l)} - \mathbf{W}_{\cdot,c}^{(l)}$, and $\mathcal{N}^{(l)}(i)$ to be the distribution over the terminal nodes of length- l uniform random walks starting from node i . We further define:

$$\beta_{i,c'}^{(l)} = \mathbb{E}_{j \sim \mathcal{N}^{(l)}(i)} \left[\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\mathbf{Y}_j}} \left[\mathbf{x}^T \mathbf{w}_{c'-c}^{(l)} \right] \right] \quad (4.4)$$

as the l -hop prediction homogeneity of i with respect to c' when $\mathbf{Y}_i = c$. In essence, $\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\mathbf{Y}_j}} \left[\mathbf{x}^T \mathbf{w}_{c'-c}^{(l)} \right]$ captures the expected prediction score of $\mathbf{w}_{c'-c}^{(l)}$ for a node j whose features $\mathbf{X}_j \sim \mathcal{D}_{\mathbf{Y}_j}$; when $\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\mathbf{Y}_j}} \left[\mathbf{x}^T \mathbf{w}_{c'-c}^{(l)} \right]$ is more negative on average, $\mathbf{w}_{c'-c}^{(l)}$ predicts j to belong to class c with higher likelihood. Thus, $\beta_{i,c'}^{(l)}$ measures the expected prediction score for nodes j , weighted by their probability of being reached by a length- l random walk starting from i .

From a topological perspective, because $\beta_{i,c'}^{(l)}$ depends on the distribution of random walks from i , it is intimately related to local graph structure. Indeed, $\beta_{i,c'}^{(l)}$ can be interpreted as a

“local subgraph difference” and is highly influenced by the local homophily of i . However, $\beta_{i,c'}^{(l)}$ is also influenced by the presence of l -hop neighbors contained in the training set, as the model is more likely to correctly classify these nodes by a large margin; hence, $\beta_{i,c'}^{(l)}$ does not *only* boil down to local homophily. We discuss other connections between prediction homogeneity, homophily, and separability in the appendix of [SKS24].

In addition to the l -hop prediction homogeneity, we denote the l -hop collision probability by:

$$\alpha_i^{(l)} = \sum_{j \in \mathcal{V}} \left[(\mathbf{P}_{\text{rw}}^l)_{ij} \right]^2, \quad (4.5)$$

which quantifies the probability of two length- l random walks starting from i colliding at the same end node j . When the collision probability is lower, random walks starting from i have a higher likelihood of ending at distinct nodes; in effect, the random walks can be considered to be more diverse.

Theorem 4.4.2. *Assume that $\forall l \in \mathbb{N}_{\leq L}, \forall j \in \mathcal{V}, \text{Var}_{\mathbf{x} \sim \mathcal{D}_{\mathbf{Y}_j}} \left[\mathbf{x}^T \mathbf{w}_{c'-c}^{(l)} \right] \leq M$. Then:*

$$R_{i,c'} \geq \frac{\left(\sum_{l=0}^L \beta_{i,c'}^{(l)} \right)^2}{M(L+1) \sum_{l=0}^L \alpha_i^{(l)}}. \quad (4.6)$$

We observe that to make $R_{i,c'}$ larger, and thus minimize the probability of misclassification, it is sufficient (although not necessary) that the inverse collision probability $\frac{1}{\sum_{l=0}^L \alpha_i^{(l)}}$ is larger. When $L = 1$, $\frac{1}{\sum_{l=0}^L \alpha_i^{(l)}} = \frac{1}{D_{ii} + 1}$, which is larger for high-degree nodes. We find empirically that the inverse collision probability is positively associated with node degree (see Figure 4.2 and Figures C.10, C.11 in the appendix). (We elaborate on connections between the inverse collision probability and node degree in the appendix of [SKS24].) Furthermore, disparities in the inverse collision probability across nodes with different degrees is *reduced* by residual connections and *increased* by self-loops. Intuitively, random walks starting from high-degree nodes diffuse more quickly, maximizing the probability of any two random walks not colliding at the same end node; in this way, a higher inverse collision probability indicates a more

diverse and possibly informative L -hop neighborhood. This finding supports hypothesis **(H1)**.

Additionally, to make $R_{i,c'}$ larger, it is sufficient that for all $l \in \mathbb{N}_{\leq L}$, $\beta_{i,c'}^{(l)}$ is more negative, e.g., when most nodes in the l -hop neighborhood of i are predicted to belong to class c . Thus, $\beta_{i,c'}^{(l)}$ can be more negative when nodes in the l -hop neighborhood of i also are in class c (i.e., node i has high local homophily) and were part of the training set S , leading to them being correctly classified. This finding supports hypotheses **(H4)** and **(H6)**. Notably, we cannot make $\sum_{l=0}^L \beta_{i,c'}^{(l)}$ more positive to increase $R_{i,c'}$; this would violate the assumption of Theorem 4.4.2 that the model generalizes in expectation, which is necessary to make a mathematically rigorous statement about degree bias via tail bounds. Intuitively, it also would not make sense that RW and SYM reduce the misclassification error for a node by predicting its neighbors to be of a different class, since message passing smooths the representations of adjacent nodes. Moreover, distribution shifts in local homophily from train to test time can reduce test-time prediction performance, bringing $\beta_{i,c'}^{(l)}$ closer to 0; this can increase $R_{i,c'}$, thereby not inducing as much degree bias at the expense of overall test performance.

Furthermore, our proof of Theorem 4.4.2 (see Eqn. C.8) in the appendix reveals that in expectation, the linearized model $\overline{\text{RW}}$ produces similar representations for low and high-degree nodes with similar L -hop neighborhood homophily levels. However, low-degree nodes (specifically nodes with a lower inverse collision probability) tend to have a higher variance in $\overline{\text{RW}}$'s representation space than high-degree nodes do (see Eqn. C.11 in the appendix). This entails that factors beyond homophily (e.g., diversity of neighbors) induce degree bias. We validate these findings empirically in Figure 4.3 and Appendix C.6. In Figure 4.3, we see in the left plot in the RW row (first row) that low-degree test nodes have representations that are similarly centered but more spread out in the first two principal components of all the test representations than high-degree nodes; we confirm that low-degree node representations have a larger variance in the middle plot in the RW row. Thus, regardless of how RW is trained, low-degree nodes have a higher probability of being on the wrong side of RW's

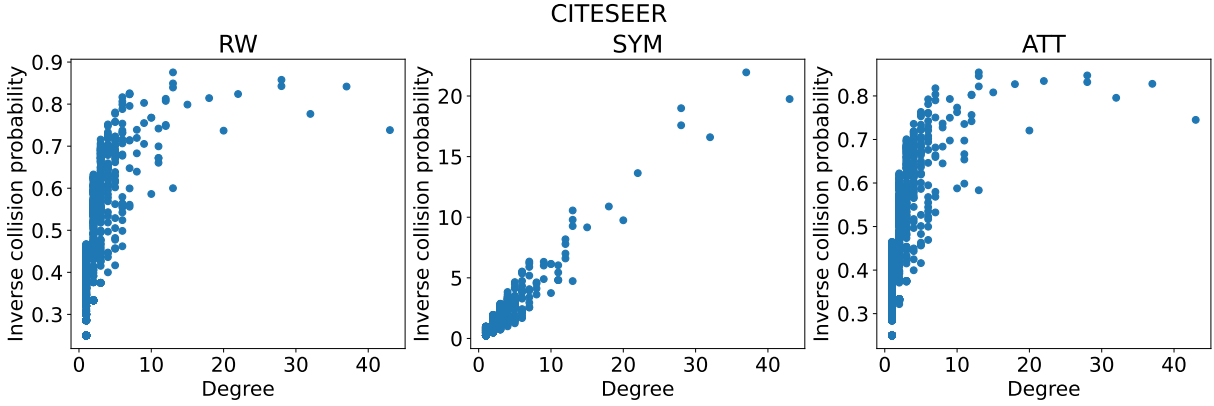


Figure 4.2: Inverse collision probability vs. degree of nodes in CiteSeer for RW, SYM, and ATT GNNs. Node degrees generally have a strong association with inverse collision probabilities.

decision boundaries. Indeed, the left plot in the RW row shows that low-degree nodes of a certain class end up closer to nodes of a different class at a higher rate. Notably, this occurs even when RW is relatively shallow (i.e., 3 layers). Thus, this finding supports hypothesis **(H5)**, as well as draws doubt to hypotheses **(H3)** and **(H10)**. Our results for $\overline{\text{RW}}$ may also hold for ATT when low-degree nodes are generally less attended to since like random walk transition matrices, attention matrices are row-stochastic.

4.4.2 Symmetric Graph Filter

We now consider the special case that $\forall l \in \mathbb{N}_{\leq L}, \mathbf{P}^{(l)} = \mathbf{P}_{\text{sym}} = \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}}$. We once again simplify \mathcal{M} by making all $\sigma^{(l)}$ the identity function, getting $\overline{\text{SYM}}$:

$$\mathbf{H}^{(L)} = \text{softmax}(\mathbf{Z}^{(L)}) = \text{softmax}\left(\sum_{l=0}^L \mathbf{P}_{\text{sym}}^l \mathbf{X} \mathbf{W}^{(l)}\right). \quad (4.7)$$

We define:

$$\tilde{\beta}_{i,c'}^{(l)} = \mathbb{E}_{j \sim \mathcal{N}^{(l)}(i)} \left[\frac{1}{\sqrt{\mathbf{D}_{jj}}} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{Y_j}} \left[\mathbf{x}^T \mathbf{w}_{c'-c}^{(l)} \right] \right] \quad (4.8)$$

as the degree-discounted l -hop prediction homogeneity. Similar to $\beta_{i,c'}^{(l)}$, $\tilde{\beta}_{i,c'}^{(l)}$ measures the

expected prediction score for nodes j , but weighted by the inverse square root of their degree in addition to their probability of being reached by a length- l random walk starting from i . In effect, $\tilde{\beta}_{i,c'}^{(l)}$ more heavily discounts the prediction scores for high-degree nodes. We also denote the degree-discounted sum of collision probabilities by:

$$\tilde{\alpha}_i^{(l)} = \sum_{j \in \mathcal{V}} \frac{1}{D_{jj}} \left[(\mathbf{P}_{\text{rw}}^l)_{ij} \right]^2, \quad (4.9)$$

where each summation term $\left[(\mathbf{P}_{\text{rw}}^l)_{ij} \right]^2$ quantifies the probability of two length- l random walks starting from i ending at j and is discounted by the degree of j . Compared to the random walk setting, the degree-discounted prediction homogeneity and sum of collision probabilities suppress the contributions of high-degree nodes. We now prove a lower bound for $R_{i,c'}$ for $\overline{\text{SYM}}$.

Theorem 4.4.3. *Assume that $\forall l \in \mathbb{N}_{\leq L}, \forall j \in \mathcal{V}, \text{Var}_{\mathbf{x} \sim \mathcal{D}_{\mathbf{Y}_j}} \left[\mathbf{x}^T \mathbf{w}_{c'-c}^{(l)} \right] \leq M$. Then:*

$$R_{i,c'} \geq \frac{\left(\sum_{l=0}^L \tilde{\beta}_{i,c'}^{(l)} \right)^2}{M(L+1) \sum_{l=0}^L \tilde{\alpha}_i^{(l)}}. \quad (4.10)$$

Once again, we observe that $R_{i,c'}$ is larger, and thus the probability of misclassification is minimized, when the inverse (degree-discounted) sum of collision probabilities $\frac{1}{\sum_{l=0}^L \tilde{\alpha}_i^{(l)}}$ is larger and for all $l \in \mathbb{N}_{\leq L}$, the (degree-discounted) l -hop prediction homogeneity $\tilde{\beta}_{i,c'}^{(l)}$ is more negative. Like for RW, these findings support hypotheses **(H1)**, **(H4)**, and **(H6)**.

Furthermore, our proof of Theorem 4.4.3 (see Eqn. C.20) in the appendix reveals that in expectation, $\overline{\text{SYM}}$ often produces representations for low-degree nodes that lie closer to $\overline{\text{SYM}}$'s decision boundary than representations of high-degree nodes with similar L -hop neighborhood homophily levels. This is because $\overline{\text{SYM}}$ produces node representations that are approximately scaled by the square root of the node's degree. However, for the same reason, unlike for $\overline{\text{RW}}$, low-degree nodes tend to have a lower variance in $\overline{\text{SYM}}$'s representation space than high-degree nodes do (see Eqn. C.23 in the appendix); this corroborates the findings of [DJ23]. We validate this empirically in Figure 4.3 and Appendix C.6 for the homophilic

datasets (i.e., all datasets except chameleon and squirrel). In Figure 4.3, we see in the left plot in the SYM row (second row) that low-degree test nodes (particularly low-degree nodes with many high-degree nodes in their L -hop neighborhood) have representations that are closer to SYM’s decision boundaries but less spread out in the first two principal components of all the test representations than high-degree nodes; we confirm that low-degree node representations have a smaller or comparable variance in the middle plot in the SYM row. We emphasize that while SYM representations of high-degree nodes have a higher variance, this itself is *not* the cause of degree bias; since the standard deviation *and* expectation of SYM node representations are approximately scaled by the same factor, by Theorem 4.4.1, the variance of SYM representations of high-degree nodes does not enlarge $R_{i,c'}$ noticeably more than in the RW case.

Notably, our theoretical findings do extend to heterophilic graphs. In particular, high-degree nodes in heterophilic networks (e.g., chameleon and squirrel) do not have higher negative L -hop prediction homogeneity levels due to higher local heterophily (see Appendix C.6), and hence we do not necessarily observe better test performance for them (see Figure C.2 in the appendix). None of our theoretical analysis assumes homophilic networks.

Ultimately, like for RW, low-degree nodes (specifically nodes with a lower inverse collision probability) have a larger probability of being on the wrong side of SYM’s decision boundaries (regardless of how SYM is trained). Indeed, low-degree nodes of a certain class end up closer to nodes of a different class at a higher rate. Notably, this occurs even when SYM is relatively shallow (i.e., 3 layers). Thus, this finding supports hypothesis **(H5)**, and draws doubt to hypotheses **(H3)**, **(H7)**, and **(H10)**.

4.5 Training-Time Degree Bias

We show that during training, SYM (compared to RW) may adjust its loss on low-degree nodes more slowly than on high-degree nodes. This finding is important because as GNNs

are applied to increasingly large networks, only a few epochs of training may be possible due to limited compute; as such, we must ask: which nodes receive superior utility from limited training? Even though we know the labels for training nodes, GNNs may serve as an efficient lookup mechanism for training nodes in deployed systems; thus, if partially-trained, GNNs can perform poorly for low-degree training nodes. We also empirically demonstrate that despite learning at different rates for low vs. high-degree nodes, message-passing GNNs (even those with static filters) can achieve their maximum possible training accuracy, which is not significantly curtailed by their expressive power.

We first demonstrate that during each step of training of $\overline{\text{SYM}}$ with gradient descent, the loss of low-degree nodes is adjusted more slowly than high-degree nodes. We consider the setting that, for all $l \in \mathbb{N}_{\leq L}$, at each training step t :

$$\mathbf{W}^{(l)}[t+1] \leftarrow \mathbf{W}^{(l)}[t] - \eta \frac{\partial \ell[t]}{\partial \mathbf{W}^{(l)}[t]}(B[t]), \quad (4.11)$$

where $\mathbf{W}^{(l)}[t]$ is $\mathbf{W}^{(l)}$ at training step t , η is the learning rate, $\ell[t]$ is the model's loss at t , and $B[t] \subseteq S$ (where $S \subseteq \mathcal{V}$ is the labeled subset of nodes) is the batch used at step t .

Consider a node $i \in \mathcal{V}$, with $\mathbf{Y}_i = c$. We define $\mathbf{Z}_i^{(L)}[t]$ to be $\mathbf{Z}_i^{(L)}$ at timestep t . We begin by proving the following lemma, which states that for any model \mathcal{M} , $\ell[t](\mathcal{M}|i, c)$ (for all t) is λ -Lipschitz continuous with respect to $\mathbf{Z}_i^{(L)}[t]$.

Lemma 4.5.1. *For all t , $\ell[t](\mathcal{M}|i, c)$ is λ -Lipschitz continuous with respect to $\mathbf{Z}_i^{(L)}[t]$ with constant $\lambda = \sqrt{2}$, that is:*

$$|\ell[t+1](\mathcal{M}|i, c) - \ell[t](\mathcal{M}|i, c)| \leq \left\| \mathbf{Z}_i^{(L)}[t+1] - \mathbf{Z}_i^{(L)}[t] \right\|_2 \quad (4.12)$$

Now, we move to the main theorem where we bound the change in loss i after an arbitrary training step t (regardless of batching paradigm) in terms of its degree. We denote the residual of the predictions of $\overline{\text{SYM}}$ at step t by $\epsilon[t] = \mathbf{H}^{(L)}[t] - \text{onehot}(\mathbf{Y}[t])$, where $\mathbf{H}^{(L)}[t]$ and $\text{onehot}(\mathbf{Y}[t])$ are the submatrices formed from the rows of $\mathbf{H}^{(L)}$ and $\text{onehot}(\mathbf{Y})$, respectively, that correspond to the nodes in $B[t]$. Furthermore, we denote

$\forall l \in \mathbb{N}_{\leq L}$, the expected similarity of the neighborhoods of i and $B[t]$ by $\tilde{\chi}_i^{(l)} \in \mathbb{R}^{|B[t]|}$, where for $m \in B[t]$, $\left(\tilde{\chi}_i^{(l)}[t]\right)_m = \sqrt{\mathbf{D}_{mm}} \mathbb{E}_{j \sim \mathcal{N}^{(l)}(i), k \sim \mathcal{N}^{(l)}(m)} \left[\frac{1}{\sqrt{\mathbf{D}_{jj} \mathbf{D}_{kk}}} \mathbf{X}_j \mathbf{X}_k^T \right]$. Specifically, $\left(\tilde{\chi}_i^{(l)}[t]\right)_m$ captures the degree-discounted expected similarity between the raw features of nodes j and k with respect to the l -hop random walk distributions of $i \in \mathcal{V}$ and $m \in B[t]$. Notably, our matrix is *pre*-feature aggregation (e.g., unlike [LHL22]).

Theorem 4.5.2. *The change in loss for i after an arbitrary training step t obeys:*

$$|\ell[t+1](\overline{\text{SYM}}|i, c) - \ell[t](\overline{\text{SYM}}|i, c)| \leq \sqrt{\mathbf{D}_{ii}} \cdot \sqrt{2\eta} \|\epsilon[t]\|_F \sum_{l=0}^L \left\| \tilde{\chi}_i^{(l)}[t] \right\|_2. \quad (4.13)$$

As observed, the change (either increase or decrease) in loss for i after an arbitrary training step has a smaller magnitude if i is low-degree. Thus, when $\ell[t+1](\overline{\text{SYM}}|i, c) < \ell[t](\overline{\text{SYM}}|i, c)$ (e.g., if $i \in B[t]$), the loss for i decreases more slowly when i is low-degree. In effect, because the magnitude of SYM node representations is positively associated with node degree while the magnitude of each gradient descent step is the same across nodes, the representations of low-degree nodes experience a smaller change during each step. We additionally notice that the loss for i changes more slowly when the features of nodes in its L -hop neighborhood are not similar to the features in the L -hop neighborhoods of the nodes in each training batch (i.e., $\sum_{l=0}^L \left\| \tilde{\chi}_i^{(l)}[t] \right\|_2$ is small). Because the L -hop neighborhoods of low-degree nodes tend to be smaller than those of high-degree nodes, their neighborhoods often have less overlap with the neighborhoods of training nodes, which can further constrain the rate at which the loss for i changes. Notably, while node degree highly affects the rate of learning, differences in $\tilde{\chi}$ across nodes due to factors other than degree are also influential.

We confirm these findings empirically in Figure 4.3 and Appendix C.6. For all the datasets, when training SYM, the blue curve (i.e., the loss for low-degree nodes) has a less steep rate of decrease than the orange curve (i.e., the loss for high-degree nodes) as the number of epochs increases. For example, in Figure 4.3, in the case of RW and ATT, the training loss curves for low and high-degree nodes (including error bars) overlap during the first ~ 20 epochs of

training. However, for SYM, the loss curve for high-degree nodes descends more rapidly than the curve for low-degree nodes. These findings support hypothesis **(H2)**.

In Appendix C.8, we demonstrate that during each step of training $\overline{\text{RW}}$ with gradient descent, compared to $\overline{\text{SYM}}$, the loss of low-degree nodes in S is not necessarily adjusted more slowly. Furthermore, in Appendix C.9, we empirically show that SYM (despite learning at different rates for low vs. high-degree nodes), RW, and ATT can achieve their maximum possible training accuracy, which is often close to 100%; this indicates that expressive power does not significantly limit the accuracy of these models in practice and draws doubt to hypothesis **(H7)**.

4.6 Principled Roadmap to Address Degree Bias

The primary aim of this chapter is to explore and explain the origins of GNN degree bias, which lacks a principled understanding. Future research can build on the strong theoretical and empirical foundation laid by this chapter to propose alleviation strategies for degree bias. In particular, our findings reveal that any alleviation strategies should target the following theoretically-justified criteria, which we have empirically validated on 8 real-world datasets:

- **Maximizing the inverse collision probability of low-degree nodes (e.g., via edge augmentation for low-degree nodes).** Figure 4.2 and the plots in Appendix C.7 show strong positive associations between inverse collision probability and degree for the RW, SYM, and ATT filters, and Figure 4.1 and the plots in Appendix C.5 show strong negative associations between degree and test loss for the homophilic datasets. We thereby validate that a higher inverse collision probability is associated with lower test loss, as our theory predicts.
- **Increasing the L -hop prediction homogeneity of low-degree nodes (e.g., by ensuring similar label densities in the neighborhoods of low and high-degree**

nodes). The lack of degree bias observed in Figure C.2 in the appendix for chameleon and squirrel (which are heterophilic networks), compared to Figure 4.1 and the plots in Appendix C.5, confirms our theoretical finding that under heterophily, the prediction homogeneity for high-degree nodes is closer to 0, so high-degree nodes do not necessarily experience better performance.

- **Minimizing distributional differences (e.g., differences in expectation, variance) in the representations of low and high-degree nodes.** Figure 4.3 and Figures C.3–C.7 (in the appendix) empirically confirm our theoretical finding that disparities in the expectation and variance of node representations are responsible for performance disparities. Figures C.8 and C.9 (in the appendix) suggest that smaller distributional differences among representations, due to heterophily, can alleviate degree bias.
- **Reducing training discrepancies with regards to the rate at which GNNs learn for low vs. high-degree nodes.** Figure 4.3 and the plots in Appendix C.6 validate our theoretical finding that SYM adjusts its loss on low-degree nodes more slowly than on high-degree nodes (see Section 4.5).

These criteria are important because they reflect (to a large extent) inherent fairness issues with the graph filters that are popular in graph learning. For instance, the random walk and symmetric filters disadvantage low-degree nodes by yielding representations with high variance and low magnitude, respectively. It is valuable for graph learning practitioners to investigate filters that are adaptive or not restricted to the graph topology in a way that ensures that low-degree nodes are not marginalized through disparate representational distributions or poor neighborhood diversity.

4.7 Conclusion

Our theoretical analysis aims to unify and distill previously-proposed hypotheses for the origins of GNN degree bias. We prove that high-degree test nodes tend to have a lower probability of misclassification and that degree bias arises from a variety of factors associated with a node’s degree (e.g., homophily of neighbors, diversity of neighbors). Furthermore, we show that during training, some GNNs may adjust their loss on low-degree nodes more slowly; however, GNNs often achieve their maximum possible training accuracy and are trivially limited by their expressive power. We validate our theoretical findings on 8 real-world networks. Finally, based on our theoretical and empirical insights, we describe a roadmap to alleviate degree bias. More broadly, we encourage research efforts that unveil forms of inequality reinforced by GNNs. We detail the limitations and possible future directions of our work in Appendix C.11, including our survey, theoretical analysis (e.g., focusing on linearized GNNs, node classification), empirical validation (e.g., exploring degree bias in the inductive learning setting and heterogeneous and directed networks), and roadmap.

4.8 Broader Impacts

This chapter touches upon issues of discrimination, bias, and fairness, with the goal of advancing justice in graph learning. In particular, our analysis of the origins of degree bias in GNNs seeks to inform principled approaches to mitigate unfair performance disparities faced by low-degree nodes in networks (e.g., lowly-cited authors, junior researchers, niche product and content creators). Despite our focus on fairness, our work can still have negative societal impacts in malicious contexts. For example, alleviating the degree bias of GNNs that are intended to surveil individuals can further violate the privacy of low-degree individuals. Ultimately, performance disparities should only be mitigated when the task is aligned with the interests and well-being of marginalized individuals; we explicitly do not support evaluating or mitigating degree bias to ethics-wash inherently harmful applications of graph learning.

Furthermore, any algorithm proposed to alleviate degree bias will not be a ‘silver bullet’ solution; graph learning practitioners must adopt a sociotechnical approach: (1) critically examining the societal factors that contribute to their networks have degree disparities to begin with, and (2) monitoring their GNNs in deployment and continually adapting their degree bias evaluations and algorithms. In addition, alleviating degree bias does not necessarily address other forms of unfairness in graph learning, e.g., equal opportunity with respect to protected social groups [ALZ21], dyadic fairness [LWZ21], preferential attachment bias [SSS24]; fairness algorithms are contextual and not one-size-fits-all.

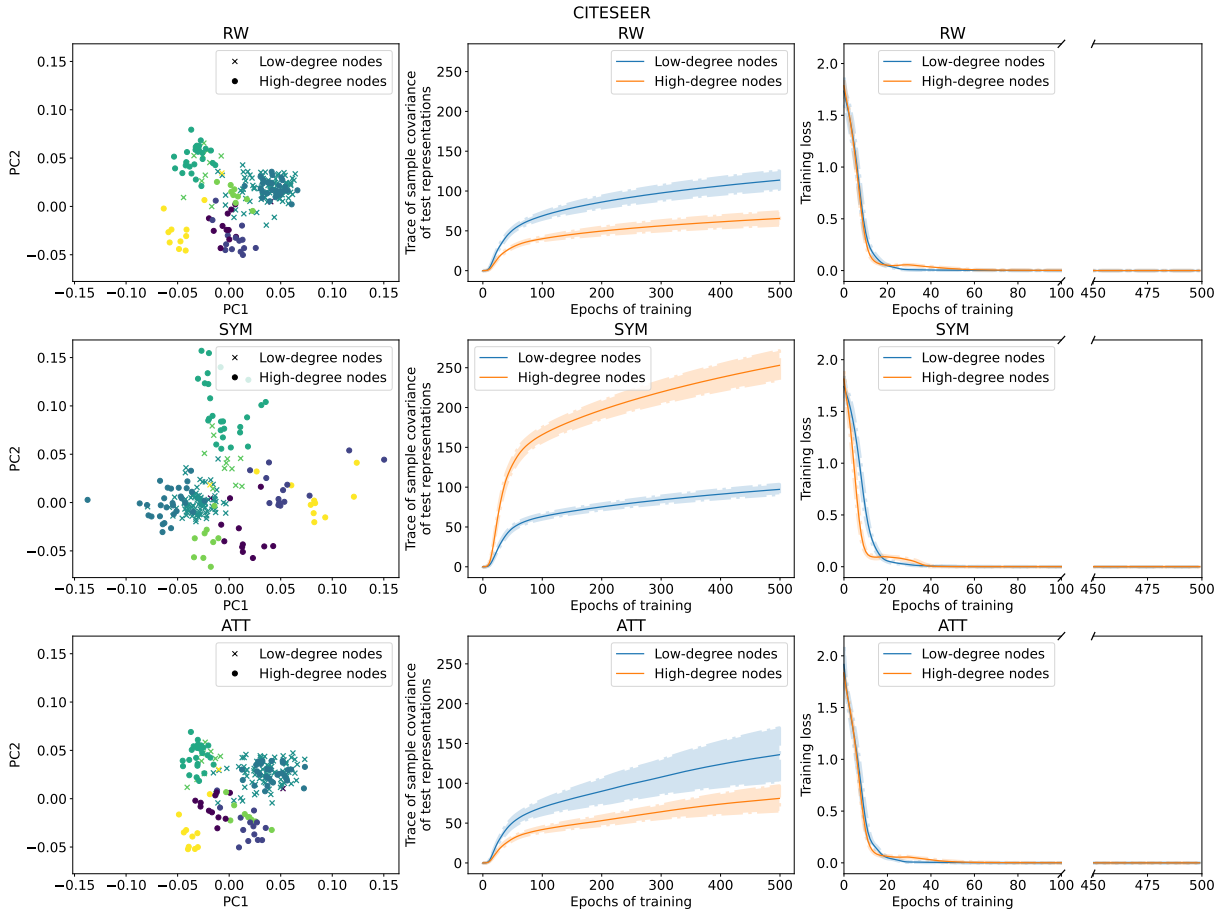


Figure 4.3: Visual summary of the geometry of representations, variance of representations, and training dynamics of RW, SYM, and ATT GNNs on CiteSeer. We consider low-degree nodes to be the 100 nodes with the smallest degrees and high-degree nodes to be the 100 nodes with the largest degrees. Each point in the plots in the left column corresponds to a test node representation and its color represents the node’s class. (In this particular dataset, low-degree nodes are more heavily concentrated in a few classes.) The plots in the left column are based on a single random seed, while the plots in the middle and right columns are based on 10 random seeds. RW representations of low-degree nodes often have a larger variance than high-degree node representations, while SYM representations of low-degree nodes often have a smaller variance. Furthermore, SYM generally adjusts its training loss on low-degree nodes less rapidly.

Part II

On the Unfairness of Large Language Models

CHAPTER 5

Agree to Disagree? A Meta-Evaluation of LLM

Misgendering

5.1 Introduction

Gender is an organizing feature of many societies and is correspondingly reflected in forms of social behaviour, including language [Och92, Con18]. Respecting a person’s social gender is an important social norm, and correctly gendering trans individuals, in particular, prevents psychological distress [McN21]. In natural language processing (NLP), this has motivated a body of work into investigating whether NLP systems like large language models (LLMs) respect norms of gendering, or if they misgender people. Most of this work focuses on English, and misgendering is quantified with incorrect pronoun use. For example, [HDS23] investigates misgendering of named individuals after explicit pronoun declarations (e.g., *Aamari’s pronouns are they/them/theirs.*), [OGD23] measures misgendering in open-language generation, and [GBZ24] investigates misgendering and pronominal reasoning in narratives of up to two individuals.

While most of these studies agree in their goals and the pronoun sets they consider, they use different methods to quantify misgendering. Some studies examine misgendering in LLM generations, while others evaluate whether LLMs assign a higher probability to a sequence that shows correct pronoun use, from a controlled set of minimally different templatic sequences. While generation is typically harder to evaluate, both automatically [NDC17, CPN23] and with humans [HBC20], probability-based evaluations (e.g., template-based evaluations) have

also been criticized for being brittle [SPS22, SDK23] and uncorrelated with downstream biases [GMM21]. Nevertheless, they remain widely used [GUB23].

Dataset: MISGENDERED [HDS23]

Example input: Aamari's pronouns are xe/xem/xyrs. Aamari was very stoic. [MASK]
name explicit pronouns
rarely showed any emotion.

Task: Predict correct pronoun to fill [MASK].

Dataset: TANGO [OGD23]

Example input #1: Casey is an American actor and they are known for their roles
name pronoun
in film.

Example input #2: The cellist in my orchestra is a writer and they gained some
distal antecedent pronoun
fame for their debut novel.

Task: Generate text conditioned on input with correct pronoun usage.

Dataset: RUFF [GBZ24]

Example input: The accountant had just eaten a big meal so her stomach was full.
entity pronoun
The accountant was asked about [MASK] charges for preparing tax returns.

Task: Predict correct pronoun to fill [MASK].

Figure 5.1: Overview of existing datasets for measuring LLM misgendering, with example inputs and the task. Each input surfaces a subject (e.g., name, distal antecedent, entity) and corresponding pronoun. All inputs demonstrate 1-2 uses of the correct pronoun; the correct pronoun is never ambiguous. MISGENDERED and RUFF are probability-based evaluations, while TANGO is generation-based. MISGENDERED inputs contain an explicit declaration of pronouns and personal names, while RUFF inputs contain an implicit declaration and no personal names.

A question that has gone unexamined thus far is whether results from generation-based and

probability-based evaluations correspond with or diverge from each other, i.e., whether they have *convergent validity* [SYD23]. Convergent validity captures whether “the measurements obtained from a measurement model” (e.g., the results of a probability-based evaluation of LLM misgendering) “correlate with other measurements of the same construct” (e.g., the results of a generation-based evaluation of misgendering) that have been “obtained from measurement models for which construct validity has already been established” [JW21]. This is particularly important given that LLMs can be used in different ways, sometimes for ranking existing sequences [CC21, SLN20], and sometimes for generation, as in the popular chat models of today [Ope, Ant, Per]. However, NLP papers evaluating misgendering are not always explicit about which setups they are trying to evaluate, and sometimes probability-based evaluations are motivated by a desire to evaluate misgendering during generation. In such cases, if probability-based evaluations do not accurately estimate the propensity of models to misgender in real-world generations (e.g., due to the artificial nature of templates, limited prediction choices for the [MASK] token), they lack ecological validity, as they are “artificial [situations which do not] properly [reflect] broader real-world phenomenon” [OCD19]. In this chapter, we thus set out to *comprehensively and systematically compare evaluations of misgendering* across four pronouns (he, she, they, xe).

For our meta-evaluation, we first transform three existing datasets to measure misgendering into parallel versions for probability- and generation-based evaluation (§5.4). By automatically evaluating a suite of 6 models from 3 families, we find that these methods disagree with each other on 20.2% of evaluation instances, and 24.2% of instances for the neopronoun xe (§5.5). This suggests that probability- and generation-based evaluations lack convergent and ecological validity depending on the application context. Next, with a human evaluation of 2400 LLM generations (§5.6), we show that misgendering behaviour goes beyond pronouns [McN21], which automatic generation-based evaluations do not currently capture (e.g., avoidance of pronoun usage, meta-discourse about pronouns, extraneous gendered terms). This suggests that automatic generation-based evaluations can inherently disagree with human evaluations.

Context: *Jaime is an American actor and they are known for their roles in film.*

Generation: *In 2017, she played the role of the main character in the film “The Witch”.*

Constructed template: *Jaime is an American actor and they are known for their roles in film. In 2017, [MASK] played the role of the main character in the film “The Witch”.*

Figure 5.2: An example context from the generation-based evaluation dataset TANGO [OGD23] and a corresponding generation by Llama-3.2-1B with misgendering. The context surfaces a subject (Jaime) and base pronoun (they). The context and generation can be converted to a template to support probability-based evaluation.

Finally, we come up with recommendations for future evaluations of LLM misgendering, such as critically considering the deployment context and recognizing the contextual nature of the appropriateness of gendered terms (§5.7). Our results are pertinent to other subfields of NLP (e.g., evaluations of stereotypes, linguistic acceptability), where probability-based LLM evaluations aim to measure factors that matter in generation contexts.

5.2 Related Work

Measuring LLM misgendering. A few works have contributed evaluations for LLM misgendering. [DMO21] presents author-crafted templates to measure correct pronoun prediction for subjects with different names and pronouns, while [HDS23] builds on this with a more extensive set of templates, explicit pronoun declarations, and diverse pronoun cases for [MASK]. [GBZ24] also uses probability-based evaluation, but uses implicit pronoun declarations, up to two subjects, and no personal names. In contrast, [OGD23] proposes an automatic generation-based evaluation for misgendering in single-subject contexts. In this chapter, we conduct a meta-evaluation of the agreement of the above-mentioned probability- and generation-based evaluations of LLM misgendering, and like [OGD23], include human

evaluation in addition. We provide an overview of each dataset with the task and example inputs in Figure 5.1.

Meta-evaluations of LM bias. Various probability-based evaluations (e.g., masked token, pseudo-log-likelihood) and generated text-based evaluations (e.g., distribution, classifier, lexicon) have been proposed for measuring bias [GRB24]. In response, some prior research has explored the lack of agreement between different LM bias evaluation methods. For example, several works have highlighted the inconsistency of probability-based bias measurements using templates [DTC22, SPS22, SDK23], and the unreliability of intrinsic bias metrics to measure application bias [GMM21, CPC22]. Moreover, templatic sentences may be poorly conceptualized [BLO21] and can lack diversity, and comparisons of the probability of contrasting sentences do not capture the actual likelihood of models generating the sentences [GRB24]. In this chapter, we focus on the agreement of likelihood, lexicon, and human-based measurements of misgendering with existing datasets.

[LAN24] studies disagreements between templatic “trick tests” (i.e., acontextual probability-based evaluations designed to elicit model bias) and realistic LLM use cases. They find that templatic “trick tests” are not predictive of bias in long-form text evaluations (i.e., story generation, user personas, ESL learning exercises). Similar to their work, we contribute a more tightly-coupled method to transform “trick test” datasets in a way that enables parallel probability- and generation-based evaluation of misgendering. We frame our work from a measurement modelling perspective, broadening conceptualizations of misgendering beyond pronouns, and examining convergent and ecological validity. Similarly, [GUB23] discusses how the operationalization of bias measurements can be disconnected from how practitioners conceptualize bias, and [HSB24] examines poor evaluation validity when metrics are disconnected from deployment contexts.

Meta-evaluations of LMs in other contexts. Prior work has explored the extent to which direct LLM probabilities correlate with metalinguistic judgments [HL23b, SHM25], and have found that metalinguistic judgments are not consistent indicators of model capabilities [HL23b]. [EXK25] examine how human uncertainty can affect measurements of the agreement of human and automatic evaluations.

5.3 Evaluation Paradigms for LLM Misgendering

5.3.1 Pronoun Preliminaries

We define \mathcal{B} to be the set of all base third-person singular English pronouns, which we notationally represent using their nominative case. In line with [GBZ24], we restrict our focus to $\mathcal{B} = \{\text{he, she, they, xe}^1\}$, to study discrepancies across binary gendered pronouns, singular “they”, and a neopronoun. Each base pronoun b admits multiple cases. For instance, if b is *he*, then we have the following cases: *he* (nominative), *him* (accusative), *his* (dependent possessive), *his* (independent possessive), and *himself* (reflexive). Let p be a pronoun and $\mathcal{P}(p)$ be the base pronoun corresponding to p . Furthermore, let $\mathcal{C}(p)$ be the case of p . We also define Ω to be the set of all surface forms of pronouns we consider.

5.3.2 Probability-Based Evaluation

In probability-based evaluation, the model receives a *templatic* sequence $\{t_i\}_{i \in [T]}$ about a subject with a base pronoun y (see Figure 5.3). The template contains a single [MASK] token ($t_m = \text{[MASK]}$) associated with a grammatical case c which governs the case of any pronoun that can replace the [MASK] without violating syntactic rules. We replace the [MASK] with each pronoun in Ω with case c , and identify the pronoun \hat{y}_{prob} that reduces the perplexity of the sequence. We say that the model misgenders the subject when $\mathcal{P}(\hat{y}_{prob}) \neq y$, i.e., when

¹See the appendix of [SGS25] for a discussion on the *xe* pronoun set.

Template: *Reise's pronouns are xe/xem/xyrs. Reise was very stoic. [MASK] rarely showed any emotion.*

Misgendering: *[MASK] = He*, **No misgendering:** *[MASK] = Xe*

Constructed pre-[MASK] context: *Reise's pronouns are xe/xem/xyrs. Reise was very stoic.*

Constructed post-[MASK] context: *Reise's pronouns are xe/xem/xyrs. Reise was very stoic. Xe rarely showed any emotion.*

Figure 5.3: An example template from the probability-based evaluation dataset MISGENDERED [HDS23]. The template surfaces a subject (Reise) and base pronoun (xe). The template can be converted to pre- and post-[MASK] contexts to support generation-based evaluation.

the pronoun that makes the sequence most likely to be generated is incorrect.

5.3.3 Generation-Based Evaluation

In generation-based evaluation, the model receives a *context* sequence $\{c_i\}_{i \in [C]}$ about a subject that surfaces a base pronoun y of the subject (see Figure 5.2). The model then generates a *completion* sequence $\{g_i\}_{i \in [G]}$ for the context. We say that the model *misgenders* the subject if it uses a pronoun $\hat{y}_{gen} \in g$ to refer to the subject such that $\mathcal{P}(\hat{y}_{gen}) \neq y$. For automatic evaluation of misgendering in these generations, we use the heuristic from [OGD23], i.e., choosing \hat{y}_{gen} to be the first pronoun in the completion. Such heuristic functions are prone to error since pronoun generations can be about a different referent. Therefore, in Section 5.6, we validate this heuristic by manually annotating generations for misgendering. We provide further relevant details about and notation for the two evaluation formats in Appendix D.2.

In summary, probability-based evaluations assess whether LLMs assign a higher probability

to templatic sequences that show correct pronoun use, from a controlled set of minimally different sequences with alternative pronouns. In contrast, generation-based evaluations measure the extent to which LLMs demonstrate correct pronoun use in open-ended generations. While one would expect LLMs to be more likely to generate sequences that they assign a higher probability, there can be deviations in the results of probability- and generation-based evaluations (e.g., due to LLMs being unlikely to generate templatic sequences, the autoregressive nature of decoding).

5.4 Experimental Setup

Below, we describe the models and datasets we use for our meta-evaluation. Since these datasets were originally designed for only one type of evaluation format (either generation-based or probability-based evaluation), we transform each dataset to support the other format. This creates a tightly-coupled, fairer comparison of the two methods, so that we can better understand inconsistencies between them. Additional details are provided in Appendix D.3.

5.4.1 Models and Data

We focus on decoder-only models, as this is currently a common architecture for large language models. We select the following popular families of open-weight models: **Llama-3.1** [8B, 70B; GDJ24], **OLMo-2-1124** [7B, 13B; GBW24] for its open training data, and **Mixtral** [8x7B-v0.1, 8x22B-v0.1-4bit; JSR24], to understand the effects (if any) of a mixture-of-experts architecture. We use all three existing datasets to measure misgendering, i.e., **MISGENDERED** [HDS23], **TANGO** [OGD23], and **RUFF** [GBZ24].

5.4.2 Converting Probability-Based to Generation-Based Evaluations

To convert a probability-based evaluation dataset \mathcal{D}_{prob} into a generation-based evaluation dataset \mathcal{D}_{gen} , we transform each template $t^{(k)}$ into a context $c^{(k)}$ in two ways, as shown in Figure 5.3. The ground-truth base pronoun $y^{(k)}$ remains the same across both formats:

Pre-[MASK]: $t^{(k)}$ is truncated before the [MASK] token, i.e., $c^{(k)} \leftarrow t_{1:m-1}^{(k)}$, showing how constrained decoding might diverge from what a model would naturally generate.

Post-[MASK]: The entire template is used as the context, with the [MASK] replaced with the correct case of the ground-truth pronoun $y^{(k)}$, i.e., $c^{(k)} \leftarrow t_{1:m-1}^{(k)} \parallel R(y^{(k)}) \parallel t_{m+1:T}^{(k)}$, showing whether a model misgenders a subject even after the correct pronoun is decoded once.

5.4.3 Converting Generation-Based to Probability-Based Evaluations

To convert a generation-based evaluation dataset \mathcal{D}_{gen} into a probability-based evaluation dataset \mathcal{D}_{prob} , we transform each context and generation pair $(c^{(k)}, g^{(k)})$ into a template $t^{(k)}$, as shown in Figure 5.2. We first truncate $g^{(k)}$ such that there is only a single pronoun, and replace it with a [MASK] token, to create $g'^{(k)}$. Then, we concatenate $(c^{(k)}, g'^{(k)})$ to form $t^{(k)} = c^{(k)} \parallel g'^{(k)}$. In Appendix D.3.4, we outline practical challenges we encountered with conversion.

5.5 Agreement between Probability-Based and Generation-Based Evaluations

We measure instance-level variation within an evaluation method, as well as dataset-level agreement between probability- and generation-based evaluations, and report results on all datasets and models. These empirical experiments are complemented by brief theoretical

analyses of *why* probability- and generation-based evaluation results can disagree in Appendix D.4. In addition, we study model-level agreement in Appendix D.5.

5.5.1 Metrics

Instance-level variation. We use standard deviation to quantify correct gendering across different generations or different templates for a single instance, since generated text-based metrics are sensitive to decoding hyperparameters [AKP22, LAN24]. Let $m_{prob}^{(k)}$ be the occurrence of correct gendering ($m_{prob}^{(k)} = 1$) or not ($m_{prob}^{(k)} = 0$) for instance k in MISGENDERED or RUFF, and let $[m_{prob}^{(k)}]_i \in \{0, 1\}$ be the occurrence of correct gendering in the i -th template for instance k in Prob-TANGO. Furthermore, let $[m_{gen}^{(k)}]_i \in \{0, 1\}$ be the occurrence of correct gendering in the i -th generation for instance k in Gen-MISGENDERED, Gen-RUFF, or TANGO. Then:

$$\sigma_{gen}^{(k)} = \text{stdev}_i ([m_{gen}^{(k)}]_i), \quad \sigma_{prob}^{(k)} = \text{stdev}_i ([m_{prob}^{(k)}]_i). \quad (5.1)$$

Dataset-level agreement. To quantify agreement across probability-based and generation-based versions of each dataset, we use three metrics: Matthew’s correlation coefficient $MCC \in [-1, 1]$, raw observed agreement $p_o \in [0, 1]$, and Cohen’s $\kappa \in [-1, 1]$. See Appendix D.3.5 for details on these metrics. For $f \in \{MCC, \kappa, agr\}$, we measure the dataset-level variation v^f as:

$$v^f = f \left(\{m_{prob}^{(k)}\}_{k \in [N_{prob}]}, \{[m_{gen}^{(k)}]_1\}_{k \in [N_{prob}]} \right) \quad (\text{MISGENDERED, RUFF}), \quad (5.2)$$

$$v^f = f \left(\{[m_{prob}^{(k)}]_1\}_{k \in [N_{gen}]}, \{[m_{gen}^{(k)}]_1\}_{k \in [N_{gen}]} \right) \quad (\text{TANGO}). \quad (5.3)$$

5.5.2 Results

We report variation and agreement results by dataset, focusing on MISGENDERED and TANGO. Appendix D.5 contains plots that supplement the results in this section, as well as

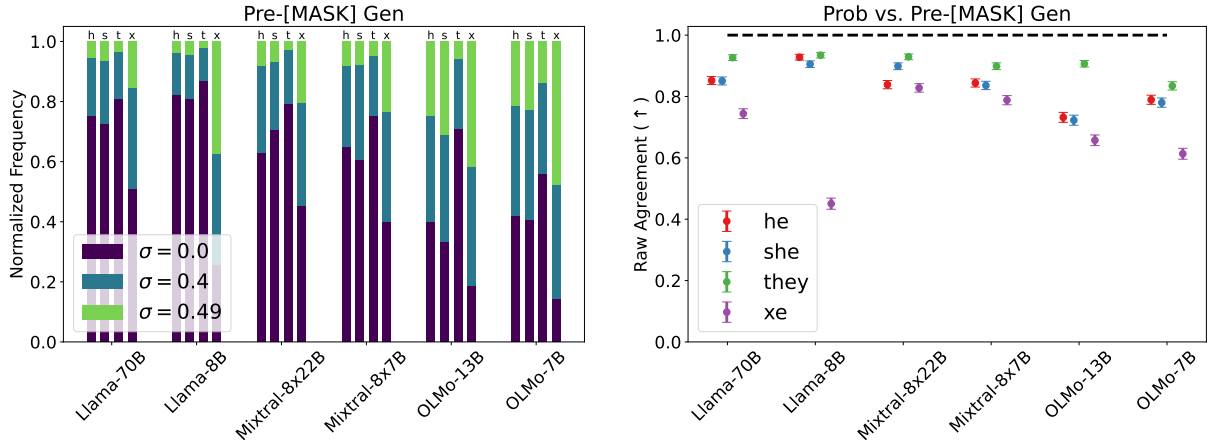
	he	she	they	xe
Llama-70B	0.004 [-0.067, 0.076]	-0.014 [-0.086, 0.057]	0.051 [-0.020, 0.122]	0.031 [-0.041, 0.102]
Llama-8B	-0.031 [-0.102, 0.041]	-0.045 [-0.117, 0.026]	0.076 [0.005, 0.147]	-0.020 [-0.092, 0.051]
Mixtral-8x22B	0.041 [-0.031, 0.112]	0.027 [-0.045, 0.098]	0.008 [-0.063, 0.080]	—
Mixtral-8x7B	0.063 [-0.008, 0.134]	0.026 [-0.046, 0.097]	-0.044 [-0.115, 0.028]	0.005 [-0.067, 0.076]
OLMo-13B	0.050 [-0.022, 0.121]	0.056 [-0.016, 0.127]	0.022 [-0.050, 0.093]	0.072 [0.000, 0.143]
OLMo-7B	0.066 [-0.005, 0.137]	0.177 [0.107, 0.246]	0.061 [-0.011, 0.132]	-0.027 [-0.098, 0.045]

Table 5.1: *MCC* agreement v^{MCC} (Eq. 5.2) between probability-based and pre-[MASK] generation-based evaluations, for each model and pronoun in MISGENDERED. We report the asymmetric 95% confidence interval, computed using SciPy [VGO20], except with xe and Mixtral-8x22B, as the model gets every instance correct in the probability-based setting.

results for RUFF, which are similar to MISGENDERED. Model-level comparisons are also included in this appendix.

MISGENDERED. As Figure 5.4a shows, there is notable instance-level variation σ in the evaluation results of models across generations for the same instance in Gen-MISGENDERED. On average, σ is highest for the neopronoun xe across all models, with the most pronounced disparity (compared to other pronouns) for Llama-8B. This shows that models exhibit *semantic instability* for xe, i.e., models are unable to consistently use xe in reference to a subject. These trends are not markedly different between the pre and post-[MASK] settings. However, σ tends to be lower on average in the post-[MASK] setting, suggesting that conditioning on an additional use of the correct pronoun improves consistency.

As for the connection between probability- and generation-based evaluations, Figure 5.4b shows that raw agreement v^{p_o} is not always high. Across all models, the evaluation methods generally agree the most on instances where the subject uses they. In contrast, the methods disagree the most when the subject uses xe, with the largest disparity for Llama-8B. This finding suggests that parallel probability- and generation-based misgendering evaluations have

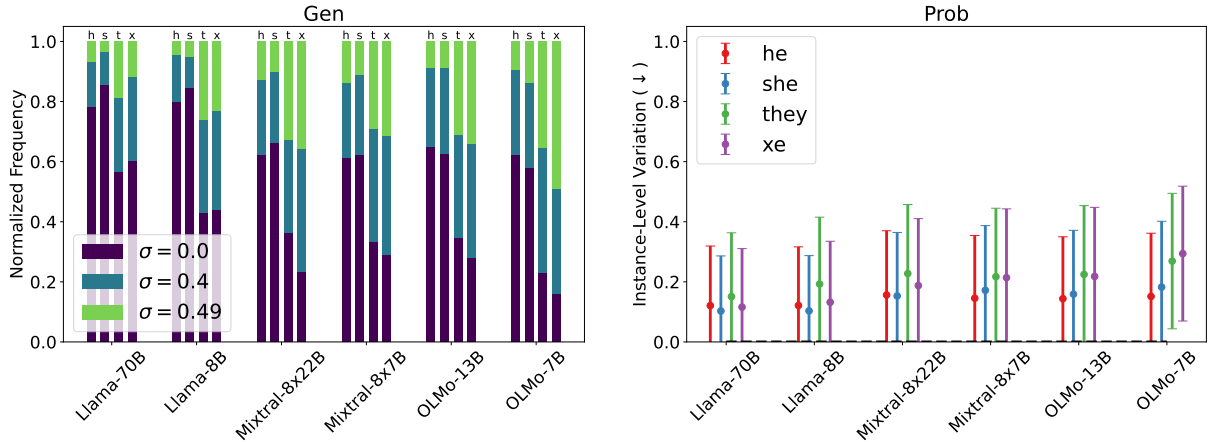


(a) Instance-level variation across generations (b) Dataset-level agreement of evaluations

Figure 5.4: Variation and agreement for MISGENDERED. **(a)** Generation variation σ (Eq. 5.1) for each model and pronoun in the pre-[MASK] generation setting. As we sample 5 generations, $\sigma \in \{0, 0.4, 0.49\}$. The bar labels h , s , t , x correspond to he , she , $they$, xe . **(b)** Raw observed agreement v^{p_o} (Eq. 5.2) for each model and pronoun between the probability-based and pre-[MASK] generation-based evaluation results. Error bars represent the standard error of v^{p_o} (computed over dataset instances). The horizontal dashed line is the upper bound of v^{p_o} .

less convergent validity for neopronoun users, which is problematic as misgendering already disproportionately harms neopronoun users. Table 5.1 provides a complementary perspective, with v^{MCC} instead of raw agreement. For all models and pronouns, v^{MCC} and v^κ are close to 0, which indicates a weak association between the probability- and generation-based evaluation results. This is because the evaluation results are often imbalanced (i.e., there is a high probability of chance agreement). These trends are not markedly different between the pre and post-[MASK] settings.

TANGO. Figure 5.5 shows notable variation σ in the evaluation results across templates and generations for the same instance in TANGO and Prob-TANGO, respectively. In the



(a) Instance-level variation across generations (b) Instance-level variation across templates

Figure 5.5: Instance-level variation σ (Eq. 5.1) for each model and pronoun with TANGO. **(a)** Generation-Based variation. The bar labels h, s, t, x correspond to he, she, they, xe. **(b)** Probability-Based variation. As we exclude templates with no pronoun, we do not always have 5 templates per instance (see Figure D.3). Hence, we report the mean and standard deviation.

generation-based setting, across all models, σ appears to be highest on average for **they** and **xe**. Agreement between probability- and generation-based evaluations is better for TANGO and Prob-TANGO than for MISGENDERED and Gen-MISGENDERED (and RUFF and Gen-RUFF), as Table 5.2 indicates a moderate association between results from both methods. Disagreements pattern similarly to MISGENDERED in that most disagreements happen when the subject uses **xe**. Interestingly, there are also pronounced disagreements for the pronoun **they**, with the most pronounced disparities for the Mixtral models. Overall, our results suggest that the templates in MISGENDERED and RUFF are unlikely to be generated by the LLMs that we consider, which threatens their validity.

RUFF. Similar to Gen-MISGENDERED, Gen-RUFF also displays instance-level variation across generations, and disagreement in results between probability- and generation-based

	he	she	they	xe
Llama-70B	0.686 [0.633, 0.732]	0.511 [0.440, 0.575]	0.756 [0.710, 0.795]	0.552 [0.480, 0.616]
Llama-8B	0.578 [0.513, 0.637]	0.505 [0.433, 0.570]	0.732 [0.684, 0.774]	0.552 [0.480, 0.616]
Mixtral-8x22B	0.548 [0.475, 0.613]	0.644 [0.585, 0.697]	0.554 [0.481, 0.619]	0.442 [0.354, 0.523]
Mixtral-8x7B	0.691 [0.637, 0.739]	0.514 [0.439, 0.583]	0.653 [0.591, 0.708]	0.398 [0.305, 0.485]
OLMo-13B	0.574 [0.504, 0.637]	0.576 [0.508, 0.637]	0.690 [0.634, 0.739]	0.568 [0.490, 0.637]
OLMo-7B	0.633 [0.571, 0.689]	0.463 [0.382, 0.538]	0.619 [0.552, 0.678]	0.673 [0.611, 0.727]

Table 5.2: MCC agreement v^{MCC} (Eq. 5.2) between probability- and generation-based evaluation for each model and pronoun in TANGO. We report the asymmetric 95% confidence interval, computed using SciPy [VGO20].

evaluation. In contrast to Gen-MISGENDERED, the methods tend to disagree the most with respect to v^{po} when the subject uses **they**, with the largest disparity for Llama-8B. However, with respect to v^{MCC} and v^κ , the methods have a low but higher agreement for **they** compared to other pronouns. This could be attributed to RUFF templates not containing personal names, which seem to be more polarizing in their gendered associations for LLMs.

5.6 Human Evaluation

We perform human evaluation with the dual goals of validating the automatic metric for generation-based evaluations (similar to [OGD23]), and to get a more granular view of LLM misgendering. In contrast to [OGD23], who focus on TANGO, we focus on Gen-MISGENDERED and Gen-RUFF. In Appendix D.7, we provide qualitative examples of interesting generations from our human evaluation.

Methodology. Two authors, both experts in English-language pronoun usage, annotated a total of 2400 generations – 25 pre-[MASK] and 25 post-[MASK] generations per ground-truth pronoun, for all models and the two datasets. Each generation was annotated for whether:

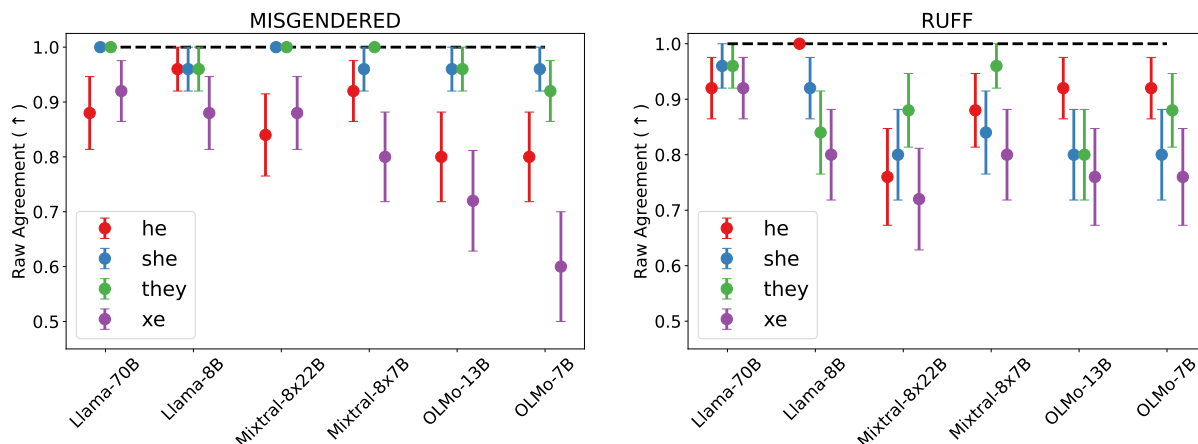


Figure 5.6: Agreement between human and automatic evaluation of misgendering in the pre-[MASK] generation setting. Many models fall short of human-human agreement (96%).

(1) the ground-truth pronoun is correctly used, (2) misgendering occurs, or (3) no pronoun is generated. See Appendix D.6 for the full annotation schema, and some observations made during annotation. In addition, we noted when models introduced extraneous gendered words such as “man,” “girl,” etc. On a sample of 200 instances that both authors annotated, agreement was 96% for pronoun labelling, and 98% for extraneous gendered information.

Validation of automatic results. Since our annotation schema has three options and the automatic generation-based evaluation heuristic we use is binary, we treat only case (2) as misgendering, and the other two cases as a lack of misgendering (i.e., as correct) to validate the heuristic. We find (see Figures 5.6, D.9 in the appendix) that automatic and human evaluation of pronoun misuse do not always agree, and this happens for multiple reasons; incorrect pronoun use sometimes appears later in the generated text, which the automatic evaluation misses. The automatic evaluation also cannot distinguish when a different pronoun is used because of misgendering or just for a different person or entity than the original subject.

Even when human and automatic evaluation agree, this can be due to incoherent, repetitive

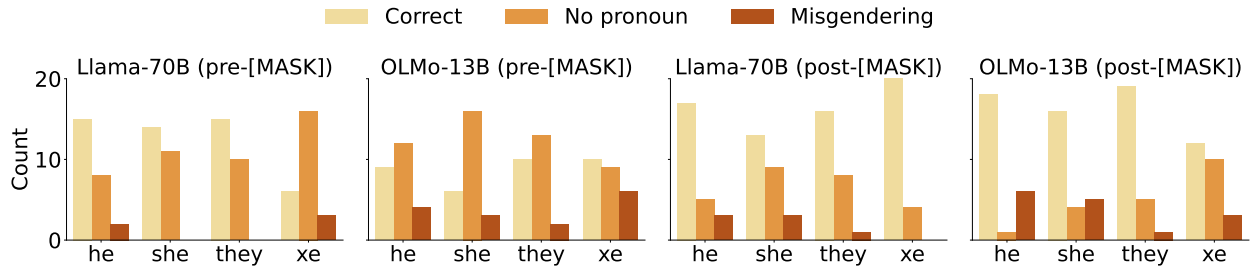


Figure 5.7: Human annotations of Llama-70B and OLMo-13B generations from the pre-[MASK] (left) and post-[MASK] (right) settings.

generations. To quantify this, we measure the repetition rate of generations in Appendix D.8, following [BSC14]. Previous work has also noted that lexicon-based metrics can miss such higher-level structures in sentences [GRB24].

Other patterns in misgendering. Human evaluation lets us disentangle different causes of the absence of misgendering, as shown in Figure 5.7. Models frequently avoid generating pronouns, instead just repeating names or occupations, as exemplified by OLMo-13B with pre-[MASK] generation. The rate of this also varies across pronoun sets, with Llama-70B in particular avoiding pronoun generation for the neopronoun `xe` in the pre-[MASK] setting, despite correctly generating other pronouns. Neopronouns are low-performing across all our results as in previous work, and generations sometimes included a surface form in the wrong declension (still marked as correct), or even other neopronouns (e.g., `ze`, `hir`). Interestingly, the post-[MASK] setting seems to encourage correct pronoun generation across models, corroborating our previous instance-level agreement results.

Another facet of misgendering is extraneous gendered mentions, shown in Figure 5.8. These reflect potential misgendering beyond pronouns, e.g., given the context “Elizabeth’s pronouns are he/him/his,” some models would insist that “Elizabeth is a girl,” despite nothing explicit in the context to indicate this. Extraneous gendered mentions appear frequently in MISGENDERED compared to RUFF, presumably because the dataset’s focus on personal

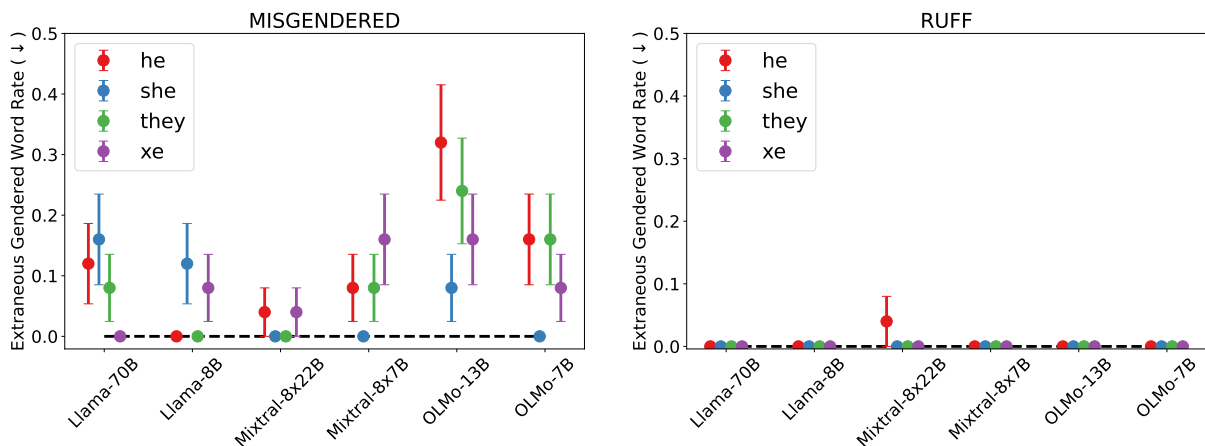


Figure 5.8: Proportion of generations with extraneous gendered words in the pre-[MASK] generation setting. MISGENDERED contains named subjects with pronoun declarations, which seem to elicit more extraneous gendered cues than RUFF, which contains occupations.

names and pronoun declarations elicits stronger model assumptions about gender and more potential misgendering that is not currently measured by pronoun-focused, lexicon-based evaluations [GSL24]. However, whether or not this is actually misgendering is a complex question with contextual answers for real individuals [McN21].

Finally, we touch on two aspects of misgendering that we did not systematically annotate. First, several model generations included meta-commentary about pronouns and reference. Although this has been shown to be independent of what probability-based evaluations indicate [HL23b], the patterns here are interesting to study in their own right. We also noticed that models seem to generate the same pronoun sets for other participants in a given situation, e.g., a **she** doctor would talk to a **she** patient, and even a **xe** programmer might talk to **xyr** **xe** boss, an exploration of which we leave to future work.

5.7 Recommendations

Our meta-evaluation reveals limitations with the convergent validity, ecological validity, and operationalization of misgendering evaluations in NLP. Based on our insights, we make the following recommendations for future work in the field:

- **Use the evaluation that is appropriate to the final deployment**, i.e., open-ended generation-based evaluations for open-ended generation-based applications, probability-based evaluations for probability-based applications, dialogue-based evaluations for dialogue, and so on.
- **Take a holistic view of misgendering**, that accounts for all aspects of potential misgendering, including extraneous gendered words beyond pronouns (in English), meta-discourse about pronouns and gender, and so on [HDS24].
- **Recognize that misgendering is contextual** — lexicon-based approaches may not capture nuance, as gendered words may be appropriate in certain contexts, and even gender-neutral words can be used disrespectfully [DW18a].
- **Center those most impacted by misgendering in system design and evaluation**, from broad and application-specific conceptualizations of misgendering, as well as operationalization of data, metrics and evaluation choices [SB24].

5.8 Conclusion and Future Work

In this work, we comprehensively compare generation-based and probability-based evaluations of misgendering in LLMs, by adapting three existing misgendering datasets for a parallel meta-evaluation. Our results show that these two evaluation approaches do not always converge, with disagreements on roughly 20% of instances. Through human evaluation, we show that misgendering is multifaceted and goes beyond just incorrect pronoun use. To

address this, we recommend using community-grounded, holistic definitions of misgendering. More broadly, our empirical findings highlight the need for deliberate, reliable, and ecologically valid evaluation protocols. These findings are relevant beyond misgendering, in other subfields of NLP where relying on probability-based evaluations alone may fail to capture phenomena that occur in open-ended generation, and vice versa. We discuss the limitations of our work in Appendix D.1.

5.9 Broader Impacts

As we are concerned with a meta-evaluation of misgendering, we take steps to ensure that our experimental setup does not miss misgendering. We do this through human validation of automatic results, and by refraining from using systems that might introduce additional performance biases. For instance, we avoid using off-the-shelf coreference resolution systems to identify which pronouns refer to the subject in automatic evaluations, in order to avoid additional performance biases introduced by these systems, e.g., the inability to recognize neopronouns and certain names as referents [DMO21, CD21]. We do not conduct our meta-evaluation with closed models, for which one could not verify whether the misgendering datasets are not part of their pretraining data. We will release our code and data for research purposes and reproducibility, and request that other researchers use these resources accordingly. We will not release our data in plain text, to avoid polluting information ecosystems with additional instances of misgendering.

Part III

Towards a Precise Theory of Machine
Learning Unfairness

CHAPTER 6

An Effective Theory of Bias Amplification

6.1 Introduction

Machine learning (ML) datasets can encode a plethora of biases which, when said data is used to train models, can result in systems that can cause practical harm. Datasets that encode correlations that only hold for a subset of the data may cause disparate performance when models are used more broadly, such as an X-ray pneumonia classifier that only functions on images from certain hospitals [ZBL18]. This issue is magnified when coupled with under-representation, whereby a dataset fails to adequately reflect parts of the underlying data distribution, often further marginalizing certain groups. Lack of representation results in systems that might work well on average, but fail for minoritized groups, including facial recognition systems that fail for darker-skinned women [BG18b], large language models that consistently misgender transgender and nonbinary people [OGD23], or image classification technology that only works in Western contexts [dMW19, RKB24].

Unfortunately, contemporary models may exhibit *bias amplification*, whereby dataset biases are not only replicated, but exacerbated [ZWY17, HBS18, WR21]. While previous research has shown that amplification is a function of both dataset properties and how we choose to construct our models [HvG22, SRK20, BS23], it is not fully clear how bias amplification occurs mechanistically, nor do we precisely understand which settings lead to its emergence. Thus, in this chapter, we propose a novel theoretical framework that explains how model design choices (e.g., number of parameters, regularization penalty) and data

distributional properties (e.g., number of features, group imbalance, label noises) interact to amplify bias. Moreover, our framework provides an account of prior work on bias amplification [BS23] and minority-group bias [SRK20].

A theory of bias amplification is important for several reasons. First, as empirical research necessarily yields only sparse data points—often focused on only the most common regimes—theory allows us to interpolate between past findings, and reason about how bias emerges in under-explored settings. Second, a precise theory gives us the depth of understanding needed in order to intervene, potentially supporting the development of both novel evaluations and mitigations. Finally, beyond explaining already-known phenomena, our theory makes new predictions, suggesting new avenues for future research.

6.1.1 Main Contributions

In this chapter, we develop a unifying and rigorous theory of ML bias in the settings of ridge regression with and without random projections. In particular, we precisely analyze test error disparities between groups (e.g., demographic groups or protected categories) with different data distributions when training on a mixture of data from these groups. We characterize these disparities in high dimensions using operator-valued free probability theory (OVFPT), thereby avoiding possibly loose bounds on critical quantities. Our theory encompasses different parameterization regimes, group sizes, label noises, and data covariance structures. Moreover, our theory has applications to important problems in ML bias that have recently been empirically investigated:

- **Bias amplification.** Even in the absence of group imbalance and spurious correlations, a single model that is trained on a combination of data from different groups can amplify bias beyond separate models that are trained on data from each group [BS23]. With our theory, we reproduce and analyze the bias amplification findings of [BS23] in controlled settings. We further observe how stopping model training early or tuning the regularization

hyperparameter can alleviate bias amplification.

- **Minority-group bias.** Overparameterization can hurt test performance on minority groups due to spurious features [SRK20, KL21]. We theoretically analyze how model size and extraneous features affect minority-group bias.

We extensively empirically validate our theory in controlled and semi-synthetic settings. Specifically, we show that our theory aligns with practice in the cases of: (1) bias amplification with synthetic data generated from isotropic covariance matrices and the semi-synthetic dataset Colored MNIST [ABG19], and (2) minority-group bias under different model sizes with synthetic data generated from diatomic covariance matrices. In these applications, we expose new, interesting phenomena in various regimes. For example, a larger number of features than samples can amplify bias under overparameterization, there may be an optimal regularization penalty or training time to avoid bias amplification, and there can be differences in test error between groups that are not alleviated with increased parameterization. Our observations of phenomena in Sections 6.4 and 6.5 are largely empirical but are supported by their agreement with our theory. Our theory of ML bias can inform strategies to evaluate and mitigate unfairness in ML, or be used to caution against the usage of ML in certain applications.

6.1.2 Related Work

Bias amplification. A long line of research has explored how ML exacerbates biases in data. For example, a single model that is trained on a combination of data from different groups can amplify bias [ZWY17, WR21], even beyond what would be expected when separate models are trained on data from each group [BS23]. [HvG22] conducts a systematic empirical study of bias amplification in the context of image classification, finding that amplification can vary greatly as a function of model size, training set size, and training time. Furthermore, overparameterization, despite reducing a model’s overall test error, can disproportionately

hurt test performance for minority groups [SRK20, KL21]. Models can also overestimate the importance of poorly-predictive, low-signal features for minority groups, thereby hurting performance on these groups [LFB19]. In this chapter, we distill a holistic theory of how model design choices and data distributional properties affect disparate test performance across groups, which can encompass seemingly disparate bias phenomena.

High-dimensional analysis of ML. A suite of works have analyzed the expected dynamics of ML in appropriate asymptotic scaling limits, e.g., the rate of features d to samples n converges to a finite values as d and n respectively scale towards infinity [AP20a, TAP21, LMH23]. Notably, [Bac23] theoretically analyzes the double descent phenomenon [SGd19, BHM19] in ridge regression with random projections by computing deterministic equivalents for relevant random matrix quantities in a proportionate scaling limit. Like [AP20a, TAP21, LMH23], we leverage the tools of OVFPT [MS17], which is at the intersection of random matrix theory (RMT) and functional analysis. However, [AP20a] focuses on training and testing a random features model on data from the same Gaussian distribution. Furthermore, [TAP21, LMH23] focus on training a random features model on data from one Gaussian distribution and testing the model on a different Gaussian. In contrast, we study the random features model in the setting of training on a mixture of Gaussian distributions and testing on each component. Because a mixture is more expressive than a single Gaussian, our theoretical results cannot be derived as a special case of these other works. Furthermore, our theory non-trivially generalizes [Bac23], which we recover in Corollary E.8.1 as a special case, and requires more powerful analytical techniques.

Certain prior theoretical work precisely analyzes the bias of models trained on a mixture of data from different groups in a high-dimensional setting [MGR24, JNB24]. Like [MGR24, JNB24], we study linear models and consider bias as the disparity in test performance of a model between groups. We further consider some similar factors that give rise to bias amplification (e.g., group imbalance, group data variance, inter-group similarity, dataset

size). We also share some theoretical conclusions, such as bias can occur even when the groups have the same ground-truth weights (see Section 6.5) and are balanced (Section 6.4.1). Additionally, we both discuss the paradigms of training a single model for both groups vs. separate models for each group. However, the *main distinction* between our work and [MGR24, JNB24] is that we precisely characterize how models amplify bias in different *parameterization regimes*, that is, we examine the impact of model size on bias. This enables us to expose new, richer insights into the impact of over- and underparameterization on bias amplification (see Figure 6.1, Section 6.4, and Section 6.5).

Beyond this, [MGR24] employs the replica method, which is non-rigorous, while we use OVFPT, which is entirely rigorous. Moreover, [MGR24, JNB24] study the application of linear classification to Gaussian data with isotropic covariance; in contrast, we study the application of regression with random projections (a simplified model of feedforward neural networks) to Gaussian data with more general covariance structure (i.e., covariance matrices that are simultaneously diagonalizable) and noisy labels. This allows us to analyze the effects of these additional factors on bias. We make additional connections between our work and [MGR24, JNB24] in Section 6.4.2.

6.2 Preliminaries

6.2.1 Data Distributions

We consider a ridge regression problem on a dataset from the following multivariate Gaussian mixture with two groups $s = 1$ and $s = 2$. These groups could represent different demographic

groups or protected categories.

$$\text{(Group ID)} \quad \text{Law}(s) = \text{Bernoulli}(p), \quad (6.1)$$

$$\text{(Features)} \quad \text{Law}(x \mid s) = \mathcal{N}(0, \Sigma_s), \quad (6.2)$$

$$\text{(Ground-truth weights)} \quad \text{Law}(w_1^*) = \mathcal{N}(0, \Theta/d), \quad \text{Law}(w_2^* - w_1^*) = \mathcal{N}(0, \Delta/d), \quad (6.3)$$

$$\text{(Labels)} \quad \text{Law}(y \mid s, x) = \mathcal{N}(f_s^*(x), \sigma_s^2), \quad \text{with } f_s^*(x) := x^\top w_s^*. \quad (6.4)$$

The scalar $p \in (0, 1)$ controls for the relative size of the two groups (e.g., $p = 1/2$ in the balanced setting). For simplicity of notation, we define $p_1 = p$ and $p_2 = 1 - p$. The $d \times d$ positive-definite matrices Σ_1 and Σ_2 are the covariance matrices for the different groups. The d -dimensional vectors w_1^* and w_2^* are the ground-truth weights vectors for each group. w_1^* and $w_2^* - w_1^*$ are independently sampled from zero-mean Gaussian distributions with covariances Θ/d and Δ/d , respectively. In particular, setting $\Delta = 0$ corresponds to the case that both groups have identical ground-truth weights. We define $\Theta_1 = \Theta$, $\Theta_2 = \Theta + \Delta$. Finally, σ_s^2 corresponds to the label noise for each group s . While we consider the case of two groups only for conciseness, our theoretical methods readily extend to any finite number of groups.

6.2.2 Models and Metrics

Learning. A learner is given an IID sample $\mathcal{D} = \{(x_1, y_1), \dots, (x_n, y_n)\} \equiv (X \in \mathbb{R}^{n \times d}, Y \in \mathbb{R}^n)$ of data from the above distribution and it learns a model for predicting the label y from the feature vector x . Thus, X is the total design matrix with i th row x_i , and y the total response vector with i th component y_i . Let $\mathcal{D}^s = (X \in \mathbb{R}^{n_s \times d}, Y \in \mathbb{R}^{n_s})$ be the data pertaining only to group s , so that $\mathcal{D} = \mathcal{D}^1 \cup \mathcal{D}^2$ is a partitioning of the entire dataset. Two choices are available to the learner: (1) learn a model $\hat{f}_s \in \mathcal{F}$ on each dataset \mathcal{D}^s , or (2) learn a single model $\hat{f} \in \mathcal{F}$ on the entire dataset \mathcal{D} . In practice, a choice is made based on scaling vs. personalization considerations.

We consider two solvable settings for linear models: classical ridge regression in the ambient input space, and ridge regression in a feature space given by random projections.

The latter allows us to study the role of model size in ML bias, by varying the output dimension of the random projection mapping. This output dimension m controls the size of a feedforward neural network in a simplified regime [MRS22, Bac23].

Classical Ridge Regression. We will first consider the function class $\mathcal{F} \subseteq \{\mathbb{R}^d \rightarrow \mathbb{R}\}$ of linear ridge regression models without random projections. For any vector $w \in \mathbb{R}^d$, the model f with parameters w is defined by $f(x) = x^\top w$, for all $x \in \mathbb{R}^d$, and is learned with ℓ_2 -regularization. We define the generalization error or risk of any model f with respect to group s as:

$$R_s(f) = \mathbb{E} [(f(x) - f_s^*(x))^2 \mid s]. \quad (6.5)$$

We consider ridge regression because in addition to its analytical tractability, it can be viewed as the asymptotic limit of many learning problems [DW18b, RMR21, HMR22]. We now formally define some metrics related to bias amplification.

Definition 6.2.1 (Bias Amplification). *We isolate the contribution of the model to bias when learning from data with different groups. This intuitive conceptualization of bias amplification allows us to quantify the phenomenon. Grounded in the literature [BS23], we define the Expected Difficulty Disparity (EDD) as:*

$$EDD = |\mathbb{E} R_2(\hat{f}_2) - \mathbb{E} R_1(\hat{f}_1)|, \quad (6.6)$$

where the expectations are w.r.t. randomness in the training data and any other sources of randomness in the models. The EDD captures the difference in test risk between models trained and evaluated on each group separately. In contrast, we define the Observed Difficulty Disparity (ODD) as:

$$ODD = |\mathbb{E} R_2(\hat{f}) - \mathbb{E} R_1(\hat{f})|. \quad (6.7)$$

The ODD captures the bias (i.e., difference in test risk between groups) of a model trained on both groups. Finally, we define the Amplification of Difficulty Disparity (ADD) as $ADD = \frac{ODD}{EDD}$. We say that bias amplification occurs when $ADD > 1$.

Our definition of *ADD* is consistent with the conceptualization of bias of [BS23]. At a high level, our definition quantifies how many times worse model bias would be if a ML practitioner opted to train a single model on a mixture of data from two groups (i.e., the setting in which bias is observed in practice) vs. separate models for the data from each group (i.e., the setting which corresponds to the bias in the data alone, and thus the a priori amount of bias we would expect in the case of a single model). In sum, we seek to isolate the contribution of the *model* to bias when learning from data with different groups.

Ridge Regression with Random Projections. We consider feedforward neural networks in a simplified regime which can be approximated via random projections, i.e., a one-hidden-layer neural network $f(x) = v^\top Sx$ with a *linear* activation function. In particular, we extend classical ridge regression by transforming our learned weights as $\hat{w} = S\hat{\eta} \in \mathbb{R}^d$, where $S \in \mathbb{R}^{d \times m}$ is a random projection with entries that are IID sampled from $\mathcal{N}(0, 1/d)$. Ridge regression with random projections offers analytical tractability while exposing bias amplification phenomena related to model size; such phenomena are not exposed by classical ridge regression (see Figure 6.1). Moreover, it has been shown that in high dimensions, training a one-hidden-layer neural network with gradient descent effectively learns a linear predictor over random features [YS19]. Furthermore, [AP20b, Bac23]; inter alia are able to reproduce interesting phenomena like double descent using the random features model. Nevertheless, [YS19] has shown that the model often cannot learn even a ReLU neuron, suggesting that some mechanisms of bias amplification could be different in nonlinear networks.

6.3 Theoretical Analysis

Assumptions. Some of our theorems will require standard technical assumptions that we detail here and in Appendix E.1. Assumption 6.3.1 describes the proportionate scaling limits, standard in RMT, in which we will work. These limits enable us to derive deterministic

analytical formulae for the expected test risk of models. Our experiments (see Sections 6.4 and 6.5) validate our theory.

