

UCLA

UCLA Electronic Theses and Dissertations

Title

Mathematical modeling of epidemics and adversarial learning in distributed systems

Permalink

<https://escholarship.org/uc/item/1j15n9k3>

Author

Li, Xia

Publication Date

2022

Peer reviewed|Thesis/dissertation

UNIVERSITY OF CALIFORNIA

Los Angeles

Mathematical modeling of epidemics and adversarial learning in distributed systems

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

by

Xia Li

2022

© Copyright by

Xia Li

2022

ABSTRACT OF THE DISSERTATION

Mathematical modeling of epidemics and adversarial learning in distributed systems

by

Xia Li

Doctor of Philosophy in Mathematics

University of California, Los Angeles, 2022

Professor Andrea L. Bertozzi, Co-Chair

Professor Deanna M. Hunter, Co-Chair

The COVID-19 epidemic has had a major global impact on humanity and the economy. Analyzing the effect of the COVID-19 pandemic can provide guidance for future pandemics. This dissertation studies three aspects of the epidemics during modern times. In the first part of the thesis, we study different aspects of pandemics along with mathematical models to address these aspects. Epidemics affect small communities in a different way than large urban centers. In chapter 2, we develop a mathematical model for finite-size effects using a stochastic compartmental susceptible–infected–recovered (SIR) model with a martingale formulation. The deterministic part coincides with the classical SIR model and we provide an upper bound for the stochastic part. Through analysis of the stochastic component depending on varying population size, we provide a theoretical explanation of finite size effects. Our theory is supported by numerical simulations of theoretical infinitesimal variance. In chapter 3, we propose a coupled model of policy-making and epidemic dynamics based on the SIR model and an optimization scheme. In chapter 4, we propose a hierarchical non-negative matrix factorization (NMF) scheme to classify the literature on COVID-19. We discover eight major latent topics and 52 granular subtopics in the body of

literature, related to vaccines, genetic structure and modeling of the disease and patient studies, as well as related diseases and virology.

Modern day machine-learning algorithms often operate in a distributed manner and are known to be vulnerable to adversarial attacks. Developing large-scale distributed methods that are robust to the presence of adversarial or corrupted workers is an important part of making such methods practical for real-world problems. In chapter 5, we propose novel methods that guarantees convergence and identify adversarial workers in highly hostile systems. The algorithm utilizes simple statistics (mode) to guarantee convergence and is capable of identifying the adversarial workers. Additionally, the efficiency of the proposed methods is shown in simulations in the presence of adversaries. The results demonstrate the great capability of such methods to tolerate different levels of adversary rates and to identify the adversarial workers with high accuracy.

The dissertation of Xia Li is approved.

Frederic R. Paik Schoenberg

Mason Alexander Porter

Deanna M. Hunter, Committee Co-Chair

Andrea L. Bertozzi, Committee Co-Chair

University of California, Los Angeles

2022

To my parents

TABLE OF CONTENTS

1	Introduction	1
1.1	Mathematical modeling of epidemics	1
1.2	Modeling of the COVID-19 literature	5
1.3	Adversarial learning in distributed systems	7
2	A Stochastic Compartmental Susceptible–Infected–Recovered (SIR) Model to Analyze Finite-size Effects in COVID-19 Case Studies	9
2.1	Background	9
2.2	Stochastic SIR-IPC model	12
2.2.1	Overview and notation	12
2.2.2	Small-time-interval probabilities for compartment fractions	14
2.3	Main results using a martingale approach	15
2.3.1	Martingale formulation of the stochastic SIR–IPC model	15
2.3.2	Estimates of the martingale variances	17
2.3.3	Deterministic analogue of the stochastic SIR model	19
2.3.4	Numerical simulations	19
2.4	Mathematical analysis of the <i>finite-size effects</i>	22
2.4.1	Scaling property of the stochastic component	22
2.4.2	Numerical simulations	24
2.5	Case studies of a small county and a cruise ship	28
2.6	Discussion	30

3	Multi-regional Policy-making and the Epidemics	33
3.1	Background	33
3.2	Policy model using optimal control	35
3.3	Models	36
3.3.1	The policy-incorporated SIR model	37
3.3.2	The policy model	37
3.3.3	The Policy model for a single region	39
3.3.4	The policy model for multi-layer multiple regions	40
3.4	Algorithms	42
3.4.1	Single region model	43
3.4.2	Multiple-layers multiple-regions model	44
3.5	Simulations	45
3.5.1	Optimal policy in France	45
3.5.2	Case study—2nd wave in Los Angeles	47
3.5.3	Multiple regions	51
3.6	Discussion and future work	53
4	COVID-19 Literature Topic-Based Search via Hierarchical nonnegative matrix factorization (NMF)	55
4.1	Background	55
4.1.1	Contributions	56
4.1.2	Related work	57
4.2	Data description	58
4.3	Hierarchical NMF for topic detection	58

4.3.1	NMF for topic detection	59
4.3.2	Hierarchical NMF	59
4.4	Discussion of results	60
4.4.1	Topic visualization	61
4.4.2	Discussion of topics	63
4.4.3	Topic coherence	64
4.4.4	Topic similarity	66
4.5	Implementation	69
4.5.1	Determining number of topics in each layer	70
4.5.2	Implementation of hierarchical NMF	71
4.6	Conclusion and future work	73
4.7	Appendix	73
4.7.1	Keywords removed	73
4.7.2	Topic keywords	74
5	Adversarial Learning in Distributed Systems	81
5.1	Background	81
5.1.1	Contribution	83
5.1.2	Related work	83
5.2	Method	84
5.3	Theoretical results	85
5.3.1	Mode distribution	85
5.3.2	Convergence without block-list	89
5.3.3	Block-list method	95

5.4	Simulations	97
5.5	Conclusion and future work	102
5.6	Appendix	103
5.6.1	Single row convergence without block-list	103
5.6.2	Finding the optimal d_0	107
5.6.3	Other proofs	109
6	Conclusion	110
	References	112

LIST OF FIGURES

1.1	An illustration of NMF.	6
2.1	Plots of infected compartment fractions $\mathbf{i}_N(t)$ for the stochastic SIR and $s(t), i(t)$, and $r(t)$ for deterministic SIR. For both models, $p_1 = 0.5, p_2 = 0.5, \gamma = 1, T = 23, s_0 = \mathbf{s}_{0N} = 0.96, i_0 = \mathbf{i}_{0N} = 0.04, r_0 = \mathbf{r}_{0N} = 0$. In Figs. 2.1b and 2.1a, 2.1d and 2.1c, 2.1f and 2.1e, and 2.1h and 2.1g, $\beta = 0.95, 1.1, 1.2, 1.3$, respectively. For the SIR-IPC model displayed in Figs. 2.1b, 2.1d, 2.1f, and 2.1h, $N = 10^{2.5}, 10^3, 10^{3.5}, 10^4$ in every panel.	21
2.2	Examples of the infinitesimal standard deviation for $\sigma_N^{(l)}(t), l = 1, 2, 3$. In both cases, $p_1 = 0.5, p_2 = 0.5, \gamma = 1, T = 23, s_0 = \mathbf{s}_{0N} = 0.96, i_0 = \mathbf{i}_{0N} = 0.04, r_0 = \mathbf{r}_{0N} = 0$. In the cases with $\mathcal{R}_0 = 0.95 < 1$ in Figs 2.2a, 2.2c, and 2.2e, $\beta = 0.95$. The blue, orange, green, and red lines show results with the same realization from simulations with the corresponding colors in Fig 2.1b. In the cases with $\mathcal{R}_0 = 1.3 > 1$ in Figs 2.2b, 2.2d, and 2.2f, $\beta = 0.95$. The blue, orange, green, and red lines show results with the same realization from simulations with the corresponding colors in Fig 2.1h.	26
2.3	Comparisons of the log–log plot of the theoretical and empirical scaling for $\overline{\mathcal{V}_N^{(l)}}$, $l = 1, 2, 3$, as in (2.38) - (2.40). The vertical bars are the error bars of the empirical scaling of the infinitesimal variance from the same realization as in Figs. 2.1b, and 2.1h. We set $[T_1, T_2]$ as $[0, 10], [1, 11], [2, 12], \dots$ and $[13, 23]$. The minimum and maximum values of $y = \log\left(\overline{\mathcal{V}_N^{(l)}}\right), l = 1, 2, 3$, taken over all such intervals are set as the upper and lower bounds of the error bars. The straight lines with slope -1 show the theoretical scaling with the x -intercepts as $\log(p_1\beta i_0 s_0), \log(p_1\beta i_0 s_0 + p_2\gamma i_0)$, and $\log(p_2\gamma i_0)$, for $l = 1, 2, 3$. In all figures, $p_1 = 0.5, p_2 = 0.5, \gamma = 1, T = 23, s_0 = \mathbf{s}_{0N} = 0.96, i_0 = \mathbf{i}_{0N} = 0.04$, and $r_0 = \mathbf{r}_{0N} = 0$. In Figs. 2.3a, 2.3c, and 2.3e, $\beta = 0.95$. In Figs. 2.3a, 2.3c, and 2.3e, $\beta = 1.3$	27

2.4	Comparison of field data (solid black line) for daily confirmed case percentages with 30 realisations of the stochastic SIR model for Churchill County, NV and the Diamond Princess Cruise Ship.	29
2.5	10-day moving average of the daily increased confirmed cases in 3 counties in California. Ventura’s peak is around end of July; San Luis Obispo’s peak is around mid August.	32
3.1	Timeline of COVID-19 Restrictions in France. Note the discrete nature of the restrictions, both in terms of the small number of categories and the fixed time intervals of enforcement.	37
3.2	An example of a three-layer hierarchical structure.	41
3.3	Different cost functions vs policy intensity α	43
3.4	Optimal policy and the SIR model of France from March 17 to May 11 2020.	47
3.5	The fraction of the infected and ‘stay-at-home’ policy over time in Los Angeles, San Francisco, and Orange County.	48
3.6	Optimal policy in Los Angeles with the basic reproduction number $R_0 = 2.5, 2.15$ and the fraction of the initial recovered population $r_0 = 0.1, 0.2, 0.3$	50
3.7	An example of three dependent counties without and with interventions. With intervention, for all counties, the coefficients for the implementation cost $\kappa = \frac{1}{2}$ and the coefficients for the impact cost $\eta = \frac{1}{2}$. The minimal policy time interval $\Delta t = 7$	52
3.8	An example with 3 counties and a governing state. For all counties, the coefficients for the implementation cost $\kappa = \frac{1}{6}$, the coefficients for the impact cost $\eta = \frac{1}{6}$ and the coefficients for the impact cost $\eta = \frac{1}{2}$ and the coefficients for the non-compliance cost $1 - \kappa - \eta = \frac{2}{3}$. For the state, the coefficients for the implementation cost $\kappa = \frac{1}{3}$, the coefficients for the impact cost $\eta = \frac{2}{3}$. The minimal policy time interval $\Delta t = 7$	53

4.1	Sunburst Diagram of the complete hierarchical structure. The top three relevant words per topic are shown. The area of each region is proportional to the number of articles in that topic. See appendix for the keywords associated with the third layer. The inner circle numeric labels are corresponding to topic number in Figure 4.2	62
4.2	Part of Topics from HNMF and related topic coherence: The first row shows the the key words for the topics in the first layer, the second row shows the subtopics of Topic 7 and the subtopics of Topic 7-1 is showed in row 3. Corresponding <i>topic coherence</i> score (see Section 4.4.3 for more details) is underneath each word cloud.	63
4.3	Topic similarity for all the topics from HNMF measured by WDM . A dark color indicates the topics are dissimilar, while a light color indicates high similarity. Note that the topics are listed from first layer to third layer from top to bottom or right to left on the vertical and horizontal axes, respectively.	68
4.4	Topic similarity between Topic 7 and its subtopics measured by WDM: Topic 7 has high topic similarity with its five subtopics (7-1, 7-2, 7-3, 7-4, 7-5) and the five topics have low similarity between themselves.	69
4.5	Topic similarity between Topic 7-1 and its subtopics measured by WDM: Topic 7-1 has high topic similarity with its four subtopics (7-1-1, 7-1-2, 7-1-3, 7-1-4) and the four topics have low similarity between themselves.	69
4.6	Plot of marginal increment in proportion of variance explained by adding another cluster to split X . It is determined that the ideal number of clusters/topics likely lies in the range $[7, 11]$, as this is where the plot starts to level off.	71
4.7	Box plot of LSS_k : Topic number 8 is the “best” as it has the highest median lss (least seed similarity) score and should be expected to yield consistent results with random seeds.	72

5.1	Effects of the number of used rows d_0 on convergence: $N_r = 20, n_r = 4, k = 3, \ e\ _\infty = 10^{-3}$. The error norms were averaged over 50 trials (the solid lines) with 90% percentiles (the shaded areas).	98
5.2	Effects of different data sizes d_0 on convergence: $N_r = 20, n_r = 4, k = 3$, and $\ e\ _\infty = 500$. The error norms were averaged over 50 trials (the solid lines) with 90% percentiles (the shaded areas).	99
5.3	Effects of the number of used workers n_r : $k = 3, d_0 = 6, N_r = 20, \ e\ _\infty = 5 \times 10^2$.	100
5.4	Effects of the number of categories k . $N_r = 20, n_r = 4, d_0 = 4, \ e\ _\infty = 5 \times 10^2$	101
5.5	Effects of the adversarial rate p , $d_0 = 3, N_r = 20, n_r = 4$, and $k = 3$. Squared error norms were averaged over 50 trials with the 90% percentiles	101
5.6	Effects of number of used row d_0 using the Breast Cancer Wisconsin data set, $N_r = 10, n_r = 4, k = 3, \ e\ _\infty = 500$.	102

LIST OF TABLES

3.1	CDC stay-at-home policies. There are 6 levels of policies and we map the levels linearly onto the interval $[0, 1]$ for simplicity. The numerical value on the left is used to graph actual policies over time in Fig. 3.5b.	48
4.1	The coherence scores based on both the C and C_V metric for each of the 8 topics in the first layer of the tree.	66
4.2	The top 10 keywords associated with each of the 8 topics in the first layer of the tree	75
4.3	The top 10 keywords associated with each of the subtopics of Topic 1 in the 2^{nd} layer of the tree and the C_V coherence score	75
4.4	The top 10 keywords associated with each of the subtopics of Topic 2 in the 2^{nd} layer of the tree and the C_V coherence score	76
4.5	The top 10 keywords associated with each of the subtopics of Topic 3 in the 2^{nd} layer of the tree and the C_V coherence score	76
4.6	The top 10 keywords associated with each of the subtopics of Topic 4 in the 2^{nd} layer of the tree and the C_V coherence score	77
4.7	The top 10 keywords associated with each of the subtopics of Topic 5 in the 2^{nd} layer of the tree and the C_V coherence score	77
4.8	The top 10 keywords associated with each of the subtopics of Topic 6 in the 2^{nd} layer of the tree and the C_V coherence score	78
4.9	The top 10 keywords associated with each of the subtopics of Topic 7 in the 2^{nd} layer of the tree and the C_V coherence score	78
4.10	The top 10 keywords associated with each of the subtopics of Topic 8 in the 2^{nd} layer of the tree and the C_V coherence score	79

4.11	The top 10 keywords associated with each of the subtopics of Topic 7-1 in the 3 rd layer of the tree and the C_V coherence score	80
4.12	The top 10 keywords associated with each of the subtopics of Topic 7-4 in the 3 rd layer of the tree and the C_V coherence score	80
5.1	Notation table	87
5.2	Total number of workers $N = 10$, number of error categories $k = 3$	95
5.3	Conditional probability of being in block-list.	96
5.4	Accuracy of the method with the block-list when number of iterations to update the block-list $S = 200, 500, 1000, 2000$, $k = 3$, $N_r = 20$, and $n_r = 4$	102
5.5	Total number of workers $N = 100$, number of chosen workers $n = 5$	105
5.6	Total number of workers $N = 100$, number of error categories $k = 5$	105

ACKNOWLEDGMENTS

I would like to thank my advisors Prof. Andrea Bertozzi and Prof. Deanna Needell (Deanna M. Hunter) for their supervision and support during my studies. I thank Andrea Bertozzi for giving me opportunities to work on various social-science-related projects where I found my passions. I am thankful to the members of Deanna Needell's research group for creating such an inspiring and supportive environment.

I would like to thank my committee members Mason Porter and Frederic (Rick) Paik Schoenberg for their valuable suggestions and insights. I would also like to thank all my collaborators for their hard work and enthusiasm. I would like to pay special regards to P. Jeffrey Brantingham and Longxiu Huang. I thank P. Jeffrey Brantingham for bringing humanity aspects to every project we collaborated and inspired me with wonderful ideas. I thank Longxiu Huang for her guidance, patient and support. I would like to thank Maida Bassili, Brenda Buenrostro, Martha Contreras, and Corinne Smith for providing administrative support.

This work is supported by NSF grant DMS-2027438, ARO MURI Grant W911NF1810208, NSF grant DMS-2027277, Simons Math + X Investigator Award number 510776, NSF BIGDATA DMS #1740325 and NSF DMS #2011140.

Thank you to all of my friends. I would like to thank to Kaiyan Peng, Joel Barnett, Dominic Yang, and Siting Liu for friendship, support and academic discussion. You made my Ph.D. journey much more fun. I thank Bon-soon Lin for helping with the basic qualifying exam. I thank Bohan Chen, Jacob Moorman, Yoni Dukler, Bohyun Kim, Yurun Ge, Ziheng Ge, Xuchen Han and many other Ph.D. students at UCLA for their help and fun times. I would like to also acknowledge the support and great love from my family, and my friends Grace, Henglin, Jordan, Jinglei and Paige. I love you all.

VITA

2014–2018 B.S. Mathematics and Applied Mathematics, Sichuan University.

2018–2020 M.A. Mathematics, UCLA.

2018–2020 Teaching Assistant, Mathematics Department, UCLA.

2020–Present Graduate Student Researcher, Mathematics Department, UCLA.

PUBLICATIONS

X. Li, C. Wang, H. Li, and A. Bertozzi, A Martingale Formulation for Stochastic Compartmental Susceptible–Infected–Resistant (SIR) Models to Analyze Finite-size Effects in COVID-19 Case Studies, *Networks & Heterogeneous Media*, 17(3), pp. 311-331, 2022.

X. Li, L. Huang, D. Needell, Distributed Randomized Kaczmarz for Adversarial Workers, *Proceeding 53rd Asilomar Conference on Signals, Systems and Computers*, Pacific Grove, CA, Nov. 2021.

R. Grotheer, K. Ha, Y. Huang, P. Li, X. Li, L. Huang, D. Needell, E. Rebrova, A. Kryshchenko, O. Kryshchenko, COVID-19 Literature Topic-Based Search via Hierarchical NMF, In *Proceedings of the 1st Workshop on NLP for COVID-19 (Part 2) at EMNLP 2020*, Online. Association for Computational Linguistics.

M. Alaverdian, W. Gilroy, V. Kirgios, X. Li, C. Matuk, D. Mckenzie, T. Ruangkriengsin, A. Bertozzi, and J. Brantingham, Who killed Lilly Kane? A Case Study in Applying Knowledge Graphs to Crime Fiction (Proc. GTA³ 4.0 at IEEE International Conference on Big Data, Atlanta, GA, 2020)

CHAPTER 1

Introduction

1.1 Mathematical modeling of epidemics

The novel strain of coronavirus, SARS-CoV-2, was first identified in Wuhan (Hubei) in December 2019, and soon after that, the number of infections grew exponentially. Despite the measures taken to contain the outbreak, a COVID-19 pandemic ensued and had spread worldwide by early April 2020. The COVID-19 pandemic provides the first instance in 100 years to study the worldwide dynamics of a novel virus outbreak. Early on in the pandemic, we observed many localized outbreaks on cruise ships, skilled nursing facilities, and through other congregating mechanisms such as conferences, parties, and places of worship. Several groups developed models to inform public health, a notable one being the March 16, 2020 report by Imperial College [\[FLN20\]](#) that forecast over 2 million deaths in the US and over 500,000 deaths in the UK if non-pharmaceutical interventions (NPIs) were not implemented. Within days, lockdowns were enforced in both countries. Over a year since that time, we have observed a variety of patterns and some of them did not follow the trends predicted by the early reports. The most obvious reason for these differences is the reduction of the reproductive number due to NPIs. An additional important factor is the departure from simple mixture models that results from isolating populations because of NPIs or simply because of geography. One way to model such effects is to consider an agent-based model with a finite population size in which the role of stochasticity is more pronounced due to the population size. This is the focus of our work in [chapter 2](#).

Mathematical modeling has an important role in studying infectious diseases on both long and

short timescales—for modern viruses such as AIDS, Severe Acute Respiratory Syndrome (SARS), Zika, and the novel coronavirus disease 2019 (COVID-19). It helps to understand how diseases spread and are controlled, providing insights on decision-making. For epidemic modeling, the standard approach involves compartmental models in which the population is divided into compartments representing individuals in one of several states, e.g. the susceptible (S), infectious (I), recovered (R) and exposed (E). This mathematical framework leads to a variety of epidemic models, such as SIR, SIS, SEIR, SIRS (see e.g. [LTH18, LMR19, KCM05, BRQ20]). Compartmental modeling has been applied to epidemic study of COVID-19, HIV, etc [TT94, Llo01, BFM20, Bai75, Bar60]. These models are generalized to agent-based forms and the most granular model among them being the stochastic models. The classic compartmental models converge to an endemic equilibrium while their stochastic counterparts usually do not. Other characteristics of the stochastic epidemic models include the probability of an outbreak, the final size distribution of an epidemic and the expected duration of an epidemic. Due to these properties of the stochastic models, in chapter 2, we adopt a martingale approach of a stochastic SIR model. We quantitatively analyze the finite-size effects under small populations with the martingale formula applicable to all population scales. We decompose the process into a deterministic component and a stochastic component, corresponding to the infinitesimal means and variances of the process, respectively. The deterministic component leads to a fluid-limit-type continuum analogue, coinciding with the deterministic SIR model (rescaled over the population size N) as in [KMW27]. Our simulations show that as N decreases, our stochastic SIR Model deviates further from the deterministic SIR model, and finite-size effects arise. The stochastic component is bounded by quantities of the same order of magnitude of $1/N$. We find that the stochastic component scales as $1/\sqrt{N}$. Furthermore, the smaller the population size is, the larger the deviation of the stochastic simulations from the deterministic ones are. We provide a theoretical estimation for this scaling, building on the analysis of continuous-time Markov processes. The output of numerical experiments of theoretical infinitesimal variances supports our theory.