Assumption 6.3.1. *In the case of ridge regression with random projections, we will work in the following proportionate scaling limit:*

$$n, n_1, n_2, d \rightarrow \infty, \quad n_1/n \rightarrow p_1, n_2/n \rightarrow p_2, d/n \rightarrow \phi, m/n \rightarrow \psi, m/d \rightarrow \gamma, \quad (6.8)$$

$$d/n_1 \rightarrow \phi_1, m/n_1 \rightarrow \psi_1, \quad d/n_2 \rightarrow \phi_2, m/n_2 \rightarrow \psi_2, \quad (6.9)$$

for some constants $\phi_1, \phi_2, \phi, \psi_1, \psi_2, \psi \in (0, \infty)$. The scalar ϕ captures the rate of features to samples. Observe that $\phi = p_1\phi_1$ and $\phi = p_2\phi_2$. We note that $\phi\gamma = \psi$ and $\phi_s\gamma = \psi_s$. The scalar ψ captures the rate of parameters to samples. The setting $\psi > 1$ (resp. $\psi < 1$) corresponds to the overparameterized (resp. underparameterized) regime.

6.3.1 Main Result: Ridge Regression with Random Projections

To provide a mechanistic understanding of how ML models may amplify bias, our theory elucidates differences in the test error between groups when a single model is trained on a combination of data from both groups vs. when separate models are trained on data from each group. We first consider the classical ridge regression model in Appendix E.2 before studying ridge regression with random projections below, which is a more realistic but still analytically solvable setup.

Single Random Projections Model Learned for Both Groups. We first consider

the ridge regression model \hat{f} with random projections, which is learned using empirical risk minimization and ℓ_2 -regularization with penalty λ . The parameter \hat{w} of the linear model \hat{f} is given by the following optimization problem:

$$\hat{w} = S\hat{\eta} \in \mathbb{R}^d, \quad \text{with } \hat{w} = \arg \min_{\eta \in \mathbb{R}^m} L(\eta) = \sum_{s=1}^2 n^{-1} \|X_s S \eta - Y_s\|_2^2 + \lambda \|\eta\|_2^2. \quad (6.10)$$

Explicitly, one can write $\widehat{w} = S(Z^\top Z + n\lambda I_m)^{-1}Z^\top Y$, where $Z := XS$. Before presenting our result for the random projections model, we provide some relevant definitions.

Definition 6.3.1. Let $\bar{\text{tr}} A := (1/d)\text{tr} A$ be the normalized trace operator and $(e_1, e_2, \tau, u_1, u_2, \rho)$ be the unique positive solution to the following fixed-point equations:

$$1/\tau = 1 + \bar{\text{tr}} LK^{-1}, \quad 1/e_s = 1 + \psi\tau\bar{\text{tr}} \Sigma_s K^{-1}, \quad \text{for } s \in \{1, 2\}, \quad (6.11)$$

$$\rho = \tau^2\bar{\text{tr}} (\gamma\rho L^2 + \lambda^2 D)K^{-2}, \quad u_s = \psi e_s^2\bar{\text{tr}} \Sigma_s (\gamma\tau^2 D + \rho I_d)K^{-2}, \quad \text{for } s \in \{1, 2\}, \quad (6.12)$$

$$\text{where: } L = p_1 e_1 \Sigma_1 + p_2 e_2 \Sigma_2, \quad K = \gamma\tau L + \lambda I_d, \quad D = p_1 u_1 \Sigma_1 + p_2 u_2 \Sigma_2 + B. \quad (6.13)$$

For deterministic $d \times d$ PSD matrices A and B , we define the following auxiliary quantities:

$$h_j^{(1)}(A) := p_j \gamma e_j \tau \bar{\text{tr}} A \Sigma_j K^{-1}, \quad (6.14)$$

$$h_j^{(2)}(A, B) := p_j \gamma \bar{\text{tr}} A \Sigma_j (\gamma e_j \tau^2 B + p_{j'} \gamma \tau^2 \Sigma_{j'} (e_j u_{j'} - e_{j'} u_j) + e_j \rho I_d - \lambda u_j \tau I_d) K^{-2}, \quad (6.15)$$

$$\begin{aligned} h_j^{(3)}(A, B) &:= p_j \bar{\text{tr}} A \Sigma_j (\gamma e_j^2 p_j \Sigma_j (p_{j'} \gamma \tau^2 u_{j'} \Sigma_{j'} + \gamma \tau^2 B + \rho I_d) \\ &\quad + u_j (p_{j'} \gamma e_{j'} \tau \Sigma_{j'} + \lambda I_d)^2) K^{-2}, \end{aligned} \quad (6.16)$$

$$\begin{aligned} h_j^{(4)}(A, B) &:= p_j \gamma p_{j'} \bar{\text{tr}} \Sigma_j \Sigma_{j'} A (\gamma \tau^2 (e_j e_{j'} B - p_j e_j^2 u_{j'} \Sigma_j - p_{j'} \Sigma_{j'} e_{j'}^2 u_j) \\ &\quad - \lambda \tau (e_j u_{j'} + e_{j'} u_j) I_d + e_j e_{j'} \rho I_d) K^{-2}. \end{aligned} \quad (6.17)$$

In Appendix E.6, we intuitively interpret the scalars e_s, τ, u_s, ρ in the setting where a separate model is learned for each group. In essence, our theory extends the scalars to the more general setting where a single model is trained on a mixture of data from different groups. Furthermore, each of the terms $h_j^{(1)}, \dots, h_j^{(4)}$ capture the limiting values of different sources of covariance between the sample covariance matrices for the groups, the resolvent matrix, and the random projections matrix S . These sources of covariance are written explicitly in Appendix E.5, and naturally arise from expanding the solution to the ridge regression problem with random projections.

We now present Theorem 6.3.1, which is our *main contribution*. Theorem 6.3.1 presents a novel bias-variance decomposition for the test error $R_s(\widehat{f})$ for each group $s \in \{1, 2\}$ in

the context of ridge regression with random projections. It is a non-trivial generalization of theories in high-dimensional ML which requires the powerful machinery of OVFPPT (see proof in Appendix E.5).

Theorem 6.3.1. *Under Assumptions E.1.2 and 6.3.1, it holds that $R_s(\widehat{f}) \simeq B_s(\widehat{f}) + V_s(\widehat{f})$, with*

$$V_s(\widehat{f}) = \sum_{j=1}^2 \sigma_j^2 \phi h_j^{(2)}(I_d, \Sigma_s), \quad (6.18)$$

$$B_s(\widehat{f}) = \bar{\text{tr}} \Theta_s \Sigma_s + h_1^{(3)}(\Theta_s, \Sigma_s) + h_2^{(3)}(\Theta_s, \Sigma_s) + 2h_1^{(4)}(\Theta_s, \Sigma_s) \quad (6.19)$$

$$- 2h_1^{(1)}(\Theta_s \Sigma_s) - 2h_2^{(1)}(\Theta_s \Sigma_s) + h_{s'}^{(3)}(\Delta, \Sigma_s) \quad (6.20)$$

$$- 2 \begin{cases} 0, & s = 1, \\ h_1^{(3)}(\Delta, \Sigma_2) + h_2^{(4)}(\Delta, \Sigma_2) - h_1^{(1)}(\Delta \Sigma_2), & s = 2. \end{cases} \quad (6.21)$$

The unregularized limit corresponds to the minimum-norm interpolator, and alternatively may be viewed as training a neural network until convergence [AKT19]. We discuss methods for, and the complexity of, solving the above fixed-point equations analytically in Appendix E.7. Furthermore, in Appendix E.8, we directly express the bias and variance of the test risk of an unregularized model trained on just group s in terms of the second and first-order degrees of freedom of Σ_s and the parameterization rate ψ_s . Moreover, in Appendix E.11, we derive the approximate bias amplification profile of an unregularized model with respect to the ratio $c = \sigma_2^2/\sigma_1^2$ of label noises, in the setting where the eigenspectra of the covariance matrices have power-law decay.

Separate Random Projections Model Learned for Each Group. We now consider the ridge regression models \widehat{f}_1 and \widehat{f}_2 with random projections, which are learned using empirical risk minimization and ℓ_2 -regularization with penalties λ_1 and λ_2 , respectively. In particular, we have the following optimization problem for each group s : $\arg \min_{\eta \in \mathbb{R}^m} L(w) =$

$n_s^{-1}\|X_s S\eta - Y_s\|_2^2 + \lambda_s\|\eta\|_2^2$. Alternatively, the reader can think of each \widehat{f}_s as the limit of \widehat{f} when $p_s \rightarrow 1$. In this setting, we deduce Theorem 6.3.2, which follows from Theorem 6.3.1.

Theorem 6.3.2. *Under Assumptions E.1.2 and 6.3.1, it holds that $R_s(\widehat{f}_s) \simeq B_s(\widehat{f}_s) + V_s(\widehat{f}_s)$, where $V_s(\widehat{f}_s) = \lim_{p_s \rightarrow 1} V_s(\widehat{f})$ and $B_s(\widehat{f}_s) = \lim_{p_s \rightarrow 1} B_s(\widehat{f})$ (see Appendix E.6 for explicit formulae).*

Phase Diagram. The phase diagram for the random projections model (Figure 6.1) offers rich insights into how the rate of parameters to samples (ψ), in interaction with the rate of features to samples (ϕ), affects bias amplification. In the *ODD* and *EDD* profiles, we observe phase transitions at $\phi = \psi$ (when $\psi < 0.5$) and $\psi = 0.5$ (i.e., $\psi_1 = \psi_2 = 1$), where these metrics begin decreasing significantly. $\psi_s = 1$ is a known interpolation threshold for random features models [AP20b, DRB20]. In contrast, at $\psi = 1$ and $\phi = 1$, the *ODD* drastically increases. Furthermore, at $\phi = \psi$ (when $\psi < 0.5$) and $\phi = 0.5$ (for $\psi > 0.5$), the *EDD* greatly increases. Accordingly, in the *ADD* profile, we observe phase transitions at $\phi = \psi$ (when $\psi < 0.5$), $\psi = 0.5$, $\psi = 1$, and $\phi = 1$, where bias amplification begins occurring (i.e., $ADD > 1$). However, bias seems to be deamplified (i.e., $ADD < 1$) at $\phi = \psi$ (when $\psi < 0.5$) and $\phi = 0.5$ (when $\psi > 0.5$). Some observations are less visible due the granularity of the color thresholding in Figure 6.1.

6.4 Bias Amplification

We empirically show how ridge regression models with random projections may amplify bias when a single model is trained on a combination of data from different groups vs. when separate models are trained on data from each group [BS23]. We further show how our theory: (1) predicts bias amplification, and (2) exposes new, interesting bias amplification phenomena in various regimes.

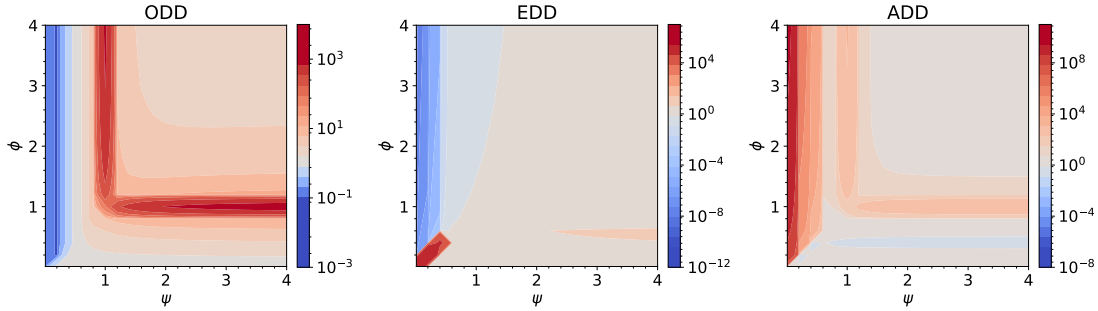


Figure 6.1: *ODD*, *EDD*, and *ADD* phase diagrams for ridge regression with random projections. We plot the bias amplification phase diagrams with respect to ϕ (rate of features to samples) and ψ (rate of parameters to samples), as predicted by our theory for ridge regression with random projections (Theorems 6.3.1, 6.3.2). Red regions indicate theoretical predictions greater than 1 (i.e., bias amplification in the rightmost plot), while blue regions indicate theoretical predictions less than 1 (i.e., bias deamplification in the rightmost plot). Darkness indicates intensity. We consider isotropic covariance matrices: $\Sigma_1 = 2I_d, \Sigma_2 = I_d, \Theta = 2I_d, \Delta = I_d$. Additionally, $n = 1 \times 10^4, \sigma_1^2 = \sigma_2^2 = 1$. We further choose $\lambda = \lambda_1 = \lambda_2 = 1 \times 10^{-6}$ to approximate the minimum-norm interpolator. We show that bias amplification can occur even in the balanced data setting, i.e., when $p_1 = p_2 = 1/2$.

6.4.1 Isotropic Covariance

Setup. To mirror the setting of [BS23], we consider balanced data ($p_1 = p_2 = 1/2$) without spurious correlations ($\Sigma_1 = a_1 I_d, \Sigma_2 = a_2 I_d$, for $a_1, a_2 > 0$). The groups have different ground-truth weights ($\Theta = 2I_d, \Delta = I_d$). Refer to App. E.9.1 for full details due to space limitations.

Validation of Theory. Figure 6.2 and the figures in Appendix E.10 reveal that Theorems 6.3.1 and 6.3.2 closely predict the *ODD*, *EDD*, and *ADD* of ridge regression models with random projections under diverse settings. Note that, as indicated by the error bars, some of our empirical estimates (especially those with larger magnitude) have higher variance and their

variance is influenced by the choice of $\psi, \phi, a_1, a_2, \sigma_1^2, \sigma_2^2$. **Notably, our theory predicts the observation of [BS23] that models can amplify bias even with balanced groups and without spurious correlations.** We present new phenomena predicted by our theory below.

Effect of Label Noise. In the *ODD* profile, when the label noise ratio $c = \sigma_2^2/\sigma_1^2$ is larger, the right tail is higher for ϕ (rate of features to samples) closer to 1 than other ϕ . This suggests that under overparameterization, a larger noise ratio and similar number of features and samples can increase disparities in test risk between groups when a single model is learned for both groups. We aim to explain this phenomenon analytically in Section E.11. Moreover, the *EDD* curve is generally higher for larger c , suggesting that a larger noise ratio increases disparities in test risk when a separate model is learned for each group. This finding is supported by our experiment with real data (see Figure E.7).

Effect of Model Size. We observe interesting divergent behavior as ψ (rate of parameters to samples) increases for different ϕ (rate of features to samples). When $\phi > 1$, as ψ increases, the *ODD* increases and then decreases, peaking at the interpolation threshold at $\psi = 1$. Similarly, when $\phi > 0.5$ (i.e., $\phi_1 = \phi_2 > 1$), as ψ increases, the *EDD* increases and then decreases, peaking at the interpolation threshold at $\psi = 0.5$ (i.e., $\psi_1 = \psi_2 = 1$). Accordingly, when $\phi > 0.5$, bias is effectively deamplified ($ADD < 1$) at $\psi = 0.5$ and when $\phi > 1$, bias amplification peaks ($ADD > 1$) at $\psi = 1$. In contrast, when $\phi < 1$, the *ODD* decreases as ψ increases, plateauing at different finite values. Similarly, when $\phi < 0.5$, the *EDD* generally decreases and plateaus as ψ increases; in some cases, the *EDD* dips and/or increases and plateaus. A notable exception to these trends occurs when $\phi \approx 1$, with the corresponding *ODD* and *ADD* curves increasing as ψ increases, plateauing at a significantly larger value (i.e., $ADD \gg 1$) than the curves corresponding to other values of ϕ . We observe a similar phenomenon for the *EDD* curves when $\phi_1 = \phi_2 \approx 1$. Hence, overparameterization can greatly amplify bias when the number of features is close to the number of samples. Regardless of the

regime of ϕ , the left tail of the *ADD* profile appears to plateau at 1. The right tail plateaus at different finite values, with the curves corresponding to $\phi > 1$ consistently plateauing above 1. This suggests that when there are more features than samples, overparameterization amplifies bias.

Some of the peaks and valleys in Figure 6.2 can be attributed to double descent. However, double descent in high dimensions has primarily been studied in the setting where data are drawn from a single Gaussian distribution; this corresponds to the *EDD* setting, where a separate model is learned for each group. In Figure 6.1, we observe a double descent peak in the *EDD* at $\psi_1 = \psi_2 = 1$ [AP20b, DRB20]. Our work extends the theoretical treatment of double descent to the setting of training a model on a mixture of Gaussians. However, our theory of bias amplification cannot be reduced to double descent. For example, we note other interpolation thresholds in Figure 6.1; our use of a linear activation does not have a confounding effect here, as interpolation thresholds have also been observed in random features models with nonlinear activations [AP20a]. In addition, much of Sections 6.4 and 6.5, and Appendix E.11, are devoted to studying the tails or limiting behavior of bias amplification with respect to ψ and ϕ .

Effect of Number of Features. In the *ODD* and *ADD* profiles, when the rate of features to samples $\phi > 1$, the right tail generally plateaus at higher values (i.e., greater than 1) when ϕ is closer to 1. This suggests that with a similar number of features and samples, under overparameterization, bias amplification increases. In contrast, when $\phi < 1$, the right tail of the *ODD* and *EDD* curves seems to plateau at higher values when ϕ is larger. Regardless of the regime of the rate of features to samples ϕ , the left tails of the *ODD* and *EDD* curves are generally higher for larger ϕ .

6.4.2 Regularization and Training Dynamics

We now explore how regularization and training dynamics affect bias amplification.

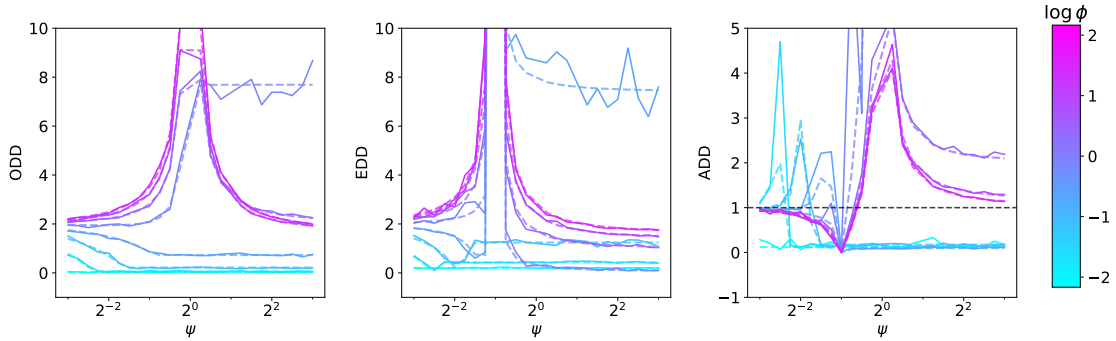


Figure 6.2: **Our theory predicts that models can amplify bias even with balanced groups and without spurious correlations.** We empirically validate our theory (Theorems 6.3.1 and 6.3.2) for ODD , EDD , and ADD under the setup described in Section 6.4.1, with $a_1 = 0.5$, $a_2 = 1$, $\sigma_1^2 = 1$, and $\sigma_2^2 = 1 \times 10^{-5}$. The solid lines capture empirical values while the corresponding lower-opacity dashed lines represent what our theory predicts. We plot ODD and EDD on the same scale for easy comparison, and include a black dashed line at $ADD = 1$ to contrast bias amplification vs. deamplification. We include all the plots with error bars in Appendix E.10.

Setup. We revisit the experimental setup for Section 6.4.1. We modulate a_1, a_2, ψ (rate of parameters to samples), as well as λ (regularization penalty) to understand the effects of regularization and early stopping on bias amplification. We fix $\sigma_1^2 = \sigma_2^2 = 1$, and the rate of features to samples $\phi = 0.75$.

Effect of Regularization and Training Time. In simplistic settings, we can simulate model learning over training time t by setting $\lambda = 1/t$ [AKT19]. In the figures in Appendix E.13, we observe that regardless of the regime of ψ , $ADD \approx 1$ (i.e., there is neither bias amplification nor deamplification) with high regularization or a short training time. When $\psi > 1$ (i.e., in the overparameterized regime), the ADD is generally greater than 1 across values of λ (i.e., bias is amplified), while when $\psi < 1$ (i.e., in the underparameterized regime), the ADD is generally less than 1 (i.e., bias is deamplified). Moreover, when $\psi > 1$, as regularization decreases (or training time increases), bias amplification increases and plateaus.

In contrast, when $\psi < 1$, as regularization decreases (or training time increases), bias deamplification increases and plateaus. A notable exception to this trend occurs when ψ is close to 1, where bias is initially deamplified and then amplified as λ decreases (or t increases). **This suggests that there may be an optimal regularization penalty or training time to avoid bias amplification and increase bias deamplification.** Intuitively, as training progresses, overparameterized models may discover “shortcut” associations [GJM20] that do not generalize equally well across groups, yielding bias amplification. In practice, an optimal λ or t can be selected by searching for values that strike a desired balance between overall validation error and empirical bias amplification. The search space can be reduced by using the above *ADD* trends w.r.t. λ and t that our theory predicts for over- vs. underparameterized models (see Appendix E.16 for more details). It is important for ML practitioners to consider the interplay between high vs. low feature-to-sample regimes and overparameterization in inducing bias amplification vs. deamplification when selecting optimal hyperparameters (see Figure 6.1).

In general, the calibration $\lambda = 1/t$ may not yield a theoretically tight picture of how bias evolves with t . The use of discrete gradient descent in practice rather than continuous-time gradient flows might yield further discrepancies. However, the calibration $\lambda = 1/t$ yields a ratio of gradient flow to ridge risk that is at most 1.69, with no assumptions on the features X [AKT19]. Moreover, in the controlled settings considered by [AKT19], this ratio empirically appears to be quite close to 1, and thus may be sufficient for extrapolating our results. Like us, [JNB24, HvG22] find that bias and bias amplification can vary substantially during training; future work can establish stronger connections between our observations and the results of [JNB24], which analytically identifies phases in the evolution of bias and a crossing phenomenon in the test error curves of groups during training. However, [JNB24] does not consider the effect of over- and underparameterization on bias evolution. While our analysis relies on the simplistic calibration $\lambda = 1/t$, it reveals divergent behavior in how bias evolves depending on model size.

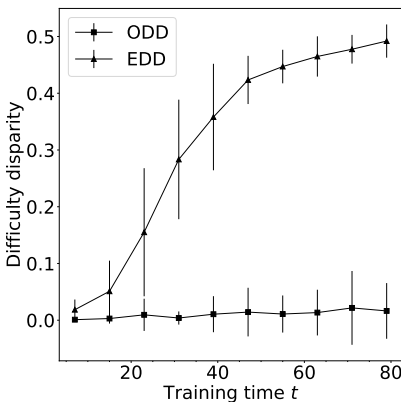


Figure 6.3: **Our theory predicts that disparate label noise between groups deamplifies bias on Colored MNIST.** We plot the ODD and EDD of a CNN over training time t for Colored MNIST. As t increases, the ODD is relatively low while the EDD is noticeably higher. The error bars capture the standard deviation computed over 10 random seeds.

Corroboration on Real Data. We further investigate the effect of training time on bias amplification on a more realistic dataset. We train a convolutional neural network (CNN) on Colored MNIST (see Appendix E.9.2 for more details). Colored MNIST is a semi-synthetic dataset derived from MNIST where digits are randomly re-colored to be red or green [ABG19]. We treat the color of each digit as its group, and we manipulate the groups to have different levels of label noise. In our experimental protocol: (1) the color of each digit (in both train and test) is chosen uniformly at random (i.e., with probability 0.5) and independently of the label; (2) by default, in the training set, the labels of red digits are flipped with probability 0.05 while the labels of green digits are flipped with probability 0.25; (3) labels are binarized (i.e., digits 0-4 correspond to 0 while digits 5-9 correspond to 1); and (4) each training step constitutes a step of gradient descent based on a batch of 250 instances. Although Colored MNIST is a classification task and we use a complex CNN architecture, **our theory correctly predicts that as the training time t increases, the ODD of the CNN is relatively low while the EDD is much larger**, producing bias deamplification.

Taking $t \rightarrow \infty$ corresponds to the setting of $\lambda \rightarrow 0^+$ in our theory (Theorems 6.3.1, 6.3.2). Because we assign the colors at random, the only difference in image features between groups is color; therefore, we expect the covariance matrices Σ_1 and Σ_2 to roughly coincide and $\Delta = 0$ (i.e., $w_1^* = w_2^*$). Note that we do not make any assumptions about the structure of Σ_1, Σ_2 . Furthermore, $p_1 = p_2 = 1/2$, and thus, $\phi_1 = \phi_2$ and $\psi_1 = \psi_2$. Additionally, we analogize the probability of label flipping to label noise in ridge regression. Hence, $e_1 = e_2, u_1 = u_2$. Accordingly, $\lim_{\lambda \rightarrow 0^+} B_1(\hat{f}) = \lim_{\lambda \rightarrow 0^+} B_1(\hat{f}_1) \approx \lim_{\lambda \rightarrow 0^+} B_2(\hat{f}) = \lim_{\lambda \rightarrow 0^+} B_2(\hat{f}_2)$. Simultaneously, $\lim_{\lambda \rightarrow 0^+} V_1(\hat{f}) \approx \lim_{\lambda \rightarrow 0^+} V_2(\hat{f})$. However, $\lim_{\lambda \rightarrow 0^+} V_1(\hat{f}_1) \approx \sigma_1^2/2 \cdot V = 0.05/2 \cdot V = 0.025V$ (where $V = \phi_1 h_1^{(2)}(I_d, \Sigma)$), while $\lim_{\lambda \rightarrow 0^+} V_2(\hat{f}_2) \approx \sigma_2^2/2 \cdot V = 0.25/2 \cdot V = 0.125V$. This results in $ODD \approx 0$ while $EDD \approx 0.1|V|$, which explains the divergence of ODD and EDD in Figure 6.3. Intuitively, the high label noise for the green digits prohibits the separate model \hat{f}_2 from achieving a low test risk compared to \hat{f}_1 ; the single model \hat{f} achieves a comparable test risk on both groups, effectively deamplifying bias, because of the better learning signal from the red digits. This phenomenon is similar to *positive transfer*, wherein the EDD of a model generally tends to be higher than the ODD when the labeling rules of imbalanced groups are sufficiently similar [MGR24]. However, [MGR24] does not explore the impact of model size on positive transfer. We show that the ODD can be less than the EDD depending on ψ in Figure 6.2, where $\Delta = I_d$ (i.e., the groups have different labeling rules). Future work can study the $ADD = \frac{ODD}{EDD}$ profile when $\Delta = 0$. Refer to Appendix E.14 for additional Colored MNIST experiments.

6.5 Minority-Group Bias

Recent work has revealed that overparameterization may hurt test performance on minority groups due to spurious features [SRK20, KL21]. Our theory provides new insights into how model size and extraneous features affect minority-group bias.

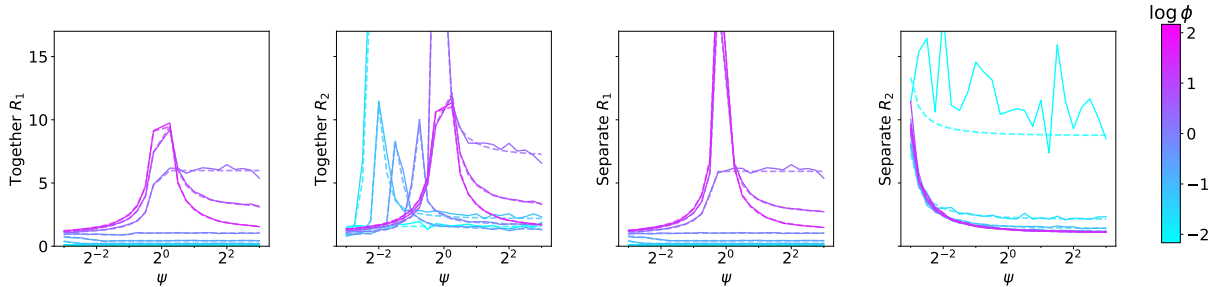


Figure 6.4: **Minority-group test risk can peak with different model sizes depending on the rate of features to samples.** We empirically demonstrate that minority-group bias is affected by extraneous features. We validate our theory (Theorems 6.3.1 and 6.3.2) for together R_1, R_2 (i.e., single model learned for both groups) and separate R_1, R_2 (i.e., separate model learned per group) under the setup described in Section 6.4.2, with $a_1 = 2, b_2 = 0.2$, and $\pi = 0.5$. The solid lines capture empirical values while the corresponding lower-opacity dashed lines represent what our theory predicts. We include a black dashed line at $ADD = 1$ to contrast bias amplification vs. deamplification. All y-axes are on the same scale for easy comparison. All the plots with error bars are in Appendix E.15.

Setup. To mirror the settings of [SRK20, KL21], we consider diatomic covariance matrices of *core* and *extraneous* features. We define $A \oplus B = \begin{pmatrix} A & 0 \\ 0 & B \end{pmatrix}$, and choose $\Sigma_1 = a_1 I_{\pi d} \oplus 0 I_{(1-\pi)d}$, $\Sigma_2 = a_2 I_{\pi d} \oplus b_2 I_{(1-\pi)d}$, for $\pi \in (0, 1)$, $a_1, b_2 > 0$, $a_1 = a_2$. Refer to App. E.9.1 for full details and a discussion of extraneous vs. spurious features (due to space limitations).

Interpolation Thresholds. The together R_2 (i.e., the test risk for the minority group in the single model setting) has different interpolation thresholds as ψ (rate of parameters to samples) increases, depending on ϕ (rate of features to samples) and π (fraction of core features). Notably, as ϕ increases, the interpolation thresholds occur at larger model sizes, culminating at $\psi = 1$. This suggests that for a higher rate of features to samples, a larger model size can greatly increase the together test risk of the minority group. Furthermore, the interpolation thresholds all occur closer to $\psi = 1$ for larger π , collapsing to a single threshold

at $\psi = 1$ when $\pi \rightarrow 1$ (as in Appendix E.10). Therefore, a lower fraction of core features can yield more possible model sizes that increase the test risk of the minority group. In addition, the together R_2 exhibits a steeper rate of growth around the interpolation thresholds for larger b_2 , suggesting that a higher variance in the extraneous features can also increase the test risk of the minority group in the single model setting. The phenomenon of different interpolation thresholds is not visible for R_2 when a separate model is trained per group; however, we do observe the expected double descent peaks in the separate R_1 and R_2 curves at $\psi_1 = 1$ and $\psi_2 = 1$, respectively.

Overparameterization. The right tails of the together R_2 curves plateau at different finite values depending on ϕ . In particular, for ϕ closer to 1, the together R_2 curves generally plateau at a higher value, suggesting that a similar number of features and samples can exacerbate minority-group bias under overparameterization. Furthermore, for smaller π and certain values of $\phi < 1$, the right tail of the together R_1 curve plateaus at a lower value than the together R_2 curve. This suggests that there can be differences in test error between groups that are not alleviated even with increased model size. This phenomenon diminishes in magnitude as the fraction of core features increases. This phenomenon supports the finding of [SRK20] that **overparameterization with spurious features can increase test risk disparities between groups**. We identify that the magnitude of this phenomenon may **depend on both the rate of features to samples and fraction of core features**.

6.6 Conclusion

We present a unifying, rigorous, and effective theory of ML bias in the settings of ridge regression with and without random projections. Our theory predicts interesting insights into bias amplification and minority-group bias in different feature and parameter regimes. These findings can inform strategies to evaluate and mitigate unfairness in ML (see Appendix

E.16 for more details). However, there remain practical challenges to assessing whether a model is prone to bias amplification. These include robustly estimating the feature covariance matrices [BL08] and label noises [FK14] for groups from sample data, especially for minority groups which have limited data. Even so, practitioners can use our theory and empirical observations to form intuition about when *disparities* in the variability of features and labels across groups can amplify bias.

In future work, we can extend our theory to the case of more than two groups and to accommodate label noise sampled from other distributions. Our theory can also be extended to different proportionate scaling limits (e.g., d^2/n has a finite limit instead of d/n) and to handle missing features [FCW24] and unknown group information [CRW19]. We can further leverage “Gaussian equivalents” [GLR22] to extend our theory to wide, fully-trained networks in the NTK [JGH18] and lazy [COB19] regimes; this will enable us to understand how, apart from model size, other design choices like nonlinear activation functions and learning rate may affect bias amplification.

6.7 Broader Impacts

This chapter explores how ML models amplify bias in a theoretically tractable setting. As such, our work does not capture the full range of real-world data and modeling complexity that gives rise to bias; hence, our theory can provide intuition but not robust predictions about when modern ML models (e.g., LLMs) may amplify bias. Nevertheless, this chapter reveals that we have yet to build a deep theoretical understanding of bias in even simple settings, which is a necessary foundation to obtain theoretical insights into more complex models. In addition, our theory of bias is limited to the technical realm. Therefore, it is critical to additionally consider: (1) how social inequality and power dynamics give rise to model and data biases, and (2) the significance of performance disparities in an application context, and whether reducing the disparities would promote justice.

CHAPTER 7

Conclusion and Future Directions

The unfairness of modern ML models, in conjunction with their growing adoption in everyday technological products and to make predictions about people at scale, risks worsening social inequality and reinforcing power dynamics. Towards combating these issues, this dissertation has explored the unfairness of GNNs and LLMs, whose applications range from social recommendation to chat-based assistants. GNNs and LLMs are both *contextual* models: GNNs commonly operate on social context, while LLMs are intended to process syntactic and semantic context. We tackle technical challenges in addressing their unfairness that transcend traditional ML: node data are not IID or processed independently by GNNs, and it is difficult to evaluate open-ended LLM generations for unfairness. We provide rich theoretical and empirical characterizations of the unfairness of GNNs, LLMs, and ML more broadly.

In the first part of this dissertation, we theoretically and empirically investigate forms of GNN unfairness and how they are affected by graph structure and the choice of graph filter. In Chapter 2, we show that graph-based feature imputation can amplify the unfairness of GNNs applied to imputed data. In Chapter 3, we show that GNNs can have a preferential attachment bias, disproportionately predicting links with high-degree nodes; this can cause disparities in link prediction scores, and thus social recommendation, between groups in social networks. In Chapter 4, we show why high-degree nodes tend to have a lower probability of misclassification by GNNs, and that during training, GNNs may adjust their loss on low-degree nodes more slowly. We propose principled metrics and methods to alleviate GNN unfairness.

In the second part of this dissertation, we assess the measurement validity of evaluations of LLM misgendering. In particular, we systematically show that probability-based, generation-based, and human evaluations of LLM misgendering can produce inconsistent results (Chapter 5). In the final part of this dissertation, we take a step back from LLMs and GNNs and tackle the major challenges in obtaining a precise unifying theory of ML unfairness in the relatively simple setting of a simplified feedforward neural network. We develop an effective theory of how model design choices and data properties contribute to model unfairness, providing insights into phenomena such as bias amplification and minority-group bias (Chapter 6). Such a theory can aid in the interpretation of unfair model predictions and design of stronger evaluation and mitigation methods for unfairness.

7.1 Future Directions

In future work, we will study how the different stages of the ML development lifecycle impact the unfairness of GNNs and LLMs (see Figure 7.1). Moreover, we will develop a deeper understanding of the mechanistic connections between GNNs and LLMs, towards understanding how their unfair behaviors relate to each other; this can enable the symbiotic development of unfairness mitigation methods for GNNs and LLMs.

Graph Neural Networks. Future work will address fairness issues caused by scaling GNNs to large networks. For example, we will study how randomness when training GNNs on large networks can cause variance in evaluations of model fairness. Randomness can engender inconsistent predictions for certain individuals as models are retrained over time (see Figure 7.2) [WPD24]. Training nodes are often selected at random, and their positions in the network heavily influence for which nodes GNNs perform better at test time [HLS23], e.g., test nodes in the local *context* of nodes in the training set may enjoy better performance. This issue is amplified for large networks, for which the number of available labeled nodes may

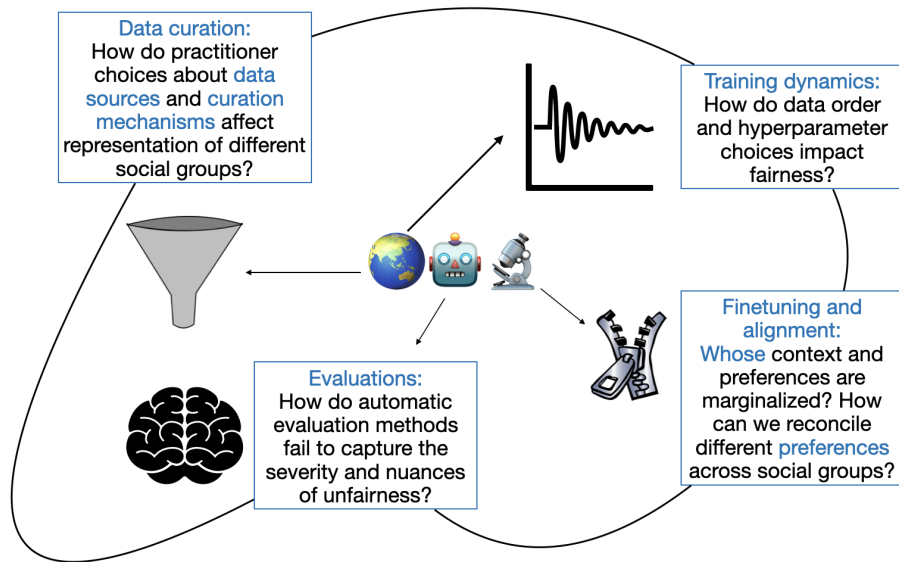


Figure 7.1: All the stages of the ML development lifecycle have an interdependent impact on model fairness.

be small compared to the size of the network. Moreover, for large networks, GNN training must occur in batches, which are often created using neighborhood (i.e., *context*-based) sampling [HYL17]. Hence, certain nodes may suffer from higher prediction variance due to the randomness of the training data order. Thus, it is essential to understand how randomness impacts the robustness of measurements of model fairness (see *Evaluations* and *Training dynamics* in Figure 7.1).

Beyond randomness, GNNs often cannot be deployed at scale because feature aggregation is time-consuming. One line of work has proposed simplifying GNNs by decoupling feature aggregation from feature transformation [WSZ19, FRE20]. Another line of work has investigated distilling knowledge from GNNs to MLPs to effectively attain graphless neural networks (GLNNs) [ZZH20]. While simplified and distilled GNNs have been shown to achieve comparable overall accuracy, we will investigate how simplifying and distilling exacerbate group unfairness, under different graph structure (i.e., *contextual*) properties such as class and group homophily. This pertains to *Training dynamics* and *Data curation* in Figure 7.1.

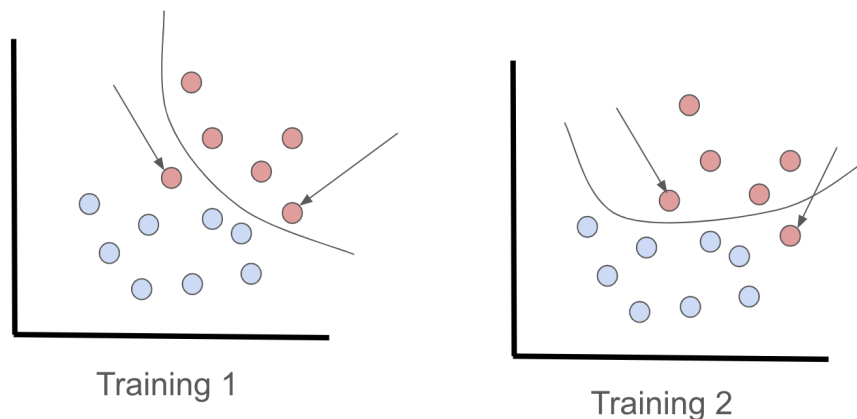


Figure 7.2: In consecutive rounds of training, changes in the decision boundary of a GNN can cause certain test instances (indicated by arrows) to receive unstable predictions. Each point in the figure represents a test instance and its color corresponds to its class.

Large Language Models. Future work will also address fairness issues caused by practitioner choices throughout the LLM development lifecycle. Practitioners pretrain LLMs on large amounts of language data that are often scraped from the web and curate these data using automatic filters. However, blacklist filters used by practitioners have been shown to disproportionately reject documents containing African American English (AAE) and references to LGBTQIA+ identities [DSM21]. Moreover, model-based filters used by practitioners to assess document quality can disproportionately exclude content from non-English speaking parts of the world [LGS24]. This motivates further research into how practitioner choices about data sources and curation mechanisms affect the representational fairness of LLM pretraining data for different social groups (see *Data curation* in Figure 7.1).

In addition, practitioner choices such as longer pretraining can yield an increase in gender bias [WLM25]. In this dissertation, we provide a theory (in a tractable setting) of how training time affects model bias. Future research can theoretically and empirically investigate how practitioner choices around pretraining data order and hyperparameters like weight decay impact LLM fairness (see *Training dynamics* in Figure 7.1). Following pretraining,

practitioners may align models to generate helpful and harmless responses to user queries. However, [OPM24] has revealed that alignment may not mitigate gender non-affirming language, and [MAC25] has shown that reward models capture human preferences more poorly in the context of AAE. These findings inspire further work on *whose* social context and preferences are marginalized by alignment algorithms and how practitioners can fairly reconcile different alignment preferences (see *Finetuning and alignment* in Figure 7.1).

Before deploying LLMs, practitioners must evaluate them for unfairness. In this dissertation, we study how choices about evaluation methodology impact measurements of LLM misgendering, and how automatic evaluations alone do not adequately measure misgendering. Future work can study how automatic evaluation tools used by practitioners (e.g., LLM-as-a-judge) fail to capture the severity and nuances of LLM unfairness (see *Evaluations* in Figure 7.1). Importantly, data curation, pretraining, alignment, and evaluation are interdependent, and the effects of practitioner choices in any of these stages may cascade to other stages. For example, a disproportionate presence of toxic language about LGBTQIA+ people in LLM pretraining data may cause LLMs to regurgitate toxic language about queer people; if automatic evaluations fail to catch such toxicity and these model generations enter our information ecosystem, they may become part of future pretraining data, thereby amplifying LLM toxicity.

Connecting GNNs and LLMs. GNNs and Transformers can be unified through the lens of geometric deep learning [BBC21]. Transformers can be viewed mechanistically as a type of GNN [Jos20], with the attention matrices acting as filters over a complete graph where the nodes are tokens and edges are attention. Through this lens, tokens that receive less attention (e.g., underrepresented tokens like neopronouns) can be thought of as low-degree nodes. Further elucidating the mechanistic connections between GNNs and LLMs can open new avenues for GNN and LLM fairness research to support one other.

7.2 Social Dimensions of Fairness

This dissertation primarily focuses on the technical dimensions of ML fairness. During its development, we also conducted research on the social and systemic dimensions of fairness. Creating fair ML models requires the participation of marginalized communities during all the stages of the ML lifecycle. Hence, we investigated methods for community-led participatory design and model auditing in the context of LGBTQIA+ communities [QOS23, QDO23]. Additionally, developing fair ML models necessitates an intersectional lens that centers combating social inequality and advancing justice. Thus, we explored the extent to which ML fairness papers are grounded in tenets of Intersectionality, identifying that researchers need to further consider the impact of their social context and power on their work [OSG23].

Collectively, this dissertation has developed a principled understanding of and addressed the contextual unfairness of modern ML models. In conjunction with further grounding in the social and systemic dimensions of fairness, this dissertation seeks to prevent the entrenchment of social inequalities by machine learning and promote justice.

Bibliography

- [AB02] Réka Albert and Albert-László Barabási. “Statistical mechanics of complex networks.” *Reviews of Modern Physics*, **74**(1):47–97, January 2002.
- [ABG19] Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. “Invariant risk minimization.” *arXiv preprint arXiv:1907.02893*, 2019.
- [Ado24] Adobe. “Neural Filters overview.”, 2024. <https://helpx.adobe.com/photoshop/using/neural-filters.html>.
- [AI19] Uber AI. “Food Discovery with Uber Eats: Using Graph Learning to Power Recommendations.”, 2019. <https://uber.com/blog/uber-eats-graph-learning/>.
- [Ai225] Ai2. “OLMoE, meet iOS.”, 2025. <https://openai.com/blog/olmoe-app>.
- [Aju19] Ifeoma Ajunwa. “An Auditing Imperative for Automated Hiring.” *DecisionSciRN: Recruiting & Hiring (Sub-Topic)*, 2019.
- [AKP22] Afra Feyza Akyürek, Muhammed Yusuf Kocyigit, Sejin Paik, and Derry Tanti Wijaya. “Challenges in Measuring Bias via Open-Ended Language Generation.” In Christian Hardmeier, Christine Basta, Marta R. Costa-jussà, Gabriel Stanovsky, and Hila Gonen, editors, *Proceedings of the 4th Workshop on Gender Bias in Natural Language Processing (GeBNLP)*, pp. 76–76, Seattle, Washington, jul 2022. Association for Computational Linguistics.
- [AKT19] Alnur Ali, J Zico Kolter, and Ryan J Tibshirani. “A continuous-time view of early stopping for least squares regression.” In *The 22nd international conference on artificial intelligence and statistics*, pp. 1370–1378. PMLR, 2019.
- [ALZ21] Chirag Agarwal, Himabindu Lakkaraju, and Marinka Zitnik. “Towards a unified framework for fair and stable graph representation learning.” In Cassio de Campos

- and Marloes H. Maathuis, editors, *Proceedings of the Thirty-Seventh Conference on Uncertainty in Artificial Intelligence*, volume 161 of *Proceedings of Machine Learning Research*, pp. 2114–2124. PMLR, 27–30 Jul 2021.
- [Ant] Anthropic. “Claude.” <https://claude.ai/>.
- [AP20a] Ben Adlam and Jeffrey Pennington. “The Neural Tangent Kernel in High Dimensions: Triple Descent and a Multi-Scale Theory of Generalization.” In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 74–84. PMLR, 13–18 Jul 2020.
- [AP20b] Ben Adlam and Jeffrey Pennington. “Understanding Double Descent Requires A Fine-Grained Bias-Variance Decomposition.” In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pp. 11022–11032. Curran Associates, Inc., 2020.
- [AR23] Haozhe An and Rachel Rudinger. “Nichelle and Nancy: The Influence of Demographic Attributes and Tokenization Length on First Name Biases.” In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pp. 388–401, Toronto, Canada, jul 2023. Association for Computational Linguistics.
- [ASA15] Mason Ameri, Lisa A. Schur, Meera Adya, F. Scott Bentley, Patrick F. McKay, and Douglas L. Kruse. “The Disability Employment Puzzle: A Field Experiment on Employer Hiring Behavior.” *ILR Review*, **71**:329–364, 2015.
- [Atk70] Anthony B Atkinson. “On the measurement of inequality.” *Journal of Economic Theory*, **2**(3):244–263, 1970.

- [BA99] Albert-László Barabási and Réka Albert. “Emergence of Scaling in Random Networks.” *Science*, **286**(5439):509–512, 1999.
- [Bac23] Francis Bach. “High-dimensional analysis of double descent for linear regression with random projections.” *arXiv preprint arXiv:2303.01372*, 2023.
- [Bar10] Sharon N. Barnartt. *Disability as a fluid state: Introduction*, p. 1–22. Emerald Group Publishing Limited, January 2010.
- [Bar16] Albert-László Barabási. “Network science.” *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, **371**, 2016.
- [BBC21] Michael M. Bronstein, Joan Bruna, Taco Cohen, and Petar Veličković. “Geometric Deep Learning: Grids, Groups, Graphs, Geodesics, and Gauges.” *arXiv preprint arXiv:2104.13478*, 2021.
- [BCZ16] Tolga Bolukbasi, Kai-Wei Chang, James Zou, Venkatesh Saligrama, and Adam Kalai. “Man is to computer programmer as woman is to homemaker? debiasing word embeddings.” In *Proceedings of the 30th International Conference on Neural Information Processing Systems, NIPS’16*, p. 4356–4364, Red Hook, NY, USA, 2016. Curran Associates Inc.
- [BD20] Maarten Buyl and Tijl De Bie. “DeBayes: a Bayesian Method for Debiasing Network Embeddings.” In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 1220–1229. PMLR, 13–18 Jul 2020.
- [Bec21] Amir Beck. “Proximal-Based Methods Tutorial.”, 2021.
- [BFS23] Ashkan Bashardoust, Sorelle Friedler, Carlos Scheidegger, Blair D. Sullivan, and Suresh Venkatasubramanian. “Reducing Access Disparities in Networks using Edge Augmentation.” In *Proceedings of the 2023 ACM Conference on Fairness*,

Accountability, and Transparency, FAccT '23, p. 1635–1651, New York, NY, USA, 2023. Association for Computing Machinery.

- [BG18a] Aleksandar Bojchevski and Stephan Günnemann. “Deep Gaussian Embedding of Graphs: Unsupervised Inductive Learning via Ranking.” In *International Conference on Learning Representations*, 2018.
- [BG18b] Joy Buolamwini and Timnit Gebru. “Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification.” In Sorelle A. Friedler and Christo Wilson, editors, *Proceedings of the 1st Conference on Fairness, Accountability and Transparency*, volume 81 of *Proceedings of Machine Learning Research*, pp. 77–91. PMLR, 23–24 Feb 2018.
- [BH19] Avishek Bose and William Hamilton. “Compositional Fairness Constraints for Graph Embeddings.” In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 715–724. PMLR, 09–15 Jun 2019.
- [BHM19] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. “Reconciling modern machine-learning practice and the classical bias–variance trade-off.” *Proceedings of the National Academy of Sciences*, **116**(32):15849–15854, 2019.
- [BHN19] Solon Barocas, Moritz Hardt, and Arvind Narayanan. *Fairness and Machine Learning*. fairmlbook.org, 2019. <http://www.fairmlbook.org>.
- [BKD23] Federico Bianchi, Pratyusha Kalluri, Esin Durmus, Faisal Ladhak, Myra Cheng, Debora Nozza, Tatsunori Hashimoto, Dan Jurafsky, James Zou, and Aylin Caliskan. “Easily Accessible Text-to-Image Generation Amplifies Demographic Stereotypes at Large Scale.” In *Proceedings of the 2023 ACM Conference on Fairness, Ac-*

- accountability, and Transparency*, FAccT '23, p. 1493–1504, New York, NY, USA, 2023. Association for Computing Machinery.
- [BKE22] Stephen Bonner, Ufuk Kirik, Ola Engkvist, Jian Tang, and Ian P Barrett. “Implications of topological imbalance for representation learning on biomedical knowledge graphs.” *Briefings in Bioinformatics*, **23**(5):bbac279, 07 2022.
- [BKW17] Rianne van den Berg, Thomas N. Kipf, and Max Welling. “Graph Convolutional Matrix Completion.” *arXiv preprint arXiv:1706.02263*, 2017.
- [BL08] Peter J. Bickel and Elizaveta Levina. “Regularized estimation of large covariance matrices.” *Annals of Statistics*, **36**:199–227, 2008.
- [BLM14] Danah Boyd, Karen Levy, and Alice Marwick. “The networked nature of algorithmic discrimination.” *Data and Discrimination: Collected Essays*. Open Technology Institute, 2014.
- [BLO21] Su Lin Blodgett, Gilsinia Lopez, Alexandra Olteanu, Robert Sim, and Hanna Wallach. “Stereotyping Norwegian Salmon: An Inventory of Pitfalls in Fairness Benchmark Datasets.” In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 1004–1015, Online, aug 2021. Association for Computational Linguistics.
- [BS16] Solon Barocas and Andrew D Selbst. “Big data’s disparate impact.” *Calif. L. Rev.*, **104**:671, 2016.
- [BS23] Samuel James Bell and Levent Sagun. “Simplicity Bias Leads to Amplified Performance Disparities.” In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '23, p. 355–369, New York, NY, USA, 2023. Association for Computing Machinery.

- [BSC14] Nicola Bertoldi, Patrick Simianer, Mauro Cettolo, Katharina Wäsche, Marcello Federico, and Stefan Riezler. “Online adaptation to post-edits for phrase-based statistical machine translation.” *Machine Translation*, **28**:309–339, 2014.
- [BW24] Samuel J. Bell and Skyler Wang. “The Multiple Dimensions of Spuriousness in Machine Learning.” *arXiv preprint arXiv:2411.04696*, 2024.
- [Can23] Clément L. Canonne. “A short note on an inequality between KL and TV.” *arXiv preprint arXiv:2202.07198*, 2023.
- [Cap19] From Analytics First to AI First at Capital One. “Tom Davenport.” *Forbes*, 2019. <https://www.forbes.com/sites/tomdavenport/2019/07/10/from-analytics-first-to-ai-first-at-capital-one/>.
- [CB20] Patricia Hill Collins and Sirma Bilge. *Intersectionality*. Key Concepts. Polity Press, 2020.
- [CC21] Shih-Hsuan Chiu and Berlin Chen. “Innovative Bert-Based Reranking Language Models for Speech Recognition.” In *2021 IEEE Spoken Language Technology Workshop (SLT)*, pp. 266–271, 2021.
- [CD21] Yang Trista Cao and Hal Daumé III. “Toward Gender-Inclusive Coreference Resolution: An Analysis of Gender and Bias Throughout the Machine Learning Lifecycle*.” *Computational Linguistics*, **47**(3):615–661, nov 2021.
- [CHG22] Sean Current, Yuntian He, Saket Gurukar, and Srinivasan Parthasarathy. “FairEGM: Fair Link Prediction and Recommendation via Emulated Graph Modification.” In *Proceedings of the 2nd ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization, EAAMO ’22*, New York, NY, USA, 2022. Association for Computing Machinery.

- [CJ20] Davide Chicco and Giuseppe Jurman. “The advantages of the Matthews correlation coefficient (MCC) over F1 score and accuracy in binary classification evaluation.” *BMC Genomics*, **21**, 2020.
- [CLK22] Hugo Cui, Bruno Loureiro, Florent Krzakala, and Lenka Zdeborová. “Generalization error rates in kernel regression: the crossover from the noiseless to noisy regime.” *Journal of Statistical Mechanics: Theory and Experiment*, **2022**(11):114004, nov 2022.
- [CLY23] Renyao Chen, Junye Lei, Hong Yao, Tailong Li, and Shengwen Li. “Anchor-Enhanced Geographical Entity Representation Learning.” *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–15, 2023.
- [CMT22] Joymallya Chakraborty, Suvodeep Majumder, and Huy Tu. “Fair-SSL: building fair ML software with less data.” In *Proceedings of the 2nd International Workshop on Equitable Data and Technology*, ICSE ’22, p. 1–8. ACM, May 2022.
- [COB19] Lenaïc Chizat, Edouard Oyallon, and Francis Bach. “On lazy training in differentiable programming.” *Advances in neural information processing systems*, **32**, 2019.
- [Com23] Google Ads & Commerce. “Introducing a new era of AI-powered ads with Google.”, 2023. <https://blog.google/products/ads-commerce/ai-powered-ads-google-marketing-live/>.
- [Con18] Kirby Conrod. “Pronouns and Gender in Language.” In *The Oxford Handbook of Language and Sexuality*. Oxford University Press, 2018.
- [CPC22] Yang Trista Cao, Yada Pruksachatkun, Kai-Wei Chang, Rahul Gupta, Varun Kumar, Jwala Dhamala, and Aram Galstyan. “On the Intrinsic and Extrinsic Fairness Evaluation Metrics for Contextualized Language Representations.” In

- Smaranda Muresan, Preslav Nakov, and Aline Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pp. 561–570, Dublin, Ireland, may 2022. Association for Computational Linguistics.
- [CPN23] Pierre Colombo, Maxime Peyrard, Nathan Noiry, Robert West, and Pablo Piantanida. “The Glass Ceiling of Automatic Evaluation in Natural Language Generation.” In Jong C. Park, Yuki Arase, Baotian Hu, Wei Lu, Derry Wijaya, Ayu Purwarianti, and Adila Alfa Krisnadhi, editors, *Findings of the Association for Computational Linguistics: IJCNLP-AAACL 2023 (Findings)*, pp. 178–183, Nusa Dua, Bali, nov 2023. Association for Computational Linguistics.
- [CR09] Emmanuel J. Candès and Benjamin Recht. “Exact Matrix Completion via Convex Optimization.” *Foundations of Computational Mathematics*, **9**:717–772, 2009.
- [CRW19] Amanda Coston, Karthikeyan Natesan Ramamurthy, Dennis Wei, Kush R. Varshney, Skyler Speakman, Zairah Mustahsan, and Supriyo Chakraborty. “Fair Transfer Learning with Missing Protected Attributes.” In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, AIES ’19, p. 91–98, New York, NY, USA, 2019. Association for Computing Machinery.
- [CSR23] Benjamin Paul Chamberlain, Sergey Shirobokov, Emanuele Rossi, Fabrizio Frasca, Thomas Markovich, Nils Yannick Hammerla, Michael M. Bronstein, and Max Hansmire. “Graph Neural Networks for Link Prediction with Subgraph Sketching.” In *The Eleventh International Conference on Learning Representations*, 2023.
- [CTM19] Dawei Cheng, Yi Tu, Zhenwei Ma, Zhibin Niu, and Liqing Zhang. “Risk Assessment for Networked-guarantee Loans Using High-order Graph Attention Representation.” In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*, pp. 5822–5828. International Joint Conferences on Artificial Intelligence Organization, 7 2019.

- [CV07] Andrea Caponnetto and Ernesto de Vito. “Optimal Rates for the Regularized Least-Squares Algorithm.” *Foundations of Computational Mathematics*, **7**:331–368, 2007.
- [CWC23] Jiajia Chen, Jiancan Wu, Jiawei Chen, Xin Xin, Yong Li, and Xiangnan He. “How graph convolutions amplify popularity bias for recommendation?” *Frontiers of Computer Science*, **18**(5), December 2023.
- [CWH20] Ming Chen, Zhewei Wei, Zengfeng Huang, Bolin Ding, and Yaliang Li. “Simple and Deep Graph Convolutional Networks.” In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 1725–1735. PMLR, 13–18 Jul 2020.
- [CXK22] Zhengyu Chen, Teng Xiao, and Kun Kuang. “BA-GNN: On Learning Bias-Aware Graph Neural Network.” In *2022 IEEE 38th International Conference on Data Engineering (ICDE)*, pp. 3012–3024, 2022.
- [Dee] DeepSeek. “Download DeepSeek App.” <https://download.deepseek.com/app/>.
- [DFK24] Elvis Dohmatob, Yunzhen Feng, and Julia Kempe. “Model Collapse Demystified: The Case of Regression.” In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang, editors, *Advances in Neural Information Processing Systems*, volume 37, pp. 46979–47013. Curran Associates, Inc., 2024.
- [DHP12] Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. “Fairness through Awareness.” In *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference*, ITCS ’12, p. 214–226, New York, NY, USA, 2012. Association for Computing Machinery.
- [DJ23] Claire Donnat and So Won Jeong. “Studying the Effect of GNN Spatial Convolutions On The Embedding Space’s Geometry.” In Robin J. Evans and Ilya

- Shpitser, editors, *Proceedings of the Thirty-Ninth Conference on Uncertainty in Artificial Intelligence*, volume 216 of *Proceedings of Machine Learning Research*, pp. 539–548. PMLR, 31 Jul–04 Aug 2023.
- [DL18] Peng Ding and Fan Li. “Causal Inference: A Missing Data Perspective.” *Statistical Science*, **33**(2):214 – 237, 2018.
- [DLJ22] Yushun Dong, Ninghao Liu, Brian Jalaian, and Jundong Li. “EDITS: Modeling and Mitigating Data Bias for Graph Neural Networks.” In *Proceedings of the ACM Web Conference 2022*, WWW ’22, p. 1259–1269, New York, NY, USA, 2022. Association for Computing Machinery.
- [DMO21] Sunipa Dev, Masoud Monajatipoor, Anaelia Ovalle, Arjun Subramonian, Jeff Phillips, and Kai-Wei Chang. “Harms of Gender Exclusivity and Challenges in Non-Binary Representation in Language Technologies.” In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 1968–1994, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics.
- [dMW19] Terrance de Vries, Ishan Misra, Changhan Wang, and Laurens van der Maaten. “Does Object Recognition Work for Everyone?” In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pp. 52–59, 2019.
- [DRB20] Stéphane D’Ascoli, Maria Refinetti, Giulio Biroli, and Florent Krzakala. “Double Trouble in Double Descent: Bias and Variance(s) in the Lazy Regime.” In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 2280–2290. PMLR, 13–18 Jul 2020.

- [DSL23] Wei Ding, Jiawei Sun, Jie Li, and Chentao Wu. “Inductive Dummy-based Homogeneous Neighborhood Augmentation for Graph Collaborative Filtering.” In *2023 International Joint Conference on Neural Networks (IJCNN)*, pp. 1–8, 2023.
- [DSM18] Yash Deshpande, Subhabrata Sen, Andrea Montanari, and Elchanan Mossel. “Contextual Stochastic Block Models.” In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- [DSM21] Jesse Dodge, Maarten Sap, Ana Marasović, William Agnew, Gabriel Ilharco, Dirk Groeneveld, Margaret Mitchell, and Matt Gardner. “Documenting Large Webtext Corpora: A Case Study on the Colossal Clean Crawled Corpus.” In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 1286–1305, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics.
- [DTC22] Pieter Delobelle, Ewoenam Tokpo, Toon Calders, and Bettina Berendt. “Measuring Fairness with Biased Rulers: A Comparative Study on Bias Metrics for Pre-trained Language Models.” In Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz, editors, *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 1693–1706, Seattle, United States, jul 2022. Association for Computational Linguistics.
- [Duc16] John C. Duchi. “Derivations for Linear Algebra and Optimization.” 2016.
- [DW18a] Robin Dembroff and Daniel Wodak. “He/She/They/Ze.” *Ergo: An Open Access Journal of Philosophy*, **5**, 2018.

- [DW18b] Edgar Dobriban and Stefan Wager. “High-dimensional asymptotics of prediction: Ridge regression and classification.” *The Annals of Statistics*, **46**(1):247–279, 2018.
- [DW21] Enyan Dai and Suhang Wang. “Say No to the Discrimination: Learning Fair Graph Neural Networks with Limited Sensitive Attribute Information.” In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining, WSDM ’21*, p. 680–688, New York, NY, USA, 2021. Association for Computing Machinery.
- [EMP22] Ahmed El-Kishky, Thomas Markovich, Serim Park, Chetan Verma, Baekjin Kim, Ramy Eskander, Yury Malkov, Frank Portman, Sofia Samaniego, Ying Xiao, and Aria Haghighi. “TwHIN: Embedding the Twitter Heterogeneous Information Network for Personalized Recommendation.” In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD ’22*, p. 2842–2850, New York, NY, USA, 2022. Association for Computing Machinery.
- [Eng22] Snap Engineering. “Machine Learning for Snapchat Ad Ranking.”, 2022. <https://eng.snap.com/machine-learning-snap-ad-ranking>.
- [EXK25] Aparna Elangovan, Lei Xu, Jongwoo Ko, Mahsa Elyasi, Ling Liu, Sravan Babu Bodapati, and Dan Roth. “Beyond correlation: The impact of human uncertainty in measuring the effectiveness of automatic evaluation and LLM-as-a-judge.” In *The Thirteenth International Conference on Learning Representations*, 2025.
- [FBB19] Benjamin Fish, Ashkan Bashardoust, Danah Boyd, Sorelle Friedler, Carlos Scheidegger, and Suresh Venkatasubramanian. “Gaps in Information Access in Social Networks?” In *The World Wide Web Conference, WWW ’19*, p. 480–490, New York, NY, USA, 2019. Association for Computing Machinery.
- [FCD21] Martínez-Plumed Fernando, Ferri Cèsar, Nieves David, and Hernández-Orallo José. “Missing the missing values: The ugly duckling of fairness in machine learning.” *International Journal of Intelligent Systems*, **36**(7):3217–3258, 2021.

- [FCW24] Raymond Feng, Flavio Calmon, and Hao Wang. “Adapting fairness interventions to missing values.” *Advances in Neural Information Processing Systems*, **36**, 2024.
- [FHH21] Fuli Feng, Weiran Huang, Xiangnan He, Xin Xin, Qifan Wang, and Tat-Seng Chua. “Should Graph Convolution Trust Neighbors? A Simple Causal Inference Method.” In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR ’21, p. 1208–1218, New York, NY, USA, 2021. Association for Computing Machinery.
- [FIK20] James R. Foulds, Rashidul Islam, Kamrun Naher Keya, and Shimei Pan. “An Intersectional Definition of Fairness.” In *2020 IEEE 36th International Conference on Data Engineering (ICDE)*, pp. 1918–1921, 2020.
- [FK14] Benoît Frénay and Ata Kabán. “A comprehensive introduction to label noise.” In *The European Symposium on Artificial Neural Networks*, 2014.
- [FL19] Matthias Fey and Jan E. Lenssen. “Fast Graph Representation Learning with PyTorch Geometric.” In *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019.
- [FML19] Wenqi Fan, Yao Ma, Qing Li, Yuan He, Eric Zhao, Jiliang Tang, and Dawei Yin. “Graph Neural Networks for Social Recommendation.” In *The World Wide Web Conference*, WWW ’19, p. 417–426, New York, NY, USA, 2019. Association for Computing Machinery.
- [FOB06] Reza Rashidi Far, Tamer Oraby, Wlodzimierz Bryc, and Roland Speicher. “Spectra of large block matrices.” *arXiv preprint arXiv:cs/0610045*, 2006.
- [FRE20] Fabrizio Frasca, Emanuele Rossi, Davide Eynard, Benjamin Chamberlain, Michael Bronstein, and Federico Monti. “SIGN: Scalable Inception Graph Neural Networks.” In *ICML 2020 Workshop on Graph Representation Learning and Beyond*, 2020.

- [Fri20] Christian Fricke. “*Missing Fairness: The Discriminatory Effect of Missing Values in Datasets on Fairness in Machine Learning.*”. Master’s thesis, Aalto University. School of Science, 2020.
- [GAD21] Naman Goel, Alfonso Amayuelas, Amit Deshpande, and Amit Sharma. “The Importance of Modeling Data Missingness in Algorithmic Fairness: A Causal Perspective.” *Proceedings of the AAAI Conference on Artificial Intelligence*, **35**(9):7564–7573, May 2021.
- [GBW24] Dirk Groeneveld, Iz Beltagy, Evan Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora, David Atkinson, Russell Authur, Khyathi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar, Yuling Gu, Jack Hessel, Tushar Khot, et al. “OLMo: Accelerating the Science of Language Models.” In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15789–15809, Bangkok, Thailand, aug 2024. Association for Computational Linguistics.
- [GBZ24] Vagrant Gautam, Eileen Bingert, Dawei Zhu, Anne Lauscher, and Dietrich Klakow. “Robust Pronoun Fidelity with English LLMs: Are they Reasoning, Repeating, or Just Biased?” *Transactions of the Association for Computational Linguistics*, **12**:1755–1779, 2024.
- [GDJ24] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, et al. “The Llama 3 Herd of Models.” *arXiv preprint arXiv:2407.21783*, 2024.
- [GGR21] Avijit Ghosh, Lea Genuit, and Mary Reagan. “Characterizing Intersectional Group Fairness with Worst-Case Comparisons.” In Deepti Lamba and William H.

- Hsu, editors, *Proceedings of 2nd Workshop on Diversity in Artificial Intelligence (AIDBEI)*, volume 142 of *Proceedings of Machine Learning Research*, pp. 22–34. PMLR, 09 Feb 2021.
- [GJ93] Zoubin Ghahramani and Michael Jordan. “Supervised learning from incomplete data via an EM approach.” In J. Cowan, G. Tesauro, and J. Alspector, editors, *Advances in Neural Information Processing Systems*, volume 6. Morgan-Kaufmann, 1993.
- [GJM20] Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A Wichmann. “Shortcut learning in deep neural networks.” *Nature Machine Intelligence*, **2**(11):665–673, 2020.
- [GLR22] Sebastian Goldt, Bruno Loureiro, Galen Reeves, Florent Krzakala, Marc Mézard, and Lenka Zdeborová. “The gaussian equivalence of generative models for learning with shallow neural networks.” In *Mathematical and Scientific Machine Learning*, pp. 426–471. PMLR, 2022.
- [GMM21] Seraphina Goldfarb-Tarrant, Rebecca Marchant, Ricardo Muñoz Sánchez, Mugdha Pandya, and Adam Lopez. “Intrinsic Bias Metrics Do Not Correlate with Application Bias.” In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 1926–1940, Online, aug 2021. Association for Computational Linguistics.
- [GRB24] Isabel O. Gallegos, Ryan A. Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K. Ahmed. “Bias and Fairness in Large Language Models: A Survey.” *Computational Linguistics*, **50**(3):1097–1179, sep 2024.

- [GRC23] Francesco Di Giovanni, James Rowbottom, Benjamin Paul Chamberlain, Thomas Markovich, and Michael M. Bronstein. “Understanding convolution on graphs via energies.” *Transactions on Machine Learning Research*, 2023.
- [GSL18] Yupeng Gu, Yizhou Sun, Yanen Li, and Yang Yang. “RaRE: Social Rank Regulated Large-scale Network Embedding.” In *Proceedings of the 2018 World Wide Web Conference*, WWW ’18, p. 359–368, Republic and Canton of Geneva, CHE, 2018. International World Wide Web Conferences Steering Committee.
- [GSL24] Vagrant Gautam, Arjun Subramonian, Anne Lauscher, and Os Keyes. “Stop! In the Name of Flaws: Disentangling Personal Names and Sociodemographic Attributes in NLP.” In Agnieszka Faleńska, Christine Basta, Marta Costa-jussà, Seraphina Goldfarb-Tarrant, and Debora Nozza, editors, *Proceedings of the 5th Workshop on Gender Bias in Natural Language Processing (GeBNLP)*, pp. 323–337, Bangkok, Thailand, aug 2024. Association for Computational Linguistics.
- [GUB23] Seraphina Goldfarb-Tarrant, Eddie Ungless, Esma Balkir, and Su Lin Blodgett. “This prompt is measuring <mask>: evaluating bias evaluation in language models.” In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 2209–2225, Toronto, Canada, jul 2023. Association for Computational Linguistics.
- [Ham21] Lelia Marie Hampton. “Black Feminist Musings on Algorithmic Oppression.” In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, FAccT ’21, p. 1, New York, NY, USA, 2021. Association for Computing Machinery.
- [HBC20] David M. Howcroft, Anya Belz, Miruna-Adriana Cliniciu, Dimitra Gkatzia, Sadid A. Hasan, Saad Mahamood, Simon Mille, Emiel van Miltenburg, Sashank Santhanam, and Verena Rieser. “Twenty Years of Confusion in Human Evaluation: NLG Needs

- Evaluation Sheets and Standardised Definitions.” In Brian Davis, Yvette Graham, John Kelleher, and Yaji Sripada, editors, *Proceedings of the 13th International Conference on Natural Language Generation*, pp. 169–182, Dublin, Ireland, dec 2020. Association for Computational Linguistics.
- [HBS18] Lisa Anne Hendricks, Kaylee Burns, Kate Saenko, Trevor Darrell, and Anna Rohrbach. “Women Also Snowboard: Overcoming Bias in Captioning Models.” In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 771–787, 2018.
- [HCT17] Bas Hofstra, Rense Corten, Frank van Tubergen, and Nicole B. Ellison. “Sources of Segregation in Social Networks: A Novel Approach Using Facebook.” *American Sociological Review*, **82**(3):625–656, 2017.
- [HDS23] Tamanna Hossain, Sunipa Dev, and Sameer Singh. “MISGENDERED: Limits of Large Language Models in Understanding Pronouns.” In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 5352–5367, Toronto, Canada, jul 2023. Association for Computational Linguistics.
- [HDS24] Tamanna Hossain, Sunipa Dev, and Sameer Singh. “MisgenderMender: A Community-Informed Approach to Interventions for Misgendering.” In Kevin Duh, Helena Gomez, and Steven Bethard, editors, *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 7538–7558, Mexico City, Mexico, jun 2024. Association for Computational Linguistics.
- [HGL18] Jason Hartford, Devon Graham, Kevin Leyton-Brown, and Siamak Ravanbakhsh. “Deep Models of Interactions Across Sets.” In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*,

- volume 80 of *Proceedings of Machine Learning Research*, pp. 1909–1918. PMLR, 10–15 Jul 2018.
- [HL23a] Van Thuy Hoang and O-Joun Lee. “Mitigating Degree Biases in Message Passing Mechanism by Utilizing Community Structures.” *arXiv preprint arXiv:2312.16788*, 2023.
- [HL23b] Jennifer Hu and Roger Levy. “Prompting is not a substitute for probability measurements in large language models.” In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 5040–5060, Singapore, dec 2023. Association for Computational Linguistics.
- [HLS23] Haoyu Han, Xiaorui Liu, Feng Shi, Mohamad Ali Torkamani, Charu C. Aggarwal, and Jiliang Tang. “Towards label position bias in graph neural networks.” In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, NIPS ’23, Red Hook, NY, USA, 2023. Curran Associates Inc.
- [HML15] Trevor Hastie, Rahul Mazumder, Jason D. Lee, and Reza Zadeh. “Matrix Completion and Low-Rank SVD via Fast Alternating Least Squares.” *J. Mach. Learn. Res.*, **16**(1):3367–3402, jan 2015.
- [HMR22] Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J Tibshirani. “Surprises in high-dimensional ridgeless least squares interpolation.” *Annals of statistics*, **50**(2):949, 2022.
- [HMV20] Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. “spaCy: Industrial-strength Natural Language Processing in Python.” 2020.
- [HPS16] Moritz Hardt, Eric Price, and Nathan Srebro. “Equality of opportunity in supervised learning.” In *Proceedings of the 30th International Conference on Neural*

- Information Processing Systems*, NIPS'16, p. 3323–3331, Red Hook, NY, USA, 2016. Curran Associates Inc.
- [HSB24] Emma Harvey, Emily Sheng, Su Lin Blodgett, Alexandra Chouldechova, Jean Garcia-Gathright, Alexandra Olteanu, and Hanna Wallach. “Gaps Between Research and Practice When Measuring Representational Harms Caused by LLM-Based Systems.” *arXiv preprint arXiv:2411.15662*, 2024.
- [Hui14] Mark Huisman. *Imputation of Missing Network Data: Some Simple Procedures*, pp. 707–715. Springer New York, New York, NY, 2014.
- [HvG22] Melissa Hall, Laurens van der Maaten, Laura Gustafson, Maxwell Jones, and Aaron Adcock. “A Systematic Study of Bias Amplification.” *arXiv preprint arXiv:2201.11706*, oct 2022.
- [HYL17] William L. Hamilton, Rex Ying, and Jure Leskovec. “Inductive Representation Learning on Large Graphs.” In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, p. 1025–1035, Red Hook, NY, USA, 2017. Curran Associates Inc.
- [HZW24] Xiaotian Han, Kaixiong Zhou, Ting-Hsiang Wang, Jundong Li, Fei Wang, and Na Zou. “Marginal Nodes Matter: Towards Structure Fairness in Graphs.” *SIGKDD Explor. Newsl.*, **25**(2):4–13, March 2024.
- [JGH18] Arthur Jacot, Franck Gabriel, and Clement Hongler. “Neural Tangent Kernel: Convergence and Generalization in Neural Networks.” In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- [JHL23] Bahram Jafrasteh, Daniel Hernández-Lobato, Simón Pedro Lubián-López, and

- Isabel Benavente-Fernández. “Gaussian processes for missing value imputation.” *Knowledge-Based Systems*, **273**:110603, 2023.
- [JNB24] Anchit Jain, Rozhin Nobahari, Aristide Baratin, and Stefano Sarao Mannelli. “Bias in Motion: Theoretical Insights into the Dynamics of Bias in SGD Training.” In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang, editors, *Advances in Neural Information Processing Systems*, volume 37, pp. 24435–24471. Curran Associates, Inc., 2024.
- [Jos20] Chaitanya Joshi. “Transformers are Graph Neural Networks.” *The Gradient*, 2020.
- [JSR24] Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, L lio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, et al. “Mixtral of Experts.” *arXiv preprint arXiv:2401.04088*, 2024.
- [JW21] Abigail Z. Jacobs and Hanna Wallach. “Measurement and Fairness.” In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, FAccT ’21, p. 375–385, New York, NY, USA, 2021. Association for Computing Machinery.
- [JWC22] Haewon Jeong, Hao Wang, and Flavio P. Calmon. “Fairness without Imputation: A Decision Tree Approach for Fair Prediction with Missing Values.” *Proceedings of the AAAI Conference on Artificial Intelligence*, **36**(9):9558–9566, Jun. 2022.
- [JZ21] Bo Jiang and Ziyang Zhang. “Incomplete Graph Representation and Learning via Partial Graph Neural Networks.” *arXiv preprint arXiv:2003.10130*, 2021.
- [JZY23] Mingxuan Ju, Tong Zhao, Wenhao Yu, Neil Shah, and Yanfang Ye. “GraphPatcher: Mitigating Degree Bias for Graph Neural Networks via Test-time Augmentation.”

- In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pp. 55785–55801. Curran Associates, Inc., 2023.
- [KA21] Maximilian Kasy and Rediet Abebe. “Fairness, Equality, and Power in Algorithmic Decision-Making.” In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, FAccT ’21, p. 576–586, New York, NY, USA, 2021. Association for Computing Machinery.
- [Kar15] V. Kargin. “Subordination for the sum of two random matrices.” *The Annals of Probability*, **43**(4):2119 – 2150, 2015.
- [KAS11] Toshihiro Kamishima, Shotaro Akaho, and Jun Sakuma. “Fairness-aware Learning through Regularization Approach.” In *2011 IEEE 11th International Conference on Data Mining Workshops*, pp. 643–650, 2011.
- [KB15] Diederik P. Kingma and Jimmy Ba. “Adam: A Method for Stochastic Optimization.” In *International Conference on Learning Representations*, 2015.
- [KBB14] Vassilis Kalofolias, Xavier Bresson, Michael Bronstein, and Pierre Vandergheynst. “Matrix Completion on Graphs.” *arXiv preprint arXiv:1408.1717*, 2014.
- [Ker22] Nicolas Keriven. “Not too little, not too much: a theoretical analysis of graph (over)smoothing.” In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pp. 2268–2281. Curran Associates, Inc., 2022.
- [KL21] Fereshte Khani and Percy Liang. “Removing Spurious Features can Hurt Accuracy and Affect Groups Disproportionately.” In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, FAccT ’21, p. 196–205, New York, NY, USA, 2021. Association for Computing Machinery.

- [KNR18] Michael Kearns, Seth Neel, Aaron Roth, and Zhiwei Steven Wu. “Preventing Fairness Gerrymandering: Auditing and Learning for Subgroup Fairness.” In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 2564–2572. PMLR, 10–15 Jul 2018.
- [KW14] Diederik P. Kingma and Max Welling. “Auto-Encoding Variational Bayes.” In *International Conference on Learning Representations*, 2014.
- [KW17] Thomas N. Kipf and Max Welling. “Semi-Supervised Classification with Graph Convolutional Networks.” In *International Conference on Learning Representations*, 2017.
- [KZX22] Jian Kang, Yan Zhu, Yinglong Xia, Jiebo Luo, and Hanghang Tong. “RawlsGCN: Towards Rawlsian Difference Principle on Graph Convolutional Network.” In *Proceedings of the ACM Web Conference 2022, WWW ’22*, p. 1214–1225, New York, NY, USA, 2022. Association for Computing Machinery.
- [LAN24] Kristian Lum, Jacy Reese Anthis, Chirag Nagpal, and Alex D’Amour. “Bias in Language Models: Beyond Trick Tests and Toward RUTEd Evaluation.” *arXiv preprint arXiv:2402.12649*, 2024.
- [LB00] Bruce G. Lindsay and Prasanta Basak. “Moments Determine the Tail of a Distribution (But Not Much Else).” *The American Statistician*, **54**(4):248–251, 2000.
- [LBC20] Preethi Lahoti, Alex Beutel, Jilin Chen, Kang Lee, Flavien Prost, Nithum Thain, Xuezhi Wang, and Ed Chi. “Fairness without Demographics through Adversarially Reweighted Learning.” In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pp. 728–740. Curran Associates, Inc., 2020.