Another line of modeling the epidemic is to use a self-exciting point process, also known as the Hawkes process (HP). With the excitation representing the transmission of disease, point processes can model the clustered nature of infections. Lewis et al. in [LM11] and Bertozzi et al. in [BFM20], applied Hawkes process to epidemic forecasting. Additionally, the multivariate Hawkes Process (MHP) utilizes a cross-excitation matrix which can be used to model the network structure of sub-processes. [MBX21] studied the transmission of COVID-19 among different age groups using MHP with a Bayesian approach. The basic reproduction number is a crucial component in measuring the spread of the pandemic. [BAC20] adopt a signal-processing approach to compute the effective reproduction number (ERN) from the daily count of newly detected cases. [MSS20] apply Cox Hawkes, a variant of the Hawkes process, to model the dynamic of ERN and analyze the impact of public health interventions on the transmission dynamics of COVID-19. More recently, the multivariate Hawkes process has been used to inform COVID-19-related decision-making in New Jersey (in collaboration with Facebook AI)[New]. The Hawkes process has no upper limit for the number of events that may occur. To address this of the original Hawkes process, in [RMK18], the author introduced a parameter N , denoting the finite total size of the population and modulated the event rate by the available population. The resulting model is referred to HawkesN. This model captures the long-term evolution of the pandemic and the authors showed a theoretical connection between HawkesN and the classic compartmental models [RMK18].

Modeling disease control and intervention is an important topic in modeling epidemics. In the course of battling COVID-19, policies provide guidelines that serve an important role in slowing down the spreading. Some of them include ‘safer-at-home’, ‘Maintaining 6 feet distance’ and ‘mask wearing’. In the absence of reliable pharmaceutical interventions like vaccination, these public health strategies are crucial. The timeline of COVID-19 globally and locally ([Cena, Wika]) indicates that the evolution of policies affected the evolution of the pandemic and vice versa. A majority of them leverage the SIR model and its variants. For example, in [KZR21], the authors proposed an SIR-based model that captures the effects of intervention policies on the disease spread parameters by leveraging intervention policy data along with the reported disease cases. The model

is also designed an observation mechanism to account for under-reporting by adding two new compartments. The model learns the spread, policy, and reporting parameters end-to-end directly from observed data via gradient-based training. In [BF07], the authors used an SEIR-based model and discovered that the timing of public health interventions had a profound influence on the pattern of the autumn wave of the 1918 pandemic in different cities. Bliman et al. in [BDP21] proposed an optimal control framework to find the optimal policy that minimizes the final pandemic size. In their work, an optimal policy should stop the disease as close as possible after crossing the herd immunity threshold. The authors proved the existence and the uniqueness of the solution and showed the optimal social distancing policy is a bang-bang controller (a control that switches from one extreme to the other). In their work, the admissible set of the policy function is a subset of all continuous functions. However, in real life, policies executed have different intensities and duration times. For example, in the county of Los Angeles, on March 21, 2020, social distancing was first suggested, a month after the first COVID 19 case in LA [Dep22]. Around that time, a health office order ‘*safer-at-home*’ was also released. One week later, beaches, hiking trails, dog parks, skate parks, etc., and more public sites and facilities were temporarily closed. As infected cases continued to increase, a month later, on May 1, facial coverings were suggested. Considering the policies executed have different intensities and duration times, in chapter 3, we model policies as piece-wise linear functions in time that only takes values from a finite set. We couple the policies with a policy-Incorporated SIR model for single-region case and a network styled policy-Incorporated SIR model for the multi-regional case. We further assume that policies function have a minimal policy time interval (MPTI) during which the policy stay constant in order to mimic the reaction time of decision-makings of the regions. We reproduced the results of France as in [BDP21] and discussed the second wave in the county of Los Angeles. In addition, we proposed a generalization of the policy-making SIR model with multiple-agents in a hierarchical structure and presented an example of three interacting counties with and without a governing state.

1.2 Modeling of the COVID-19 literature

Another aspect of studying COVID-19 is to provide an effective method to navigate COVID-19-related literature. As of July 21, 2020, there were over 196,630 COVID-19-related scholarly articles on PubMed, PubMed Central, bioRxiv and medRxiv preprint servers [ASM21]. There are several reasons for such navigation: to organize and coordinate the literature and to help explore research topics addressed. In chapter 4, we build an interactive search engine for COVID literature using a hierarchical non-negative matrix factorization (HNMF) method. The HNMF is a hierarchical version of a method called non-negative matrix factorization (NMF) and NMF is a classic method to classify the topics of articles. It is a matrix factorization method that decomposes a given matrix $X \in \mathbb{R}^{n \times d}$ into two low-rank, non-negative matrices $W_+ \in \mathbb{R}^{n \times r}$ and $H_+ \in \mathbb{R}^{r \times d}$ for some r to be specified (Fig. 1.1). For simplicity, we omit $+$ in the discussion. To find the W and H , we solve for the minimization problem:

$$\min_{W,H} \|X - WH\|_F^2. \quad (1.1)$$

To format the narratives into inputs for our topic modeling methods, we used a bag-of-words model, which takes in the corpus and creates a vocabulary out of each unique word in the corpus. It then models each document as a vector with length equal to the number of words in the vocabulary where entry i in the vector corresponds to how many times word i occurred in the document. Thus, the bag-of-words model gives us a matrix X with dimension $n \times d$ where n is the number of words in the vocabulary, and d is the number of documents. The parameter r represents the number of the topic, which needs to be determined. For the matrix W , the words per topic matrix, each column of the W represents a word distribution of a topic. For the matrix H , the topic per document matrix, i th column represents the topic distribution of i th document. In HNMF, NMF is first applied to the original corpus matrix X to obtain the dictionary matrix W and coding matrix H . The documents are then sorted into matrices X_1, X_2, \dots, X_k , each representing a different topic, according to the coding matrix H . Then NMF is applied to each one of the matrices X_1, X_2, \dots, X_k to obtain the

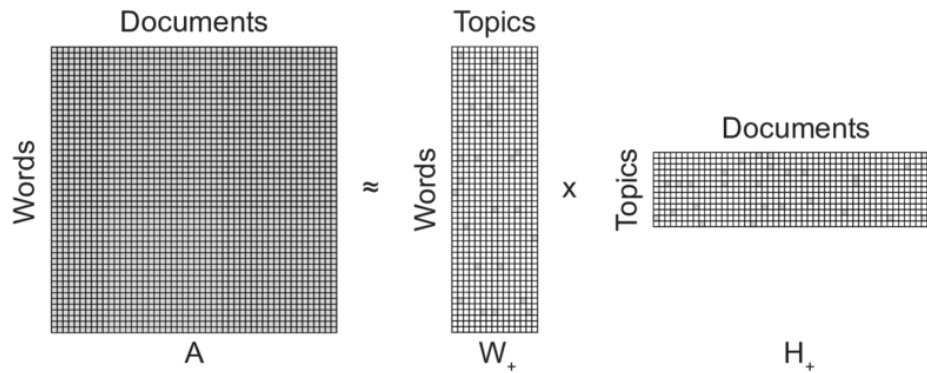


Figure 1.1: An illustration of NMF.

matrix for subtopic-super-topic relations and the corresponding coding matrix. The above process is repeated until the desired number of layers of the hierarchy is reached.

The traditional NMF method treats the detected topics as a flat structure, which limits the ability of the representation of such method. In contrast, a hierarchical NMF (HNMF) framework is able to detect supertopics, subtopics, and the relationship between them, creating a tree structure. In chapter 4, the interactive and hierarchical structure facilitates the search by researchers and is available to use through a corresponding website. The topics discovered by HNMF reveal that early research of interest to the COVID-19 research community divides into diverse areas such as research related to other corona viruses, research related to other respiratory diseases, virology and genetic research, as well as research relating to the public health response. A topic coherence metric reveals that the topics discovered are consistent and semantically meaningful, while a topic similarity metric reveals that the topics differ sufficiently from one another to allow a diversity of choice and areas of interest on the user's part.

1.3 Adversarial learning in distributed systems

The last section of this thesis is on a different topic. As machine-learning algorithms gain popularity in industrial applications, it is important to study methods to improve their robustness and protect them from adversarial attacks. There are several ways a machine learning model can be attacked, including evasion attacks [GMP18], data poisoning attacks [GFH20], model extraction [WXG21, KTP19]. Large-scale machine-learning problems are typically run on distributed systems. A typical attack in the distributed setting is the Byzantine attack [LSP82]. The Byzantine system refers to a computer system where every single component (also known as ‘worker’ or ‘node’) communicates with each other, and some send conflicting information. There have been works on stochastic optimization in a Byzantine setting. For example, in [EFG20], the authors studied the problem in an adversarial setting where, out of the m machines which allegedly compute stochastic gradients every iteration, a fraction of workers are Byzantine, and can behave arbitrarily and adversarially. They showed that the proposed algorithm is information-theoretically optimal both in terms of sampling complexity and time complexity. In [FXM14], a framework called the distributed robust learning (DRL) is proposed. This framework can reduce the computational time of traditional robust learning methods by several orders of magnitude. The authors showed that DRL not only preserves the robustness of the base robust learning method but also tolerates contaminations on a constant fraction of results from computing workers (node failures). In chapter 5, we focus on developing robust algorithms when the workers are adversarial and detect the adversarial workers.

A special application of solving optimization problems in the distributed systems is the least squares problem. In chapter 5, we consider solving the over-determined linear system $Ax = b$, where $A \in \mathbb{R}^{d_1 \times d_2}$, $b \in \mathbb{R}^{d_1}$ and $x \in \mathbb{R}^{d_2}$. This problem can be reformulated as the least squares problem $\min_x \|Ax - b\|_2^2 = 0$. One way to solving linear systems in an iterative manner is the Kaczmarz method. It was first proposed by [Kac37] which is also known under the name *Algebraic Reconstruction Technique* (ART) in computer tomography [GHJ75, HM93, Nat01] and has found various applications ranging from computer tomography to digital signal processing. Later Strohmer

et al. [SV09] proposed a randomized version of Kaczmarz, where the probability of each row being selected is set to be proportional to the Euclidean norm of the row and proves the exponential bound on the expected rate of convergence. Typically, when the dimension of the data matrix is large and cannot be loaded into memory at once, the direct methods like using the pseudo-inverse can be infeasible or expensive, then randomized methods such as the randomized Kaczmarz (RK) [SV09] are more effective. Different ways to use RK in a distributed setting can lead to different communication costs, storage overhead and possibly different convergence speeds. For example, a complete redundancy, i.e., every worker holds a complete copy of the data, is useful when a majority of workers are not reliable. However, such a scheme has a large storage overhead. Another example is partial redundancy, i.e., each worker holds part of the data, and the central server needs to aggregate their computation results to obtain an update. This is generally more effective and efficient. The RK operates on single rows of the matrix A at a time. While RK randomly selects a row of A to work with, Motzkin’s Method (MM) [Agm54] employs a greedy row selection. In [HM21], the authors proposed a hybrid algorithm based on these two algorithms: the sampling Kaczmarz-Motzkin (SKM) algorithm, which samples a random subset of β rows of the data matrix A and then greedily selects the best row of the subset. In chapter 5, we adopt a similar idea. We consider that the data is initially distributed to workers with some redundancy, and each worker holds multiple rows. At each iteration, the central worker sample a random subset of d_0 rows and a random subset of the workers who hold those rows to compute the residual in RK (Alg. 1). We utilize the simple statistic (**mode**) to avoid using results from the adversarial workers. We provide numerical and theoretical analysis to show the effectiveness of our method.

Algorithm 1 Randomized Kaczmarz Algorithm

- 1: Select a row index $i_j \in [d_2]$ with probability $p_{i_j} = \frac{\|A_{i_j}\|_2^2}{\|A\|_F^2}$
 - 2: Update $x_{j+1} = \arg \min_{x \in \mathbb{R}^{d_1}} \|x - x_j\|$ s.t. $A_{i_j} x_{j+1} = b_{i_j}$
 - 3: Repeat until convergence
-

CHAPTER 2

A Stochastic Compartmental Susceptible–Infected–Recovered (SIR) Model to Analyze Finite-size Effects in COVID-19 Case Studies

This chapter is adapted from the original paper [LWL22] that I co-authored with Chuntian Wang, Hao Li and Andrea Bertozzi. I was first author on that paper and contributed to the modeling and all of the numerical experiments in the paper.

2.1 Background

The standard approach to model epidemics involves compartmental models in which the population is divided into compartments representing individuals in one of several states, e.g. the susceptible (S), infectious (I), recovered (R) and exposed (E). This mathematical framework can lead to a variety of epidemic models, such as SIR, SIS, SEIR, SIRS ([LTH18, LMR19, KCM05, BRQ20]). Compartmental modeling has been applied to epidemic study of COVID-19, HIV, etc [TT94, Llo01, BFM20, Bai75, Bar60]. These models can have an agent-based form with the most granular being the stochastic models. Because of the granularity of the stochastic compartmental model, its fluid limit is often adopted to reduce computational complexity. By the law of large numbers, stochastic compartment models with Markov processes can be approximated by their deterministic ODE counterpart. The central limit effect arises under the diffusion scaling, and the fluid-scaled

model converges to a Gaussian diffusion of stochastic differential equations about the deterministic solution. All these approximation methods require a large population size to control variances. As the population size decreases the fluid-limit approximation acquires larger stochastic variances, and *finite-size effects* may arise. For example, empirical variances of realizations of stochastic models can be at the same magnitude of the inverse of the total population size ([Ish91, Ish93]). Since population scale within a relatively small community (like small counties, cruise ships) is very likely to fall within the regime of considerable stochastic deviations, quantitative study and analysis of *finite-size effects* is relevant to real disease statistics. To this end, we investigate stochastic variability in the stochastic SIR model driven by independent Poisson clocks (IPC), which will be referred to as the SIR-IPC Model. We adopt Poisson clocks, rather than time steps discretized with a fixed duration. A more realistic model should treat all events as occurring independently, according to their own stochastic clock, rather than occurring at regular intervals, and arriving according to the same schedule. Time steps are turned into exponentially distributed random variables with the introduction of Poisson clocks. Moreover, independent Poisson clocks allow us to use the theory of continuous-time Markov pure jump processes, e.g. a martingale approach (see e.g. [Dur96, Kun86, Lig80, Lig85, Lig10]). And this leads to a martingale formulation.

In prior literature, only the components of the martingale formulation were derived (see e.g. [All08, All11, Yan08]), while the study of the complete formulation has been lacking. Specifically, theoretical infinitesimal variances have also been studied in [BIH18, LMR19, All11, Yan08, YC19]. However, their main purpose is to explore ways to approximate stochastic compartmental models with a large population size, such as fluid limit, diffusion limit, and linear SDE approximations. With the martingale formula applicable to all population scales, we can quantitatively analyze the *finite-size effects* under small populations. It is possible to merge independent Poisson clocks to obtain systems that are driven by only one or two Poisson clocks. And this idea is worked out in [All08, All11, Yan08]. There the SIR models are driven by two Poisson clocks, one for transitions from S to I compartment, the other for those from I to R compartment. In our model, we assume that individual events each occur at a certain Poisson rate. These events include infectious contacts

of each pair of susceptible and infected individuals, and recovery of each infectious individuals. In the simplest models, births and deaths (due to natural causes) during the course of the epidemic are not taken into account. In more detailed models, births and deaths (due to natural causes) are included, and these are called “vital dynamics” in [Het00]. We assume the absence of vital dynamics and a deterministic population size parameter N . With the martingale approach, we decompose the process into a deterministic component and a stochastic component, corresponding to the infinitesimal means and variances of the process, respectively. The deterministic component leads to a fluid-limit-type continuum analogue, which coincides with the deterministic SIR model (rescaled over N) as in [KMW27]. Computer simulations show that as N decreases, the SIR-IPC model deviates further from the deterministic SIR model and *finite-size effects* arise. The stochastic component is bounded by quantities of the same order of magnitude as $1/N$. We find that the stochastic component scales as $1/\sqrt{N}$. This is not a surprise, as it is expected based on the law of large numbers that the standard deviation scales like $1/\sqrt{N}$.

Furthermore, the smaller the population size is, the larger the deviation of the stochastic simulations from the deterministic ones are. We provide a theoretical estimation for this scaling, building on the analysis of continuous-time Markov processes. The output of numerical experiments of theoretical infinitesimal variances supports our theory. Since the outbreak of COVID-19 on the Diamond Princess Cruise Ship, at least 25 other such vessels have confirmed COVID-19 cases and studies of the transmission of the disease on cruise ship have drawn consideration attention (e.g. [AKL21]).

The chapter is organized as follows. In Section 2.2, the stochastic SIR Model is introduced. The martingale formulation is derived in Section 2.3. Based on the martingale formulation, a deterministic and continuum analogue of the stochastic SIR model is derived (Section 2.3.3). Simulations are run to compare the stochastic SIR model to its fluid limit, illustrating the importance of *finite-size effects* (Section 2.3.4). In Section 2.4 the *finite-size effects* are analyzed theoretically based upon the martingale formulation. The theory is supported by simulations and field data in Section 2.5.

2.2 Stochastic SIR-IPC model

2.2.1 Overview and notation

Our setting is similar to the continuous-time-Markov-processes stochastic SIR models in [All08], [RMK18] and [Yan08]. We assume a continuous time variable $t \in [0, \infty)$, and a probability space $(\Omega, \mathcal{F}, \mathbb{P})$, and a deterministic and fixed population size N . The classical stochastic SIR model without vital dynamics consists of three compartments: susceptible, infected, and recovered individuals. For $\omega \in \Omega$, we assume that at a given time t the numbers of the compartments are $\mathbf{S}_N(\omega, t)$, $\mathbf{I}_N(\omega, t)$, and $\mathbf{R}_N(\omega, t)$, and $N = \mathbf{S}_N(\omega, t) + \mathbf{I}_N(\omega, t) + \mathbf{R}_N(\omega, t)$.

For simplicity of modeling, we assume that there is no vital dynamics, i.e., we view deaths as a subset of recovered individuals and humans' natural resistance to the disease does not introduce new susceptible people after recovery. Below we assume that all the parameters are independent of N , unless otherwise specified. We assume that the initial data are deterministic and denoted as

$$(\mathbf{S}_N(\omega, 0), \mathbf{I}_N(\omega, 0), \mathbf{R}_N(\omega, 0)) = (S_{0N}, I_{0N}, R_{0N}). \quad (2.1)$$

A Poisson clock governs the infectious contact of any pair of a susceptible and an infectious individuals. There are in total $\mathbf{S}_N(\omega, t)\mathbf{I}_N(\omega, t)$ such pairs at time t . Each individual Poisson clock advances according to a Poisson process with rate β/N . We denote this clock as the “i-clock”. Suppose that the i-clock advances at time t^- . At time t , with probability $p_1 \in [0, 1]$ the susceptible of the pair associated with the clock contracts the disease and transitions to the infectious compartment through contact. Another Poisson clock governs the recovery of the infectious individuals. Each infectious individual is assigned a Poisson clock, denoted as “r-clock”, which advances with rate γ . Suppose that the “r-clock” advances at time t^- . At time t , the infectious individual with the clock is recovered, with probability $p_2 \in [0, 1]$. All the Poisson clocks are exponentially distributed independent random variables and are independent with time increments.

Similar to the classic SIR, βp_1 denotes the transmission rate constant, and γp_2 is denotes

recovery rate constant. Thus the basic reproduction number \mathcal{R}_0 for our model is

$$\mathcal{R}_0 = \frac{\beta p_1}{\gamma p_2}. \quad (2.2)$$

Remark 2.2.1. *To facilitate the computation of the SIR-IPC model, we show the following equivalent description of the SIR-IPC model with merged Poisson clocks.*

All the i -clocks can be merged into one master i -clock to govern the arrivals of infectious contact. Once the master i -clock advances, with probability p_1 a uniformly and randomly chosen susceptible individual contracts the disease. The master i -clock advances with rate $\beta \mathbf{I}_N(\omega, t^-) \mathbf{S}_N(\omega, t^-) / N$. Likewise all the Poisson r -clocks can be merged into one master r -clock as we do for the i -clocks. Once the master r -clock advances, one uniformly and randomly chosen infectious individual is recovered, and with probability p_2 develops immunity and becomes disinfected. The master r -clock advances according to a Poisson process with rate $\gamma \mathbf{I}_N(\omega, t^-)$.