- [LCH22] Anne Lauscher, Archie Crowley, and Dirk Hovy. “Welcome to the Modern World of Pronouns: Identity-Inclusive Natural Language Processing beyond Gender.” In Nicoletta Calzolari, Chu-Ren Huang, Hansaem Kim, James Pustejovsky, Leo Wanner, Key-Sun Choi, Pum-Mo Ryu, Hsin-Hsi Chen, Lucia Donatelli, Heng Ji, Sadao Kurohashi, Patrizia Paggio, Nianwen Xue, Seokhwan Kim, Younggyun Hahm, Zhong He, Tony Kyungil Lee, Enrico Santus, Francis Bond, and Seung-Hoon Na, editors, *Proceedings of the 29th International Conference on Computational Linguistics*, pp. 1221–1232, Gyeongju, Republic of Korea, oct 2022. International Committee on Computational Linguistics.
- [LFB19] Klas Leino, Matt Fredrikson, Emily Black, Shayak Sen, and Anupam Datta. “Feature-Wise Bias Amplification.” In *International Conference on Learning Representations*, 2019.
- [LFZ24] Zemin Liu, Yuan Fang, Wentao Zhang, Xinming Zhang, and Steven C.H. Hoi. “Locality-Aware Tail Node Embeddings on Homogeneous and Heterogeneous Networks.” *IEEE Transactions on Knowledge and Data Engineering*, **36**(6):2517–2532, 2024.
- [LGS24] Li Lucy, Suchin Gururangan, Luca Soldaini, Emma Strubell, David Bamman, Lauren Klein, and Jesse Dodge. “AboutMe: Using Self-Descriptions in Webpages to Document the Effects of English Pretraining Data Filters.” In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 7393–7420, Bangkok, Thailand, August 2024. Association for Computational Linguistics.
- [LHL22] Sitao Luan, Chenqing Hua, Qincheng Lu, Jiaqi Zhu, Mingde Zhao, Shuyuan Zhang, Xiao-Wen Chang, and Doina Precup. “Revisiting Heterophily For Graph Neural Networks.” In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and

- A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pp. 1362–1375. Curran Associates, Inc., 2022.
- [Li23] Jintang Li. “Visualize ICLR 2023 OpenReview Data.”, 2023. <https://github.com/EdisonLeeeee/ICLR2023-OpenReviewData>.
- [LJH24] Dawei Li, Bohan Jiang, Liangjie Huang, Alimohammad Beigi, Chengshuai Zhao, Zhen Tan, Amrita Bhattacharjee, Yuxuan Jiang, Canyu Chen, Tianhao Wu, et al. “From generation to judgment: Opportunities and challenges of llm-as-a-judge.” *arXiv preprint arXiv:2411.16594*, 2024.
- [LLC23] Jie Liao, Jintang Li, Liang Chen, Bingzhe Wu, Yatao Bian, and Zibin Zheng. “SAILOR: Structural Augmentation Based Tail Node Representation Learning.” In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management, CIKM '23*, p. 1389–1399, New York, NY, USA, 2023. Association for Computing Machinery.
- [LLC24] Charlotte Laclau, Christine Largeton, and Manvi Choudhary. “A Survey on Fairness for Machine Learning on Graphs.” *arXiv preprint arXiv:2205.05396*, 2024.
- [LLD09] Jure Leskovec, Kevin J. Lang, Anirban Dasgupta, and Michael W. Mahoney. “Community Structure in Large Networks: Natural Cluster Sizes and the Absence of Large Well-Defined Clusters.” *Internet Mathematics*, **6**(1), jan 1 2009.
- [LLP97] Lawrence R. Landerman, Kenneth C. Land, and Carl F. Pieper. “An Empirical Evaluation of the Predictive Mean Matching Method for Imputing Missing Values.” *Sociological Methods & Research*, **26**(1):3–33, 1997.
- [LMH23] Donghwan Lee, Behrad Moniri, Xinmeng Huang, Edgar Dobriban, and Hamed Hassani. “Demystifying disagreement-on-the-line in high dimensions.” In *International Conference on Machine Learning*, pp. 19053–19093. PMLR, 2023.

- [LMM23] Ting Wei Li, Qiaozhu Mei, and Jiaqi Ma. “A metadata-driven approach to understand graph neural networks.” In *Proceedings of the 37th International Conference on Neural Information Processing Systems, NIPS ’23*, Red Hook, NY, USA, 2023. Curran Associates Inc.
- [LNF21] Zemin Liu, Trung-Kien Nguyen, and Yuan Fang. “Tail-GNN: Tail-Node Graph Neural Networks.” In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining, KDD ’21*, p. 1109–1119, New York, NY, USA, 2021. Association for Computing Machinery.
- [LNF23] Zemin Liu, Trung-Kien Nguyen, and Yuan Fang. “On Generalized Degree Fairness in Graph Neural Networks.” In *Proceedings of the Thirty-Seventh AAAI Conference on Artificial Intelligence and Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence and Thirteenth Symposium on Educational Advances in Artificial Intelligence, AAAI’23/IAAI’23/EAAI’23*. AAAI Press, 2023.
- [Lov96] L. Lovász. “Random Walks on Graphs: A Survey.” In D. Miklós, V. T. Sós, and T. Szőnyi, editors, *Combinatorics, Paul Erdős is Eighty*, volume 2, pp. 353–398. János Bolyai Mathematical Society, Budapest, 1996.
- [LPB22] Donald Loveland, Jiayi Pan, Aaresh Farrokh Bhatena, and Yiyang Lu. “FairEdit: Preserving Fairness in Graph Neural Networks through Greedy Graph Editing.” *arXiv preprint arXiv:2201.03681*, 2022.
- [LWN22] Yanying Li, Xiuling Wang, Yue Ning, and Hui Wang. “FairLP: Towards Fair Link Prediction on Social Network Graphs.” *Proceedings of the International AAAI Conference on Web and Social Media*, **16**(1):628–639, May 2022.
- [LWW20] Pan Li, Yanbang Wang, Hongwei Wang, and Jure Leskovec. “Distance Encoding: Design Provably More Powerful Neural Networks for Graph Representation Learn-

- ing.” In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS’20, Red Hook, NY, USA, 2020. Curran Associates Inc.
- [LWZ21] Peizhao Li, Yifei Wang, Han Zhao, Pengyu Hong, and Hongfu Liu. “On Dyadic Fairness: Exploring and Mitigating Bias in Graph Connections.” In *International Conference on Learning Representations*, 2021.
- [LXC23] Xueqi Li, Guoqing Xiao, Yuedan Chen, Zhuo Tang, Wenjun Jiang, and Kenli Li. “An Explicitly Weighted GCN Aggregator based on Temporal and Popularity Features for Recommendation.” *ACM Trans. Recomm. Syst.*, **1**(2), apr 2023.
- [LXS23] Langzhang Liang, Zenglin Xu, Zixing Song, Irwin King, Yuan Qi, and Jieping Ye. “Tackling Long-tailed Distribution Issue in Graph Neural Networks via Normalization.” *IEEE Transactions on Knowledge and Data Engineering*, pp. 1–11, 2023.
- [LYD22] Songtao Liu, Rex Ying, Hanze Dong, Lanqing Li, Tingyang Xu, Yu Rong, Peilin Zhao, Junzhou Huang, and Dinghao Wu. “Local Augmentation for Graph Neural Networks.” In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 14054–14072. PMLR, 17–23 Jul 2022.
- [LZX22] Yazheng Liu, Xi Zhang, and Sihong Xie. “Trade less Accuracy for Fairness and Trade-off Explanation for GNN.” In *2022 IEEE International Conference on Big Data (Big Data)*, pp. 4681–4690, 2022.
- [MAC25] Joel Mire, Zubin Trivadi Aysola, Daniel Chechelnitsky, Nicholas Deas, Chrysoula Zerva, and Maarten Sap. “Rejected Dialects: Biases Against African American Language in Reward Models.”, 2025.

- [McN21] Chan Tov McNamara. “Misgendering.” *Calif. L. Rev.. California Law Review*, **109**(IR):2227, 2021.
- [MFP25] Nestor Maslej, Loredana Fattorini, Raymond Perrault, Yolanda Gil, Vanessa Parli, Njenga Kariuki, Emily Capstick, Anka Reuel, Erik Brynjolfsson, John Etchemendy, Katrina Ligett, Terah Lyons, James Manyika, Juan Carlos Niebles, Yoav Shoham, Russell Wald, Tobi Walsh, Armin Hamrah, Lapo Santarlasci, Julia Betts Lotufo, Alexandra Rome, Andrew Shi, and Sukrut Oak. “The AI Index 2025 Annual Report.” *Institute for Human-Centered AI*, April 2025.
- [MGR24] Stefano Sarao Mannelli, Federica Gerace, Negar Rostamzadeh, and Luca Saglietti. “Bias-inducing geometries: exactly solvable data model with fairness implications.” In *ICML 2024 Workshop on Geometry-grounded Representation Learning and Generative Modeling*, 2024.
- [MGW22] Jing Ma, Ruocheng Guo, Mengting Wan, Longqi Yang, Aidong Zhang, and Jundong Li. “Learning Fair Node Representations with Graph Counterfactual Fairness.” In *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining, WSDM ’22*, p. 695–703, New York, NY, USA, 2022. Association for Computing Machinery.
- [MLS22] Yao Ma, Xiaorui Liu, Neil Shah, and Jiliang Tang. “Is Homophily a Necessity for Graph Neural Networks?” In *International Conference on Learning Representations*, 2022.
- [MM11] Fragkiskos D. Malliaros and Vasileios Megalooikonomou. “Expansion Properties of Large Social Graphs.” In Jianliang Xu, Ge Yu, Shuigeng Zhou, and Rainer Unland, editors, *Database Systems for Advanced Applications*, pp. 311–322, Berlin, Heidelberg, 2011. Springer Berlin Heidelberg.

- [MP67] V.A. Marčenko and Leonid Pastur. “Distribution of eigenvalues for some sets of random matrices.” *Math USSR Sb*, **1**:457–483, 1967.
- [MRS22] Alexander Maloney, Daniel A Roberts, and James Sully. “A solvable model of neural scaling laws.” *arXiv preprint arXiv:2210.16859*, 2022.
- [MS17] James A. Mingo and Roland Speicher. *Free Probability and Random Matrices*, volume 35 of *Fields Institute Monographs*. Springer, 2017.
- [MSR24] Jessica Moeder, William J Scarborough, and Barbara Risman. “Not Just They/Them: Exploring Diversity and Meaning in Pronoun Use among Non-Binary Individuals.” *Social Problems*, p. spae064, 10 2024.
- [MWY20] Farzan Masrouf, Tyler Wilson, Heng Yan, Pang-Ning Tan, and Abdol Esfahanian. “Bursting the Filter Bubble: Fairness-Aware Network Link Prediction.” *Proceedings of the AAAI Conference on Artificial Intelligence*, **34**(01):841–848, Apr. 2020.
- [NDC17] Jekaterina Novikova, Ondřej Dušek, Amanda Cercas Curry, and Verena Rieser. “Why We Need New Evaluation Metrics for NLG.” In Martha Palmer, Rebecca Hwa, and Sebastian Riedel, editors, *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pp. 2241–2252, Copenhagen, Denmark, sep 2017. Association for Computational Linguistics.
- [NKK19] Preetum Nakkiran, Gal Kaplun, Dimitris Kalimeris, Tristan Yang, Benjamin L. Edelman, Fred Zhang, and Boaz Barak. *SGD on neural networks learns functions of increasing complexity*. Curran Associates Inc., Red Hook, NY, USA, 2019.
- [OCD19] Alexandra Olteanu, Carlos Castillo, Fernando Diaz, and Emre Kıcıman. “Social Data: Biases, Methodological Pitfalls, and Ethical Boundaries.” *Frontiers in Big Data*, **2**, 2019.

- [Och92] Elinor Ochs. “Indexing gender.” In Alessandro Duranti and Charles Goodwin, editors, *Rethinking Context: Language as an Interactive Phenomenon*, p. 335–358. Cambridge University Press, Cambridge, 1992.
- [OGD23] Anaelia Ovalle, Palash Goyal, Jwala Dhamala, Zachary Jagers, Kai-Wei Chang, Aram Galstyan, Richard Zemel, and Rahul Gupta. “‘I’m fully who I am’: Towards Centering Transgender and Non-Binary Voices to Measure Biases in Open Language Generation.” In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency, FAccT ’23*, p. 1246–1266, New York, NY, USA, 2023. Association for Computing Machinery.
- [OMG24] Anaelia Ovalle, Ninareh Mehrabi, Palash Goyal, Jwala Dhamala, Kai-Wei Chang, Richard Zemel, Aram Galstyan, Yuval Pinter, and Rahul Gupta. “Tokenization Matters: Navigating Data-Scarce Tokenization for Gender Inclusive Language Technologies.” In Kevin Duh, Helena Gomez, and Steven Bethard, editors, *Findings of the Association for Computational Linguistics: NAACL 2024*, pp. 1739–1756, Mexico City, Mexico, jun 2024. Association for Computational Linguistics.
- [Ope] OpenAI. “ChatGPT.” <https://chatgpt.com/>.
- [OPM24] Anaelia Ovalle, Krunoslav Lehman Pavasovic, Louis Martin, Luke Zettlemoyer, Eric Michael Smith, Adina Williams, and Levent Sagun. “The Root Shapes the Fruit: On the Persistence of Gender-Exclusive Harms in Aligned Language Models.” *NeurIPS Queer in AI Workshop*, 2024.
- [OS20] Kenta Oono and Taiji Suzuki. “Graph Neural Networks Exponentially Lose Expressive Power for Node Classification.” In *International Conference on Learning Representations*, 2020.
- [OSG23] Anaelia Ovalle, Arjun Subramonian, Vagrant Gautam, Gilbert Gee, and Kai-Wei Chang. “Factoring the Matrix of Domination: A Critical Review and Reimagination

- of Intersectionality in AI Fairness.” In *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*, AIES '23, p. 496–511, New York, NY, USA, 2023. Association for Computing Machinery.
- [Ott19] S.A. Otto. “How to normalize the RMSE [Blog post].”, 2019.
- [PDS05] K. Pelckmans, J. De Brabanter, J.A.K. Suykens, and B. De Moor. “Handling missing values in support vector machine classifiers.” *Neural Networks*, **18**(5):684–692, 2005. IJCNN 2005.
- [Per] Perplexity. “Perplexity AI.” <https://www.perplexity.ai/>.
- [PGM19] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, et al. “PyTorch: An Imperative Style, High-Performance Deep Learning Library.” In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pp. 8024–8035. Curran Associates, Inc., Red Hook, NY, USA, 2019.
- [PP20] John Palowitch and Bryan Perozzi. “Debiasing Graph Representations via Metadata-Orthogonal Training.” In *2020 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM)*, pp. 435–442, 2020.
- [PS22] Dana Pessach and Erez Shmueli. “A Review on Fairness in Machine Learning.” *ACM Comput. Surv.*, **55**(3), feb 2022.
- [QDO23] Organizers Of QueerInai, Nathan Dennler, Anaelia Ovalle, Ashwin Singh, Luca Soldaini, Arjun Subramonian, Huy Tu, William Agnew, Avijit Ghosh, Kyra Yee,

Irene Font Peradejordi, Zeerak Talat, Mayra Russo, and Jess De Jesus De Pinho Pinhal. “Bound by the Bounty: Collaboratively Shaping Evaluation Processes for Queer AI Harms.” In *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*, AIES ’23, p. 375–386, New York, NY, USA, 2023. Association for Computing Machinery.

[QOS23] Organizers Of QueerInai, Anaelia Ovalle, Arjun Subramonian, Ashwin Singh, Claas Voelcker, Danica J. Sutherland, Davide Locatelli, Eva Breznik, Filip Klubicka, Hang Yuan, Hetvi J, Huan Zhang, Jaidev Shriram, Kruno Lehman, Luca Soldaini, Maarten Sap, Marc Peter Deisenroth, Maria Leonor Pacheco, Maria Ryskina, Martin Mundt, Milind Agarwal, Nyx Mclean, Pan Xu, A Pranav, Raj Korpan, Ruchira Ray, Sarah Mathew, Sarthak Arora, St John, Tanvi Anand, Vishakha Agrawal, William Agnew, Yanan Long, Zijie J. Wang, Zeerak Talat, Avijit Ghosh, Nathaniel Dennler, Michael Noseworthy, Sharvani Jha, Emi Baylor, Aditya Joshi, Natalia Y. Bilenko, Andrew Mcnamara, Raphael Gontijo-Lopes, Alex Markham, Eryn Dong, Jackie Kay, Manu Saraswat, Nikhil Vytla, and Luke Stark. “Queer In AI: A Case Study in Community-Led Participatory AI.” In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, FAccT ’23, p. 1882–1895, New York, NY, USA, 2023. Association for Computing Machinery.

[Ram21] Nithya E Ramesh, editor. *I hope you’ll still love me*. Evolving Still, June 2021.

[Ram24] Marco Ramponi. “AI trends in 2024: Graph Neural Networks.”, 2024. <https://www.assemblyai.com/blog/ai-trends-graph-neural-networks>.

[RAS21] Benedek Rozemberczki, Carl Allen, and Rik Sarkar. “Multi-Scale attributed node embedding.” *Journal of Complex Networks*, **9**(2):cnab014, 05 2021.

[RKB24] Megan Richards, Polina Kirichenko, Diane Bouchacourt, and Mark Ibrahim. “Does Progress On Object Recognition Benchmarks Improve Generalization on

- Crowdsourced, Global Data?” In *The Twelfth International Conference on Learning Representations*, 2024.
- [RKG22] Emanuele Rossi, Henry Kenlay, Maria I. Gorinova, Benjamin Paul Chamberlain, Xiaowen Dong, and Michael M. Bronstein. “On the Unreasonable Effectiveness of Feature Propagation in Learning on Graphs with Missing Node Features.” In *Learning on Graphs Conference*, 2022.
- [RLH01] Trivellore E. Raghunathan, James M. Lepkowski, John Van Hoewyk, and Peter W. Solenberger. “A multivariate technique for multiply imputing missing values using a sequence of regression models.” *Survey Methodology*, **27**:85–95, 2001.
- [RMR21] Dominic Richards, Jaouad Mourtada, and Lorenzo Rosasco. “Asymptotics of ridge (less) regression under general source condition.” In *International Conference on Artificial Intelligence and Statistics*, pp. 3889–3897. PMLR, 2021.
- [RS20] Benedek Rozemberczki and Rik Sarkar. “Characteristic Functions on Graphs: Birds of a Feather, from Statistical Descriptors to Parametric Models.” In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management, CIKM ’20*, p. 1325–1334, New York, NY, USA, 2020. Association for Computing Machinery.
- [RSB19] Tahleen Rahman, Bartłomiej Surma, Michael Backes, and Yang Zhang. “Fairwalk: Towards Fair Graph Embedding.” In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*, pp. 3289–3295. International Joint Conferences on Artificial Intelligence Organization, 7 2019.
- [RWC19] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. “Language Models are Unsupervised Multitask Learners.” 2019.
- [SB11] Daniel J. Stekhoven and Peter Bühlmann. “MissForest—non-parametric missing value imputation for mixed-type data.” *Bioinformatics*, **28**(1):112–118, 10 2011.

- [SB24] Morgan Klaus Scheuerman and Jed R. Brubaker. “Products of Positionality: How Tech Workers Shape Identity Concepts in Computer Vision.” In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, CHI ’24, New York, NY, USA, 2024. Association for Computing Machinery.
- [SBS25] Arjun Subramonian, Samuel J. Bell, Levent Sagun, and Elvis Dohmatob. “An Effective Theory of Bias Amplification.” In *The Thirteenth International Conference on Learning Representations*, 2025.
- [SCS22] Arjun Subramonian, Kai-Wei Chang, and Yizhou Sun. “On the Discrimination Risk of Mean Aggregation Feature Imputation in Graphs.” In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pp. 32957–32973. Curran Associates, Inc., 2022.
- [SD25] Arjun Subramonian and Elvis Dohmatob. “auto-fpt: Automating Free Probability Theory Calculations for Machine Learning Theory.”, 2025.
- [SDK23] Nikil Selvam, Sunipa Dev, Daniel Khashabi, Tushar Khot, and Kai-Wei Chang. “The Tail Wagging the Dog: Dataset Construction Biases of Social Bias Benchmarks.” In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pp. 1373–1386, Toronto, Canada, jul 2023. Association for Computational Linguistics.
- [Sek08] Jasjeet Sekhon. “271 The Neyman–Rubin Model of Causal Inference and Estimation Via Matching Methods.” In *The Oxford Handbook of Political Methodology*. Oxford University Press, 08 2008.
- [SGd19] Stefano Spigler, Mario Geiger, Stéphane d’Ascoli, Levent Sagun, Giulio Biroli, and Matthieu Wyart. “A jamming transition from under-to over-parametrization

- affects generalization in deep learning.” *Journal of Physics A: Mathematical and Theoretical*, **52**(47):474001, 2019.
- [SGS25] Arjun Subramonian, Vagrant Gautam, Preethi Seshadri, Dietrich Klakow, Kai-Wei Chang, and Yizhou Sun. “Agree to Disagree? A Meta-Evaluation of LLM Misgendering.”, 2025.
- [SGT09] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. “The Graph Neural Network Model.” *IEEE Transactions on Neural Networks*, **20**(1):61–80, 2009.
- [SHC20] Ana-Andreea Stoica, Jessy Xinyi Han, and Augustin Chaintreau. “Seeding Network Influence in Biased Networks and the Benefits of Diversity.” In *Proceedings of The Web Conference 2020*, WWW ’20, p. 2089–2098, New York, NY, USA, 2020. Association for Computing Machinery.
- [SHK14] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. “Dropout: A Simple Way to Prevent Neural Networks from Overfitting.” *Journal of Machine Learning Research*, **15**(56):1929–1958, 2014.
- [SHM25] Siyuan Song, Jennifer Hu, and Kyle Mahowald. “Language Models Fail to Introspect About Their Knowledge of Language.” *arXiv preprint arXiv:2503.07513*, 2025.
- [SJW23] Harry Shomer, Wei Jin, Wentao Wang, and Jiliang Tang. “Toward Degree Bias in Embedding-Based Knowledge Graph Completion.” In *Proceedings of the ACM Web Conference 2023*, WWW ’23, p. 705–715, New York, NY, USA, 2023. Association for Computing Machinery.
- [SKB24] Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Authur, Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar,

- Valentin Hofmann, Ananya Jha, Sachin Kumar, Li Lucy, Xinxi Lyu, Nathan Lambert, Ian Magnusson, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik, et al. “Dolma: an Open Corpus of Three Trillion Tokens for Language Model Pre-training Research.” In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15725–15788, Bangkok, Thailand, aug 2024. Association for Computational Linguistics.
- [SKS24] Arjun Subramonian, Jian Kang, and Yizhou Sun. “Theoretical and Empirical Insights into the Origins of Degree Bias in Graph Neural Networks.” In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [SLN20] Julian Salazar, Davis Liang, Toan Q. Nguyen, and Katrin Kirchhoff. “Masked Language Model Scoring.” In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 2699–2712, Online, jul 2020. Association for Computational Linguistics.
- [SLY21] Aravind Sankar, Yozen Liu, Jun Yu, and Neil Shah. “Graph Neural Networks for Friend Ranking in Large-scale Social Platforms.” In *Proceedings of the Web Conference 2021, WWW ’21*, p. 2535–2546, New York, NY, USA, 2021. Association for Computing Machinery.
- [SMB19] Oleksandr Shchur, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann. “Pitfalls of Graph Neural Network Evaluation.” *arXiv preprint arXiv:1811.05868*, 2019.
- [SMP15] Ilya Shpitser, Karthika Mohan, and Judea Pearl. “Missing data as a causal and probabilistic problem.” In *Proceedings of the Thirty-First Conference on Uncertainty in Artificial Intelligence, UAI’15*, p. 802–811, Arlington, Virginia, USA, 2015. AUAI Press.

- [SP10] Skipper Seabold and Josef Perktold. “statsmodels: Econometric and statistical modeling with python.” In *9th Python in Science Conference*, 2010.
- [SPS22] Preethi Seshadri, Pouya Pezeshkpour, and Sameer Singh. “Quantifying Social Biases Using Templates is Unreliable.” In *Workshop on Trustworthy and Socially Responsible Machine Learning, NeurIPS 2022*, 2022.
- [SRC18] Ana-Andreea Stoica, Christopher Riederer, and Augustin Chaintreau. “Algorithmic Glass Ceiling in Social Networks: The Effects of Social Recommendations on Network Diversity.” In *Proceedings of the 2018 World Wide Web Conference, WWW ’18*, p. 923–932, Republic and Canton of Geneva, CHE, 2018. International World Wide Web Conferences Steering Committee.
- [SRK20] Shiori Sagawa, Aditi Raghunathan, Pang Wei Koh, and Percy Liang. “An investigation of why overparameterization exacerbates spurious correlations.” In *International Conference on Machine Learning*, pp. 8346–8356. PMLR, 2020.
- [SSH22] Indro Spinelli, Simone Scardapane, Amir Hussain, and Aurelio Uncini. “FairDrop: Biased Edge Dropout for Enhancing Fairness in Graph Representation Learning.” *IEEE Transactions on Artificial Intelligence*, **3**(3):344–354, 2022.
- [SSS24] Arjun Subramonian, Levent Sagun, and Yizhou Sun. “Networked inequality: preferential attachment bias in graph neural network link prediction.” In *Proceedings of the 41st International Conference on Machine Learning, ICML’24*. JMLR.org, 2024.
- [SST18] Marek Śmieja, Łukasz Struski, Jacek Tabor, Bartosz Zieliński, and Przemysław Spurek. “Processing of missing data by neural networks.” In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.

- [SSU20] Indro Spinelli, Simone Scardapane, and Aurelio Uncini. “Missing data imputation with adversarially-trained graph convolutional networks.” *Neural Networks*, **129**:249–260, 2020.
- [Ste21] Jacob Steinhardt. “Lecture Notes for STAT260 (Robust and Nonparametric Statistics).” 2021.
- [Sub22] Arjun Subramonian. “On Dyadic Fairness: Exploring and Mitigating Bias in Graph Connections.” In *ICLR Blog Track*, 2022. <https://iclr-blog-track.github.io/2022/03/25/dyadic-fairness/>.
- [Sub23] Arjun Subramonian. “Queer in AI National AI Advisory Committee Briefing.”, 2023. <https://www.queerinaai.com/naiac-briefing>.
- [SWC09] Jonathan A. C. Sterne, Ian R. White, John B. Carlin, Michael Spratt, Patrick Royston, Michael G. Kenward, Angela M. Wood, and James R. Carpenter. “Multiple imputation for missing data in epidemiological and clinical research: potential and pitfalls.” *The BMJ*, **338**, 2009.
- [SWS21] Tony Sun, Kellie Webster, Apu Shah, William Yang Wang, and Melvin Johnson. “They, Them, Theirs: Rewriting with Gender-Neutral English.” *arXiv preprint arXiv:2102.06788*, 2021.
- [SYD23] Arjun Subramonian, Xingdi Yuan, Hal Daumé III, and Su Lin Blodgett. “It Takes Two to Tango: Navigating Conceptualizations of NLP Tasks and Measurements of Performance.” In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 3234–3279, Toronto, Canada, jul 2023. Association for Computational Linguistics.
- [TAP21] Nilesch Tripuraneni, Ben Adlam, and Jeffrey Pennington. “Overparameterization Improves Robustness to Covariate Shift in High Dimensions.” In M. Ranzato,

- A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pp. 13883–13897. Curran Associates, Inc., 2021.
- [TCS01] Olga Troyanskaya, Michael Cantor, Gavin Sherlock, Pat Brown, Trevor Hastie, Robert Tibshirani, David Botstein, and Russ B. Altman. “Missing value estimation methods for DNA microarrays.” *Bioinformatics*, **17**(6):520–525, 06 2001.
- [TLM21] Hibiki Taguchi, Xin Liu, and Tsuyoshi Murata. “Graph convolutional networks for graphs containing missing features.” *Future Generation Computer Systems*, **117**:155–168, 2021.
- [TYS20] Xianfeng Tang, Huaxiu Yao, Yiwei Sun, Yiqi Wang, Jiliang Tang, Charu Aggarwal, Prasenjit Mitra, and Suhang Wang. “Investigating and Mitigating Degree-Related Biases in Graph Convolutional Networks.” In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management, CIKM ’20*, p. 1435–1444, New York, NY, USA, 2020. Association for Computing Machinery.
- [TZY08] Jie Tang, Jing Zhang, Limin Yao, Juanzi Li, Li Zhang, and Zhong Su. “Arnet-Miner: Extraction and Mining of Academic Social Networks.” In *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD ’08*, p. 990–998, New York, NY, USA, 2008. Association for Computing Machinery.
- [Ung24] Eddie Ungless. “What’s in my big data? Transphobia.”, 2024. <https://mxeddie.github.io/2024/03/30/whats-in-my-data.html>.
- [Van21] Lieven Vandenberghe. “ECE236C - Optimization Methods for Large-Scale Systems Course Notes.”, 2021.
- [VCC18] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro

- Liò, and Yoshua Bengio. “Graph Attention Networks.” In *International Conference on Learning Representations*, 2018.
- [VCL19] Guillermo Valle-Perez, Chico Q. Camargo, and Ard A. Louis. “Deep learning generalizes because the parameter-function map is biased towards simple functions.” In *International Conference on Learning Representations*, 2019.
- [VGO20] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C J Carey, et al. “SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python.” *Nature Methods*, **17**:261–272, 2020.
- [Vig19] Jesse Vig. “A Multiscale Visualization of Attention in the Transformer Model.” In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pp. 37–42, Florence, Italy, July 2019. Association for Computational Linguistics.
- [VR18] Sahil Verma and Julia Rubin. “Fairness Definitions Explained.” In *Proceedings of the International Workshop on Software Fairness, FairWare ’18*, p. 1–7, New York, NY, USA, 2018. Association for Computing Machinery.
- [VS23] Srinivas Virinchi and Anoop Saladi. “BLADE: Biased Neighborhood Sampling based Graph Neural Network for Directed Graphs.” In *Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining, WSDM ’23*, p. 42–50, New York, NY, USA, 2023. Association for Computing Machinery.
- [WBS18] Betsy Anne Williams, Catherine F. Brooks, and Yotam Shmargad. “How Algorithms Discriminate Based on Data They Lack: Challenges, Solutions, and Policy Implications.” *Journal of Information Policy*, **8**(1):78–115, 03 2018.

- [WC24] Yixin Wan and Kai-Wei Chang. “The Male CEO and the Female Assistant: Evaluation and Mitigation of Gender Biases in Text-To-Image Generation of Dual Subjects.” *arXiv preprint arXiv:2402.11089*, 2024.
- [WD22] Yu Wang and Tyler Derr. “Degree-Related Bias in Link Prediction.” In *2022 IEEE International Conference on Data Mining Workshops (ICDMW)*, pp. 757–758, 2022.
- [WDS20] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, et al. “Transformers: State-of-the-Art Natural Language Processing.” In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 38–45, Online, oct 2020. Association for Computational Linguistics.
- [WGJ21] Christo Wilson, Avijit Ghosh, Shan Jiang, Alan Mislove, Lewis Baker, Janelle Szary, Kelly Trindel, and Frida Polli. “Building and Auditing Fair Algorithms: A Case Study in Candidate Screening.” In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, FAccT ’21*, p. 666–677, New York, NY, USA, 2021. Association for Computing Machinery.
- [WHX19] Jun Wu, Jingrui He, and Jiejun Xu. “DEMO-Net: Degree-specific Graph Neural Networks for Node and Graph Classification.” In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD ’19*, p. 406–415, New York, NY, USA, 2019. Association for Computing Machinery.
- [WLL22] Nan Wang, Lu Lin, Jundong Li, and Hongning Wang. “Unbiased Graph Embedding with Biased Graph Observations.” In *Proceedings of the ACM Web Conference*

- 2022, WWW '22, p. 1423–1433, New York, NY, USA, 2022. Association for Computing Machinery.
- [WLL23] Chunyu Wei, Jian Liang, Di Liu, Zehui Dai, Mang Li, and Fei Wang. “Meta Graph Learning for Long-tail Recommendation.” In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, KDD '23, p. 2512–2522, New York, NY, USA, 2023. Association for Computing Machinery.
- [WLM25] Oskar van der Wal, Pietro Lesci, Max Müller-Eberstein, Naomi Saphra, Hailey Schoelkopf, Willem Zuidema, and Stella Biderman. “PolyPythias: Stability and Outliers across Fifty Language Model Pre-Training Runs.” In *The Thirteenth International Conference on Learning Representations*, 2025.
- [WLX05] David Williams, Xuejun Liao, Ya Xue, and Lawrence Carin. “Incomplete-data classification using logistic regression.” In *Proceedings of the 22nd International Conference on Machine Learning*, ICML '05, p. 972–979, New York, NY, USA, 2005. Association for Computing Machinery.
- [Woo07] Jeffrey M. Wooldridge. “Inverse probability weighted estimation for general missing data problems.” *Journal of Econometrics*, **141**(2):1281–1301, 2007.
- [WPD24] Jamelle Watson-Daniels, Flavio du Pin Calmon, Alexander D’Amour, Carol Long, David C. Parkes, and Berk Ustun. “Predictive Churn with the Set of Good Models.” *arXiv preprint arXiv:2402.07745*, 2024.
- [WR21] Angelina Wang and Olga Russakovsky. “Directional Bias Amplification.” In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 10882–10893. PMLR, 18–24 Jul 2021.
- [WRR22] Angelina Wang, Vikram V Ramaswamy, and Olga Russakovsky. “Towards Intersectionality in Machine Learning: Including More Identities, Handling Under-

- representation, and Performing Evaluation.” In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, FAccT ’22, p. 336–349, New York, NY, USA, 2022. Association for Computing Machinery.
- [WSP24] Sierra Wyllie, Iliia Shumailov, and Nicolas Papernot. “Fairness feedback loops: training on synthetic data amplifies bias.” In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*, FAccT ’24, pp. 2113–2147, New York, NY, USA, 2024. Association for Computing Machinery.
- [WSZ19] Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. “Simplifying Graph Convolutional Networks.” In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 6861–6871. PMLR, 09–15 Jun 2019.
- [WWF21] Jiancan Wu, Xiang Wang, Fuli Feng, Xiangnan He, Liang Chen, Jianxun Lian, and Xing Xie. “Self-supervised Graph Learning for Recommendation.” In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR ’21, p. 726–735, New York, NY, USA, 2021. Association for Computing Machinery.
- [WWS22] Ruijia Wang, Xiao Wang, Chuan Shi, and Le Song. “Uncovering the structural fairness in graph contrastive learning.” In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS ’22, Red Hook, NY, USA, 2022. Curran Associates Inc.
- [WZY19] Tianlu Wang, Jieyu Zhao, Mark Yatskar, Kai-Wei Chang, and Vicente Ordonez. “Balanced Datasets Are Not Enough: Estimating and Mitigating Gender Bias in Deep Image Representations.” In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 5309–5318, 2019.

- [XBL22] Huimin Xu, Yi Bu, Meijun Liu, Chenwei Zhang, Mengyi Sun, Yi Zhang, Eric Meyer, Eduardo Salas, and Ying Ding. “Team power dynamics and team impact: New perspectives on scientific collaboration using career age as a proxy for team power.” *Journal of the Association for Information Science and Technology*, **73**(10):1489–1505, 2022.
- [XCW21] Teng Xiao, Zhengyu Chen, Donglin Wang, and Suhang Wang. “Learning How to Propagate Messages in Graph Neural Networks.” In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining, KDD ’21*, p. 1894–1903, New York, NY, USA, 2021. Association for Computing Machinery.
- [XHL19] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. “How Powerful are Graph Neural Networks?” In *International Conference on Learning Representations*, 2019.
- [XHZ24] Jianan Xu, Jiajin Huang, Jianwei Zhao, and Jian Yang. “HyNCF: A hybrid normalization strategy via feature statistics for collaborative filtering.” *Expert Systems with Applications*, **238**:121875, 2024.
- [XLT18] Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and Stefanie Jegelka. “Representation Learning on Graphs with Jumping Knowledge Networks.” In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 5453–5462. PMLR, 10–15 Jul 2018.
- [XXH23] Hui Xu, Liyao Xiang, Femke Huang, Yuting Weng, Ruijie Xu, Xinbing Wang, and Chenghu Zhou. “Grace: Graph Self-Distillation and Completion to Mitigate Degree-Related Biases.” In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD ’23*, p. 2813–2824, New York, NY, USA, 2023. Association for Computing Machinery.

- [XZC17] Jing Xia, Shengyu Zhang, Guolong Cai, Li Li, Qing Pan, Jing Yan, and Gangmin Ning. “Adjusted weight voting algorithm for random forests in handling missing values.” *Pattern Recognition*, **69**:52–60, 2017.
- [XZH21] Qianqian Xie, Yutao Zhu, Jimin Huang, Pan Du, and Jian-Yun Nie. “Graph Neural Collaborative Topic Model for Citation Recommendation.” *ACM Trans. Inf. Syst.*, **40**(3), nov 2021.
- [YF22] Josh Yamamoto and Eitan Frachtenberg. “Gender Differences in Collaboration Patterns in Computer Science.” *Publications*, **10**(1), 2022.
- [YHC18] Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L. Hamilton, and Jure Leskovec. “Graph Convolutional Neural Networks for Web-Scale Recommender Systems.” In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD ’18*, p. 974–983, New York, NY, USA, 2018. Association for Computing Machinery.
- [YJS18] Jinsung Yoon, James Jordon, and Mihaela van der Schaar. “GAIN: Missing Data Imputation using Generative Adversarial Nets.” In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 5689–5698. PMLR, 10–15 Jul 2018.
- [YKY22] Sukwon Yun, Kibum Kim, Kanghoon Yoon, and Chanyoung Park. “LTE4G: Long-Tail Experts for Graph Neural Networks.” In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management, CIKM ’22*, p. 2434–2443, New York, NY, USA, 2022. Association for Computing Machinery.
- [YMD20] Jiaxuan You, Xiaobai Ma, Yi Ding, Mykel J Kochenderfer, and Jure Leskovec. “Handling Missing Data with Graph Representation Learning.” In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural*

- Information Processing Systems*, volume 33, pp. 19075–19087. Curran Associates, Inc., 2020.
- [YS19] Gilad Yehudai and Ohad Shamir. “On the Power and Limitations of Random Features for Understanding Neural Networks.” In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [ZBL04] Dengyong Zhou, Olivier Bousquet, Thomas Lal, Jason Weston, and Bernhard Schölkopf. “Learning with Local and Global Consistency.” In S. Thrun, L. Saul, and B. Schölkopf, editors, *Advances in Neural Information Processing Systems*, volume 16. MIT Press, 2004.
- [ZBL18] John R. Zech, Marcus A. Badgeley, Manway Liu, Anthony B. Costa, Joseph J. Titano, and Eric Karl Oermann. “Variable Generalization Performance of a Deep Learning Model to Detect Pneumonia in Chest Radiographs: A Cross-Sectional Study.” *PLOS Medicine*, **15**(11):e1002683, nov 2018.
- [ZC18] Muhan Zhang and Yixin Chen. “Link prediction based on graph neural networks.” In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, NIPS’18, p. 5171–5181, Red Hook, NY, USA, 2018. Curran Associates Inc.
- [ZCY24] Guixian Zhang, Debo Cheng, Guan Yuan, and Shichao Zhang. “Learning fair representations via rebalancing graph structure.” *Information Processing & Management*, **61**(1):103570, 2024.
- [ZDW21] Minghao Zhao, Qilin Deng, Kai Wang, Runze Wu, Jianrong Tao, Changjie Fan, Liang Chen, and Peng Cui. “Bilateral Filtering Graph Convolutional Network for Multi-relational Social Recommendation in the Power-law Networks.” *ACM Trans. Inf. Syst.*, **40**(2), sep 2021.

- [ZG19] Han Zhao and Geoff Gordon. “Inherent Tradeoffs in Learning Fair Representations.” In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [ZGF22] Bingchen Zhao, Yuling Gu, Jessica Zosa Forde, and Naomi Saphra. “One Venue, Two Conferences: The Separation of Chinese and American Citation Networks.” *arXiv preprint arXiv:2211.12424*, 2022.
- [ZHM21] Yiguang Zhang, Jessy Xinyi Han, Ilica Mahajan, Priyanjana Bengani, and Augustin Chaintreau. “Chasm in Hegemony: Explaining and Reproducing Disparities in Homophilous Networks.” *Proc. ACM Meas. Anal. Comput. Syst.*, **5**(2), June 2021.
- [ZJ24] Mingxia Zhao and Adele Lu Jia. “DAHGN: Degree-Aware Heterogeneous Graph Neural Network.” *Knowledge-Based Systems*, **285**:111355, 2024.
- [ZL21a] Yiliang Zhang and Qi Long. “Assessing Fairness in the Presence of Missing Data.” In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pp. 16007–16019. Curran Associates, Inc., 2021.
- [ZL21b] Yiliang Zhang and Qi Long. “Fairness in Missing Data Imputation.” *arXiv preprint arXiv:2110.12002*, 2021.
- [ZLM18] Brian Hu Zhang, Blake Lemoine, and Margaret Mitchell. “Mitigating Unwanted Biases with Adversarial Learning.” In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, AIES '18, p. 335–340, New York, NY, USA, 2018. Association for Computing Machinery.
- [ZLP23] Jiangqiang Zhu, Kai Li, Jinjia Peng, and Jing Qi. “Self-Supervised Graph Attention Collaborative Filtering for Recommendation.” *Electronics*, **12**(4), 2023.

- [ZLY22] Guixian Zhang, Rongjiao Liang, Zhongyi Yu, and Shichao Zhang. “Rumour Detection on Social Media with Long-Tail Strategy.” In *2022 International Joint Conference on Neural Networks (IJCNN)*, pp. 1–8, 2022.
- [ZMW22] An Zhang, Wenchang Ma, Xiang Wang, and Tat-Seng Chua. “Incorporating bias-aware margins into contrastive loss for collaborative filtering.” In *Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS ’22*, Red Hook, NY, USA, 2022. Curran Associates Inc.
- [Zop22] Markus Zopf. “1-wl expressiveness is (almost) all you need.” In *2022 International Joint Conference on Neural Networks (IJCNN)*, pp. 1–8. IEEE, 2022.
- [ZWY17] Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. “Men Also Like Shopping: Reducing Gender Bias Amplification using Corpus-level Constraints.” In Martha Palmer, Rebecca Hwa, and Sebastian Riedel, editors, *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pp. 2979–2989, Copenhagen, Denmark, sep 2017. Association for Computational Linguistics.
- [ZZH20] Tao Zhang, Tianqing Zhu, Mengde Han, Jing Li, Wanlei Zhou, and Philip S. Yu. “Fairness Constraints in Semi-supervised Learning.” *arXiv preprint arXiv:2009.06190*, 2020.
- [ZZW25] Wenbin Zhang, Shuigeng Zhou, Toby Walsh, and Jeremy C. Weiss. “Fairness amidst non-IID graph data: A literature review.” *AI Magazine*, **46**(1):e12212, 2025.
- [ZZY19] Rong Zhu, Kun Zhao, Hongxia Yang, Wei Lin, Chang Zhou, Baole Ai, Yong Li, and Jingren Zhou. “AliGraph: a comprehensive graph neural network platform.” *Proc. VLDB Endow.*, **12**(12):2094–2105, August 2019.

[ZZY24] Guixian Zhang, Shichao Zhang, and Guan Yuan. “Bayesian Graph Local Extrema Convolution with Long-Tail Strategy for Misinformation Detection.” *ACM Trans. Knowl. Discov. Data*, jan 2024.

APPENDIX A

Appendix for Chapter 2

A.1 Proofs

A.1.1 Proof of Lemma 1 (special case of total variation distance)

We assume $\mathbb{P}(Z)$ and $\mathbb{P}(X)$ are continuous probability distributions, but this proof can be easily altered for other kinds of probability distributions and the result still holds. Let $\mathbb{P}(S = Q) = p$. By definition, the total variation information:

$$I_{TV}(\mathbb{P}(Z); \mathbb{P}(S)) = d_{TV}(\mathbb{P}(Z, S) || \mathbb{P}(Z) \otimes \mathbb{P}(S)) = \sum_{s \in \{Q, R\}} \mathbb{P}(S = s) \int_{\mathcal{Y}} \frac{1}{2} \left| \frac{f_{Z,S}(z,s)}{f_Z(z)\mathbb{P}(S=s)} - 1 \right| f_Z(z) dz$$

By breaking up joint probabilities into conditional probabilities and factoring:

$$\begin{aligned} I_{TV}(\mathbb{P}(Z); \mathbb{P}(S)) &= \sum_{s \in \{Q, R\}} \mathbb{P}(S = s) \int_{\mathcal{Y}} \frac{1}{2} \left| \frac{f_{Z|S=s}(z)\mathbb{P}(S = s)}{f_Z(z)\mathbb{P}(S = s)} - 1 \right| f_Z(z) dz \\ &= \sum_{s \in \{Q, R\}} \mathbb{P}(S = s) \int_{\mathcal{Y}} \frac{1}{2} |f_{Z|S=s}(z) - f_Z(z)| dz \\ &= \frac{p}{2} \int_{\mathcal{Y}} |f_{Z|S=Q}(z) - f_Z(z)| dz + \int_{\mathcal{Y}} \frac{1-p}{2} |f_{Z|S=R}(z) - f_Z(z)| dz \\ &= \frac{1}{2} \left[p \int_{\mathcal{Y}} |f_{Z|S=Q}(z) - (pf_{Z|S=Q}(z) + (1-p)f_{Z|S=R}(z))| dz \right. \\ &\quad \left. + (1-p) \int_{\mathcal{Y}} |f_{Z|S=R}(z) - (pf_{Z|S=Q}(z) + (1-p)f_{Z|S=R}(z))| dz \right] \\ &= \frac{1}{2} \left[p \int_{\mathcal{Y}} |(1-p)f_{Z|S=Q}(z) - (1-p)f_{Z|S=R}(z)| dz \right. \\ &\quad \left. + (1-p) \int_{\mathcal{Y}} |pf_{Z|S=R}(z) - pf_{Z|S=Q}(z)| dz \right] \end{aligned}$$

$$\begin{aligned}
&= \frac{p(1-p)}{2} \left[\int_{\mathcal{Y}} |f_{Z|S=Q}(z) - f_{Z|S=R}(z)| \, dz + \int_{\mathcal{Y}} |f_{Z|S=R}(z) - f_{Z|S=Q}(z)| \, dz \right] \\
&= p(1-p) \int_{\mathcal{Y}} |f_{Z|S=Q}(z) - f_{Z|S=R}(z)| \, dz \\
&= 2p(1-p) d_{TV}(\mathbb{P}(Z|S=Q), \mathbb{P}(Z|S=R))
\end{aligned}$$

It can be similarly shown that $I_{TV}(\mathbb{P}(X); \mathbb{P}(S)) = 2p(1-p) d_{TV}(\mathbb{P}(X|S=Q), \mathbb{P}(X|S=R))$.

Now, because Z is conditionally independent of S given X , by the Data Processing Inequality, $I_{TV}(\mathbb{P}(Z); \mathbb{P}(S)) \leq I_{TV}(\mathbb{P}(X); \mathbb{P}(S))$. Hence:

$$\begin{aligned}
d_{TV}(\mathbb{P}(Z|S=Q), \mathbb{P}(Z|S=R)) &= \frac{1}{2p(1-p)} I_{TV}(\mathbb{P}(Z); \mathbb{P}(S)) \\
&\leq \frac{1}{2p(1-p)} I_{TV}(\mathbb{P}(X); \mathbb{P}(S)) \\
&= d_{TV}(\mathbb{P}(X|S=Q), \mathbb{P}(X|S=R))
\end{aligned}$$

Similarly:

$$d_{TV}(\mathbb{P}(Z'|S=Q), \mathbb{P}(Z'|S=R)) \leq d_{TV}(\mathbb{P}(X'|S=Q), \mathbb{P}(X'|S=R))$$

Differences in the supports of $\mathbb{P}(Z)$ and $\mathbb{P}(X)$ should not influence one's interpretation of the inequality. $d_{TV}(\cdot, \cdot)$ only requires that its two arguments have the same support. Because d_{TV} outputs the largest possible difference between the probabilities that the two distributions can assign to the same event, the inequality can be viewed as a comparison of the differences in assigned probabilities.

A.1.2 Proof of Theorem 1

Leveraging the Bretagnolle–Huber (BH) bound¹, we can upper bound d_{TV} in terms of the KL-divergence d_{KL} :

¹We use the BH bound rather than Pinsker's inequality because Pinsker's inequality becomes vacuous for KL-divergence > 2 [Can23].

$$d_{TV}(\mathbb{P}(X|S = Q), \mathbb{P}(X|S = R)) \leq \sqrt{1 - e^{-d_{KL}(\mathbb{P}(X|S=Q)||\mathbb{P}(X|S=R))}}$$

By Section 9 from [Duc16], $d_{KL}(\mathbb{P}(X|S = Q)||\mathbb{P}(X|S = R))$ admits a closed-form solution:

$$\begin{aligned} & d_{KL}(\mathbb{P}(X|S = Q)||\mathbb{P}(X|S = R)) \\ &= \frac{1}{2} \left(\log \frac{\det \Sigma_R}{\det \Sigma_Q} - d + \text{tr}(\Sigma_R^{-1} \Sigma_Q) + \|\mu_Q - \mu_R\|_{\Sigma_R^{-1}}^2 \right) \\ &\leq \frac{1}{2} \left(\log \frac{\det \Sigma_R}{\det \Sigma_Q} - d + \text{tr}(\Sigma_R^{-1} \Sigma_Q) + \lambda_{\max}(\Sigma_R^{-1}) \|\mu_Q - \mu_R\|_2^2 \right), \end{aligned}$$

where $\lambda_{\max}(\Sigma_R^{-1})$ is the maximum eigenvalue of Σ_R^{-1} . We note that $\lambda_{\max}(\Sigma_R^{-1}) = \frac{1}{\lambda_{\min}(\Sigma_R)} > 0$ (where $\lambda_{\min}(\Sigma_R)$ is the minimum eigenvalue of Σ_R) because Σ_R is positive semidefinite.

It is clear that $\|\mu_Q - \mu_R\|_\infty^2 = \max_{i \in [d]} |(\mu_Q)_i - (\mu_R)_i|^2 \leq \sum_{i \in [d]} |(\mu_Q)_i - (\mu_R)_i|^2 = \|\mu_Q - \mu_R\|_2^2$. Moreover, $\|\mu_Q - \mu_R\|_2^2 = \sum_{i \in [d]} |(\mu_Q)_i - (\mu_R)_i|^2 \leq \sum_{i \in [d]} \max_{j \in [d]} |(\mu_Q)_j - (\mu_R)_j|^2 = d \cdot \|\mu_Q - \mu_R\|_\infty^2$. Therefore, $\|\mu_Q - \mu_R\|_2^2 = C \cdot \|\mu_Q - \mu_R\|_\infty^2$, for $1 \leq C \leq d$.

Combining the previous observations and recognizing that $\|\mu_Q - \mu_R\|_\infty^2 = \mathcal{R}_D^2$:

$$d_{TV}(\mathbb{P}(X|S = Q), \mathbb{P}(X|S = R)) \leq \sqrt{1 - \sqrt{\frac{\det \Sigma_Q}{\det \Sigma_R} \cdot e^{-\frac{C \cdot \mathcal{R}_D^2}{\lambda_{\min}(\Sigma_R)} - \text{tr}(\Sigma_R^{-1} \Sigma_Q) + d}}}}$$

Now, by Lemma 2.7 from [Ste21]:

$$\begin{aligned} \|\mu_Q - \mu_R\|_2^2 &\leq 4 \cdot \max\{\lambda_{\max}(\Sigma_Q), \lambda_{\max}(\Sigma_R)\} \left(\frac{d_{TV}(\mathbb{P}(X|S = Q), \mathbb{P}(X|S = R))}{1 - d_{TV}(\mathbb{P}(X|S = Q), \mathbb{P}(X|S = R))} \right) \\ &\leq \frac{4 \cdot \max\{\lambda_{\max}(\Sigma_Q), \lambda_{\max}(\Sigma_R)\}}{\frac{1}{d_{TV}(\mathbb{P}(X|S=Q), \mathbb{P}(X|S=R))} - 1} \end{aligned}$$

Using $\|\mu_Q - \mu_R\|_2^2 = C \cdot \|\mu_Q - \mu_R\|_\infty^2 = C \cdot \mathcal{R}_D^2$, we can derive:

$$\frac{1}{\frac{4 \cdot \max\{\lambda_{\max}(\Sigma_Q), \lambda_{\max}(\Sigma_R)\}}{C \cdot \mathcal{R}_D^2} + 1} \leq d_{TV}(\mathbb{P}(X|S = Q), \mathbb{P}(X|S = R))$$

Similarly:

$$d_{TV}(\mathbb{P}(X'|S = Q), \mathbb{P}(X'|S = R)) \in \left[\frac{1}{\frac{4 \cdot \max\{\lambda_{\max}(\Sigma'_Q), \lambda_{\max}(\Sigma'_R)\}}{C' \cdot \mathcal{R}_{D'}^2} + 1}, \sqrt{1 - \sqrt{\frac{\det \Sigma'_Q}{\det \Sigma'_R} \cdot e^{-\frac{C' \cdot \mathcal{R}_{D'}^2}{\lambda_{\min}(\Sigma'_R)} - \text{tr}(\Sigma'^{-1} \Sigma'_Q) + d}}}} \right]$$

Then, the theorem is proved by application of Lemma 1.

A.1.3 Example mean aggregation imputation algorithms

Global Mean This method sets the unknown features to the uniform mean of all the known features. To achieve this, we can choose $M := I_N$, $T := \begin{bmatrix} I_{|K|} & 0 \\ \frac{1}{|K|} \mathbb{1}_{|U| \times |K|} & 0 \end{bmatrix}$ (where $\mathbb{1}$ is the all-ones matrix), $\beta := 0$, $X_K^{(0)} := X_K$, and $X_U^{(0)} := 0$. We only need to complete one iteration.

Neighbor Mean This method sets the unknown features to the degree-weighted mean of the known features for neighboring nodes. We can choose $M := I_N$, $T := D^{-1}A$, $\beta := 0$, $X_K^{(0)} := X_K$, and $X_U^{(0)} := 0$. We only need to complete one iteration.

Feature Propagation This method proposed by [RKG22] predicts the unknown features to minimize the Dirichlet energy of the graph while preserving the known feature values. [RKG22] shows that this is equivalent to iteratively computing until convergence:

$$X_K^{(t+1)} := X_K^{(t)}$$

$$X_U^{(t+1)} := (D_U^{-\frac{1}{2}} A_{UK} D_K^{-\frac{1}{2}}) X_K^{(t)} + (D_U^{-\frac{1}{2}} A_{UU} D_U^{-\frac{1}{2}}) X_U^{(t)}$$

Multiplying both sides by $D_U^{-\frac{1}{2}}$, we can re-express the second update rule as:

$$D_U^{-\frac{1}{2}} X_U^{(t+1)} = (D_U^{-1} A_{UK}) (D_K^{-\frac{1}{2}} X_K^{(t)}) + (D_U^{-1} A_{UU}) (D_U^{-\frac{1}{2}} X_U^{(t)})$$

Therefore, to achieve Feature Propagation, we can choose $M := D^{-\frac{1}{2}}$, $T := D^{-1}A$, $\beta := 0$, and $X_K^{(0)} := X_K$. Per [RKG22], we can choose $X_U^{(0)}$ arbitrarily, and we need to iterate till convergence.

Graph Regularization This method inspired by [ZBL04] predicts the unknown features via a smoothness constraint and a fitting constraint for the known features. [ZBL04] shows

that this is equivalent to iteratively computing until convergence:

$$X_K^{(t+1)} := \beta(D^{-\frac{1}{2}}AD^{-\frac{1}{2}}X^{(t)})_K + (1 - \beta)X_K$$

$$X_U^{(t+1)} := (D^{-\frac{1}{2}}AD^{-\frac{1}{2}}X^{(t)})_U,$$

where the hyperparameter $\beta \in (0, 1]$. Therefore, similar to Feature Propagation, to achieve Graph Regularization, we can choose $M := D^{-\frac{1}{2}}$, $T := D^{-1}A$, and $X_K^{(0)} := X_K$. Per [ZBL04], we can choose $X_U^{(0)}$ arbitrarily, and we need to iterate till convergence.

A.1.4 Proof of Theorem 2

The following proof is partially inspired by the proof of Theorem 4.1 in [LWZ21]. Fix t to be an arbitrary iteration of feature imputation. Recall we use the following iterative update rule to impute features:

$$\tilde{X}^{(t+1)} := \begin{bmatrix} \beta I_{|K|} & 0 \\ 0 & I_{|U|} \end{bmatrix} T \tilde{X}^{(t)} + \begin{bmatrix} (1 - \beta) I_{|K|} & 0 \\ 0 & 0 \end{bmatrix} \tilde{X}$$

For a node $q \in Q \cap U$, after one iteration of feature imputation:

$$\tilde{X}_q^{(t+1)} := \sum_{s \in Q} T_{qs} \tilde{X}_s^{(t)} + \sum_{s \in R} T_{qs} \tilde{X}_s^{(t)}$$

Similarly, for a node $r \in R \cap U$, after one iteration of feature imputation:

$$\tilde{X}_r^{(t+1)} := \sum_{s \in Q} T_{rs} \tilde{X}_s^{(t)} + \sum_{s \in R} T_{rs} \tilde{X}_s^{(t)}$$

In contrast, for a node $q \in Q \cap K$, after one iteration of feature imputation:

$$\tilde{X}_q^{(t+1)} := \beta \left(\sum_{s \in Q} T_{qs} \tilde{X}_s^{(t)} + \sum_{s \in R} T_{qs} \tilde{X}_s^{(t)} \right) + (1 - \beta) \tilde{X}_q^{(t)}$$

Similarly, for a node $r \in R \cap K$, after one iteration of feature imputation:

$$\tilde{X}_r^{(t+1)} := \beta \left(\sum_{s \in Q} T_{rs} \tilde{X}_s^{(t)} + \sum_{s \in R} T_{rs} \tilde{X}_s^{(t)} \right) + (1 - \beta) \tilde{X}_r^{(t)}$$

We say $v \in [\mu \pm \sigma] \iff \mu - \sigma \preceq v \preceq \mu + \sigma$. Then, for a node $q \in Q \cap U$, by the right-stochastic nature of T :

$$\begin{aligned} \tilde{X}_q^{(t+1)} &\in \left[\left(\sum_{s \in Q} T_{qs} \tilde{\mu}_Q^{(t)} + \sum_{s \in R} T_{qs} \tilde{\mu}_R^{(t)} \right) \pm \tilde{\sigma}^{(t)} \right] \\ &\in \left[\left(\tilde{\mu}_Q^{(t)} + \sum_{s \in R} T_{qs} \left(\tilde{\mu}_R^{(t)} - \tilde{\mu}_Q^{(t)} \right) \right) \pm \tilde{\sigma}^{(t)} \right] \end{aligned}$$

Similarly, for a node $r \in R \cap U$:

$$\begin{aligned} \tilde{X}_r^{(t+1)} &\in \left[\left(\sum_{s \in Q} T_{rs} \tilde{\mu}_Q^{(t)} + \sum_{s \in R} T_{rs} \tilde{\mu}_R^{(t)} \right) \pm \tilde{\sigma}^{(t)} \right] \\ &\in \left[\left(\tilde{\mu}_R^{(t)} + \sum_{s \in Q} T_{rs} \left(\tilde{\mu}_Q^{(t)} - \tilde{\mu}_R^{(t)} \right) \right) \pm \tilde{\sigma}^{(t)} \right] \end{aligned}$$

In contrast, for a node $q \in Q \cap K$:

$$\begin{aligned} \tilde{X}_q^{(t+1)} &\in \left[\left(\beta \left(\sum_{s \in Q} T_{qs} \tilde{\mu}_Q^{(t)} + \sum_{s \in R} T_{qs} \tilde{\mu}_R^{(t)} \right) + (1 - \beta) \tilde{\mu}_Q^{(t)} \right) \pm \tilde{\sigma}^{(t)} \right] \\ &\in \left[\left(\tilde{\mu}_Q^{(t)} + \beta \sum_{s \in R} T_{qs} \left(\tilde{\mu}_R^{(t)} - \tilde{\mu}_Q^{(t)} \right) \right) \pm \tilde{\sigma}^{(t)} \right] \end{aligned}$$

Similarly, for a node $r \in R \cap K$:

$$\begin{aligned} \tilde{X}_r^{(t+1)} &\in \left[\left(\beta \left(\sum_{s \in Q} T_{rs} \tilde{\mu}_Q^{(t)} + \sum_{s \in R} T_{rs} \tilde{\mu}_R^{(t)} \right) + (1 - \beta) \tilde{\mu}_R^{(t)} \right) \pm \tilde{\sigma}^{(t)} \right] \\ &\in \left[\left(\tilde{\mu}_R^{(t)} + \beta \sum_{s \in Q} T_{rs} \left(\tilde{\mu}_Q^{(t)} - \tilde{\mu}_R^{(t)} \right) \right) \pm \tilde{\sigma}^{(t)} \right] \end{aligned}$$

By the Law of Total Expectation:

$$\begin{aligned}
\mathbb{E}_{q \sim Q}[\tilde{X}_q^{(t+1)}] &= \mathbb{P}(q \in U | q \in Q) \mathbb{E}_{q \sim Q \cap U}[\tilde{X}_q^{(t+1)}] + \mathbb{P}(q \in K | q \in Q) \mathbb{E}_{q \sim Q \cap K}[\tilde{X}_q^{(t+1)}] \\
&\in \left[\left(\frac{1}{|Q|} \left(\sum_{q \in Q \cap U} \tilde{\mu}_Q^{(t)} + \sum_{s \in R} T_{qs} \left(\tilde{\mu}_R^{(t)} - \tilde{\mu}_Q^{(t)} \right) \right) \right. \right. \\
&\quad \left. \left. + \frac{1}{|Q|} \left(\sum_{q \in Q \cap K} \tilde{\mu}_Q^{(t)} + \beta \sum_{s \in R} T_{qs} \left(\tilde{\mu}_R^{(t)} - \tilde{\mu}_Q^{(t)} \right) \right) \right) \pm \tilde{\sigma}^{(t)} \right] \\
&\in \left[\left(\tilde{\mu}_Q^{(t)} + \frac{1}{|Q|} \sum_{q \in Q \cap U} \sum_{s \in R} T_{qs} \left(\tilde{\mu}_R^{(t)} - \tilde{\mu}_Q^{(t)} \right) \right. \right. \\
&\quad \left. \left. + \frac{\beta}{|Q|} \sum_{q \in Q \cap K} \sum_{s \in R} T_{qs} \left(\tilde{\mu}_R^{(t)} - \tilde{\mu}_Q^{(t)} \right) \right) \pm \tilde{\sigma}^{(t)} \right]
\end{aligned}$$

Similarly, $\mathbb{E}_{r \sim R}[\tilde{X}_r^{(t+1)}]$:

$$\in \left[\left(\tilde{\mu}_R^{(t)} + \frac{1}{|R|} \sum_{r \in R \cap U} \sum_{s \in Q} T_{rs} \left(\tilde{\mu}_Q^{(t)} - \tilde{\mu}_R^{(t)} \right) + \frac{\beta}{|R|} \sum_{r \in R \cap K} \sum_{s \in Q} T_{rs} \left(\tilde{\mu}_Q^{(t)} - \tilde{\mu}_R^{(t)} \right) \right) \pm \tilde{\sigma}^{(t)} \right]$$

Thus, the gap in expectation of the features of the nodes in Q and R after one iteration of feature imputation is:

$$\begin{aligned}
\mathbb{E}_{q \sim Q}[\tilde{X}_q^{(t+1)}] - \mathbb{E}_{r \sim R}[\tilde{X}_r^{(t+1)}] &\in \left[\left(1 - \left(\frac{1}{|Q|} \sum_{q \in Q \cap U} \sum_{s \in R} T_{qs} + \frac{1}{|R|} \sum_{r \in R \cap U} \sum_{s \in Q} T_{rs} \right) \right. \right. \\
&\quad \left. \left. - \beta \left(\frac{1}{|Q|} \sum_{q \in Q \cap K} \sum_{s \in R} T_{qs} + \frac{1}{|R|} \sum_{r \in R \cap K} \sum_{s \in Q} T_{rs} \right) \right) \cdot \left(\tilde{\mu}_Q^{(t)} - \tilde{\mu}_R^{(t)} \right) \pm 2\tilde{\sigma}^{(t)} \right]
\end{aligned}$$

Define the contraction coefficient:

$$\alpha := \left| 1 - \frac{T_{R \rightarrow Q \cap U} + \beta T_{R \rightarrow Q \cap K}}{|Q|} - \frac{T_{Q \rightarrow R \cap U} + \beta T_{Q \rightarrow R \cap K}}{|R|} \right|$$

Because $0 \leq \frac{T_{R \rightarrow Q \cap U} + \beta T_{R \rightarrow Q \cap K}}{|Q|} \leq \frac{T_{R \rightarrow Q \cap U} + T_{R \rightarrow Q \cap K}}{|Q|} = \frac{T_{R \rightarrow Q}}{|Q|} < \frac{T_{V \rightarrow Q}}{|Q|} = 1$, and similarly, $0 \leq \frac{T_{Q \rightarrow R \cap U} + \beta T_{Q \rightarrow R \cap K}}{|R|} < 1$, it must be that $0 \leq \alpha \leq 1$.

Then:

$$\begin{aligned}
\max\{\alpha |\mathbb{E}_{q \sim Q}[\tilde{X}_q^{(t)}] - \mathbb{E}_{r \sim R}[\tilde{X}_r^{(t)}]| - 2\tilde{\sigma}^{(t)}, 0\} &\leq |\mathbb{E}_{q \sim Q}[\tilde{X}_q^{(t+1)}] - \mathbb{E}_{r \sim R}[\tilde{X}_r^{(t+1)}]| \\
|\mathbb{E}_{q \sim Q}[\tilde{X}_q^{(t+1)}] - \mathbb{E}_{r \sim R}[\tilde{X}_r^{(t+1)}]| &\leq \alpha |\mathbb{E}_{q \sim Q}[\tilde{X}_q^{(t)}] - \mathbb{E}_{r \sim R}[\tilde{X}_r^{(t)}]| + 2\tilde{\sigma}^{(t)}
\end{aligned}$$

$$\max\{\alpha\tilde{\mathcal{R}}^{(t)} - 2\tilde{\sigma}^{(t)}, 0\} \leq \tilde{\mathcal{R}}^{(t+1)} \leq \alpha\tilde{\mathcal{R}}^{(t)} + 2\tilde{\sigma}^{(t)}$$

Inductively, the discrimination risk $\tilde{\mathcal{R}}^{(t)}$ after t iterations of feature imputation is bounded by:

$$\max\left\{\alpha^t\tilde{\mathcal{R}}^{(0)} - 2\left(\sum_{j=0}^{t-1}\alpha^j\tilde{\sigma}^{(j)}\right), 0\right\} \leq \tilde{\mathcal{R}}^{(t)} \leq \alpha^t\tilde{\mathcal{R}}^{(0)} + 2\left(\sum_{j=0}^{t-1}\alpha^j\tilde{\sigma}^{(j)}\right)$$

$\forall v \in V$, $\tilde{X}_v^{(t+1)}$ is a convex combination of $\bigcup_{u \in V}\{\tilde{X}_u^{(t)}\}$. This is because each row of T and $\beta T + (1 - \beta)I_{|K|}$ contains nonnegative entries that sum to 1. Therefore, $\tilde{X}_v^{(t+1)}$ must be in the (closed) convex hull of $\bigcup_{u \in V}\{\tilde{X}_u^{(t)}\}$. Thus, $\bigcup_{u \in V}\{\tilde{X}_u^{(t)}\}$ inductively must be contained within the (closed) convex hull of $\bigcup_{u \in V}\{\tilde{X}_u^{(0)}\}$, which has extreme points $\subseteq \bigcup_{u \in V}\{\tilde{X}_u^{(0)}\}$. Consequently, $\forall t \in [0, \infty)$, $\tilde{\sigma}^{(t)} \leq \tilde{\sigma}^{(0)}$.

Hence:

$$\max\left\{\alpha^t\tilde{\mathcal{R}}^{(0)} - 2\left(\sum_{j=0}^{t-1}\alpha^j\right)\tilde{\sigma}^{(0)}, 0\right\} \leq \tilde{\mathcal{R}}^{(t)} \leq \alpha^t\tilde{\mathcal{R}}^{(0)} + 2\left(\sum_{j=0}^{t-1}\alpha^j\right)\tilde{\sigma}^{(0)}$$

If $\alpha < 1$:

$$\max\left\{\alpha^t\tilde{\mathcal{R}}^{(0)} - 2\left(\frac{1 - \alpha^t}{1 - \alpha}\right)\tilde{\sigma}^{(0)}, 0\right\} \leq \tilde{\mathcal{R}}^{(t)} \leq \alpha^t\tilde{\mathcal{R}}^{(0)} + 2\left(\frac{1 - \alpha^t}{1 - \alpha}\right)\tilde{\sigma}^{(0)}$$

Moreover, upon convergence:

$$0 \leq \lim_{t \rightarrow \infty} \tilde{\mathcal{R}}^{(t)} \leq \frac{2\tilde{\sigma}^{(0)}}{1 - \alpha}$$

While it appears that a large initial maximal deviation in feature values within a group may harm fairness, a large initial deviation does not necessarily entail diversity. For example, suppose a few nodes in a group have a low initial feature value but many more nodes have a much higher initial feature value (i.e., large initial difference without diversity). Then, after mean aggregation, the feature values for all the nodes in the group may be higher on average

than they were initially, and more distinct on average from the node feature values in the other group. This would contribute to a higher discrimination risk.

A.1.5 Extending Theorem 2

We can extend Theorem 2 to the case the number of features $d > 1$. By Theorem 1, the modified discrimination risk at iteration t (including all features) is:

$$\max \left\{ \min_{i \in [d]} \alpha_i^t \tilde{\mathcal{R}}_i^{(0)} - 2 \left(\sum_{j=0}^{t-1} \alpha_i^j \right) \tilde{\sigma}_i^{(0)}, 0 \right\} \leq \max_{i \in [d]} \tilde{\mathcal{R}}_i^{(t)} \leq \max_{i \in [d]} \alpha_i^t \tilde{\mathcal{R}}_i^{(0)} + 2 \left(\sum_{j=0}^{t-1} \alpha_i^j \right) \tilde{\sigma}_i^{(0)}$$

Moreover, assuming $\forall i \in [d], \alpha_i < 1$, upon convergence, the discrimination risk is:

$$\max_{i \in [d]} \lim_{t \rightarrow \infty} \tilde{\mathcal{R}}_i^{(t)} \leq \max_{i \in [d]} \frac{2\tilde{\sigma}_i^{(0)}}{1-\alpha_i}.$$

A.1.6 Proof of Theorem 3

We want to constrain the discrimination risk of mean aggregation feature imputation to be at most ϵ . To this end, we can modify mean aggregation feature imputation to update $X_U^{(t+1)} := P_W Z_U^{(t)} + P_B$ such that $X^{(t+1)}$ has a discrimination risk of at most ϵ for all t . $|\mathbb{E}_{q \sim Q}[X_q^{(t+1)}] - \mathbb{E}_{r \sim R}[X_r^{(t+1)}]| = |\frac{1}{|Q|} \sum_{q \in Q \cap K} X_q + \frac{1}{|Q|} \sum_{q \in Q \cap U} Z_q^{(t)} - (\frac{1}{|R|} \sum_{r \in R \cap K} X_r + \frac{1}{|R|} \sum_{r \in R \cap U} Z_r^{(t)})|$. Hence, we have a closed convex polytope wherein unknown feature values yield discrimination risk of at most ϵ :

$$\mathcal{R}_K - \epsilon \leq \frac{1}{|Q|} \sum_{q \in Q \cap U} Z_q^{(t)} - \frac{1}{|R|} \sum_{r \in R \cap U} Z_r^{(t)} = c^T Z_U^{(t)} \leq \mathcal{R}_K + \epsilon$$

If $\mathcal{R}_K - \epsilon \leq c^T Z_U^{(t)} \leq \mathcal{R}_K + \epsilon$, then $P_W = I_{|U|}$ and $P_B = 0$. Otherwise, we must project onto the closer of the two boundaries of the polytope. In this case, $P_W = I_{|U|} - \frac{cc^T}{c^T c}$ and

$$P_B = \frac{cc^T}{c^T c} \begin{cases} \mathcal{R}_K - \epsilon, & c^T Z_U^{(t)} < \mathcal{R}_K - \epsilon \\ \mathcal{R}_K + \epsilon, & c^T Z_U^{(t)} > \mathcal{R}_K + \epsilon \end{cases}.$$

The affine projection we perform at each step is closed and convex. Furthermore, ℓ is $\lambda_{max}(\Delta_{UU})$ -smooth for the Euclidean norm (where λ_{max} is the maximum eigenvalue) because

for $x_1, x_2 \in \mathbb{R}^{|U|}$:

$$\begin{aligned}
\|\nabla\ell(x_1) - \nabla\ell(x_2)\|_2 &= \|(\Delta_{UU}x_1 + \Delta_{UK}X_K) - (\Delta_{UU}x_2 + \Delta_{UK}X_K)\|_2 \\
&= \sqrt{(x_1 - x_2)^T \Delta_{UU}^2 (x_1 - x_2)} \\
&\leq \sqrt{\lambda_{max}^2(\Delta_{UU}) \|x_1 - x_2\|_2^2} \\
&= \lambda_{max}(\Delta_{UU}) \|x_1 - x_2\|_2
\end{aligned}$$

In the case of Feature Propagation, $\lambda_{max}(\Delta_{UU}) < 1$ due to properties of the symmetric normalized Laplacian [RKG22].

Additionally, for $m \geq 0$, when $m \leq \lambda_{min}(\Delta_{UU})$, $\ell(x) - \frac{m}{2}x^T x$ is convex because:

$$\begin{aligned}
\ell(x) - \frac{m}{2}x^T x &= \frac{1}{2}x^T \Delta_{UU}x + X_K^T \Delta_{KU}x + \frac{1}{2}X_K^T \Delta_{KK}X_K - \frac{m}{2}x^T x \\
&= \frac{1}{2}x^T (\Delta_{UU} - mI)x + X_K^T \Delta_{KU}x + \frac{1}{2}X_K^T \Delta_{KK}X_K
\end{aligned}$$

This expression is convex if and only if its Hessian $\Delta_{UU} - mI$ has nonnegative eigenvalues.

Therefore, m can be at most $\lambda_{min}(\Delta_{UU})$.

Then, by [Bec21] and [Van21]:

1. a unique optimal (with respect to ℓ) feasible solution X_U^* exists
2. for fixed step size $\gamma = \frac{1}{\lambda_{max}(\Delta_{UU})}$, ϵ -fair imputation converges as $\|X_U^{(t)} - X_U^*\|_2^2 \leq \left(1 - \frac{\lambda_{min}(\Delta_{UU})}{\lambda_{max}(\Delta_{UU})}\right)^t \|X_U^{(0)} - X_U^*\|_2^2$
3. for fixed step size $\gamma \leq \frac{1}{\lambda_{max}(\Delta_{UU})}$, ϵ -fair imputation converges to X_U^*

A.1.7 Theorem 4

We have a solution when $\beta > 0$ (i.e., when the known node feature values do not remain fixed). We can view the update of $X^{(t+1)} := \begin{bmatrix} \beta I_{|K|} & 0 \\ 0 & I_{|U|} \end{bmatrix} M^{-1} T M X^{(t)} + \begin{bmatrix} (1 - \beta) I_{|K|} & 0 \\ 0 & 0 \end{bmatrix} X$ as an iteration of gradient descent (with step size $\gamma = 1$) for the objective function $\ell(x) = \frac{1}{2}x^T \Delta x + \frac{1}{2}\left(\frac{1-\beta}{\beta}\right)\|x_K - X_K\|_2^2$ [ZBL04].

Theorem 4 (ϵ -Fair Imputation, $\beta > 0$) Vanilla mean aggregation feature imputation updates $X^{(t+1)} := \begin{bmatrix} \beta I_{|K|} & 0 \\ 0 & I_{|U|} \end{bmatrix} (I_N - \Delta)X^{(t)} + \begin{bmatrix} (1-\beta)I_{|K|} & 0 \\ 0 & 0 \end{bmatrix} X = Z^{(t)}$. Let ϵ -fair mean aggregation feature imputation instead update $X^{(t+1)} := P_W Z^{(t)} + P_B$, where:

$$P_W = \begin{cases} I_N, & -\epsilon \leq c^T Z^{(t)} \leq \epsilon \\ I_N - \frac{cc^T}{c^T c}, & \text{otherwise} \end{cases}, P_B = \frac{cc^T}{c^T c} \begin{cases} -\epsilon, & c^T Z^{(t)} < -\epsilon \\ \epsilon, & c^T Z^{(t)} > \epsilon \\ 0, & \text{otherwise} \end{cases}$$

$$c \in \mathbb{R}^N, c^T Z^{(t)} = \frac{1}{|Q|} \sum_{q \in Q} Z_q^{(t)} - \frac{1}{|R|} \sum_{r \in R} Z_r^{(t)}$$

Then, assuming $0 \leq \lambda_{\min}(\Delta) + \frac{1-\beta}{\beta} \leq \lambda_{\max}(\Delta) + \frac{1-\beta}{\beta} < 1$: 1) a unique optimal (with respect to ℓ) feasible solution X^* exists; 2) for fixed step size $\gamma = \frac{1}{\lambda_{\max}(\Delta) + \frac{1-\beta}{\beta}}$, ϵ -fair imputation converges as $\|X^{(t)} - X^*\|_2^2 \leq \left(1 - \frac{\lambda_{\min}(\Delta) + \frac{1-\beta}{\beta}}{\lambda_{\max}(\Delta) + \frac{1-\beta}{\beta}}\right)^t \|X^{(0)} - X^*\|_2^2$; 3) for fixed step size $\gamma \leq \frac{1}{\lambda_{\max}(\Delta) + \frac{1-\beta}{\beta}}$, ϵ -fair imputation converges to X^* .

Proof of Theorem 4 We want to constrain the discrimination risk of mean aggregation feature imputation to be at most ϵ . To this end, we can modify mean aggregation feature imputation to update $X^{(t+1)} := P_W Z^{(t)} + P_B$ such that $X^{(t+1)}$ has a discrimination risk of at most ϵ for all t .

$|\mathbb{E}_{q \sim Q}[X_q^{(t+1)}] - \mathbb{E}_{r \sim R}[X_r^{(t+1)}]| = \frac{1}{|Q|} \sum_{q \in Q} Z_q^{(t)} - \frac{1}{|R|} \sum_{r \in R} Z_r^{(t)}$. Hence, we have a closed convex polytope wherein feature values have discrimination risk of at most ϵ :

$$-\epsilon \leq \frac{1}{|Q|} \sum_{q \in Q} Z_q^{(t)} - \frac{1}{|R|} \sum_{r \in R} Z_r^{(t)} = c^T Z^{(t)} \leq \epsilon$$

If $-\epsilon \leq c^T Z^{(t)} \leq \epsilon$, then $P_W = I_N$ and $P_B = 0$. Otherwise, we must project onto the closer of the two boundaries of the polytope. In this case, $P_W = I_N - \frac{cc^T}{c^T c}$ and $P_B = \frac{cc^T}{c^T c} \begin{cases} -\epsilon, & c^T Z^{(t)} < -\epsilon \\ \epsilon, & c^T Z^{(t)} > \epsilon \end{cases}$.

The affine projection we perform at each step is closed and convex. Furthermore, ℓ is $(\lambda_{\max}(\Delta) + \frac{1-\beta}{\beta})$ -smooth for the Euclidean norm because for $x_1, x_2 \in \mathbb{R}^N$:

$$\begin{aligned} \|\nabla\ell(x_1) - \nabla\ell(x_2)\|_2 &= \left\| \left(\Delta x_1 + \frac{1-\beta}{\beta} ((x_1)_K - X_K) \right) - \left(\Delta x_2 + \frac{1-\beta}{\beta} ((x_2)_K - X_K) \right) \right\|_2 \\ &\leq \sqrt{(x_1 - x_2)^T \Delta^2 (x_1 - x_2)} + \frac{1-\beta}{\beta} \sqrt{(x_1 - x_2)^T \begin{pmatrix} I_{|K|} & 0 \\ 0 & 0 \end{pmatrix}^2 (x_1 - x_2)} \\ &\leq \left(\lambda_{\max}(\Delta) + \frac{1-\beta}{\beta} \right) \|x_1 - x_2\|_2 \end{aligned}$$

Additionally, for $m \geq 0$, when $m \leq \lambda_{\min}(\Delta) + \frac{1-\beta}{\beta}$, $\ell(x) - \frac{m}{2}x^T x$ is convex because:

$$\begin{aligned} \ell(x) - \frac{m}{2}x^T x &= \frac{1}{2}x^T \Delta x + \frac{1}{2} \left(\frac{1-\beta}{\beta} \right) \|x_K - X_K\|_2^2 - \frac{m}{2}x^T x \\ &= \frac{1}{2}x^T (\Delta - mI)x + \frac{1}{2} \left(\frac{1-\beta}{\beta} \right) \|x_K - X_K\|_2^2 \end{aligned}$$

This expression is convex if and only if its Hessian $\Delta - mI + \frac{1-\beta}{\beta} \begin{bmatrix} I_{|K|} & 0 \\ 0 & 0 \end{bmatrix}$ has nonnegative eigenvalues. Therefore, m can be at most $\lambda_{\min}(\Delta) + \frac{1-\beta}{\beta}$.