In general, theory regarding the merging and splitting of Poisson processes (see e.g. [Dur99]) implies that independent Poisson clocks can be treated as one merged Poisson process; and with probability in proportion to the rate of each Poisson clock, the Poisson clocks compete to advance first in the merged Poisson process. More precisely, suppose that $X_1(t), X_2(t), \dots, X_N(t)$ denote the numbers of arrivals corresponding to independent Poisson processes with arriving rates $\lambda_1, \dots, \lambda_N$, respectively. Then $\sum_{n=1}^N X_n(t)$ has the same distribution as the number of arrivals corresponding to one Poisson process with arriving rate $\sum_{n=1}^N \lambda_n(t)$. Furthermore, any particular arrival of the merged process has probability $\lambda_n / \sum_{n=1}^N \lambda_n(t)$ of originating from the n -th process, for $n = 1, 2, \dots, N$, independent of all other arrivals and their origins. The variance of the merged process at time t can be derived as follows:

$$\text{Var} \left(\sum_{n=1}^N X_n(t) \right) = \sum_{n=1}^N \text{Var} (X_n(t)) = \sum_{n=1}^N \lambda_n t = t \sum_{n=1}^N \lambda_n. \quad (2.3)$$

When $\lambda_n \equiv \lambda$ for all n , we have

$$\text{Var} \left(\sum_{n=1}^N X_n(t) \right) = tN\lambda. \quad (2.4)$$

2.2.2 Small-time-interval probabilities for compartment fractions

We denote the fractions of each compartment as

$$(\mathbf{s}_N(\omega, t), \mathbf{i}_N(\omega, t), \mathbf{r}_N(\omega, t)) := \left(\frac{\mathbf{S}_N}{N}(\omega, t), \frac{\mathbf{I}_N}{N}(\omega, t), \frac{\mathbf{R}_N}{N}(\omega, t) \right), \quad (2.5)$$

with the initial condition of population fractions denoted as

$$(\mathbf{s}_N(\omega, 0), \mathbf{i}_N(\omega, 0), \mathbf{r}_N(\omega, 0)) := (\mathbf{s}_{0N}, \mathbf{i}_{0N}, \mathbf{r}_{0N}). \quad (2.6)$$

This together with (2.1) implies that

$$(\mathbf{s}_{0N}, \mathbf{i}_{0N}, \mathbf{r}_{0N}) = \left(\frac{S_{0N}}{N}, \frac{I_{0N}}{N}, \frac{R_{0N}}{N} \right). \quad (2.7)$$

With similar idea as in e.g. [All08], we derive the small-time-interval probabilities for compartment fractions. For Δt a short time interval, we have

$$\begin{aligned} \mathbb{P}(\Delta \mathbf{s}_N(\omega, t), \Delta \mathbf{i}_N(\omega, t)) &= (k, j) | (\mathbf{s}_N(\omega, t), \mathbf{i}_N(\omega, t) = (s, i)) \\ &= \begin{cases} p_1 \beta N i s \Delta t + o_1(\Delta t), & (k, j) = \left(-\frac{1}{N}, \frac{1}{N}\right), \\ p_2 \gamma N i \Delta t + o_2(\Delta t), & (k, j) = \left(0, -\frac{1}{N}\right), \\ 1 - p_1 \beta N i s \Delta t - p_2 \gamma N i \Delta t + o_3(\Delta t), & (k, j) = (0, 0), \\ 0, & \text{otherwise,} \end{cases} \end{aligned} \quad (2.8)$$

where $\Delta \mathbf{s}_N(\omega, t) = \mathbf{s}_N(\omega, t + \Delta t) - \mathbf{s}_N(\omega, t)$, $\Delta \mathbf{i}_N(\omega, t) = \mathbf{i}_N(\omega, t + \Delta t) - \mathbf{i}_N(\omega, t)$ and $o_l(\Delta t)$ are functions of Δt that satisfy $\lim_{\Delta t \rightarrow 0} o_l(\Delta t)/\Delta t = 0$, $l = 1, 2, 3$.

Remark 2.2.2. *Although the set up of the Markov pure jump process implies (2.8), the converse is not true. For example, the set of small-time-interval equations do not uniquely determine the rates of the Poisson clocks. For the same reason, combining the parameters p_1 and β as one parameter leads to a different stochastic SIR model. For one thing, the variance of a Bernoulli random variable with parameter p_1 is $p_1(1 - p_1)$, which is not linear with p_1 .*

2.3 Main results using a martingale approach

The martingale formulation of a Markov pure jump process characterizes the process as the sum of an integral part involving the infinitesimal mean and a martingale part involving the infinitesimal variance.

2.3.1 Martingale formulation of the stochastic SIR-IPC model

For every t , we define $\mathbf{s}_N(t) := \{\mathbf{s}_N(\omega, t) : \omega \in \Omega\}$. In a similar way, we can define the stochastic processes $\mathbf{i}_N(t)$ and $\mathbf{r}_N(t)$ associated with the stochastic SIR-IPC model. As $(\mathbf{s}_N(t), \mathbf{i}_N(t), \mathbf{r}_N(t))$ is a Markov pure jump process with state space \mathbb{R}_+^3 , a martingale formulation (see e.g. [Dur96, Dur99, KL99, Lig10, SV06]) can be derived. We first introduce the notion of blow-up time (see e.g. [HPS72]).

Definition 2.3.1. *It is possible that for a certain realization of SIR-IPC Model, the Poisson clocks may generate time increments that add up to be finite, say time τ and $\tau < \infty$. In this case the corresponding stochastic process is defined only for $0 \leq t < \tau$. By the time $t = \tau$ the Poisson clocks will have made an infinite number of advances. We say that the stochastic process **explodes**, and τ is called the **blow-up time**.*

This is an uncommon occurrence and it is a different issue from the process extinguishing itself, in which the infected population goes to zero in finite time. The latter is very real and of interest for finite-size epidemics.

In the previous literature (e.g. [YC19, Yan08, All08, All11, AA03, BIH18, Ish93, Ish05]), there are works that consider the infinitesimal mean and variance of stochastic SIR models but typically only in the cases when the total population N increases to infinity. By contrast, we are interested in the dynamics of the SIR-IPC model for moderate size N in which the finite-size effects matter. Thus we derive the full martingale formulation encompassing the infinitesimal mean and variance of the SIR-IPC Model below.

Theorem 2.3.2. *Before the possible blow-up time, $(\mathbf{s}_N(t), \mathbf{i}_N(t), \mathbf{r}_N(t))$ can be written as*

$$\begin{cases} \mathbf{s}_N(t) = \mathbf{s}_{0N} + \int_0^t \mathcal{G}_N^{(1)}(w) dw + \mathcal{M}_N^{(1)}(t), \\ \mathbf{i}_N(t) = \mathbf{i}_{0N} + \int_0^t \mathcal{G}_N^{(2)}(w) dw + \mathcal{M}_N^{(2)}(t), \\ \mathbf{r}_N(t) = \mathbf{r}_{0N} + \int_0^t \mathcal{G}_N^{(3)}(w) dw + \mathcal{M}_N^{(3)}(t), \end{cases} \quad (2.9)$$

where $\mathcal{M}_N^{(l)}(t) = \mathcal{M}^{(l)}(\mathbf{s}_N(t), \mathbf{i}_N(t), \mathbf{r}_N(t))$, $l = 1, 2, 3$, are martingales that start at $t = 0$ as zeros, and $\mathcal{G}_N^{(l)}(t) = \mathcal{G}^{(l)}(\mathbf{s}_N(t), \mathbf{i}_N(t), \mathbf{r}_N(t))$, $l = 1, 2, 3$, are the infinitesimal means for $\mathbf{s}_N(t)$, $\mathbf{i}_N(t)$, and $\mathbf{r}_N(t)$, respectively, and

$$\mathcal{G}_N^{(1)}(t) = -p_1 \beta \mathbf{i}_N(t) \mathbf{s}_N(t), \quad (2.10)$$

$$\mathcal{G}_N^{(2)}(t) = p_1 \beta \mathbf{i}_N(t) \mathbf{s}_N(t) - p_2 \gamma \mathbf{i}_N(t), \quad (2.11)$$

$$\mathcal{G}_N^{(3)}(t) = p_2 \gamma \mathbf{i}_N(t). \quad (2.12)$$

The variances of $\mathcal{M}_N^{(l)}(t)$, $l = 1, 2, 3$, can be characterized in the following way:

$$\text{Var} \left(\mathcal{M}_N^{(l)}(t) \right) = \int_0^t \mathbb{E} \left[\mathcal{V}_N^{(l)}(w) \right] dw, \quad l = 1, 2, 3, \quad (2.13)$$

where $\mathcal{V}_N^{(l)}(t) = \mathcal{V}_N^{(l)}(\mathbf{s}_N(t), \mathbf{i}_N(t), \mathbf{r}_N(t))$, $l = 1, 2, 3$, are the infinitesimal variances for $\mathbf{s}_N(t)$,

$\mathbf{i}_N(t)$, and $\mathbf{r}_N(t)$ respectively, and

$$\mathcal{V}_N^{(1)}(t) = \frac{1}{N} p_1 \beta \mathbf{i}_N(t) \mathbf{s}_N(t), \quad (2.14)$$

$$\mathcal{V}_N^{(2)}(t) = \frac{1}{N} p_1 \beta \mathbf{i}_N(t) \mathbf{s}_N(t) + \frac{1}{N} p_2 \gamma \mathbf{i}_N(t), \quad (2.15)$$

$$\mathcal{V}_N^{(3)}(t) = \frac{1}{N} p_2 \gamma \mathbf{i}_N(t). \quad (2.16)$$

Proof of Theorem 2.3.2. To prove Theorem 2.3.2, we compute the infinitesimal means and variances for the Markov pure jump process $(\mathbf{s}_N(t), \mathbf{i}_N(t), \mathbf{r}_N(t))$ for a fixed N , using the methods in e.g., [App09, CW14, Dur02, HWY92, KT81, Kun86, M82, MP80, PZ07, Pro05].

As $\mathcal{G}_N^{(l)}(t)$, $l = 1, 2, 3$, are the infinitesimal means for $\mathbf{s}_N(t)$, $\mathbf{i}_N(t)$, and $\mathbf{r}_N(t)$, respectively, from (2.8) we have (2.10) – (2.12). As $\mathcal{V}_N^{(l)}(t)$, $l = 1, 2, 3$, are the infinitesimal variance of $\mathbf{s}_N(t)$, $\mathbf{i}_N(t)$, and $\mathbf{r}_N(t)$, respectively, we have

$$\mathcal{V}_N^{(1)}(t) \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} \mathbb{E} [\Delta(\mathbf{s}_N(t^-))^2 | ((\mathbf{s}_N(t^-), \mathbf{i}_N(t^-), \mathbf{r}_N(t^-)))] . \quad (2.17)$$

From (2.17) we obtain (2.14). In a similar way we obtain (2.15)–(2.16). With the infinitesimal means and variances obtained, we apply Theorem 1.6, [Dur96] or Theorem 3.32, [Lig10], to obtain (2.9), and apply Exercise 3.8.12 of [Bic02] to obtain (2.13). The proof of Theorem 2.3.2 is completed. \square

2.3.2 Estimates of the martingale variances

Here we derive upper bounds for variances of martingales of the martingale formulation (2.9).

Theorem 2.3.3. *Before the possible blow-up time, we have the following estimates:*

$$\text{Var} \left(\mathcal{M}_N^{(1)}(t) \right) \leq \frac{1}{N} \mathbf{s}_{0N}, \quad (2.18)$$

$$\text{Var} \left(\mathcal{M}_N^{(2)}(t) \right) \leq \frac{1}{N} [\mathbf{s}_{0N} + p_2 \gamma (\mathbf{i}_{0N} + \mathbf{s}_{0N}) t], \quad (2.19)$$

$$\text{Var} \left(\mathcal{M}_N^{(3)}(t) \right) \leq \frac{1}{N} p_2 \gamma (\mathbf{i}_{0N} + \mathbf{s}_{0N}) t. \quad (2.20)$$

Proof of Theorem 2.3.3. Taking expectation on both sides of the first equation of (2.9), we have

$$\mathbb{E} [\mathbf{s}_N(t)] = \mathbf{s}_{0N} + \int_0^t \mathbb{E} \left[\mathcal{G}_N^{(1)}(w) \right] dw. \quad (2.21)$$

This together with (2.10), (2.13) and (2.14) implies

$$\text{Var} \left(\mathcal{M}_N^{(1)}(t) \right) = \frac{1}{N} \int_0^t \mathbb{E} [p_1 \beta \mathbf{i}_N(w) \mathbf{s}_N(w)] dw = \frac{1}{N} (\mathbf{s}_{0N} - \mathbb{E} [\mathbf{s}_N(t)]). \quad (2.22)$$

As $\mathbb{E} [\mathbf{s}_N(t)] \geq 0$, we can bound the right-hand-side of (2.22) from above by \mathbf{s}_{0N}/N , and obtain (2.18).

From (2.13) and (2.16) we infer

$$\begin{aligned} \text{Var} \left(\mathcal{M}_N^{(3)}(t) \right) &= \frac{1}{N} p_2 \gamma \int_0^t \mathbb{E} [\mathbf{i}_N(w)] dw \\ &\leq \frac{1}{N} p_2 \gamma \int_0^t (\mathbf{i}_{0N} + \mathbf{s}_{0N}) dw \\ &= \frac{1}{N} p_2 \gamma (\mathbf{i}_{0N} + \mathbf{s}_{0N}) t. \end{aligned} \quad (2.23)$$

From (2.23) we obtain (2.20).

From (2.13)–(2.16), we obtain

$$\text{Var} \left(\mathcal{M}_N^{(2)}(t) \right) = \text{Var} \left(\mathcal{M}_N^{(1)}(t) \right) + \text{Var} \left(\mathcal{M}_N^{(3)}(t) \right). \quad (2.24)$$

This together with (2.18) and (2.20) implies (2.19). \square

Theorem 2.3.3 implies that the variances of the martingales of $\mathbf{s}_N(t)$, $\mathbf{i}_N(t)$, $\mathbf{r}_N(t)$ are bounded by the initial states of each compartment and the total population N . Moreover, they are bounded by quantities with an equal or lower order of magnitude than $1/N$.

2.3.3 Deterministic analogue of the stochastic SIR model

With a similar derivation of the fluid dynamic limit of Markov pure jump processes [EK86, GPV88, KL99, KOV89, SV06, TV03, Var95, Var00], we can find a deterministic analogue of the (renormalized) stochastic SIR model when N increases to infinity. The analysis is based on the martingale formulation (2.9), and the conclusion of Theorem 2.3.3, where the martingale variances are shown to have an equal or lower order of magnitude than $1/N$.

And we see that for N large, it is reasonable to set the infinitesimal mean vector as the generator of the deterministic flow of a deterministic analogue of the stochastic SIR model. And by (2.10), (2.11) and (2.12) we obtain

$$\left\{ \begin{array}{l} \frac{\partial s(t)}{\partial t} = -p_1 \beta i(t) s(t), \\ \frac{\partial i(t)}{\partial t} = p_1 \beta i(t) s(t) - p_2 \gamma i(t), \\ \frac{\partial r(t)}{\partial t} = p_2 \gamma i(t), \\ s(0) = s_0, \quad i(0) = i_0, \quad r(0) = r_0, \end{array} \right. \quad (2.25)$$

where $s(t)$, $i(t)$, and $r(t)$, $t \in (0, \infty)$ are the deterministic versions of $\mathbf{s}_N(t)$, $\mathbf{i}_N(t)$, and $\mathbf{r}_N(t)$, respectively, and $s_0 = S_0/N$, $i_0 = I_0/N$, $r_0 = R_0/N$. We note that the ODEs in Eqs. (2.25) do not explicitly include the population size N , and thus, are independent of N . Note that the deterministic analogue of the stochastic SIR model (2.25) are the same as the classical deterministic SIR model (without vital dynamics) as in e.g. [KMW27] rescaled by N .

2.3.4 Numerical simulations

We compare the stochastic and deterministic analogue of the stochastic SIR model through simulations. For the stochastic case, we use the classical Gillespie algorithm for continuous-time Markov

processes (see e.g. [Gil76]). Roughly speaking, we first simulate the sojourn times using exponential distribution generators, and then simulate the sample paths of the embedded discrete-time Markov process as described in Section 2.2.1. The algorithm ensures that independent Poisson clocks advance randomly without a prescribed diagram. Namely, we do not require that the i -clocks and r -clocks advance sequentially with a pre-determined arrangement. In this way, we are able to produce a number of sample paths corresponding to random realizations of sequences of events.

We focus on the cases when the basic reproduction number \mathcal{R}_0 is near the self-sustaining level of $\mathcal{R}_0 = 1$, as is typical in a scenario where public health measures trade off against economic and educational needs of the population¹. Specifically, we consider the cases when $\mathcal{R}_0 = 0.95, 1.1, 1.2$, and 1.3 , respectively. Parameters of the simulations are recorded in the figure caption.

In all the cases, we show output of the SIR-IPC model simulations with $N = 10^{2.5}, 10^3, 10^{3.5}$, and 10^4 , represented by the blue, orange, green, and red lines, respectively (Figs. 2.1b, 2.1d, 2.1f, 2.1h). For the case when $\mathcal{R}_0 = 0.95$, we infer from (2.25) that the parameters and data used to obtain the blue, orange, green, and red lines in Fig. 2.1b give rise to the same solution $(s(t), i(t), r(t))$ to the deterministic SIR Model. Therefore, we only display the deterministic output once in Fig. 2.1a. This same arrangement also applies to other cases with different values of \mathcal{R}_0 , which are displayed in Figs. 2.1c, 2.1e, and 2.1g, respectively. In all the figures, black, magenta, and cyan lines represent $s(t), i(t)$, and $r(t)$ associated with the deterministic simulation, respectively.

The *finite-size effects* are exhibited. As the population size decreases, deviations of dynamics of the stochastic SIR model from its continuum equation increase. The same simulation output is also observed over other random paths.

¹As a specific example we note that the dynamic reproductive number has been estimated weekly in Los Angeles County using a Bayesian SEIR model applied to hospital demand data—during the period July 2020–April 2021 it remained within a 25% window of one [BBC21].

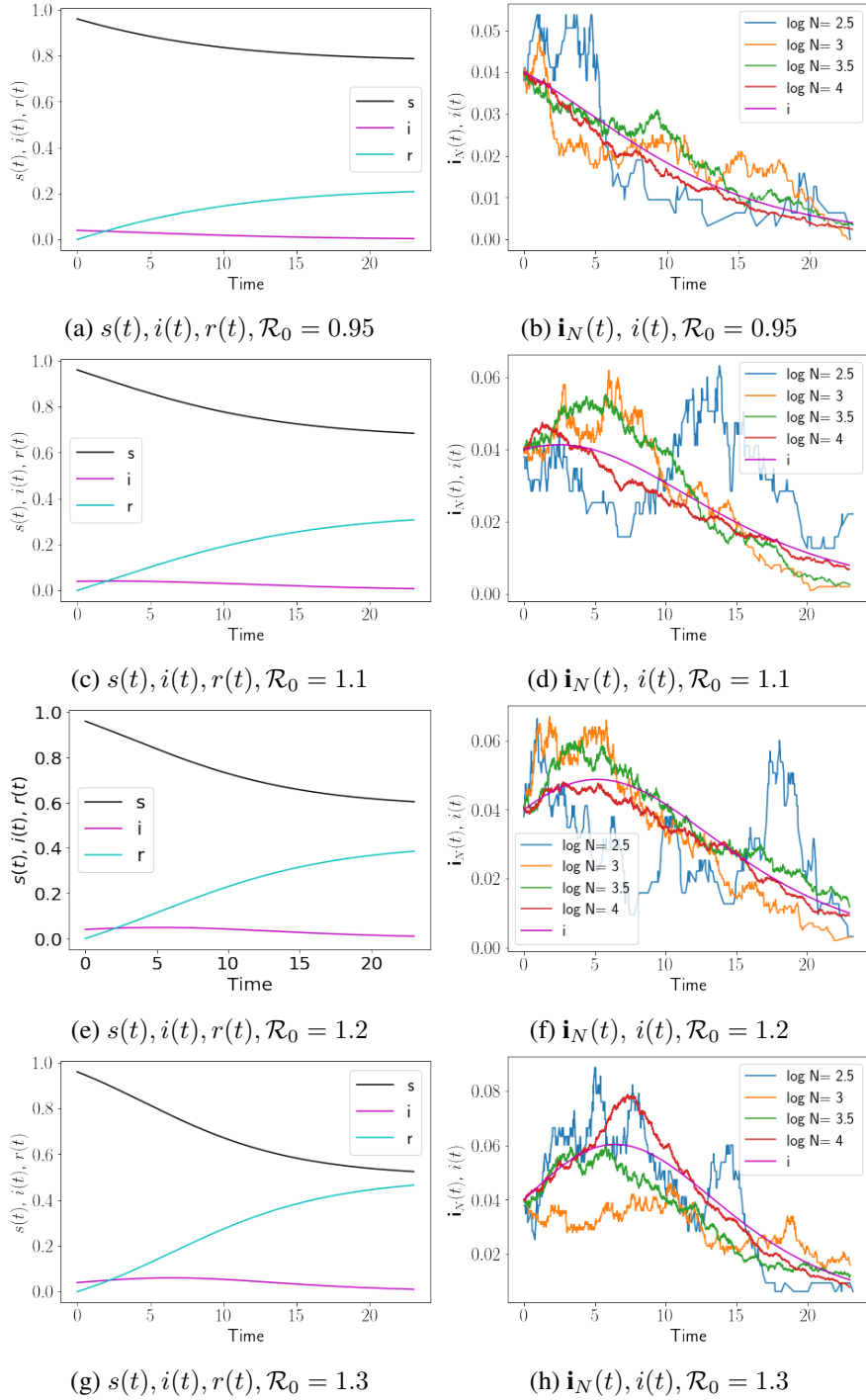


Figure 2.1: Plots of infected compartment fractions $\mathbf{i}_N(t)$ for the stochastic SIR and $s(t), i(t)$, and $r(t)$ for deterministic SIR. For both models, $p_1 = 0.5, p_2 = 0.5, \gamma = 1, T = 23, s_0 = \mathbf{s}_{0N} = 0.96, i_0 = \mathbf{i}_{0N} = 0.04, r_0 = \mathbf{r}_{0N} = 0$. In Figs. 2.1b and 2.1a, 2.1d and 2.1c, 2.1f and 2.1e, and 2.1h and 2.1g, $\beta = 0.95, 1.1, 1.2, 1.3$, respectively. For the SIR-IPC model displayed in Figs. 2.1b, 2.1d, 2.1f, and 2.1h, $N = 10^{2.5}, 10^3, 10^{3.5}, 10^4$ in every panel.

2.4 Mathematical analysis of the *finite-size effects*

In Theorem 2.3.3, a global upper bound of the same order of magnitude of $1/N$ is derived and the result indicates the order of the magnitude of the martingale infinitesimal variances when $N \rightarrow \infty$. To further understand the martingale variances when N is relatively smaller, in this section, we analyze the *finite-size effects*, based on the martingale formulation, and simulations are run which support our theoretical conclusion.

2.4.1 Scaling property of the stochastic component

We analyze the deterministic and stochastic component of the martingale formulation with varying population size N . We fix the initial condition of the compartment fractions and denote them as

$$(\mathbf{s}_{0N}, \mathbf{i}_{0N}, \mathbf{r}_{0N}) \equiv (\mathbf{s}_0, \mathbf{i}_0, \mathbf{r}_0). \quad (2.26)$$

We analyze the martingale formulation of the stochastic SIR model with varying population size N . Applying (2.9) and (2.13) to $(\mathbf{s}_N(t), \mathbf{i}_N(t), \mathbf{r}_N(t))$ over a small time step Δt , we obtain

$$\begin{cases} \mathbf{s}_N(t + \Delta t) = \mathbf{s}_N(t) + \mathcal{G}_N^{(1)}(t)\Delta t + \mathcal{M}_N^{(1)}(t + \Delta t) - \mathcal{M}_N^{(1)}(t), \\ \mathbf{i}_N(t + \Delta t) = \mathbf{i}_N(t) + \mathcal{G}_N^{(2)}(t)\Delta t + \mathcal{M}_N^{(2)}(t + \Delta t) - \mathcal{M}_N^{(2)}(t), \\ \mathbf{r}_N(t + \Delta t) = \mathbf{r}_N(t) + \mathcal{G}_N^{(3)}(t)\Delta t + \mathcal{M}_N^{(3)}(t + \Delta t) - \mathcal{M}_N^{(3)}(t). \end{cases} \quad (2.27)$$

By (2.13) and additivity of the variance in time for martingales, we have

$$\sqrt{\text{Var}(\Delta \mathcal{M}_N^{(l)}(t))} \cong \sqrt{\mathbb{E}[\mathcal{V}_N^{(l)}(t)] \Delta t}, \quad l = 1, 2, 3, \quad (2.28)$$

where $\Delta \mathcal{M}_N^{(l)}(t) = \mathcal{M}_N^{(l)}(t + \Delta t) - \mathcal{M}_N^{(l)}(t)$, $l = 1, 2, 3$. This together with (2.27) implies that the infinitesimal variances are the key to estimate the infinitesimal standard deviation of the stochastic component and the deviation of the trajectories of the evolution of the model from its deterministic

component.

We perform estimates at the first time step. At time zero, from (2.14)–(2.16) we infer

$$V_N^{(1)}(0) = \frac{1}{N} p_1 \beta \mathbf{i}_0 \mathbf{s}_0, \quad (2.29)$$

$$\mathcal{V}_N^{(2)}(0) = \frac{1}{N} (p_1 \beta \mathbf{i}_0 \mathbf{s}_0 + p_2 \gamma \mathbf{i}_0), \quad (2.30)$$

$$\mathcal{V}_N^{(3)}(0) = \frac{1}{N} p_2 \gamma \mathbf{i}_0, \quad (2.31)$$

which implies that the infinitesimal variances for the attractiveness are each inversely proportional to N :

$$\mathcal{V}_N^{(l)}(0) \propto \frac{1}{N}, l = 1, 2, 3. \quad (2.32)$$

This together with (2.28) implies that at the first time step we have

$$\text{Var} \left(\mathcal{M}_N^{(l)}(\Delta t) \right) \propto \frac{1}{N}, l = 1, 2, 3. \quad (2.33)$$

From (2.33) and (2.27) we infer that at the first time step a smaller value of N leads to a larger deviation of the trajectory of $(\mathbf{s}_N(\omega, t), \mathbf{i}_N(\omega, t), \mathbf{r}_N(\omega, t))$ from its deterministic component. This explains the *finite-size effects* at the first time step. This suggests that smaller value of N leading to a larger deviation remains to be true at an arbitrary later time, namely,

$$\mathcal{V}_N^{(l)}(t) > \mathcal{V}_{\tilde{N}}^{(l)}(t), \text{ for } 0 < N < \tilde{N} \text{ and } t > 0, l = 1, 2, 3, \quad (2.34)$$

which leads to a theory of the *finite-size effects* at an arbitrary later time.

Next we estimate the right-hand-side of formulas (2.14)–(2.16), so that we can estimate $\mathcal{V}_N^{(l)}(t)$, $l = 1, 2, 3$, at later times for $t > 0$.

In the prior literature of stochastic partial differential equations (see e.g. [Br, CGS13, Rao99, Uch08, WTG11]), usually Itô calculus is applied, as the infinitesimal variances are multiplied by Wiener processes, and the population or number of particles is required to be large, which does not

fit with the cases on which our focus here. We conjecture that formulas (2.29)–(2.31) can be written for an arbitrary $t > 0$ with the same leading order:

$$\mathcal{V}_N^{(1)}(t) \propto \frac{1}{N} p_1 \beta i_0 s_0, \quad (2.35)$$

$$\mathcal{V}_N^{(2)}(t) \propto \frac{1}{N} (p_1 \beta i_0 s_0 + p_2 \gamma i_0), \quad (2.36)$$

$$\mathcal{V}_N^{(3)}(t) \propto \frac{1}{N} p_2 \gamma i_0. \quad (2.37)$$

If our conjecture above is true, then we can fix a time period $[T_1, T_2]$, and integrate of both sides of (2.35)–(2.37) over the time period $[T_1, T_2]$. By taking the average we obtain

$$\overline{\mathcal{V}_N^{(1)}} \propto \frac{1}{N} p_1 \beta i_0 s_0, \quad (2.38)$$

$$\overline{\mathcal{V}_N^{(2)}} \propto \frac{1}{N} (p_1 \beta i_0 s_0 + p_2 \gamma i_0), \quad (2.39)$$

$$\overline{\mathcal{V}_N^{(3)}} \propto \frac{1}{N} p_2 \gamma i_0, \quad (2.40)$$

where $\overline{\mathcal{V}_N^{(l)}}$, $l = 1, 2, 3$ denotes the time average of the infinitesimal variance over this time period:

$$\overline{\mathcal{V}_N^{(l)}} := \frac{1}{T_2 - T_1} \int_{T_1}^{T_2} \mathcal{V}_N^{(l)}(t) dt, \quad l = 1, 2, 3. \quad (2.41)$$

Here we slightly abuse the notation and omit the dependence of the time average over T_1 and T_2 .

2.4.2 Numerical simulations

To check the validity of (2.34), we perform direct simulations of the infinitesimal standard deviations, namely

$$\sigma_N^{(l)}(t) = \sqrt{\mathcal{V}_N^{(l)}(t)}, \quad l = 1, 2, 3. \quad (2.42)$$

Example output can be found in Figure. 2.2. Figs. 2.2a, 2.2c, and 2.2e show results of $\sigma_N^{(l)}(t)$, $l = 1, 2, 3$, in the case with $\mathcal{R}_0 = 0.95 < 1$. The blue, orange, green and red lines show results

with the same realizations from simulations with the corresponding colors in Fig. 2.1b. Figs. 2.2b, 2.2d, and 2.2f show results of $\sigma_N^{(l)}(t)$, $l = 1, 2, 3$, in the case with $\mathcal{R}_0 = 1.3 > 1$. The blue, orange, green and red lines show results with the same realizations from simulations with the corresponding colors in Fig. 2.1h.

The outputs of the simulations agree with (2.34). The same simulation results are also observed over other random paths. We validate our conjecture of (2.38)–(2.40) through the following numerical simulation.

Fig. 2.3 shows the log–log plot with error bars for (2.38) - (2.40). The lines show the theoretical scaling with slope as -1 and the x -intercepts as $\log(p_1\beta i_0 s_0)$, $\log(p_1\beta i_0 s_0 + p_2\gamma i_0)$, and $\log(p_2\gamma i_0)$, respectively, and the error bars show the true scaling with the x -coordinate and y -coordinate as follows:

$$(x, y) = \left(\log N, \log \left(\overline{\mathcal{V}_N^{(l)}} \right) \right), \quad l = 1, 2, 3, \quad (2.43)$$

for $N = 10^{2.5}, 10^3, 10^{3.5}, 10^4$, and $l = 1, 2, 3$. Here $[T_1, T_2]$ are chosen as $[0, 10]$, $[1, 11]$, $[2, 12]$, ..., and $[13, 23]$. The minimum and maximum values of $y = \log \left(\overline{\mathcal{V}_N^{(l)}} \right)$, $l = 1, 2, 3$, taken over all such intervals are set as the upper and lower bounds of the error bars. Figs. 2.3a, 2.3c, and 2.3e show results in the cases when $\mathcal{R}_0 = 0.95$ for $l = 1, 2$, and 3 , respectively. The error bars with horizontal x -axis as 2.5, 3, 3.5, and 4 show results with the same realization from simulations of the blue, orange, green, and red lines in Fig. 2.1b, respectively. Figs. 2.3b, 2.3d, and 2.3f show results in the cases when $\mathcal{R}_0 = 0.13$ for $l = 1, 2$, and 3 , respectively. The error bars with x -axis as 2.5, 3, 3.5, and 4, show results with the same realization from simulations of the blue, orange, green, and red lines in Fig. 2.1h, respectively.

The output shows that the error bars include the straight lines representing the theory and scale with the line. The same simulation results are also observed over other random paths. These results support the validity of equations (2.38)–(2.40) and our theory for the *finite-size effects* based on the martingale formulation.

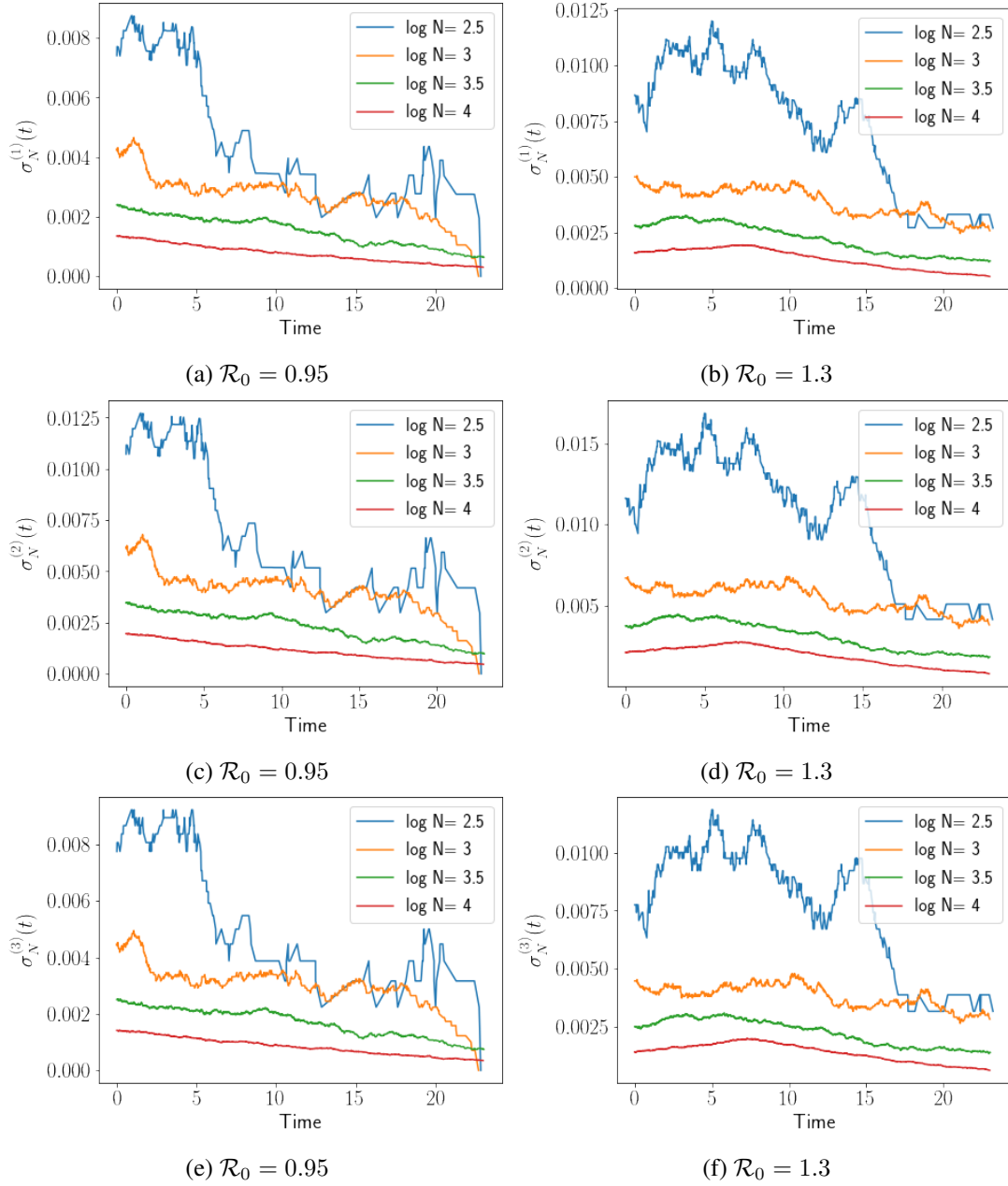


Figure 2.2: Examples of the infinitesimal standard deviation for $\sigma_N^{(l)}(t)$, $l = 1, 2, 3$. In both cases, $p_1 = 0.5$, $p_2 = 0.5$, $\gamma = 1$, $T = 23$, $s_0 = \mathbf{s}_{0N} = 0.96$, $i_0 = \mathbf{i}_{0N} = 0.04$, $r_0 = \mathbf{r}_{0N} = 0$. In the cases with $\mathcal{R}_0 = 0.95 < 1$ in Figs 2.2a, 2.2c, and 2.2e, $\beta = 0.95$. The blue, orange, green, and red lines show results with the same realization from simulations with the corresponding colors in Fig 2.1b. In the cases with $\mathcal{R}_0 = 1.3 > 1$ in Figs 2.2b, 2.2d, and 2.2f, $\beta = 0.95$. The blue, orange, green, and red lines show results with the same realization from simulations with the corresponding colors in Fig 2.1h.

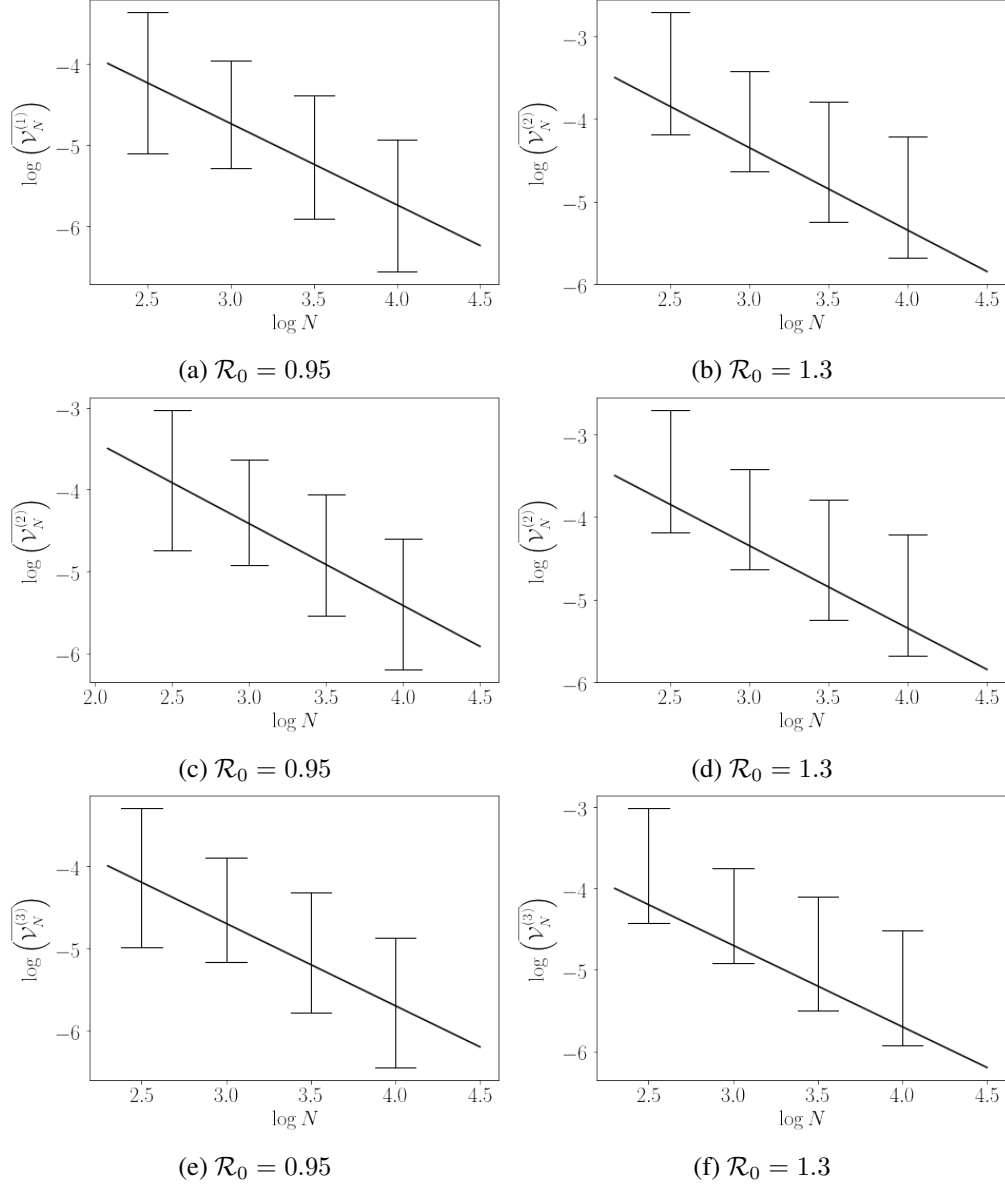


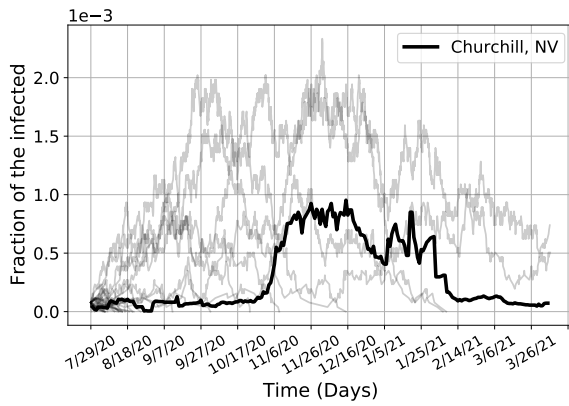
Figure 2.3: Comparisons of the log–log plot of the theoretical and empirical scaling for $\overline{\mathcal{V}_N^{(l)}}$, $l = 1, 2, 3$, as in (2.38) – (2.40). The vertical bars are the error bars of the empirical scaling of the infinitesimal variance from the same realization as in Figs. 2.1b, and 2.1h. We set $[T_1, T_2]$ as $[0, 10]$, $[1, 11]$, $[2, 12]$, \dots and $[13, 23]$. The minimum and maximum values of $y = \log(\overline{\mathcal{V}_N^{(l)}})$, $l = 1, 2, 3$, taken over all such intervals are set as the upper and lower bounds of the error bars. The straight lines with slope -1 show the theoretical scaling with the x -intercepts as $\log(p_1\beta i_0 s_0)$, $\log(p_1\beta i_0 s_0 + p_2\gamma i_0)$, and $\log(p_2\gamma i_0)$, for $l = 1, 2, 3$. In all figures, $p_1 = 0.5$, $p_2 = 0.5$, $\gamma = 1$, $T = 23$, $s_0 = \mathbf{s}_{0N} = 0.96$, $i_0 = \mathbf{i}_{0N} = 0.04$, and $r_0 = \mathbf{r}_{0N} = 0$. In Figs. 2.3a, 2.3c, and 2.3e, $\beta = 0.95$. In Figs. 2.3b, 2.3d, and 2.3f, $\beta = 1.3$.

2.5 Case studies of a small county and a cruise ship

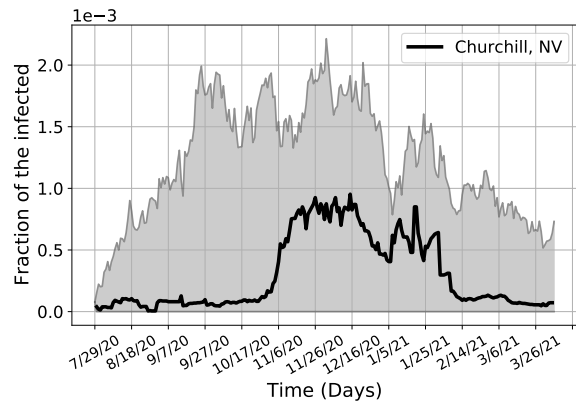
In the following experiments, we set parameters to mimic real world cases of a small county or a cruise ship. We will display plots for the fractional of compartment and the infinitesimal standard deviations, and analyze the role of stochastic variances in each case. The infinitesimal variances represent chance variations that lead to deviations from the deterministic SIR model. Through simulations of infinitesimal variances as in (2.14)–(2.16), we can produce ranges that encompass a majority of the random paths. Realistically, our findings suggest that for small populations, the estimating probabilities of rare events, e.g. the early-die-out or early-outbreak, are important.

We study the outbreak of COVID-19 in Churchill County, Nevada from July, 2020 to March, 2021 and on the Diamond Princess cruise ship from February to March, 2020. Since early testing on the Diamond Princess was done by sampling from the population due to limited availability of testing, there is an under-reporting issue in the infected counts. The data of Churchill County was collected from [DDG20] and the reported number of confirmed cases of Princess Diamond, and the quarantine process were retrieved from the Princess Cruise website of the Carnival Cooperation [Cru21] and the official website of Ministry of Health, Labor and Welfare, Japan (Ministry of Health, Labor and Welfare, Japan [HW20]).

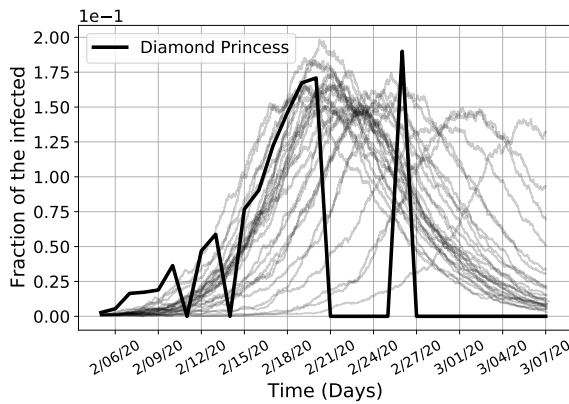
We plot the daily confirmed cases with 30 realiations from the SIR-IPC models in Fig. 2.4. For Fig. 2.4a, the total population is set to be Churchill County’s population $N = 25715$ and the basic reproduction number $\mathcal{R}_0 = 1.007$ using linear regression optimized by a grid search. For Fig. 2.4c the total population N is the total number of passengers and crew on Diamond Princess and $N = 3711$. Cruise ships carry a large number of people in confined spaces with relative homogeneous mixing, amplifying an already highly transmissible disease. Following the estimation of the dynamic basic reproduction number in [RSW20], we take an average of these estimates from Feb 5 to Feb 26, resulting in $\mathcal{R}_0 = 2.02$. With these smaller populations, the infinitesimal standard deviations are non-negligible in the stochastic SIR model. Public health officials thus need to be aware of the variability of possible outcomes due to finite-size effects.