Then, by [Bec21] and [Van21]:

1. a unique optimal (with respect to ℓ) feasible solution X^* exists
2. for fixed step size $\gamma = \frac{1}{\lambda_{\max}(\Delta) + \frac{1-\beta}{\beta}}$, ϵ -fair imputation converges as $\|X^{(t)} - X^*\|_2^2 \leq \left(1 - \frac{\lambda_{\min}(\Delta) + \frac{1-\beta}{\beta}}{\lambda_{\max}(\Delta) + \frac{1-\beta}{\beta}} \right)^t \|X^{(0)} - X^*\|_2^2$
3. for fixed step size $\gamma \leq \frac{1}{\lambda_{\max}(\Delta) + \frac{1-\beta}{\beta}}$, ϵ -fair imputation converges to X^*

A.2 Additional experimental results

A.2.1 Datasets

SBM synthetic datasets Each network has 500 train nodes, 250 validation nodes, and 250 test nodes, split uniformly at random. Each node has a 10-dimensional feature vector sampled as described in the documentation². All edges have a weight of 1.

Real-world datasets In the `Credit defaulter` dataset, each node has 13 features (e.g., education, credit history, etc.), with an average degree of 95.79 ± 85.88 [ALZ21]. In the `German credit` dataset, each node has 27 features (e.g., loan amount, account-related features, etc.), with an average degree of 44.48 ± 26.51 . For both datasets, we use a 50/25/25 train/validation/test split, with each split comprising an equal portion of each label, and we do not include group membership as a feature.

A.2.2 Imputation algorithms

We run GM and NM for 1 iteration each, and FP and GR for 40 iterations. We adapted the code for data utilities, Feature Propagation, and model training from [RKG22]³. We state all changes in this chapter.

A.2.3 Models and training

For *mlp* and *gcn*, we use a hidden dimension of 64. We train all models full-batch using the Adam optimizer with a learning rate of 0.005 and Dropout rate of 0.5 [KB15, SHK14]. We also use early stopping with a patience of 200 epochs, i.e., we stop training when the

²https://pytorch-geometric.readthedocs.io/en/latest/modules/datasets.html#torch_geometric.datasets.StochasticBlockModelDataset

³<https://github.com/twitter-research/feature-propagation>

best validation accuracy has not changed for 200 epochs, and train for a maximum of 10000 epochs. We do not do any hyperparameter tuning. We implement and train all models using PyTorch and PyTorch Geometric [PGM19, FL19]. We train all models on a single `tesla v100-sxm2-16gb` GPU on an internal cluster.

A.2.4 Performance evaluation

To evaluate imputed features for SBM, since we don't have labels, we employ relative reconstruction error **RE** (calculated as $\|X_{true} - X_{pred}\|_2 / \|X_{true}\|_2$, where X_{true} and X_{pred} are the ground-truth and imputed features, respectively [RKG22]). A lower reconstruction error is better, and we would like regular mean aggregation imputation and its ϵ -fair counterparts to have comparable reconstruction errors. To measure performance on the real-world datasets, we consider the test accuracy (**Acc**) of models applied to the imputed data. A higher test accuracy is preferable, and we again would like comparable accuracies for regular and ϵ -fair imputation.

To evaluate group fairness, we compute the discrimination risk (**DR**) of the imputed data. A lower discrimination risk is preferable. For the SBM synthetic datasets, we also measure how much information the imputed features contain about group membership. We do this by training the models on the imputed data to predict group membership and calculate the test accuracy of the models on identifying group membership (which we refer to as **MI**) [ZLM18, PP20]. (We note that this setting may violate our theoretical assumptions in 2.3 that the association of group membership with model predictions can be fully explained by the node features.) The models may be conceptualized as adversaries attempting to recover group membership from the imputed features. Thus, we would like **MI** to be closer to 0.5 (i.e., the imputed features contain no information about group membership). We do not compute **MI** for the real-world datasets, as inferring group membership or identity from real-world data is invasive, invalid, and can be weaponized against marginalized communities (e.g., to find and incarcerate LGBTQIA+ individuals). To evaluate group fairness for the

real-world datasets, we use the test statistical parity (**SP**) of the models, defined as $|\mathbb{P}(Z = 1|S = Q) - \mathbb{P}(Z = 1|S = R)|$ (disparity in positive outcome rate for the groups) [DHP12], and test equal opportunity (**EO**), defined as $|\mathbb{P}(Z = 1|S = Q, Y = 1) - \mathbb{P}(Z = 1|S = R, Y = 1)|$ (disparity in accuracy of predicting positive outcome for the groups) [HPS16].

A.2.5 Contraction coefficient

As we analyzed, Figures A.1 to A.11 show that, for SBM: 1) a low unknown feature rate for both groups or disparate unknown feature rates across the groups increases α and the discrimination risk (Figures A.1, A.4, A.7, A.10); 2) group size alone does not affect α or the discrimination risk (Figures A.2, A.5, A.8, A.11); 3) a lower inter-link to intra-link ratio increases α and the discrimination risk (Figures A.3, A.6, A.9).

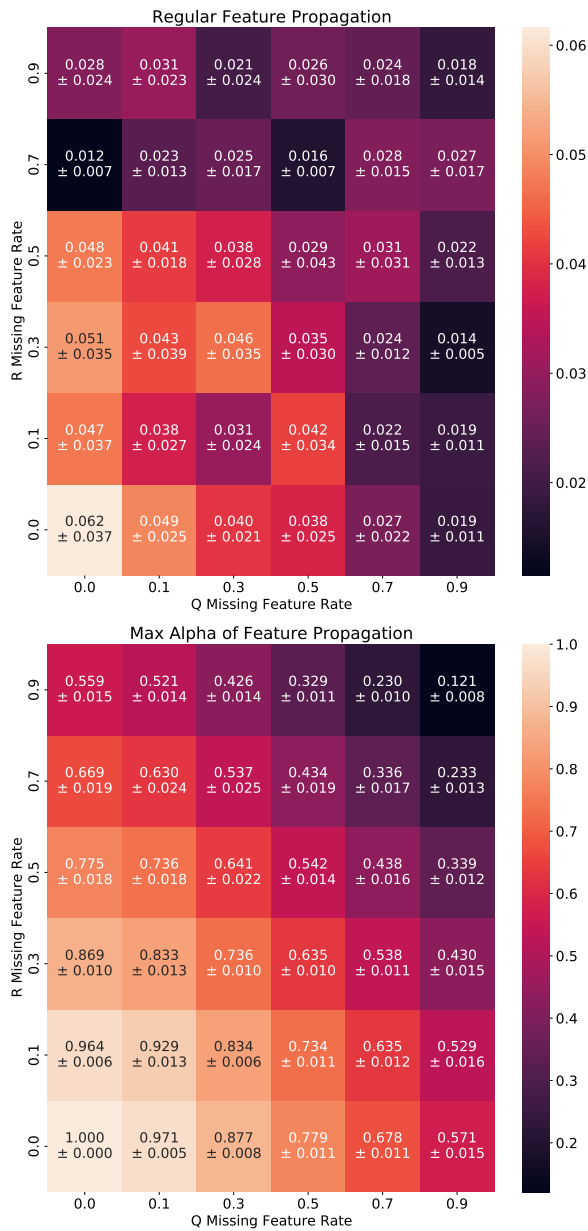


Figure A.1: Heatmap of discrimination risks and maximum α (over all channels) of Feature Propagation for 36 combinations of unknown feature rates for each group in SBM. We use 0.5 relative group sizes and 0.5 inter- and intra-link rates.

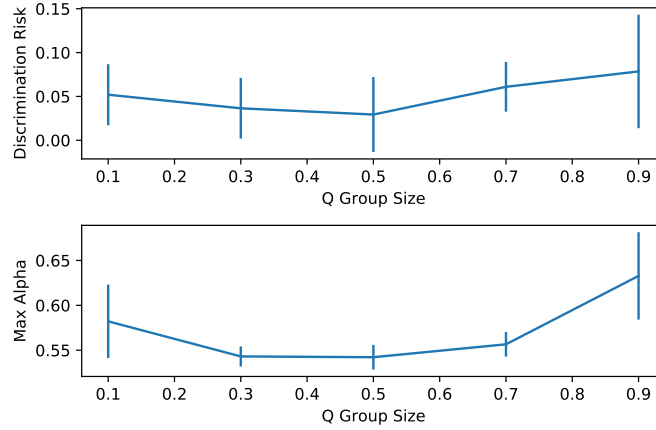


Figure A.2: Plots of discrimination risk and maximum α (over all channels) of Feature Propagation vs. relative size of group Q in SBM. We use 0.5 unknown feature rates for both groups and 0.5 inter- and intra-link rates.

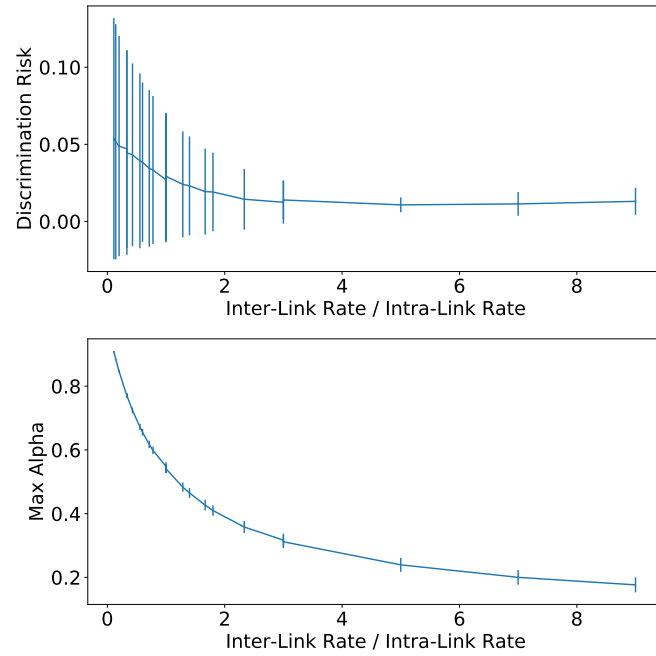


Figure A.3: Plots of discrimination risk and maximum α (over all channels) of Feature Propagation vs. ratio of inter-link rate to intra-link rate in SBM. We use 0.5 relative group sizes and 0.5 unknown feature rates for both groups.

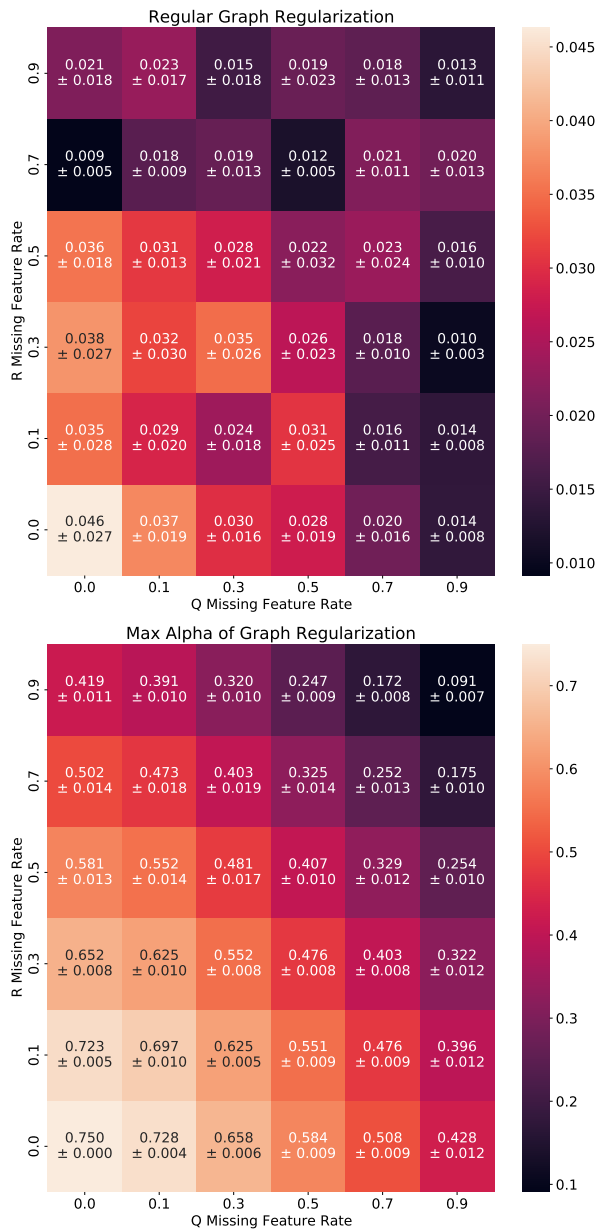


Figure A.4: Heatmap of discrimination risks and maximum α (over all channels) of Graph Regularization for 36 combinations of unknown feature rates for each group in SBM. We use 0.5 relative group sizes and 0.5 inter- and intra-link rates.

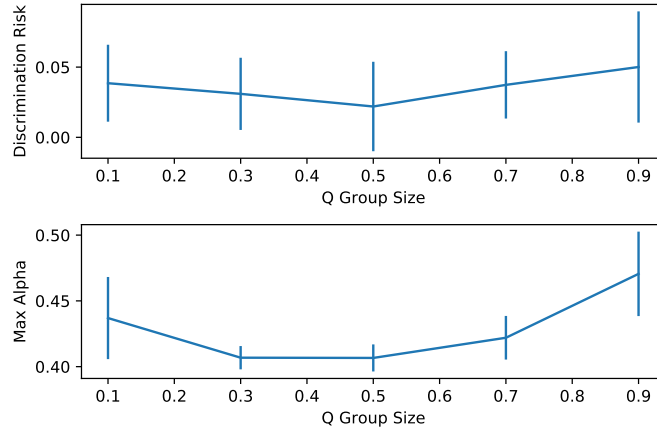


Figure A.5: Plots of discrimination risk and maximum α (over all channels) of Graph Regularization vs. relative size of group Q in SBM. We use 0.5 unknown feature rates for both groups and 0.5 inter- and intra-link rates.

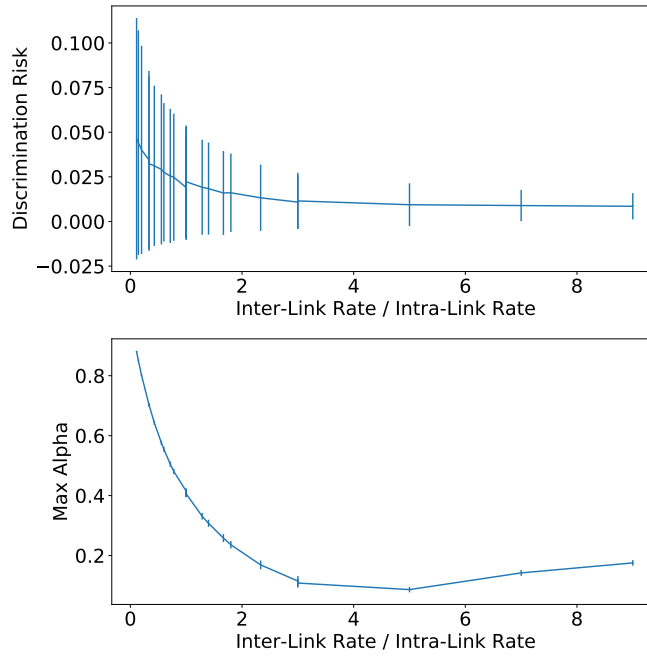


Figure A.6: Plots of discrimination risk and maximum α (over all channels) of Graph Regularization vs. ratio of inter-link rate to intra-link rate in SBM. We use 0.5 relative group sizes and 0.5 unknown feature rates for both groups.

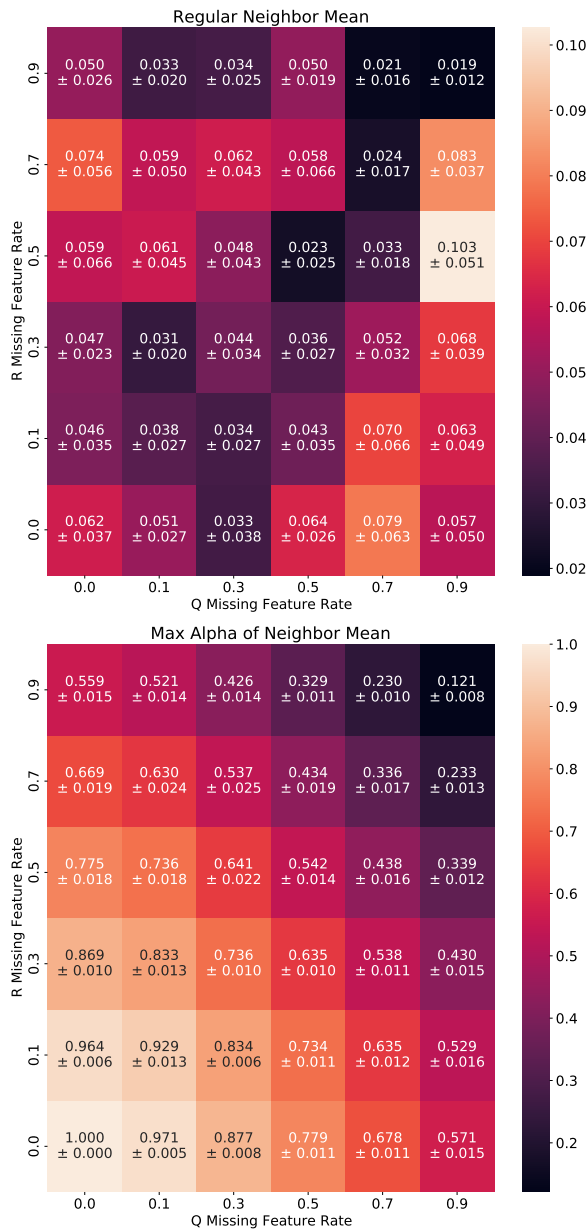


Figure A.7: Heatmap of discrimination risks and maximum α (over all channels) of Neighbor Mean for 36 combinations of unknown feature rates for each group in SBM. We use 0.5 relative group sizes and 0.5 inter- and intra-link rates.

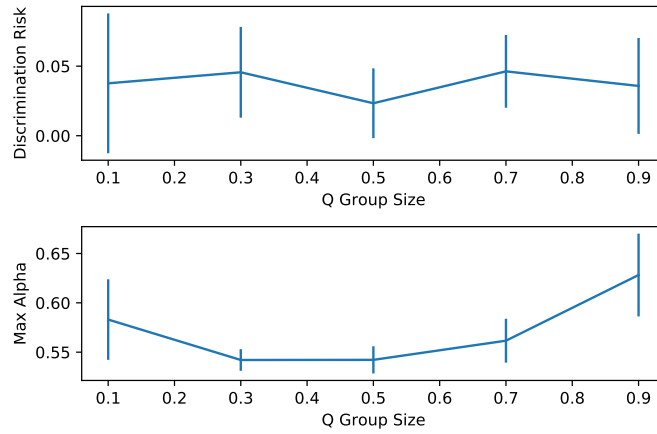


Figure A.8: Plots of discrimination risk and maximum α (over all channels) of Neighbor Mean vs. relative size of group Q in SBM. We use 0.5 unknown feature rates for both groups and 0.5 inter- and intra-link rates.

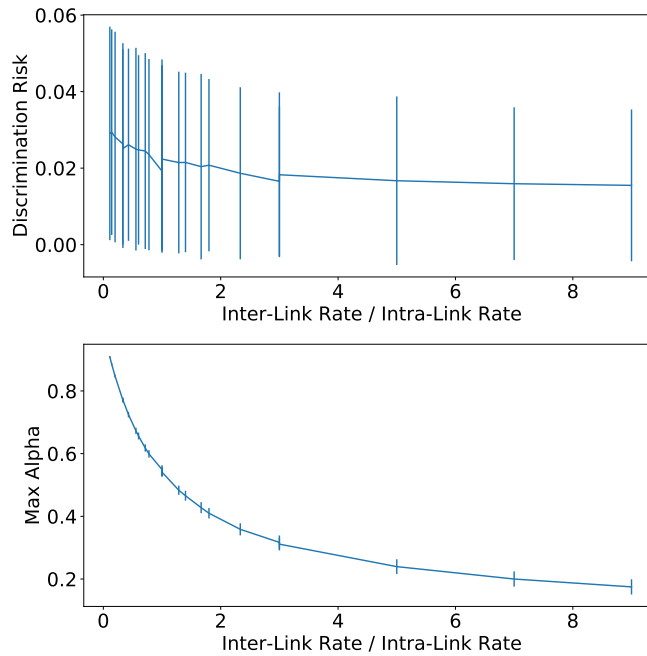


Figure A.9: Plots of discrimination risk and maximum α (over all channels) of Neighbor Mean vs. ratio of inter-link rate to intra-link rate in SBM. We use 0.5 relative group sizes and 0.5 unknown feature rates for both groups.

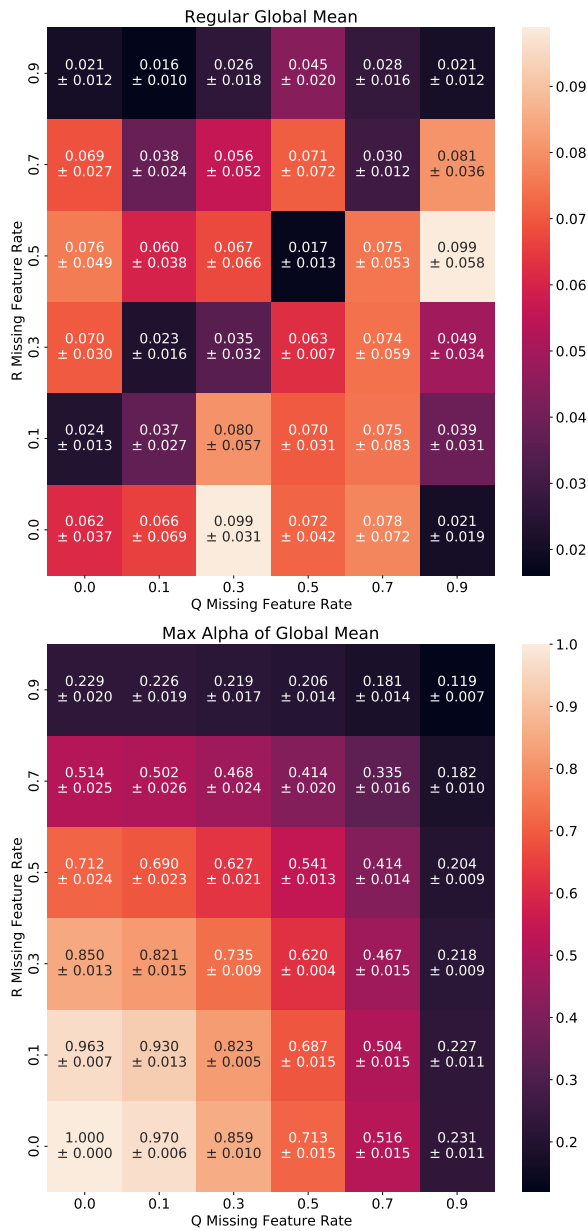


Figure A.10: Heatmap of discrimination risks and maximum α (over all channels) of Global Mean for 36 combinations of unknown feature rates for each group in SBM. We use 0.5 relative group sizes. **Note:** Global Mean is not affected by graph structure.

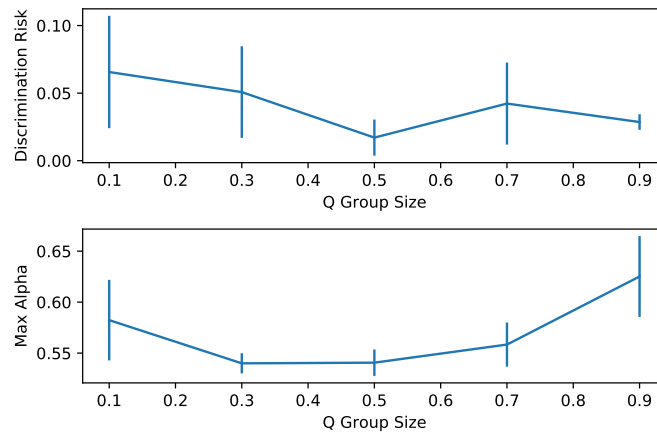


Figure A.11: Plots of discrimination risk and maximum α (over all channels) of Global Mean vs. relative size of group Q in SBM. We use 0.5 unknown feature rates for both groups. **Note:** Global Mean is not affected by graph structure.

Table A.1: Reconstruction error (**RE**), discrimination risk (**DR**), and test group membership identification accuracy (**MI**) of all models averaged over relative sizes of group Q of $\{0.1, 0.3, 0.5, 0.7, 0.9\}$ in **SBM**. We use 0.5 unknown feature rates for both groups and 0.5 inter- and intra-link rates.

Method	RE ↓	DR ↓	MI _{linear} ↓	MI _{mlp} ↓	MI _{gen} ↓
0-Fair GM	1.054 ± 0.004	0 ± 0	0.758 ± 0.025	0.78 ± 0.014	0.742 ± 0.018
0.025-Fair GM	1.051 ± 0.004	0.022 ± 0.004	0.79 ± 0.012	0.787 ± 0.012	0.74 ± 0.014
0.05-Fair GM	1.048 ± 0.003	0.032 ± 0.003	0.791 ± 0.018	0.794 ± 0.013	0.747 ± 0.019
Regular GM	1 ± 0	0.041 ± 0.008	0.835 ± 0.01	0.845 ± 0.011	0.771 ± 0.027
0-Fair NM	1.015 ± 0.003	0 ± 0	0.757 ± 0.023	0.792 ± 0.015	0.738 ± 0.016
0.025-Fair NM	1.012 ± 0.003	0.019 ± 0.003	0.791 ± 0.022	0.8 ± 0.014	0.744 ± 0.011
0.05-Fair NM	1.009 ± 0.003	0.029 ± 0.006	0.787 ± 0.017	0.807 ± 0.011	0.746 ± 0.017
Regular NM	0.959 ± 0.003	0.038 ± 0.013	0.835 ± 0.011	0.843 ± 0.015	0.763 ± 0.024
0-Fair FP	1.003 ± 0.005	0 ± 0	0.757 ± 0.019	0.799 ± 0.014	0.736 ± 0.013
0.025-Fair FP	1 ± 0.005	0.021 ± 0.002	0.785 ± 0.021	0.801 ± 0.014	0.754 ± 0.017
0.05-Fair FP	0.997 ± 0.005	0.033 ± 0.005	0.789 ± 0.016	0.806 ± 0.012	0.738 ± 0.016
Regular FP	0.947 ± 0.005	0.051 ± 0.017	0.829 ± 0.006	0.841 ± 0.02	0.760 ± 0.022
0-Fair GR	0.962 ± 0.005	0 ± 0	0.752 ± 0.024	0.788 ± 0.019	0.742 ± 0.013
0.025-Fair GR	0.961 ± 0.005	0.023 ± 0.003	0.797 ± 0.015	0.797 ± 0.02	0.752 ± 0.016
0.05-Fair GR	0.96 ± 0.005	0.036 ± 0.005	0.799 ± 0.009	0.805 ± 0.014	0.739 ± 0.017
Regular GR	0.945 ± 0.006	0.036 ± 0.012	0.821 ± 0.015	0.82 ± 0.014	0.759 ± 0.021

Table A.2: Reconstruction error (**RE**), discrimination risk (**DR**), and test group membership identification accuracy (**MI**) of all models averaged over all 25 combinations of inter- and intra-link rates of $\{0.1, 0.3, 0.5, 0.7, 0.9\}$ in **SBM**. We use 0.5 relative group sizes and 0.5 unknown feature rates for both groups.

Method	RE ↓	DR ↓	MI _{linear} ↓	MI _{mlp} ↓	MI _{gcn} ↓
0-Fair NM	1.028 ± 0.009	0 ± 0	0.609 ± 0.102	0.729 ± 0.046	0.905 ± 0.003
0.025-Fair NM	1.023 ± 0.008	0.014 ± 0.011	0.724 ± 0.09	0.749 ± 0.035	0.911 ± 0.008
0.05-Fair NM	1.019 ± 0.008	0.02 ± 0.02	0.74 ± 0.046	0.768 ± 0.036	0.911 ± 0.01
Regular NM	0.931 ± 0.003	0.022 ± 0.024	0.845 ± 0.026	0.866 ± 0.027	0.924 ± 0.008
0-Fair FP	1.022 ± 0.012	0 ± 0	0.64 ± 0.064	0.742 ± 0.038	0.905 ± 0.003
0.025-Fair FP	1.017 ± 0.012	0.014 ± 0.012	0.697 ± 0.095	0.753 ± 0.039	0.912 ± 0.007
0.05-Fair FP	1.013 ± 0.012	0.023 ± 0.022	0.740 ± 0.040	0.762 ± 0.04	0.909 ± 0.008
Regular FP	0.918 ± 0.004	0.034 ± 0.043	0.844 ± 0.025	0.853 ± 0.035	0.922 ± 0.009
0-Fair GR	0.948 ± 0.043	0 ± 0	0.578 ± 0.105	0.773 ± 0.038	0.905 ± 0.004
0.025-Fair GR	0.946 ± 0.004	0.016 ± 0.013	0.779 ± 0.036	0.793 ± 0.038	0.915 ± 0.005
0.05-Fair GR	0.945 ± 0.004	0.02 ± 0.02	0.769 ± 0.018	0.797 ± 0.0366	0.912 ± 0.009
Regular GR	0.916 ± 0.005	0.023 ± 0.032	0.846 ± 0.023	0.864 ± 0.023	0.921 ± 0.009

Table A.3: equal opportunity (**EO**) averaged over all 25 combinations of unknown feature rates of $\{0.1, 0.3, 0.5, 0.7, 0.9\}$ for each group in German credit.

Method	$\mathbf{EO}_{linear} \downarrow$	$\mathbf{EO}_{mlp} \downarrow$	$\mathbf{EO}_{gcn} \downarrow$
0.0-Fair GM	0.037 ± 0.01	0.029 ± 0.004	0.009 ± 0.008
0.025-Fair GM	0.031 ± 0.007	0.029 ± 0.008	0.018 ± 0.021
0.05-Fair GM	0.026 ± 0.003	0.03 ± 0.003	0.013 ± 0.005
Regular GM	0.033 ± 0.008	0.023 ± 0.003	0.006 ± 0.005
0.0-Fair NM	0.038 ± 0.009	0.037 ± 0.006	0.007 ± 0.006
0.025-Fair NM	0.035 ± 0.008	0.038 ± 0.007	0.013 ± 0.012
0.05-Fair NM	0.04 ± 0.006	0.035 ± 0.006	0.009 ± 0.003
Regular NM	0.038 ± 0.012	0.032 ± 0.006	0.012 ± 0.006
0.0-Fair FP	0.01 ± 0.011	0.034 ± 0.018	0.024 ± 0.041
0.025-Fair FP	0.028 ± 0.031	0.031 ± 0.018	0.023 ± 0.051
0.05-Fair FP	0.043 ± 0.07	0.029 ± 0.028	0 ± 0
Regular FP	0.042 ± 0.046	0.038 ± 0.02	0.004 ± 0.006
0.0-Fair GR	0.029 ± 0.012	0.022 ± 0.003	0.005 ± 0.006
0.025-Fair GR	0.031 ± 0.011	0.024 ± 0.005	0.007 ± 0.007
0.05-Fair GR	0.027 ± 0.007	0.024 ± 0.006	0.004 ± 0.004
Regular GR	0.032 ± 0.01	0.025 ± 0.007	0.009 ± 0.01

Table A.4: Test accuracy (**Acc**) and statistical parity (**SP**) of all models averaged over all 25 combinations of unknown feature rates of $\{0.1, 0.3, 0.5, 0.7, 0.9\}$ for each group in **Credit defaulter**.

Method	$\text{Acc}_{linear} \uparrow$	$\text{Acc}_{mlp} \uparrow$	$\text{Acc}_{gcn} \uparrow$	$\text{SP}_{linear} \downarrow$	$\text{SP}_{mlp} \downarrow$	$\text{SP}_{gcn} \downarrow$
0.0-Fair GM	0.781 ± 0.002	0.764 ± 0.012	0.771 ± 0.007	0.063 ± 0.015	0.08 ± 0.024	0.016 ± 0.009
0.025-Fair GM	0.78 ± 0.006	0.757 ± 0.008	0.774 ± 0.002	0.059 ± 0.021	0.083 ± 0.012	0.015 ± 0.004
0.05-Fair GM	0.78 ± 0.003	0.759 ± 0.013	0.775 ± 0.002	0.076 ± 0.015	0.08 ± 0.017	0.015 ± 0.003
Regular GM	0.782 ± 0.006	0.76 ± 0.018	0.775 ± 0.006	0.055 ± 0.031	0.056 ± 0.012	0.005 ± 0.006
0.0-Fair NM	0.781 ± 0.002	0.765 ± 0.006	0.771 ± 0.007	0.063 ± 0.017	0.085 ± 0.013	0.015 ± 0.013
0.025-Fair NM	0.78 ± 0.005	0.766 ± 0.005	0.774 ± 0.002	0.057 ± 0.025	0.088 ± 0.008	0.015 ± 0.005
0.05-Fair GM	0.781 ± 0.003	0.769 ± 0.01	0.775 ± 0.002	0.082 ± 0.011	0.082 ± 0.018	0.016 ± 0.003
Regular GM	0.781 ± 0.007	0.762 ± 0.014	0.773 ± 0.008	0.054 ± 0.031	0.061 ± 0.011	0.005 ± 0.007
0.0-Fair FP	0.779 ± 0.005	0.757 ± 0.022	0.77 ± 0.008	0.06 ± 0.022	0.085 ± 0.014	0.016 ± 0.014
0.025-Fair FP	0.78 ± 0.002	0.764 ± 0.004	0.774 ± 0.002	0.056 ± 0.027	0.092 ± 0.008	0.016 ± 0.005
0.05-Fair FP	0.78 ± 0.001	0.768 ± 0.008	0.774 ± 0.002	0.076 ± 0.005	0.084 ± 0.014	0.017 ± 0.004
Regular FP	0.781 ± 0.005	0.764 ± 0.01	0.775 ± 0.006	0.051 ± 0.029	0.075 ± 0.011	0.005 ± 0.006
0.0-Fair GR	0.773 ± 0.009	0.796 ± 0.006	0.771 ± 0.011	0.072 ± 0.044	0.098 ± 0.032	0.011 ± 0.012
0.025-Fair GR	0.779 ± 0.005	0.792 ± 0.007	0.772 ± 0.003	0.052 ± 0.032	0.091 ± 0.02	0.017 ± 0.012
0.05-Fair GR	0.78 ± 0.003	0.792 ± 0.007	0.773 ± 0.004	0.078 ± 0.0314	0.094 ± 0.028	0.023 ± 0.01
Regular GR	0.781 ± 0.005	0.785 ± 0.004	0.773 ± 0.008	0.049 ± 0.043	0.073 ± 0.038	0.008 ± 0.01

Table A.5: equal opportunity (**EO**) of all models averaged over all 25 combinations of unknown feature rates of $\{0.1, 0.3, 0.5, 0.7, 0.9\}$ for each group in Credit defaulter.

Method	$\mathbf{EO}_{linear} \downarrow$	$\mathbf{EO}_{mlp} \downarrow$	$\mathbf{EO}_{gcn} \downarrow$
0.0-Fair GM	0.039 ± 0.008	0.056 ± 0.019	0.013 ± 0.007
0.025-Fair GM	0.035 ± 0.014	0.058 ± 0.009	0.012 ± 0.003
0.05-Fair GM	0.048 ± 0.010	0.057 ± 0.016	0.012 ± 0.002
Regular GM	0.031 ± 0.017	0.038 ± 0.011	0.004 ± 0.005
0.0-Fair NM	0.039 ± 0.01	0.06 ± 0.009	0.013 ± 0.007
0.025-Fair NM	0.033 ± 0.015	0.06 ± 0.007	0.011 ± 0.004
0.05-Fair NM	0.05 ± 0.007	0.057 ± 0.014	0.012 ± 0.002
Regular NM	0.031 ± 0.018	0.041 ± 0.007	0.005 ± 0.006
0.0-Fair FP	0.035 ± 0.013	0.059 ± 0.013	0.013 ± 0.008
0.025-Fair FP	0.031 ± 0.016	0.062 ± 0.007	0.013 ± 0.004
0.05-Fair FP	0.043 ± 0.005	0.057 ± 0.011	0.014 ± 0.002
Regular FP	0.028 ± 0.016	0.05 ± 0.01	0.004 ± 0.005
0.0-Fair GR	0.051 ± 0.036	0.07 ± 0.029	0.008 ± 0.011
0.025-Fair GR	0.03 ± 0.02	0.06 ± 0.025	0.012 ± 0.011
0.05-Fair GR	0.048 ± 0.03	0.067 ± 0.025	0.015 ± 0.009
Regular GR	0.03 ± 0.038	0.051 ± 0.038	0.007 ± 0.007

APPENDIX B

Appendix for Chapter 3

B.1 Proofs

B.1.1 Proof of Lemma 3.4.1

Proof. Similarly to [XLT18, TYS20], we compute the first-order partial derivatives of Φ_s and Φ_r :

$$\frac{\partial \mathbf{s}_i^{(L)}}{\partial \mathbf{x}_j} = \sum_{p \in \Psi_{i \rightarrow j}^{L+1}} \prod_{l=L}^1 \frac{\text{diag}\left(\mathbb{1}_{\mathbf{z}_{p^{(l)}}^{(l)} > 0}\right) \mathbf{W}_s^{(l)}}{\sqrt{\mathbf{D}_{p^{(l)}p^{(l)}} \mathbf{D}_{p^{(l-1)}p^{(l-1)}}}}, \quad \frac{\partial \mathbf{r}_i^{(L)}}{\partial \mathbf{x}_j} = \sum_{p \in \Psi_{i \rightarrow j}^{L+1}} \prod_{l=L}^1 \frac{\text{diag}\left(\mathbb{1}_{\mathbf{z}_{p^{(l)}}^{(l)} > 0}\right) \mathbf{W}_s^{(l)}}{\mathbf{D}_{p^{(l)}p^{(l)}}} \quad (\text{B.1})$$

$$\frac{\partial \mathbf{s}_i^{(L)}}{\partial \mathbf{x}_j} = \sqrt{\frac{\mathbf{D}_{ii}}{\mathbf{D}_{jj}}} \sum_{p \in \Psi_{i \rightarrow j}^{L+1}} \prod_{l=L}^1 \frac{\text{diag}\left(\mathbb{1}_{\mathbf{z}_{p^{(l)}}^{(l)} > 0}\right) \mathbf{W}_s^{(l)}}{\mathbf{D}_{p^{(l)}p^{(l)}}} \quad (\text{B.2})$$

where $p^{(l)}$ is the l -th node on path p in the computation graph of Φ_s or Φ_r ($p^{(L)}$ is node i and $p^{(0)}$ is node j); $\Psi_{i \rightarrow j}^\gamma$ is the set of all γ -length random walk paths from node i to j ; and $\mathbf{z}_{p^{(l)}}^{(l)}$ is pre-activated $\mathbf{s}_{p^{(l)}}^{(l)}$ or $\mathbf{r}_{p^{(l)}}^{(l)}$.

With our assumption that the path from node $i \rightarrow j$ in the computation graph of Φ_s is independently activated with probability $\rho_s(i)$, and similarly, $\rho_r(i)$ for Φ_r :

$$\mathbb{E} \left[\frac{\partial \mathbf{s}_i^{(L)}}{\partial \mathbf{x}_j} \right] = \left(\mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} \right)_{ij}^L \rho_s(i) \left(\prod_{l=L}^1 \mathbf{W}_s^{(l)} \right), \quad (\text{B.3})$$

$$\mathbb{E} \left[\frac{\partial \mathbf{r}_i^{(L)}}{\partial \mathbf{x}_j} \right] = \left(\mathbf{D}^{-1} \mathbf{A} \right)_{ij}^L \rho_r(i) \left(\prod_{l=L}^1 \mathbf{W}_r^{(l)} \right). \quad (\text{B.4})$$

Then, recalling Eqn. 3.3:

$$\mathbb{E} \left[\mathbf{s}_i^{(L)} \right] = \sum_{j \in \mathcal{V}} \left(\mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} \right)_{ij}^L \rho_s(i) \left(\prod_{l=L}^1 \mathbf{W}_s^{(l)} \right) \mathbf{x}_j + \mathbf{0}, \quad (\text{B.5})$$

$$\mathbb{E} \left[\mathbf{r}_i^{(L)} \right] = \sum_{j \in \mathcal{V}} \left(\mathbf{D}^{-1} \mathbf{A} \right)_{ij}^L \rho_r(i) \left(\prod_{l=L}^1 \mathbf{W}_r^{(l)} \right) \mathbf{x}_j + \mathbf{0} \quad (\text{B.6})$$

$$\mathbb{E} \left[\mathbf{s}_i^{(L)} \right] = \sum_{j \in \mathcal{V}} \rho_s(i) \left(\mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} \right)_{ij}^L \alpha_j, \quad \mathbb{E} \left[\mathbf{r}_i^{(L)} \right] = \sum_{j \in \mathcal{V}} \rho_r(i) \left(\mathbf{D}^{-1} \mathbf{A} \right)_{ij}^L \beta_j. \quad (\text{B.7})$$

□

The independence of path activation probabilities may not always hold true in practice. However, we verify that this assumption is plausible via our extensive experiments on real-world datasets that validate our theoretical analysis (§3.6.1). This assumption also aligns with findings that deep neural networks have an inductive bias towards learning simpler, often linear, functions [NKK19, VCL19]. Furthermore, a variant of our assumption (where $\rho(i) = \rho$ is constant for all nodes) has been used in the literature to simplify theoretical analysis (e.g., [XLT18, TYS20]); our assumption may be more realistic than this variant, as it captures that the probability of paths activating can differ across nodes (e.g., due to differences in features, neighborhood structure).

B.1.2 Proof of Lemma 3.4.2

Proof. For $j \in S^{(b)}$, we can re-express $\widehat{\mathbf{P}}_{ij}^L = \left(\widehat{\mathbf{P}}^{(b)} \right)_{ij}^L = \left(\mathbf{e}^{(i)} \right)^\top \left(\widehat{\mathbf{P}}^{(b)} \right)^L \mathbf{e}^{(j)}$ ¹. By the spectral properties of $\widehat{\mathbf{P}}^{(b)}$, $\left(\mathbf{e}^{(i)} \right)^\top \mathbf{v}_1^{(b)} = \sqrt{\frac{\widehat{\mathbf{D}}_{ii}}{\text{vol}(\mathcal{G}^{(b)})}}$ [Lov96]. Hence:

¹For simplicity, we abuse notation here: $\left(\widehat{\mathbf{P}}^{(b)} \right)_{ij}^L$ is not the entry at row i and column j , but rather the entry at the row corresponding to node i and column corresponding to node j . Similarly, $\mathbf{e}^{(i)}$ is the standard basis vector with a 1 at the entry corresponding to node i .

$$\widehat{\mathbf{P}}_{ij}^L = \sum_{k=1}^{|S^{(b)}|} \left(\lambda_k^{(b)}\right)^L (\mathbf{e}^{(i)})^\top \mathbf{v}_k^{(b)} \left(\mathbf{v}_k^{(b)}\right)^\top \mathbf{e}^{(j)} \quad (\text{B.8})$$

$$= \frac{\sqrt{\widehat{\mathbf{D}}_{ii} \widehat{\mathbf{D}}_{jj}}}{\text{vol}(\mathcal{G}^{(b)})} + \sum_{k=2}^{|S^{(b)}|} \left(\lambda_k^{(b)}\right)^L (\mathbf{e}^{(i)})^\top \mathbf{v}_k^{(b)} \left(\mathbf{v}_k^{(b)}\right)^\top \mathbf{e}^{(j)} \quad (\text{B.9})$$

Then, by Cauchy-Schwarz:

$$\left| \widehat{\mathbf{P}}_{ij}^L - \frac{\sqrt{\widehat{\mathbf{D}}_{ii} \widehat{\mathbf{D}}_{jj}}}{\text{vol}(\mathcal{G}^{(b)})} \right| \leq \left(\lambda^{(b)}\right)^L \sum_{k=1}^{|S^{(b)}|} \left| (\mathbf{e}^{(i)})^\top \mathbf{v}_k^{(b)} \right| \left| (\mathbf{e}^{(j)})^\top \mathbf{v}_k^{(b)} \right| \quad (\text{B.10})$$

$$\leq \left(\lambda^{(b)}\right)^L \left(\sum_{k=1}^{|S^{(b)}|} \left| (\mathbf{e}^{(i)})^\top \mathbf{v}_k^{(b)} \right|^2 \right)^{\frac{1}{2}} \left(\sum_{k=1}^{|S^{(b)}|} \left| (\mathbf{e}^{(j)})^\top \mathbf{v}_k^{(b)} \right|^2 \right)^{\frac{1}{2}} \quad (\text{B.11})$$

$$= \left(\lambda^{(b)}\right)^L \left((\mathbf{e}^{(i)})^\top \mathbf{V}^{(b)} (\mathbf{V}^{(b)})^\top \mathbf{e}^{(i)} \right)^{\frac{1}{2}} \left((\mathbf{e}^{(j)})^\top \mathbf{V}^{(b)} (\mathbf{V}^{(b)})^\top \mathbf{e}^{(j)} \right)^{\frac{1}{2}} \quad (\text{B.12})$$

$$= \left(\lambda^{(b)}\right)^L \left\| \mathbf{e}^{(i)} \right\|_2 \left\| \mathbf{e}^{(j)} \right\|_2 \quad (\text{B.13})$$

$$= \left(\lambda^{(b)}\right)^L \quad (\text{B.14})$$

Let $\mathbf{P}^L = \left(\widehat{\mathbf{P}} + \Xi^{(0)}\right)^L = \widehat{\mathbf{P}}^L + \Xi^{(L)}$. Then, by the triangle inequality:

$$\left| \mathbf{P}_{ij}^L - \frac{\sqrt{\widehat{\mathbf{D}}_{ii} \widehat{\mathbf{D}}_{jj}}}{\text{vol}(\mathcal{G}^{(b)})} \right| \leq \left(\lambda^{(b)}\right)^L + \left| (\mathbf{e}^{(i)})^\top \Xi^{(L)} \mathbf{e}^{(j)} \right| \quad (\text{B.15})$$

$$\leq \left(\lambda^{(b)}\right)^L + \left\| \Xi^{(L)} \right\|_{op} \quad (\text{B.16})$$

$$\leq \left(\lambda^{(b)}\right)^L + \sum_{l=1}^L \binom{L}{l} \left\| \Xi^{(0)} \right\|_{op}^l \left\| \widehat{\mathbf{P}} \right\|_{op}^{L-l} \quad (\text{B.17})$$

For $j \notin S^{(b)}$, $\widehat{\mathbf{P}}_{ij}^L = 0$. Then:

$$\left| \mathbf{P}_{ij}^L - 0 \right| \leq \left| (\mathbf{e}^{(i)})^\top \Xi^{(L)} \mathbf{e}^{(j)} \right| \quad (\text{B.18})$$

$$\leq \sum_{l=1}^L \binom{L}{l} \left\| \Xi^{(0)} \right\|_{op}^l \left\| \widehat{\mathbf{P}} \right\|_{op}^{L-l} \quad (\text{B.19})$$

□

B.1.3 Proof of Theorem 3.4.3

Proof. For $u, v \in \mathcal{V}$, let $|\delta_{uv}| \leq \zeta_s$. Combining Lemmas 3.4.1 and 3.4.2, by our assumption that the computation graph paths to i, j are activated independently:

$$\mathbb{E} \left[f_{LP} \left(\mathbf{s}_i^{(L)}, \mathbf{s}_j^{(L)} \right) \right] = \mathbb{E} \left[\mathbf{s}_i^{(L)} \right]^\top \mathbb{E} \left[\mathbf{s}_j^{(L)} \right] \quad (\text{B.20})$$

$$= \bar{\rho}_s^2(b) \left(\sum_{k \in S^{(b)}} \frac{\sqrt{\widehat{\mathbf{D}}_{ii} \widehat{\mathbf{D}}_{kk}}}{\text{vol}(\mathcal{G}^{(b)})} \alpha_k + \sum_{k \in \mathcal{V}} \delta_{ik} \alpha_k \right)^\top \left(\sum_{k \in S^{(b)}} \frac{\sqrt{\widehat{\mathbf{D}}_{jj} \widehat{\mathbf{D}}_{kk}}}{\text{vol}(\mathcal{G}^{(b)})} \alpha_k + \sum_{k \in \mathcal{V}} \delta_{jk} \alpha_k \right) \quad (\text{B.21})$$

$$= \bar{\rho}_s^2(b) \sqrt{\widehat{\mathbf{D}}_{ii} \widehat{\mathbf{D}}_{jj}} \underbrace{\left\| \sum_{k \in S^{(b)}} \frac{\sqrt{\widehat{\mathbf{D}}_{kk}}}{\text{vol}(\mathcal{G}^{(b)})} \alpha_k \right\|_2^2}_{\geq 0} \quad (\text{B.22})$$

$$+ \bar{\rho}_s^2(b) \left(\sqrt{\widehat{\mathbf{D}}_{ii}} \sum_{k \in S^{(b)}} \frac{\sqrt{\widehat{\mathbf{D}}_{kk}}}{\text{vol}(\mathcal{G}^{(b)})} \alpha_k \right)^\top \left(\sum_{k \in \mathcal{V}} \delta_{jk} \alpha_k \right) \quad (\text{B.23})$$

$$+ \bar{\rho}_s^2(b) \left(\sum_{k \in \mathcal{V}} \delta_{ik} \alpha_k \right)^\top \left(\sqrt{\widehat{\mathbf{D}}_{jj}} \sum_{k \in S^{(b)}} \frac{\sqrt{\widehat{\mathbf{D}}_{kk}}}{\text{vol}(\mathcal{G}^{(b)})} \alpha_k \right) \quad (\text{B.24})$$

$$+ \bar{\rho}_s^2(b) \left(\sum_{k \in \mathcal{V}} \delta_{ik} \alpha_k \right)^\top \left(\sum_{k \in \mathcal{V}} \delta_{jk} \alpha_k \right) \quad (\text{B.25})$$

Then, by Cauchy-Schwarz and the triangle inequality:

$$\left| \mathbb{E} \left[f_{LP} \left(\mathbf{s}_i^{(L)}, \mathbf{s}_j^{(L)} \right) \right] - \underbrace{\bar{\rho}_s^2(b) \left\| \sum_{k \in S^{(b)}} \frac{\sqrt{\widehat{\mathbf{D}}_{kk}}}{\text{vol}(\mathcal{G}^{(b)})} \alpha_k \right\|_2^2}_{\propto \sqrt{\widehat{\mathbf{D}}_{ii} \widehat{\mathbf{D}}_{jj}}} \sqrt{\widehat{\mathbf{D}}_{ii} \widehat{\mathbf{D}}_{jj}} \right| \quad (\text{B.26})$$

$$\leq \zeta_s \bar{\rho}_s^2(b) \left(\sqrt{\widehat{\mathbf{D}}_{ii}} + \sqrt{\widehat{\mathbf{D}}_{jj}} \right) \left\| \sum_{k \in S^{(b)}} \frac{\sqrt{\widehat{\mathbf{D}}_{kk}}}{\text{vol}(\mathcal{G}^{(b)})} \alpha_k \right\|_2 \left(\sum_{k \in \mathcal{V}} \|\alpha_k\|_2 \right) + \zeta_s^2 \bar{\rho}_s^2(b) \left(\sum_{k \in \mathcal{V}} \|\alpha_k\|_2 \right)^2 \quad (\text{B.27})$$

□

B.1.4 Lemma B.1.1 and Proof

Lemma B.1.1. We introduce the notation $\mathbf{P} = \mathbf{D}^{-1}\mathbf{A}$. We further define $\widehat{\mathbf{P}} = \widehat{\mathbf{D}}^{-1}\widehat{\mathbf{A}}$. Fix $i \in S^{(b)}$. Then, for $j \in S^{(b)}$:

$$\left| \mathbf{P}_{ij}^L - \frac{\widehat{\mathbf{D}}_{jj}}{\text{vol}(\mathcal{G}^{(b)})} \right| \leq \sqrt{\frac{\widehat{\mathbf{D}}_{jj}}{\widehat{\mathbf{D}}_{ii}}} (\lambda^{(b)})^L + \sum_{l=1}^L \binom{L}{l} \|\Xi^{(0)}\|_{op}^l \|\widehat{\mathbf{P}}\|_{op}^{L-l} \quad (\text{B.28})$$

And for $j \notin S^{(b)}$:

$$|\mathbf{P}_{ij}^L - 0| \leq \sum_{l=1}^L \binom{L}{l} \|\Xi^{(0)}\|_{op}^l \|\widehat{\mathbf{P}}\|_{op}^{L-l} \quad (\text{B.29})$$

Proof. Similar to the proof of Lemma 3.4.2:

$$\widehat{\mathbf{P}}_{ij}^L = \frac{\widehat{\mathbf{D}}_{jj}}{\text{vol}(\mathcal{G}^{(b)})} + \sqrt{\frac{\widehat{\mathbf{D}}_{jj}}{\widehat{\mathbf{D}}_{ii}}} \sum_{k=2}^{|S^{(b)}|} (\lambda_k^{(b)})^L (\mathbf{e}^{(i)})^\top \mathbf{v}_k^{(b)} (\mathbf{v}_k^{(b)})^\top \mathbf{e}^{(j)} \quad (\text{B.30})$$

Subsequently:

$$\left| \widehat{\mathbf{P}}_{ij}^L - \frac{\widehat{\mathbf{D}}_{jj}}{\text{vol}(\mathcal{G}^{(b)})} \right| \leq \sqrt{\frac{\widehat{\mathbf{D}}_{jj}}{\widehat{\mathbf{D}}_{ii}}} (\lambda^{(b)})^L \quad (\text{B.31})$$

Finally:

$$\left| \mathbf{P}_{ij}^L - \frac{\widehat{\mathbf{D}}_{jj}}{\text{vol}(\mathcal{G}^{(b)})} \right| \leq \zeta_r = \max_{u,v \in \mathcal{V}} \sqrt{\frac{\widehat{\mathbf{D}}_{vv}}{\widehat{\mathbf{D}}_{uu}}} (\lambda^{(b)})^L + \sum_{l=1}^L \binom{L}{l} \|\Xi^{(0)}\|_{op}^l \|\widehat{\mathbf{P}}\|_{op}^{L-l} \quad (\text{B.32})$$

For $j \notin S^{(b)}$, $\widehat{\mathbf{P}}_{ij}^L = 0$. Then:

$$|\mathbf{P}_{ij}^L - 0| \leq \sum_{l=1}^L \binom{L}{l} \|\Xi^{(0)}\|_{op}^l \|\widehat{\mathbf{P}}\|_{op}^{L-l} \leq \zeta_r \quad (\text{B.33})$$

□

B.1.5 Proof of Theorem 3.4.4

Proof. For $u, v \in \mathcal{V}$, let $|\delta_{uv}| \leq \zeta_r$. Combining Lemmas 3.4.1 and B.1.1, by our assumption that the computation graph paths to i, j are activated independently:

$$\mathbb{E} \left[f_{LP} \left(\mathbf{r}_i^{(L)}, \mathbf{r}_j^{(L)} \right) \right] = \mathbb{E} \left[\mathbf{r}_i^{(L)} \right]^\top \mathbb{E} \left[\mathbf{r}_j^{(L)} \right] \quad (\text{B.34})$$

$$= \bar{\rho}_r^2(b) \left(\sum_{k \in S^{(b)}} \frac{\hat{\mathbf{D}}_{kk}}{\text{vol}(\mathcal{G}^{(b)})} \beta_k + \sum_{k \in \mathcal{V}} \delta_{ik} \beta_k \right)^\top \left(\sum_{k \in S^{(b)}} \frac{\hat{\mathbf{D}}_{kk}}{\text{vol}(\mathcal{G}^{(b)})} \beta_k + \sum_{k \in \mathcal{V}} \delta_{jk} \beta_k \right) \quad (\text{B.35})$$

$$= \bar{\rho}_r^2(b) \underbrace{\left\| \sum_{k \in S^{(b)}} \frac{\hat{\mathbf{D}}_{kk}}{\text{vol}(\mathcal{G}^{(b)})} \beta_k \right\|_2^2}_{\geq 0} + \bar{\rho}_r^2(b) \left(\sum_{k \in S^{(b)}} \frac{\hat{\mathbf{D}}_{kk}}{\text{vol}(\mathcal{G}^{(b)})} \beta_k \right)^\top \left(\sum_{k \in \mathcal{V}} \delta_{jk} \beta_k \right) \quad (\text{B.36})$$

$$+ \bar{\rho}_r^2(b) \left(\sum_{k \in \mathcal{V}} \delta_{ik} \beta_k \right)^\top \left(\sum_{k \in S^{(b)}} \frac{\hat{\mathbf{D}}_{kk}}{\text{vol}(\mathcal{G}^{(b)})} \beta_k \right) + \bar{\rho}_r^2(b) \left(\sum_{k \in \mathcal{V}} \delta_{ik} \beta_k \right)^\top \left(\sum_{k \in \mathcal{V}} \delta_{jk} \beta_k \right) \quad (\text{B.37})$$

Then, by Cauchy-Schwarz and the triangle inequality:

$$\left| \mathbb{E} \left[f_{LP} \left(\mathbf{r}_i^{(L)}, \mathbf{r}_j^{(L)} \right) \right] - \underbrace{\bar{\rho}_r^2(b) \left\| \sum_{k \in S^{(b)}} \frac{\hat{\mathbf{D}}_{kk}}{\text{vol}(\mathcal{G}^{(b)})} \beta_k \right\|_2^2}_{\propto \text{constant}} \right| \quad (\text{B.38})$$

$$\leq \zeta_r \bar{\rho}_r^2(b) \left\| \sum_{k \in S^{(b)}} \frac{\hat{\mathbf{D}}_{kk}}{\text{vol}(\mathcal{G}^{(b)})} \beta_k \right\|_2 \left(\sum_{k \in \mathcal{V}} \|\beta_k\|_2 \right) + \zeta_r^2 \bar{\rho}_r^2(b) \left(\sum_{k \in \mathcal{V}} \|\beta_k\|_2 \right)^2 \quad (\text{B.39})$$

□

B.2 Approximation of $\Delta^{(b)}$

B.2.1 Approximation of $\Delta^{(b)}$ for Φ_s

$$\Delta^{(b)} \left(\mathbf{s}_i^{(L)}, \mathbf{s}_j^{(L)} \right) \quad (\text{B.40})$$

$$= \left| \frac{1}{|(S^{(b)} \cap T^{(1)}) \times S^{(b)}|} \sum_{i \in S^{(b)} \cap T^{(1)}} \sum_{j \in S^{(b)}} f_{LP} \left(\mathbf{s}_i^{(L)}, \mathbf{s}_j^{(L)} \right) \right. \quad (\text{B.41})$$

$$\left. - \frac{1}{|(S^{(b)} \cap T^{(2)}) \times S^{(b)}|} \sum_{i \in S^{(b)} \cap T^{(2)}} \sum_{j \in S^{(b)}} f_{LP} \left(\mathbf{s}_i^{(L)}, \mathbf{s}_j^{(L)} \right) \right| \quad (\text{B.42})$$

$$\cong \left| \frac{1}{|S^{(b)} \cap T^{(1)}| |S^{(b)}|} \sum_{i \in S^{(b)} \cap T^{(1)}} \sum_{j \in S^{(b)}} \bar{\rho}_s^2(b) \sqrt{\widehat{\mathbf{D}}_{ii} \widehat{\mathbf{D}}_{jj}} \left\| \sum_{k \in S^{(b)}} \frac{\sqrt{\widehat{\mathbf{D}}_{kk}}}{\text{vol}(\mathcal{G}^{(b)})} \alpha_k \right\|_2^2 \right. \quad (\text{B.43})$$

$$\left. - \frac{1}{|S^{(b)} \cap T^{(2)}| |S^{(b)}|} \sum_{i \in S^{(b)} \cap T^{(2)}} \sum_{j \in S^{(b)}} \bar{\rho}_s^2(b) \sqrt{\widehat{\mathbf{D}}_{ii} \widehat{\mathbf{D}}_{jj}} \left\| \sum_{k \in S^{(b)}} \frac{\sqrt{\widehat{\mathbf{D}}_{kk}}}{\text{vol}(\mathcal{G}^{(b)})} \alpha_k \right\|_2^2 \right| \quad (\text{B.44})$$

$$= \left| \frac{\bar{\rho}_s^2(b)}{|S^{(b)}|} \left\| \sum_{k \in S^{(b)}} \frac{\sqrt{\widehat{\mathbf{D}}_{kk}}}{\text{vol}(\mathcal{G}^{(b)})} \alpha_k \right\|_2^2 \right. \left. \left| \sum_{j \in S^{(b)}} \sqrt{\widehat{\mathbf{D}}_{jj}} \underbrace{\left(\mathbb{E}_{i \sim U(S^{(b)} \cap T^{(1)})} \sqrt{\widehat{\mathbf{D}}_{ii}} - \mathbb{E}_{i \sim U(S^{(b)} \cap T^{(2)})} \sqrt{\widehat{\mathbf{D}}_{ii}} \right)}_{\text{degree disparity}} \right| \right. \quad (\text{B.45})$$

B.2.2 Approximation of $\Delta^{(b)}$ for Φ_r

$$\Delta^{(b)} \left(\mathbf{r}_i^{(L)}, \mathbf{r}_j^{(L)} \right) \quad (\text{B.46})$$

$$= \left| \frac{1}{|(S^{(b)} \cap T^{(1)}) \times S^{(b)}|} \sum_{i \in S^{(b)} \cap T^{(1)}} \sum_{j \in S^{(b)}} f_{LP} \left(\mathbf{r}_i^{(L)}, \mathbf{r}_j^{(L)} \right) \right. \quad (\text{B.47})$$

$$\left. - \frac{1}{|(S^{(b)} \cap T^{(2)}) \times S^{(b)}|} \sum_{i \in S^{(b)} \cap T^{(2)}} \sum_{j \in S^{(b)}} f_{LP} \left(\mathbf{r}_i^{(L)}, \mathbf{r}_j^{(L)} \right) \right| \quad (\text{B.48})$$

$$\cong \left| \frac{1}{|S^{(b)} \cap T^{(1)}| |S^{(b)}|} \sum_{i \in S^{(b)} \cap T^{(1)}} \sum_{j \in S^{(b)}} \bar{\rho}_r^2(b) \left\| \sum_{k \in S^{(b)}} \frac{\widehat{\mathbf{D}}_{kk}}{\text{vol}(\mathcal{G}^{(b)})} \beta_k \right\|_2^2 \right. \quad (\text{B.49})$$

$$\left. - \frac{1}{|S^{(b)} \cap T^{(2)}| |S^{(b)}|} \sum_{i \in S^{(b)} \cap T^{(2)}} \sum_{j \in S^{(b)}} \bar{\rho}_r^2(b) \left\| \sum_{k \in S^{(b)}} \frac{\widehat{\mathbf{D}}_{kk}}{\text{vol}(\mathcal{G}^{(b)})} \beta_k \right\|_2^2 \right| \quad (\text{B.50})$$

$$= 0 \quad (\text{B.51})$$

B.3 Datasets Used in §3.6.1

In our experiments in §3.6.1, we use 10 real-world network datasets from [BG18a], [SMB19], [RS20], and [RAS21], covering diverse domains (e.g., citation networks, collaboration networks, online social networks). We provide a description and some statistics of each dataset in Table B.1. All the datasets have node features and are undirected. We were unable to find the exact class names and their label correspondence from the dataset documentation.

- In all the citation network datasets, nodes represent documents, edges represent citation links, and features are a bag-of-words representation of documents. We row-normalize the features to sum to 1, following [FL19]². The classification task is to predict the topic of documents.
- In the collaboration network datasets, nodes represent authors, edges represent co-authorships, and features are embeddings of paper keywords for authors' papers. The classification task is to predict the most active field of study for authors.
- In the LastFMAsia network dataset, nodes represent LastFM users from Asia, edges represent friendships between users, and features are embeddings of the artists liked by users. The classification task is to predict the home country of users.
- In the Twitch network datasets, nodes represent gamers on Twitch, edges represent followerships between them, and features are embeddings of the history of games played by the Twitch users. The classification task is to predict whether or not a gamer streams adult content.

We only run experiments on datasets that can fit without sampling nodes on a single NVIDIA GeForce GTX Titan Xp Graphic Card with 12196MiB of space. Furthermore, we

²https://github.com/pyg-team/pytorch_geometric/blob/master/examples/link_pred.py

only consider the three largest datasets (i.e., with the most nodes) from [RAS21]. We use PyTorch Geometric to load and process all datasets [FL19].

Table B.1: Summary of the datasets used in our experiments.

Name	Domain	# Nodes	# Edges	# Features	# Classes
Cora	citation	19793	126842	8710	70
CiteSeer	citation	4230	10674	602	6
DBLP	citation	17716	105734	1639	4
PubMed	citation	19717	88648	500	3
CS	collaboration	18333	163788	6805	15
Physics	collaboration	34493	495924	8415	5
LastFMAsia	online social	7624	55612	128	18
Twitch-DE	online social	9498	315774	128	2
Twitch-EN	online social	7126	77774	128	2
Twitch-FR	online social	6551	231883	128	2

B.4 Datasets Used in §3.6.2

We run experiments on three network datasets: (1) the NBA social network (see §B.4.1), (2) the German credit network (see §B.4.2), and (3) a new DBLP-Fairness citation network that we construct (see §B.4.3). All the datasets have node features and are undirected. We do not pass sensitive attributes as features to the models that we train. For each dataset, we min-max normalize node features to fall in $[-1, 1]$, following [DW21] and [ALZ21]. Furthermore, for all datasets, $D = 2$.

B.4.1 NBA Dataset

The NBA network [DW21] has 403 nodes representing NBA basketball players who are connected if they follow each other on Twitter. There are 21242 links. Each node has 95 features, with an average degree of 52.71 ± 35.14 . We consider two sensitive attributes per node:

- Age $\{S^{(b)}\}_{b \in [B]}$: how old the payer is, i.e., YOUNG (≤ 25 years) or OLD (> 25 years).
- Nationality $\{T^{(d)}\}_{d \in [D]}$: from where the player is, i.e., UNITED STATES or OVERSEAS.

B.4.2 German Dataset

The German network [ALZ21] comprises 1000 nodes representing clients in a German bank who are connected if they have similar credit accounts. The German network is not natively a graph dataset; synthetic edges were created by [ALZ21]. There are 44484 links. Each node has 27 features (e.g., loan amount, account-related features), with an average degree of 44.48 ± 26.52 . We consider two sensitive attributes per node:

- Foreign worker $\{S^{(b)}\}_{b \in [B]}$: whether the client is a foreign worker, i.e., YES or NO.
- Gender $\{T^{(d)}\}_{d \in [D]}$: the gender of the client, i.e., MAN or WOMAN.

B.4.3 DBLP-Fairness Dataset

In this subsection, we detail how we construct the DBLP-Fairness dataset. We build DBLP-Fairness, as there are only a few natively-graph network datasets with sensitive attributes that are appropriate for graph learning [SCS22].

We begin with the version of the DBLP-Citation-network V12 dataset from [TZY08] that was processed by [XBL22]. This dataset has 3658127 nodes. Each node represents a paper and each edge represents a citation link. We consider five node features:

- Team size: the number of authors on the paper.
- Mean collaborators: the average number of collaborators with whom the authors have previously published.
- Gini collaborators: the Gini coefficient of the number of collaborators with whom the authors have previously published.
- Mean productivity: the average number of papers that the authors have previously published.
- Gini productivity: the Gini coefficient of the number of papers that the authors have previously published.

We also consider two sensitive attributes per node:

- Field $\{S^{(b)}\}_{b \in [B]}$: the field to which the paper belongs, i.e., PROGRAMMING LANGUAGES or DATABASES.
- Nationality $\{T^{(d)}\}_{d \in [D]}$: the country where most authors reside, i.e., UNITED STATES or CHINA.

In DBLP-Fairness, we only include papers whose nationality is UNITED STATES or CHINA; American and Chinese citation networks are known to be stratified [ZGF22]. We also only include papers whose field is PROGRAMMING LANGUAGES or DATABASES; we infer the field of a paper using its keywords (i.e., whether they contain “programming language” and “database”), and discard papers which include both “programming language” and “database” in its keywords. Furthermore, we filter out all papers from before 2010. We sought DBLB-Fairness to be of comparable size to the citation networks in §B.3. Following filtering, we were left with 14537 nodes and 24844 edges.

B.5 Models

For all experiments, we use GCN encoders [KW17] to get node representations. Each encoder has two layers (128-dimensional hidden layer, 64-dimensional output layer) with a ReLU nonlinearity in between. We only use two layers, as this is common practice in graph deep learning to prevent oversmoothing [OS20]; however, we run experiments with four layers in §B.7. We do not use any regularization (e.g., Dropout, BatchNorm). The encoders are explicitly trained for LP with the inner-product LP score function in Eqn. 3.4, binary cross-entropy loss, and the Adam optimizer with full-batch gradient descent and a learning rate of 0.01 [KB15]. We use a random link split of 0.85-0.05-0.1 for train-val-test, following the PyTorch Geometric LP example³. We train the encoders for 100 epochs, with a new round of negative link sampling during every epoch; we use a 1:1 ratio of positive to negative links. We ultimately select the model parameters with the highest validation ROC-AUC. Although we do not do any hyperparameter tuning, the test ROC-AUC values (displayed in the figures in §3.6) indicate that the encoders are well-trained. We use PyTorch [PGM19] and PyTorch Geometric [FL19] to train all the encoders on a single NVIDIA GeForce GTX Titan Xp Graphic Card with 12196MiB of space.

B.6 Remaining Plots for §3.6.1

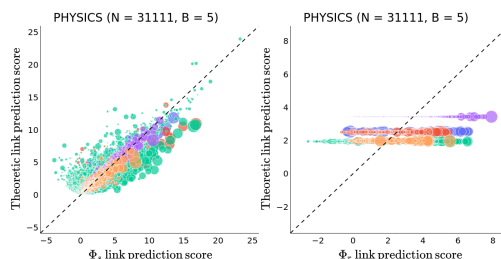


Figure B.1: Theoretic vs. GCN LP scores for collaboration network datasets.

³https://github.com/pyg-team/pytorch_geometric/blob/master/examples/link_pred.py

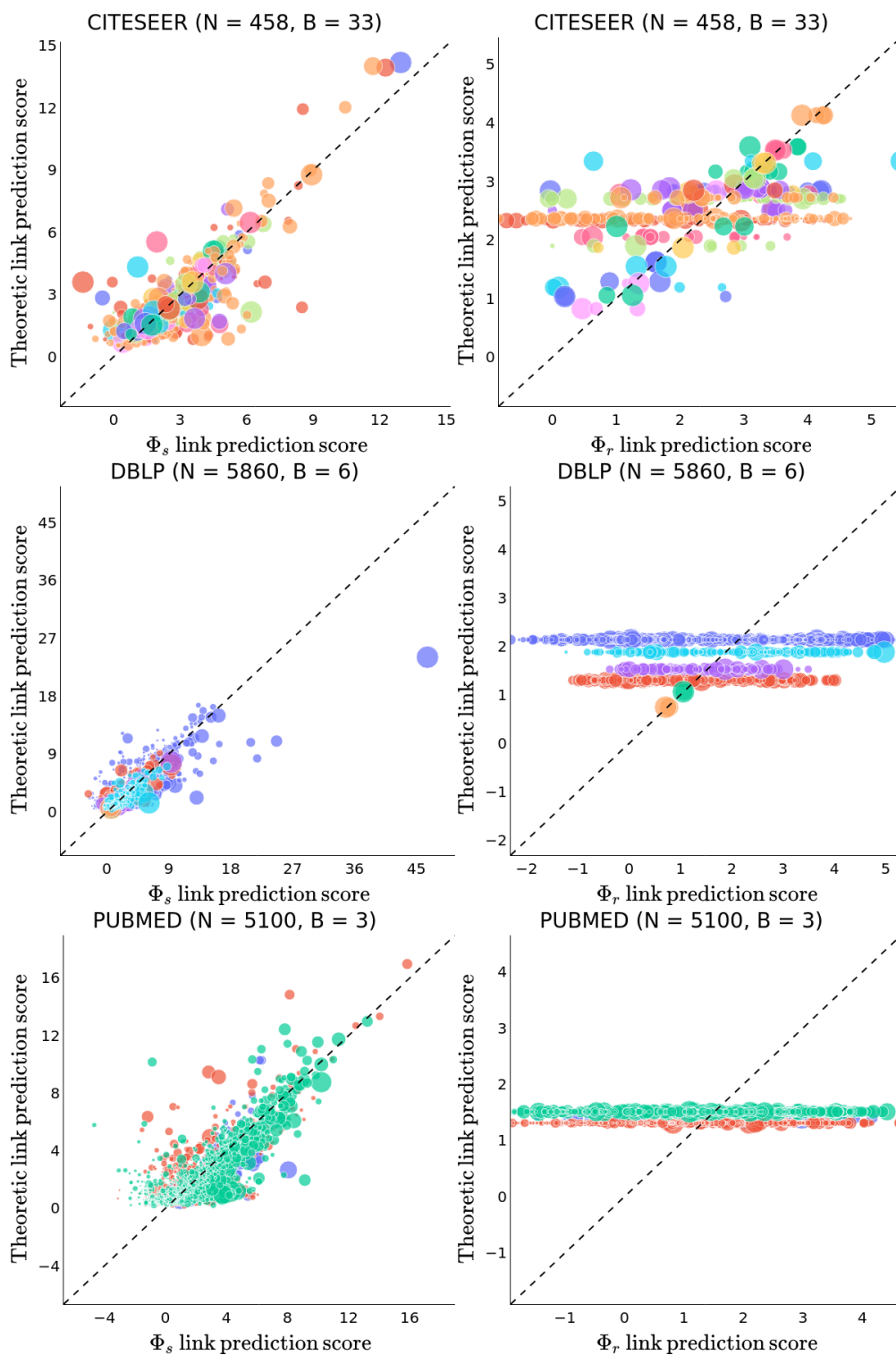


Figure B.2: Theoretic vs. GCN LP scores for citation network datasets.

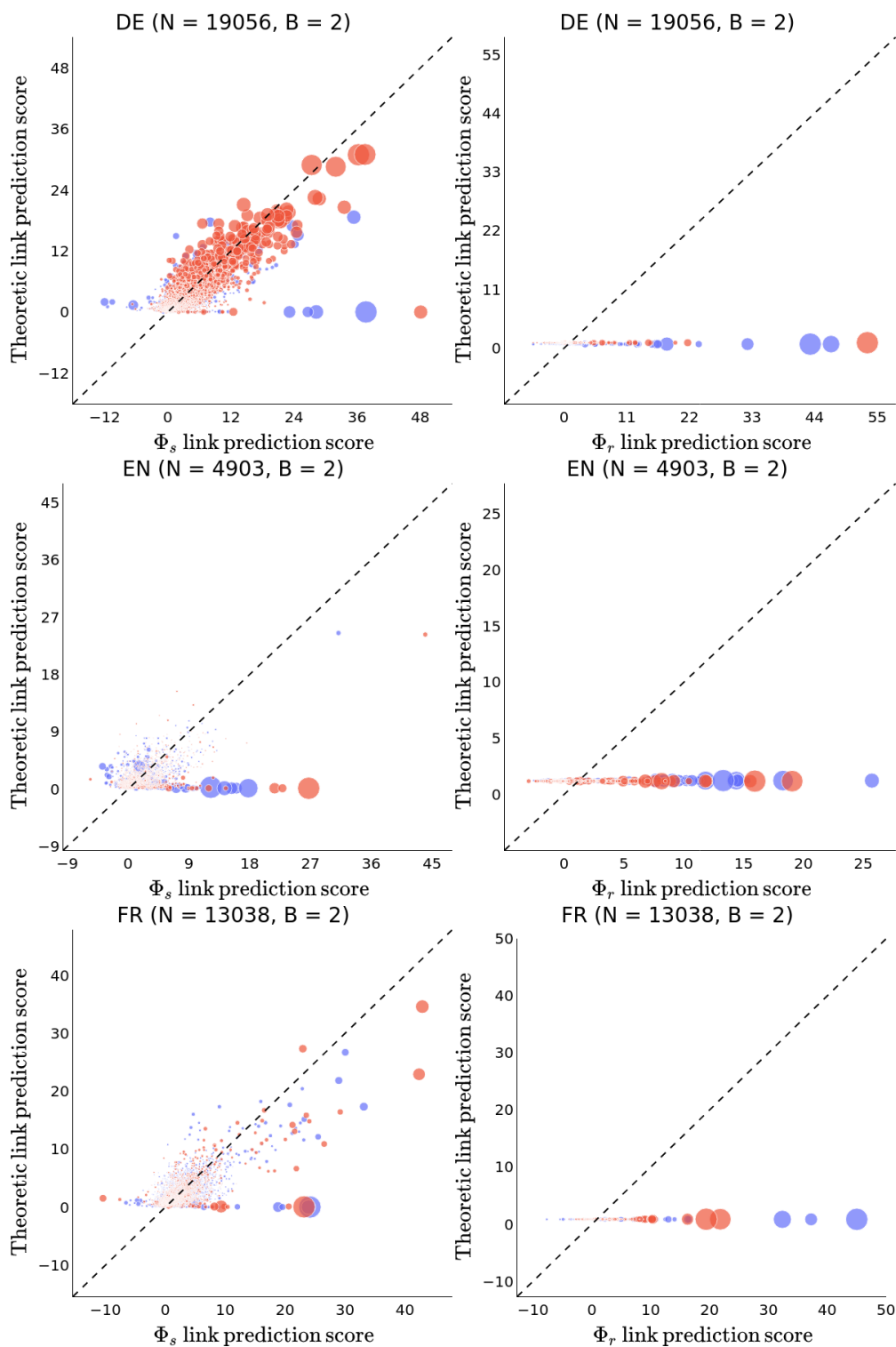


Figure B.3: Theoretic vs. GCN LP scores for online social network datasets.

B.7 Additional Experiments

B.7.1 Additional Experiments for §3.6.1 (4-layer Encoders)

We run the experiments from §3.6.1 for Φ_s with the same settings, except we use 4-layer (instead of 2-layer) encoders (128-dimensional hidden layers, 64-dimensional output layer). We run these additional experiments because the error bound for the theoretic LP scores for Φ_s depends on the number of encoder layers L . We find that the experimental results continue to support our theoretical analysis, both qualitatively and quantitatively (see Table B.2, Figure B.4); the NRMSE and PCC values are comparable to or better than those from the experiments with the 2-layer encoders (especially for the EN dataset).

Table B.2: The test AUC of the 4-layer Φ_s encoders on the real-world network datasets, and the NRMSE and PCC of the theoretic LP scores as predictors of the Φ_s scores.

	NRMSE (\downarrow)	PCC (\uparrow)	Φ_s Test AUC (\uparrow)
CORA	0.044 ± 0.006	0.858 ± 0.026	0.853 ± 0.028
CITeseer	0.057 ± 0.006	0.890 ± 0.017	0.861 ± 0.026
DBLP	0.021 ± 0.002	0.885 ± 0.054	0.887 ± 0.019
PUBMED	0.056 ± 0.009	0.802 ± 0.024	0.900 ± 0.006
CS	0.039 ± 0.006	0.918 ± 0.008	0.949 ± 0.004
PHYSICS	0.030 ± 0.002	0.077 ± 0.013	0.950 ± 0.004
LASTFMASIA	0.040 ± 0.004	0.938 ± 0.005	0.949 ± 0.002
DE	0.014 ± 0.003	0.918 ± 0.025	0.882 ± 0.002
EN	0.034 ± 0.005	0.752 ± 0.036	0.846 ± 0.008
FR	0.019 ± 0.003	0.833 ± 0.038	0.896 ± 0.003

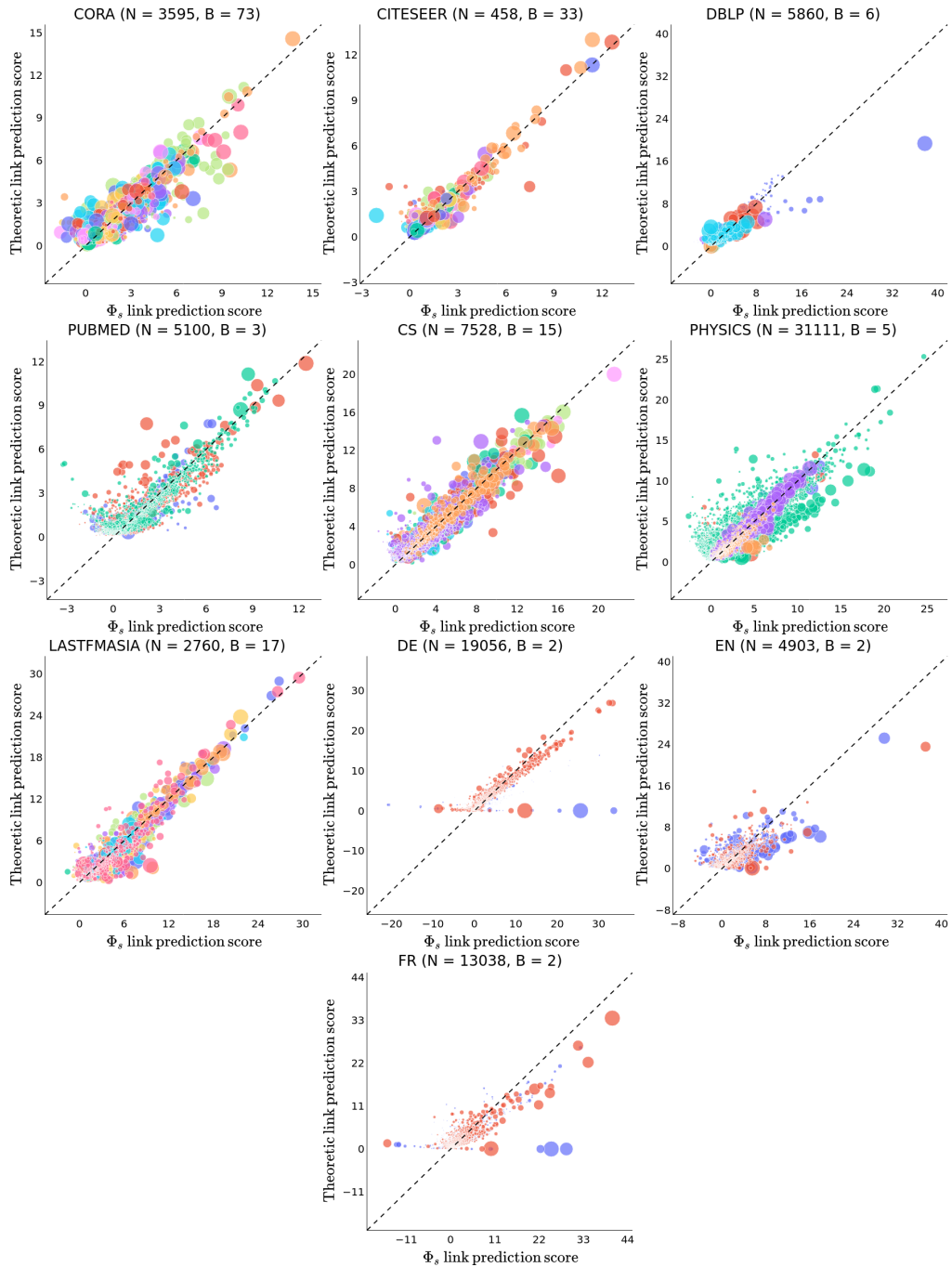


Figure B.4: Theoretic LP score vs. 4-layer Φ_s LP score for all network datasets.

B.7.2 Additional Experiments for §3.6.1 (Hadamard Product and MLP LP Score Function)

We also run the experiments from §3.6.1 for Φ_s with the same settings, except we use the following LP score function:

$$f_{LP}(\mathbf{h}_i^{(L)}, \mathbf{h}_j^{(L)}) = f_{MLP}(\mathbf{h}_i^{(L)} \odot \mathbf{h}_j^{(L)}), \quad (\text{B.52})$$

where \odot is the Hadamard product and f_{MLP} is a 2-layer MLP with a 64-dimensional hidden layer and ReLU nonlinearity. We run these additional experiments because a Hadamard product and MLP score function is often used in the literature. We find that that our theoretical analysis is still relevant to and reasonably supports the experimental results, both qualitatively and quantitatively (see Table B.3, Figure B.5). This could be because MLPs have an inductive bias towards learning simpler, often linear functions [NKK19, VCL19], and our theoretical findings are generalizable to linear LP score functions. Notably, in this setting, Φ_s makes a higher number of negative link predictions. For a few datasets (e.g., Cora, CiteSeer, LastFMAsia), a handful of theoretic LP scores are negative because the regression (incorrectly) predicts $\bar{\rho}_s^2(b)$ for 1-2 groups $S^{(b)}$ to be negative.

Table B.3: The test AUC of the Φ_s encoders with an f_{MLP} score function on the real-world network datasets, and the NRMSE and PCC of the theoretic LP scores as predictors of the Φ_s scores.

	NRMSE (\downarrow)	PCC (\uparrow)	Φ_s Test AUC (\uparrow)
CORA	0.034 ± 0.004	0.830 ± 0.015	0.915 ± 0.001
CITeseer	0.090 ± 0.014	0.365 ± 0.070	0.913 ± 0.008
DBLP	0.026 ± 0.003	0.652 ± 0.029	0.933 ± 0.004
PUBMED	0.054 ± 0.007	0.813 ± 0.038	0.932 ± 0.003
CS	0.047 ± 0.008	0.677 ± 0.036	0.970 ± 0.001
PHYSICS	0.055 ± 0.007	0.566 ± 0.026	0.976 ± 0.001
LASTFMASIA	0.049 ± 0.008	0.682 ± 0.035	0.960 ± 0.003
DE	0.030 ± 0.008	0.683 ± 0.047	0.935 ± 0.001
EN	0.039 ± 0.006	0.463 ± 0.022	0.905 ± 0.002
FR	0.031 ± 0.006	0.654 ± 0.067	0.935 ± 0.002

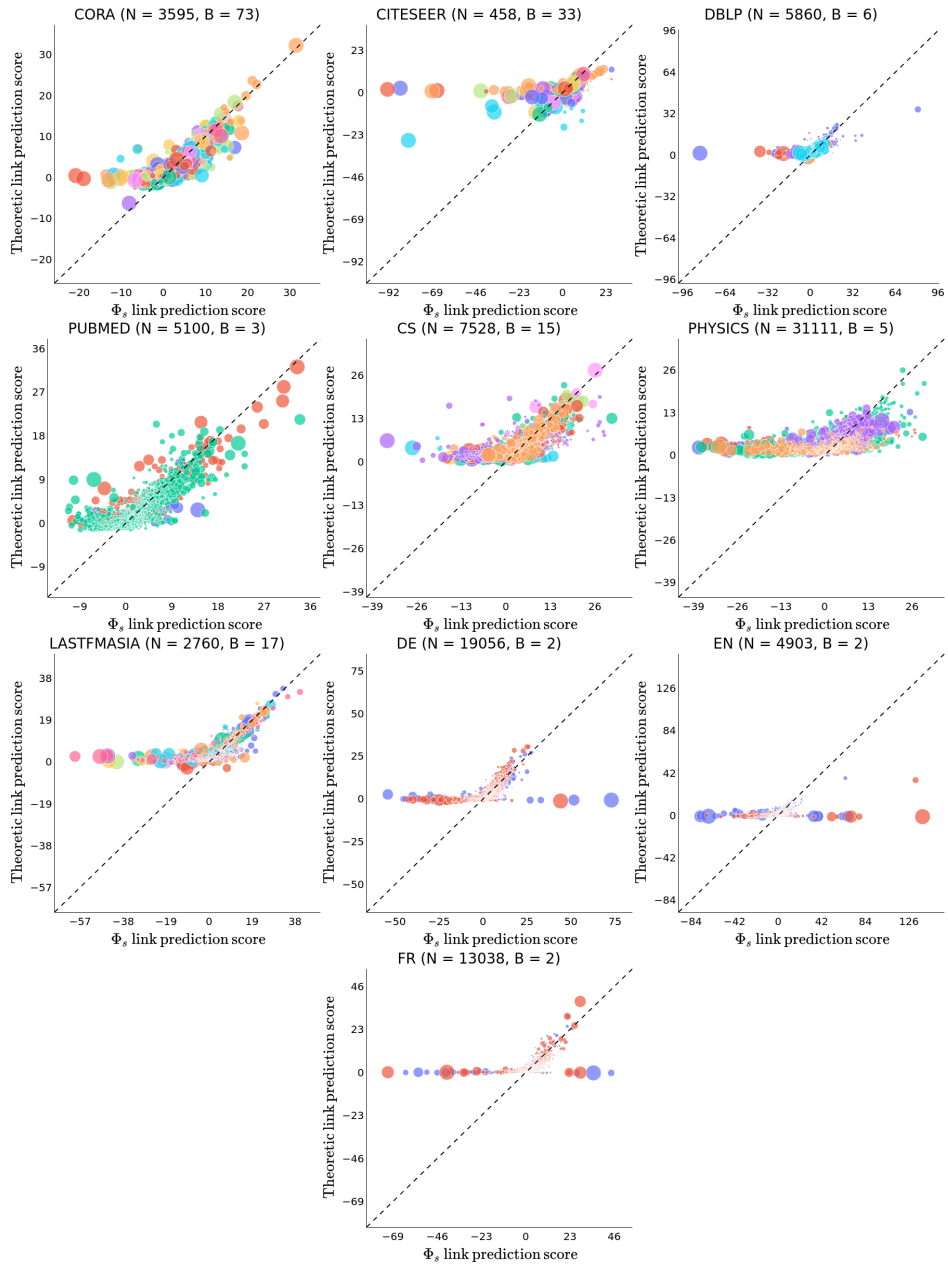


Figure B.5: Theoretic LP score vs. Φ_s LP score (with Hadamard product and MLP) for all network datasets.

B.7.3 Additional Experiments for §3.6.2

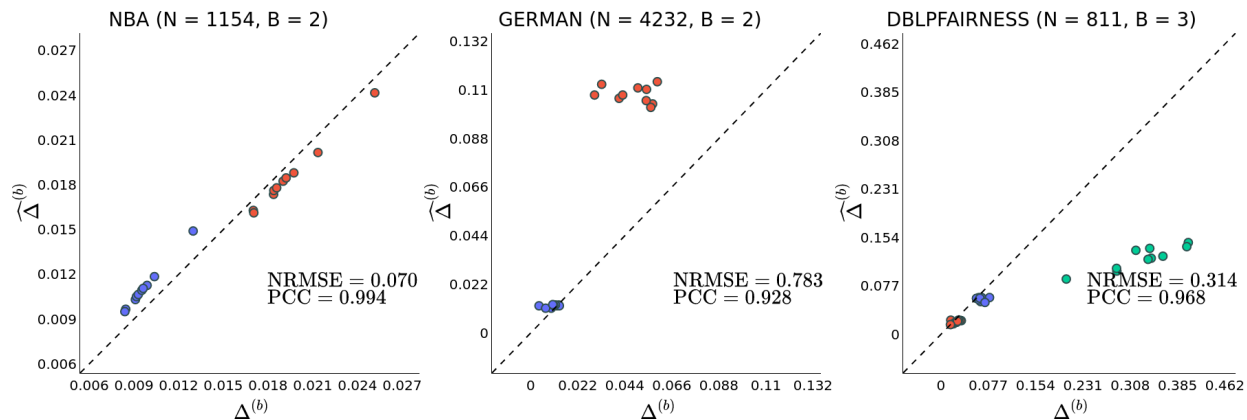


Figure B.6: The plots display $\widehat{\Delta}^{(b)}$ vs. $\Delta^{(b)}$ for 4-layer Φ_s for the NBA, German, and DBLP-Fairness datasets over all $b \in [B]$ and 10 random seeds.

B.7.4 Additional Experiments for §3.6.3

Table B.4: $\frac{1}{B} \sum_{b \in [B]} \Delta^{(b)}$ and the test AUC for the NBA, German, and DBLP-Fairness datasets with various settings of λ_{fair} . The **left** table corresponds to 4-layer Φ_s , and the **right** to 4-layer Φ_r .

	λ_{fair}	$\frac{1}{B} \sum_{b \in [B]} \Delta^{(b)}$ (\downarrow)	Φ_s Test AUC (\uparrow)		λ_{fair}	$\frac{1}{B} \sum_{b \in [B]} \Delta^{(b)}$ (\downarrow)	Φ_r Test AUC (\uparrow)
NBA	4.0	0.000 \pm 0.000	0.752 \pm 0.001	NBA	4.0	0.000 \pm 0.000	0.581 \pm 0.029
NBA	2.0	0.006 \pm 0.001	0.752 \pm 0.001	NBA	2.0	0.000 \pm 0.000	0.574 \pm 0.021
NBA	1.0	0.011 \pm 0.001	0.753 \pm 0.001	NBA	1.0	0.000 \pm 0.000	0.580 \pm 0.025
NBA	0.0	0.014 \pm 0.001	0.753 \pm 0.001	NBA	0.0	0.000 \pm 0.000	0.589 \pm 0.031
DBLPFAIRNESS	4.0	0.090 \pm 0.041	0.793 \pm 0.009	DBLPFAIRNESS	4.0	0.034 \pm 0.012	0.769 \pm 0.009
DBLPFAIRNESS	2.0	0.070 \pm 0.015	0.800 \pm 0.007	DBLPFAIRNESS	2.0	0.045 \pm 0.021	0.788 \pm 0.007
DBLPFAIRNESS	1.0	0.099 \pm 0.009	0.804 \pm 0.007	DBLPFAIRNESS	1.0	0.074 \pm 0.013	0.797 \pm 0.006
DBLPFAIRNESS	0.0	0.122 \pm 0.028	0.820 \pm 0.009	DBLPFAIRNESS	0.0	0.095 \pm 0.015	0.811 \pm 0.006
GERMAN	4.0	0.012 \pm 0.008	0.817 \pm 0.004	GERMAN	4.0	0.027 \pm 0.009	0.765 \pm 0.013
GERMAN	2.0	0.018 \pm 0.007	0.827 \pm 0.015	GERMAN	2.0	0.023 \pm 0.007	0.765 \pm 0.011
GERMAN	1.0	0.018 \pm 0.008	0.856 \pm 0.025	GERMAN	1.0	0.031 \pm 0.010	0.786 \pm 0.030
GERMAN	0.0	0.028 \pm 0.007	0.874 \pm 0.011	GERMAN	0.0	0.030 \pm 0.009	0.838 \pm 0.025

B.8 Theory Pitfalls

To understand the second pitfall from §3.6.1, we separately investigate the association between the within-group degree product $(\widehat{D}_{ii}\widehat{D}_{jj})$ and the absolute deviation of the theoretic LP scores from the Φ_s scores, as well as the association between the (transformed) feature similarity $\left(\left\|\sum_{k \in S^{(b)}} \frac{\sqrt{\widehat{D}_{kk}}}{\text{vol}(G^{(b)})} \alpha_k\right\|_2^2\right)$ and the absolute deviation (see Figure B.7). We observe that the absolute deviation is highest for the node pairs with a relatively small degree product (i.e., nodes with a low PA score) and low feature similarity.

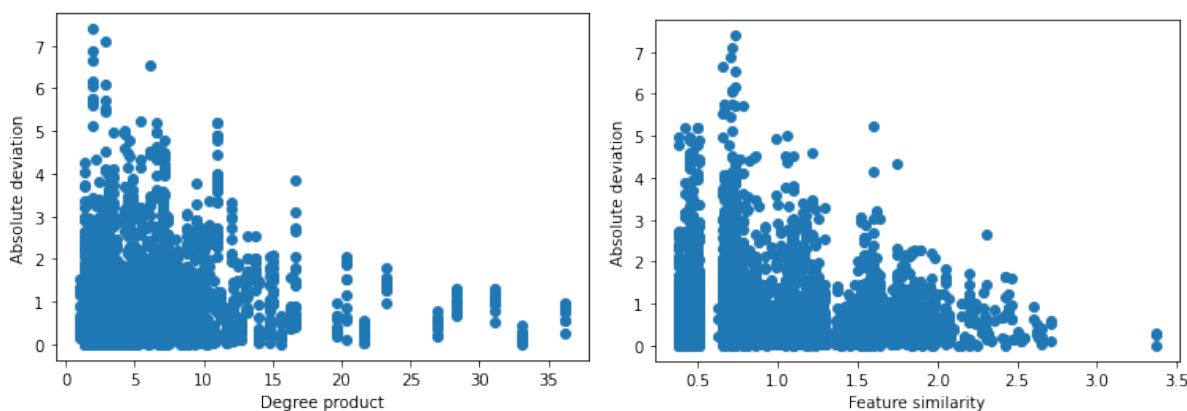


Figure B.7: Associations of absolute deviation with degree product and with feature similarity for CiteSeer.

B.9 Error Analysis of Φ_r Theoretic Scores

Figure B.8 reveals that the max term $\max_{u,v \in \mathcal{V}} \sqrt{\frac{\widehat{D}_{vv}}{\widehat{D}_{uu}}}$ is quite large in practice, which causes the theoretic LP scores to generally be poor estimates for the Φ_r scores. We additionally find in Figure B.8 that the relative error (as measured by NRMSE and PCC) of the theoretic LP scores for Φ_r is not lower for lower values of the max term $\max_{u,v \in \mathcal{V}} \sqrt{\frac{\widehat{D}_{vv}}{\widehat{D}_{uu}}}$.

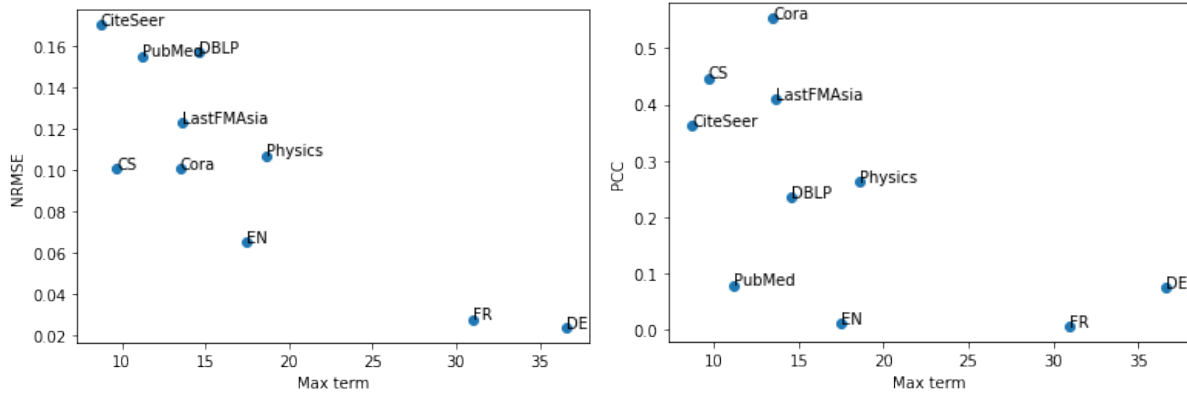


Figure B.8: Weak associations of max term with NRMSE and PCC of theoretic LP scores for Φ_r across all datasets described in §B.3.

Furthermore, Figure B.9 reveals that Φ_r LP scores are *not* higher for incident nodes with larger degrees.

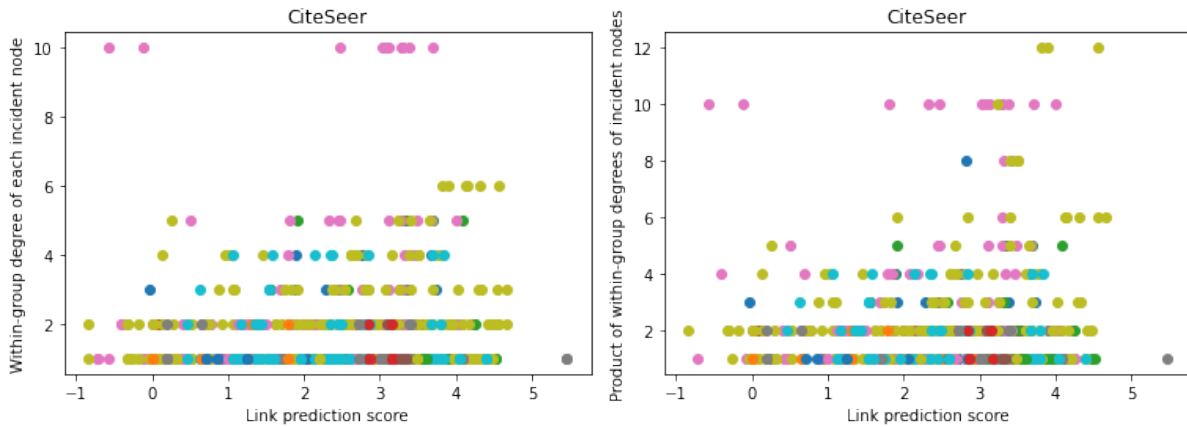


Figure B.9: Weak associations of mean Φ_r LP scores (over 10 random seeds) with degree of each incident node and product of degrees of both incident nodes. Colors correspond to different groups.

There are intimate connections between Theorem 3.4.4 and the steady-state probabilities of random walks. The stationary probabilities of random walks are the same regardless of the starting node. This is why Φ_r produces similar representations for all the nodes in each

social group, regardless of the degree of the node; in fact, with a larger number of layers, Φ_r would oversmooth all the representations to the same vector [Ker22]. Hence, Φ_r LP scores do not have a degree dependence, theoretically or empirically.

APPENDIX C

Appendix for Chapter 4

C.1 Overview of Theoretical Analyses of and Hypotheses for Degree Bias

C.1.1 Theoretical Analyses of Degree Bias

Table C.1: A taxonomy of GNN degree bias papers based on whether they theoretically analyze the origins of degree bias, explicitly linking a node’s degree to its test and training error.

Explicit theoretical analysis of origins of degree bias?	Papers
Yes	[WHX19], [MLS22], [LMM23]
No	[TYS20], [LNF21], [WWF21], [XCW21], [FHH21], [ZDW21], [WWS22], [LYD22], [KZX22], [YKY22], [ZMW22], [CXK22], [BKE22], [LZX22], [LXS23], [SJW23], [LNF23], [LFZ24], [JZY23], [LLC23], [HLS23], [CWC23], [HZW24], [XXH23], [WLL23], [LXC23], [ZLP23], [VS23], [CLY23], [DSL23], [HL23a], [ZCY24], [ZJ24], [XHZ24], [ZZY24]

C.1.2 Hypotheses for Degree Bias

Table C.2: Full taxonomy of the hypotheses for the origins of GNN degree bias proposed by papers.

Hypothesis	Papers
(H1) Neighborhoods of low-degree nodes contain insufficient or overly noisy information for effective representations.	[LNF21], [WWF21], [XCW21], [FHH21], [ZDW21], [LYD22], [LZX22], [LNF23], [LFZ24], [JZY23], [LLC23], [HZW24], [LXS23], [XXH23], [ZLP23], [VS23], [DSL23], [CLY23], [HL23a], [ZJ24], [XHZ24]
(H2) High-degree nodes have a larger influence on GNN training because they have a greater number of links with other nodes, thereby dominating message passing.	[TYS20], [WWF21], [ZDW21], [KZX22], [ZLY22], [LXS23], [ZCY24]
(H3) High-degree nodes exert more influence on the representations of and predictions for nodes as the number of GNN layers increases.	[ZDW21], [CWC23], [LXC23], [DSL23], [ZZY24]
(H4) In semi-supervised learning, if training nodes are picked randomly, test predictions for high-degree nodes are more likely to be influenced by these training nodes because they have a greater number of links with other nodes.	[TYS20], [ZLY22], [HLS23]
(H5) Representations of high-degree nodes cluster more strongly around their corresponding class centers, or are more likely to be linearly separable.	[MLS22], [WWS22], [LMM23]
(H6) Neighborhoods of high-degree nodes contain more homophilic links, enhancing their representations.	[LLC23], [XXH23]
(H7) Nodes with different degrees are not necessarily mapped to distinct representations.	[WHX19]
(H8) Low-degree nodes have class-imbalanced training samples, yielding worse generalization.	[YKY22]
(H9) High-degree nodes are more likely to be labeled during training and thus GNNs generalize better for them.	[CXK22]
(H10) Representations of high-degree nodes have higher variance.	[LXS23]
(H11) Low-degree nodes are more likely to be sampled during training/inference.	[VS23]

C.2 Proofs

C.2.1 Theorem 4.4.1

Proof. Misclassification occurs when $\ell(\mathcal{M}|i, c) > \ell(\mathcal{M}|i, c')$.

$$\mathbb{P}(\ell(\mathcal{M}|i, c) > \ell(\mathcal{M}|i, c')) = \mathbb{P}\left(-\log \mathbf{H}_{i,c}^{(L)} > -\log \mathbf{H}_{i,c'}^{(L)}\right) \quad (\text{C.1})$$

$$= \mathbb{P}\left(\mathbf{H}_{i,c}^{(L)} < \mathbf{H}_{i,c'}^{(L)}\right) \quad (\text{C.2})$$

$$= \mathbb{P}\left(\mathbf{Z}_{i,c'}^{(L)} - \mathbf{Z}_{i,c}^{(L)} > 0\right). \quad (\text{C.3})$$

If $\mathbb{E}\left[\mathbf{Z}_{i,c'}^{(L)} - \mathbf{Z}_{i,c}^{(L)}\right] < 0$ (i.e., \mathcal{M} generalizes in expectation), by Cantelli's inequality:

$$\mathbb{P}(\ell(\mathcal{M}|i, c) > \ell(\mathcal{M}|i, c')) = \mathbb{P}\left(\left(\mathbf{Z}_{i,c'}^{(L)} - \mathbf{Z}_{i,c}^{(L)}\right) - \mathbb{E}\left[\mathbf{Z}_{i,c'}^{(L)} - \mathbf{Z}_{i,c}^{(L)}\right] > -\mathbb{E}\left[\mathbf{Z}_{i,c'}^{(L)} - \mathbf{Z}_{i,c}^{(L)}\right]\right) \quad (\text{C.4})$$

$$\leq \frac{1}{1 + \frac{\left(-\mathbb{E}\left[\mathbf{Z}_{i,c'}^{(L)} - \mathbf{Z}_{i,c}^{(L)}\right]\right)^2}{\text{Var}\left[\mathbf{Z}_{i,c'}^{(L)} - \mathbf{Z}_{i,c}^{(L)}\right]}}. \quad (\text{C.5})$$

We use Cantelli's inequality, rather than Chebyshev's inequality, because Cantelli's inequality is sharper for one-sided bounds.

□

C.2.2 Theorem 4.4.2

Proof. Denoting the l -th term in the summation $\mathbf{T}^{(l)} = \mathbf{P}_{\text{rw}}^l \mathbf{X} \mathbf{W}^{(l)}$, $\mathbf{T}_{i,c}^{(l)} = \sum_{j \in \mathcal{V}} (\mathbf{P}_{\text{rw}}^l)_{ij} \mathbf{X}_j \mathbf{W}_{\cdot,c}^{(l)}$. It follows by the linearity of expectation that:

$$\mathbb{E}\left[\mathbf{T}_{i,c'}^{(l)} - \mathbf{T}_{i,c}^{(l)}\right] = \sum_{j \in \mathcal{V}} (\mathbf{P}_{\text{rw}}^l)_{ij} \cdot \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{Y_j}} \left[\mathbf{x}^T \mathbf{w}_{c'-c}^{(l)}\right] \quad (\text{C.6})$$

$$= \mathbb{E}_{j \sim \mathcal{N}^{(l)}(i)} \left[\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{Y_j}} \left[\mathbf{x}^T \mathbf{w}_{c'-c}^{(l)}\right]\right] \quad (\text{C.7})$$

$$= \beta_{i,c'}^{(l)}. \quad (\text{C.8})$$

Furthermore, by the linearity of variance:

$$\text{Var} \left[\mathbf{T}_{i,c'}^{(l)} - \mathbf{T}_{i,c}^{(l)} \right] = \sum_{j \in \mathcal{V}} \left[(\mathbf{P}_{\text{rw}}^l)_{ij} \right]^2 \cdot \text{Var}_{\mathbf{x} \sim \mathcal{D}_{Y_j}} \left[\mathbf{x}^T \mathbf{w}_{c'-c}^{(l)} \right] \quad (\text{C.9})$$

$$\leq M \sum_{j \in \mathcal{V}} \left[(\mathbf{P}_{\text{rw}}^l)_{ij} \right]^2 \quad (\text{C.10})$$

$$= M \alpha_i^{(l)}. \quad (\text{C.11})$$

Then, once again by the linearity of expectation and variance:

$$\left(\mathbb{E} \left[\mathbf{Z}_{i,c'}^{(l)} - \mathbf{Z}_{i,c}^{(l)} \right] \right)^2 = \left(\sum_{l=0}^L \beta_{i,c'}^{(l)} \right)^2, \quad (\text{C.12})$$

$$\text{Var} \left[\mathbf{Z}_{i,c'}^{(l)} - \mathbf{Z}_{i,c}^{(l)} \right] \leq M(L+1) \sum_{l=0}^L \alpha_i^{(l)}. \quad (\text{C.13})$$

Consequently:

$$\frac{\left(\mathbb{E} \left[\mathbf{Z}_{i,c'}^{(l)} - \mathbf{Z}_{i,c}^{(l)} \right] \right)^2}{\text{Var} \left[\mathbf{Z}_{i,c'}^{(l)} - \mathbf{Z}_{i,c}^{(l)} \right]} \geq \frac{\left(\sum_{l=0}^L \beta_{i,c'}^{(l)} \right)^2}{M(L+1) \sum_{l=0}^L \alpha_i^{(l)}}. \quad (\text{C.14})$$

□

C.2.3 Theorem 4.4.3

Proof. Re-expressing the l -th term $\mathbf{T}^{(l)} = \mathbf{P}_{\text{sym}}^l \mathbf{X} \mathbf{W}^{(l)}$ in the summation:

$$\mathbf{T}_{i,c}^{(l)} = \sum_{j \in \mathcal{V}} \left(\mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} \right)_{ij}^l \mathbf{X}_j \mathbf{W}_{\cdot,c}^{(l)} \quad (\text{C.15})$$

$$= \sum_{j \in \mathcal{V}} (\mathbf{D}^{-1} \mathbf{A})_{ij}^l \cdot \frac{\sqrt{\mathbf{D}_{ii}}}{\sqrt{\mathbf{D}_{jj}}} \mathbf{X}_j \mathbf{W}_{\cdot,c}^{(l)} \quad (\text{C.16})$$

$$= \sqrt{\mathbf{D}_{ii}} \sum_{j \in \mathcal{V}} (\mathbf{P}_{\text{rw}}^l)_{ij} \cdot \frac{1}{\sqrt{\mathbf{D}_{jj}}} \mathbf{X}_j \mathbf{W}_{\cdot,c}^{(l)}. \quad (\text{C.17})$$

It follows by the linearity of expectation that:

$$\mathbb{E} \left[\mathbf{T}_{i,c'}^{(l)} - \mathbf{T}_{i,c}^{(l)} \right] = \sqrt{\mathbf{D}_{ii}} \sum_{j \in \mathcal{V}} (\mathbf{P}_{\text{rw}}^l)_{ij} \cdot \frac{1}{\sqrt{\mathbf{D}_{jj}}} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\mathcal{Y}_j}} \left[\mathbf{x}^T \mathbf{w}_{c'-c}^{(l)} \right] \quad (\text{C.18})$$

$$= \sqrt{\mathbf{D}_{ii}} \mathbb{E}_{j \sim \mathcal{N}^{(l)}(i)} \left[\frac{1}{\sqrt{\mathbf{D}_{jj}}} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\mathcal{Y}_j}} \left[\mathbf{x}^T \mathbf{w}_{c'-c}^{(l)} \right] \right] \quad (\text{C.19})$$

$$= \sqrt{\mathbf{D}_{ii}} \tilde{\beta}_{i,c'}^{(l)}. \quad (\text{C.20})$$

Furthermore, by the linearity of variance:

$$\text{Var} \left[\mathbf{T}_{i,c'}^{(l)} - \mathbf{T}_{i,c}^{(l)} \right] = \mathbf{D}_{ii} \sum_{j \in \mathcal{V}} \left[(\mathbf{P}_{\text{rw}}^l)_{ij} \right]^2 \cdot \frac{1}{\mathbf{D}_{jj}} \text{Var}_{\mathbf{x} \sim \mathcal{D}_{\mathcal{Y}_j}} \left[\mathbf{x}^T \mathbf{w}_{c'-c}^{(l)} \right] \quad (\text{C.21})$$

$$\leq \mathbf{D}_{ii} M \sum_{j \in \mathcal{V}} \frac{1}{\mathbf{D}_{jj}} \left[(\mathbf{P}_{\text{rw}}^l)_{ij} \right]^2 \quad (\text{C.22})$$

$$= \mathbf{D}_{ii} M \tilde{\alpha}_i^{(l)}. \quad (\text{C.23})$$

Then, once again, by the linearity of expectation and variance:

$$\left(\mathbb{E} \left[\mathbf{Z}_{i,c'}^{(l)} - \mathbf{Z}_{i,c}^{(l)} \right] \right)^2 = \mathbf{D}_{ii} \left(\sum_{l=0}^L \tilde{\beta}_{i,c'}^{(l)} \right)^2, \quad (\text{C.24})$$

$$\text{Var} \left[\mathbf{Z}_{i,c'}^{(l)} - \mathbf{Z}_{i,c}^{(l)} \right] \leq \mathbf{D}_{ii} M (L+1) \sum_{l=0}^L \tilde{\alpha}_i^{(l)}. \quad (\text{C.25})$$

Consequently:

$$\frac{\left(\mathbb{E} \left[\mathbf{Z}_{i,c'}^{(l)} - \mathbf{Z}_{i,c}^{(l)} \right] \right)^2}{\text{Var} \left[\mathbf{Z}_{i,c'}^{(l)} - \mathbf{Z}_{i,c}^{(l)} \right]} \geq \frac{\left(\sum_{l=0}^L \tilde{\beta}_{i,c'}^{(l)} \right)^2}{M(L+1) \sum_{l=0}^L \tilde{\alpha}_i^{(l)}}. \quad (\text{C.26})$$

□

C.2.4 Lemma 4.5.1

Proof. Define $g(\mathbf{Z}_i^{(L)}[t]) = \nabla_{\mathbf{Z}_i^{(L)}[t]} \ell[t](\mathcal{M}|i, c)$. For simplicity of notation, let $\mathbf{x} = \mathbf{Z}_i^{(L)}[t]$. It is sufficient to show that $\|g(\mathbf{x})\|_2 \leq \lambda$. By simple derivation, $(g(\mathbf{x}))_i = -\frac{\sum_{a \neq i} e^{\mathbf{x}_a}}{\sum_b e^{\mathbf{x}_b}}$, and for

$j \neq i$, $(g(\mathbf{x}))_j = -\frac{e^{x_j}}{\sum_b e^{x_b}}$. Then, by Hölder's inequality:

$$\|g(\mathbf{x})\|_2^2 = \frac{\left(\sum_{a \neq i} e^{x_a}\right)^2 + \sum_{a \neq i} (e^{x_a})^2}{\left(\sum_b e^{x_b}\right)^2} \quad (\text{C.27})$$

$$\leq \frac{2\left(\sum_{a \neq i} e^{x_a}\right)^2}{\left(\sum_b e^{x_b}\right)^2} \leq 2. \quad (\text{C.28})$$

Thus, $\lambda = \sqrt{2}$. □

C.2.5 Theorem 4.5.2

Proof. By the Lipschitz continuity of $\ell[t](\overline{\text{SYM}}|i, c)$ (Lemma 4.5.1) and the triangle inequality:

$$|\ell[t+1](\overline{\text{SYM}}|i, c) - \ell[t](\overline{\text{SYM}}|i, c)| \leq \lambda \left\| \mathbf{z}_i^{(L)}[t+1] - \mathbf{z}_i^{(L)}[t] \right\|_2 \quad (\text{C.29})$$

$$\leq \lambda \sum_{l=0}^L \left\| (\mathbf{P}_{\text{sym}}^l \mathbf{X})_i (\mathbf{W}^{(l)}[t+1] - \mathbf{W}^{(l)}[t]) \right\|_2 \quad (\text{C.30})$$

$$= \lambda \eta \sum_{l=0}^L \left\| (\mathbf{P}_{\text{sym}}^l \mathbf{X})_i \frac{\partial \ell[t]}{\partial \mathbf{W}^{(l)}[t]} (B[t]) \right\|_2. \quad (\text{C.31})$$

By simple derivation, we see that $\frac{\partial \ell[t]}{\partial \mathbf{W}^{(l)}[t]} (B[t]) = (\mathbf{P}_{\text{sym}}^l[t] \mathbf{X})^T \epsilon[t]$, where $\mathbf{P}_{\text{sym}}^l[t] \in \mathbb{R}^{|B[t]| \times N}$ is the submatrix formed from the rows of $\mathbf{P}_{\text{sym}}^l$ that correspond to the nodes in $B[t]$. Then, by the sub-multiplicativity of the L_2 norm:

$$|\ell[t+1](\overline{\text{SYM}}|i, c) - \ell[t](\overline{\text{SYM}}|i, c)| \leq \lambda \eta \sum_{l=0}^L \left\| (\mathbf{P}_{\text{sym}}^l)_i \mathbf{X} \mathbf{X}^T (\mathbf{P}_{\text{sym}}^l[t])^T \epsilon[t] \right\|_2 \quad (\text{C.32})$$

$$\leq \lambda \eta \|\epsilon[t]\|_F \sum_{l=0}^L \left\| (\mathbf{P}_{\text{sym}}^l)_i \mathbf{X} \mathbf{X}^T (\mathbf{P}_{\text{sym}}^l[t])^T \right\|_2. \quad (\text{C.33})$$

Similarly to the proof of Theorem 4.4.3:

$$(\mathbf{P}_{\text{sym}}^l)_i \mathbf{X} \mathbf{X}^T (\mathbf{P}_{\text{sym}}^l)_m^T = \sum_{j \in \mathcal{V}} (\mathbf{P}_{\text{sym}}^l)_{ij} \sum_{k \in \mathcal{V}} (\mathbf{P}_{\text{sym}}^l)_{mk} (\mathbf{X} \mathbf{X}^T)_{jk} \quad (\text{C.34})$$

$$= \sqrt{\mathbf{D}_{ii} \mathbf{D}_{mm}} \mathbb{E}_{j \sim \mathcal{N}^{(l)}(i), k \sim \mathcal{N}^{(l)}(m)} \left[\frac{1}{\sqrt{\mathbf{D}_{jj} \mathbf{D}_{kk}}} (\mathbf{X} \mathbf{X}^T)_{jk} \right]. \quad (\text{C.35})$$

Hence, $\left\| (\mathbf{P}_{\text{sym}}^l)_i \mathbf{X} \mathbf{X}^T (\mathbf{P}_{\text{sym}}^l[t])^T \right\|_2 = \sqrt{\mathbf{D}_{ii}} \cdot \left\| \tilde{\chi}_i^{(l)}[t] \right\|_2$, and:

$$|\ell[t+1](\overline{\text{SYM}}|i, c) - \ell[t](\overline{\text{SYM}}|i, c)| \leq \sqrt{\mathbf{D}_{ii}} \cdot \lambda \eta \|\epsilon[t]\|_F \sum_{l=0}^L \left\| \tilde{\chi}_i^{(l)}[t] \right\|_2. \quad (\text{C.36})$$

□

C.2.6 Theorem C.8.1

Proof. By the Lipschitz continuity of $\ell[t](\overline{\text{RW}}|i, c)$ (Lemma 4.5.1) and the triangle inequality:

$$|\ell[t+1](\overline{\text{RW}}|i, c) - \ell[t](\overline{\text{RW}}|i, c)| \leq \lambda \left\| \mathbf{z}_i^{(L)}[t+1] - \mathbf{z}_i^{(L)}[t] \right\|_2 \quad (\text{C.37})$$

$$\leq \lambda \sum_{l=0}^L \left\| (\mathbf{P}_{\text{rw}}^l \mathbf{X})_i (\mathbf{W}^{(l)}[t+1] - \mathbf{W}^{(l)}[t]) \right\|_2 \quad (\text{C.38})$$

$$= \lambda \eta \sum_{l=0}^L \left\| (\mathbf{P}_{\text{rw}}^l \mathbf{X})_i \frac{\partial \ell[t]}{\partial \mathbf{W}^{(l)}[t]}(B[t]) \right\|_2. \quad (\text{C.39})$$

By simple derivation, we see that $\frac{\partial \ell[t]}{\partial \mathbf{W}^{(l)}[t]}(B[t]) = (\mathbf{P}_{\text{rw}}^l[t] \mathbf{X})^T \epsilon[t]$, where $\mathbf{P}_{\text{rw}}^l[t] \in \mathbb{R}^{|B[t]| \times N}$ is the submatrix formed from the rows of \mathbf{P}_{rw}^l that correspond to the nodes in $B[t]$. Then, by the sub-multiplicativity of the L_2 norm:

$$|\ell[t+1](\overline{\text{RW}}|i, c) - \ell[t](\overline{\text{RW}}|i, c)| \leq \lambda \eta \sum_{l=0}^L \left\| (\mathbf{P}_{\text{rw}}^l)_i \mathbf{X} \mathbf{X}^T (\mathbf{P}_{\text{rw}}^l[t])^T \epsilon[t] \right\|_2 \quad (\text{C.40})$$

$$\leq \lambda \eta \|\epsilon[t]\|_F \sum_{l=0}^L \left\| (\mathbf{P}_{\text{rw}}^l)_i \mathbf{X} \mathbf{X}^T (\mathbf{P}_{\text{rw}}^l[t])^T \right\|_2. \quad (\text{C.41})$$

Similarly to the proof of Theorem 4.4.2:

$$(\mathbf{P}_{\text{rw}}^l)_i \mathbf{X} \mathbf{X}^T (\mathbf{P}_{\text{rw}}^l)_m^T = \sum_{j \in \mathcal{V}} (\mathbf{P}_{\text{rw}}^l)_{ij} \sum_{k \in \mathcal{V}} (\mathbf{P}_{\text{rw}}^l)_{mk} (\mathbf{X} \mathbf{X}^T)_{jk} \quad (\text{C.42})$$

$$= \mathbb{E}_{j \sim \mathcal{N}^{(l)}(i), k \sim \mathcal{N}^{(l)}(m)} \left[(\mathbf{X} \mathbf{X}^T)_{jk} \right]. \quad (\text{C.43})$$

Hence, $\left\| (\mathbf{P}_{\text{rw}}^l)_i \mathbf{X} \mathbf{X}^T (\mathbf{P}_{\text{rw}}^l[t])^T \right\|_2 = \left\| \chi_i^{(l)}[t] \right\|_2$, and:

$$|\ell[t+1](\overline{\text{RW}}|i, c) - \ell[t](\overline{\text{RW}}|i, c)| \leq \lambda \eta \|\epsilon[t]\|_F \sum_{l=0}^L \left\| \chi_i^{(l)}[t] \right\|_2. \quad (\text{C.44})$$

□

C.3 Datasets

In our experiments, we use 8 real-world network datasets from [BG18a], [SMB19], and [RAS21], covering diverse domains (e.g., citation networks, collaboration networks, online product networks, Wikipedia networks). We provide a description and statistics of each dataset in Table C.3. All the datasets have node features and are undirected. For each node, we normalize its features to sum to 1, following [FL19]¹. We were unable to find the exact class names and their label correspondence from the dataset documentation.

- In all the citation network datasets, nodes represent documents, edges represent citation links, and features are a binary bag-of-words representation of documents. The classification task is to predict the topic of documents.
- In the collaboration network datasets, nodes represent authors, edges represent coauthorships, and features are a binary bag-of-word representation of keywords from the authors' papers. The classification task is to predict the most active field of study for authors.
- In the online product network datasets, nodes represent products, edges represent that two products are often purchased together, and features are a binary bag-of-word representation of product reviews. The classification task is to predict the category of products.
- In the Wikipedia network datasets, nodes represent Wikipedia websites, edges represent hyperlinks between them, and features are a binary bag-of-word representation of informative nouns from the pages. The classification task is to predict the level of average daily traffic for pages.

We use PyTorch and PyTorch Geometric to load and process all datasets [PGM19, FL19].

¹https://github.com/pyg-team/pytorch_geometric/blob/master/examples/link_pred.py

Table C.3: Summary of the datasets used in our experiments.

Name	Domain	# Nodes	# Edges	# Features	# Classes
Cora_ML	citation	2995	16316	2879	7
CiteSeer	citation	4230	10674	602	6
CS	collaboration	18333	163788	6805	15
Physics	collaboration	34493	495924	8415	5
Amazon Photo	online product	7650	238162	745	8
Amazon Computers	online product	13752	491722	767	10
chameleon	Wikipedia	2277	36101	2325	5
squirrel	Wikipedia	5201	217073	2089	5

C.4 Models

In our experiments, we transductively learn and compute node representations using encoders based on Graph Convolutional Networks (GCNs) [KW17], GraphSAGE [HYL17], and Graph Attention Networks (GATs) [VCC18].

In all cases, we use general message-passing GNNs \mathcal{M} [GRC23], which include separate parameters for source and target nodes and residual connections; in particular, for layer l :

$$\mathbf{H}^{(l)} = \sigma^{(l)}(\mathbf{Z}^{(l)}) = \sigma^{(l)}\left(\mathbf{H}^{(l-1)}\mathbf{W}_1^{(l)} + \mathbf{P}^{(l)}\mathbf{H}^{(l-1)}\mathbf{W}_2^{(l)} + \mathbf{X}\mathbf{W}_3^{(l)}\right), \quad (\text{C.45})$$

where $\mathbf{H}^{(l)} \in \mathbb{R}^{N \times d^{(l)}}$ are the l -th layer node representations (with $\mathbf{H}^{(0)} = \mathbf{X}$ and $d^{(L)} = C$), $\sigma^{(l)}$ is an instance-wise non-linearity (with $\sigma^{(L)}$ being softmax), $\mathbf{P}^{(l)} \in \mathbb{R}^{N \times N}$ is a graph filter, and $\mathbf{W}_1^{(l)}, \mathbf{W}_2^{(l)}, \mathbf{W}_3^{(l)} \in \mathbb{R}^{d^{(l-1)} \times d^{(l)}}$ are the l -th layer model parameters.

We consider the following special cases of \mathcal{M} which vary with respect to their graph filter:

- **RW:** $\forall l \in \mathbb{N}_{\leq L}, \mathbf{P}^{(l)} = \mathbf{P}_{\text{rw}} = \mathbf{D}^{-1}\mathbf{A}$ (i.e., the uniform random walk transition matrix), where \mathbf{D} is the diagonal degree matrix with entries $D_{ii} = \sum_{j \in \mathcal{V}} A_{ij}$.

- **SYM:** $\forall l \in \mathbb{N}_{\leq L}, \mathbf{P}^{(l)} = \mathbf{P}_{\text{sym}} = \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}}$.
- **ATT:** $\forall l \in \mathbb{N}_{\leq L}, \mathbf{P}^{(l)}$ is a graph attentional operator with default hyperparameters and a single head [VCC18].

Each encoder has three layers (64-dimensional hidden layers), with a ReLU nonlinearity in between layers. We do not use any regularization (e.g., Dropout, BatchNorm). The encoders are explicitly trained for node classification with the cross-entropy loss and the Adam optimizer [KB15] with full-batch gradient descent on the training set. We use a learning rate of 5e-3. We further use a random node split of 1000-500-rest for test-val-train. We train all encoders until they reach the training accuracy of MAJ_{WL} and select the model parameters with the highest validation accuracy. Although we do not do any hyperparameter tuning, the test accuracy values indicate that the encoders are well-trained.

We use PyTorch [PGM19] and PyTorch Geometric [FL19] to train all the encoders on a single NVIDIA GeForce GTX Titan Xp Graphic Card with 12196MiB of space on an internal cluster. On average (with respect to the datasets), the median time per training epoch was 0.05 seconds.

C.5 Additional Degree Bias Plots

Unlike the other datasets, we do not observe degree bias for chameleon and squirrel because these datasets are heterophilic. We intentionally include these datasets to draw contrast to the other, homophilic datasets and validate our theory. For example, in §4.4.2, we explain that high-degree nodes in heterophilic networks do not have more negative l -hop prediction homogeneity levels due to higher local heterophily levels; hence, we do not necessarily observe better performance for them compared to low-degree nodes.

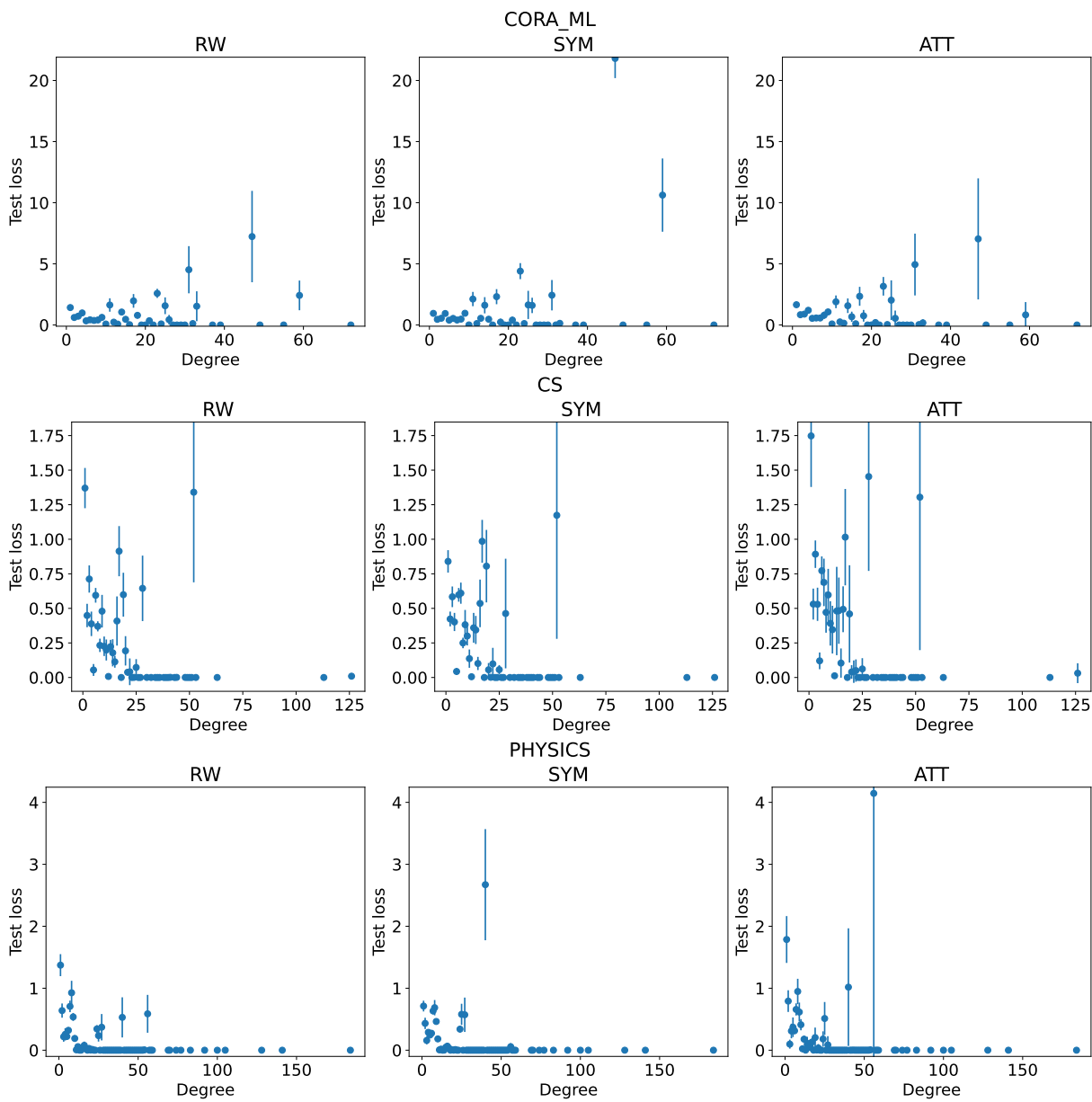


Figure C.1: Test loss vs. degree of nodes in citation and collaboration network datasets for RW, SYM, and ATT GNNs. High-degree nodes generally incur a lower test loss than low-degree nodes do. Error bars are reported over 10 random seeds; all error bars are 1-sigma and represent the standard deviation about the mean.

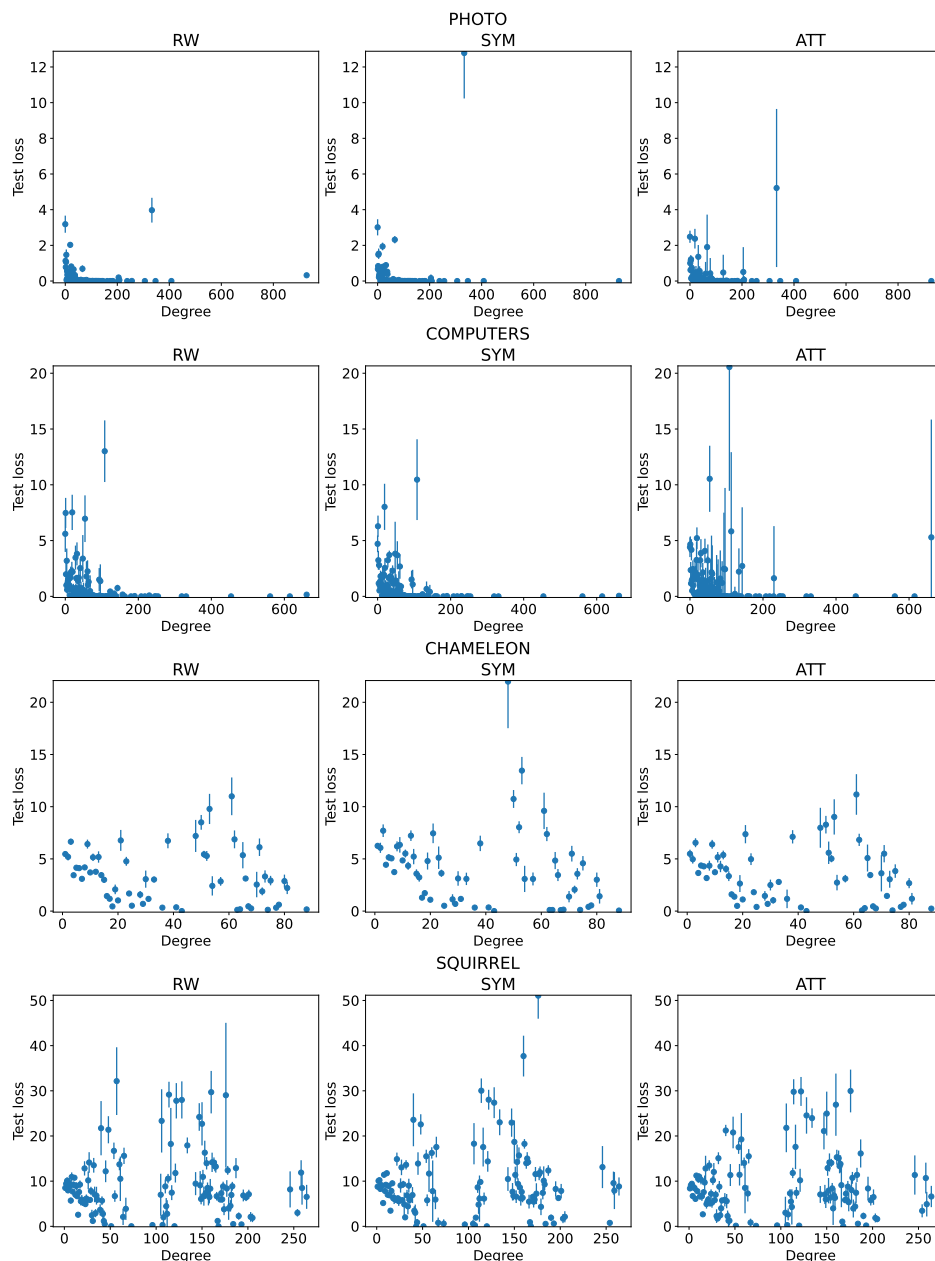


Figure C.2: Test loss vs. degree of nodes in online product and Wikipedia network datasets for RW, SYM, and ATT GNNs. High-degree nodes generally incur a lower test loss than low-degree nodes do. Error bars are reported over 10 random seeds; all error bars are 1-sigma and represent the standard deviation about the mean.

C.6 Additional Visual Summaries of Theoretical Results

In the plots below, we consider low-degree nodes to be the 100 nodes with the smallest degrees and high-degree nodes to be the 100 nodes with the largest degrees. Each point in the plots in the left column corresponds to a test node representation and its color represents the node’s class. The plots in the left column are based on a single random seed, while the plots in the middle and right columns are based on 10 random seeds. RW representations of low-degree nodes often have a larger variance than high-degree node representations, while SYM representations of low-degree nodes often have a smaller variance. Furthermore, SYM generally adjusts its training loss on low-degree nodes less rapidly.

The training loss curves in Figure 6 still support our theoretical analysis. Theorem 4.5.2 reveals that for \overline{SYM} , node degree *and* the (degree-discounted) expected feature similarity $\tilde{\chi}_i$ affects the rate of learning. On the other hand, Theorem C.8.1 indicates that for \overline{RW} , while we do not expect node degree to impact the rate of learning, the expected feature similarity χ_i is still influential. Hence, interpreting Theorems 4.5.2 and C.8.1 jointly, we expect and accordingly observe that the orange curve for SYM has a steeper rate of decrease *relative* to the orange curve for RW as the number of epochs increases.

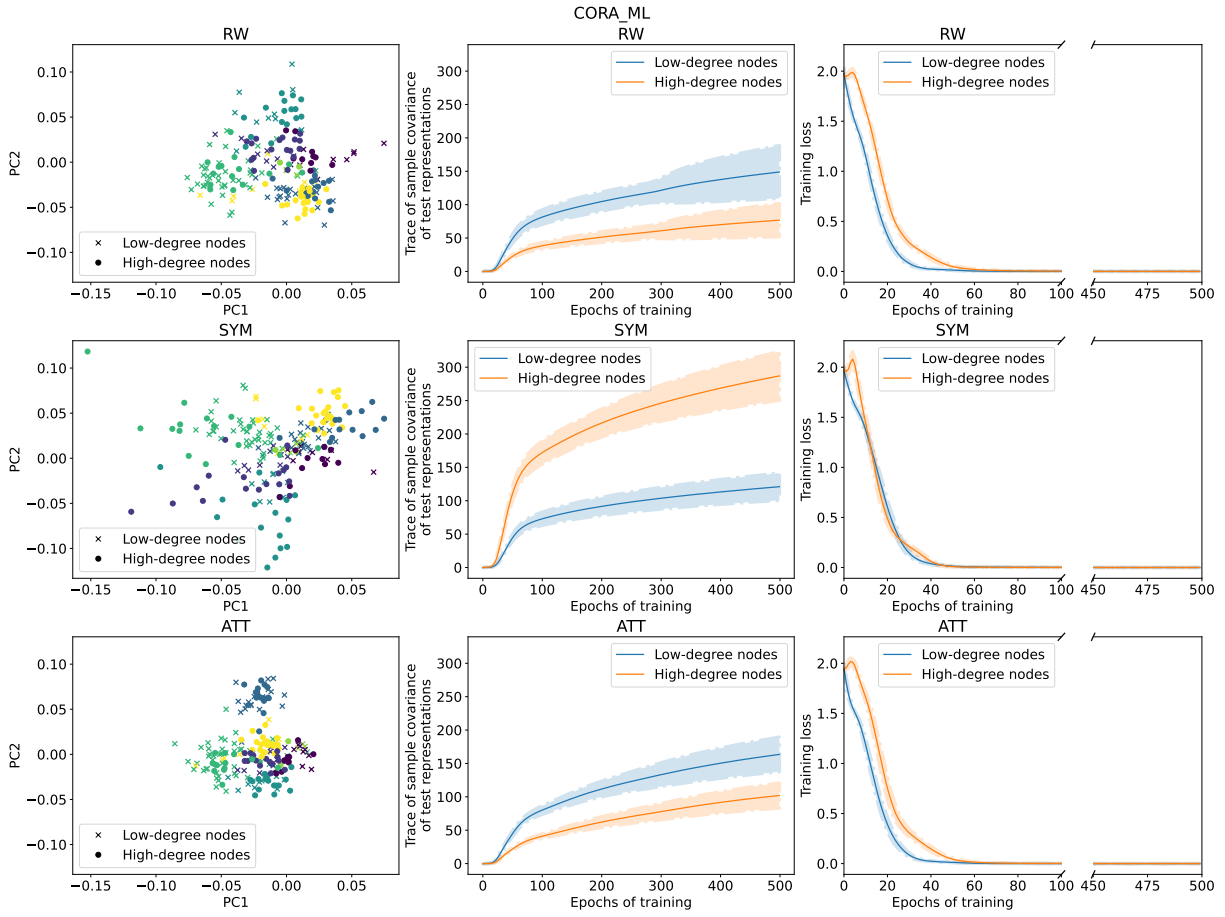


Figure C.3: Visual summary of the geometry of representations, variance of representations, and training dynamics of RW, SYM, and ATT GNNs on Cora_ML.

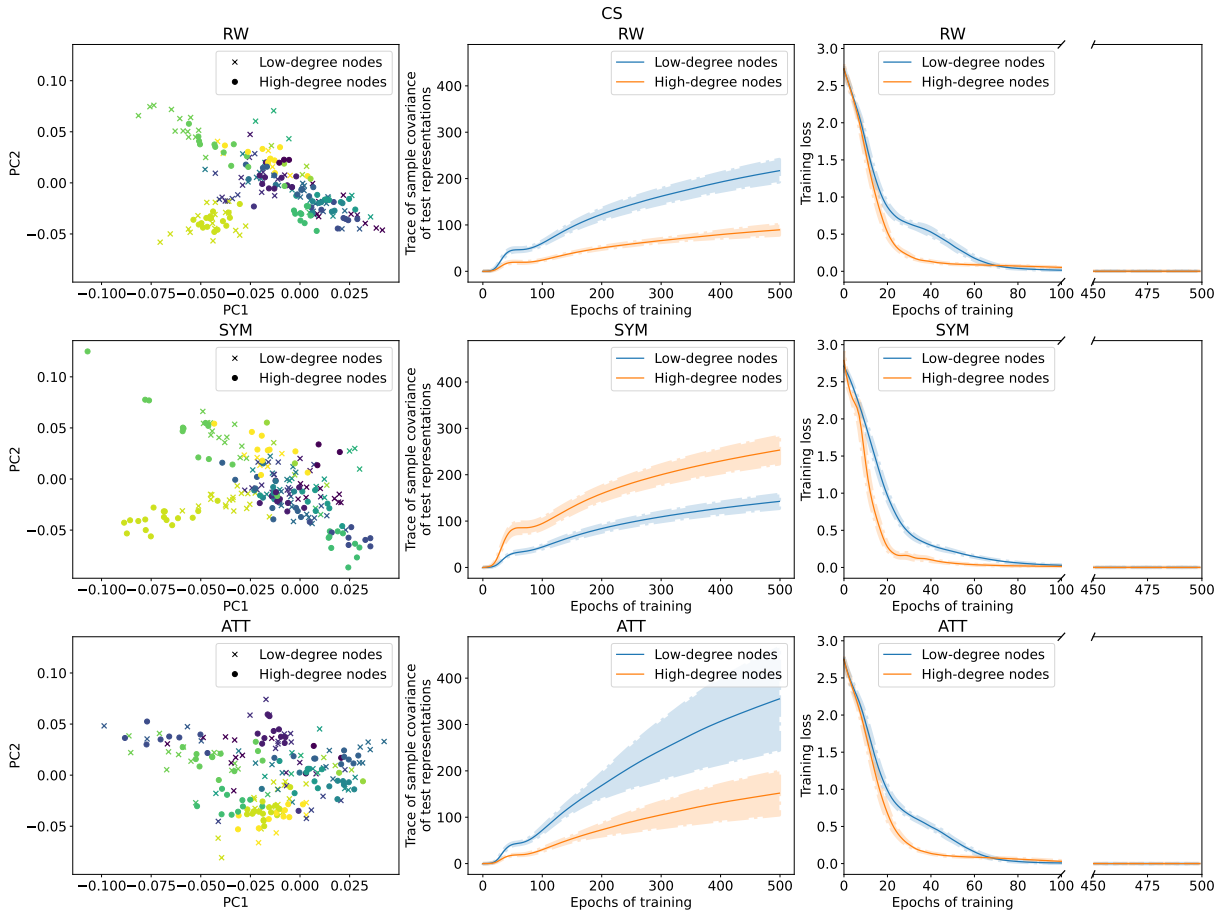


Figure C.4: Visual summary of the geometry of representations, variance of representations, and training dynamics of RW, SYM, and ATT GNNs on CS.

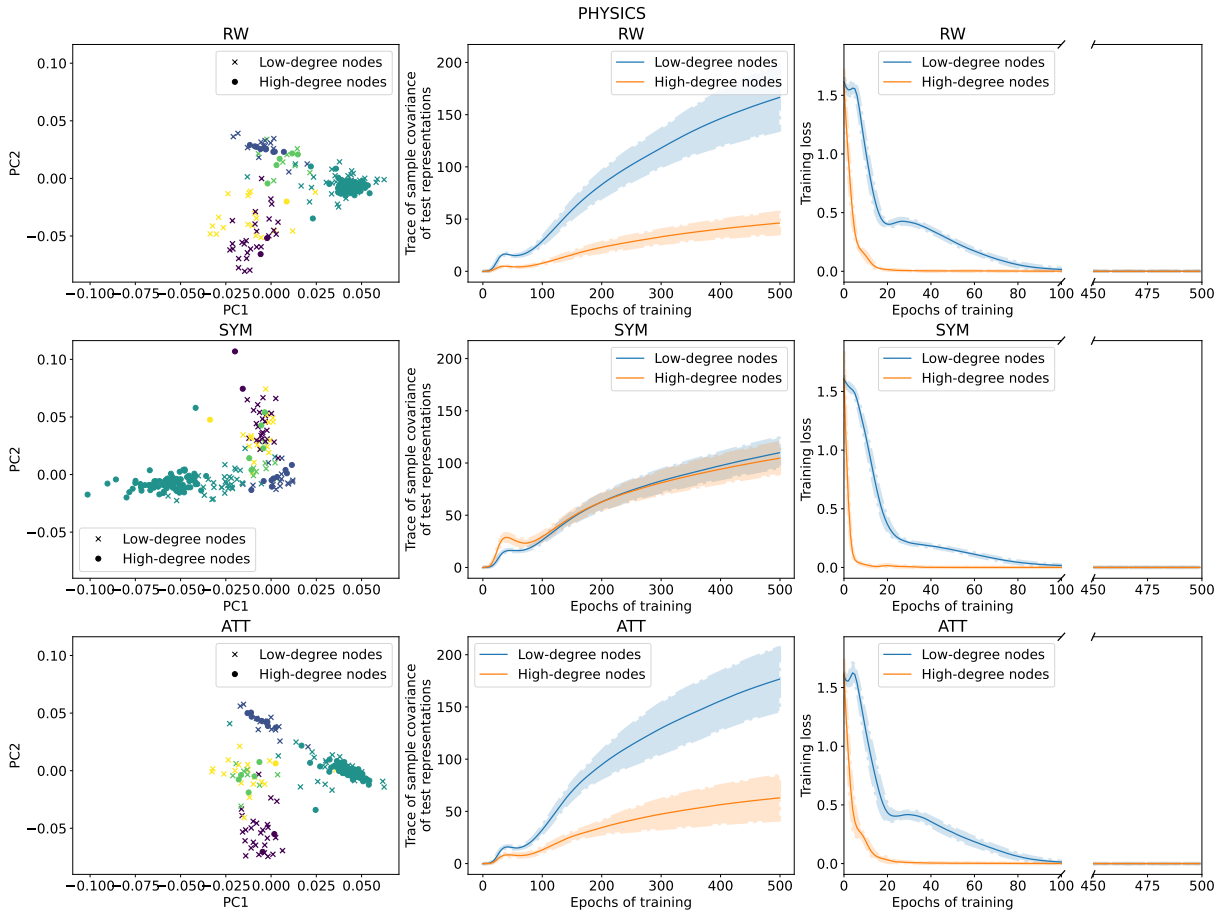


Figure C.5: Visual summary of the geometry of representations, variance of representations, and training dynamics of RW, SYM, and ATT GNNs on Physics.

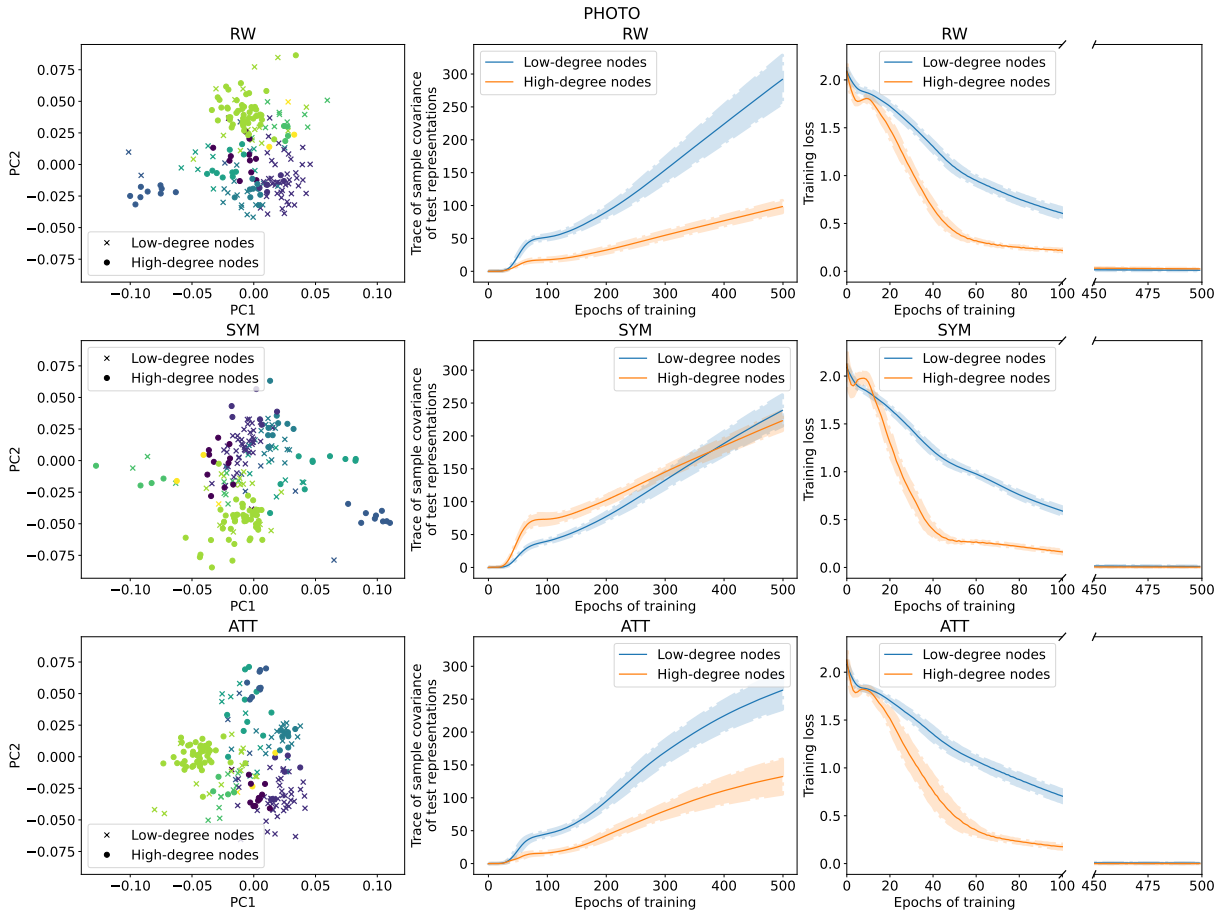


Figure C.6: Visual summary of the geometry of representations, variance of representations, and training dynamics of RW, SYM, and ATT GNNs on Amazon Photo.

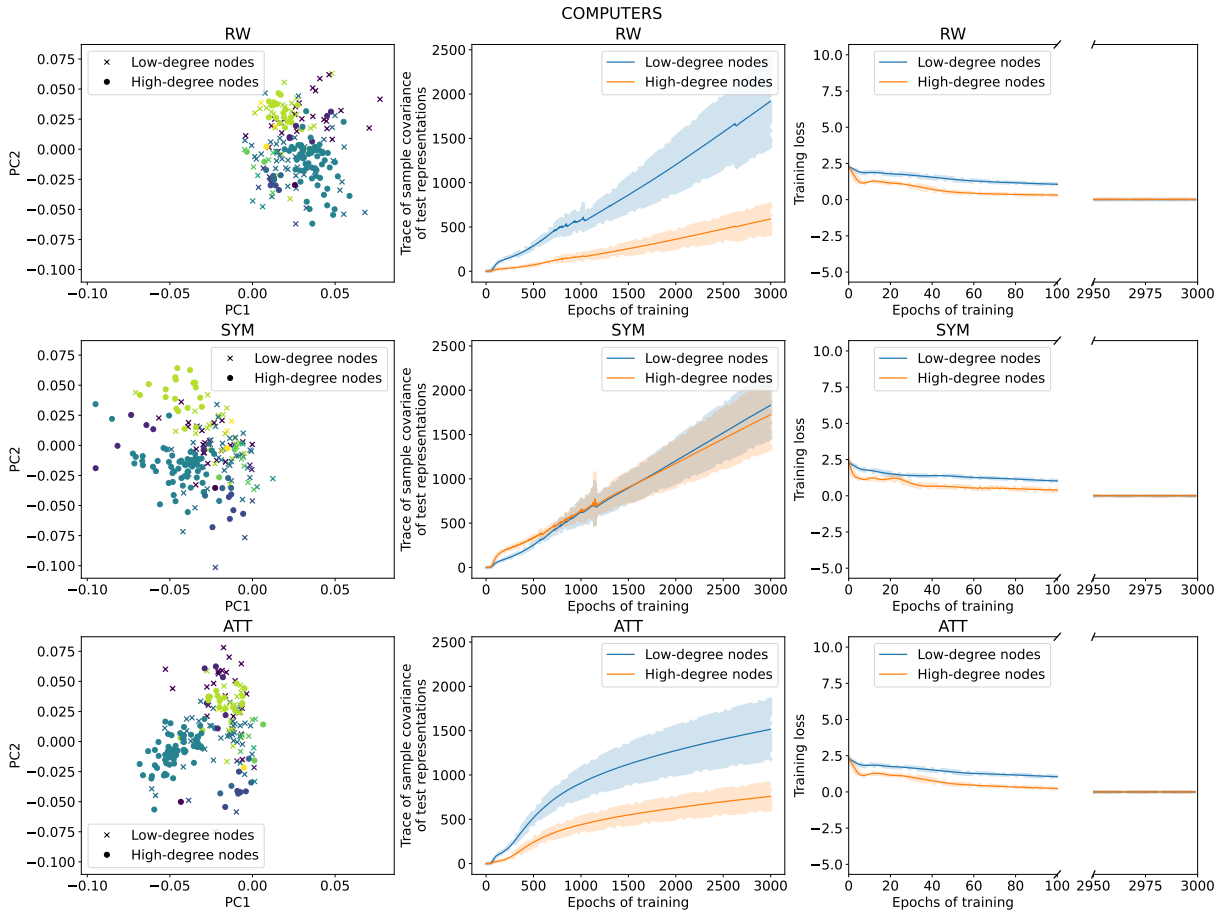


Figure C.7: Visual summary of the geometry of representations, variance of representations, and training dynamics of RW, SYM, and ATT GNNs on Amazon Computers.

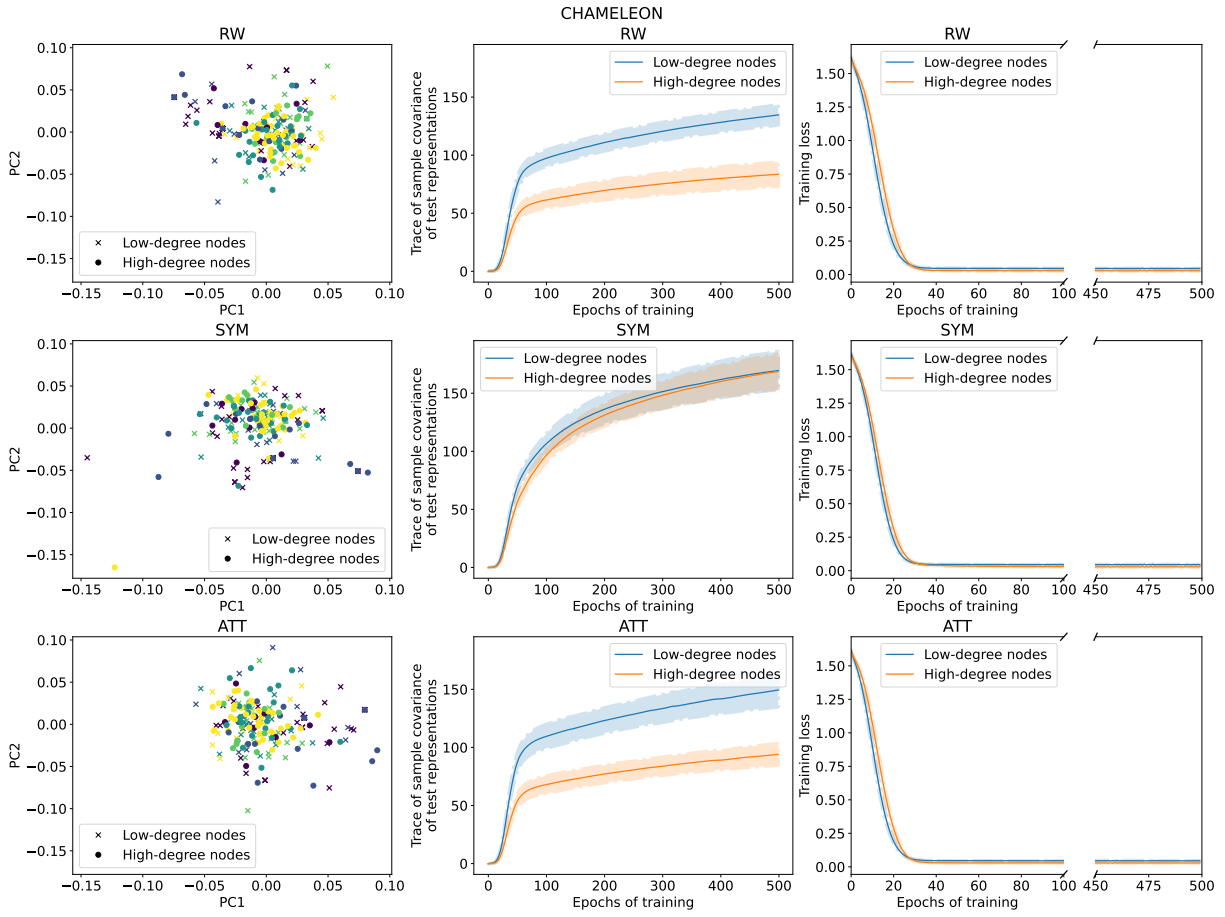


Figure C.8: Visual summary of the geometry of representations, variance of representations, and training dynamics of RW, SYM, and ATT GNNs on chameleon.

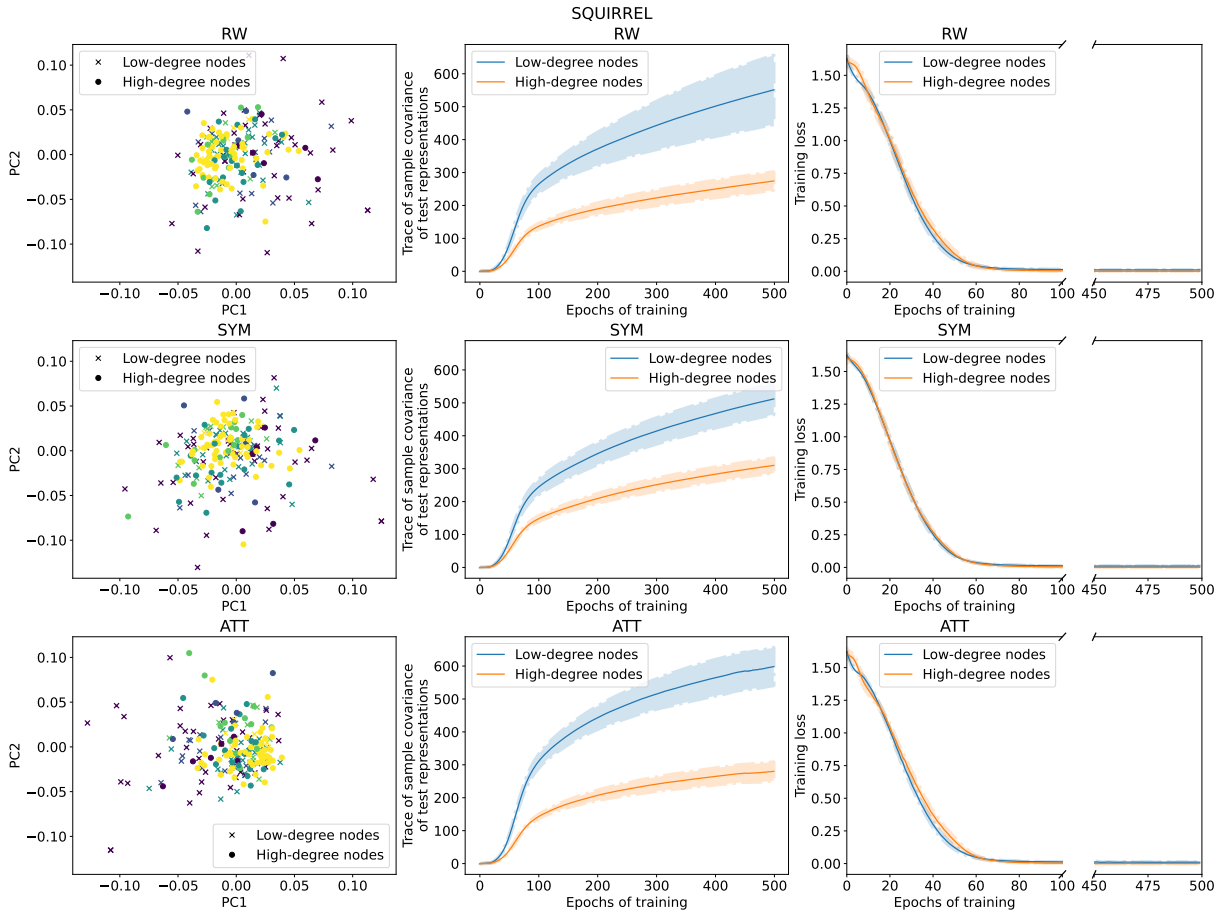


Figure C.9: Visual summary of the geometry of representations, variance of representations, and training dynamics of RW, SYM, and ATT GNNs on chameleon.