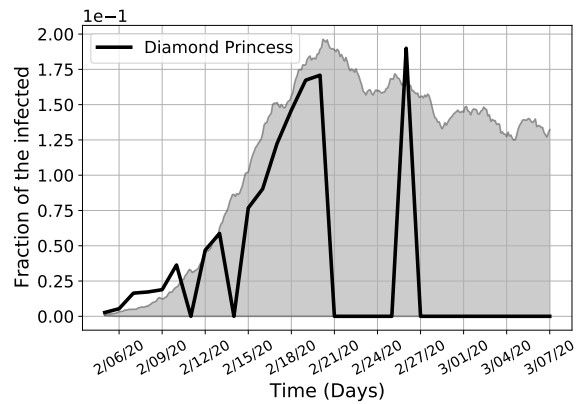
(a) Churchill. $R_0 = 1.007$, $N = 25715$



(b) The range of the Churchill simulated data.



(c) Diamond Princess. $R_0 = 2.02$, $N = 3711$



(d) The range of the Diamond Princess simulated data.

Figure 2.4: Comparison of field data (solid black line) for daily confirmed case percentages with 30 realisations of the stochastic SIR model for Churchill County, NV and the Diamond Princess Cruise Ship.

2.6 Discussion

Compartmental models are a powerful tool to predict and control infectious diseases. In this chapter, we study theoretically and numerically the *finite-size effects* arising in stochastic compartmental models, where individual realizations of the models deviate from their fluid limit. We apply compound Poisson processes to the classical SIR compartmental models without vital dynamics. The result is a continuous-time Markov pure jump process, and a martingale approach can be applied to this process. The process is expressed as the sum of a deterministic and a stochastic component, which provides us with a tool to study both the statistical and stochastic features of the process. The deterministic part coincides with the classical deterministic SIR model. By providing a bound of the variances of the martingale, we show that the fluid limit is indeed the deterministic SIR model.

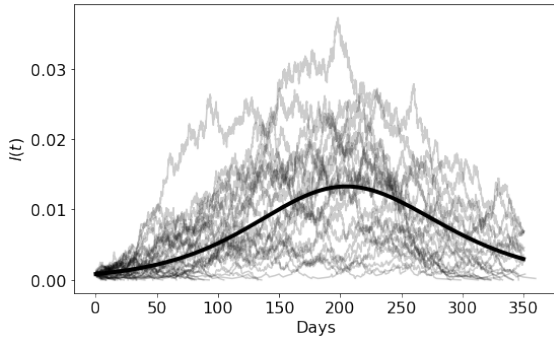
However, a small population size leads to stochastic fluctuations that deviate from the deterministic SIR model. Numerical and theoretical analysis of these *finite-size effects* have not been well-studied. We found a theoretical explanation for the *finite-size effects* by observing that the stochastic component of the martingale formulation scales as the inverse of the square root of the population size. A larger variance both in the outbreak size and its temporal behavior arises as population size decreases. This scaling property is verified at time zero with equilibrium initial data. Direct numerical simulations of the theoretical infinitesimal variance support our theory as output shows that the dependence on population size remains to be true at any time. Here we simulate with fixed initial fractions of compartments and different total population sizes. To the best of our knowledge, this is the first time that simulations of theoretical infinitesimal variances for stochastic compartmental models are implemented. In previous works only empirical variations are simulated (see e.g. [RMK18], [Ish91]), and the focus was on the effects of statistical fluctuations on the basic reproduction number assuming varying initial compartment fractions [AG18]. We also simulate *finite-size effects* with small populations and analyze with real data. All the simulations support our theory. Our results exhibit the danger of fitting data collected during an outbreak to deterministic

counterparts of the stochastic compartmental models, especially for small populations.

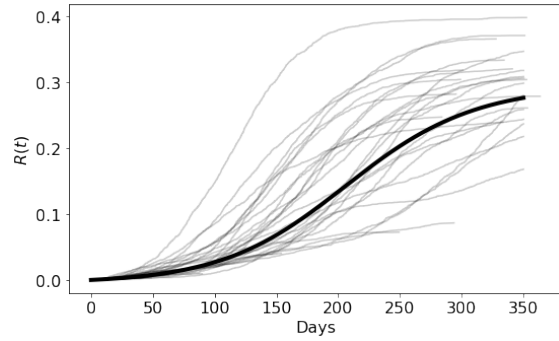
finite-size effects are observed in Markov processes for population dynamics in many fields ([CTW19, GWW19, Ham18, HMZ18, IMM18, Jim18, KWT18, KW18, PNW19, TW18, TWW18, War18]) but there have been few quantitative studies to date about *finite-size effects* in epidemics. The methodology developed here may be broadly useful for quantitative social and natural sciences may provide a mathematical and theoretical framework that may contribute to epidemic policy for public health agencies.

So far, the analysis is purely theoretical. It remains an open problem to integrate the real-world data to the model and provide insights from the data to control the disease. One example is that in order to control the spread, governments should know when and where to take certain measures to reduce the transmission rate β .

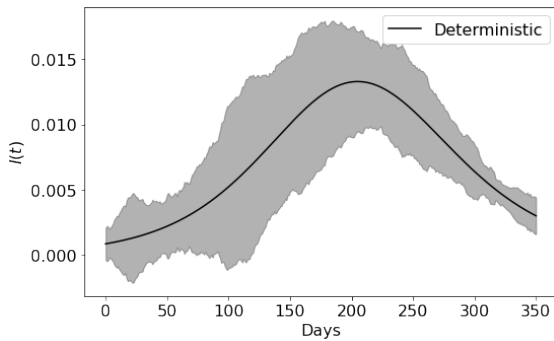
From Fig. 2.5c, with a smaller population, the range of synthetic data can be wide. Additionally, the process with a smaller population also has more paths that die out early in the process (see Fig. 2.2a, 2.2c, 2.2e, and 2.5a). Our work provides a guide for authorities of smaller populations like cruise ships and small towns to estimate risk over time in order to prepare for outbreaks. It is important to bear in mind that the broader variations in the pandemic caused by the smaller population lead to a wide outcome when it comes to estimating risk. In the past hundreds of years of human history, there have been several infectious diseases, including SARS, and MERS. However, only few of them, e.g. COVID-19 and 1918 influenza escalated into a pandemic. It is important to understand when, why and how a disease dies out and a future direction is to quantitatively study the die-out event and its relations to finite size population and basic reproduction number \mathcal{R}_0 .



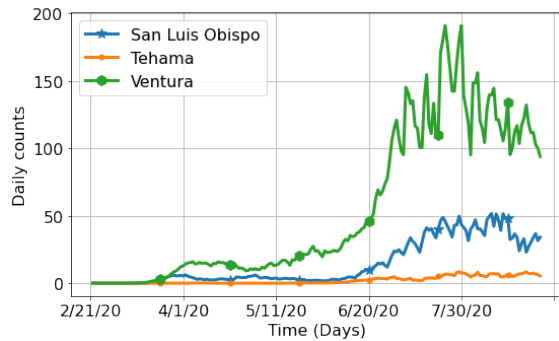
(a) Deterministic $i(t)$ with 50 simulations of stochastic $\mathbf{i}(t)$. The parameters are estimated using daily infected counts and total population of San Luis Obispo county, CA and $\beta = 0.13, \gamma = 0.11, N = 4777$.



(b) Deterministic $r(t)$ with 50 simulations of stochastic $\mathbf{r}(t)$.



(c) Deterministic $i(t)$ with its theoretical Confidence band with the width being $6\sqrt{\mathcal{V}_I(t)}$.



(d) 10-day moving average of the daily confirmed cases in 3 counties in California. Ventura's peak is around end of July, San Luis Obispo's peak is around mid August.

Figure 2.5: 10-day moving average of the daily increased confirmed cases in 3 counties in California. Ventura's peak is around end of July; San Luis Obispo's peak is around mid August.

CHAPTER 3

Multi-regional Policy-making and the Epidemics

This chapter is based on collaborative work with Andrea Bertozzi, P. Jeffrey Brantingham and Yevgeniy Vorobeychik. The problem and approach was suggested by the collaborators. I developed the numerical methods and all the simulations as well as fine-tuning the details of the models.

3.1 Background

In the course of battling COVID-19, public health policies determine enforced non-pharmaceutical interventions to slow down or halt the spread of the pandemic. Some of them include ‘safer-at-home’, ‘maintaining 6 feet distance’ and ‘mask wearing’ which were considered crucial during the early stage prior to the availability of vaccines. The timeline of COVID-19 globally and locally ([Cena, Wika]) indicates that the evolution of policy affects the evolution of the pandemic and vice versa. For example, in the county of Los Angeles, physical distancing was first mandated [Dep22] on March 21, 2020, about a month after the first reported COVID-19 case in LA. Around that time, the mayor’s office released ‘safer-at-home’ [Los20]. One week later, beaches, hiking trails, dog parks, skate parks, etc., and more public sites and facilities were temporarily closed. As infected cases continued to increase, a month later, on May 1st, facial coverings were suggested. In hindsight, a natural question to ask is whether the policies that were enforced were done so in an optimal way. What can we learn from past examples and by using mathematical modeling to understand the interplay between policy and spread of disease. This chapter introduces a policy model coupled with the susceptible–infected–recovered (SIR) epidemic model to study interactions between policy-making and the dynamics of epidemics. There have been several studies on the

relationship between policies and epidemics [BKM02, BF07, PHB18, BEC06]. In study analyzing data from 16 US cities during the 1918 pandemic [BF07], Bootsma and Ferguson analyzed specific outcomes related to the impact of the delay of lockdown policies on the total deaths and also on the appearance of second waves of outbreaks due to reopening too early. That analysis was done fitting to an SEIR model. For optimal control, they considered the simple SIR and the end-state of the pandemic, noting that there exists an optimal control level with fewer deaths and no second wave. More recently, Bliman et al. [BDP21] develop a theoretical study of the optimal control of a classical SIR outbreak. In their work, they do not consider the possibility of vaccines or pharmaceutical interventions. Focusing exclusively on NPIs, they ask the question of how to design an optimal policy that achieves an end state as close as possible to the herd immunity threshold. This is the same problem considered briefly in a section of [BF07]. Bliman et al. prove the existence and the uniqueness of the solution and showed the optimal social distancing policy is a bang-bang controller [BD21], generalizing the results of [BDP21] by modeling without prescribing the onset of the policy. Their model assumes a policy that can change continuously in time, however this is impractical in real-life situations. As observed during the COVID-19 pandemic, Policy-makers must provide easy-to-follow policies, with a small number of different intensity levels (see Fig. 3.5b). Moreover, policies can not change frequently in time or they can not be easily followed. A practical implementation requires a minimum time duration for a particular stage of the policy. This can be modeled as a piece-wise constant function of time with a minimum time interval for each policy stage. With this idea in mind, we model policies as piece-wise linear functions in time and aim to find an optimal and most practical policy among all piece-wise policies. Another practical issue is the trade off between decreased infections and negative impact on other aspects of society such as remote learning for young students, lack of employment for people in certain job sections, and lack of services provided to the public. The prior models do not consider these important issues.

In this work, we modify the model in [BDP21] to take into account all of the practical issues described above. The chapter is organized in the following way: we first introduce the work in [BDP21] and reproduce the results using our methods. We discuss different optimal policies

resulting from different parameters. Next, we discuss a case study of the 2nd wave (November 2020–May 2021) in the county of Los Angeles, California and a simulated case with multiple regions. Lastly, we study a case of three counties with and without a governing state as an example of the multi-layer multi-regional model

3.2 Policy model using optimal control

A policy function is a continuous function that has a range of $[0, 1]$. As the numerical value increases, the strictness of the policy decreases. The Numerical value 0 denotes a total lockdown and 1 denotes no control. We assume a policy $u(t)$ directly influences the level of a lockdown, which affects the rate of the population transport from compartment S to I . We use the following policy-incorporated SIR:

$$\left\{ \begin{array}{l} \frac{dS(t)}{dt} = -u(t)\beta\frac{I(t)S(t)}{N}, \\ \frac{dI(t)}{dt} = u(t)\beta\frac{I(t)S(t)}{N} - \gamma I(t), \\ \frac{dR(t)}{dt} = \gamma I(t), \\ S(0) = S_0, \quad I(0) = I_0, \quad R(0) = R_0. \end{array} \right. \quad (3.1)$$

In [BDP21], a policy $u(t)$ is assumed to belong to the admissible set $\mathcal{U}_{\alpha_{\max}, T_0}$ defined by

$$\{u \in \mathcal{L}^\infty([0, +\infty]), \alpha_{\max} \leq u(t) \leq 1 \text{ if } t \in [0, T_0], u(t) = 1 \text{ if } t > T_0\}.$$

The constant T_0 characterizes the duration of the policy, and α_{\max} its maximal intensity. In [BDP21], Theorem 2.1 states that no finite time intervention is able to stop the epidemics before or exactly at the herd immunity. However, one may stop arbitrarily close to the latter by allowing sufficiently long intervention, provided that the intensity is sufficiently strong. Therefore, stopping S arbitrarily close to the herd immunity threshold is only possible by sufficiently long intervention of a strong enough intensity. To determine the closest state to this threshold attainable by control of maximal

intensity α_{\max} on the interval $[0, T_0]$, one is led to consider the following optimal control problem:

$$\sup_{u \in \mathcal{U}_{\alpha_{\max}, T_0}} S_{\infty}(u). \quad (3.2)$$

Furthermore, they prove, in theory, the existence and uniqueness of the optimal solution to problem 3.2 and that the solution is a bang-bang controller (a control that switches from one extreme to the other). More specifically, they have the following theorem:

Theorem 3.2.1. (Theorem 2.1 in [BDP21]) *Let $\alpha_{\max} \in [0, 1)$ and $T_0 > 0$. Problem 3.2 admits a unique solution u^* . Furthermore,*

- (i) *the maximal value $S_{\infty, \alpha_{\max}, T_0}^* := \{\max S_{\infty}(u) : u \in \mathcal{U}_{\alpha_{\max}, T_0}\}$ is non-increasing with respect to α_{\max} and non-decreasing with respect to T_0 .*
- (ii) *there exists a unique $T_0^* \in [0, T_0)$ such that $u^* = u_{T_0^*} := \mathbb{1}_{[0, T_0^*]} + \mathbb{1}_{[T_0^*, T_0]} + \mathbb{1}_{[T_0, +\infty)}$ (in particular, the optimal control is bang-bang).*

3.3 Models

We use the same policy-incorporated SIR model for the epidemic dynamic as in [BDP21]. Instead of minimizing the final epidemic size alone, we adopt a similar policy-making process as in [JML21] by using a cost function that takes into account the cost of implementing the policy, the impact of the infection and a penalty for being non-compliant. We also consider the practical implementation constraints, namely that the policy can only be implemented using a finite number of discrete levels of control and with a minimal time interval of during for each level. As an example, consider the policy implementation in France during the year 2020 and 2021 shown in Fig. 3.1 [Wikb]. This discrete model is not only practical in a use-case setting but also allows for ease of numerical computation of the optimal policy by searching through a discrete set of values rather than a continuum of policies.

Timeline of measures

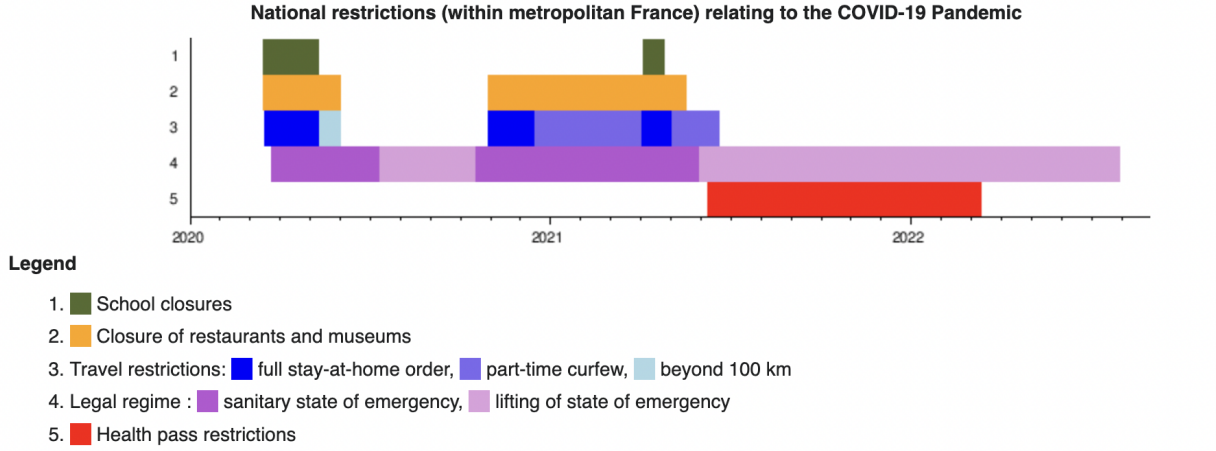


Figure 3.1: Timeline of COVID-19 Restrictions in France. Note the discrete nature of the restrictions, both in terms of the small number of categories and the fixed time intervals of enforcement.

3.3.1 The policy-incorporated SIR model

To model the evolution of the pandemic, we discretize the system of ODE using forward Euler's method with a time step of 1:

$$\left\{ \begin{array}{l} S(t) = S(t-1) - \alpha\beta \frac{I(t-1)S(t-1)}{N}, \\ I(t) = I(t-1) + \alpha\beta \frac{I(t-1)S(t-1)}{N} - \gamma I(t-1), \\ R(t) = R(t-1) + \gamma I(t-1), \\ S(0) = S_0, \quad I(0) = I_0, \quad R(0) = R_0. \end{array} \right. \quad (3.3)$$

Eqn. 3.3 can be seen as a first order approximation of the system of ODE.

3.3.2 The policy model

Policy function Instead of assuming policy functions to be continuous, we consider a more realistic set of policies by assuming policies have different stages. Therefore, we consider policy

functions from a subset of the admissible set $\mathcal{U}_{\alpha_{\max}, T_0}$ in [BDP21]: we define the minimal policy time interval (MPTI) as the minimal duration time during which a policy remains unchanged. This notion assumes that there is a minimal duration time of different stages of a policy. In addition to $u \in \mathcal{U}_{\alpha_{\max}, T_0}$, we assume that every policy u has a minimal policy time interval Δt and in our simulations, the duration of each stage is a multiple of the MPTI. We denote this subset of policy functions as $\mathcal{U}_{\alpha_{\max}, T_0}^{\Delta t}$. In the past, many public health agencies enforced policies for time periods that corresponded to the work week (e.g. seven days) or multiples of this (e.g. one month). For the purpose of this chapter, we assume the MPTI is an integer multiple of seven days. We additionally assume that policy functions take value from a finite number of intensity levels $A [\alpha_{\max}, 1]$, corresponding to different stages of the policy. In the simulations, we use $A = \{\alpha_{\max}, \frac{\alpha_{\max}+1}{2}, 1\}$. As a result, policy functions we consider are in a special form of piece-wise functions.

Cost function At time t , let $u(t) = \alpha$. The cost at time t is defined by:

$$c(\alpha) = \kappa c^{\text{implementation}}(\alpha) + \eta c^{\text{impact}}(\alpha) + (1 - \kappa - \eta) c^{\text{non-compliance}}(\alpha) \quad (3.4)$$

The cost function is a linear combination of three parts:

- (i) the impact cost, which represents the impact of the epidemic on economy and medical system, etc.
- (ii) the implementation cost, which represents the psychological and economic costs of a lock-down.
- (iii) the non-compliance cost, which is a penalty imposed by a policy-maker upon an agent within its jurisdiction for deviating from its recommendation (e.g., a fine or litigation costs).

The implementation cost is a non-increasing function of α and the impact cost function is a non-decreasing function of α . The coefficients $\kappa, \eta \in [0, 1]$. The cost from time t_1 to t_2 is defined as the

averaged integral of the cost function over a total time period T :

$$c_{t_1 t_2}(u) = \frac{1}{T} \int_{t_1}^{t_2} c(\alpha(t)) dt. \quad (3.5)$$

There are different ways to parameterize the cost function. In this chapter, the cost function is parameterized in the following way:

$$c_{t_1 t_2}(u) = \kappa \left(1 - \frac{\int_{t_1}^{t_2} u(t) dt}{T}\right) + \eta R_{t_2}(u) + (1 - \kappa - \eta) \frac{1}{T} \int_{t_1}^{t_2} (u(t) - \pi(u(t)))^2 dt, \quad (3.6)$$

where $R_{t_2}(u)$ is the fraction of the recovered population at time t_2 if policy u is adopted during $[t_1, t_2]$ and $\pi(u)$ is the policy of the agent one level above. The parameterization of the implementation cost and the non-compliance cost are adopted from [JML21]. The impact cost is parameterized as the recovered population at time t_2 to approximate the impact on the medical system since a fraction of the recovered represents the hospitalized population. Assume the cost function over time interval $[t_1, t_2]$ is constant α , the cost can be simplified as:

$$c_{t_1 t_2}(u) = \kappa \frac{t_2 - t_1}{T} (1 - \alpha) + \eta R_{t_2}(\alpha) + (1 - \kappa - \eta) \frac{t_2 - t_1}{T} (\alpha - \pi(\alpha))^2. \quad (3.7)$$

An example of cost functions of different weights using the above parameterization is shown in Fig. 3.3.

3.3.3 The Policy model for a single region

In our simulation for a single region, we use a averaged total cost over a time period T as the following:

$$\begin{aligned} c_{\text{total}}(u) &= \frac{1}{T} \int_0^T c(\alpha(t)) dt \\ &= \kappa \left(1 - \frac{\int_0^T \alpha(t) dt}{T}\right) + \eta R_T(\alpha) + (1 - \kappa - \eta) \int_0^T (\alpha(t) - \pi(\alpha(t)))^2 dt. \end{aligned} \quad (3.8)$$

If at time T , the SIR model has reached the equilibrium, we can use $R_T(\alpha)$ to approximate R_∞ —the fraction of the final size of the recovered population. To find the optimal policy, we solve for the following optimization problem:

$$u(t) = \arg \min_{u'} c_{\text{total}}(u') \quad (3.9)$$

3.3.4 The policy model for multi-layer multiple regions

In our framework, the policy-making goes top-down across layers. For example, in a three-layer model, the federal government makes the decision first, followed by the states, and counties (see Fig. 3.2). Within a layer, the regions make decision in a game-like situation. There are several ways to model this ‘game’. One type of game is that the counties unilaterally make the best decision given the policies of other counties are fixed, whereas in the Pareto game [Deb54], the counties make decisions jointly. An optimal outcome is said be the Pareto efficient if there is no outcome that can increases at least one player’s payoff without decreasing anyone else’s.

There are several differences between our work and [JML21]’s. First, their model is based on Nash equilibrium, where agents make decisions with other agents’ possible actions in mind. We use the idea of ‘learning in game’ [FL98]. We assume that the agents gradually evolve towards the best decisions instead of being optimal instantly. In practice, each region in the game assume other regions’ policies stay the same when optimizing for its cost function. Second, we focus on the dynamics, instead of a snapshot in time.

Network SIR In practice, the counties can hardly be treated as independent. People travel across counties, e.g. work and visiting families and friends. The majority of the literature ([KMS17]) of network style SIR focus on the single people as a node and study the effects of interpersonal network on the pandemics. For example, [MSH12] empirically study how well various centrality measures perform at identifying which nodes in a network will be the best spreaders of disease. [TKH20] explains why most COVID-19 infection curves are linear after the first peak in the context

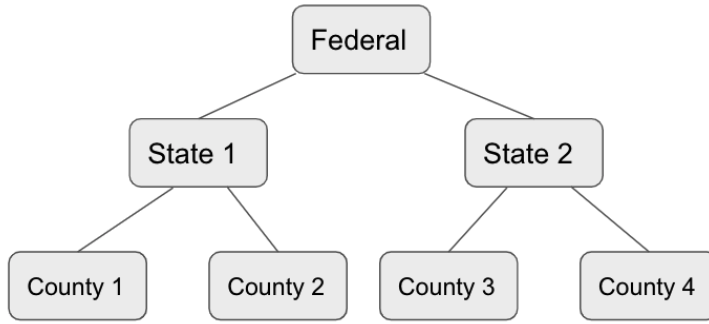


Figure 3.2: An example of a three-layer hierarchical structure.

of the contact network using the SIR model. There are a few works that study the interplay between different geographical regions rather than the interpersonal contact network. In [GKG21], a kernel-modulated SIR model was introduced to model the spread of COVID-19 across counties. The kernel is based on the spatial proximity between regions. Metapopulation epidemic models are based on the spatial structure of the environment, and the detailed knowledge of transportation infrastructures and movement patterns. The metapopulation dynamics of infectious diseases has generated a wealth of models and results considering both mechanistic approaches taking explicitly into account the movement of individuals ([GEG03, KR02, RWL06]). For example, in [RWL06], the authors proposed a multi-regional compartmental model using medical geography theory (central place theory) and studied the effect of the travel of individuals (especially the infected and exposed ones) between regions on the global spread of severe acute respiratory syndrome (SARS). Another way to account for the interplay between regions is to use a cross excitation matrix [YLB19]. This scheme assumes the a uniform mixing of the population across regions and the infected population in one region can trigger the infection in another. The entries of the matrix records the pair-wise cross excitation from one region to another. In this chapter, we assume uniform mixing in the population and use an excitation matrix $K = \{K_{aa'}\}$ to model the travel and infections across counties. Our

network-style SIR is the following:

$$\left\{ \begin{array}{l} \frac{dS_a(t)}{dt} = -\alpha_a \beta \sum_{a'} K_{aa'} \frac{I_{a'}(t) S_a(t)}{N_a}, \\ \frac{dI_a(t)}{dt} = \alpha_a \beta \sum_{a'} K_{aa'} \frac{I_{a'}(t) S_a(t)}{N_a} - \gamma I_a(t), \\ \frac{dR_a(t)}{dt} = \gamma I_a(t), \\ S(0) = S_0, \quad I(0) = I_0, \quad R(0) = R_0. \end{array} \right. \quad (3.10)$$

For any county a , the rate of change from S_a to I_a triggered by $I_{a'}$ depends on $K_{aa'}$, the current fraction of the susceptible S_a in county a and the current fraction of the infected $I_{a'}$ in county a' . Note that $K_{aa} = 1$. When $K = I$, the network SIR is the independent SIR.

Cost function Consider the i -th time interval $[i\Delta t, (i+1)\Delta t]$ and $u(t) = \alpha$ for $t \in [i\Delta t, (i+1)\Delta t]$. Regions adopt the following cost function:

$$c_{i\Delta t, (i+1)\Delta t}(\alpha) = \kappa(1 - \alpha)\Delta t/T + \eta R_T(\alpha) + (1 - \kappa - \eta)(\alpha - \pi(\alpha))^2 \Delta t/T.$$

For the top-layer regions, there is no non-compliance cost and the cost function is

$$c_{\Delta t}(\alpha) = \kappa(1 - \alpha)\Delta t/T + \eta R_T(\alpha),$$

where $\kappa + \eta = 1$.

3.4 Algorithms

We discretize time by MPTI Δt and the policy intensity into multiple levels. Let T be the total time and A be the set of possible policy intensities (e.g., $A = \{0, 0.5, 1\}$).

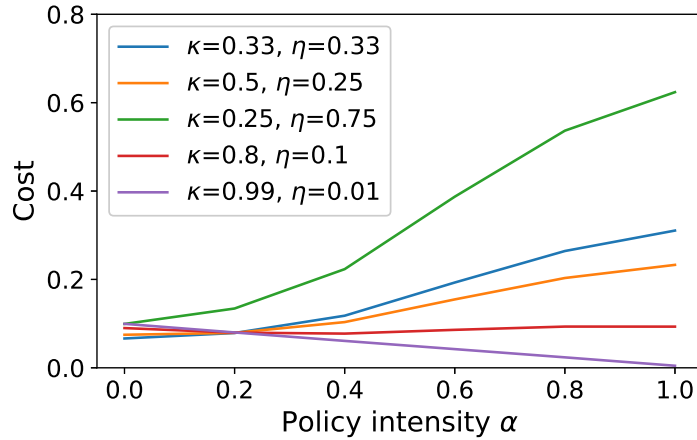


Figure 3.3: Different cost functions vs policy intensity α .

3.4.1 Single region model

We search for all the policies that lead $S_{\text{final}} \in (S_{\text{herd}} - \epsilon, S_{\text{herd}} + \epsilon)$, for some small ϵ using a depth-first search algorithm [Tar71]. The depth-first search algorithm stores the cost up to current time interval and reuse the this result to get the total cost for each policy function through backtracking. Let $N = \frac{T}{\Delta t}$ and N denote the number of stages of a policy. In total, there are $|A|^N$ policies. We initialize the minimal cost c_{min} to be 9999. For n -th time interval ($n < N$), we choose a value from the set intensity levels A that is not used before, calculate the cost for the policy intensity, add it to the previous cost, and move to $(n + 1)$ -th time interval. If the end time interval is reached, check if $S_{\text{final}} \in (S_{\text{herd}} - \epsilon, S_{\text{herd}} + \epsilon)$. If $S_{\text{final}} \in (S_{\text{herd}} - \epsilon, S_{\text{herd}} + \epsilon)$, calculate the cost for the final time interval and add it to the previous cost to get the current total cost c . If the total cost c is smaller than c_{min} , we update c_{min} with the total cost c , and the optimal policy u_{opt} with u . Next, we go back to the previous time interval and repeat the same procedure. After searching over all policies, the policy with the lowest cost is the optimal policy (Alg. 2).

Algorithm 2 SINGLE-REGION POLICY SIR

```
1: Input: Time  $T$ , intensity levels  $A$ , minimal policy time interval  $\Delta t$ , policy end time  $T_0$ ,  
   Tol  $\epsilon$   
2: Initialize county policies, minimal cost  $c_{\min} = 9999$ , current cost  $c = 0$   
3:  $N = \frac{T}{\Delta t}$ ,  $n = 1$   
4: if  $n == N$  then  
5:   for intensity level  $\alpha \in A$  do  
6:     calculate  $S_{\text{final}}$  using the intensity level  $\alpha$   
7:     if  $S_{\text{final}} \in (S_{\text{herd}} - \epsilon, S_{\text{herd}} + \epsilon)$  then  
8:       calculate the cost  $c_{\text{temp}} = C(\alpha)$  for  $N$ -th time interval,  $c += c_{\text{temp}}$   
9:       if  $c \leq c_{\min}$  then  
10:         $c_{\min} = c$ ,  $u_{\text{opt}} = u$   
11:       end if  
12:      $c -= c_{\text{temp}}$   
13:   end if  
14: end for  
15: else  
16:   for intensity level  $\alpha \in A$  do  
17:     calculate the cost  $c_{\text{temp}} = C(\alpha)$  for the  $n$ -th time interval,  $c += c_{\text{temp}}$   
18:      $n += 1$   
19:     repeat line 4–18 until  $n = N$   
20:      $c -= c_{\text{temp}}$   
21:   end for  
22: end if  
23: return  $c_{\min}, u_{\text{opt}}$ 
```

3.4.2 Multiple-layers multiple-regions model

The single region algorithm minimizes over all admissible piece-wise functions, while the multiple-region algorithm only minimizes over every time interval. We use an example of a two-layer model (states and counties) to illustrate the algorithm. At n -th time interval, we first determine the optimal policy intensity that minimizes the cost for each state $C_{n\Delta t, (n+1)\Delta t}^s$ for the period $[n\Delta t, (n+1)\Delta t]$ unilaterally, i.e., assuming other states follow their previous policies. Next, we choose the optimal policy intensity for the counties in the same manner, except the cost functions $C_{n\Delta t, (n+1)\Delta t}^a$ include the non-compliance cost. The full details of the two-layer model is in Alg. 3

Algorithm 3 GAME POLICY SIR

```
1: Input: Time  $T$ , excitation matrix  $K$ , intensity levels  $A$ , time interval  $\Delta t$ 
2: Initialize state, county policies
3: Number of policy stages  $N = \frac{T}{\Delta t}$ ,  $n = 1$ 
4: while  $n \leq N$  do
5:    $t = n\Delta t$ 
6:   while  $t < T$  do
7:     for every state  $s$  do
8:       for every county  $a$  in state  $s$  do
9:         update  $S_a, I_a, R_a$  according to the current policy  $\alpha_a$  and the excitation matrix  $K$ :
10:         $S_a(t) = S_a(t-1) - \alpha_a \beta \sum_{a'} K_{aa'} \frac{I_{a'}(t-1)S_a(t-1)}{N_a}$ 
11:         $I_a(t) = I_a(t-1) + \alpha_a \beta \sum_{a'} K_{aa'} \frac{I_{a'}(t-1)S_a(t-1)}{N_a} - \gamma I_a(t-1)$ 
12:         $R_a(t) = R_a(t-1) + \gamma I_a(t-1)$ 
13:       end for
14:     end for
15:      $t += 1$ 
16:   end while
17:   for every state  $s$  do
18:      $\alpha_s = \arg \min_{\alpha' \in A} C_{n\Delta t, (n+1)\Delta t}^s(\alpha')$ 
19:     for every county  $a$  in state  $s$  do
20:        $\alpha_a = \arg \min_{\alpha' \in A} C_{n\Delta t, (n+1)\Delta t}^a(\alpha')$ 
21:     end for
22:   end for
23:    $n += 1$ 
24: end while
```

3.5 Simulations

In this section, we present the results for both single-region and multiple-region cases. We first compare the results of our discretized algorithm with the results from [BDP21]. Next, we study the second wave (November 2020–May 2021) in the county of Los Angeles. In addition, we present a 3-county example of the multiple regions game and a 3-county example with a state.

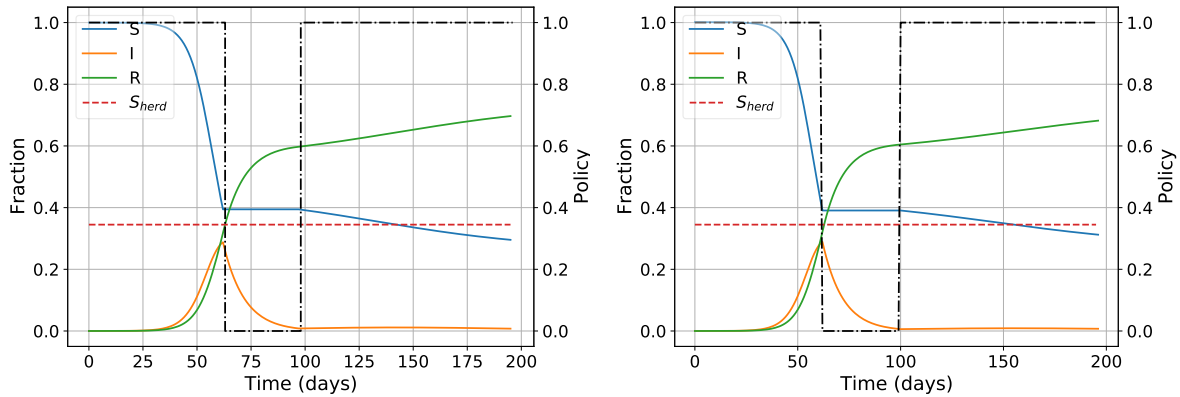
3.5.1 Optimal policy in France

We compare the results from [BDP21] to our model with the same cost function but only three possible levels of policy intensity α . As in [BDP21], the general cost function (3.8) reduces to the

impact cost and is parameterized as the final epidemic size R_∞ . In their work, an assumption is made that the paths considered all reach herd immunity. Therefore, in our search for the optimal policy, we filtered out those cases that do not reach herd immunity. Note that without this constraint, the optimal solution is the strictest policy starting from the beginning of the pandemic, resulting in the least number of infections. For ease of computation, we consider three levels of policy intensity: 0, 0.5, 1 and fixed time intervals for the MPTI. We use the same set of parameters for the SIR model as in [BDP21]: $N = 6.7 \times 10^7$, $I_0 = 10^3$, $S_0 = N - I_0$, $\mathcal{R}_0 = 2.9$. Following [BDP21], we choose the policy end time T_0 as close as possible to 100, thus setting $T_0 = 98$ since the time interval needs to be a multiple of the MPTI of seven days. Our algorithm produces the result in Fig. 3.4a which we compare to the result from [BDP21], shown in Fig. 3.4b. Both solutions exhibit a bang-bang controller. The solution using our model starts the control on day 63 rather than day 61.9 it is a multiple of seven. Slightly more people are infected with the policy that is forced to use seven day intervals rather than a finer gradation in time.

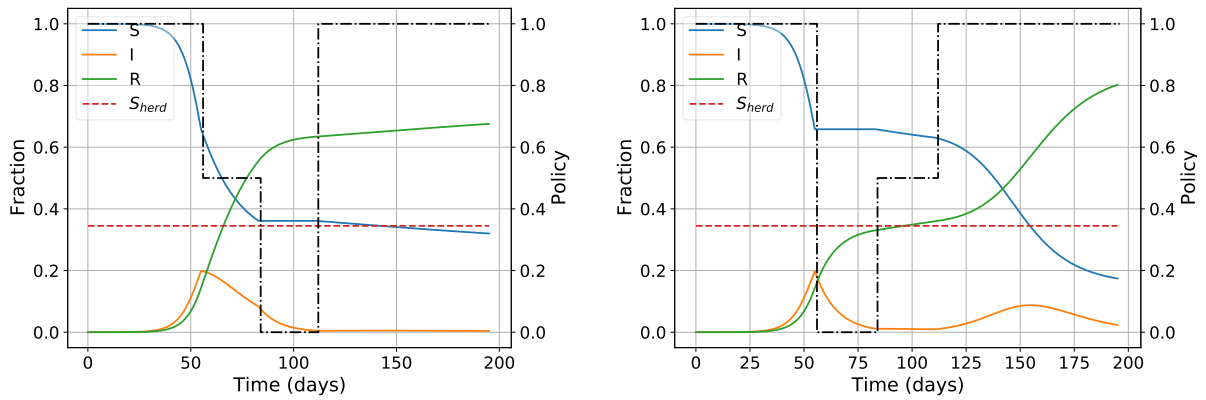
Using a larger minimal policy time interval of 28 days and $T_0 = 112$, the optimal solution is no longer a bang-bang controller, shown in Fig. 3.4c with a larger $S_\infty = 0.32$. The optimal policy starts with a looser policy and then a stricter one. Interestingly, in practice, during COVID-19 it was common for policies to start with the strictest restrictions followed by partial opening [Wikb, Dep22]. Thus it is interesting to contrast the optimal policy with a policy in which the two stages are flipped in time, see Fig. 3.4d. The flipped policy is not the optimal solution—it results in a larger pandemic size and a second wave of infections, as was often seen during the first two years of the COVID-19 pandemic. Nevertheless, the policy in Fig. 3.4d, while infecting more people, divides this population into two waves which could decrease daily hospital demand over the course of the outbreak. Our policy model does optimize for hospital demand. Since many public health agencies considered hospital demand when making policy decisions, it might be interesting to consider it in further studies. The model considered here is idealized and meant to provide some insights for future research on policy rather than to directly create policy change. However future policy-makers may want to review these results in light of the additional infections a policy might

produce.



(a) Optimal policy with the minimal policy time interval $\Delta t = 7$ days, $S_\infty = 0.296$

(b) Optimal policy in [BDP21], $S_\infty = 0.311$



(c) Optimal policy with the minimal policy time interval $\Delta t = 28$ days, $S_\infty = 0.32$

(d) Flipped policy from panel (c), $S_\infty = 0.174$

Figure 3.4: Optimal policy and the SIR model of France from March 17 to May 11 2020.

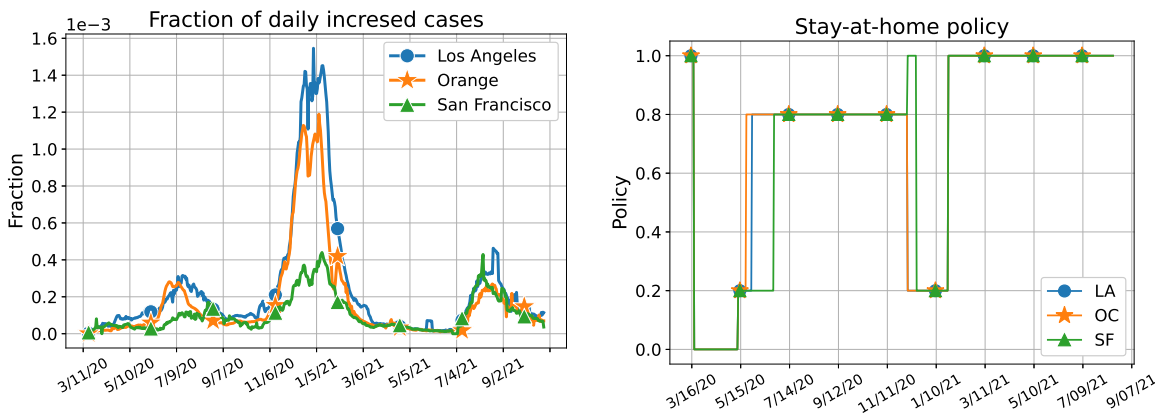
3.5.2 Case study—2nd wave in Los Angeles

We first present the infection in three counties in California and their corresponding ‘stay-at-home’ policy from Mar 2020 to Sept 2021. Fig. 3.5a shows the 7-day rolling average of the fraction of the daily increased infected cases based on the data from [DDG20] in 3 counties with the largest population density in California, namely, San Francisco (SF), Orange county (OC), and Los Angeles (LA). There were 3 major outbreaks during the given time interval. For the first and the second

wave, orange county and LA county have similar situations while SF stayed more contained. Due to the massive travel in the holiday season, the second wave has a much larger size than the first one. In [Cenb], the US Centers for Disease Control and Prevention describes six levels of ‘stay-at-home’ policy. The intensity of the policy decreases as the numerical value increases. The exact descriptions of the five levels of policies and their numerical representation are shown in Table 3.1. Fig. 3.5b shows the change of the intensity of the ‘stay-at-home’ policy during the same period. The policy during the first wave was *proactive*, whereas the policy for the second wave was more *reactive*. Society is vigilant when the virus first shows up, but the fatigue in policy-making and compliance with the masking wearing and stay-at-home order increases the response time.

Numerical value	‘Stay-at-home’ policy
0	Mandatory for all individuals
0.2	Mandatory only for all individuals in certain areas of the jurisdiction
0.4	Mandatory only for at-risk individuals in the jurisdiction
0.6	Mandatory only for at-risk individuals in certain areas of the jurisdiction
0.8	Advisory/Recommendation
1	No order for individuals to stay home

Table 3.1: CDC stay-at-home policies. There are 6 levels of policies and we map the levels linearly onto the interval $[0, 1]$ for simplicity. The numerical value on the left is used to graph actual policies over time in Fig. 3.5b.



(a) The fraction of the daily increase of the infected with a 7-day rolling average.

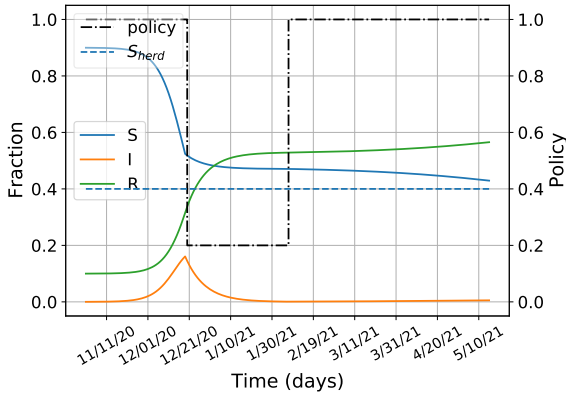
(b) Stay-at-home Policy

Figure 3.5: The fraction of the infected and ‘stay-at-home’ policy over time in Los Angeles, San Francisco, and Orange County.