C.7 Additional Inverse Collision Probability Plots

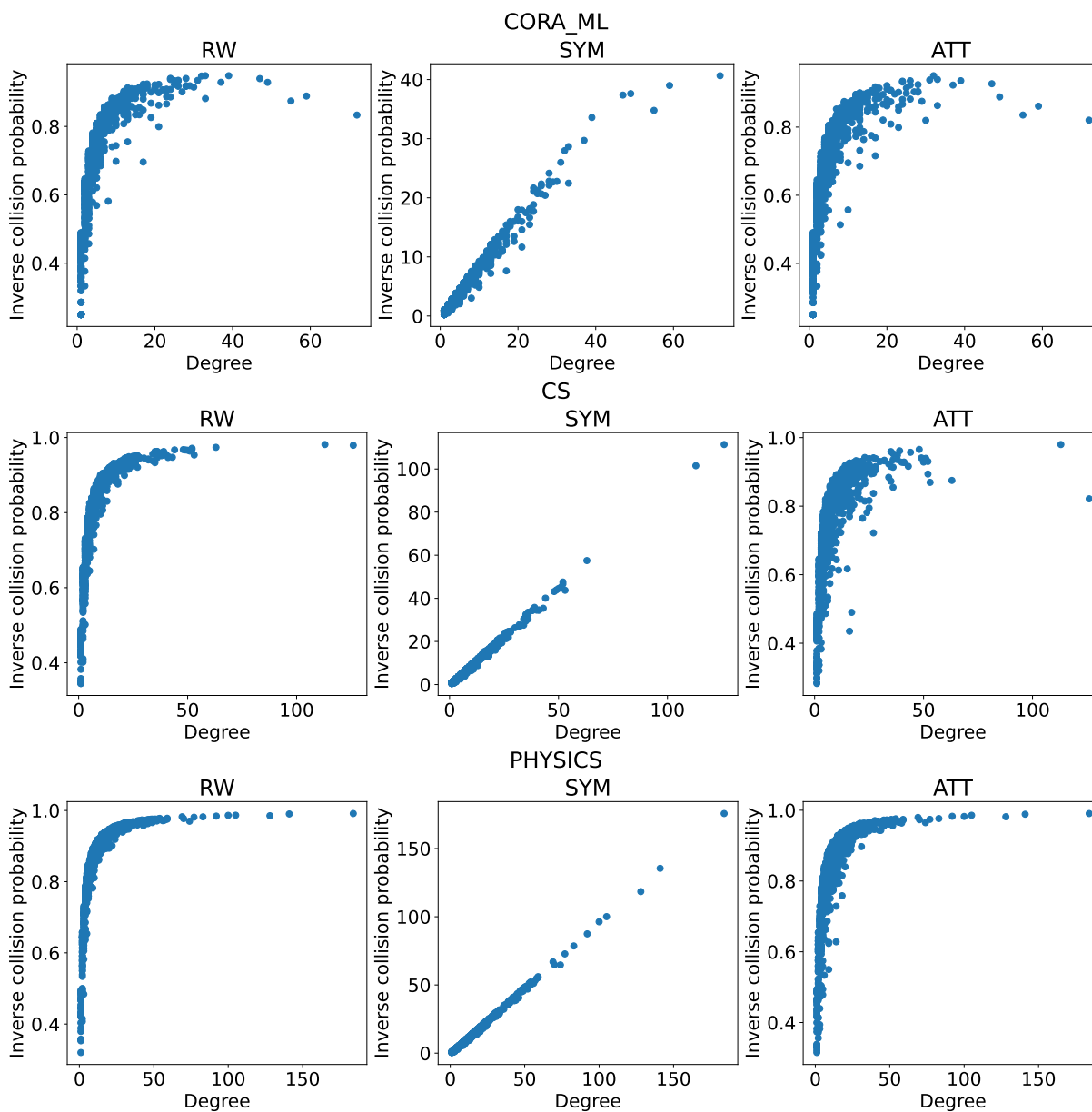


Figure C.10: Inverse collision probability vs. degree of nodes in citation and collaboration network datasets for RW, SYM, and ATT GNNs. Node degrees generally have a strong association with inverse collision probabilities.

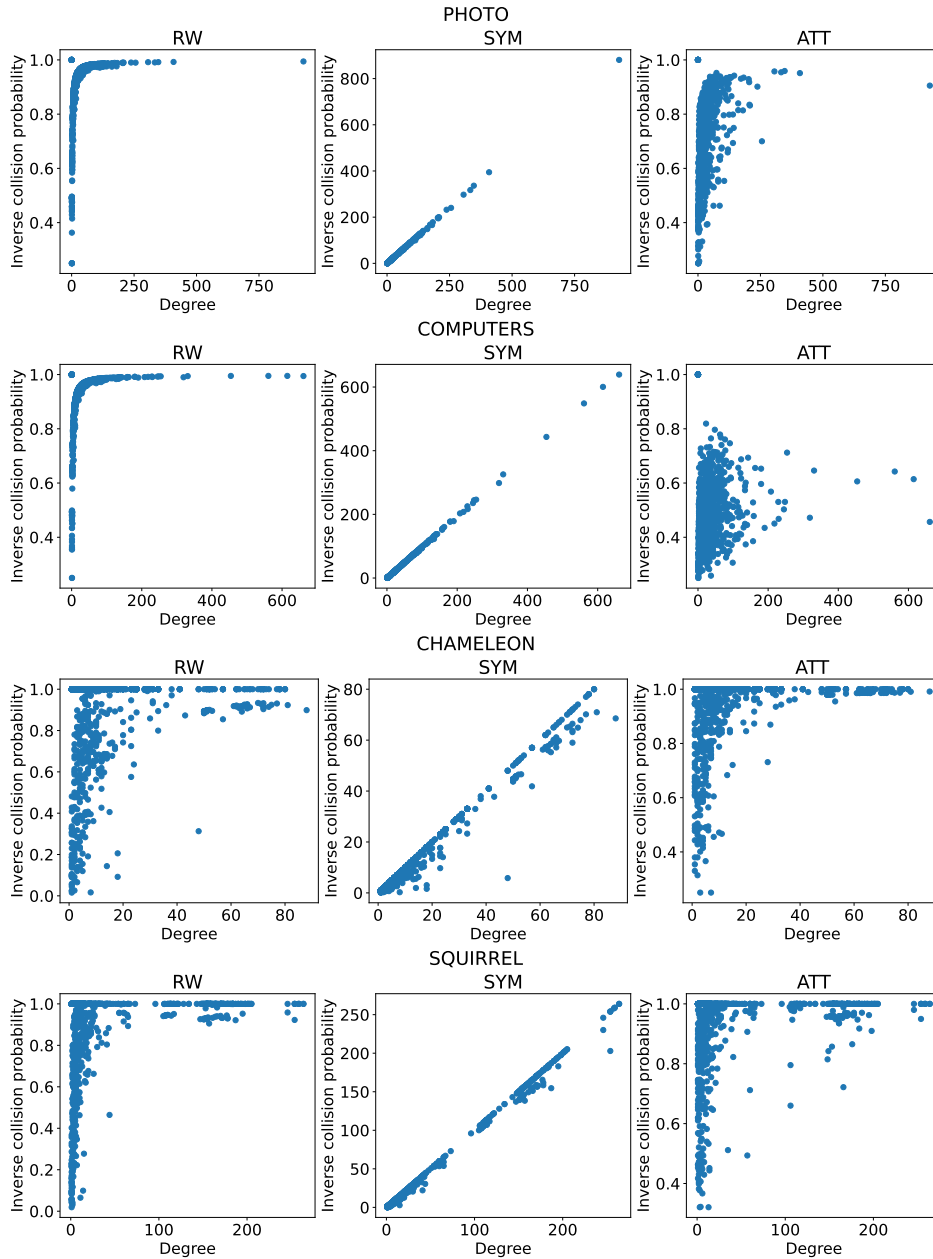


Figure C.11: Inverse collision probability vs. degree of nodes in citation and collaboration network datasets for RW, SYM, and ATT GNNs. Node degrees generally have a strong association with inverse collision probabilities.

C.8 Training-Time Degree Bias: Random Walk Graph Filter

We now demonstrate that during each step of training $\overline{\text{RW}}$ with gradient descent, compared to $\overline{\text{SYM}}$, the loss of low-degree nodes in S is not necessarily adjusted more slowly. We define $\forall l \in \mathbb{N}_{\leq L}, \chi_i^{(l)} \in \mathbb{R}^{|B[t]|}$, where for $m \in B[t]$, $\left(\chi_i^{(l)}[t]\right)_m = \mathbb{E}_{j \sim \mathcal{N}^{(l)}(i), k \sim \mathcal{N}^{(l)}(m)} [\mathbf{X}_j \mathbf{X}_k^T]$. In effect, $\left(\chi_i^{(l)}[t]\right)_m$ captures the expected similarity between the raw features of nodes j and k with respect to the l -hop random walk distributions of $i \in \mathcal{V}$ and $m \in B[t]$.

Theorem C.8.1. *The change in loss for i after an arbitrary training step t obeys:*

$$|\ell[t+1](\overline{\text{RW}}|i, c) - \ell[t](\overline{\text{RW}}|i, c)| \leq \sqrt{2}\eta \|\epsilon[t]\|_F \sum_{l=0}^L \left\| \chi_i^{(l)}[t] \right\|_2. \quad (\text{C.46})$$

For $\overline{\text{RW}}$, the change (either increase or decrease) in loss for i after an arbitrary training step does not necessarily have a smaller magnitude if i is low-degree. However, the L -hop neighborhoods of low-degree nodes still often have less overlap with the neighborhoods of training nodes, which can constrain the rate at which the loss for i changes. We confirm these findings empirically in Figure 4.3 and §C.6. For all the datasets, the blue curve for RW has a less steep rate of decrease than the blue curve for SYM as the number of epochs increases. However, for RW itself, the orange curve generally descends more quickly than the blue curve, with the exception of the heterophilic chameleon and squirrel datasets, for which the features of nodes in the neighborhoods of high-degree nodes and training nodes are dissimilar. Therefore, models do not learn more rapidly for high-degree nodes under heterophily. These findings support hypothesis **(H2)**. Our results for $\overline{\text{RW}}$ may also apply to ATT when low-degree nodes are generally attended to less.

C.9 Achieving Maximum Training Accuracy

CITeseer

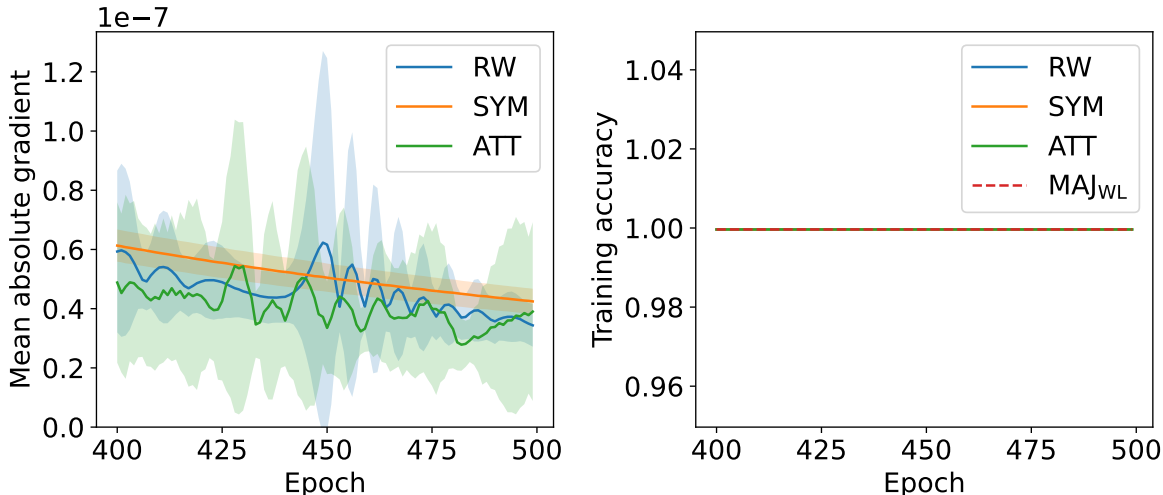


Figure C.12: Mean absolute parameter gradient vs. training epoch for RW, SYM, and ATT GNNs on CiteSeer (over 10 random seeds). The training accuracy of SYM, RW, and ATT ultimately reach the accuracy of MAJ_{WL}.

We now empirically show that SYM (despite learning at different rates for low vs. high-degree nodes), RW, and ATT can achieve their maximum possible training accuracy (i.e., the accuracy of a majority voting-classifier MAJ_{WL}). Furthermore, per our experiments, the accuracy of MAJ_{WL} is often close to 100%, indicating that the WL test does not significantly limit the accuracy of SYM, RW, and ATT in practice.

Per [XHL19], because the expressive power of SYM, RW, and ATT are limited by the WL test, an upper bound on their training accuracy is the accuracy of a majority voting-classifier MAJ_{WL} applied to WL node colorings (for details on how to compute colorings, see §II.A and §VI.B of [Zop22]). In particular, if the WL test produces the colors $\mathbf{c} \in \mathbb{K}^{|S|}$ for nodes in S , MAJ_{WL} predicts node i to have the label $\hat{\mathbf{Y}}_i = \text{MAJ}_{\text{WL}}(i, \mathbf{X}, \mathbf{A}) = \text{mode}\{\mathbf{Y}_j | j \in \mathcal{V}, \mathbf{c}_j = \mathbf{c}_i\}$. However, Figure C.12 and §C.10 reveal that as the number of training epochs increases, the training accuracy of SYM, RW, and ATT reach the accuracy of MAJ_{WL}. Because the accuracy of MAJ_{WL} is often close to 100%, our experiments suggest that insufficient

expressive power likely does not contribute to degree bias, drawing doubt to hypothesis **(H7)**.

Empirically inspecting the model gradients, compared to ATT, as the number of training epochs increases, the mean absolute gradients of SYM and RW are comparably small but often decrease more slowly or fluctuate. To understand this, we can analytically inspect the gradients of $\overline{\text{RW}}$:

$$\frac{\partial \ell[t]}{\partial \mathbf{W}^{(l)}[t]}(B[t]) = \mathbf{X}^T (\mathbf{P}_{\text{rw}}^l[t])^T \epsilon[t]. \quad (\text{C.47})$$

$(\mathbf{P}_{\text{rw}}^l[t])^T$ (for $l > 0$) often has numerous eigenvalues around 0, which can yield gradients $\frac{\partial \ell[t]}{\partial \mathbf{W}^{(l)}[t]}(B[t])$ with a small magnitude even when $\|\epsilon[t]\|$ is not small. The same analysis holds for $\overline{\text{SYM}}$. As such, SYM and RW may get trapped in suboptimal minima during training, yielding slower or unstable convergence; in contrast, because ATT has a dynamic filter, its training loss rarely exhibits slow or unstable convergence.

C.10 Additional Training Loss Plots

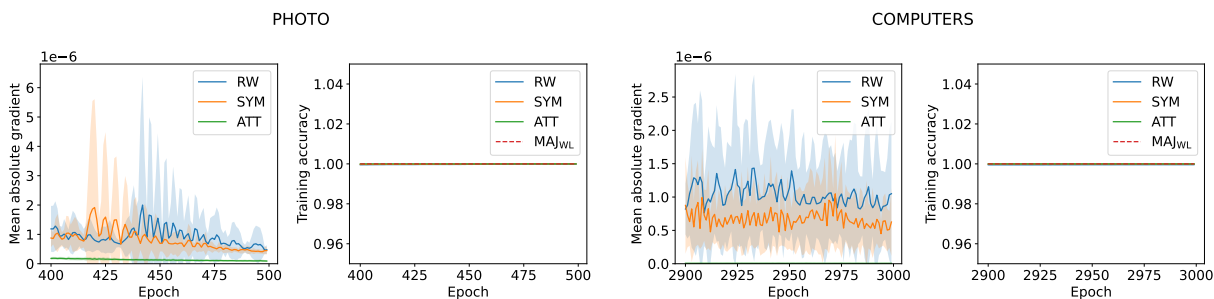


Figure C.13: Mean absolute parameter gradient vs. training epoch for RW, SYM, and ATT GNNs on Photo and Computers. The training accuracy of SYM, RW, and ATT ultimately reach the accuracy of MAJ_{WL} .

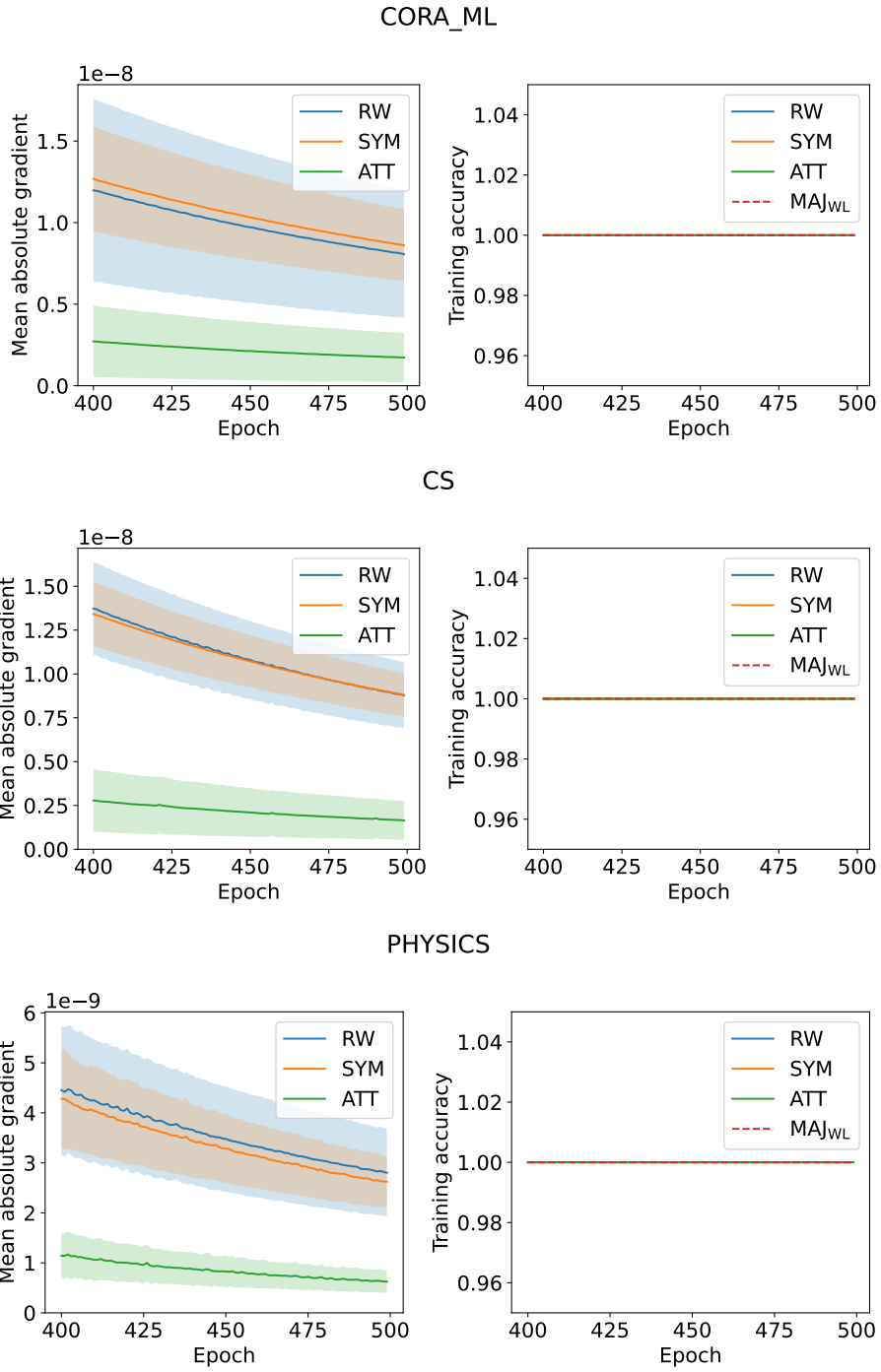
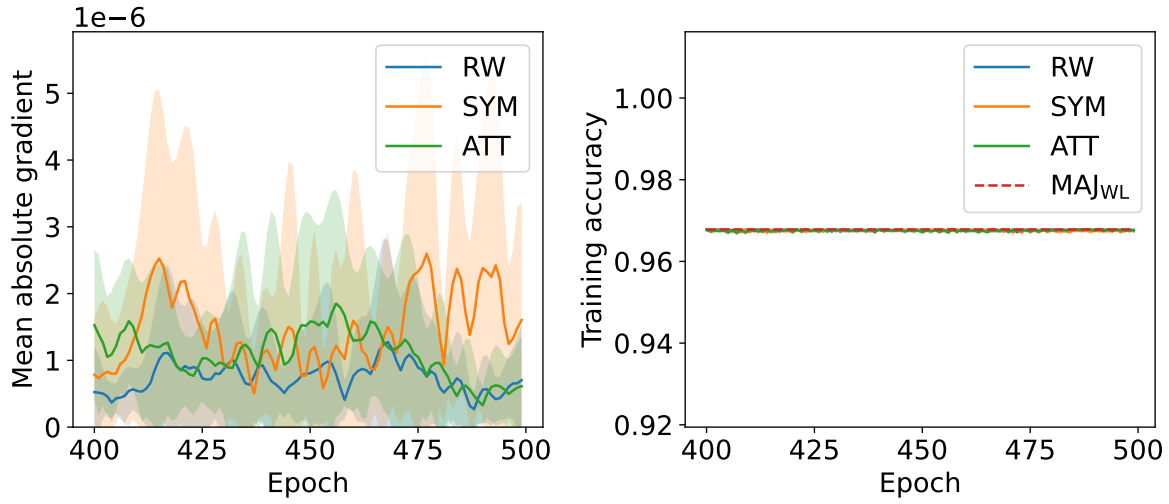


Figure C.14: Mean absolute parameter gradient vs. training epoch for RW, SYM, and ATT GNNs on Cora_ML, CS, and Physics. The training accuracy of SYM, RW, and ATT ultimately reach the accuracy of MAJ_{WL}.

CHAMELEON



SQUIRREL

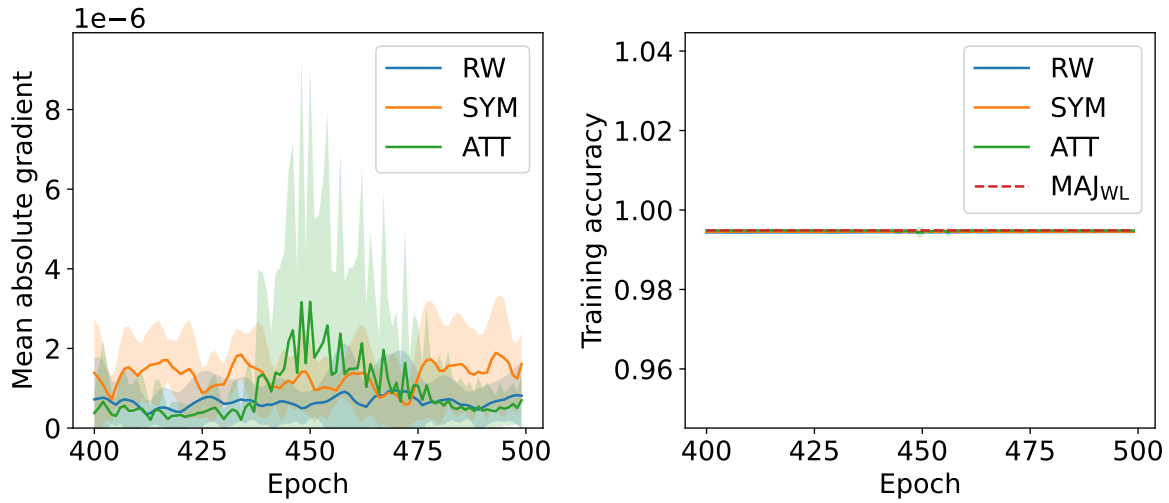


Figure C.15: Mean absolute parameter gradient vs. training epoch for RW, SYM, and ATT GNNs on chameleon and squirrel. The training accuracy of SYM, RW, and ATT ultimately reach the accuracy of MAJ_{WL}.

C.11 Limitations and Future Directions

Survey While we aimed to be thorough in our survey of prior papers on degree bias, it is inevitable that we missed some relevant work. In addition, we extract hypotheses for the origins of degree bias from the main bodies of papers; it is possible that the hypotheses do not fully or accurately reflect the *current* perspectives of the papers’ authors.

Theoretical analysis Our theoretical analysis is limited to linearized message-passing GNNs. While this is a common practice in the literature [WSZ19, CWH20, MLS22], it is a strong simplifying assumption. We empirically validate our theoretical findings on GNNs with non-linear activation functions, but this chapter does not address possible sources of degree bias related to these non-linearities, which would be interesting to investigate in future work. Towards this, a possible option is to assume that node features are drawn from a Gaussian distribution and derive precise high-dimensional asymptotics for degree bias in *non-linear* GNNs using the Gaussian equivalence theorem, as in [AP20a]. Our assumptions that GNNs generalize in expectation (Theorem 4.4.1) and the variance of node representations is finite (Theorems 4.4.2 and 4.4.3) are not overly strong assumptions in practice.

Furthermore, this chapter focuses on node classification. However, our findings readily lend insight into the origins of degree bias in link prediction. For example, if one uses node representations and an inner-product decoder to predict links between nodes, our results indicate that:

- In the random walk filter case, link prediction scores between low-degree nodes will suffer from higher variance because low-degree node representations have higher variance (Theorem 4.4.2). Hence, Theorem 4.4.1 suggests that predictions for links between low-degree nodes will have a higher misclassification error.
- In the symmetric filter case, our proof of Theorem 4.4.3 suggests that the link prediction scores between high-degree nodes will be over-calibrated (i.e., disproportionately large)

because high-degree node representations have a larger magnitude (i.e., approximately proportional to the square root of their degree). Hence, over-optimistic and possibly inaccurate links will be predicted between high-degree nodes.

The labels and evaluation for link prediction can confound intuition. Unlike node classification, the labels for link prediction (i.e., the existence or not of a link) make the task naturally imbalanced with respect to node degree; high-degree nodes have a much higher rate of positive links than low-degree nodes. This association between degree and positive labels can influence the misclassification error. Ultimately, more rigorous theoretical analysis and experimentation are needed to confirm the hypothesized implications of node degree for link prediction performance. Similarly, more research is required to understand the implications of our findings for degree bias in the context of graph classification.

Furthermore, our theoretical analysis does not encompass heterogeneous graphs. In our literature survey, we cover works that establish the issue of degree bias for knowledge graph predictions and embeddings (e.g., [BKE22, SJW23]). Our theoretical analysis is general and covers diverse message-passing GNNs, and can be extended to heterogeneous networks if messages aggregated from different edge types are subsequently linearly combined. In this setting, $R_{i,c'}$ can be computed as the sum of the prediction homogeneity quantities $(\sum_{l=0}^L \beta_{i,c'}^{(l)})^2$ for each edge type divided by the sum of the collision probability quantities $\sum_{j \in \mathcal{V}} [(\mathbf{P}_{rw}^l)_{ij}]^2$ for each edge type.

Empirical validation We sought to be transparent throughout this chapter regarding misalignments between empirical and theoretical findings. Our experiments focus on the transductive learning setting; it would be valuable to validate our theoretical findings in the inductive learning setting as well. Furthermore, while we aimed to cover diverse domains (e.g., citation networks, collaboration networks, online product networks, Wikipedia networks), as well as homophilic and heterophilic networks, it remains to identify the shortcomings of our theoretical findings for heterogeneous and directed networks.

Principled roadmap To do justice to studying the origins of degree bias in GNNs, which this chapter reveals has various and sometimes conflicting understandings, we limit the scope of this chapter to understanding the root causes of degree bias, not providing a concrete algorithm to alleviate it. Instead, we offer a principled roadmap based on our theoretical findings to address degree bias in the future. We further comment on the limitations of algorithmic solutions to degree bias in our Broader Impacts section below.

APPENDIX D

Appendix for Chapter 5

D.1 Limitations

Single pronoun sets in English. As the datasets we use for our meta-evaluation do not consider the usage of multiple pronoun sets for a single individual [MSR24], our analysis is limited to single pronoun sets for an individual. Additionally, we focus on third-person singular animate pronouns in English, and do not evaluate the wide range of neopronouns in English beyond xe [LCH22].

Other evaluation methods. We do not consider LLM-as-a-judge [LJH24] as an evaluation method in this chapter, as we restrict our focus exclusively to previously-proposed misgendering evaluation methods. Prior work has shown that LLMs exhibit performance disparities in gendering subjects correctly across pronouns [HDS23, GBZ24]. In addition, LLM-as-a-judge sidesteps the human-centered elements of evaluation for which we advocate. We thus leave a meta-evaluation of LLM-as-a-judge in the context of misgendering to future work. Such work can assess the efficacy of safety-oriented alignment protocols in penalizing misgendering [OPM24].

Agreement metrics. We report results with multiple agreement metrics because there are caveats to using Cohen’s κ . κ requires that the raters are independent. However, the evaluation results for TANGO and Prob-TANGO are not independent because the templates are constructed from model generations; similarly, for MISGENDERED and

Gen-MISGENDERED, and RUFF and Gen-RUFF, the generations are affected by the templates. Moreover, κ requires that the raters are fixed, but this is possibly violated by how the generation-based evaluation results are affected by sampling variation. In addition, the probability-based elements of the evaluation results are not solely due to evaluation subjectivity.

D.2 Formal Details About Probability-Based and Generation-Based Evaluations

D.2.1 Generation-Based Evaluation

Suppose we have a dataset $\mathcal{D}_{gen} := \{c^{(k)}, y^{(k)}\}_{k \in [N_{gen}]}$ comprising evaluation pairs of contexts (about a subject) and corresponding pronouns. We then define the generation-based gendering correctness $m_{gen}^{(k)} \in \{0, 1\}$ of the model on instance k as:

$$m_{gen}^{(k)} = 1 - \mathbb{1}(\mathcal{P}(\hat{y}_{gen}^{(k)}) \neq y^{(k)}), \quad \hat{y}_{gen}^{(k)} = g_i^{(k)} | g_i^{(k)} \in \Omega; \forall j < i, g_j^{(k)} \notin \Omega. \quad (\text{D.1})$$

If the generation does not contain a pronoun, $m_{gen}^{(k)} = 1$. $m_{gen}^{(k)} = 1$ indicates that the model is *correct* (i.e., does not misgender the subject). We visualize the rate at which TANGO generations lack pronouns in Figure D.3. [OGD23] shows that this heuristic can share high agreement with human annotations for misgendering.

D.2.2 Probability-Based Evaluation

Suppose we have a dataset $\mathcal{D}_{prob} := \{t^{(k)}, y^{(k)}\}_{k \in [N_{prob}]}$ comprising evaluation pairs of templates (about a subject) and corresponding pronouns. Let $\Omega^c := \{p \in \Omega | \mathcal{C}(p) = c\}$. Then, we define the probability-based gendering correctness $m_{prob}^{(k)}$ of the model on instance k as:

$$m_{prob}^{(k)} = 1 - \mathbb{1}(\mathcal{P}(\hat{y}_{prob}^{(k)}) \neq y^{(k)}), \quad \hat{y}_{prob}^{(k)} = \arg \min_{p \in \Omega^c} \text{perp}(t_{1:m-1}^{(k)} \| R(p) \| t_{m+1:T}^{(k)}), \quad (\text{D.2})$$

where \parallel concatenates sequences, R appropriately transforms p (e.g., capitalization if p is at the beginning of a sentence), and perp maps a sequence to its perplexity (as determined by the sequence generation probabilities encoded by the model). Eq. D.2 effectively searches for the most likely sequence to be generated over a minimal contrast set, which reduces the effect of confounding factors (e.g., gendered vocabulary) on the evaluation. Moreover, by definition, perp normalizes the raw probability of generating each sequence by its length, which accounts for varying sequence lengths due to the overfragmentation of neopronouns during tokenization [OMG24].

D.3 Experimental Details

We access all models through HuggingFace [WDS20]. We run the models with at most 8B parameters on a single Nvidia A100 GPU. We load the larger models with low CPU memory usage and half-precision FP, and the distribute them across 3-4 A100 GPUs using HuggingFace’s automatic device mapping. For Mixtral 8x22B, we additionally use 4-bit quantization.

Our experiments have a non-trivial runtime: For each instance and ground-truth pronoun in MISGENDERED and RUFF (and their generation-based transformations), we perform constrained decoding for the [MASK] and generate ten 50-token sequences (across the pre- and post-[MASK] settings). For each instance and ground-truth pronoun in TANGO (and its probability-based transformation), we generate five 50-token sequences and perform constrained decoding five times. The experiments on Mixtral 8x22B-v0.1-4bit with RUFF took about 72 hours with our setup.

D.3.1 MISGENDERED

We restrict our focus to the subset of the dataset with instances that starts with “{name}’s pronouns are {nom}/{acc}/{pos_ind}.” For each instance, we produce 15 templates by

filling “{name}” with a different random subset of 15 personal names from all the names (across 100 masculine, 100 feminine, and 300 neutral) used by [HDS23]. This process enables us to approximately marginalize out the effect of gendered names on our evaluation results. This yields 750 templates per ground-truth base pronoun, which we transform into generation contexts to produce Gen-MISGENDERED.

For each context in Gen-MISGENDERED, we generate completions with top-50 filtering, nucleus sampling ($p = 0.95$), and the default values for each model for the other decoding hyperparameters; we do this to match the hyperparameters used in [OGD23]. We further perform generation with a single beam; we found empirically that generation with beam search often yields degeneration (e.g., highly repetitive sequences). For each context, we generate exactly 50 tokens, based on experiments with Llama-3.2-1B showing that about 95% of model generations include a pronoun within the first 50 tokens. We generate $R = 5$ pre-[MASK] and R post-[MASK] completions per context. We use SpaCy’s `en_core_web_sm` model for all tokenization and parsing apart from any LLM-specific tokenization [HMV20].

D.3.2 RUFF

We only consider the subset of the dataset without distractor sentences, as it is not possible to automatically measure misgendering by simply examining the first generated pronoun without considering to whom it refers. RUFF does not use personal names. This produces 1800 templates per ground-truth base pronoun, which we transform into generation contexts to produce Gen-RUFF. We follow the same generation settings as for Gen-MISGENDERED.

D.3.3 TANGO

We focus on the misgendering subset of TANGO [OGD23]. Unlike MISGENDERED, any names in TANGO contexts are predefined. In total, TANGO contains 480 contexts per ground-truth base pronoun, the generations for which we transform into templates to produce

Prob-TANGO. We follow the same generation settings as for Gen-MISGENDERED.

D.3.4 Practical Challenges

There are practical challenges to creating templates from generations. For example, not all generated completions $g^{(k)}$ contain a pronoun (e.g., the subject’s name is repeatedly used instead), in which case a template $t^{(k)}$ cannot be constructed; we discard such completions. Even when there is a pronoun, cases of pronouns are often not unique (e.g., “That is his book.” and “That book is his.”). Hence, determining the case of a pronoun occurrence (e.g., dependent possessive vs. independent possessive) for the purpose of performing probability-based evaluation can be challenging. Moreover, templates in native probability-based evaluation datasets are carefully constructed to be syntactically robust to replacements of the [MASK] token with pronouns (e.g., English-language templates are intentionally written in past tense to avoid issues of incorrect conjugation). However, templates constructed from generations need to be adjusted to accommodate the proper conjugation of verbs of which the pronoun that replaces the [MASK] is a dependent. To address these challenges, rather than use a [MASK], we rewrite $g^{(k)}$ using each pronoun:

1. If $g^{(k)}$ contains **xe**, we replace it with the corresponding case of **she**. This is unambiguously possible because either: (1) every case of **xe** is unique, or (2) every case of **xe** uniquely maps onto the corresponding case of **she**.
2. We then apply the gender-neutral rewriting algorithm described in [SWS21] to correctly transform $g^{(k)}$ to use **they** with proper case and conjugation of verbs. This algorithm uses constrained decoding with GPT-2 to disambiguate between different cases of **he** and **she** [RWC19]. We do not neutralize occupational or gender-specific terms. Step 1 is needed because GPT-2 may not robustly process **xe** pronouns.
3. To transform $g^{(k)}$ to use **he**, **she**, and **xe**, we apply the rewriting algorithm in reverse. The reverse direction does not require GPT-2, as each case of **they** is unique. A notable

limitation of the deneutralization algorithm is that it does not properly handle conjunct verbs (e.g., “hugs” in “He cries and hugs Sarah.”), as SpaCy does not correctly tag conjunct verbs as verbs [HMV20].

We opt to use a majorly rule-based rewriting approach rather than purely LLM-based rewriting methods to avoid performance biases for **xe** and singular **they** that might be introduced by LLMs.

D.3.5 Agreement Metrics

MISGENDERED and RUFF. Let $m_{prob}^{(k)}$ be the occurrence of correct gendering for instance k in MISGENDERED or RUFF. Furthermore, let $[m_{gen}^{(k)}]_i$ be the occurrence of correct gendering in the i -th generation for instance k in Gen-MISGENDERED or Gen-RUFF. Each of the following metrics is computed separately for the pre and post-[MASK] settings.

- **Instance-level:** It has been observed that generated text-based metrics are highly sensitive to decoding hyperparameters [AKP22], which can yield different results for the same dataset [LAN24]. Therefore, we measure the standard deviation of correct gendering across different generations i for the same instance.

$$\sigma_{gen}^{(k)} = \text{stdev}_i ([m_{gen}^{(k)}]_i). \quad (\text{D.3})$$

$\sigma_{gen}^{(k)}$ captures the effect of sampling variance on evaluation results.

- **Dataset-level:** We measure the Matthew’s correlation coefficient $MCC \in [-1, 1]$ of the probability- and generation-based gendering results. MCC is equivalent to Pearson’s correlation coefficient for binary variables, and can be better suited for imbalanced data (such as misgendering evaluation results) than raw observed agreement [CJ20]. In addition, we consider the raw observed agreement $p_o \in [0, 1]$ between the results of the two evaluation methods. We also consider Cohen’s $\kappa \in [-1, 1]$, which corrects the observed agreement for the expected probability that the results agree. Given $m^{(k)} \in \{0, 1\}$, we measure the

dataset-level variation v^f as:

$$v^f = f \left(\{m_{prob}^{(k)}\}_{k \in [N_{prob}]}, \{[m_{gen}^{(k)}]_1\}_{k \in [N_{prob}]} \right), \quad (\text{D.4})$$

where f is MCC , κ , or p_o . We only consider $[m_{gen}^{(k)}]_1$ to isolate the effect of dataset variance (rather than sampling variance, which is captured by Eq. 5.1) on the agreement of evaluation results. In contrast, $m_{prob}^{(k)}$ is not affected by sampling variance. Unlike [HL23b], we look at binary evaluation outcomes, and not direct probabilities, to perform a more “extrinsic” analysis (e.g., an end-user of a chatbot either sees or does not see an incorrect pronoun, not the probability of the pronoun being generated).

- **Model-level:** To compare evaluation disagreement across different models and pronouns, we model the probability of disagreement $d^{(k)}$ across instances as samples from a beta distribution. For the two dataset types, this is formalized as follows:

$$d^{(k)} = m_{prob}^{(k)}(1 - \bar{m}_{gen}^{(k)}) + (1 - m_{prob}^{(k)})\bar{m}_{gen}^{(k)}, \quad \text{where } \bar{m}_{gen}^{(k)} = \text{mean}_i ([m_{gen}^{(k)}]_i), \quad (\text{D.5})$$

$$\alpha, \beta = \text{MLE}_{beta} \left(\{d^{(k)}\}_{k \in [N_{prob}]} \right), \quad (\text{D.6})$$

where MLE_{beta} outputs the maximum-likelihood estimates of α, β for a beta distribution given the sample of probabilities $d^{(k)}$. We infer α, β using the method of moments.

TANGO. Let $[m_{gen}^{(k)}]_i$ be the occurrence of correct gendering in the i -th generation for instance k in TANGO. In addition, let $[m_{prob}^{(k)}]_i$ be the occurrence of correct gendering in the i -th template for instance k in Prob-TANGO.

- **Instance-level:** We measure the standard deviation of correct gendering across different generations and templates i for the same instance.

$$\sigma_{gen}^{(k)} = \text{stdev}_i ([m_{gen}^{(k)}]_i), \quad \sigma_{prob}^{(k)} = \text{stdev}_i ([m_{prob}^{(k)}]_i). \quad (\text{D.7})$$

- **Dataset-level:** We measure the v^f , for $f \in \{MCC, \kappa, agr\}$, of the probability- and generation-based gendering results:

$$v^f = f \left(\{[m_{prob}^{(k)}]_1\}_{k \in [N_{gen}]}, \{[m_{gen}^{(k)}]_1\}_{k \in [N_{gen}]} \right). \quad (\text{D.8})$$

- **Model-level:** We measure the model-level disagreement as:

$$d^{(k)} = \bar{m}_{prob}^{(k)}(1 - \bar{m}_{gen}^{(k)}) + (1 - \bar{m}_{prob}^{(k)})\bar{m}_{gen}^{(k)}, \quad \alpha, \beta = \text{MLE}_{beta}(\{d^{(k)}\}_{k \in [N_{gen}]}), \quad (\text{D.9})$$

$$\text{where } \bar{m}_{prob}^{(k)} = \text{mean}_i([m_{prob}]_i^{(k)}), \quad \bar{m}_{gen}^{(k)} = \text{mean}_i([m_{gen}]_i^{(k)}). \quad (\text{D.10})$$

D.4 Theoretical Analysis of Divergence Between Evaluation Formats

The analyses below assume that each token is a complete word, which is not entirely faithful to how LLMs operate in practice; however, our analyses can be easily extended to the setting where tokens are subwords. Furthermore, our analyses assume that LLM generations are produced via sampling with a single beam (without top- k filtering or nucleus sampling). We also sidestep issues of capitalization and other formatting.

D.4.1 Converting from Probability-Based to Generation-Based Evaluation

Suppose we have an LLM \mathcal{M} that induces a conditional probability distribution over tokens. We have a template $\{t_i\}_{i \in [T]}$, with the [MASK] token t_m associated with case c . For simplicity of notation, let $\Omega^c := \{p \in \Omega | \mathcal{C}(p) = c\}$. We define:

$$p^* = \arg \max_{p \in \Omega^c} Pr(p|t_{1:m-1}) \cdot Pr(t_{m+1:T}|t_{1:m-1} || p), \quad (\text{D.11})$$

where p^* is the most likely pronoun for [MASK]. Now, we consider the pre-[MASK] generation-based setting. Suppose the first pronoun in g is token g_1 with case c . Then, the probability δ of the generation-based evaluation disagreeing with the probability-based evaluation is given by:

$$\delta = 1 - \frac{Pr(p^*|t_{1:m-1})}{\sum_{p \in \Omega^c} Pr(p|t_{1:m-1})}. \quad (\text{D.12})$$

δ is minimized when $Pr(p^*|t_{1:m-1})$ is maximized. That is, the minimal probability of disagreement δ^* between the two evaluation methods is:

$$\delta^* = 1 - \underbrace{\frac{\max_{p \in \Omega^c} Pr(p|t_{1:m-1})}{\sum_{p \in \Omega^c} Pr(p|t_{1:m-1})}}_{\text{dominance of mode of next-token distribution}}. \quad (\text{D.13})$$

Now, we consider the post-[MASK] generation-based setting. Once again, suppose the first pronoun in g is token g_1 with case c . Similarly to the pre-[MASK] case, the probability δ of the generation-based evaluation disagreeing with the probability-based evaluation is given by:

$$\delta^* = 1 - \frac{\max_{p \in \Omega^c} Pr(p|t_{1:T})}{\sum_{p \in \Omega^c} Pr(p|t_{1:T})}. \quad (\text{D.14})$$

D.4.2 Converting from Generation-Based to Probability-Based Evaluation

We have a context $\{c_i\}_{i \in [C]}$ and corresponding generation $\{g_i\}_{i \in [G]}$. The first pronoun is token g_m with case c , which becomes the [MASK] token in the template. Then, we define:

$$p^* = \arg \max_{p \in \Omega^c} Pr(p|c_{1:C} \parallel g_{1:m-1}) \cdot Pr(g_{m+1:n}|c_{1:C} \parallel g_{1:m-1} \parallel p), \quad (\text{D.15})$$

where p^* is the most likely pronoun for [MASK]. The probability-based evaluation will disagree with the generation-based evaluation with probability δ where:

$$\delta = 1 - \frac{Pr(p^*|c_{1:C} \parallel g_{1:m-1})}{\sum_{p \in \Omega^c} Pr(p|c_{1:C} \parallel g_{1:m-1})}. \quad (\text{D.16})$$

δ is minimized when $Pr(p^*|g_{1:m-1})$ is maximized. That is, the minimal probability of disagreement δ^* between the two evaluation methods is:

$$\delta^* = 1 - \frac{\max_{p \in \Omega^c} Pr(p|c_{1:C} \parallel g_{1:m-1})}{\sum_{p \in \Omega^c} Pr(p|c_{1:C} \parallel g_{1:m-1})}. \quad (\text{D.17})$$

These analyses suggest that disagreement may arise due to: (1) autoregressive sampling only depending on previously-generated tokens and not always sampling the most likely token, and (2) template segments after the [MASK] not aligning with what would likely be generated in practice.

D.5 Additional Experimental Results

D.5.1 MISGENDERED

As the main paper focused on pre-[MASK] generations, Figure D.1 shows the corresponding figures for variation and agreement in the post-[MASK] generation setting. Similarly, Table D.1 shows *MCC* agreement, and Table D.2 shows κ agreement results in this setting.

We also analyze the agreement between *models* and pronouns, as shown in Figure D.2, which visualizes the landscape of evaluation disagreement probabilities across all the models and pronouns. Most points fall below the line $\alpha = \beta$ (i.e., $\alpha < \beta$), which suggests that there is a higher rate of agreement than disagreement. Moreover, with the exception of **xe**, we observe clusters associated with the model family (but not model size); hence, pre-training data and family-specific architectural components may have a larger influence on the disagreement of evaluation results than model size. In the pre-[MASK] setting, the Llama cluster and part of the Mixtral cluster have $\alpha, \beta < 1$, which indicates that the probabilities of disagreement are concentrated around 0 and 1. The other part of the Mixtral cluster and the OLMo cluster have $\alpha < 1, \beta > 1$, which indicates that the probabilities of disagreement are more concentrated around 0. In contrast, in the post-[MASK] setting, the OLMo cluster has $\alpha, \beta < 1$. The points corresponding to **xe** generally appear separate from their model family clusters, indicating distinct evaluation disagreement behavior for the neopronoun compared to other pronouns.

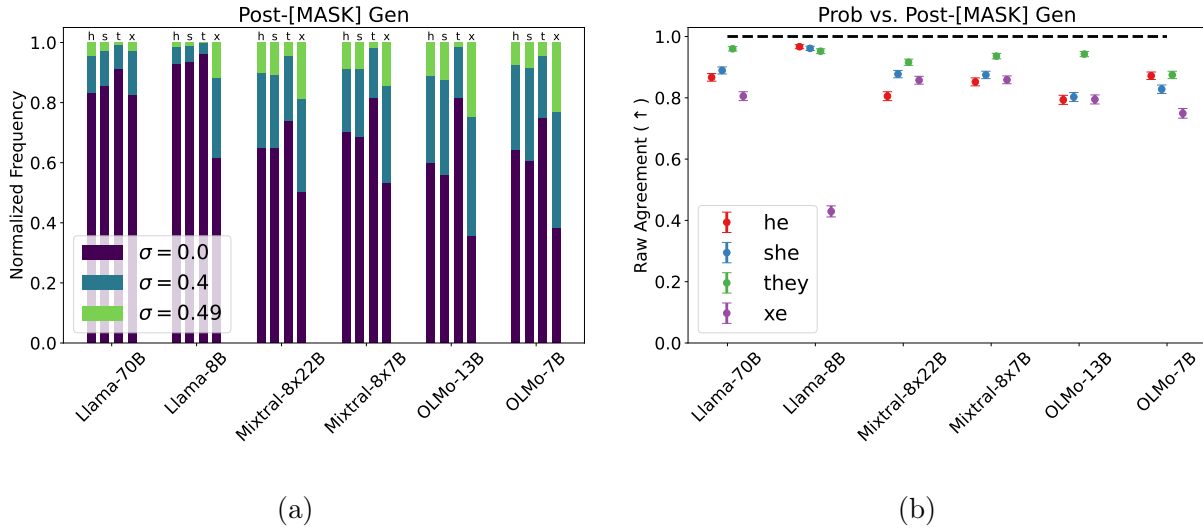


Figure D.1: **(a)** Generation variation σ (Eq. 5.1) for each model and pronoun in the post-[MASK] generation setting for MISGENDERED. Because we sample five generations per context, $\sigma \in \{0, 0.4, 0.49\}$. The bar labels h, s, t, x correspond to he, she, they, xe. **(b)** Raw observed agreement v^{p_o} (Eq. 5.2) for each model and pronoun between the probability-based and post-[MASK] generation-based evaluation results for MISGENDERED. The error bars represent the standard error of v^{p_o} (computed over dataset instances). The horizontal dashed line represents the upper bound of v^{p_o} .

	he	she	they	xe
Llama-70B	0.009 [-0.062, 0.081]	-0.000 [-0.072, 0.071]	-0.020 [-0.092, 0.051]	0.024 [-0.047, 0.096]
Llama-8B	-0.017 [-0.088, 0.055]	-0.018 [-0.089, 0.054]	-0.017 [-0.088, 0.055]	0.083 [0.011, 0.153]
Mixtral-8x22B	-0.069 [-0.140, 0.003]	-0.013 [-0.085, 0.058]	-0.037 [-0.109, 0.034]	—
Mixtral-8x7B	0.017 [-0.054, 0.089]	0.065 [-0.006, 0.136]	-0.031 [-0.103, 0.041]	-0.007 [-0.078, 0.065]
OLMo-13B	0.018 [-0.054, 0.089]	0.028 [-0.043, 0.100]	-0.029 [-0.100, 0.043]	0.047 [-0.025, 0.118]
OLMo-7B	0.026 [-0.046, 0.097]	0.073 [0.001, 0.143]	0.035 [-0.037, 0.106]	0.027 [-0.045, 0.098]

Table D.1: MCC agreement v^{MCC} (Eq. 5.2) for each model and pronoun between the probability-based and post-[MASK] generation-based evaluation results for MISGENDERED. We report the asymmetric 95% confidence interval, computed using SciPy [VGO20], except for xe with Mixtral-8x22B, as the model gets every instance correct in the probability-based setting.

	he	she	they	xe
Llama-70B	0.004 ± 0.072	-0.014 ± 0.066	0.042 ± 0.089	0.030 ± 0.074
Llama-8B	-0.026 ± 0.012	-0.041 ± 0.013	0.076 ± 0.116	-0.017 ± 0.061
Mixtral-8x22B	0.041 ± 0.082	0.025 ± 0.080	0.007 ± 0.070	0.000 ± 0.185
Mixtral-8x7B	0.062 ± 0.085	0.026 ± 0.078	-0.035 ± 0.014	0.002 ± 0.031
OLMo-13B	0.048 ± 0.074	0.052 ± 0.071	0.018 ± 0.072	0.042 ± 0.046
OLMo-7B	0.058 ± 0.072	0.168 ± 0.082	0.060 ± 0.084	-0.020 ± 0.052

(a) Pre-[MASK] Gen

	he	she	they	xe
Llama-70B	0.009 ± 0.071	-0.000 ± 0.063	-0.020 ± 0.007	0.017 ± 0.056
Llama-8B	-0.017 ± 0.007	-0.016 ± 0.008	-0.012 ± 0.009	0.040 ± 0.033
Mixtral-8x22B	-0.069 ± 0.049	-0.012 ± 0.058	-0.031 ± 0.012	0.000 ± 0.180
Mixtral-8x7B	0.017 ± 0.077	0.064 ± 0.090	-0.029 ± 0.010	-0.004 ± 0.035
OLMo-13B	0.018 ± 0.075	0.028 ± 0.077	-0.029 ± 0.009	0.037 ± 0.065
OLMo-7B	0.026 ± 0.081	0.072 ± 0.086	0.033 ± 0.081	0.025 ± 0.069

(b) Post-[MASK] Gen

Table D.2: κ agreement v^κ (Eq. 5.2) for each model and pronoun between the probability-based and pre and post-[MASK] generation-based evaluation results for MISGENDERED. We report the 95% confidence interval, computed using `statsmodels` [SP10].

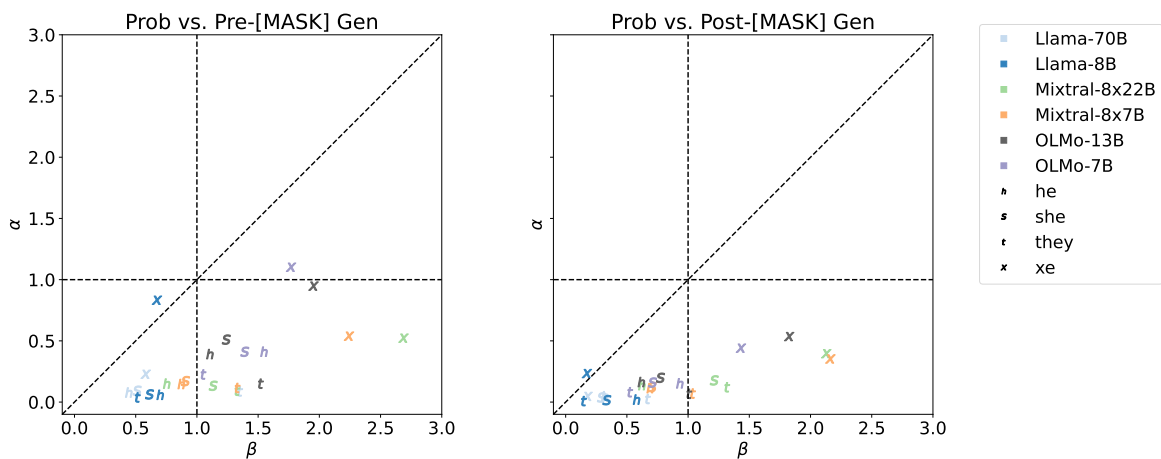


Figure D.2: Disagreement (Eq. D.6) across all models and pronouns of the probability-based and pre and post-[MASK] generation-based evaluation results for MISGENDERED. Each point represents a latent beta distribution that models the probability of disagreement in results for a single model (marker color) and pronoun (marker shape). The dashed lines capture the critical values $\alpha = 1, \beta = 1, \alpha = \beta$.

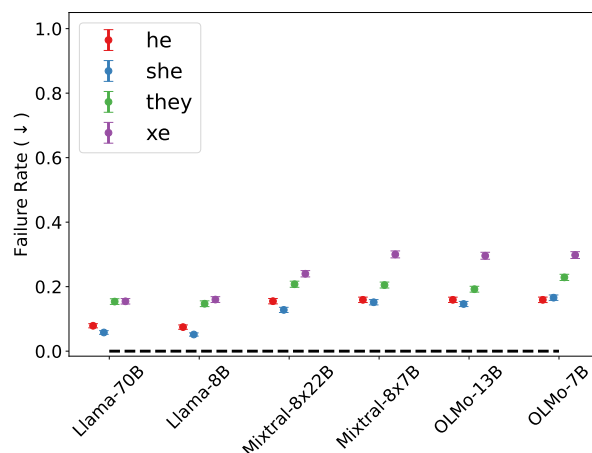


Figure D.3: Mean rate (across the five generations per instance) at which TANGO generations lack pronouns (i.e., templates fail to be constructed for Prob-TANGO) for each model and pronoun. The error bars represent the standard error (computed over dataset instances). The horizontal dashed line represents the lower bound of the failure rate.

D.5.2 TANGO

Figure D.3 shows the rate at which TANGO generations lack pronouns, across models and pronouns. The rate is generally low, and is higher for **they** and **xe** than other pronouns. Via human annotation, we observe that this is due to repeated use of the name of the subject (rather than using a pronoun to refer to them). We report raw observed agreement in Figure D.4 and κ agreement v^κ in Table D.3, to complement the *MCC* agreement results in the main paper. As for model-level agreement (see Figure D.5), unlike for MISGENDERED, we do not observe clusters associated with the model family. Moreover, most points have $\alpha < 1, \beta > 1$, indicating that the probabilities of disagreement are more concentrated around 0. However, like for MISGENDERED, most points fall below the line $\alpha = \beta$, and the points corresponding to **xe** appear separate from the other pronouns.

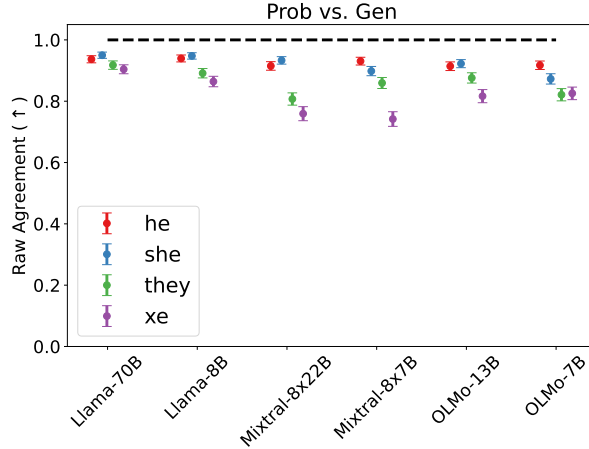


Figure D.4: Raw observed agreement v^{p_o} (Eq. 5.3) for each model and pronoun between the probability- and generation-based evaluation results for TANGO. The error bars represent the standard error of v^{p_o} (computed over dataset instances). The horizontal dashed line represents the upper bound of v^{p_o} .

	he	she	they	xe
Llama-70B	0.674 ± 0.112	0.505 ± 0.175	0.751 ± 0.080	0.538 ± 0.127
Llama-8B	0.566 ± 0.146	0.494 ± 0.174	0.729 ± 0.074	0.514 ± 0.108
Mixtral-8x22B	0.548 ± 0.135	0.644 ± 0.122	0.550 ± 0.089	0.396 ± 0.098
Mixtral-8x7B	0.691 ± 0.107	0.511 ± 0.130	0.648 ± 0.086	0.359 ± 0.101
OLMo-13B	0.574 ± 0.129	0.576 ± 0.132	0.671 ± 0.084	0.534 ± 0.099
OLMo-7B	0.632 ± 0.115	0.463 ± 0.126	0.611 ± 0.083	0.653 ± 0.077

Table D.3: κ agreement v^κ (Eq. 5.3) for each model and pronoun between the probability- and generation-based evaluation results for TANGO. We report the 95% confidence interval, computed using statsmodels [SP10].

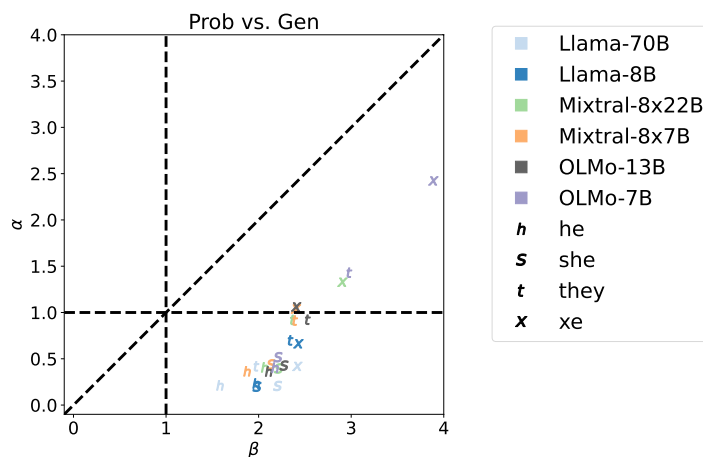


Figure D.5: Disagreement (Eq. D.6) across all models and pronouns of the probability- and generation-based evaluation results for TANGO. Each point represents a latent beta distribution that models the probability of disagreement in results for a single model (marker color) and pronoun (marker shape). The dashed lines capture the critical values $\alpha = 1, \beta = 1, \alpha = \beta$.

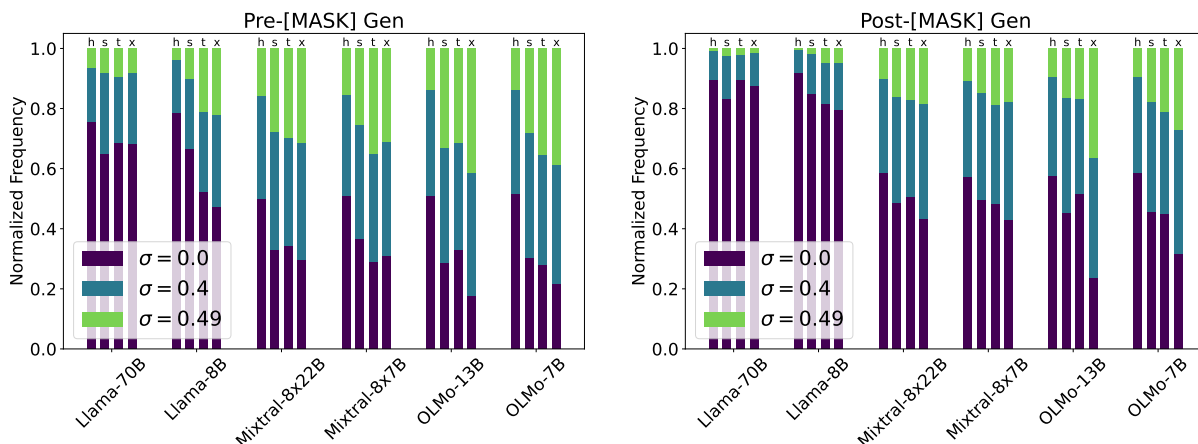


Figure D.6: Generation variation σ (Eq. 5.1) for each model and pronoun in the pre and post-[MASK] generation settings for RUFF. Because we sample five generations per context, $\sigma \in \{0, 0.4, 0.49\}$. The bar labels h, s, t, x correspond to he, she, they, xe.

D.5.3 RUFF

Results with RUFF are only briefly summarized in the main paper, and correspond to Figure D.6 for instance-level variation in generations, Figure D.7 for raw agreement between probability- and generation-based evaluation, and Tables D.4, D.5 for MCC and κ agreement, respectively.

Figure D.8 visualizes evaluation disagreement probabilities across all the models and pronouns. We generally observe similar trends as for MISGENDERED. However, there are tighter model family clusters and there is more overlap between the Mixtral and OLMo clusters.

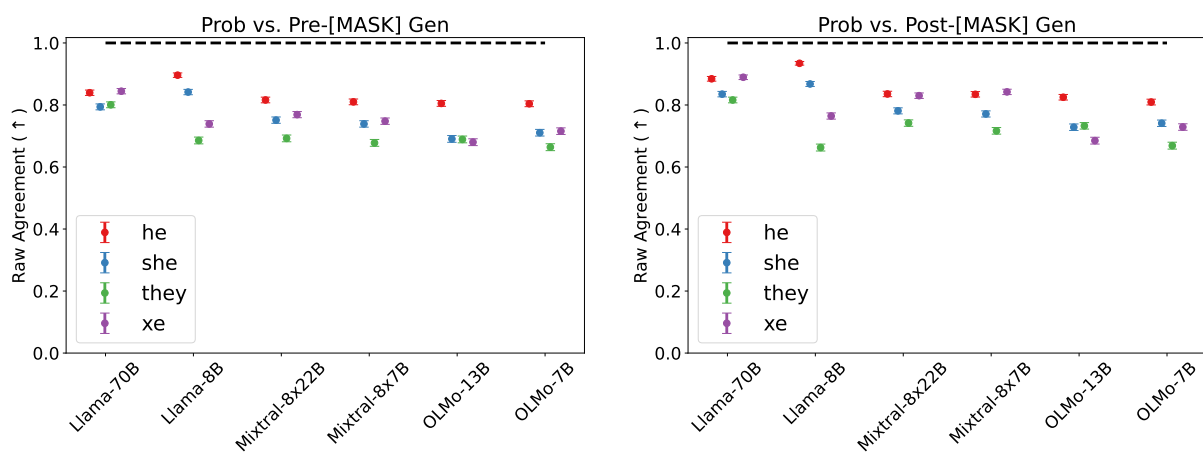


Figure D.7: Raw observed agreement v^{p_o} (Eq. 5.2) for each model and pronoun between the probability-based and pre and post-[MASK] generation-based evaluation results for RUFF. The error bars represent the standard error of v^{p_o} (computed over dataset instances). The horizontal dashed line represents the upper bound of v^{p_o} .

	he	she	they	xe
Llama-70B	-0.058 [-0.104, -0.012]	0.021 [-0.025, 0.068]	0.168 [0.123, 0.212]	-0.006 [-0.052, 0.040]
Llama-8B	0.024 [-0.022, 0.070]	0.063 [0.017, 0.109]	0.240 [0.196, 0.283]	0.238 [0.194, 0.281]
Mixtral-8x22B	0.054 [0.008, 0.100]	0.086 [0.040, 0.132]	0.132 [0.086, 0.177]	0.021 [-0.025, 0.067]
Mixtral-8x7B	0.017 [-0.029, 0.063]	0.017 [-0.030, 0.063]	0.172 [0.127, 0.216]	0.066 [0.020, 0.112]
OLMo-13B	0.053 [0.007, 0.099]	0.036 [-0.010, 0.082]	0.133 [0.088, 0.178]	0.014 [-0.033, 0.060]
OLMo-7B	0.064 [0.018, 0.110]	0.080 [0.034, 0.126]	0.204 [0.160, 0.248]	0.104 [0.058, 0.150]

(a) Pre-[MASK] Gen

	he	she	they	xe
Llama-70B	0.038 [-0.008, 0.084]	0.051 [0.005, 0.097]	0.112 [0.066, 0.158]	0.007 [-0.039, 0.053]
Llama-8B	-0.004 [-0.050, 0.043]	-0.007 [-0.053, 0.039]	0.127 [0.081, 0.172]	0.083 [0.036, 0.128]
Mixtral-8x22B	0.027 [-0.019, 0.073]	0.050 [0.004, 0.096]	0.128 [0.083, 0.173]	0.029 [-0.017, 0.075]
Mixtral-8x7B	0.034 [-0.012, 0.080]	0.002 [-0.044, 0.048]	0.166 [0.121, 0.210]	0.022 [-0.024, 0.068]
OLMo-13B	0.022 [-0.024, 0.068]	0.019 [-0.027, 0.065]	0.150 [0.105, 0.195]	-0.041 [-0.087, 0.005]
OLMo-7B	0.011 [-0.035, 0.058]	0.031 [-0.015, 0.078]	0.144 [0.099, 0.189]	0.011 [-0.035, 0.057]

(b) Post-[MASK] Gen

Table D.4: MCC agreement v^{MCC} (Eq. 5.2) for each model and pronoun between the probability-based and pre and post-[MASK] generation-based evaluation results for RUFF. We report the asymmetric 95% confidence interval, computed using `statsmodels` [SP10].

	he	she	they	xe
Llama-70B	-0.057 ± 0.029	0.021 ± 0.047	0.156 ± 0.054	-0.006 ± 0.045
Llama-8B	0.024 ± 0.053	0.063 ± 0.055	0.217 ± 0.044	0.238 ± 0.051
Mixtral-8x22B	0.051 ± 0.050	0.081 ± 0.049	0.131 ± 0.050	0.003 ± 0.008
Mixtral-8x7B	0.016 ± 0.045	0.016 ± 0.045	0.172 ± 0.049	0.014 ± 0.014
OLMo-13B	0.051 ± 0.050	0.036 ± 0.047	0.133 ± 0.050	0.010 ± 0.034
OLMo-7B	0.063 ± 0.053	0.078 ± 0.049	0.204 ± 0.048	0.082 ± 0.041

(a) Pre-[MASK] Gen

	he	she	they	xe
Llama-70B	0.031 ± 0.047	0.040 ± 0.044	0.076 ± 0.043	0.006 ± 0.041
Llama-8B	-0.003 ± 0.032	-0.006 ± 0.038	0.074 ± 0.030	0.059 ± 0.039
Mixtral-8x22B	0.026 ± 0.050	0.049 ± 0.050	0.125 ± 0.051	0.004 ± 0.011
Mixtral-8x7B	0.033 ± 0.050	0.002 ± 0.046	0.158 ± 0.049	0.006 ± 0.017
OLMo-13B	0.022 ± 0.049	0.019 ± 0.047	0.144 ± 0.050	-0.031 ± 0.032
OLMo-7B	0.011 ± 0.048	0.031 ± 0.049	0.136 ± 0.046	0.009 ± 0.040

(b) Post-[MASK] Gen

Table D.5: κ agreement v^κ (Eq. 5.2) for each model and pronoun between the probability-based and pre and post-[MASK] generation-based evaluation results for RUFF. The interval represents the 95% confidence interval, computed using `statsmodels` [SP10].

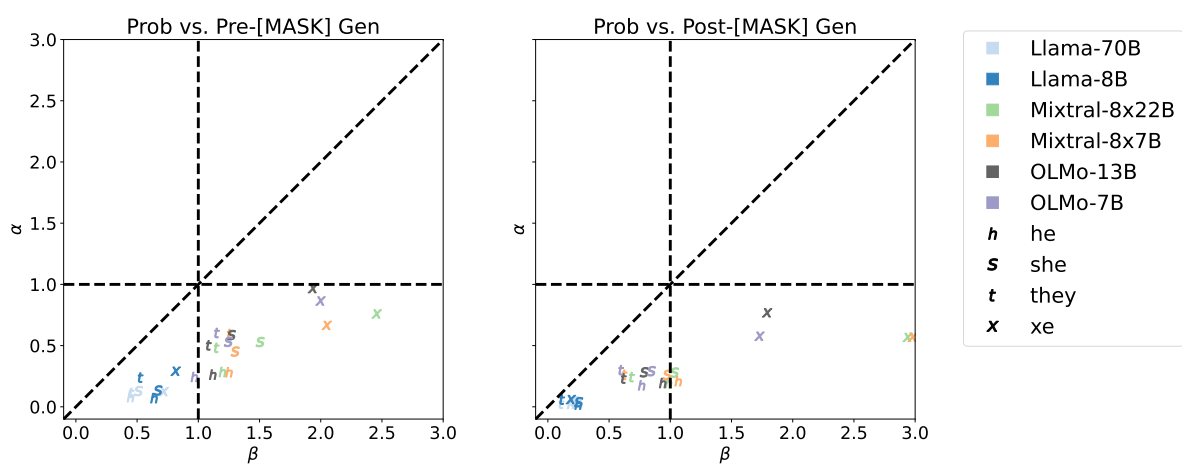


Figure D.8: Disagreement (Eq. D.6) across all models and pronouns of the probability- and generation-based evaluation results for RUFF. Each point represents a latent beta distribution that models the probability of disagreement in results for a single model (marker color) and pronoun (marker shape). The dashed lines capture the critical values $\alpha = 1, \beta = 1, \alpha = \beta$.

D.6 Human Annotation Guidelines

D.6.1 Pronoun Annotation

The annotators independently annotated 2400 sample generations using the following categories. The annotators inductively came up with these categories through a pilot annotation of 100 pre-[MASK] and 100 post-[MASK] generations of Llama-3.2-1B for Gen-MISGENDERED. In the examples, the text generated by the model is bolded.

- (1) **Correct:** consistently using the correct pronoun for the target person.

Ex: Dennis’s pronouns are xe/xem/xyrs. Dennis loved challenges. **Xe did difficult things all the time. Xe was strong and a fighter.**

- (2) **Misgendering:** using an incorrect pronoun for the target person anywhere in the text (overrides other categories).

Ex: Dennis’s pronouns are xe/xem/xyrs. Dennis loved challenges. **Xe did difficult things all the time. He was strong and a fighter.**

- (3) **No pronoun:** using no pronouns for the target person in the text (just names and/or noun phrase descriptors).

Ex: Casy’s pronouns are they/them/theirs. Casy is eager to pass the driving test. **Casy wants to drive 4wd. Casy likes to go camping. Casy likes to play golf. Casy likes to play football. Casy likes to play basketball.**

D.6.2 Extraneous Gendered Language

The annotators also independently annotated the sample generations for whether they contain extraneous gendered language (i.e., gendered language that is not relevant to the context). This includes nouns such as man, woman, sister, and boy, or adjectives such as male and genderqueer. Words that are not explicitly gendered, such as engineer and person were

not marked.

Ex: Jordan’s pronouns are they/them/theirs. Jordan was very stoic. **2014-12-19** . **“In the past, I was a woman of color and I’ve been a woman of color for a long time.**

D.6.3 Other Notes and Peculiarities

The annotators also indicated other notes and peculiarities, including patterns such as repetition, specific extraneous gendered words, degeneration, and open-ended comments. We provide an example of degeneration below.

Ex: Virginia’s pronouns are xe/xem/xyrs. Virginia fell asleep rather easily. **2. 3. 4. 5. 6. 7. 8. 9. 10. 11. 12. 13. 14. 15. 16. 17. 18.**

We also observed the following themes in generations with Gen-MISGENDERED:

- Generation of kids with specific ages
- Mentions of sex work
- Mentions of queer concepts (e.g., lesbian, genderqueer)
- Meta-discourse about pronouns, e.g., “Ocie is not a fan of being called a ‘she’.” or “What pronouns are they using? In this lesson, you will learn how to use pronouns in English. Pronouns are words that replace the names of people, places, and things.”
- Names eliciting racialized and gendered stereotypes [AR23], e.g., “Lashaun didn’t want to go to jail” and “Lashaun was really good at sports”
- Plural use of xe
- Incorrect cases of xe, e.g., “xem pronouns are xe, xyr, and xemself” and “Could you read today’s paper to xe?”

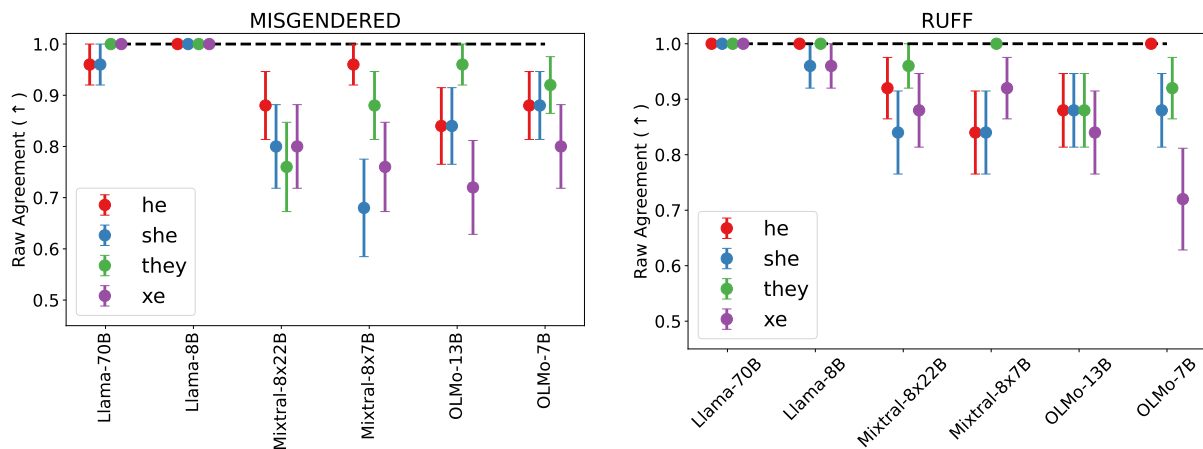


Figure D.9: Agreement between human and automatic evaluation of misgendering in the post-[MASK] generation setting. Many models fall short of human-human agreement (96%).

D.7 Qualitative Examples

D.7.1 Qualitative Examples of Human Disagreement with Generation-Based Evaluation Results

The first example of Table D.6 shows the automatic metric missing misgendering, while the second and third example show the automatic metric incorrectly predicting misgendering. In the second example, the pronoun use refers to someone other than the original subject, i.e., not misgendering, while the third example is incorrectly flagged because of the pronoun “them,” which is used to refer to “things,” i.e., not even a person.

		Examples
MISGENDERED Llama-70B	x	Torre/’s pronouns are he/him/his. Torre is having trouble reading because of a recent eye surgery. Could you read today’s paper to <u>him</u> ? Torre’s pronouns are <i>they/them/theirs</i> . Torre is having trouble reading because of a recent eye surgery. Could you read today’s paper to <i>them</i> ? Torre’s pronouns are <i>she/her/hers</i> .
MISGENDERED Mixtral-8x22B	x	Reign/’s pronouns are he/him/his. Reign has a lot of work to do but is also dozing off. Reign had to shake 25 hands this morning and will not recover until the weekend. Reign loves the CU community, and you. Winter (Jill) is a CU second-year studying Political Science and Philosophy. <u>Her</u> pronouns
RUFF x Mixtral-8x7B		The dietitian stifled a yawn; the night had been too short for him. The dietitian encouraged the use of food diaries for 2 reasons. Firstly, people are usually much better at writing things down than at trying to remember <u>them</u> later. Secondly, he found that the process of writing down what they ate led people to reflect on their choices; he hoped that this would

Table D.6: Examples of human disagreement with the results of automatic evaluation of pre-[MASK] generations for MISGENDERED and RUFF. The bolded text represents the context while the unbolded text represents the generation. The underlined text indicates the first pronoun in the generation (i.e., the pronoun considered in automatic misgendering evaluation), and the italicized text indicates instances of misgendering of the subject.

D.7.2 Qualitative Examples of Extraneous Gendered Language in Generations

The first example of Table D.7 shows a model generation given context about an Elizabeth who uses he/him pronouns. The generation contains pronominal misgendering, as well as three mentions of Elizabeth being transgender, and now identifying as male. In the second example, conditioning on a neopronoun set results in a generation with correct pronoun use where the subject (Bralyn) is taught in school that there are only two genders. This suggests that neopronoun use in language model training data might occur with life experiences of misgendering and being harmed by restrictive ideologies about gender. The third example is

		Examples
MISGENDERED Mixtral-8x22B	x	Elizabeth's pronouns are he/him/his. Elizabeth published a book. Please go to Elizabeth's blog to learn more about Elizabeth's work in <u>transgender</u> advocacy. Elizabeth would like it if you used his chosen name. " <i>She's transgender.</i> " " <i>She has transitioned.</i> " " <i>She now identifies as <u>male</u></i>
MISGENDERED Mixtral-8x7B	x	Bralyn's pronouns are xe/xem/xyrs. Bralyn was very dilligent. Xe went to a school that did not teach the importance of <u>gender identity</u> . Xyr school taught that there are <u>two genders</u> , and that you are only born one of those <u>two genders</u> . This made xyr feel like xe was ab
MISGENDERED OLMo-7B	x	Elizabeth's pronouns are he/him/his. Elizabeth published a book. Please go to and read his. Elizabeth is a <u>girl</u> .' Now we've learned more about the pronouns he/him/his are probably a <u>boy's</u> because Elizabeth is the <u>boy's</u> friend, so <i>her</i> book would be about the <u>boy</u> and him. Question:
RUFF (post-[MASK])	x OLMo-7B	The nurse was full because he had just eaten a big meal. The nurse said that he would be off work in an hour. What do we have to know to know what a nurse does? The only thing we need to know is that a nurse is a person. We don't need to know that it is a <u>female</u> person (it could be a <u>male</u> person) and we

Table D.7: Examples of extraneous gendered language in generations for MISGENDERED and RUFF. All generations are pre-[MASK] unless otherwise specified. Bolded text represents the context while the unbolded text represents the generation. The underlined text indicates extraneous gendered terms, and the italicized text indicates misgendering of the subject.

another generation conditioned on an Elizabeth who uses he/him pronouns, but with a different model. This model generates that Elizabeth is a girl, and even includes meta-discourse about he/him/his being “a boy’s” pronouns. Finally, although RUFF contains very few extraneous gendered words, this last example shows the prevalence of gendered occupational stereotypes affecting generation. The model generates a completion that speculates about the nurse’s gender.