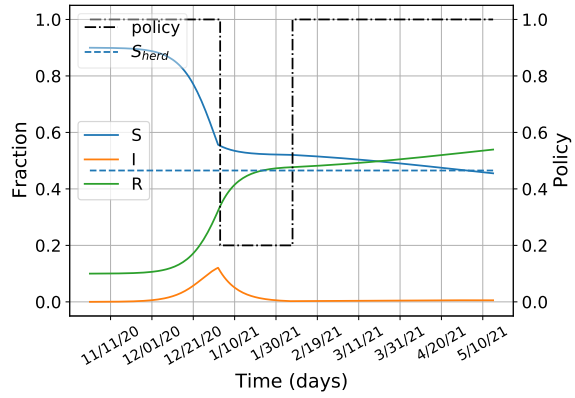
During the second wave, with a relatively strict policy, the regions all stayed below herd immunity. With vaccination available in early 2021, the pandemic in all three regions tapered off.

Now we consider a counterfactual study of how the pandemic would evolve if herd immunity is reached during the second wave, controlled by our policy model, using parameters measured from the Los Angeles data. We choose to study the period of the second wave for several reasons. First, the reporting scheme of COVID-19 is improved compared to the first wave and the data during the second wave is more accurate. In addition, with the experience and knowledge gained from the first wave, it's more likely for the authorities to make optimal decisions. Given that there was no complete lockdown during the second wave, we consider the policy intensity levels $A = \{0.2, 0.6, 1\}$, and use the minimal policy time interval $\Delta t = 7$. We choose 0.2 as our maximal policy intensity because full lockdown was not desirable during this period. We choose a second policy level of 0.6 because it is the average of 0.2 and 1. In all simulations we optimize for final pandemic size and compare the optimal controls found.

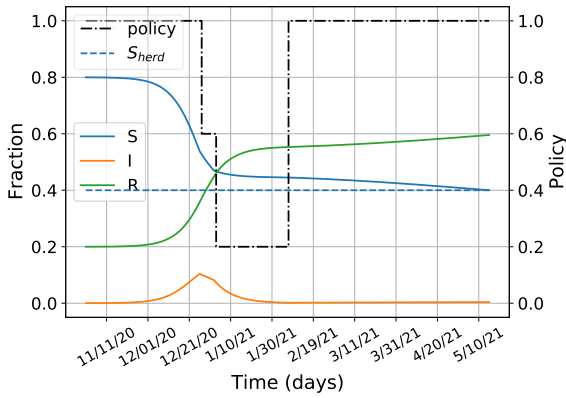
In Fig. 3.6, the left column (Figs. 3.6a, 3.6c, 3.6e) is the simulated SIR with the optimal policy when the basic reproduction number $R_0 = 2.5$ and the initial recovered $r_0 = 0.1, 0.2, 0.3$. The right column (Figs. 3.6b, 3.6d, 3.6f) is the simulated SIR with the optimal policy when the reproduction number $R_0 = 2.15$ and the initial recovered $r_0 = 0.1, 0.2, 0.3$. This value $R_0 = 2.5$ is estimated from the early COVID-19 infected data ([DDG20]) and the $R_0 = 2.15$ is using the infected data from October to early November 2021 ([DDG20]), prior to the second wave. All optimal policies have a bang-bang-like shape. The policy started approximately around the peak of the infected curve, and the resulting dynamics approach herd immunity. For larger values of r_0 , we expect that a shorter period of high intensity policy is needed to reach herd immunity and our results confirm this. Once enough of the population is infected and recovered, a shorter control policy is needed.



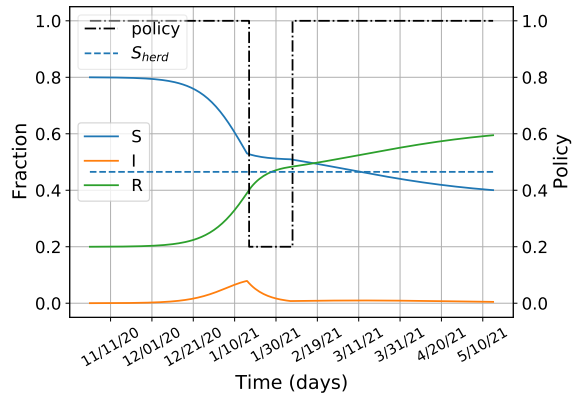
(a) $R_0 = 2.5, r_0 = 0.1.$



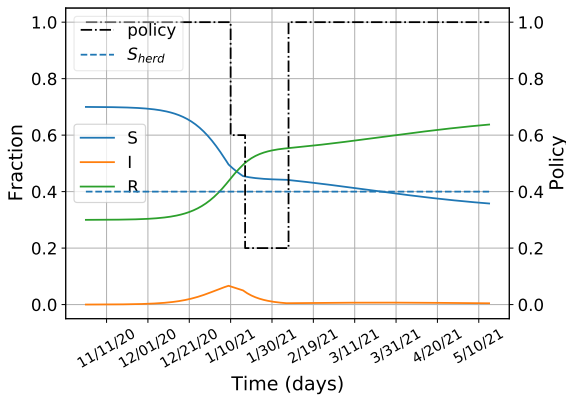
(b) $R_0 = 2.15, r_0 = 0.1.$



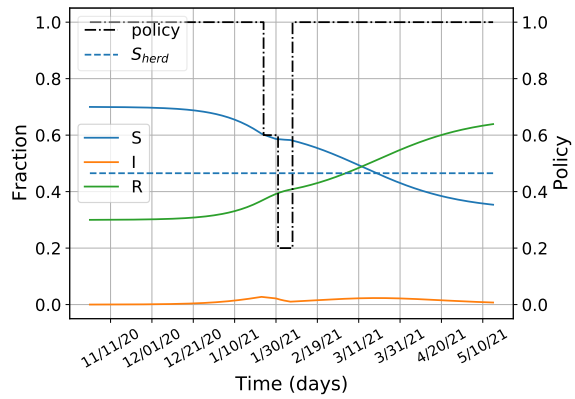
(c) $R_0 = 2.5, r_0 = 0.2$



(d) $R_0 = 2.15, r_0 = 0.2$



(e) $R_0 = 2.5, r_0 = 0.3$



(f) $R_0 = 2.15, r_0 = 0.3.$

Figure 3.6: Optimal policy in Los Angeles with the basic reproduction number $R_0 = 2.5, 2.15$ and the fraction of the initial recovered population $r_0 = 0.1, 0.2, 0.3$.

3.5.3 Multiple regions

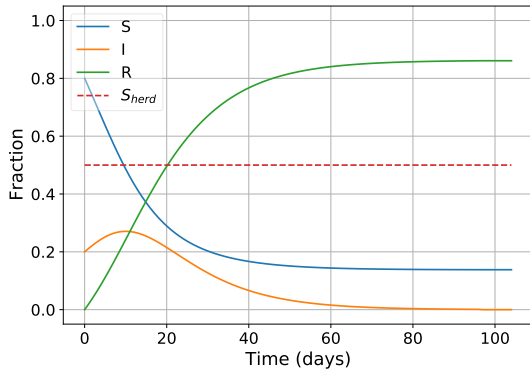
It is interesting to extend this model to the case of multiple regions. First, we discuss when one layer exists, i.e., only counties. The game between the counties is through cross excitation of infection among the counties. Next, we study the case when a governing state is added.

We consider three interacting counties with the excitation matrix K :

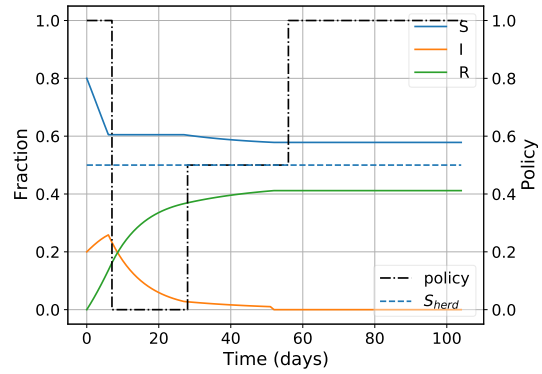
$$K = \begin{bmatrix} 1 & 0 & 0 \\ 0.1 & 1 & 0 \\ 0 & 0.1 & 1 \end{bmatrix}$$

Counties 1, 2, 3 have initial fractions of the infected population as $i_0 = 0.2, 0.1, 0.1$, respectively. This implies that county 1 has a bigger outbreak initially, and part of the infection in county 2 is excited from county 1 and part of the infection in county 3 is excited from county 2. The cost functions for all counties consist of an implementation cost and an impact cost with equal weights ($\eta = \kappa = 1/2$). The minimal policy time interval Δ is set to be 7.

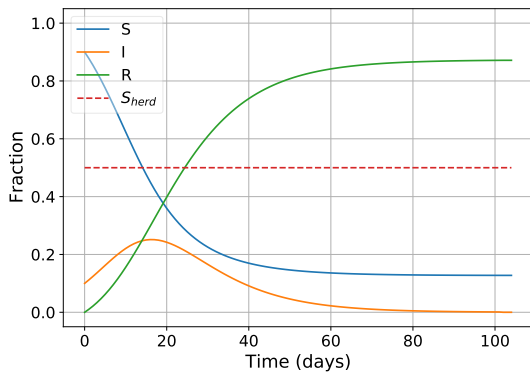
The left column (Figs. 3.7a, 3.7c, 3.7e) are the simulations for the counties without any intervention and the right column (Figs. 3.7b, 3.7d, 3.7f) are the simulations with interventions. Without intervention, we see propagation of waves of infection from county 1 to county 2 and then to county 3. All of the counties reached herd immunity eventually. With interventions, the policies started on day 7 and for county 2 and 3, the infected curves decrease before reaching their peaks. With the control, county 1 contained the pandemic and the final S_∞ is close herd immunity level S_{herd} . With a fewer infected population to begin with, county 2 and 3 contained the pandemic before reaching herd immunity. Fig. 3.8d shows the results of adding a governing state on top of the county layer. We keep the ratio of the weights for the implementation cost and the impact cost to be 1:1, the same as in the no-state case in Fig. 3.7. The state has slightly different weights, with the ratio of the weights for the implementation cost and the impact cost being 1:2. Compared to Fig. 3.7, by adding a state, the three counties ended up with the same policy. In this case, the



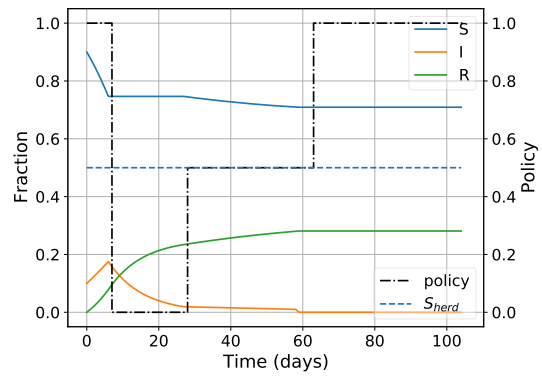
(a) County 1. No intervention.



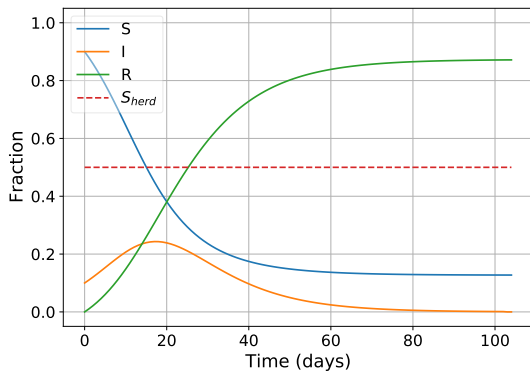
(b) County 1.



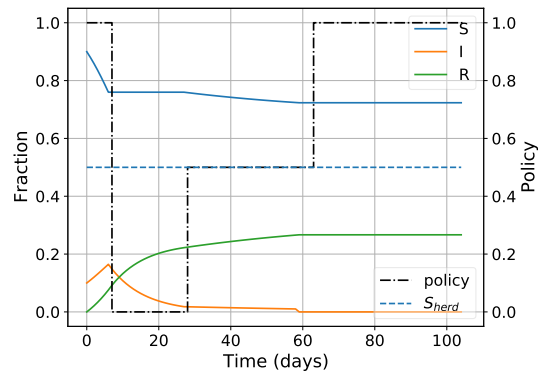
(c) County 2. No intervention.



(d) County 2.



(e) County 3. No intervention.



(f) County 3.

Figure 3.7: An example of three dependent counties without and with interventions. With intervention, for all counties, the coefficients for the implementation cost $\kappa = \frac{1}{2}$ and the coefficients for the impact cost $\eta = \frac{1}{2}$. The minimal policy time interval $\Delta t = 7$.

noncompliance cost results in each county choosing the same policy as the state rather than different policies.

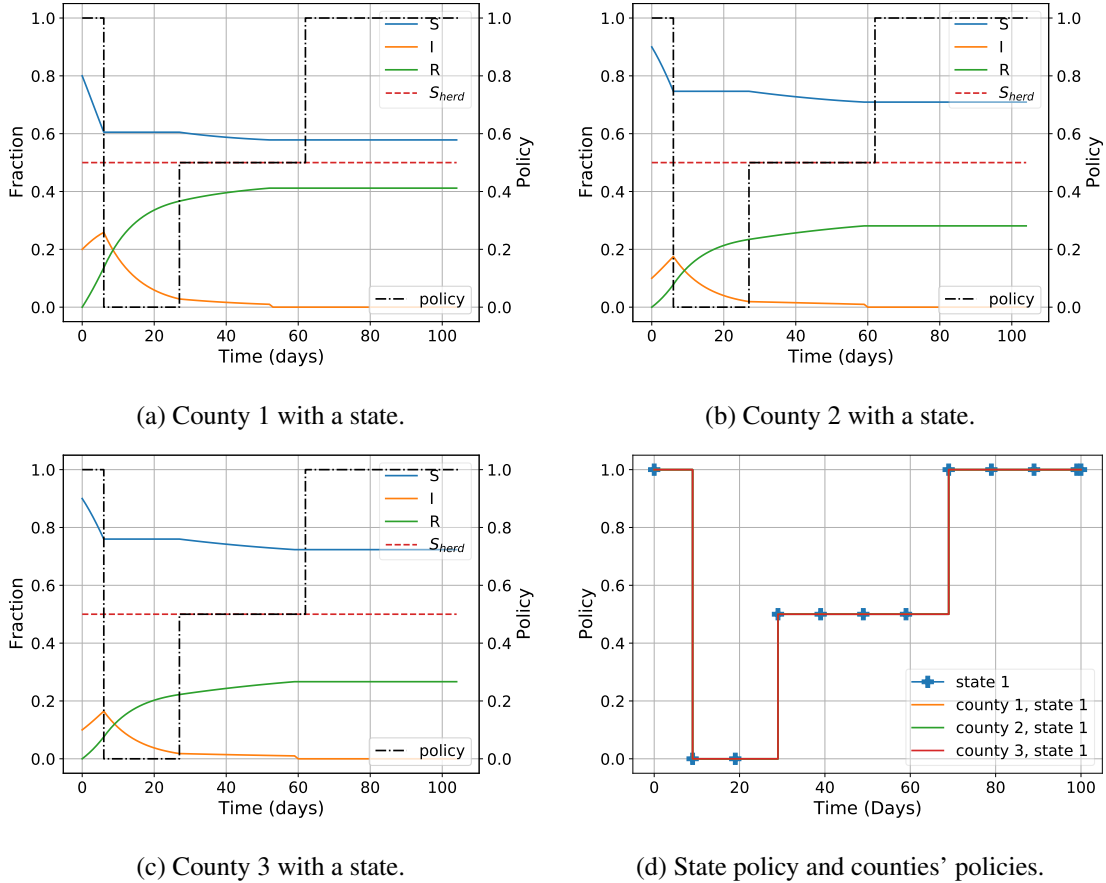


Figure 3.8: An example with 3 counties and a governing state. For all counties, the coefficients for the implementation cost $\kappa = \frac{1}{6}$, the coefficients for the impact cost $\eta = \frac{1}{6}$ and the coefficients for the impact cost $\eta = \frac{1}{2}$ and the coefficients for the non-compliance cost $1 - \kappa - \eta = \frac{2}{3}$. For the state, the coefficients for the implementation cost $\kappa = \frac{1}{3}$, the coefficients for the impact cost $\eta = \frac{2}{3}$. The minimal policy time interval $\Delta t = 7$.

3.6 Discussion and future work

We propose a policy-making model coupled with the SIR model to study a single region and game-like interactions between multiple regions. The model demonstrates its ability to model real-life situations with different sets of parameters in both one region and multiple regions scenarios. One

can extend the model to a hierarchical structure by building multiple layers of the multiple regions model and study the cross-layer effects.

In the search for the optimal policy, we use a naive depth-first search algorithm for the one-region model. One can speed up the algorithm by removing some of the obvious non-optimal paths.

In our model, the policy intensity α is a heuristic representation of the lockdown, social distancing and mask policy. It remains to be discussed how other policies, for example, vaccination policies, affects the spreading in the different stages of the pandemic. The model ignores some of the important features like the limitation of the hospital capacity, which could be added as constraints when minimizing the cost function. Fig. 3.5b shows the policy for the first wave is proactive while the one for the second wave is reactive. This implies that as the pandemic goes on, the fatigue of policy-making start to show up. So far, the model does not have the capability of modeling this fatigue. In the future, one could consider an adaptive term in the cost function to model it.

CHAPTER 4

COVID-19 Literature Topic-Based Search via Hierarchical nonnegative matrix factorization (NMF)

This chapter is adapted from an original paper [[GHH20](#)] that I co-authored with Rachel Grotheer, Longxiu Huang, Yihuan Huang, Alona Kryshchenko, Oleksandr Kryshchenko, Pengyu Li, Elizaveta Rebrova, Kyung Ha, and Deanna Needell. The problem and approach were suggested by Deanna Needell . I preprocessed the data set, developed and implemented the numerical method for HNMF and analysed topic similarities with the rest of the collaborators.

4.1 Background

The appearance of the novel SARS-CoV-2 virus on the global scale has generated demand for rapid research into the virus and the disease it causes, COVID-19. However, the literature about coronaviruses such as SARS-CoV-2 is vast and difficult to sift through. This chapter describes an attempt to organize existing literature on coronaviruses, other pandemics, and early research on the current COVID-19 outbreak in response to the call to action issued by the White House Office of Science and Technology policy [[SP](#)] and posted on the Semantic Scholar [[Sch20](#)] and [[Kag20](#)] websites. The original dataset posted on that site is augmented by adding articles drawn from other databases in order to make the final interactive organizational structure more robust for researchers.

Our primary goal is to create a framework for a topic-based search of papers within this dataset that is helpful to those investigating the novel coronavirus, SARS-CoV-2, and the global COVID-19

pandemic. In order to discover the latent topics present in the collection of scholarly articles, as well as to organize them into a hierarchical tree structure that allows for an interactive search, we use a modified hierarchical nonnegative matrix factorization (HNMF) approach. A website¹ that allows users to walk through the topic tree based on the top keywords associated with each topic is created using this hierarchical organization of the papers.

4.1.1 Contributions

Our methods help make sense of a vast and rapidly growing body of COVID-19 related literature. The main contributions of this chapter are as follows:

- A diverse dataset of COVID-19 related scientific literature is compiled, consisting of articles with full-text available drawn from several online collections.
- A tree-like soft² cluster structure is created of all the papers in the dataset based on inherent relations between their topics using hierarchical NMF.
- The best number of topics for each layer is defined as the number that produces the most consistent clustering of the dataset with random initializations of NMF algorithm. A variance analysis method is used to identify the best number of topics on each layer.
- The effectiveness of the method is measured by exploring the coherence of each topic and dissimilarity between the topics.
- The discovered topics and distribution of articles into each of the topics are discussed, revealing major areas of interest and research in the early months of the pandemic, as well as how existing epidemic literature can be effectively organized to allow efficient comparison to COVID-19 related research.

¹<http://covid-19-literature-clustering.net/>

²*Soft* here means that clusters can intersect, as one paper could belong to more than one topic.

- The theoretical results are complemented with an interactive website:

<http://covid-19-literature-clustering.net/>

4.1.2 Related work

Some relevant works that motivate our approach are briefly reviewed. NMF was first proposed for document clustering [XLG03], and since then many variants of the NMF method have been proposed and applied to help organize various types of data [LS99, Buc08, KCP15]. In particular, there exist several recent papers that use NMF to find a hierarchy of topics in a set of documents. For example, [KP13] apply a rank-2 NMF to the recursive splitting of a text corpus and also provide an efficient on-the-fly stopping criterion. In [GHM19], the authors discuss a different version of HNMF, when the hierarchy of topics is generated by aggregation of the topics (rather than splitting). The first application of NMF produces the initial set of the most refined topics, and the subsequent NMF iterations find supertopics in which the previous set of topics can be summarized. This approach is referred as a *bottom-to-top* viewpoint, and the former as a *top-to-bottom*. Approaches that utilize tools from neural networks such as back propagation to improve the topic representations have also been developed recently [TBZ16, RHW15, SNT17, GHM19]. In [TCL18] a hierarchical online non-negative matrix factorization method (HONMF) is proposed to generate topic hierarchies from data streams. The proposed method can dynamically adjust the topic hierarchy to adapt to the emerging, evolving and fading process of the topics. This work most closely aligns with what we present here, and although we do not consider the online setting, our method can easily be adapted to such.

Finally, several authors have sought to address the issue of interpretability of topics discovered by NMF, especially in datasets comprised of text documents. For example, in [ASN17], the authors apply NMF to the documents using a word embedding model, *Word2Vec* [MSC13], that focuses on the semantic relationship between words. We make use of this embedding to analyze the usefulness of the topics generated by examining their semantic similarity.

4.2 Data description

The dataset used is compiled from 4 different databases that contain scholarly articles related to COVID-19, various coronavirus diseases, other infectious diseases, and epidemiology [Sch20, DP, Bio20b, bio20a]. From each of these databases, only articles written in English that have a complete abstract and text body available are included. Punctuation and words on the NLTK English stopwords list [BKL09] are removed from the text body and abstract of each article. An initial application of NMF generated a topic consisting primarily of words that, upon further investigation, were found to be part of the copyright and publishing information present at the top of articles primarily drawn from the bioRxiv database, and not the content of the articles themselves. Therefore, we also remove the top 30 keywords of that topic from the corpus. See Appendix 4.7.1 for the list of these removed keywords. Finally, the articles are lemmatized and each word in the text body and abstract is represented by a TD-IDF embedding [SB88].