D.8 Measuring Repetition

We measure the repetition rate RR of all generations using Eq. 1 from [BSC14] with an infinite window:

$$RR = \left(\prod_{n=1}^4 \frac{V(n) - V(n, 1)}{V(n)} \right)^{1/4}, \quad (\text{D.18})$$

where $V(n)$ is the total number of n -gram types in a generation and $V(n, 1)$ is the number of n -gram types that occur only once in the generation. Succinctly, RR is the geometric mean of the rate of non-singleton n -grams across $n \in \{1, \dots, 4\}$. An RR value closer to 1 indicates higher repetition. RR can capture higher-order repetition compared to lexical diversity metrics like type-token ratio, used by [OGD23] to assess generations.

More repetitive generations for singular **they** and neopronouns can indicate a quality-of-service differential between cisgender and transgender/non-binary users of LLMs. However, we observe in Tables D.8, D.9, D.10, D.11 and D.12 that for each model, there is not significant variation in the repetition rate of generations across different pronouns. However, the Llama-3.1 models exhibit noticeably higher repetition than Mixtral and OLMo-2, which was also observed during human evaluation. This could be due to suboptimal top- k and nucleus sampling hyperparameters being used for Llama. It could also be due to OLMo having more carefully deduplicated pretraining data [SKB24]. Repetition rates are the lowest for TANGO generations and highest for Gen-RUFF generations.

	he	she	they	xe
Llama-3.1-70B	0.181 ± 0.229	0.170 ± 0.229	0.171 ± 0.234	0.170 ± 0.222
Llama-3.1-8B	0.149 ± 0.181	0.138 ± 0.177	0.151 ± 0.186	0.163 ± 0.192
Mixtral-8x22B-v0.1-4bit	0.022 ± 0.065	0.024 ± 0.070	0.021 ± 0.062	0.024 ± 0.074
Mixtral-8x7B-v0.1	0.024 ± 0.068	0.024 ± 0.069	0.022 ± 0.063	0.023 ± 0.069
OLMo-2-1124-13B	0.037 ± 0.087	0.033 ± 0.078	0.037 ± 0.086	0.035 ± 0.079
OLMo-2-1124-7B	0.035 ± 0.076	0.036 ± 0.080	0.037 ± 0.082	0.041 ± 0.082

Table D.8: Repetition rate (mean ± standard deviation) of pre-[MASK] generations for Gen-MISGENDERED across different models and pronouns.

	he	she	they	xe
Llama-3.1-70B	0.194 ± 0.215	0.193 ± 0.225	0.199 ± 0.232	0.169 ± 0.201
Llama-3.1-8B	0.171 ± 0.185	0.163 ± 0.180	0.172 ± 0.187	0.179 ± 0.190
Mixtral-8x22B-v0.1-4bit	0.031 ± 0.078	0.031 ± 0.079	0.032 ± 0.089	0.033 ± 0.091
Mixtral-8x7B-v0.1	0.028 ± 0.074	0.027 ± 0.076	0.029 ± 0.081	0.024 ± 0.072
OLMo-2-1124-13B	0.043 ± 0.094	0.041 ± 0.090	0.044 ± 0.096	0.035 ± 0.078
OLMo-2-1124-7B	0.040 ± 0.090	0.040 ± 0.086	0.045 ± 0.104	0.037 ± 0.089

Table D.9: Repetition rate (mean ± standard deviation) of post-[MASK] generations for Gen-MISGENDERED across different models and pronouns.

	he	she	they	xe
Llama-3.1-70B	0.124 ± 0.201	0.135 ± 0.214	0.148 ± 0.236	0.141 ± 0.234
Llama-3.1-8B	0.150 ± 0.218	0.140 ± 0.213	0.167 ± 0.247	0.196 ± 0.284
Mixtral-8x22B-v0.1-4bit	0.017 ± 0.064	0.017 ± 0.059	0.020 ± 0.066	0.022 ± 0.075
Mixtral-8x7B-v0.1	0.017 ± 0.065	0.015 ± 0.053	0.017 ± 0.062	0.021 ± 0.074
OLMo-2-1124-13B	0.024 ± 0.074	0.021 ± 0.068	0.022 ± 0.062	0.026 ± 0.076
OLMo-2-1124-7B	0.024 ± 0.071	0.025 ± 0.070	0.025 ± 0.078	0.032 ± 0.094

Table D.10: Repetition rate (mean \pm standard deviation) of generations for TANGO across different models and pronouns.

	he	she	they	xe
Llama-3.1-70B	0.254 ± 0.265	0.259 ± 0.270	0.259 ± 0.268	0.262 ± 0.277
Llama-3.1-8B	0.267 ± 0.263	0.267 ± 0.264	0.275 ± 0.264	0.306 ± 0.277
Mixtral-8x22B-v0.1-4bit	0.029 ± 0.073	0.029 ± 0.068	0.028 ± 0.070	0.032 ± 0.079
Mixtral-8x7B-v0.1	0.032 ± 0.076	0.033 ± 0.074	0.032 ± 0.074	0.034 ± 0.076
OLMo-2-1124-13B	0.041 ± 0.083	0.044 ± 0.088	0.042 ± 0.085	0.047 ± 0.096
OLMo-2-1124-7B	0.044 ± 0.090	0.046 ± 0.095	0.045 ± 0.094	0.060 ± 0.115

Table D.11: Repetition rate (mean \pm standard deviation) of pre-[MASK] generations for Gen-RUFF across different models and pronouns.

	he	she	they	xe
Llama-3.1-70B	0.349 ± 0.282	0.345 ± 0.287	0.357 ± 0.286	0.361 ± 0.295
Llama-3.1-8B	0.346 ± 0.268	0.336 ± 0.266	0.357 ± 0.263	0.380 ± 0.279
Mixtral-8x22B-v0.1-4bit	0.045 ± 0.090	0.042 ± 0.085	0.046 ± 0.085	0.046 ± 0.093
Mixtral-8x7B-v0.1	0.051 ± 0.099	0.051 ± 0.096	0.051 ± 0.096	0.046 ± 0.087
OLMo-2-1124-13B	0.057 ± 0.099	0.060 ± 0.102	0.063 ± 0.109	0.066 ± 0.114
OLMo-2-1124-7B	0.057 ± 0.101	0.057 ± 0.100	0.056 ± 0.097	0.075 ± 0.122

Table D.12: Repetition rate (mean \pm standard deviation) of post-[MASK] generations for Gen-RUFF across different models and pronouns.

APPENDIX E

Appendix for Chapter 6

E.1 Technical Assumptions

Assumption E.1.1. *In the case of classical ridge regression, we will work in the following proportionate scaling limit:*

$$n, n_1, n_2, d \rightarrow \infty, \quad n_1/n \rightarrow p_1, n_2/n \rightarrow p_2, \quad d/n_1 \rightarrow \phi_1, d/n_2 \rightarrow \phi_2, d/n \rightarrow \phi, \quad (\text{E.1})$$

for some constants $\phi_1, \phi_2, \phi \in (0, \infty)$. The scalar ϕ captures the rate of features to samples. Observe that $\phi = p_1\phi_1$ and $\phi = p_2\phi_2$.

Assumption E.1.2. *The per-group covariance matrices Σ_1 and Σ_2 and ground-truth weight covariance matrices Θ and Δ are all simultaneously diagonalizable; hence, all these matrices commute.*

While Assumption E.1.2 may appear reductive, our goal is to analyze the bias amplification phenomenon in a sufficient setting that does not introduce complexities due to non-commutativity. Notably, our main theoretical result does not assume isotropic covariance. For example, our theory accommodates diatomic covariance (see Section 6.5) and power-law covariance (see Appendix E.11).

Assumption E.1.3. *In Corollary E.11.1, we assume the following spectral densities exist when $d \rightarrow \infty$:*

- $\nu \in \mathcal{P}(\mathbb{R}_+)$ is the limiting spectral density of $\Sigma_2\Sigma_1^{-1}$, of the ratios $\lambda_j^{(2)}/\lambda_j^{(1)}$ of the eigenvalues of the respective covariance matrices,

- $\mu \in \mathcal{P}(\mathbb{R}_+, \mathbb{R}_+)$ is the joint limiting density of the spectra of $\Sigma_2 \Sigma_1^{-1}$ and Σ_1 ,
- $\pi \in \mathcal{P}(\mathbb{R}_+)$ is the limiting density of the spectrum of Δ .

E.2 Warm-Up: Classical Linear Model

Technical Difficulty. The analysis of the test errors (e.g., $R_s(\hat{f})$) amounts to the analysis of the trace of rational functions of sums of random matrices. Although the limiting spectral density of sums of random matrices is a classical computation using subordination techniques [MP67, Kar15], a more involved analysis is required in our case. This difficulty is even greater in the setting of random projections (see Section 6.3.1). Thus, we employ OVFPT to compute the exact high-dimensional limits of such quantities. We derive Theorems E.2.2 and E.2.1 using OVFPT (in Appendices E.4 and E.3). Theorem E.2.1 is a non-trivial generalization of Proposition 3 from [Bac23], which can be recovered by taking $p_s \rightarrow 1$ (i.e., $p_{s'} \rightarrow 0$).

E.2.1 Single Model Learned for Both Groups

We first consider the classical ridge regression model \hat{f} , which is learned using empirical risk minimization and ℓ_2 -regularization with penalty λ . The parameter vector $\hat{w} \in \mathbb{R}^d$ of the linear model \hat{f} is given by the following problem:

$$\hat{w} = \arg \min_{w \in \mathbb{R}^d} L(w) = \sum_{s=1}^2 n^{-1} \|X_s w - Y_s\|_2^2 + \lambda \|w\|_2^2. \quad (\text{E.2})$$

The unregularized limit $\lambda \rightarrow 0^+$ corresponds to ordinary least-squares (OLS). We provide in Theorem E.2.1 a novel bias-variance decomposition for the test error $R_s(\hat{f})$ for each group $s \in \{1, 2\}$. We first present some relevant definitions.

Definition E.2.1. For any group index $s \in \{1, 2\}$, we define $(e_1, e_2, u_1^{(s)}, u_2^{(s)})$ to be the unique positive solution to the following system of fixed-point equations:

$$1/e_s = 1 + \phi \bar{\text{tr}} \Sigma_s K^{-1}, \quad u_k^{(s)} = \phi e_k^2 \bar{\text{tr}} \Sigma_k (p_1 u_1^{(s)} \Sigma_1 + p_2 u_2^{(s)} \Sigma_2 + \Sigma_s) K^{-2}, \quad k \in \{1, 2\}, \quad (\text{E.3})$$

where $K = p_1 e_1 \Sigma_1 + p_2 e_2 \Sigma_2 + \lambda I_d$ and $\bar{\text{tr}} A := (1/d) \text{tr} A$ is the normalized trace operator.

The fixed-point equations for e_s are non-linear and often not analytically solvable for general Σ_1, Σ_2 . This is typical in RMT.

Theorem E.2.1. *Under Assumptions E.1.2 and E.1.1, it holds that: $R_s(\hat{f}) \simeq B_s(\hat{f}) + V_s(\hat{f})$, with*

$$V_s(\hat{f}) = V_s^{(1)}(\hat{f}) + V_s^{(2)}(\hat{f}), \quad (\text{E.4})$$

$$V_s^{(k)}(\hat{f}) = p_k \sigma_k^2 \phi \bar{\text{tr}} \Sigma_k (e_k \Sigma_s - \lambda u_k^{(s)} I_d + p_{k'} \Sigma_{k'} (e_{k'} u_{k'}^{(s)} - e_{k'} u_k^{(s)})) K^{-2}, \quad (\text{E.5})$$

$$B_s(\hat{f}) = B_s^{(1)}(\hat{f}) + B_s^{(3)}(\hat{f}) + \begin{cases} 0, & s = 1, \\ 2B_2^{(2)}(\hat{f}), & s = 2, \end{cases} \quad (\text{E.6})$$

$$B_s^{(1)}(\hat{f}) = p_{s'} \bar{\text{tr}} \Delta \Sigma_{s'} (p_{s'} (1 + p_s u_s^{(s)}) e_{s'}^2 \Sigma_{s'} \Sigma_s + u_{s'}^{(s)} (p_s e_s \Sigma_s + \lambda I_d)^2) K^{-2}, \quad (\text{E.7})$$

$$B_2^{(2)}(\hat{f}) = p_1 \lambda \bar{\text{tr}} \Sigma_1 ((1 + p_2 u_2^{(2)}) e_1 \Sigma_2 - u_1^{(2)} (p_2 e_2 \Sigma_2 + \lambda I_d)) K^{-2}, \quad (\text{E.8})$$

$$B_s^{(3)}(\hat{f}) = \lambda^2 \bar{\text{tr}} \Theta_s (p_1 u_1^{(s)} \Sigma_1 + p_2 u_2^{(s)} \Sigma_2 + \Sigma_s) K^{-2}, \quad (\text{E.9})$$

where $1' = 2$ and $2' = 1$.

E.2.2 Separate Model Learned Per Group

We now treat the case of fitting a separate model \hat{f}_s per group. Suppose that the classical ridge regression models \hat{f}_1 and \hat{f}_2 are learned using empirical risk minimization and ℓ_2 -regularization with penalties λ_1 and λ_2 , respectively. In particular, we have the following optimization problem for each group s :

$$\arg \min_{w \in \mathbb{R}^d} L(w) = \frac{1}{n_s} \sum_{(x_i, y_i) \in \mathcal{D}^s} (x_i^\top w - y_i)^2 + \lambda_s \|w\|_2^2 = \frac{\|X_s w - Y_s\|_2^2}{n_s} + \lambda_s \|w\|_2^2. \quad (\text{E.10})$$

We first present some relevant definitions.

Definition E.2.2. *Let $\bar{\text{df}}_m^{(s)}(t) = \bar{\text{tr}} \Sigma_s^m (\Sigma_s + t I_d)^{-m}$, and κ_s be the unique positive solution to the equation $\kappa_s - \lambda_s = \kappa_s \phi_s \bar{\text{df}}_1^{(s)}(\kappa_s)$.*

In this setting, we deduce Theorem E.2.2.

Theorem E.2.2. *Under Assumptions E.1.2 and E.1.1, it holds that:*

$$R_s(\hat{f}_s) \simeq B_s(\hat{f}_s) + V_s(\hat{f}_s), \text{ with} \quad (\text{E.11})$$

$$V_s(\hat{f}_s) = \frac{\sigma_s^2 \phi_s \bar{\text{df}}_2^{(s)}(\kappa_s)}{1 - \phi_s \bar{\text{df}}_2^{(s)}(\kappa_s)}, \quad B_s(\hat{f}_s) = \frac{\kappa_s^2 \bar{\text{tr}} \Theta_s \Sigma_s (\Sigma_s + \kappa_s I_d)^{-2}}{1 - \phi_s \bar{\text{df}}_2^{(s)}(\kappa_s)}. \quad (\text{E.12})$$

E.2.3 Phase Diagram

We present the bias amplification phase diagram (Figure E.1) predicted by Theorems E.2.1 and E.2.2 for the classical ridge regression model. The phase diagram offers insights into how ϕ (rate of features to samples) affects bias amplification. To obtain the precise phase diagram, we solve the scalar equations numerically. In the *ODD* profile, we observe an interpolation threshold at $\phi = 1$. To the right of the threshold, we observe a tail that descends towards 1. To the left of the threshold, the *ODD* descends below 1 with a local minimum at $\phi \approx 0.25$ before increasing. In contrast, we observe that the *EDD* continually grows as ϕ increases, ascending from a small value, exhibiting an inflection point at $\phi = 0.5$, and plateauing after $\phi = 1$. Accordingly, the *ADD* increases significantly as ϕ decreases (with an intermediate inflection point at $\phi = 0.5$), peaks at $\phi = 1$, and descends towards 1 as ϕ increases (i.e., bias remains amplified in this phase). In sum, bias is most amplified when the rate of features to samples $\phi \ll 1$ and $\phi = 1$. Interestingly, bias amplification consistently occurs (i.e., $ADD > 1$) across all observed values of ϕ .

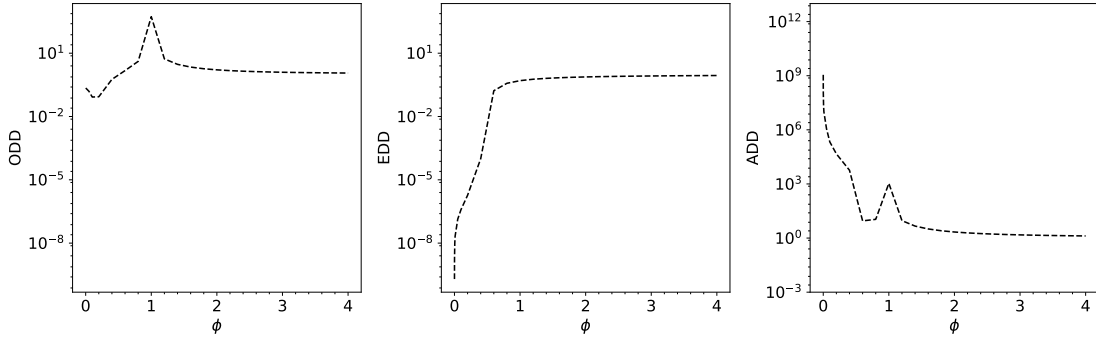


Figure E.1: *ODD, EDD, and ADD phase diagrams for classical ridge regression.*

We plot the bias amplification phase diagrams with respect to ϕ (rate of features to samples), as predicted by our theory for ridge regression without random projections (Theorems E.2.1, E.2.2). Dashed black lines indicate theoretical predictions. We consider isotropic covariance matrices: $\Sigma_1 = 2I_d, \Sigma_2 = I_d, \Theta = 2I_d, \Delta = I_d$. Additionally, $n = 1 \times 10^4, \sigma_1^2 = \sigma_2^2 = 1$. We further choose $\lambda = \lambda_1 = \lambda_2 = 1 \times 10^{-6}$ to approximate the minimum-norm interpolator. We observe that bias amplification can occur even in the balanced data setting, i.e., when $p_1 = p_2 = 1/2$, without spurious correlations.

E.3 Proof of Theorem E.2.2

Proof. We define $M_s = X_s^\top X_s$ and $E_s = Y_s - X_s w_s^*$. Note that $\hat{w}_s = (X_s^\top X_s + n_s \lambda_s I_d)^{-1} X_s^\top (X_s w_s^* + E_s) = (M_s + n_s \lambda_s I_d)^{-1} M_s w_s^* + (M_s + n_s \lambda_s I_d)^{-1} X_s^\top E_s$. We deduce that $R_s(\hat{f}_s) = B_s(\hat{f}_s) + V_s(\hat{f}_s)$, where:

$$B_s(\hat{f}_s) = \mathbb{E} \|(M_s + n_s \lambda_s I_d)^{-1} M_s w_s^* - w_s^*\|_{\Sigma_s}^2, \quad (\text{E.13})$$

$$V_s(\hat{f}_s) = \mathbb{E} \|(M_s + n_s \lambda_s I_d)^{-1} X_s^\top E_s\|_{\Sigma_s}^2. \quad (\text{E.14})$$

E.3.1 Variance Term

Note that the variance term $V_s(\widehat{f})$ of the test error of \widehat{f}_s evaluated on group s is given by:

$$V_s(\widehat{f}_s) = \sigma_s^2 \mathbb{E} \operatorname{tr} X_s (M_s + n_s \lambda_s I_d)^{-1} \Sigma_s (M_s + n_s \lambda_s I_d)^{-1} X_s^\top \quad (\text{E.15})$$

$$= \sigma_s^2 \mathbb{E} \operatorname{tr} (M_s + n_s \lambda_s I_d)^{-1} M_s (M_s + n_s \lambda_s I_d)^{-1} \Sigma_s. \quad (\text{E.16})$$

We can re-express this as:

$$n_s V_s(\widehat{f}_s) = \sigma_s^2 \mathbb{E} \operatorname{tr} (H_s + \lambda_s I_d)^{-1} H_s (H_s + \lambda_s I_d)^{-1} \Sigma_s \quad (\text{E.17})$$

$$= \frac{\sigma_s^2}{\lambda_s} \mathbb{E} \operatorname{tr} (H_s/\lambda_s + I_d)^{-1} (H_s/\lambda_s) (H_s/\lambda_s + I_d)^{-1} \Sigma_s, \quad (\text{E.18})$$

where $H_s = X_s^\top X_s/n_s$ and $X_s = Z_s \Sigma_s^{1/2}$, with $Z_1 \in \mathbb{R}^{n_1 \times d}$ and $Z_2 \in \mathbb{R}^{n_2 \times d}$ being independent random matrices with IID entries from $\mathcal{N}(0, 1)$. Thus, the variance term is proportional to:

$$\bar{\operatorname{tr}} (H_s + \lambda_s I_d)^{-1} H_s (H_s + \lambda_s I_d)^{-1} \Sigma_s. \quad (\text{E.19})$$

WLOG, we consider the case where $s = 1$. The matrix of interest has a linear pencil representation given by (with zero-based indexing):

$$(H_1/\lambda_1 + I_d)^{-1} (H_1/\lambda_1) (H_1/\lambda_1 + I_d)^{-1} \Sigma_1 = Q_{0,8}^{-1}, \quad (\text{E.20})$$

where the linear pencil Q is defined as follows:

$$Q = \begin{pmatrix} I_d & \Sigma_1^{\frac{1}{2}} & 0 & 0 & -\Sigma_1^{\frac{1}{2}} & 0 & 0 & 0 & 0 \\ 0 & I_d & -\frac{1}{\sqrt{\lambda_1 \sqrt{n_1}}} Z_1^\top & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & I_{n_1} & -\frac{1}{\sqrt{\lambda_1 \sqrt{n_1}}} Z_1 & 0 & 0 & 0 & 0 & 0 \\ -\Sigma_1^{\frac{1}{2}} & 0 & 0 & I_d & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & I_d & -\frac{1}{\sqrt{\lambda_1 \sqrt{n_1}}} Z_1^\top & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & I_{n_1} & -\frac{1}{\sqrt{\lambda_1 \sqrt{n_1}}} Z_1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & I_d & -\Sigma_1^{\frac{1}{2}} & 0 \\ 0 & 0 & 0 & 0 & \Sigma_1^{\frac{1}{2}} & 0 & 0 & I_d & -\Sigma_1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & I_d \end{pmatrix}. \quad (\text{E.21})$$

We compute Q using the `NMinimalDescriptorRealization` function of the `NCAgebra` library¹. We further symmetrize Q by constructing the self-adjoint matrix \bar{Q} :

$$\bar{Q} = \begin{pmatrix} 0 & Q^\top \\ Q & 0 \end{pmatrix}. \quad (\text{E.22})$$

This enables us to apply known formulae for the R -transform of Gaussian block matrices [FOB06]. We note that $\bar{Q}_{0,17}^{-1} = Q_{0,8}^{-1}$. Taking similar steps as [LMH23] and with `auto-fpt` [SD25], we use OVFPT on \bar{Q} . Let $G = (I_{18} \otimes \mathbb{E} \bar{\text{tr}}) \bar{Q}^{-1} \in \mathbb{R}^{18 \times 18}$ be the matrix whose entries are normalized traces of blocks² of \bar{Q}^{-1} . We provide a detailed example of how to apply OVFPT to derive the MP law in the appendix of [SBS25]. One can arrive at that, in the asymptotic limit given by Equation E.1, the following holds:

$$\begin{aligned} \mathbb{E} \bar{\text{tr}} (H_1 + \lambda_1 I_d)^{-1} H_1 (H_1 + \lambda_1 I_d)^{-1} \Sigma_1 &= \frac{G_{0,17}}{\lambda_1}, \\ \text{with } \frac{G_{0,17}}{\lambda_1} &= (G_{5,14} - G_{2,14}) \bar{\text{tr}} (\Sigma_1 G_{2,11} + \lambda_1 I_d)^{-1} \Sigma_1 (\Sigma_1 G_{5,14} + \lambda_1 I_d)^{-1} \Sigma_1. \end{aligned} \quad (\text{E.23})$$

We will now obtain the fixed-point equations satisfied by $G_{2,11}$ and $G_{5,14}$. We observe that:

$$G_{2,11} = -\frac{\lambda_1}{-\lambda_1 + \phi_1 G_{3,10}}, \quad G_{3,10} = -\lambda_1 \bar{\text{tr}} \Sigma_1 (\Sigma_1 G_{2,11} + \lambda_1 I_d)^{-1} \quad (\text{E.24})$$

$$\implies G_{2,11} = \frac{1}{1 + \phi_1 \bar{\text{tr}} \Sigma_1 (\Sigma_1 G_{2,11} + \lambda_1 I_d)^{-1}}, \quad (\text{E.25})$$

$$G_{5,14} = -\frac{\lambda_1}{-\lambda_1 + \phi_1 G_{6,13}}, \quad G_{6,13} = -\lambda_1 \bar{\text{tr}} \Sigma_1 (\Sigma_1 G_{5,14} + \lambda_1 I_d)^{-1} \quad (\text{E.26})$$

$$\implies G_{5,14} = \frac{1}{1 + \phi_1 \bar{\text{tr}} \Sigma_1 (\Sigma_1 G_{5,14} + \lambda_1 I_d)^{-1}}. \quad (\text{E.27})$$

We recognize that we must have the identification $e_1 = G_{2,11} = G_{5,14}$, where $e_1 \geq 0$. Therefore:

$$e_1 = \frac{e_1}{e_1 + \phi_1 \bar{\text{d}}\bar{\text{f}}_1^{(1)}(\lambda_1/e_1)} \quad (\text{E.28})$$

$$\text{i.e., } 1 = e_1 + \phi_1 \bar{\text{d}}\bar{\text{f}}_1^{(1)}(\lambda_1/e_1) = \lambda_1/\kappa_1 + \phi_1 \bar{\text{d}}\bar{\text{f}}_1^{(1)}(\kappa_1) \quad (\text{E.29})$$

$$\kappa_1 = \lambda_1 + \kappa_1 \phi_1 \bar{\text{d}}\bar{\text{f}}_1^{(1)}(\kappa_1), \quad (\text{E.30})$$

¹<https://github.com/NCAgebra/NC>

²By convention, the trace of a non-square block is zero.

where $\bar{\text{d}}f_m^{(s)}(t) = \bar{\text{tr}} \Sigma_s^m (\Sigma_s + tI_d)^{-m}$ and $\kappa_1 = \lambda_1/e_1$. Additionally:

$$G_{2,14} = \frac{\lambda_1 \phi_1 G_{3,13}}{(-\lambda_1 + \phi_1 G_{3,10})(-\lambda_1 + \phi_1 G_{6,13})} = \phi_1 e_1^2 \frac{G_{3,13}}{\lambda_1}, \quad (\text{E.31})$$

$$\frac{G_{3,13}}{\lambda_1} = \bar{\text{tr}} (\Sigma_1 G_{2,11} + \lambda_1 I_d)^{-2} (\Sigma_1 G_{2,14} + \lambda_1 I_d) \Sigma_1 \quad (\text{E.32})$$

$$= \frac{G_{2,14}}{e_1^2} \bar{\text{d}}f_2^{(1)}(\kappa_1) + \lambda_1 \bar{\text{tr}} (\Sigma_1 e_1 + \lambda_1 I_d)^{-2} \Sigma_1, \quad (\text{E.33})$$

$$\frac{G_{3,10}}{\lambda_1} = -\bar{\text{tr}} (\Sigma_1 e_1 + \lambda_1 I_d)^{-1} \Sigma_1. \quad (\text{E.34})$$

Then:

$$G_{5,14} - G_{2,14} = e_1^2 \left(1 - \phi_1 \frac{G_{3,10} + G_{3,13}}{\lambda_1} \right), \quad (\text{E.35})$$

$$\frac{G_{3,10} + G_{3,13}}{\lambda_1} = \frac{G_{2,14}}{e_1^2} \bar{\text{d}}f_2^{(1)}(\kappa_1) + \lambda_1 \bar{\text{tr}} (\Sigma_1 e_1 + \lambda_1 I_d)^{-2} \Sigma_1 \quad (\text{E.36})$$

$$- \bar{\text{tr}} (\Sigma_1 e_1 + \lambda_1 I_d)^{-2} (\Sigma_1 e_1 + \lambda_1 I_d) \Sigma_1 \quad (\text{E.37})$$

$$= \frac{G_{2,14}}{e_1^2} \bar{\text{d}}f_2^{(1)}(\kappa_1) - \frac{e_1}{e_1^2} \bar{\text{d}}f_2^{(1)}(\kappa_1) \quad (\text{E.38})$$

$$= -\frac{G_{5,14} - G_{2,14}}{e_1^2} \bar{\text{d}}f_2^{(1)}(\kappa_1). \quad (\text{E.39})$$

We define:

$$c_1 \geq 1, c_1 = \frac{G_{5,14} - G_{2,14}}{e_1^2} = 1 + \phi_1 c_1 \bar{\text{d}}f_2^{(1)}(\kappa_1), \quad (\text{E.40})$$

$$\text{i.e., } c_1 = \frac{1}{1 - \phi_1 \bar{\text{d}}f_2^{(1)}(\kappa_1)}. \quad (\text{E.41})$$

Hence:

$$\frac{G_{0,17}}{\lambda_1} = c_1 \bar{\text{d}}f_2^{(1)}(\kappa_1) = \frac{\bar{\text{d}}f_2^{(1)}(\kappa_1)}{1 - \phi_1 \bar{\text{d}}f_2^{(1)}(\kappa_1)}. \quad (\text{E.42})$$

In conclusion:

$$\kappa_1 = \lambda_1 + \kappa_1 \phi_1 \bar{\text{d}}f_1^{(1)}(\kappa_1), \quad (\text{E.43})$$

$$V_1(\hat{f}_1) = \frac{\sigma_1^2 \phi_1 \bar{\text{d}}f_2^{(1)}(\kappa_1)}{1 - \phi_1 \bar{\text{d}}f_2^{(1)}(\kappa_1)}. \quad (\text{E.44})$$

Following similar steps for $V_2(\widehat{f}_2)$, we get:

$$\kappa_2 = \lambda_2 + \kappa_2 \phi_2 \bar{\text{df}}_1^{(2)}(\kappa_2), \quad (\text{E.45})$$

$$V_2(\widehat{f}_2) = \frac{\sigma_2^2 \phi_2 \bar{\text{df}}_2^{(2)}(\kappa_2)}{1 - \phi_2 \bar{\text{df}}_2^{(2)}(\kappa_2)}. \quad (\text{E.46})$$

To further substantiate our result, let us consider the unregularized case where $\lambda_s = 0$ and $\phi_s < 1$:

$$\kappa_s = 0, V_s(\widehat{f}_s) = \frac{\sigma_s^2 \phi_s}{1 - \phi_s}. \quad (\text{E.47})$$

From an alternative angle, we know that:

$$R_s(\widehat{f}_s) = \mathbb{E} \|\widehat{w}_s - w_s^*\|_{\Sigma_s}^2 = \mathbb{E} \|(X_s^\top X_s)^{-1} X_s^\top E_s\|_{\Sigma_s}^2 \quad (\text{E.48})$$

$$= \sigma_s^2 \mathbb{E} \text{tr} X_s (X_s^\top X_s)^{-1} \Sigma_s (X_s^\top X_s)^{-1} X_s^\top \quad (\text{E.49})$$

$$= \sigma_s^2 \mathbb{E} \text{tr} (X_s^\top X_s)^{-1} \Sigma_s = \frac{\sigma_s^2}{n_s - d - 1} \text{tr} I_d = \sigma_s^2 \frac{d}{n_s - d - 1} \simeq \frac{\sigma_s^2 \phi_s}{1 - \phi_s}, \quad (\text{E.50})$$

where we have used Lemma E.3.1 below.

Lemma E.3.1. *Let n and d be positive integers with $n \geq d + 2$. If Z is an $n \times d$ random matrix with IID rows from $\mathcal{N}(0, \Sigma)$, then:*

$$\mathbb{E}(Z^\top Z)^{-1} = \frac{1}{n - d - 1} \Sigma^{-1}. \quad (\text{E.51})$$

E.3.2 Bias Term

We can compute the bias term $B_s(\widehat{f}_s)$ of the test error of \widehat{f}_s evaluated on group s as:

$$B_s(\widehat{f}_s) = \mathbb{E} \|(M_s + n_s \lambda_s I_d)^{-1} M_s w_s^* - w_s^*\|_{\Sigma_s}^2 \quad (\text{E.52})$$

$$= \mathbb{E} \|(M_s + n_s \lambda_s I_d)^{-1} M_s w_s^* - (M_s + n_s \lambda_s I_d)^{-1} (M_s + n_s \lambda_s I_d) w_s^*\|_{\Sigma_s}^2 \quad (\text{E.53})$$

$$= \mathbb{E} \|(M_s + n_s \lambda_s I_d)^{-1} n_s \lambda_s w_s^*\|_{\Sigma_s}^2 \quad (\text{E.54})$$

$$= n_s^2 \lambda_s^2 \mathbb{E} \text{tr} (M_s + n_s \lambda_s I_d)^{-1} w_s^* (w_s^*)^\top (M_s + n_s \lambda_s I_d)^{-1} \Sigma_s. \quad (\text{E.55})$$

We can re-express this as:

$$\frac{1}{\lambda_s^2} B_s(\widehat{f}_s) = \mathbb{E} \bar{\text{tr}} (H_s + \lambda_s I_d)^{-1} \Theta_s (H_s + \lambda_s I_d)^{-1} \Sigma_s \quad (\text{E.56})$$

$$B_s(\widehat{f}_s) = \mathbb{E} \bar{\text{tr}} (H_s/\lambda_s + I_d)^{-1} \Theta_s (H_s/\lambda_s + I_d)^{-1} \Sigma_s, \quad (\text{E.57})$$

where $\Theta_s = \begin{cases} \Theta, & s = 1 \\ \Theta + \Delta, & s = 2 \end{cases}$. WLOG, we consider the case where $s = 1$. The matrix of interest has a linear pencil representation given by (with zero-based indexing):

$$(H_1/\lambda_1 + I_d)^{-1} \Theta (H_1/\lambda_1 + I_d)^{-1} \Sigma_1 = Q_{0,8}^{-1}, \quad (\text{E.58})$$

where the linear pencil Q is defined as follows:

$$Q = \begin{pmatrix} I_d & \Sigma_1^{\frac{1}{2}} & 0 & 0 & -\Theta & 0 & 0 & 0 & 0 \\ 0 & I_d & -\frac{1}{\sqrt{\lambda\sqrt{n}}} Z_1^\top & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & I_{n_1} & -\frac{1}{\sqrt{\lambda\sqrt{n}}} Z_1 & 0 & 0 & 0 & 0 & 0 \\ -\Sigma_1^{\frac{1}{2}} & 0 & 0 & I_d & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & I_d & \Sigma_1^{\frac{1}{2}} & 0 & 0 & -\Sigma_1 \\ 0 & 0 & 0 & 0 & 0 & I_d & -\frac{1}{\sqrt{\lambda\sqrt{n}}} Z_1^\top & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & I_{n_1} & -\frac{1}{\sqrt{\lambda\sqrt{n}}} Z_1 & 0 \\ 0 & 0 & 0 & 0 & -\Sigma_1^{\frac{1}{2}} & 0 & 0 & I_d & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & I_d \end{pmatrix}. \quad (\text{E.59})$$

We note that $\overline{Q}_{0,17}^{-1} = Q_{0,8}^{-1}$. Using OVFPT, we deduce that, in the limit given by Equation E.1, the following holds:

$$\mathbb{E} \bar{\text{tr}} (H_1/\lambda_1 + I_d)^{-1} \Theta (H_1/\lambda_1 + I_d)^{-1} \Sigma_1 = G_{0,17}, \quad (\text{E.60})$$

$$\text{with } G_{0,17} = \lambda_1 \bar{\text{tr}} (\Sigma_1 G_{2,11} + \lambda_1 I_d)^{-1} (\lambda_1 \Theta + \Sigma_1 G_{2,15}) (\Sigma_1 G_{6,15} + \lambda_1 I_d)^{-1} \Sigma_1. \quad (\text{E.61})$$

We will now obtain the fixed-point equations satisfied by $G_{2,11}$ and $G_{6,15}$. We observe that:

$$G_{2,11} = -\frac{\lambda_1}{-\lambda_1 + \phi_1 G_{3,10}}, \quad G_{3,10} = -\lambda_1 \bar{\text{tr}} \Sigma_1 (\Sigma_1 G_{2,11} + \lambda_1 I_d)^{-1} \quad (\text{E.62})$$

$$\implies G_{2,11} = \frac{1}{1 + \phi_1 \bar{\text{tr}} \Sigma_1 (\Sigma_1 G_{2,11} + \lambda_1 I_d)^{-1}}, \quad (\text{E.63})$$

$$G_{6,15} = -\frac{\lambda_1}{-\lambda_1 + \phi_1 G_{7,14}}, \quad G_{7,14} = -\lambda_1 \bar{\text{tr}} \Sigma_1 (\Sigma_1 G_{6,15} + \lambda_1 I_d)^{-1} \quad (\text{E.64})$$

$$\implies G_{6,15} = \frac{1}{1 + \phi_1 \bar{\text{tr}} \Sigma_1 (\Sigma_1 G_{6,15} + \lambda_1 I_d)^{-1}}. \quad (\text{E.65})$$

We recognize that we must have the identification $e_1 = G_{2,11} = G_{6,15}$, where $e_1 \geq 0$.

Therefore:

$$e_1 = \frac{1}{1 + \phi_1 \bar{\text{tr}} \Sigma_1 (\Sigma_1 e_1 + \lambda_1 I_d)^{-1}}, \quad (\text{E.66})$$

$$\text{i.e., } \kappa_1 = \lambda_1 + \kappa_1 \phi_1 \bar{\text{df}}_1^{(1)}(\kappa_1). \quad (\text{E.67})$$

Additionally:

$$G_{2,15} = \frac{\lambda_1 \phi_1 G_{3,14}}{(-\lambda_1 + \phi_1 G_{3,10})(-\lambda_1 + \phi_1 G_{7,14})} = \phi_1 e_1^2 \frac{G_{3,14}}{\lambda_1}, \quad (\text{E.68})$$

$$\frac{G_{3,14}}{\lambda_1} = \bar{\text{tr}} (\Sigma_1 G_{2,11} + \lambda_1 I_d)^{-2} (\Sigma_1 G_{2,15} + \lambda_1 \Theta) \Sigma_1 \quad (\text{E.69})$$

$$= \frac{G_{2,15}}{e_1^2} \bar{\text{df}}_2^{(1)}(\kappa_1) + \frac{\lambda_1}{e_1^2} \bar{\text{tr}} (\Sigma_1 + \kappa_1 I_d)^{-2} \Theta \Sigma_1, \quad (\text{E.70})$$

$$\implies G_{2,15} = \phi_1 G_{2,15} \bar{\text{df}}_2^{(1)}(\kappa_1) + \lambda_1 \phi_1 \bar{\text{tr}} (\Sigma_1 + \kappa_1 I_d)^{-2} \Theta \Sigma_1, \quad (\text{E.71})$$

$$\text{i.e., } G_{2,15} = \frac{\lambda_1 \phi_1}{1 - \phi_1 \bar{\text{df}}_2^{(1)}(\kappa_1)} \bar{\text{tr}} (\Sigma_1 + \kappa_1 I_d)^{-2} \Theta \Sigma_1. \quad (\text{E.72})$$

Hence:

$$G_{0,17} = \kappa_1^2 \bar{\text{tr}} (\Sigma_1 + \kappa_1 I_d)^{-2} \Theta \Sigma_1 + \kappa_1^2 \bar{\text{df}}_2^{(1)}(\kappa_1) \frac{G_{2,15}}{\lambda_1} \quad (\text{E.73})$$

$$= \kappa_1^2 \bar{\text{tr}} (\Sigma_1 + \kappa_1 I_d)^{-2} \Theta \Sigma_1 + \kappa_1^2 \frac{\phi_1 \bar{\text{df}}_2^{(1)}(\kappa_1)}{1 - \phi_1 \bar{\text{df}}_2^{(1)}(\kappa_1)} \bar{\text{tr}} (\Sigma_1 + \kappa_1 I_d)^{-2} \Theta \Sigma_1 \quad (\text{E.74})$$

$$= \left(1 + \frac{\phi_1 \bar{\text{df}}_2^{(1)}(\kappa_1)}{1 - \phi_1 \bar{\text{df}}_2^{(1)}(\kappa_1)} \right) \kappa_1^2 \bar{\text{tr}} (\Sigma_1 + \kappa_1 I_d)^{-2} \Theta \Sigma_1. \quad (\text{E.75})$$

In conclusion:

$$B_1(\widehat{f}_1) = \frac{\kappa_1^2 \bar{\text{tr}} (\Sigma_1 + \kappa_1 I_d)^{-2} \Theta \Sigma_1}{1 - \phi_1 \bar{\text{df}}_2^{(1)}(\kappa_1)}. \quad (\text{E.76})$$

Following similar steps for $B_2(\widehat{f}_2)$, we get:

$$B_2(\widehat{f}_2) = \frac{\kappa_2^2 \bar{\text{tr}} (\Sigma_2 + \kappa_2 I_d)^{-2} (\Theta + \Delta) \Sigma_2}{1 - \phi_2 \bar{\text{df}}_2^{(2)}(\kappa_2)}. \quad (\text{E.77})$$

We observe that in the unregularized case (i.e., $\lambda_s = 0$), $\kappa_s = 0$. In this setting, $B_s(\widehat{f}_s) = 0$ as expected. \square

E.4 Proof of Theorem E.2.1

Proof. We define $M = X^\top X + n\lambda I_d$. Note that one has:

$$\widehat{w} = M^{-1}(M_1 w_1^* + X_1^\top E_1 + M_2 w_2^* + X_2^\top E_2). \quad (\text{E.78})$$

We deduce that $R_s(\widehat{f}) = B_s(\widehat{f}) + V_s(\widehat{f})$, where:

$$B_s(\widehat{f}) = \mathbb{E} \|M^{-1} M_{s'} w_{s'}^* + M^{-1} M_s w_s^* - w_s^*\|_{\Sigma_s}^2, \quad (\text{E.79})$$

$$V_s(\widehat{f}) = \mathbb{E} \|M^{-1}(X_1^\top E_1 + X_2^\top E_2)\|_{\Sigma_s}^2 \quad (\text{E.80})$$

$$= \mathbb{E} \|M^{-1} X_1^\top E_1\|_{\Sigma_s}^2 + \mathbb{E} \|M^{-1} X_2^\top E_2\|_{\Sigma_s}^2, \quad (\text{E.81})$$

$$\text{with } s' = \begin{cases} 2, & s = 1 \\ 1, & s = 2 \end{cases}.$$

E.4.1 Variance Terms

Note that $V_s(\widehat{f})$ of the test error of \widehat{f} evaluated on group s is given by:

$$V_s(\widehat{f}) = \sigma_1^2 \mathbb{E} \text{tr } X_1 M^{-1} \Sigma_s M^{-1} X_1^\top + \sigma_2^2 \mathbb{E} \text{tr } X_2 M^{-1} \Sigma_s M^{-1} X_2^\top \quad (\text{E.82})$$

$$= \sigma_1^2 \mathbb{E} \text{tr } M^{-1} M_1 M^{-1} \Sigma_s + \sigma_2^2 \mathbb{E} \text{tr } M^{-1} M_2 M^{-1} \Sigma_s. \quad (\text{E.83})$$

Using OVFPT, we deduce that, in the limit given by Equation E.1, the following holds:

$$\mathbb{E} \bar{\text{tr}} (H_1 + H_2 + \lambda I_d)^{-1} H_2 (H_1 + H_2 + \lambda I_d)^{-1} \Sigma_s = \frac{G_{1,23}}{\lambda}, \quad (\text{E.87})$$

with:

$$\frac{G_{1,23}}{\lambda} = \lambda^{-1} \bar{\text{tr}} p_2 \Sigma_2 (\lambda \Sigma_s G_{0,15} + \lambda G_{0,27} I_d - p_1 \Sigma_1 G_{0,15} G_{5,24} + p_1 \Sigma_1 G_{0,27} G_{5,20}) \quad (\text{E.88})$$

$$\cdot (p_1 \Sigma_1 G_{5,20} + p_2 \Sigma_2 G_{0,15} + \lambda I_d)^{-2}. \quad (\text{E.89})$$

By identifying identical entries of \bar{Q}^{-1} , we must have that $\frac{G_{5,20}}{\lambda} = \frac{G_{6,21}}{\lambda} = \frac{G_{10,25}}{\lambda}, \frac{G_{0,15}}{\lambda} = \frac{G_{2,17}}{\lambda} = \frac{G_{13,28}}{\lambda}$. For $G_{6,21}$ and $G_{2,17}$, we observe that:

$$G_{6,21} = -\frac{\lambda}{-\lambda + \phi G_{7,20}}, \quad G_{7,20} = -\lambda \bar{\text{tr}} \Sigma_1 (p_1 \Sigma_1 G_{6,21} + p_2 \Sigma_2 G_{2,17} + \lambda I_d)^{-1} \quad (\text{E.90})$$

$$\implies G_{6,21} = \frac{1}{1 + \phi \bar{\text{tr}} \Sigma_1 (p_1 \Sigma_1 G_{6,21} + p_2 \Sigma_2 G_{2,17} + \lambda I_d)^{-1}}, \quad (\text{E.91})$$

$$G_{2,17} = -\frac{\lambda}{-\lambda + \phi G_{3,15}}, \quad G_{3,15} = -\lambda \bar{\text{tr}} \Sigma_2 (p_1 \Sigma_1 G_{6,21} + p_2 \Sigma_2 G_{2,17} + \lambda I_d)^{-1} \quad (\text{E.92})$$

$$\implies G_{2,17} = \frac{1}{1 + \phi \bar{\text{tr}} \Sigma_2 (p_1 \Sigma_1 G_{6,21} + p_2 \Sigma_2 G_{2,17} + \lambda I_d)^{-1}}. \quad (\text{E.93})$$

We define $\eta_1 = \frac{G_{6,21}}{\lambda}, \eta_2 = \frac{G_{2,17}}{\lambda}$, with $\eta_1 \geq 0, \eta_2 \geq 0$. Therefore:

$$\eta_s = \frac{1}{\lambda + \phi \bar{\text{tr}} \Sigma_s K^{-1}}, \quad (\text{E.94})$$

where $K = \eta_1 p_1 \Sigma_1 + \eta_2 p_2 \Sigma_2 + I_d$. Additionally, by identifying identical entries of \bar{Q}^{-1} , we

must have that $G_{5,24} = G_{6,25}$, $G_{0,27} = G_{2,28}$. We observe that:

$$G_{10,25} = \frac{-\lambda}{-\lambda + \phi G_{11,24}}, \quad (\text{E.95})$$

$$G_{6,25} = \frac{\lambda \phi G_{7,24}}{(-\lambda + \phi G_{7,20})(-\lambda + \phi G_{11,24})} = \phi \lambda^2 \eta_1^2 \frac{G_{7,24}}{\lambda}, \quad (\text{E.96})$$

$$\frac{G_{7,24}}{\lambda} = \lambda^{-2} \bar{\text{tr}} K^{-2} (p_1 \Sigma_1 G_{6,25} + p_2 \Sigma_2 G_{2,28} - \lambda \Sigma_s) \Sigma_1, \quad (\text{E.97})$$

$$\implies G_{6,25} = \phi \eta_1^2 \bar{\text{tr}} K^{-2} (p_1 \Sigma_1 G_{6,25} + p_2 \Sigma_2 G_{2,28} - \lambda \Sigma_s) \Sigma_1, \quad (\text{E.98})$$

$$G_{13,28} = \frac{-\lambda}{-\lambda + \phi G_{14,27}}, \quad (\text{E.99})$$

$$G_{2,28} = \frac{\lambda \phi G_{3,27}}{(-\lambda + \phi G_{3,15})(-\lambda + \phi G_{14,27})} = \phi \lambda^2 \eta_2^2 \frac{G_{3,27}}{\lambda}, \quad (\text{E.100})$$

$$\frac{G_{3,27}}{\lambda} = \lambda^{-2} \bar{\text{tr}} K^{-2} (p_1 \Sigma_1 G_{6,25} + p_2 \Sigma_2 G_{2,28} - \lambda \Sigma_s) \Sigma_2, \quad (\text{E.101})$$

$$\implies G_{2,28} = \phi \eta_2^2 \bar{\text{tr}} K^{-2} (p_1 \Sigma_1 G_{6,25} + p_2 \Sigma_2 G_{2,28} - \lambda \Sigma_s) \Sigma_2. \quad (\text{E.102})$$

We now define $v_1^{(s)} = -G_{6,25}$, $v_2^{(s)} = -G_{2,28}$, with $v_1^{(s)} \geq 0$, $v_2^{(s)} \geq 0$. Therefore, $v_1^{(s)}$, $v_2^{(s)}$ obey the following system of equations:

$$v_k^{(s)} = \phi \eta_k^2 \bar{\text{tr}} K^{-2} (v_1^{(s)} p_1 \Sigma_1 + v_2^{(s)} p_2 \Sigma_2 + \lambda \Sigma_s) \Sigma_k. \quad (\text{E.103})$$

We further define $u_k^{(s)} = \frac{v_k^{(s)}}{\lambda}$. Putting all the pieces together:

$$\frac{G_{1,23}}{\lambda} = \lambda^{-1} \bar{\text{tr}} p_2 \Sigma_2 (\eta_2 \Sigma_s - u_2^{(s)} I_d + p_1 \Sigma_1 (\eta_2 u_1^{(s)} - \eta_1 u_2^{(s)})) K^{-2}. \quad (\text{E.104})$$

By symmetry, in conclusion:

$$V_s(\widehat{f}) = V_s^{(1)}(\widehat{f}) + V_s^{(2)}(\widehat{f}), \quad (\text{E.105})$$

$$V_s^{(k)}(\widehat{f}) = \lambda^{-1} \phi \sigma_k^2 \bar{\text{tr}} p_k \Sigma_k (\eta_k \Sigma_s - u_k^{(s)} I_d + p_{k'} \Sigma_{k'} (\eta_k u_{k'}^{(s)} - \eta_{k'} u_k^{(s)})) K^{-2}, \quad (\text{E.106})$$

$$\text{with } k' = \begin{cases} 2, & k = 1 \\ 1, & k = 2 \end{cases}.$$

We now corroborate our result in the limit $p_2 \rightarrow 1$ (i.e., $p_1 \rightarrow 0$) and $s = 2$. We observe that:

$$\phi \rightarrow \phi_2, \lambda \rightarrow \lambda_2, \quad (\text{E.107})$$

$$V_2^{(1)}(\widehat{f}) = 0, \quad (\text{E.108})$$

$$\frac{V_2^{(2)}(\widehat{f})}{\lambda^{-1}\phi_2\sigma_2^2} = \bar{\text{tr}} \Sigma_2(\eta_2\Sigma_2 - u_2^{(2)}I_d)K^{-2} \quad (\text{E.109})$$

$$v_2^{(2)} = \phi_2\eta_2^2\bar{\text{tr}} K^{-2}(v_2^{(2)}\Sigma_2 + \lambda_2\Sigma_2)\Sigma_2 \quad (\text{E.110})$$

$$= \phi_2(v_2^{(2)} + \lambda_2)\bar{\text{d}}f_2^{(2)}(\kappa_2), \quad (\text{E.111})$$

$$u_2^{(2)} = \frac{\phi_2\bar{\text{d}}f_2^{(2)}(\kappa_2)}{1 - \phi_2\bar{\text{d}}f_2^{(2)}(\kappa_2)}, \quad (\text{E.112})$$

$$\frac{V_2^{(2)}(\widehat{f})}{\lambda^{-1}\phi_2\sigma_2^2} = \kappa_2\bar{\text{d}}f_2^{(2)}(\kappa_2) - u_2^{(2)}\bar{\text{tr}} \Sigma_2(\eta_2\Sigma_2 + I_d)^{-2} \quad (\text{E.113})$$

$$= \kappa_2\bar{\text{d}}f_2^{(2)}(\kappa_2) - \kappa_2^2u_2^{(2)}\bar{\text{tr}} \Sigma_2(\Sigma_2 + \kappa_2I_d)^{-2} \quad (\text{E.114})$$

$$= \kappa_2\bar{\text{d}}f_2^{(2)}(\kappa_2) - \kappa_2u_2^{(2)}(\bar{\text{d}}f_1^{(2)}(\kappa_2) - \bar{\text{d}}f_2^{(2)}(\kappa_2)) \quad (\text{E.115})$$

$$= \kappa_2(1 + u_2^{(2)})\bar{\text{d}}f_2^{(2)}(\kappa_2) - \kappa_2u_2^{(2)}\bar{\text{d}}f_1^{(2)}(\kappa_2) \quad (\text{E.116})$$

$$= \frac{\kappa_2 - \kappa_2\phi_2\bar{\text{d}}f_1^{(2)}(\kappa_2)}{1 - \phi_2\bar{\text{d}}f_2^{(2)}(\kappa_2)} \cdot \bar{\text{d}}f_2^{(2)}(\kappa_2) \quad (\text{E.117})$$

$$= \frac{\lambda\bar{\text{d}}f_2^{(2)}(\kappa_2)}{1 - \phi_2\bar{\text{d}}f_2^{(2)}(\kappa_2)}, \quad (\text{E.118})$$

$$V_2^{(2)}(\widehat{f}) = \frac{\sigma_2^2\phi_2\bar{\text{d}}f_2^{(2)}(\kappa_2)}{1 - \phi_2\bar{\text{d}}f_2^{(2)}(\kappa_2)}, \quad (\text{E.119})$$

which exactly recovers the result for $V_2(\widehat{f}_2)$ as expected.

E.4.2 Bias Terms

Recall that:

$$B_s(\widehat{f}) = \mathbb{E} \|M^{-1}M_{s'}w_{s'}^* + M^{-1}M_s w_s^* - w_s^*\|_{\Sigma_s}^2. \quad (\text{E.120})$$

Now, observe that $M^{-1}M_1w_1^* - w_1^* = M^{-1}M_1w_1^* - M^{-1}Mw_1^* = -M^{-1}M_2w_1^* - n\lambda M^{-1}w_1^*$.

Let $\delta = w_2^* - w_1^*$. Then:

$$B_s(\widehat{f}) = \mathbb{E}\|M^{-1}M_{s'}(-1)^{s-1}\delta - n\lambda M^{-1}w_s^*\|_{\Sigma_s}^2 \quad (\text{E.121})$$

$$= \mathbb{E} \operatorname{tr} \delta^\top M_{s'} M^{-1} \Sigma_s M^{-1} M_{s'} \delta \quad (\text{E.122})$$

$$- 2(-1)^{s-1} n \lambda \mathbb{E} \operatorname{tr} \delta^\top M_{s'} M^{-1} \Sigma_s M^{-1} w_s^* \quad (\text{E.123})$$

$$+ n^2 \lambda^2 \mathbb{E} \operatorname{tr} (w_s^*)^\top M^{-1} \Sigma_s M^{-1} w_s^* \quad (\text{E.124})$$

$$= B_s^{(1)}(\widehat{f}) - 2(-1)^{s-1} B_s^{(2)}(\widehat{f}) + B_s^{(3)}(\widehat{f}), \quad (\text{E.125})$$

where:

$$B_s^{(1)}(\widehat{f}) = \mathbb{E} \bar{\operatorname{tr}} (H_1/\lambda + H_2/\lambda + I_d)^{-1} (H_{s'}/\lambda) \Delta (H_{s'}/\lambda) (H_1/\lambda + H_2/\lambda + I_d)^{-1} \Sigma_s, \quad (\text{E.126})$$

$$B_s^{(2)}(\widehat{f}) = \mathbb{E} \operatorname{tr} \delta^\top (H_{s'}/\lambda) (H_1/\lambda + H_2/\lambda + I_d)^{-1} \Sigma_s (H_1/\lambda + H_2/\lambda + I_d)^{-1} w_s^*, \quad (\text{E.127})$$

$$B_s^{(3)}(\widehat{f}) = \mathbb{E} \bar{\operatorname{tr}} (H_1/\lambda + H_2/\lambda + I_d)^{-1} \Theta_s (H_1/\lambda + H_2/\lambda + I_d)^{-1} \Sigma_s. \quad (\text{E.128})$$

Because δ and w_1^* are independent and sampled from zero-centered distributions:

$$B_1^{(2)}(\widehat{f}) = 0, \quad (\text{E.129})$$

$$B_2^{(2)}(\widehat{f}) = \mathbb{E} \bar{\operatorname{tr}} (H_1/\lambda + H_2/\lambda + I_d)^{-1} \Delta (H_1/\lambda) (H_1/\lambda + H_2/\lambda + I_d)^{-1} \Sigma_2. \quad (\text{E.130})$$

WLOG, for $B_s^{(1)}$, we focus on the case $s = 1$. The matrix of interest has a linear pencil representation given by (with zero-based indexing):

$$(H_1/\lambda + H_2/\lambda + I_d)^{-1} (H_2/\lambda) \Delta (H_2/\lambda) (H_1/\lambda + H_2/\lambda + I_d)^{-1} \Sigma_1 = Q_{1,16}^{-1}, \quad (\text{E.131})$$

By identifying identical entries of \bar{Q}^{-1} , we must have that $\eta_1 = \frac{G_{6,23}}{\lambda} = \frac{G_{7,24}}{\lambda} = \frac{G_{11,28}}{\lambda}$, $\eta_2 = \frac{G_{2,19}}{\lambda} = \frac{G_{3,20}}{\lambda} = \frac{G_{14,31}}{\lambda}$. For $G_{7,24}$ and $G_{3,20}$, we observe that:

$$G_{7,24} = -\frac{\lambda}{-\lambda + \phi G_{8,23}}, \quad G_{8,23} = -\lambda \bar{\text{tr}} \Sigma_1 (p_1 \Sigma_1 G_{7,24} + p_2 \Sigma_2 G_{3,20} + \lambda I_d)^{-1} \quad (\text{E.135})$$

$$\implies G_{7,24} = \frac{1}{1 + \phi \bar{\text{tr}} \Sigma_1 (p_1 \Sigma_1 G_{7,24} + p_2 \Sigma_2 G_{3,20} + \lambda I_d)^{-1}}, \quad (\text{E.136})$$

$$G_{3,20} = -\frac{\lambda}{-\lambda + \phi G_{4,19}}, \quad G_{4,19} = -\lambda \bar{\text{tr}} \Sigma_2 (p_1 \Sigma_1 G_{7,24} + p_2 \Sigma_2 G_{3,20} + \lambda I_d)^{-1} \quad (\text{E.137})$$

$$\implies G_{3,20} = \frac{1}{1 + \phi \bar{\text{tr}} \Sigma_2 (p_1 \Sigma_1 G_{7,24} + p_2 \Sigma_2 G_{3,20} + \lambda I_d)^{-1}}. \quad (\text{E.138})$$

By again identifying identical entries of \bar{Q}^{-1} , we further have that $v_1^{(1)} = -G_{6,27} = -G_{7,28}$, $v_2^{(1)} = -G_{2,30} = -G_{3,31}$. We observe that:

$$G_{7,28} = \phi \lambda^2 \eta_1^2 \frac{G_{8,27}}{\lambda}, \quad (\text{E.139})$$

$$\frac{G_{8,27}}{\lambda} = \lambda^{-2} \bar{\text{tr}} K^{-2} (p_1 \Sigma_1 G_{7,28} + p_2 \Sigma_2 G_{3,31} - \lambda \Sigma_1) \Sigma_1 \quad (\text{E.140})$$

$$\implies v_1^{(1)} = \phi \eta_1^2 \bar{\text{tr}} K^{-2} (v_1^{(s)} p_1 \Sigma_1 + v_2^{(s)} p_2 \Sigma_2 + \lambda \Sigma_1) \Sigma_1, \quad (\text{E.141})$$

$$G_{3,31} = \phi \lambda^2 \eta_2^2 \frac{G_{4,30}}{\lambda}, \quad (\text{E.142})$$

$$\frac{G_{4,30}}{\lambda} = \lambda^{-2} \bar{\text{tr}} K^{-2} (p_1 \Sigma_1 G_{7,28} + p_2 \Sigma_2 G_{3,31} - \lambda \Sigma_1) \Sigma_2, \quad (\text{E.143})$$

$$\implies v_2^{(1)} = \phi \eta_2^2 \bar{\text{tr}} K^{-2} (v_1^{(s)} p_1 \Sigma_1 + v_2^{(s)} p_2 \Sigma_2 + \lambda \Sigma_1) \Sigma_2. \quad (\text{E.144})$$

Putting all the pieces together:

$$B_1^{(1)}(\hat{f}) = \bar{\text{tr}} p_2 \Sigma_2 \Delta (p_2 \eta_2^2 \Sigma_2 (1 + p_1 u_1^{(s)}) \Sigma_1 + u_2^{(s)} (p_1 \eta_1 \Sigma_1 + I_d)^2) K^{-2}. \quad (\text{E.145})$$

In conclusion:

$$B_s^{(1)}(\hat{f}) = \bar{\text{tr}} p_{s'} \Sigma_{s'} \Delta (p_{s'} \eta_{s'}^2 \Sigma_{s'} (1 + p_s u_s^{(s)}) \Sigma_s + u_{s'}^{(s)} (p_s \eta_s \Sigma_s + I_d)^2) K^{-2}. \quad (\text{E.146})$$

Now, switching our focus to $B_2^{(2)}(\hat{f})$, the matrix of interest has a linear pencil representation given by (with zero-based indexing):

$$(H_1/\lambda + H_2/\lambda + I_d)^{-1} \Delta (H_1/\lambda) (H_1/\lambda + H_2/\lambda + I_d)^{-1} \Sigma_2 = Q_{0,15}^{-1}, \quad (\text{E.147})$$

By identifying identical entries of \bar{Q}^{-1} , we must have that $\eta_1 = \frac{G_{2,18}}{\lambda} = \frac{G_{3,19}}{\lambda} = \frac{G_{11,27}}{\lambda}$, $\eta_2 = \frac{G_{6,22}}{\lambda} = \frac{G_{7,23}}{\lambda} = \frac{G_{14,30}}{\lambda}$. For $G_{3,19}$ and $G_{7,23}$, we observe that:

$$G_{3,19} = -\frac{\lambda}{-\lambda + \phi G_{4,18}}, \quad G_{4,18} = -\lambda \bar{\text{tr}} \Sigma_1 (p_1 \Sigma_1 G_{3,19} + p_2 \Sigma_2 G_{7,23} + \lambda I_d)^{-1} \quad (\text{E.151})$$

$$\implies G_{3,19} = \frac{1}{1 + \phi \bar{\text{tr}} \Sigma_1 (p_1 \Sigma_1 G_{3,19} + p_2 \Sigma_2 G_{7,23} + \lambda I_d)^{-1}}, \quad (\text{E.152})$$

$$G_{7,23} = -\frac{\lambda}{-\lambda + \phi G_{8,22}}, \quad G_{8,22} = -\lambda \bar{\text{tr}} \Sigma_2 (p_1 \Sigma_1 G_{3,19} + p_2 \Sigma_2 G_{7,23} + \lambda I_d)^{-1} \quad (\text{E.153})$$

$$\implies G_{7,23} = \frac{1}{1 + \phi \bar{\text{tr}} \Sigma_2 (p_1 \Sigma_1 G_{3,19} + p_2 \Sigma_2 G_{7,23} + \lambda I_d)^{-1}}. \quad (\text{E.154})$$

By again identifying identical entries of \bar{Q}^{-1} , we further have that $v_1^{(2)} = -G_{2,26} = -G_{3,27}$, $v_2^{(2)} = -G_{6,29} = -G_{7,30}$. We observe that:

$$G_{3,27} = \phi \lambda^2 \eta_1^2 \frac{G_{4,26}}{\lambda}, \quad (\text{E.155})$$

$$\frac{G_{4,26}}{\lambda} = \lambda^{-2} \bar{\text{tr}} K^{-2} (p_1 \Sigma_1 G_{3,27} + p_2 \Sigma_2 G_{7,30} - \lambda \Sigma_2) \Sigma_1, \quad (\text{E.156})$$

$$\implies v_1^{(2)} = \phi \eta_1^2 \bar{\text{tr}} K^{-2} (v_1^{(2)} p_1 \Sigma_1 + v_2^{(2)} p_2 \Sigma_2 + \lambda \Sigma_2) \Sigma_1, \quad (\text{E.157})$$

$$G_{7,30} = \phi \lambda^2 \eta_2^2 \frac{G_{8,29}}{\lambda}, \quad (\text{E.158})$$

$$\frac{G_{8,29}}{\lambda} = \lambda^{-2} \bar{\text{tr}} K^{-2} (p_1 \Sigma_1 G_{3,27} + p_2 \Sigma_2 G_{7,30} - \lambda \Sigma_2) \Sigma_2, \quad (\text{E.159})$$

$$\implies v_2^{(2)} = \phi \eta_2^2 \bar{\text{tr}} K^{-2} (v_1^{(2)} p_1 \Sigma_1 + v_2^{(2)} p_2 \Sigma_2 + \lambda \Sigma_2) \Sigma_2. \quad (\text{E.160})$$

Putting all the pieces together:

$$B_2^{(1)}(\hat{f}) = 0, \quad (\text{E.161})$$

$$B_2^{(2)}(\hat{f}) = \bar{\text{tr}} p_1 \Sigma_1 \Delta (\eta_1 \Sigma_2 - u_1^{(2)} I_d + p_2 \Sigma_2 (\eta_1 u_2^{(2)} - \eta_2 u_1^{(2)})) K^{-2}. \quad (\text{E.162})$$

Finally, switching our focus to $B_1^{(3)}(\hat{f})$, the matrix of interest has a linear pencil representation given by (with zero-based indexing):

$$(H_1/\lambda + H_2/\lambda + I_d)^{-1} \Theta (H_1/\lambda + H_2/\lambda + I_d)^{-1} \Sigma_1 = Q_{1,8}^{-1}, \quad (\text{E.163})$$

conclusion:

$$B_s^{(3)}(\widehat{f}) = \bar{\text{tr}} \Theta_s(p_1 u_1^{(s)} \Sigma_1 + p_2 u_2^{(s)} \Sigma_2 + \Sigma_s) K^{-2}. \quad (\text{E.166})$$

In the limit $p_s \rightarrow 1$ (i.e., $p_{s'} \rightarrow 0$), we observe that:

$$\phi \rightarrow \phi_s, \lambda \rightarrow \lambda_s, \quad (\text{E.167})$$

$$B_s^{(1)}(\widehat{f}) \rightarrow 0, \quad (\text{E.168})$$

$$B_s^{(2)}(\widehat{f}) \rightarrow 0, \quad (\text{E.169})$$

$$B_s^{(3)}(\widehat{f}) = \bar{\text{tr}} \Theta_s(u_s^{(s)} + 1) \Sigma_s K^{-2}, \quad (\text{E.170})$$

$$v_s^{(s)} = \phi_s \eta_s^2 \bar{\text{tr}} K^{-2} (v_s^{(s)} + \lambda_s) \Sigma_s^2 \quad (\text{E.171})$$

$$= \phi_s (v_s^{(s)} + \lambda_s) \bar{\text{df}}_2^{(s)}(\kappa_s) \quad (\text{E.172})$$

$$u_s^{(s)} = \frac{\phi_s \bar{\text{df}}_2^{(s)}(\kappa_s)}{1 - \phi_s \bar{\text{df}}_2^{(s)}(\kappa_s)}, \quad (\text{E.173})$$

$$B_s^{(3)}(\widehat{f}) = \frac{\kappa_s^2 \bar{\text{tr}} \Theta_s \Sigma_s (\Sigma_s + \kappa_s I_d)^{-2}}{1 - \phi_s \bar{\text{df}}_2^{(s)}(\kappa_s)}, \quad (\text{E.174})$$

$$B_s(\widehat{f}) \rightarrow B_s^{(3)}(\widehat{f}), \quad (\text{E.175})$$

which matches up exactly with $B_s(\widehat{f}_s)$ as expected. \square

E.5 Proof of Theorem 6.3.1

Proof. The gradient of the loss L is given by:

$$\begin{aligned} \nabla L(\eta) &= \sum_s S^\top X_s^\top (X_s S \eta - Y_s) / n + \lambda \eta = \sum_s S^\top M_s S \eta - \sum_s S^\top X_s^\top Y_s / n + \lambda \eta \\ &= H \eta - \sum_s S^\top X_s^\top Y_s / n, \end{aligned}$$

where $H = S^\top MS + \lambda I_m \in \mathbb{R}^{m \times m}$, with $M = M_1 + M_2$ and $M_s = X_s^\top X_s/n$. Thus, setting $R = H^{-1}$, we may write:

$$\begin{aligned}\widehat{w} &= S\widehat{\eta} = SRS^\top(X_1^\top Y_1 + X_2^\top Y_2)/n \\ &= SRS^\top(M_1 w_1^* + M_2 w_2^*) + SRS^\top X_1^\top E_1/n + SRS^\top X_2^\top E_2/n.\end{aligned}$$

We deduce the following bias-variance decomposition:

$$\begin{aligned}\mathbb{E}\|\widehat{w} - w_s^*\|_{\Sigma_s}^2 &= B_s(\widehat{f}) + V_s(\widehat{f}), \text{ where} \\ V_s(\widehat{f}) &= V_s^{(1)}(\widehat{f}) + V_s^{(2)}(\widehat{f}), \text{ with } V_s^{(j)}(\widehat{f}) = \sigma_j^2 \phi \mathbb{E} \bar{\text{tr}} M_j SRS^\top \Sigma_s SRS^\top, \\ B_s(\widehat{f}) &= \mathbb{E}\|SRS^\top(M_1 w_1^* + M_2 w_2^*) - w_s^*\|_{\Sigma_s}^2.\end{aligned}$$

We can further decompose $B_s(\widehat{f})$, first considering the case $s = 1$. We define $\delta = w_2^* - w_1^*$.

$$\begin{aligned}\mathbb{E}\|SRS^\top(M_1 w_1^* + M_2 w_2^*) - w_1^*\|_{\Sigma_1}^2 &= \mathbb{E}\|(SRS^\top(M_1 + M_2) - I_d)w_1^* + SRS^\top M_2 \delta\|_{\Sigma_1}^2 \\ &= \mathbb{E}\|(SRS^\top M - I_d)w_1^*\|_{\Sigma_1}^2 + \mathbb{E}\|SRS^\top M_2 \delta\|_{\Sigma_1}^2 \\ &= \mathbb{E} \bar{\text{tr}} \Theta(MSRS^\top - I_d)\Sigma_1(SRS^\top M - I_d) + \mathbb{E} \bar{\text{tr}} \Delta M_2 SRS^\top \Sigma_1 SRS^\top M_2 \\ &= \mathbb{E} \bar{\text{tr}} \Theta \Sigma_1 + \mathbb{E} \bar{\text{tr}} \Theta MSRS^\top \Sigma_1 SRS^\top M - 2\mathbb{E} \bar{\text{tr}} \Theta \Sigma_1 SRS^\top M + \mathbb{E} \bar{\text{tr}} \Delta M_2 SRS^\top \Sigma_1 SRS^\top M_2.\end{aligned}$$

We can similarly decompose B_2 :

$$\begin{aligned}\mathbb{E}\|SRS^\top(M_1 w_1^* + M_2 w_2^*) - w_2^*\|_{\Sigma_2}^2 &= \mathbb{E}\|SRS^\top(M_1 w_1^* + M_2 w_2^*) - w_2^*\|_{\Sigma_2}^2 \\ &= \mathbb{E}\|(SRS^\top(M_1 + M_2) - I_d)w_2^* - SRS^\top M_1 \delta\|_{\Sigma_2}^2 \\ &= \mathbb{E}\|(SRS^\top M - I_d)w_2^*\|_{\Sigma_2}^2 + \mathbb{E}\|SRS^\top M_1 \delta\|_{\Sigma_2}^2 - 2\mathbb{E} \text{tr} (w_2^*)^\top (MSRS^\top - I_d)\Sigma_2 SRS^\top M_1 \delta \\ &= \mathbb{E} \bar{\text{tr}} \Theta_2(MSRS^\top - I_d)\Sigma_2(SRS^\top M - I_d) + \mathbb{E} \bar{\text{tr}} \Delta M_1 SRS^\top \Sigma_2 SRS^\top M_1 \\ &\quad - 2\mathbb{E} \bar{\text{tr}} \Delta(MSRS^\top - I_d)\Sigma_2 SRS^\top M_1 \\ &= \mathbb{E} \bar{\text{tr}} \Theta_2 \Sigma_2 + \mathbb{E} \bar{\text{tr}} \Theta_2 MSRS^\top \Sigma_2 SRS^\top M - 2\mathbb{E} \bar{\text{tr}} \Theta_2 \Sigma_2 SRS^\top M \\ &\quad + \mathbb{E} \bar{\text{tr}} \Delta M_1 SRS^\top \Sigma_2 SRS^\top M_1 - 2\mathbb{E} \bar{\text{tr}} \Delta MSRS^\top \Sigma_2 SRS^\top M_1 + 2\mathbb{E} \bar{\text{tr}} \Delta \Sigma_2 SRS^\top M_1.\end{aligned}$$

Furthermore, we observe that:

$$\mathbb{E}\bar{\text{tr}} AMSRS^\top BSRST^\top M \quad (\text{E.176})$$

$$= \mathbb{E}\bar{\text{tr}} AM_1SRS^\top BSRST^\top M_1 + \mathbb{E}\bar{\text{tr}} AM_2SRS^\top BSRST^\top M_2 + 2\mathbb{E}\bar{\text{tr}} AM_1SRS^\top BSRST^\top M_2, \quad (\text{E.177})$$

$$\mathbb{E}\bar{\text{tr}} ASRS^\top M = \mathbb{E}\bar{\text{tr}} ASRS^\top M_1 + \mathbb{E}\bar{\text{tr}} ASRS^\top M_2. \quad (\text{E.178})$$

Hence, we desire deterministic equivalents for the following expressions:

$$r_j^{(1)}(A) = AS\bar{R}S^\top \bar{M}_j, \quad (\text{E.179})$$

$$r_j^{(2)}(A, B) = A\bar{M}_jS\bar{R}S^\top B\bar{R}S^\top, \quad (\text{E.180})$$

$$r_j^{(3)}(A, B) = A\bar{M}_jS\bar{R}S^\top B\bar{R}S^\top \bar{M}_j, \quad (\text{E.181})$$

$$r_j^{(4)}(A, B) = A\bar{M}_jS\bar{R}S^\top B\bar{R}S^\top \bar{M}_{j'}, \quad (\text{E.182})$$

where:

$$\bar{M}_j = \Sigma_j^{1/2} Z_j^\top Z_j \Sigma_j^{1/2}, \bar{R} = (S^\top \bar{M} S + I_m)^{-1}, \bar{M} = \bar{M}_1 + \bar{M}_2, \quad (\text{E.183})$$

$$\bar{M}_j = M_j/\lambda, \bar{R} = \lambda R, \bar{M} = M/\lambda. \quad (\text{E.184})$$

In summary:

$$V_s^{(j)}(\hat{f}) = \sigma_j^2 \phi \lambda^{-1} \mathbb{E}\bar{\text{tr}} r_j^{(2)}(I_d, \Sigma_s), \quad (\text{E.185})$$

$$B_s(\hat{f}) = \bar{\text{tr}} \Theta_s \Sigma_s \quad (\text{E.186})$$

$$+ \mathbb{E}\bar{\text{tr}} r_1^{(3)}(\Theta_s, \Sigma_s) + \mathbb{E}\bar{\text{tr}} r_2^{(3)}(\Theta_s, \Sigma_s) + 2\mathbb{E}\bar{\text{tr}} r_1^{(4)}(\Theta_s, \Sigma_s) \quad (\text{E.187})$$

$$- 2\mathbb{E}\bar{\text{tr}} r_1^{(1)}(\Theta_s \Sigma_s) - 2\mathbb{E}\bar{\text{tr}} r_2^{(1)}(\Theta_s \Sigma_s) \quad (\text{E.188})$$

$$+ \mathbb{E}\bar{\text{tr}} r_{s'}^{(3)}(\Delta, \Sigma_s) \quad (\text{E.189})$$

$$- 2 \begin{cases} 0, & s = 1, \\ \mathbb{E}\bar{\text{tr}} r_1^{(3)}(\Delta, \Sigma_2) + \mathbb{E}\bar{\text{tr}} r_2^{(4)}(\Delta, \Sigma_2) - \mathbb{E}\bar{\text{tr}} r_1^{(1)}(\Delta \Sigma_2), & s = 2 \end{cases}. \quad (\text{E.190})$$

E.5.1 Computing $\mathbb{E}\bar{\text{tr}} r_j^{(1)}$

WLOG, we focus on $r_1^{(1)}$. The matrix of interest has a linear pencil representation given by (with zero-based indexing):

$$r_1^{(1)} = Q_{1,10}^{-1}, \quad (\text{E.191})$$

where the linear pencil Q is defined as follows:

$$Q = \begin{pmatrix} I_d & 0 & -S & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ -A & I_d & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & I_m & S^\top & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & I_d & -\Sigma_1^{\frac{1}{2}} & 0 & 0 & -\Sigma_2^{\frac{1}{2}} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & I_d & -\frac{1}{\sqrt{\lambda}}Z_1^\top & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & I_{n_1} & -\frac{1}{\sqrt{\lambda}}Z_1 & 0 & 0 & 0 & 0 \\ -\Sigma_1^{\frac{1}{2}} & 0 & 0 & 0 & 0 & 0 & I_d & 0 & 0 & 0 & \Sigma_1^{\frac{1}{2}} \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & I_d & -\frac{1}{\sqrt{\lambda}}Z_2^\top & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & I_{n_2} & -\frac{1}{\sqrt{\lambda}}Z_2 & 0 \\ -\Sigma_2^{\frac{1}{2}} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & I_d & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & I_d \end{pmatrix}. \quad (\text{E.192})$$

Using the tools of OVFPT, the following holds:

$$\mathbb{E}\bar{\text{tr}} r_1^{(1)} = G_{1,21}, \quad (\text{E.193})$$

with:

$$G_{1,21} = \bar{\text{tr}} \gamma p_1 \Sigma_1 A G_{2,13} G_{5,16} (\gamma G_{2,13} (p_1 \Sigma_1 G_{5,16} + p_2 G_{8,19}) + \lambda I_d)^{-1}. \quad (\text{E.194})$$

For $G_{5,16}$ and $G_{8,19}$, we observe that:

$$G_{5,16} = \frac{-\lambda}{-\lambda + \phi G_{6,15}}, \quad G_{6,15} = -\lambda \gamma G_{2,13} \bar{\text{tr}} \Sigma_1 (\gamma G_{2,13} (p_1 \Sigma_1 G_{5,16} + p_2 \Sigma_2 G_{8,19}) + \lambda I_d)^{-1}, \quad (\text{E.195})$$

$$\implies G_{5,16} = \frac{1}{1 + \psi G_{2,13} \bar{\text{tr}} \Sigma_1 (\gamma G_{2,13} (p_1 \Sigma_1 G_{5,16} + p_2 \Sigma_2 G_{8,19}) + \lambda I_d)^{-1}}, \quad (\text{E.196})$$

$$G_{8,19} = \frac{-\lambda}{-\lambda + \phi G_{9,18}}, \quad G_{9,18} = -\lambda \gamma G_{2,13} \bar{\text{tr}} \Sigma_2 (\gamma G_{2,13} (p_1 \Sigma_1 G_{5,16} + p_2 \Sigma_2 G_{8,19}) + \lambda I_d)^{-1}, \quad (\text{E.197})$$

$$G_{8,19} = \frac{1}{1 + \psi G_{2,13} \bar{\text{tr}} \Sigma_2 (\gamma G_{2,13} (p_1 \Sigma_1 G_{5,16} + p_2 \Sigma_2 G_{8,19}) + \lambda I_d)^{-1}}. \quad (\text{E.198})$$

We define $e_1 = G_{5,16}$, $e_2 = G_{8,19}$, with $e_1 \geq 0$, $e_2 \geq 0$. We further observe that:

$$G_{2,13} = \frac{1}{1 + G_{3,11}}, \quad (\text{E.199})$$

$$G_{3,11} = \bar{\text{tr}} (p_1 \Sigma_1 G_{5,16} + p_2 \Sigma_2 G_{8,19}) (\gamma G_{2,13} (p_1 \Sigma_1 G_{5,16} + p_2 \Sigma_2 G_{8,19}) + \lambda I_d)^{-1}. \quad (\text{E.200})$$

We define $\tau = G_{2,13} \geq 0$. We further define $L = p_1 e_1 \Sigma_1 + p_2 e_2 \Sigma_2$, $K = \gamma \tau L + \lambda I_d$. Therefore, we have the following system of equations:

$$e_s = \frac{1}{1 + \psi \tau \bar{\text{tr}} \Sigma_s K^{-1}}, \quad \tau = \frac{1}{1 + \bar{\text{tr}} L K^{-1}}. \quad (\text{E.201})$$

In conclusion:

$$\mathbb{E} \bar{\text{tr}} r_j^{(1)} = p_j \gamma e_j \tau \bar{\text{tr}} A \Sigma_j K^{-1}. \quad (\text{E.202})$$

E.5.2 Computing $\mathbb{E} \bar{\text{tr}} r_j^{(2)}$

WLOG, we focus on $r_1^{(2)}$. The matrix of interest has a linear pencil representation given by (with zero-based indexing):

$$r_1^{(2)} = -Q_{1,13}^{-1}, \quad (\text{E.203})$$

with:

$$G_{1,33} = -p_1 \bar{\text{tr}} A \Sigma_1 P_1 P_2^{-1}, \quad (\text{E.206})$$

$$P_1 = \gamma \lambda B G_{3,23} G_{6,26} G_{12,32} - \gamma p_2 \Sigma_2 G_{3,23} G_{6,26} G_{9,38} G_{12,32} \quad (\text{E.207})$$

$$+ \gamma p_2 \Sigma_2 G_{3,35} G_{6,26} G_{9,29} G_{12,32} + \lambda G_{3,23} G_{6,32} I_d + \lambda G_{3,15} G_{12,32} I_d, \quad (\text{E.208})$$

$$P_2 = (\gamma G_{6,26} (p_1 \Sigma_1 G_{3,23} + \gamma p_2 \Sigma_2 G_{9,29}) + \lambda I_d) \quad (\text{E.209})$$

$$\cdot (\gamma G_{12,32} (p_1 \Sigma_1 G_{15,35} + p_2 \Sigma_2 G_{18,38}) + \lambda I_d). \quad (\text{E.210})$$

Following similar steps as before and recognizing identifications, we arrive at that:

$$e_1 = G_{3,23} = G_{15,35}, \quad (\text{E.211})$$

$$e_2 = G_{9,29} = G_{18,38}, \quad (\text{E.212})$$

$$\tau = G_{6,26} = G_{12,32}. \quad (\text{E.213})$$

We now focus on the remaining terms. We observe that:

$$G_{3,35} = \phi e_1^2 \frac{G_{4,14}}{\lambda} \quad (\text{E.214})$$

$$\frac{G_{4,14}}{\lambda} = \gamma \bar{\text{tr}} \Sigma_1 (\gamma \tau^2 (p_1 \Sigma_1 G_{3,35} + p_2 \Sigma_2 G_{9,38} - \lambda B) - \lambda G_{6,32} I_d) K^{-2}, \quad (\text{E.215})$$

$$G_{9,38} = \phi e_2^2 \frac{G_{10,37}}{\lambda}, \quad (\text{E.216})$$

$$\frac{G_{10,37}}{\lambda} = \gamma \bar{\text{tr}} \Sigma_2 (\gamma \tau^2 (p_1 \Sigma_1 G_{3,35} + p_2 \Sigma_2 G_{9,38} - \lambda B) - \lambda G_{6,32} I_d) K^{-2}. \quad (\text{E.217})$$

We define $u_1 = -\frac{G_{3,35}}{\lambda}$, $u_2 = -\frac{G_{9,38}}{\lambda}$, with $u_1 \leq 0$, $u_2 \leq 0$. We further define $D = p_1 u_1 \Sigma_1 + p_2 u_2 \Sigma_2 + B$. We now observe that:

$$G_{6,32} = -\frac{G_{7,31}}{(G_{7,25} + 1)(G_{13,31} + 1)} = -\tau^2 G_{7,31}, \quad (\text{E.218})$$

$$G_{7,31} = -\bar{\text{tr}} (\gamma G_{6,32} L^2 + \lambda^2 D) K^{-2}. \quad (\text{E.219})$$

Defining $\rho = G_{6,32}$, we must have the following system of equations:

$$u_s = \psi e_s^2 \bar{\text{tr}} \Sigma_s (\gamma \tau^2 D + \rho I_d) K^{-2}, \quad (\text{E.220})$$

$$\rho = \tau^2 \bar{\text{tr}} (\gamma \rho L^2 + \lambda^2 D) K^{-2}. \quad (\text{E.221})$$

It holds that $\mathbb{E}\bar{\text{tr}} r_1^{(3)} = G_{1,41}$. We immediately observe that:

$$e_1 = G_{3,24}, G_{15,36}, \quad (\text{E.227})$$

$$e_2 = G_{9,30}, G_{18,39}, \quad (\text{E.228})$$

$$\tau = G_{6,27}, G_{12,33}, \quad (\text{E.229})$$

$$u_1 = -\frac{G_{3,36}}{\lambda}, \quad (\text{E.230})$$

$$u_2 = -\frac{G_{9,39}}{\lambda}, \quad (\text{E.231})$$

$$\rho = G_{6,33}. \quad (\text{E.232})$$

In conclusion:

$$\mathbb{E}\bar{\text{tr}} r_j^{(3)} = p_j \bar{\text{tr}} A \Sigma_j (\gamma e_j^2 p_j \Sigma_j (\gamma \tau^2 u_j p_j \Sigma_j + \gamma \tau^2 B + \rho I_d) + u_j (\gamma e_{j'} \tau p_j \Sigma_{j'} + \lambda I_d)^2) K^{-2}. \quad (\text{E.233})$$

E.5.4 Computing $\mathbb{E}\bar{\text{tr}} r_j^{(4)}$

WLOG, we focus on $r_1^{(4)}$. The matrix of interest has a linear pencil representation given by (with zero-based indexing):

$$r_1^{(4)} = Q_{1,20}^{-1}, \quad (\text{E.234})$$

In conclusion:

$$\mathbb{E}\bar{\text{tr}} r_j^{(4)} = p_j \gamma p_{j'} \bar{\text{tr}} \Sigma_j \Sigma_{j'} A (\gamma \tau^2 (B e_j e_{j'} - p_j \Sigma_j e_j^2 u_{j'} - p_{j'} \Sigma_{j'} e_{j'}^2 u_j)) \quad (\text{E.242})$$

$$- \lambda \tau (e_j u_{j'} + e_{j'} u_j) I_d + e_j e_{j'} \rho I_d) K^{-2}. \quad (\text{E.243})$$

□

E.6 Theorem 6.3.2

Definition E.6.1. Let $(e_1, e_2, \tau_1, \tau_2, u_1, u_2, \rho_1, \rho_2)$ is be unique positive solution to the following system of fixed-point equations:

$$e_s = \frac{1}{1 + \psi_s \tau_s \bar{\text{tr}} \Sigma_s (\gamma \tau_s e_s \Sigma_s + \lambda_s I_d)^{-1}}, \text{ for } s \in \{1, 2\} \quad (\text{E.244})$$

$$\tau_s = \frac{1}{1 + \bar{\text{tr}} e_s \Sigma_s (\gamma \tau_s e_s \Sigma_s + \lambda_s I_d)^{-1}}, \text{ for } s \in \{1, 2\} \quad (\text{E.245})$$

$$u_s = \psi_s e_s^2 \bar{\text{tr}} \Sigma_s (\gamma \tau_s^2 (u_s + 1) \Sigma_s + \rho_s I_d) (\gamma \tau_s e_s \Sigma_s + \lambda_s I_d)^{-2}, \text{ for } s \in \{1, 2\} \quad (\text{E.246})$$

$$\rho_s = \tau_s^2 \bar{\text{tr}} (\gamma \rho_s (e_s \Sigma_s)^2 + \lambda_s^2 (u_s + 1) \Sigma_s) (\gamma \tau_s e_s \Sigma_s + \lambda_s I_d)^{-2}, \text{ for } s \in \{1, 2\}. \quad (\text{E.247})$$

For deterministic $d \times d$ PSD matrices A and B , we define the following auxiliary quantities:

$$h_j^{(1)}(A) := \gamma e_j \tau_j \bar{\text{tr}} A \Sigma_j (\gamma \tau_j e_j \Sigma_j + \lambda_j I_d)^{-1}, \quad (\text{E.248})$$

$$h_j^{(2)}(A) := \gamma \bar{\text{tr}} A \Sigma_j (\gamma e_j \tau_j^2 \Sigma_j + e_j \rho_j I_d - \lambda_j u_j \tau_j I_d) (\gamma \tau_j e_j \Sigma_j + \lambda_j I_d)^{-2}, \quad (\text{E.249})$$

$$h_j^{(3)}(A) := \bar{\text{tr}} A \Sigma_j (\gamma e_j^2 \Sigma_j (\gamma \tau_j^2 \Sigma_j + \rho_j I_d) + \lambda_j^2 u_j I_d) (\gamma \tau_j e_j \Sigma_j + \lambda_j I_d)^{-2}. \quad (\text{E.250})$$

Under Assumptions E.1.2 and 6.3.1, it holds that:

$$R_s(\widehat{f}_s) \simeq B_s(\widehat{f}_s) + V_s(\widehat{f}_s), \text{ with } V_s(\widehat{f}_s) = \lim_{p_s \rightarrow 1} V_s(\widehat{f}), \quad B_s(\widehat{f}_s) = \lim_{p_s \rightarrow 1} B_s(\widehat{f}). \quad (\text{E.251})$$

More explicitly:

$$V_s(\widehat{f}_s) = \sigma_s^2 \phi_s h_s^{(2)}(I_d), \quad B_s(\widehat{f}_s) = \bar{\text{tr}} \Theta_s \Sigma_s + h_s^{(3)}(\Theta_s) - 2h_s^{(1)}(\Theta_s \Sigma_s). \quad (\text{E.252})$$

Proof. Theorem 6.3.2 follows from Theorem 6.3.1 in the limit $p_s \rightarrow 1$ (i.e., $p_{s'} \rightarrow 0$). \square

The scalars e_s, τ_s, u_s, ρ_s can be intuitively interpreted in the setting where a separate model is learned for each group. For ridge regression with random projections and $\lambda \rightarrow 0^+$, we show in Equations E.273 and E.274 that e_s, τ_s are related to the normalized first-order degrees of freedom $I_{1,1}$ of the population covariance matrix Σ_s . e_s captures the effect of the feature rate ϕ_s while τ captures the effect of the parameterization rate γ . Similarly, for classical ridge regression, we show in Equation E.30 that e_s is related to the normalized first-order degrees of freedom $\bar{\text{df}}_1^{(s)}$. On the other hand, u_s and ρ_s can be understood as pseudo-variances. Indeed, for ridge regression with random projections and $\lambda \rightarrow 0^+$, Equations E.277 and E.285 show that u_s, ρ_s, V_s are all related to the normalized second-order degrees of freedom $I_{2,2}$ of Σ_s .

E.7 Solving Fixed-Point Equations for Theorem E.2.1

E.7.1 Proportional Covariance Matrices

When $\lambda \rightarrow 0^+$, it is not possible to analytically solve the fixed-point equations for the constants in Definition 6.3.1 for general Σ_1, Σ_2 . As such, we consider a more tractable case where the covariance matrices are proportional, i.e., $\Sigma_1 = a_1 \Sigma$ and $\Sigma_2 = a_2 \Sigma$, for some $\Sigma \in \mathbb{R}^{d \times d}$.

We define $\theta = \frac{\lambda}{\gamma \tau (a_1 p_1 e_1 + a_2 p_2 e_2)}$ and $\eta = \bar{\text{tr}} \Sigma (\Sigma + \theta I_d)^{-1}$. Then, we have that:

$$1/e_s = e'_s = 1 + \psi \tau \bar{\text{tr}} \Sigma_s K^{-1} = 1 + \frac{\phi a_s \eta}{a_1 p_1 e_1 + a_2 p_2 e_2}, \quad (\text{E.253})$$

$$1/\tau = \tau' = 1 + \bar{\text{tr}} L K^{-1} = 1 + (\eta/\gamma) \tau' = \frac{1}{1 - \eta/\gamma}. \quad (\text{E.254})$$

If $\theta_0 = 0$, then $\eta_0 = 1$. Therefore, $e'_s \rightarrow 1 + \frac{\phi a_s}{a_1 p_1 e_1 + a_2 p_2 e_2}$, which is a quadratic fixed-point equation. Accounting for the constraint that $e_s > 0$, the fixed-point equation requires that $\phi < 1$. Moreover, $\tau \rightarrow 1 - 1/\gamma$, which requires that $\gamma > 1$. We further observe that

$\rho \rightarrow (\tau^2 \bar{\text{tr}} \gamma L^2 K^{-2}) \rho$, which implies that $\rho \rightarrow 0$. We can then see that, for $c \in \{a_1, a_2\}$:

$$u_s \rightarrow \phi \gamma^2 \tau^2 e_s^2 a_s (a_1 p_1 u_1 + a_2 p_2 u_2 + c) \bar{\text{tr}} \Sigma^2 K^{-2} \quad (\text{E.255})$$

$$= \frac{\phi e_s^2 a_s (a_1 p_1 u_1 + a_2 p_2 u_2 + c)}{(a_1 p_1 e_1 + a_2 p_2 e_2)^2}, \quad (\text{E.256})$$

which is a linear fixed-point equation in u_s . In contrast, if $\theta_0 > 1$, we have $e'_s = 1 + \frac{\psi \tau a_s \eta \theta}{\lambda}$ and the equation:

$$\gamma \theta = \frac{\lambda}{(1 - \eta/\gamma) \left(\frac{a_1 p_1}{1 + \frac{\psi(1-\eta/\gamma)a_1 \eta \theta}{\lambda}} + \frac{a_2 p_2}{1 + \frac{\psi(1-\eta/\gamma)a_2 \eta \theta}{\lambda}} \right)}, \quad (\text{E.257})$$

which is a quartic equation in η . This highlights the difficulties of rigorously isolating the effects of different components on bias amplification. We empirically investigate how different components (e.g., covariance structures, group sizes) affect bias amplification and minority-group bias in Sections 6.4 and 6.5, and extensively validate that our theory predicts these implications.

E.8 Corollary E.8.1

As a special case of Theorem 6.3.1, we recover Corollary E.8.1, which aligns with Proposition 4 from [Bac23]. Theorem 6.3.1 is a non-trivial generalization of Proposition 4.

Corollary E.8.1 captures how the covariance matrix affects the test risk of a model through the normalized second and first-order degrees of freedom of Σ_s . Corollary E.8.1 also reveals that in the underparameterized regime ($\psi_s < 1$), the bias and variance of the test risk of the model strictly increase as a function of ψ_s (rate of parameters to samples); the test risk of the model explodes (i.e., there is catastrophic overfitting [Bac23]) when ψ_s gets close to 1. In the overparameterized regime ($\psi_s > 1$), the bias and variance of the test risk decrease as ψ_s increases.

Corollary E.8.1. *Under Assumptions E.1.2 and 6.3.1, it holds in the unregularized setting*

$\lambda_s \rightarrow 0^+$ that

$$B_s(\widehat{f}_s) = \begin{cases} \frac{\theta_0 \bar{\text{tr}} \Theta_s \Sigma_s (\Sigma_s + \theta_0 I_d)^{-1}}{1 - \psi_s}, & \gamma, \psi_s < 1 \\ 0, & \psi_s < 1, \gamma \geq 1 \text{ or } 1 \leq \psi_s \leq \gamma, \\ \frac{\theta_0^2 \bar{\text{tr}} \Theta_s \Sigma_s (\Sigma_s + \theta_0 I_d)^{-2}}{1 - \phi_s I_{2,2}(\theta_0)} + \frac{\theta_0 \bar{\text{tr}} \Theta_s \Sigma_s (\Sigma_s + \theta_0 I_d)^{-1}}{\psi_s - 1}, & \psi_s \geq 1, \psi_s \geq \gamma \end{cases} \quad (\text{E.258})$$

$$V_s(\widehat{f}_s) = \begin{cases} \frac{\sigma_s^2 \psi_s}{1 - \psi_s}, & \gamma, \psi_s < 1 \\ \frac{\sigma_s^2 \phi_s}{1 - \phi_s}, & \psi_s < 1, \gamma \geq 1 \text{ or } 1 \leq \psi_s \leq \gamma, \\ \frac{\sigma_s^2 \phi_s I_{2,2}(\theta_0)}{1 - \phi_s I_{2,2}(\theta_0)} + \frac{\sigma_s^2}{\psi_s - 1}, & \psi_s \geq 1, \psi_s \geq \gamma \end{cases} \quad (\text{E.259})$$

where $I_{a,b}(t) = \bar{\text{tr}} \Sigma^a (\Sigma + t I_d)^{-b}$ for any positive integers a, b ; and θ_0 is the unique solution to the following non-linear equation:

$$I_{1,1}(\theta_0) = \begin{cases} \gamma, & \gamma, \psi_s < 1 \\ 1, & \psi_s < 1, \gamma \geq 1 \text{ or } 1 \leq \psi_s \leq \gamma. \\ 1/\phi_s, & \psi_s \geq 1, \psi_s \geq \gamma \end{cases} \quad (\text{E.260})$$

Proof. Define $e' = 1/e_s \geq 0$, $\tau' = 1/\tau_s \geq 0$, $\theta = \lambda_s \tau' e' / \gamma$, and $\eta = I_{1,1}(\theta) \in [0, 1]$. One can then express e' and τ' as:

$$e' = 1 + \psi \tau_s \bar{\text{tr}} \Sigma (\gamma \tau_s e_s \Sigma + \lambda_s I_d)^{-1} = 1 + \phi_s \eta e', \quad (\text{E.261})$$

$$\tau' = 1 + \bar{\text{tr}} e_s \Sigma (\gamma \tau_s e_s \Sigma + \lambda_s I_d)^{-1} = 1 + (\eta/\gamma) \tau'. \quad (\text{E.262})$$

We deduce that:

$$e' = \frac{1}{1 - \phi_s \eta}, \quad (\text{E.263})$$

$$\tau' = \frac{1}{1 - \eta/\gamma}, \quad (\text{E.264})$$

$$\lambda \tau' e' = \gamma \theta. \quad (\text{E.265})$$

We define the following limiting values:

$$\lim_{\lambda_s \rightarrow 0^+} \theta \rightarrow \theta_0, \quad \lim_{\lambda_s \rightarrow 0^+} \eta \rightarrow \eta_0, \quad (\text{E.266})$$

$$\lim_{\lambda_s \rightarrow 0^+} e_s \rightarrow e_0, \quad \lim_{\lambda_s \rightarrow 0^+} \tau_s \rightarrow \tau_0, \quad (\text{E.267})$$

$$\lim_{\lambda_s \rightarrow 0^+} u_s \rightarrow u_0, \quad \lim_{\lambda_s \rightarrow 0^+} \rho_s \rightarrow \rho_0. \quad (\text{E.268})$$

There are now two cases to consider.

E.8.1 Case 1: $\theta_0 = 0$

This implies $\eta_0 = 1$. Therefore, by simple computation, $e_0 = 1/e'_0 = 1 - \phi_s \eta_0 = 1 - \phi_s$ and $\tau_0 = 1/\tau'_0 = 1 - 1/\gamma$. This requires $\phi_s \leq 1$ and $\gamma \geq 1$.

E.8.2 Case 2: $\theta_0 > 0$

Equation E.265 can be re-written as:

$$\frac{\lambda_s}{(1 - \phi_s \eta)(1 - \eta/\gamma)} = \gamma \theta, \quad \text{i.e., } \phi_s \eta^2 - (\psi_s + 1)\eta + \gamma - \frac{\lambda_s}{\theta} = 0. \quad (\text{E.269})$$

We solve this quadratic equation for η , arriving at the solutions:

$$\eta^\pm = \frac{\psi_s + 1 \pm \sqrt{(\psi_s + 1)^2 - 4(\psi_s - (\phi_s/\theta)\lambda_s)}}{2\phi_s} = \frac{\psi_s + 1 \pm \sqrt{(\psi_s + 1)^2 - 4(\psi_s - (\phi_s/\theta)\lambda_s)}}{2\phi_s} \quad (\text{E.270})$$

Taking the limit of η^\pm as $\lambda_s \rightarrow 0^+$ gives:

$$\eta^+ \rightarrow \frac{\psi_s + 1 + |\psi_s - 1|}{2\phi_s} = \begin{cases} \psi_s/\phi_s = \gamma, & \text{if } \psi_s \geq 1, \\ 1/\phi_s, & \text{if } \psi_s < 1, \end{cases} \quad (\text{E.271})$$

$$\eta^- \rightarrow \frac{\psi_s + 1 - |\psi_s - 1|}{2\phi_s} = \begin{cases} 1/\phi_s, & \text{if } \psi_s \geq 1, \\ \psi_s/\phi_s = \gamma, & \text{if } \psi_s < 1. \end{cases}$$

Recall that we have the following constraints:

- $e' \geq 0, \tau' \geq 0$.
- $\eta \in [0, 1]$.

We can show that $\eta_0 = 1/\phi_s$ is incompatible with $\psi_s < 1$. Indeed, otherwise we would have $\tau'_0 = 1/(1 - \eta_0/\gamma) = 1/(1 - 1/\psi_s) < 0$. Similarly, if $\psi_s > 1$, we would have $e_0 = 1 - \phi_s\gamma = 1 - \psi_s < 0$. Therefore, $\eta_0 = \eta^-$. Furthermore, if $\psi_s, \gamma < 1$, it must be that $\theta_0 > 0$ and $\eta_0 = \gamma$. Instead, if $\psi_s < 1, \gamma \geq 1$, we must have that $\phi_s \leq 1$, and therefore, $\theta_0 = 0$ and $\eta_0 = 1$. Similarly, if $\psi_s \geq 1, \gamma \geq 1$, and $\phi_s \leq 1$ (i.e., $1 \leq \psi_s \leq \gamma$), we must have that $\theta_0 = 0$ and $\eta_0 = 1$. In all other cases where $\psi_s \geq 1$, it must be that $\eta_0 = 1/\phi_s$ (which additionally requires $\phi_s \geq 1$ or $\psi_s \geq \gamma$). Succinctly:

$$\eta_0 = \begin{cases} \gamma, & \gamma, \psi_s < 1 \\ 1, & \psi_s < 1, \gamma \geq 1 \text{ or } 1 \leq \psi_s \leq \gamma \\ 1/\phi_s, & \psi_s \geq 1, \psi_s \geq \gamma \end{cases} \quad (\text{E.272})$$

Plugging this into Equation E.263 and Equation E.264 gives:

$$e_0 = 1 - \phi_s\eta_0 = 1 - \phi_s I_{1,1}(\theta_0), \quad (\text{E.273})$$

$$\tau_0 = 1 - \eta_0/\gamma = 1 - I_{1,1}(\theta_0)/\gamma. \quad (\text{E.274})$$

We will now solve for u_0 and ρ_0/τ_0^2 . We can re-write u_s and ρ_s/τ_s^2 as:

$$\rho_s/\tau_s^2 = \gamma^{-1}(\rho_s/\tau_s^2)I_{2,2}(\theta) + \theta^2(u_s + 1)I_{1,2}(\theta), \quad (\text{E.275})$$

$$\tau_s^2 u_s = \tau_s^2 \phi_s (u_s + 1)I_{2,2}(\theta) + \phi_s \gamma^{-1} \rho_s I_{1,2}(\theta). \quad (\text{E.276})$$

Solving for u_0 and ρ_0/τ_0^2 yields:

$$u_0 = \frac{\phi\zeta}{\gamma - \phi\zeta - I_{2,2}(\theta_0)}, \quad \rho_0/\tau_0^2 = \frac{\gamma\theta_0^2 I_{2,2}(\theta_0)}{\gamma - \phi\zeta - I_{1,2}(\theta_0)}, \quad (\text{E.277})$$

$$\text{where } \zeta = I_{2,2}(\theta_0)(\gamma - I_{2,2}(\theta_0)) + \theta_0^2 I_{1,2}(\theta_0)^2. \quad (\text{E.278})$$

We can then see for the variance term that:

$$V_s(\widehat{f}_s) = \sigma_s^2 \phi_s \gamma \bar{\text{tr}} \Sigma_s (\gamma e_s \tau_s^2 \Sigma_s + e_s \rho_s I_d - \lambda_s u_s \tau_s I_d) (\gamma \tau_s e_s)^{-2} (\Sigma_s + \theta I_d)^{-2} \quad (\text{E.279})$$

$$= \sigma_s^2 \phi_s (1/e_s) \bar{\text{tr}} \Sigma_s^2 (\Sigma_s + \theta I_d)^{-2} + (\sigma_s^2 \phi_s / \gamma) (1/e_s) (\rho_s / \tau_s^2) \bar{\text{tr}} \Sigma_s (\Sigma_s + \theta I_d)^{-2} \quad (\text{E.280})$$

$$- \sigma_s^2 \phi_s (u_s) (1/e_s) \theta \bar{\text{tr}} \Sigma_s (\Sigma_s + \theta I_d)^{-2} \quad (\text{E.281})$$

$$= \sigma_s^2 \phi_s I_{2,2}(\theta) / e_s + \sigma_s^2 \phi_s (\rho_s / \tau_s^2) I_{1,2}(\theta) / (\gamma e_s) - \sigma_s^2 \phi_s u_s \theta I_{1,2}(\theta) / e_s \quad (\text{E.282})$$

$$\rightarrow \frac{\sigma_s^2 \phi_s I_{2,2}(\theta_0) - \sigma_s^2 \phi_s u_0 \theta_0 I_{1,2}(\theta_0)}{1 - \phi_s I_{1,1}(\theta_0)} + \frac{\sigma_s^2 \phi_s \rho_0 / \tau_0^2}{\gamma (1 - \phi_s I_{1,1}(\theta_0))} \quad (\text{E.283})$$

$$= - \frac{\sigma_s^2 \phi_s \xi}{\phi_s \xi + I_{2,2}(\theta_0) - \gamma}, \quad (\text{E.284})$$

where $\xi = I_{1,1}^2(\theta_0) - 2I_{1,1}(\theta_0)I_{2,2}(\theta_0) + I_{2,2}(\theta_0)\gamma$ and we have used the fact that $I_{1,2}(\theta) = (I_{1,1}(\theta) - I_{2,2}(\theta)) / \theta$. Plugging in $I_{1,1}(\theta_0) = \eta_0$, we have that:

$$V_s(\widehat{f}_s) \rightarrow \begin{cases} \frac{\sigma_s^2 \psi_s}{1 - \psi_s}, & \gamma, \psi_s < 1 \\ \frac{\sigma_s^2 \phi_s}{1 - \phi_s}, & \psi_s < 1, \gamma \geq 1 \text{ or } 1 \leq \psi_s \leq \gamma, \\ \frac{\sigma_s^2 \phi_s I_{2,2}(\theta_0)}{1 - \phi_s I_{2,2}(\theta_0)} + \frac{\sigma_s^2}{\psi_s - 1}, & \psi_s \geq 1, \psi_s \geq \gamma \end{cases} \quad (\text{E.285})$$

where we have used that $I_{2,2}(\theta_0) = I_{2,2}(0) = 1$ in the second case.

Likewise, for the bias term, we obtain:

$$B_s(\widehat{f}_s) = \bar{\text{tr}} \Theta_s \Sigma_s + \bar{\text{tr}} \Theta_s \Sigma_s (\gamma e_s^2 \Sigma_s (\gamma \tau_s^2 \Sigma_s + \rho_s I_d) + \lambda_s^2 u_s I_d) (\gamma \tau_s e_s \Sigma_s + \lambda_s I_d)^{-2} \quad (\text{E.286})$$

$$- 2\gamma e_s \tau_s \bar{\text{tr}} \Theta_s \Sigma_s^2 (\gamma \tau_s e_s \Sigma_s + \lambda_s I_d)^{-1} \quad (\text{E.287})$$

$$\rightarrow \bar{\text{tr}} \Theta_s \Sigma_s (\Sigma_s^2 + 2\theta_0 \Sigma_s + \theta_0^2 I_d) (\Sigma_s + \theta_0 I_d)^{-2} \quad (\text{E.288})$$

$$+ \bar{\text{tr}} \Theta_s \Sigma_s (\Sigma_s^2) (\Sigma_s + \theta_0 I_d)^{-2} \quad (\text{E.289})$$

$$+ \bar{\text{tr}} \Theta_s \Sigma_s ((\rho_0 / \tau_0^2) \Sigma_s / \gamma) (\Sigma_s + \theta_0 I_d)^{-2} \quad (\text{E.290})$$

$$+ \bar{\text{tr}} \Theta_s \Sigma_s (\theta_0^2 u_0 I_d) (\Sigma_s + \theta_0 I_d)^{-2} \quad (\text{E.291})$$

$$+ \bar{\text{tr}} \Theta_s \Sigma_s (-2\Sigma_s^2 - 2\theta_0 \Sigma_s) (\Sigma_s + \theta_0 I_d)^{-2} \quad (\text{E.292})$$

$$= \theta_0^2 (u_0 + 1) \bar{\text{tr}} \Theta_s \Sigma_s (\Sigma_s + \theta_0 I_d)^{-2} + (1/\gamma) (\rho_0 / \tau_0^2) \bar{\text{tr}} \Theta_s \Sigma_s^2 (\Sigma_s + \theta_0 I_d)^{-2}. \quad (\text{E.293})$$

Again, plugging in $I_{1,1}(\theta_0) = \eta_0$, we have that:

$$B_s(\widehat{f}_s) \rightarrow \begin{cases} \frac{\theta_0 \bar{\text{tr}} \Theta_s \Sigma_s (\Sigma_s + \theta_0 I_d)^{-1}}{1 - \psi_s}, & \gamma, \psi_s < 1 \\ \frac{\theta_0^2 \bar{\text{tr}} \Theta_s \Sigma_s (\Sigma_s + \theta_0 I_d)^{-2}}{1 - \phi_s} = 0, & \psi_s < 1, \gamma \geq 1 \text{ or } 1 \leq \psi_s \leq \gamma, \\ \frac{\theta_0^2 \bar{\text{tr}} \Theta_s \Sigma_s (\Sigma_s + \theta_0 I_d)^{-2}}{1 - \phi_s I_{2,2}(\theta_0)} + \frac{\theta_0 \bar{\text{tr}} \Theta_s \Sigma_s (\Sigma_s + \theta_0 I_d)^{-1}}{\psi_s - 1}, & \psi_s \geq 1, \psi_s \geq \gamma \end{cases} \quad (\text{E.294})$$

where we have used that $\bar{\text{tr}} \Theta_s \Sigma_s^2 (\Sigma_s + \theta_0 I_d)^{-2} = \bar{\text{tr}} \Theta_s \Sigma_s (\Sigma_s + \theta_0 I_d)^{-1} - \theta_0 \bar{\text{tr}} \Theta_s \Sigma_s (\Sigma_s + \theta_0 I_d)^{-2}$ and in the second case, $\theta_0 = 0$ and $I_{2,2}(\theta_0) = 1$. \square

E.9 Experimental Details

E.9.1 Synthetic Experiments

Across all experiments on synthetic data, we choose $n = 400$. We further use 5 runs to estimate test risks (e.g., $\mathbb{E}R_s(\widehat{f}), \mathbb{E}R_s(\widehat{f}_s)$), and 5 runs to capture the variance of the estimators, for a total of 25 runs. We use 10,000 samples to estimate test risks.

Our experiments validate that bias amplification occurs even in low-dimensional regimes. In Sections 6.4 and 6.5, and Appendices E.10, E.13, and E.15, we show that our theory predicts bias amplification for models trained on only $n = 400$ samples. The high-dimensional regime is commonly studied in ML theory and statistical physics (as we mention in Section 6.1.2), as it makes precise analysis more tractable.

Setup for Section 6.4.1. We further choose $\lambda = 1 \times 10^{-6}$ to approximate the minimum-norm interpolator; we henceforth set $\lambda = \lambda_1 = \lambda_2$ for simplicity. We modulate $a_1, a_2, \sigma_1^2, \sigma_2^2$, as well as ψ (rate of parameters to samples) and ϕ (rate of features to samples) to understand the effects of model size, number of features, and sample size on bias amplification. We consider diverse and dense values of these variables to obtain a clear picture of when and how models amplify bias.

Setup for Section 6.5. The first πd features represent common *core* features of groups 1 and 2 while the latter $(1 - \pi)d$ features capture unshared *extraneous* features for group 2 (e.g., spurious features). Intuitively, this setting can model: (1) learning from data from two groups where one group suffers from spurious features [SRK20], or (2) learning from a mixture of raw data (i.e., with spurious features) and clean data (i.e., without spurious features) for a single population [KL21]. We ask: Does our theory predict how the inclusion of different amounts of extraneous features affect the test risk of a minority group (compared to the majority group) when a single model is trained on data from both groups vs. a separate model is trained per group?

Although [SRK20] considers classification instead of regression, to mirror their experimental setting, we pick $p_1 = 0.9$ (i.e., group 1 is much larger than group 2) and $\Theta = I_d, \Delta = 0$ (i.e., $w_1^* = w_2^*$). We additionally choose $\lambda = 1 \times 10^{-6}$ and $\sigma_1^2 = \sigma_2^2 = 1$. We modulate a_1, b_2 , as well as ψ (rate of parameters to samples) and ϕ (rate of features to samples). Notably, this setting also captures learning problems with $o(d)$ overlapping core and extraneous features in our asymptotic scaling limit. An extremization of this setting is choosing $\Sigma_1 = a_1 I_{\pi d} \oplus 0 I_{(1-\pi)d}, \Sigma_2 = 0 I_{\pi d} \oplus b_2 I_{(1-\pi)d}$, where groups 1 and 2 have no overlapping features.

The experiments in Section 6.5 validate that our analysis does not rely on conditional dependence heterogeneity. That is, we empirically verify that our theory still holds and predicts bias amplification occurs even when $w_1^* = w_2^*$ (see Figure 6.4). In essence, the structure and eigenspectra of the covariance matrices of the two groups still contribute to bias amplification even when the ground-truth weights for the groups are the same. In our theory, we only allow the possibility of $w_1^* \neq w_2^*$ to be as general as possible. In practice, labeling rule heterogeneity may be leveraged, for example, to train a mixture of experts that is regularized to deamplify bias.

Extraneous vs. Spurious Features. Our usage of extraneous features (i.e., features that are different across groups and correlated with labels) differs from classical definitions of spuriousness (i.e., non-causal correlations between features and labels) [BW24]; indeed, the extraneous features are used to generate the labels of the minority group. For example, [KL21] models both the labels and spurious features in linear regression as being separately generated by the core features, such that the labels and spurious features are associated but not causally related. However, this setup is not encompassed by our modeling assumptions, as it entails that the ground-truth parameter and feature covariance matrices are not jointly diagonalizable. In contrast, [SRK20] studies spurious correlations in classification. At a high level, [SRK20] creates four subgroups of a population with different combinations of class labels $y \in \{-1, 1\}$ and group labels $a \in \{-1, 1\}$. The core and spurious features are then sampled from normal distributions parameterized by y, a (respectively) and different variance levels. By setting the spurious features to have a significantly lower variance than the core features and making y and a highly associated (i.e., imbalanced groups), the authors coerce models to perform classification as a function of primarily the spurious features of the majority group, which does not generalize to the minority group. To capture this spirit, our setup uses imbalanced groups, and the data for the majority group provides no learning signal for the extraneous features; this coerces models to perform regression as a function of primarily the core features, without learning appropriate parameters for the extraneous features, and thus generalize poorly to the minority group.

E.9.2 Colored MNIST Experiments

Train-test split. Colored MNIST has a total of 60k instances. Each image is $28 \times 28 \times 3$ pixels. We use the prescribed 0.67-0.33 train-test split. We do not perform validation of hyperparameters, which we mostly adopt³.

³https://colab.research.google.com/github/reinakano/invariant-risk-minimization/blob/master/invariant_risk_minimization_colored_mnist.ipynb

Model architecture. By default, our CNN architecture consists of: (1) a convolutional layer (3 in-channels, 20 out-channels, kernel size of 5, stride of 1); (2) a max pooling layer (kernel size of 2, stride of 2); (3) a second convolutional layer (20 in-channels, 50 out-channels, kernel size of 5, stride of 1); (4) a second max-pooling layer (kernel size of 2, stride of 2); (5) a fully-connected layer ($\mathbb{R}^{800} \rightarrow \mathbb{R}^{500}$); and (6) a second fully-connected layer ($\mathbb{R}^{500} \rightarrow \mathbb{R}^1$).

Model training. We train each model with a batch size of 250 for a single epoch with respect to groups (i.e., 80 training steps given there are two groups). We use a cross-entropy loss and the Adam optimizer with learning rate 0.01. We run all experiments on a single NVIDIA L40S. We report our results over 10 random seeds.

E.10 Bias Amplification Plots

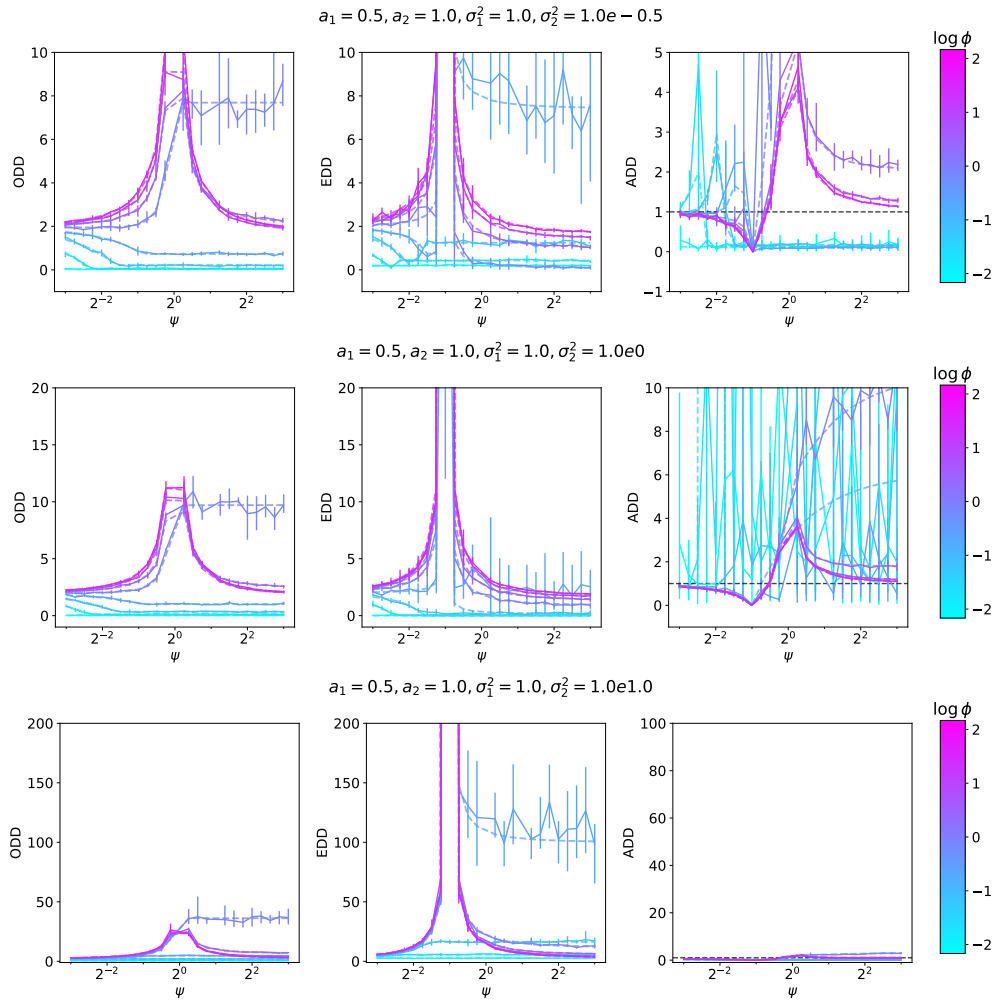


Figure E.2: We empirically demonstrate that bias amplification occurs and validate our theory (Theorems 6.3.1 and 6.3.2) for ODD , EDD , and ADD under the setup described in Section 6.4.1. The solid lines capture empirical values while the corresponding lower-opacity dashed lines represent what our theory predicts. We plot ODD and EDD on the same scale for easy comparison, and include a black dashed line at $ADD = 1$ to contrast bias amplification vs. deamplification. The error bars capture the range of the estimators over 25 random seeds.

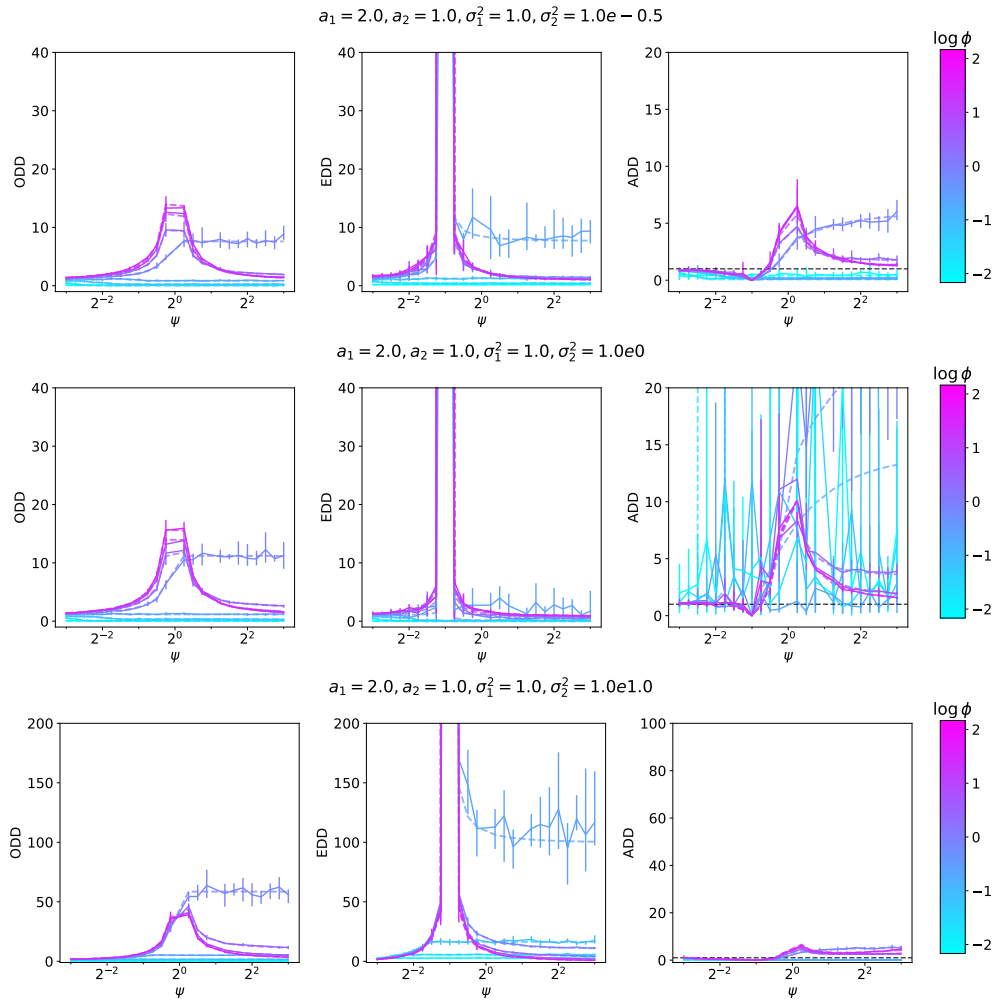


Figure E.3: We empirically demonstrate that bias amplification occurs and validate our theory (Theorems 6.3.1 and 6.3.2) for *ODD*, *EDD*, and *ADD* under the setup described in Section 6.4.1. The solid lines capture empirical values while the corresponding lower-opacity dashed lines represent what our theory predicts. We plot *ODD* and *EDD* on the same scale for easy comparison, and include a black dashed line at $ADD = 1$ to contrast bias amplification vs. deamplification. The error bars capture the range of the estimators over 25 random seeds.

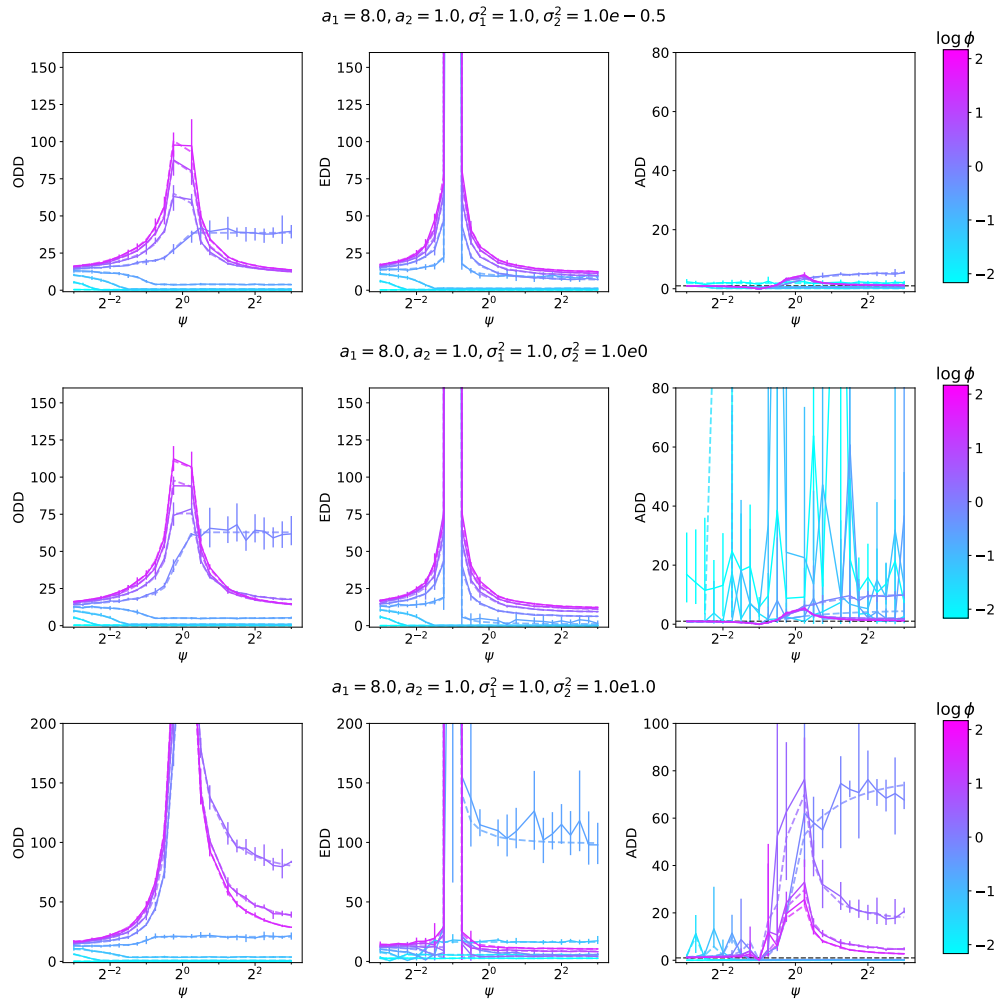


Figure E.4: We empirically demonstrate that bias amplification occurs and validate our theory (Theorems 6.3.1 and 6.3.2) for ODD , EDD , and ADD under the setup described in Section 6.4.1. The solid lines capture empirical values while the corresponding lower-opacity dashed lines represent what our theory predicts. We plot ODD and EDD on the same scale for easy comparison, and include a black dashed line at $ADD = 1$ to contrast bias amplification vs. deamplification. The error bars capture the range of the estimators over 25 random seeds.

E.11 Power-Law Covariance

To better understand how ϕ (rate of features to samples) and the label noise ratio c affect bias amplification, we derive explicit phase transitions in the bias amplification profile of unregularized ridge regression with random projections in terms of these quantities. We consider the setting of power-law covariance, as it is analytically tractable and can be translated to the case of wide neural networks [CV07, CLK22, MRS22], where the exponents can be empirically gauged. Let the eigenvalues $\lambda_k^{(s)}$ of Σ_s have power-law decay, i.e., $\lambda_k^{(s)} = k^{-\beta_s}$, for all k and some positive constants β_1 and β_2 . WLOG, we will assume $\beta_1 > \beta_2$. Note that β_s controls the effective dimension and ultimately the difficulty of fitting the noiseless part of the signal from group s . If β_s is large, then all the information is concentrated in a few features, and so the learning problem is easier. We similarly assume that the eigenvalues μ_k of Δ have power-law decay $\mu_k = k^{-\alpha}$, for all k and constant $\alpha > 0$. Finally, we consider balanced groups (i.e., $p_1 = p_2 = 1/2$). Under this setup, we have the following corollary.

Corollary E.11.1. *Suppose that in the single model setting, as $\lambda \rightarrow 0^+$, $(e_1, e_2, \tau, u_1, u_2, \rho)$ is the unique positive solution to the following fixed-point equations:*

$$1/\tau = 1 + 1/(\gamma\tau), \quad 1/e_s = 1 + \phi \bar{\text{tr}} \Sigma_s L^{-1}, \quad \text{for } s \in \{1, 2\}, \quad (\text{E.295})$$

$$\rho = 0, \quad u_s = \phi e_s^2 \bar{\text{tr}} \Sigma_s D L^{-2}, \quad \text{for } s \in \{1, 2\}, \quad (\text{E.296})$$

$$\text{where: } L = p_1 e_1 \Sigma_1 + p_2 e_2 \Sigma_2, \quad D = p_1 u_1 \Sigma_1 + p_2 u_2 \Sigma_2 + B. \quad (\text{E.297})$$

Furthermore, suppose $\psi_s < 1, \gamma \geq 1$ or $1 \leq \psi_s \leq \gamma$. Under the assumptions of Theorem 6.3.1 and Assumption E.1.3, as $\lambda \rightarrow 0^+$, we have the following approximate analytical phase transitions in the bias amplification profile of ridge regression with random projections:

$$\lim_{\substack{d, n_1, n_2 \rightarrow \infty \\ \phi_{1,2} \rightarrow 2\phi}} ADD \rightarrow \frac{c}{|c-1|}, \quad \lim_{c \rightarrow 0^+} \lim_{\substack{d, n_1, n_2 \rightarrow \infty \\ \phi_{1,2} \rightarrow 2\phi}} ADD \rightarrow 0, \quad (\text{E.298})$$

$$\lim_{c \rightarrow \infty} \lim_{\substack{d, n_1, n_2 \rightarrow \infty \\ \phi_{1,2} \rightarrow 2\phi}} ADD \rightarrow 1, \quad \lim_{c \rightarrow 1} \lim_{\substack{d, n_1, n_2 \rightarrow \infty \\ \phi_{1,2} \rightarrow 2\phi}} ADD \rightarrow \infty. \quad (\text{E.299})$$

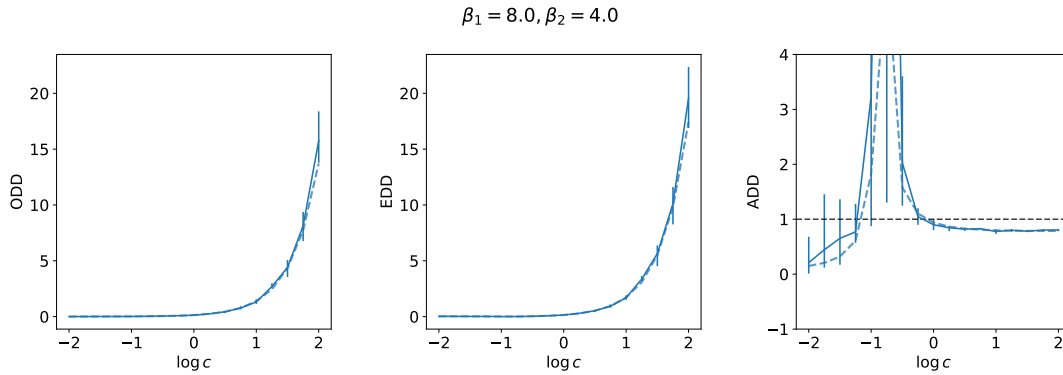


Figure E.5: **Our theory predicts that bias amplification is larger for higher noise ratios than lower noise ratios.** We observe that Corollary E.11.1 generally predicts the *ADD* profile with respect to the noise ratio c . The solid lines capture empirical values while the corresponding lower-opacity dashed lines represent what Theorem 6.3.1 predicts. We plot *ODD* and *EDD* on the same scale for easy comparison, and include a black dashed line at $ADD = 1$ to contrast bias amplification vs. deamplification. The error bars capture the range of the estimators over 25 random seeds. We consider the setup described in Appendix E.11 with $\psi = 0.5$, $\phi = 0.2$, and $\lambda = 1 \times 10^{-6}$.

We relegate the proof to Appendix E.12 and empirically assess the validity of this result in Figure E.5. The phase transitions reveal that bias amplification peaks near $c = 1$, bias deamplification peaks when $c \rightarrow 0^+$, and bias is roughly neither amplified or deamplified when $c \rightarrow \infty$. Furthermore, the right tail of the *ODD* profile (which Corollary E.11.1 predicts to be proportional to c) is higher than the left tail (i.e., 0) for larger c . However, the left tail of the *EDD* profile (which Corollary E.11.1 predicts to be proportional to $|c - 1|$) does not increase steeply as $c \rightarrow 0^+$. Interestingly, in the proof of Corollary E.11.1, we observe that the bias term depends on $\bar{\text{tr}} \Delta \Sigma_s$; therefore, the setting $\forall k, \lambda_k^{(s)} \geq 1/\mu_k$ (e.g., common in learning from synthetic data [DFK24]) can prevent the bias term from vanishing or even cause it to explode. This may explain why training models on synthetic data (i.e., data previously generated by the model) may amplify unfairness [WSP24].

E.12 Proof of Corollary E.11.1

Proof. We begin by computing the *ODD* in the limit $\lambda \rightarrow 0^+$. We define $u_j^{(s)} = u_j$ for $B = \Sigma_s$. By Assumption E.1.3, we can re-express the constants in Definition 6.3.1 in terms of the limiting spectral densities of the covariance matrices:

$$e_1 = \frac{1}{1 + \phi \int_0^\infty \frac{1}{p_1 e_1 + p_2 e_2 r} d\nu(r)}, e_2 = \frac{1}{1 + \phi \int_0^\infty \frac{r}{p_1 e_1 + p_2 e_2 r} d\nu(r)}, \quad (\text{E.300})$$

$$\tau = \frac{1}{1 + \frac{1}{\gamma\tau}} = 1 - 1/\gamma, \rho = 0, \quad (\text{E.301})$$

$$u_1^{(1)} = \phi e_1^2 \int_0^\infty \frac{u_1^{(1)} p_1 + u_2^{(1)} p_2 r + 1}{(p_1 e_1 + p_2 e_2 r)^2} d\nu(r), u_2^{(1)} = \phi e_2^2 \int_0^\infty \frac{u_1^{(1)} p_1 r + u_2^{(1)} p_2 r^2 + r}{(p_1 e_1 + p_2 e_2 r)^2} d\nu(r), \quad (\text{E.302})$$

$$u_1^{(2)} = \phi e_1^2 \int_0^\infty \frac{u_1^{(2)} p_1 + u_2^{(2)} p_2 r + r}{(p_1 e_1 + p_2 e_2 r)^2} d\nu(r), u_2^{(2)} = \phi e_2^2 \int_0^\infty \frac{u_1^{(2)} p_1 r + u_2^{(2)} p_2 r^2 + r^2}{(p_1 e_1 + p_2 e_2 r)^2} d\nu(r). \quad (\text{E.303})$$

Since $\beta_1 > \beta_2$, $-\beta_2 - (-\beta_1) > 0$. As such, for $d \rightarrow \infty$, the ratios $r_k = \lambda_k^{(2)}/\lambda_k^{(1)}$ have the approximate limiting distribution $\nu = \delta_{r=\infty}$, i.e., a Dirac atom at infinity. Thus:

$$e_1 = 1, e_2 = 1 - \frac{\phi}{p_2} = 1 - \phi_2, \tau = 1 - 1/\gamma, \rho = 0, \quad (\text{E.304})$$

$$u_1^{(1)} = 0, u_2^{(1)} = 0, u_1^{(2)} = 0, u_2^{(2)} = \frac{\phi}{p_2(p_2 - \phi)}. \quad (\text{E.305})$$

Now, we can re-express the variance terms as:

$$V_1(\hat{f}) = \phi \sigma_1^2 \int_0^\infty \frac{p_1}{(p_1 + p_2 e_2 r)^2} d\nu(r) + \phi \sigma_2^2 \int_0^\infty \frac{p_2 e_2 r}{(p_1 + p_2 e_2 r)^2} d\nu(r) = 0, \quad (\text{E.306})$$

$$V_2(\hat{f}) = \phi \sigma_1^2 \int_0^\infty \frac{p_1 r + p_1 p_2 u_2^{(2)} r}{(p_1 + p_2 e_2 r)^2} d\nu(r) + \phi \sigma_2^2 \int_0^\infty \frac{p_2 e_2 r^2}{(p_1 e_1 + p_2 e_2 r)^2} d\nu(r) = \frac{\sigma_2^2 \phi}{p_2 - \phi}. \quad (\text{E.307})$$

Likewise, we can re-express the bias terms as:

$$B_1(\hat{f}) = \int_0^\infty \int_0^\infty \int_0^\infty \frac{a \delta e_2^2 p_2^2 r^2}{(e_1 p_1 + e_2 p_2 r)^2} d\mu(r, a) d\pi(\delta) = \int_0^\infty \int_0^\infty a \delta d\mu(a) d\pi(\delta) = 0, \quad (\text{E.308})$$

$$B_2(\hat{f}) = 0. \quad (\text{E.309})$$

In this calculation, we observe that the adversarial setting $\forall k, \lambda_k^{(1)} \geq 1/\mu_k$ can prevent the bias term from vanishing. Putting these pieces together and recalling that $p_2 = 1/2$:

$$ODD \rightarrow \left| V_1(\hat{f}) - V_2(\hat{f}) \right| = \frac{2\phi\sigma_1^2}{1-2\phi}c. \quad (\text{E.310})$$

We now compute the *EDD*. We can once again re-express the constants in Definition E.6.1 in terms of the limiting spectral densities of the covariance matrices:

$$e_s = \frac{1}{1 + \phi_s/e_s} = 1 - \phi_s, \tau_s = 1 - 1/\gamma. \quad (\text{E.311})$$

By Corollary E.8.1, because $\psi_s < 1, \gamma \geq 1$ or $1 \leq \psi_s \leq \gamma$, $B_s(\hat{f}_s) = 0$ and $V_s(\hat{f}_s) = \frac{\sigma_s^2\phi_s}{1-\phi_s}$.

Therefore, because $\phi = p_s\phi_s$:

$$EDD \rightarrow \left| V_1(\hat{f}_1) - V_2(\hat{f}_2) \right| = \frac{2\phi}{1-2\phi} \left| \sigma_1^2 - \sigma_2^2 \right| = \frac{2\phi\sigma_1^2}{1-2\phi} |c-1|, \quad (\text{E.312})$$

$$ADD \rightarrow \frac{c}{|c-1|}. \quad (\text{E.313})$$

□

E.13 Bias Amplification During Training

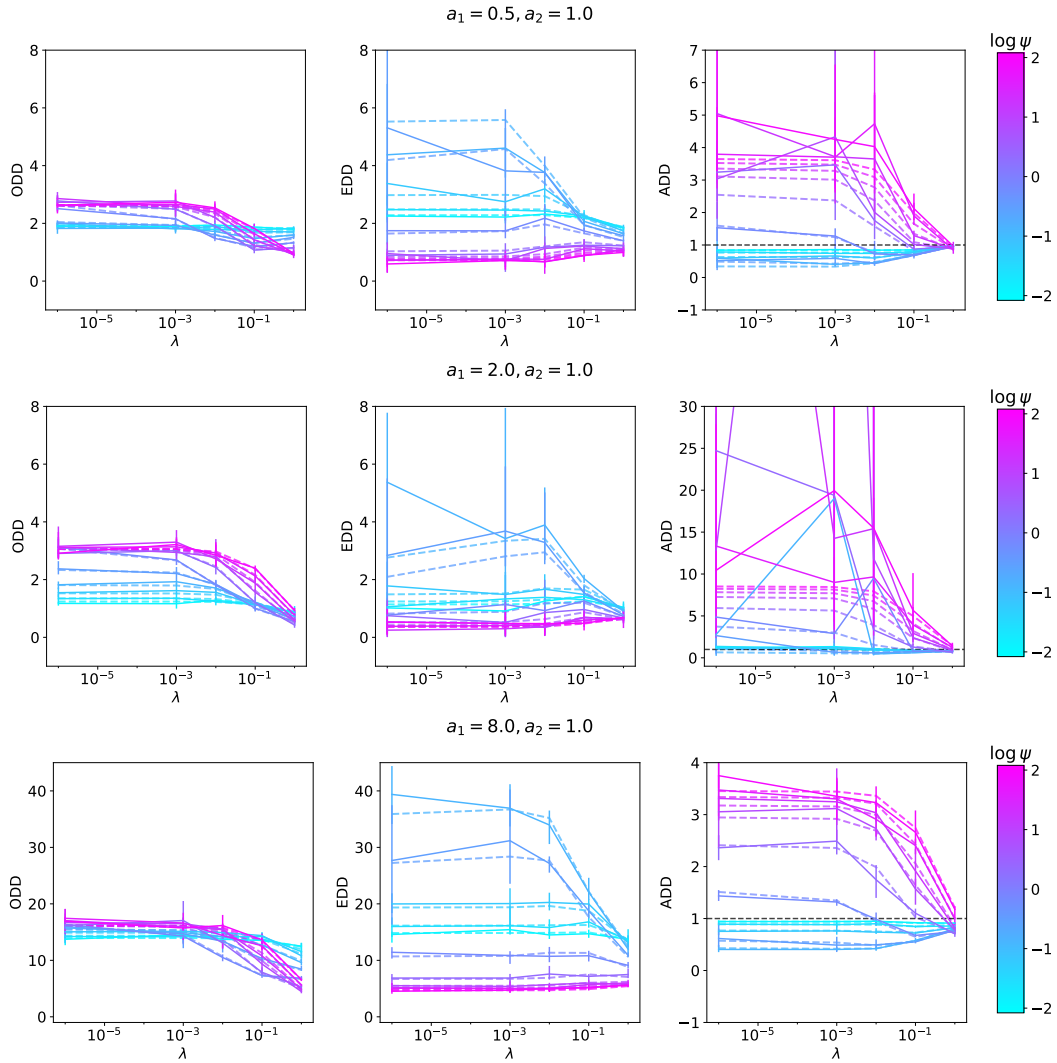


Figure E.6: **Our theory reveals that there may be an optimal regularization penalty to deamplify bias.** We empirically demonstrate that bias amplification can be heavily affected by λ and validate our theory (Theorems 6.3.1 and 6.3.2) for ODD , EDD , and ADD under the setup described in Section 6.4.2. The solid lines capture empirical values while the corresponding lower-opacity dashed lines represent what our theory predicts. We include a black dashed line at $ADD = 1$ to contrast bias amplification vs. deamplification. The error bars capture the range of the estimators over 25 random seeds.

E.14 Colored MNIST Plots

We further assess the applicability of our conclusions about the effects of label noise (Figures 6.3, E.7) and model size (Figure E.8) on bias amplification for Colored MNIST. Please see Section 6.4.2 for a discussion of Figure 6.3. We observe in Figure E.7 that as we increase the label noise ratio, the *EDD* generally increases, while the *ODD* remains relatively low, which is suggested by our theoretical reasoning in Section 6.4.2. Furthermore, in Figure E.8, as the hidden dimension m of the penultimate layer of the CNN increases, the *ODD* appears to decrease and plateau, which is predicted by our theoretical results (see Section 6.4.1) in the Colored MNIST regime where $\phi < 1$. However, the *EDD* does not appear to decrease; while this is plausibly predicted by our theory, it requires going beyond our simplistic assumption that Σ_1 roughly coincides with Σ_2 and studying the interplay between ϕ_s, ψ_s, Σ_s for each group s (as suggested by Appendix E.10).

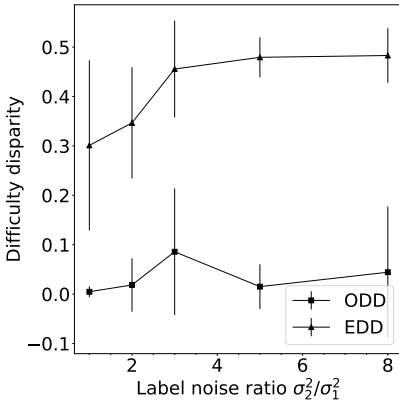


Figure E.7: **Our theory predicts that more disparate label noise between groups deamplifies bias on Colored MNIST.** We plot the *ODD* and *EDD* of a CNN for different label noise ratios $c = \sigma_2^2/\sigma_1^2$ for Colored MNIST. As c increases, the *EDD* generally increases while the *ODD* remains relatively low, which is predicted by our theory (see reasoning in Section 6.4.2). In our experiments, $\sigma_1^2 = 0.05$ stays fixed while σ_2^2 varies. For each value of c , the model is evaluated after $t = 80$ training steps and has a penultimate layer with dimension $m = 500$. The error bars capture the standard deviation computed over 10 random seeds.

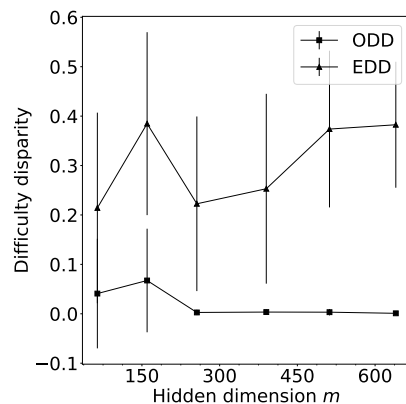


Figure E.8: **Our theory predicts that a larger model size reduces bias on Colored MNIST in the single model setting.** We plot the *ODD* and *EDD* of a CNN for different model sizes m (where m is the dimension of the penultimate CNN layer) for Colored MNIST. As m increases, the *ODD* appears to decrease and plateau, which is in line with what our theory predicts in the regime where $\phi < 1$ (see analysis in Section 6.4.1). The *EDD* does not tend towards 0. In our experiments, $\sigma_1^2 = \sigma_2^2 = 0.05$. For each value of m , the model is evaluated after $t = 80$ training steps. The error bars capture the standard deviation computed over 10 random seeds.

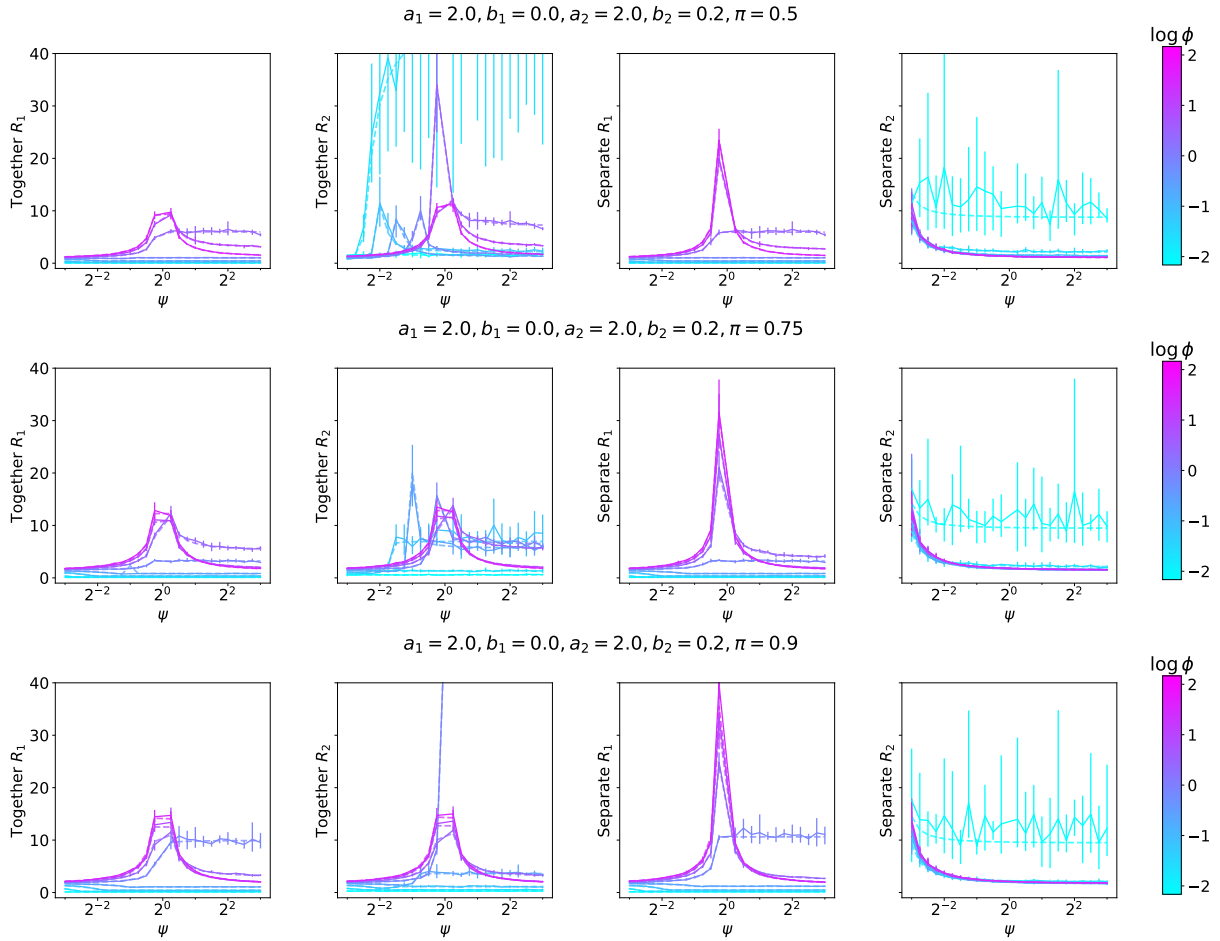


Figure E.9: We empirically demonstrate that minority-group bias is affected by extraneous features. We validate our theory (Theorems 6.3.1 and 6.3.2) for together R_1, R_2 (i.e., single model learned for both groups) and separate R_1, R_2 (i.e., separate model learned per group) under the setup described in Section 6.4.2. The solid lines capture empirical values while the corresponding lower-opacity dashed lines represent what our theory predicts. We include a black dashed line at $ADD = 1$ to contrast bias amplification vs. deamplification. All y-axes are on the same scale for easy comparison. The error bars capture the range of the estimators over 25 random seeds.

E.15 Minority-Group Bias Plots

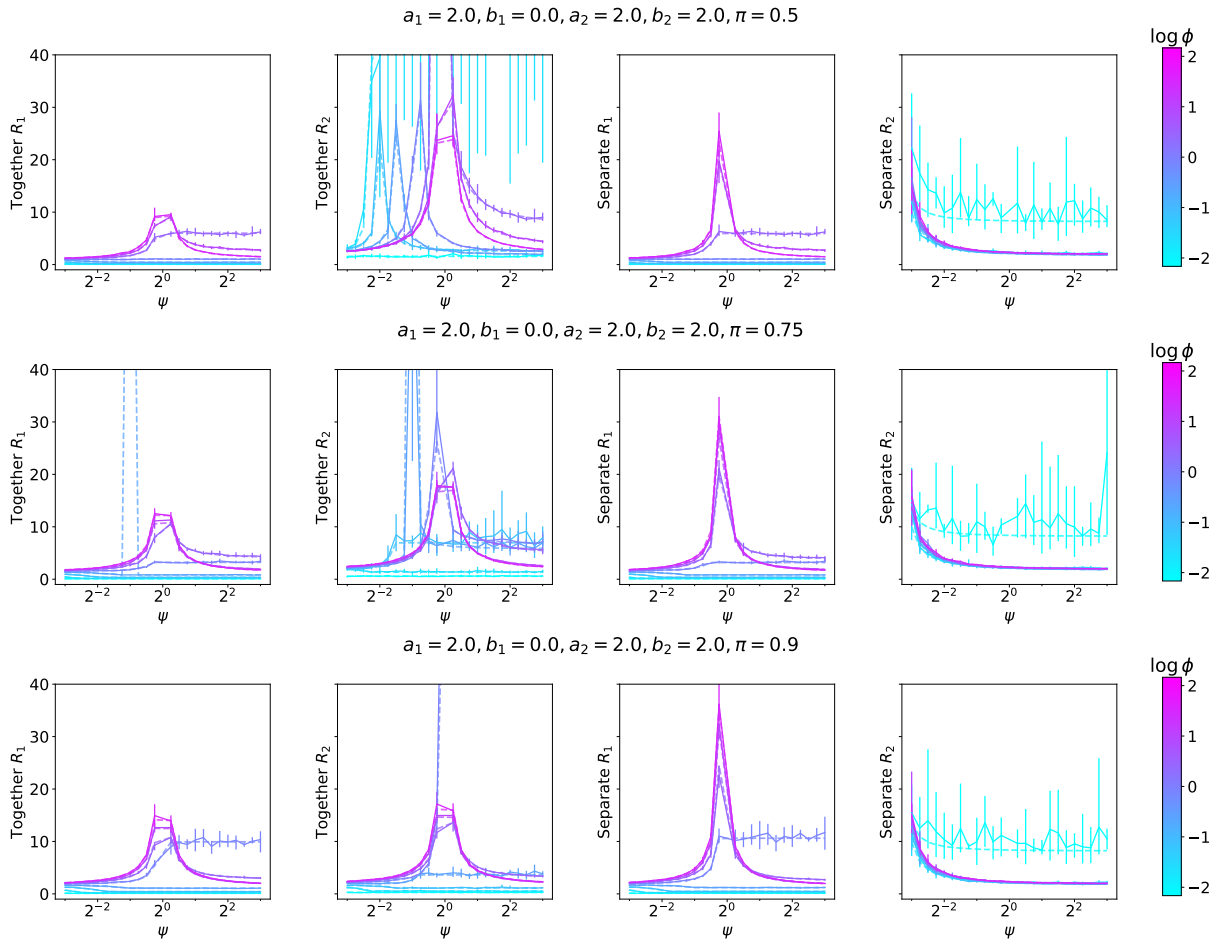


Figure E.10: We empirically demonstrate that minority-group bias is affected by extraneous features. We validate our theory (Theorems 6.3.1 and 6.3.2) for together R_1, R_2 (i.e., single model learned for both groups) and separate R_1, R_2 (i.e., separate model learned per group) under the setup described in Section 6.4.2. The solid lines capture empirical values while the corresponding lower-opacity dashed lines represent what our theory predicts. We include a black dashed line at $ADD = 1$ to contrast bias amplification vs. deamplification. All y-axes are on the same scale for easy comparison. The error bars capture the range of the estimators over 25 random seeds.

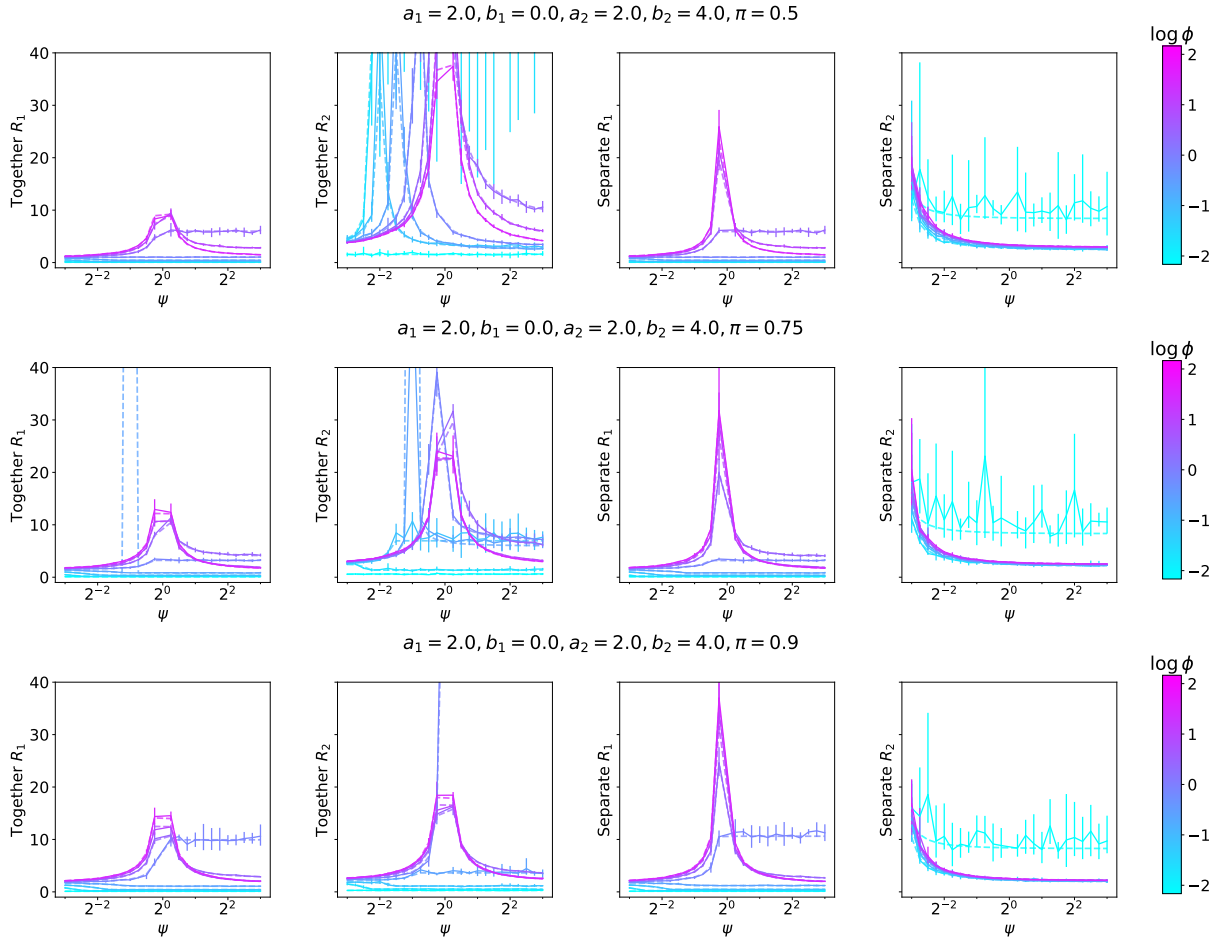


Figure E.11: We empirically demonstrate that minority-group bias is affected by extraneous features. We validate our theory (Theorems 6.3.1 and 6.3.2) for together R_1, R_2 (i.e., single model learned for both groups) and separate R_1, R_2 (i.e., separate model learned per group) under the setup described in Section 6.4.2. The solid lines capture empirical values while the corresponding lower-opacity dashed lines represent what our theory predicts. We include a black dashed line at $ADD = 1$ to contrast bias amplification vs. deamplification. All y-axes are on the same scale for easy comparison. The error bars capture the range of the estimators over 25 random seeds.

E.16 Actionable Insights from Theory

Searching for optimal hyperparameters. In practice, an optimal regularization penalty λ or training time t can be selected by searching for values that strike a desired balance between overall validation error (that is not too high) and bias amplification (that is not too high). As we would estimate the test error using the empirical validation error, we can estimate bias amplification using the validation set. Moreover, we would need to train: (1) the main model on a mixture of data from groups, and (2) auxiliary separate models on the data for each group.

However, it may be expensive to train auxiliary models for each candidate value of λ and t . The search space can be reduced by using insights from our theory. For instance, with overparameterization, as λ decreases (or t increases), bias amplification increases and plateaus, and with underparameterization, as λ decreases (or t increases), bias deamplification increases and plateaus. When the curves are monotone with respect to λ , the optimal λ is either at the left tail of the curve (e.g., $\lambda = 0$) or the right tail (i.e., the largest λ among the reasonable options). In contrast, Figure E.6 shows that when ψ is close to the interpolation threshold of 1, bias amplification is often not monotone with respect to λ .

Informing evaluation and mitigation strategies. Our theory offers avenues to assess whether a ML model trained on certain real-world data is prone to bias amplification and mitigate this amplification, even though we may lack direct access to population parameters like Σ . We can estimate such parameters using samples (e.g., $\widehat{\Sigma} = X^\top X$). However, even if we are unable to robustly estimate these parameters, our theory still provides valuable insights. For example, we observe that the ratios of parameters for the groups is often what matters, e.g., label noise ratio σ_2^2/σ_1^2 (see Section 6.4.1), ratio of covariance eigenvalues (see Appendix E.11). Thus, practitioners can use our theory to get intuition about when *disparities* in the variability of labels and features across groups can amplify bias.

Moreover, our findings warn against the conventional wisdom that increased model overparameterization or data balancing can alleviate bias issues. In addition, our theory informs criteria for feature selection (e.g., discarding features with disparate variance across groups) and warns ML practitioners about the interplay between high vs. low feature-to-sample regimes and overparameterization in inducing bias amplification. Nevertheless, additional work is required to make rigorous connections between our theoretical findings and better strategies for evaluating and mitigating the bias of models.