After processing and cleaning, the final dataset contains 25,663 articles. Most of these databases are regularly updated and one of the important future directions of this work will include developing a dynamic tree structure that pulls new articles from these databases weekly.

4.3 Hierarchical NMF for topic detection

In a vector space model, a corpus can be represented by a $d \times n$ matrix X , where d is the size of the vocabulary, and n is the number of documents. The underlying assumptions in topic modeling using NMF [BL07] are that a latent topic can be represented as a distribution over the words, and that every document is a mixture of topics, i.e. comprises a statistical distribution of topics that can be obtained by “adding up” all of the distributions of all the topics covered. In this section, we will introduce how to apply hierarchical NMF for topic detection and creation of the hierarchical tree structure. As a preliminary step, a brief introduction to using NMF for topic detection is given.

4.3.1 NMF for topic detection

In NMF, the bag-of-words model is used to format the documents into inputs. It takes in the corpus and creates a vocabulary out of each unique word in the corpus. It then models each document as a vector with length equal to the number of words in the vocabulary where entry i in the vector corresponds to how many times word i occurred in the document. Thus, the bag-of-words model gives us a matrix X with dimension $d \times n$ where n is the number of words in the vocabulary, and d is the number of documents. Then corpus matrix $X \in \mathbb{R}_{\geq 0}^{d \times n}$ is decomposed into a pair of low-rank nonnegative matrices $W \in \mathbb{R}^{d \times k}$, also known as the dictionary matrix, and $H \in \mathbb{R}^{k \times n}$, also known as the coding matrix, by solving the following optimization problem

$$\inf_{W \in \mathbb{R}_{\geq 0}^{d \times k}, H \in \mathbb{R}_{\geq 0}^{k \times n}} \|X - WH\|_F^2, \quad (4.1)$$

where $\|A\|_F^2 = \sum_{i,j} A_{ij}^2$ denotes the matrix Frobenius norm. Solving NMF using iterative optimization algorithms, has a drawback: the objective function is usually non-convex and has multiple local minima. Therefore a different random initialization of the NMF procedure will result in a different matrix factorization. More importantly, this changes the interpretation of the results, including topic vector representations (W) as well as the relevance between articles and topics (H). Another possible source of variability in the algorithm is the choice of the number of topics, k . Different combinations of initializations of W , H , and k yield different topics, leading to different article clustering results. See Section 4.5.1 for more discussion and implementation details in this vein.

4.3.2 Hierarchical NMF

The traditional NMF method treats the detected topics as a flat structure, which limits the ability of the representation of such method. In contrast, a hierarchical NMF (HNMF) framework is able to detect supertopics, subtopics, and the relationship between them, creating a tree structure. Compared with traditional NMF, HNMF improves topic interpretability. For instance, while both NMF and

HNMF may produce topics that are similar to one another, if these topics are in the bottom layer of the tree structure provided by HNMF, their associated supertopics provide additional context to help distinguish the related subtopics. Besides improving topic interpretability, HNMF also provides a more user-friendly search framework, which is suitable for building website. The hierarchy allows users to search for relevant topics more effectively, while progressively narrowing their search.

Given the complex nature of the coronavirus literature corpus, such a hierarchical approach is appealing. Thus, we apply the HNMF algorithm summarized in Algorithm 4. Note that this algorithm is similar to the one in [TCL18], which has been shown to be effective for topic detection.

In HNMF, NMF is first applied to the original corpus matrix X to obtain the dictionary matrix W and coding matrix H . The documents are then sorted into matrices X_1, X_2, \dots, X_k , each representing a different topic, according to the coding matrix H , or into the matrix X_e that temporarily holds unassigned articles. Whether the leaves need to be further divided depends on the number of the documents in each topic matrix (leaf). If the number of documents sorted into a topic is greater than a pre-specified value m , then a further division is needed. The above process is repeated until the number of documents in each leaf is less than m . More details on the implementation of the HNMF algorithm are provided in Section 4.5.2.

4.4 Discussion of results

This section begins with a discussion and visualization of the hierarchical tree structure obtained using Algorithm 4. Then in Sections 4.4.3 and 4.4.4 quantitative evidence is provided that the discovered topics are reasonable. In doing this, we seek to measure both the rationality of a given topic and the similarity between topics to evaluate whether the topics differ enough to be useful for a user.

Algorithm 4 Hierarchical NMF

- 1: **Input:** Corpus matrix X .
 - 2: $[W, H] = \text{NMF}(X, k^*)$ where topic number k^* is chosen by Algorithm 5
 - 3: assign articles to the related topics X_1, \dots, X_{k^*} according to the threshold α in H , and any remaining articles to “Extra Document” matrix X_e
 - 4: **while** # of the articles assigned to a topic $i > m$ **do**
 - 5: determine the # of sub-topics k_i^* of the topic i in X_i by Algorithm 5
 - 6: $[W_i, H_i] = \text{NMF}(X_i, k_i^*)$
 - 7: assign the documents to the topics by the a threshold α in H_i s
 - 8: assign the rest to X_e
 - 9: **end while**
 - 10: **for** article x_i in X_e **do**
 - 11: calculate cosine similarity between x_i and leaves, and assign the article to the most related leaf
 - 12: **end for**
 - 13: repeat both while and for loops until the number of the articles assigned to each topic is less than m .
-

4.4.1 Topic visualization

Implementation of Algorithm 4 on the dataset results in a hierarchical clustering of the articles into eight supertopics, each with five to six subtopics. Two of these subtopics, the first and fourth subtopics of supertopic 7, are further decomposed into a third layer of subtopics as the number of articles assigned to the first and fourth subtopics are larger than the selected m in Algorithm 4. The full hierarchical tree structure is visualized in the diagram in Figure 4.1. Each color represents one of the eight supertopics and the size of each slice is proportional to the number of articles that are clustered into that topic. It is important to note that only the top three words associated to each topic are shown due to space constraints, but in some cases extending the list of highly related words is necessary to clarify the difference between the subtopics. For reference, the top ten keywords associated with each topic and subtopic can be found in Appendix 4.7. Additionally, the five most probable words associated to each topic are displayed on the associated website to aid users in more effectively choosing the topics of personal interest.

In order to examine the structure in more depth, Figure 4.2 displays a branch of the resulting tree represented by word clouds, generated from the top five words associated with each topic. The

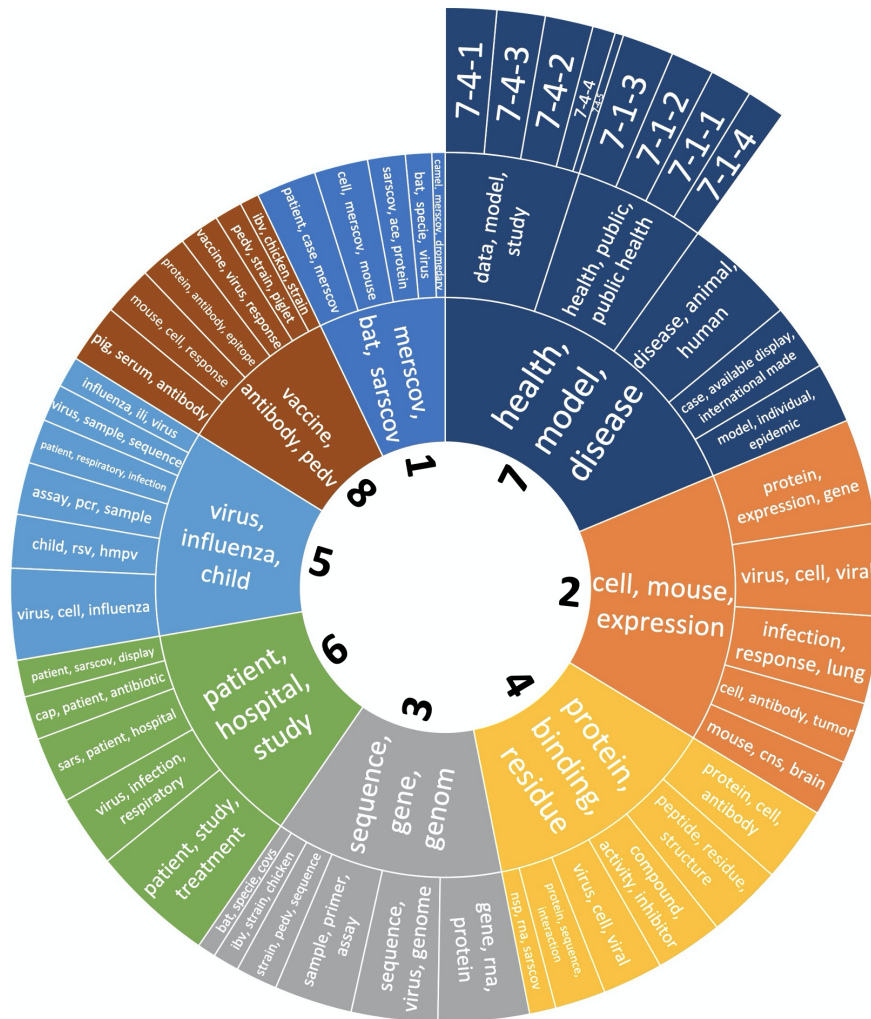


Figure 4.1: Sunburst Diagram of the complete hierarchical structure. The top three relevant words per topic are shown. The area of each region is proportional to the number of articles in that topic. See appendix for the keywords associated with the third layer. The inner circle numeric labels are corresponding to topic number in Figure 4.2

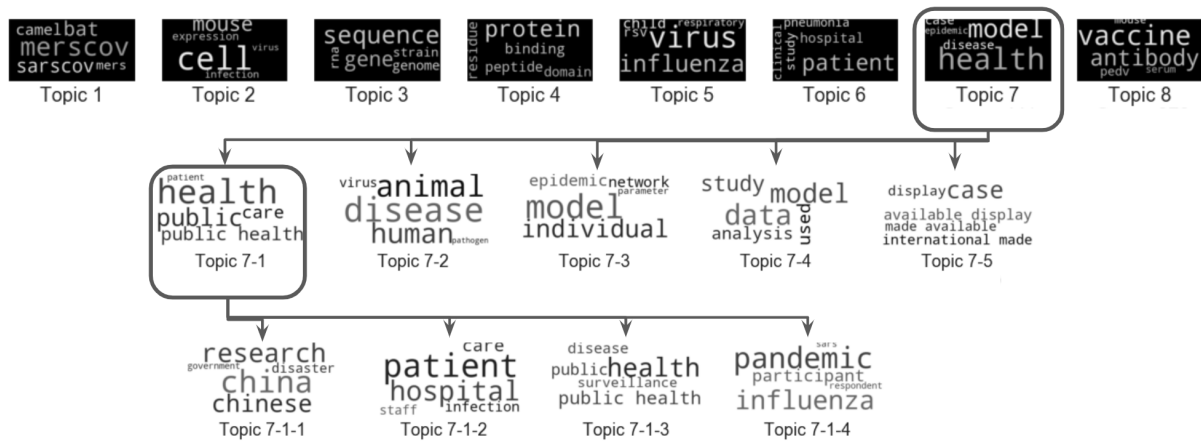


Figure 4.2: Part of Topics from HNMF and related topic coherence: The first row shows the the key words for the topics in the first layer, the second row shows the subtopics of Topic 7 and the subtopics of Topic 7-1 is showed in row 3. Corresponding *topic coherence* score (see Section 4.4.3 for more details) is underneath each word cloud.

size of the words in each word cloud cell are proportional to their weight in the corresponding W matrices, and thus, the probability they are associated with that topic. In particular, the figure follows one path down the tree structure, focusing on Topic 7 and its associated subtopics, and then continuing to the subtopics of Topic 7-1. When moving to deeper layers in the tree, the general “health” and “model” topic further differentiates into subtopics ranging from public health to animal to human transmission diseases, and data modeling. Finally, the public health subtopic leads to clusters of articles specifically related to China or hospital care, for example.

4.4.2 Discussion of topics

Perhaps not surprisingly, the topic to which the highest number of articles are assigned, Topic 7, is about the general study of the disease (with the most highly associated words being “health, model, disease, case, epidemic, outbreak, public, country, population, transmission”), further split into two additional layers of subtopics. This is the only topic that was split into a third layer, allowing a more effective differentiation between articles covering a similar topic.

Also unsurprisingly, much of the literature, which was compiled early on during the pandemic,

is clustered around the study of other coronavirus-caused diseases. Topic 8, for example, focuses on vaccine development through the lens of the Porcine Epidemic Diarrhea Virus (PEDV). Although this is a coronavirus found only in pigs, several vaccines have been developed, especially within the last seven years, when PEDV was first discovered in North America [GZ17]. Hence, it is reasonable that this topic would be of interest to current researchers looking to develop a vaccine for SARS-CoV-2. Similarly, Topic 1 focuses on coronaviruses known to infect humans, such as SARS-CoV, and MERS-CoV. Topic 4 also contains a couple of subtopics that look specifically at the genetic structure of SARS-CoV.

Other topics of interest focus on articles about diseases with related symptoms, although they may be caused by a different type of virus. For example, both Topics 5 and 6 examine literature related to respiratory illnesses such as influenza, though Topic 5 clusters articles more related to laboratory study and Topic 6 clusters articles more related to hospital studies and patient care.

Other major topics focus more on microbiology, including the genomic structure of the virus, the cellular infection and immuno-response, and cell-protein interaction. Thus, the hierarchical tree structure separates papers between macro- (public health) and micro- (biological) studies of the virus, and into papers that study related viruses. This creates a clear delineation of topics for those investigating papers, and gives insight into areas of interest for early researchers of SARS-CoV-2. This organizational structure appears to be more robust and high-level than e.g. a keyword based search or organization.

4.4.3 Topic coherence

One measure of effectiveness of the topics discovered by HNMF is *topic coherence*. Topic coherence is a quantitative measure of how well the keywords that define a topic make sense as a whole to a human observer and collectively provide a consistent interpretation of the topic.

While many topic coherence measures have been proposed, [RBH15] found that the C_V coherence metric correlates the most closely with evaluation by human experts. The C_V measure

calculates the similarity between two words w_i and w_j using the normalized pointwise mutual information (NMPI) metric defined as,

$$\text{NPMI}(w_i, w_j)^\gamma = \left(\frac{\log \frac{\mathcal{P}(w_i, w_j) + \epsilon}{\mathcal{P}(w_i) \cdot \mathcal{P}(w_j)}}{-\log(\mathcal{P}(w_i, w_j) + \epsilon)} \right)^\gamma \quad (4.2)$$

where $\mathcal{P}(w_i)$ and $\mathcal{P}(w_i, w_j)$ are probabilities defined as the number of documents in which either w_i or (w_i, w_j) , respectively, appear, divided by the total number of documents. These probabilities are calculated using a sliding Boolean window of size s that slides over a document at the rate of one word per step. The sliding window allows for the proximity of the words to be taken into account. The γ allows for more weight to be placed on higher NPMI values. After the NPMI score is calculated between each of the top N words, W' , in each topic $W = \{W_1, W_2, \dots, W_N\}$ and each of the remaining $N - 1$ words, W^* , these scores are added together to form a context vector $\vec{v}(W')$. Using the notation given by [SS17], who applied the C_V metric to topics found using latent Dirichlet allocation (LDA), we define the context vector as,

$$\vec{v}(W') = \left\{ \sum_{w_i \in W'} \text{NPMI}(w_i, w_j)^\gamma \right\}_{j=1, \dots, N} \quad (4.3)$$

Finally, the cosine similarity between all context vector pairs within $S_i = (W', W^*)$ is calculated, giving the confirmation measure ϕ_{S_i} ,

$$\phi_{S_i}(\vec{u}, \vec{w}) = \frac{\sum_{i=1}^N u_i \cdot w_i}{\|\vec{u}\|_2 \cdot \|\vec{w}\|_2} \quad (4.4)$$

which is a measure of how well word W' in topic W is supported by the word W^* relative to all the words in W .

To further support these results, we additionally calculate the coherence score defined by [MWT11]

Topic	C_V	C
1	0.68	321
2	0.63	442
3	0.68	407
4	0.56	419
5	0.61	405
6	0.66	420
7	0.63	441
8	0.59	378

Table 4.1: The coherence scores based on both the C and C_V metric for each of the 8 topics in the first layer of the tree.

for each topic. The coherence score C_i for topic i , $i = 1, \dots, k$ is given by,

$$C_i(W^{(i)}) = \sum_{p=2}^N \sum_{\ell=1}^{p-1} \log \frac{P(w_p^{(i)}, w_\ell^{(i)}) + 1}{P(w_\ell^{(i)})}. \quad (4.5)$$

The topic coherence scores for each of the topics in the first layer, using both the C_V and C metrics are in Table 4.1.

The C_V coherence metric has values between 0 and 1, with values closer to one indicating that the keywords form a topic that would be highly ranked by human expert. A positive, large coherence score using the C metric indicates the same. A coherence score that is close to 0 (for C_V) or negative (for C) indicates that a topic is less meaningful, which may occur, for example, if the associated keywords fall into two unrelated groups, or if the keywords are seemingly random and have no obvious connection. Most of our identified subtopics have coherence scores whose values suggest that they are understandable and useful to human users. The C_V scores for each of the subtopics can be found in Appendix 4.7.2.

4.4.4 Topic similarity

Another test of the usefulness of the hierarchical structure generated is to evaluate whether the topics are different enough to allow for informative choice between them. To evaluate this, we quantify topic similarity using a metric known as the Word Mover's Distance (WMD). WMD is a popular

tool for measuring distances between documents [KSK15]. WMD utilizes *Word2Vec* [MCC13], a word embedding technique, and treats each document as a set of vectors in the embedded vector space. This embedding allows the WMD metric to consider the semantic meaning of a given word, rather than just its spelling. Thus, for example, it allows for identification of synonyms as having the same meaning in a given context despite being different words, which makes it more preferable than traditional metrics such as cosine similarity or Euclidean distance. The distance between two documents A and B is defined as the minimum cumulative distance that words from document A need to travel to match exactly the words of document B. We note that while there are other state of the art semantic representations, such as BERT [DCL19] and ELMo [PNI18], and associated metrics, since the topics extracted are a bag of words with weights, the WMD with Word2Vec is sufficient for our purposes.

The topic similarity across the layers and within each layer is evaluated by computing the WMD between a topic and its associated subtopics and between the subtopics themselves, where each topic is represented by its 100 most related words. The similarities between all topics in the hierarchical structure obtained from HNMF is visualized in the heat map in Figure 4.3. As indicated by the overall dark colors, in general each topic in the tree is dissimilar from the others.

When examining the similarities between a topic and its subtopics, results show that for a given topic, its subtopics are less correlated with each other than with their parent topic. For example, in Figure 4.4, for Topic 7, the similarity scores between its subtopics are much lower than the scores between subtopics and their parent Topic 7. Similar results can be drawn for Topic 7-1 and its subtopics, as shown in Figure 4.5.

However, there are some high similarity scores between subtopics that belong to different topics, for example the light off-diagonal spot in Figure 4.3 showing the similarities between Topics 6-3 and 5-3. Examining the top ten keywords associated with each topic, we find that both topics are associated with the words “influenza”, “virus”, and “study” indicating that both topics deal with studies related to the influenza virus.

The insight into the difference between the two subtopics comes from examining supertopics

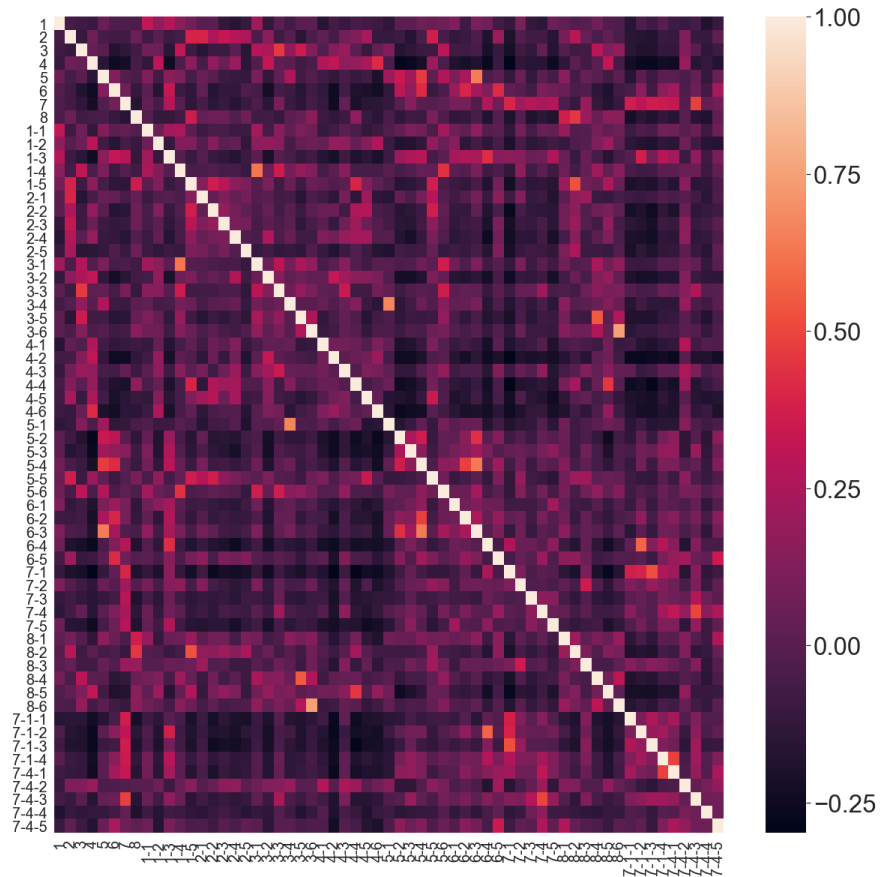


Figure 4.3: Topic similarity for all the topics from HNMF measured by WDM . A dark color indicates the topics are dissimilar, while a light color indicates high similarity. Note that the topics are listed from first layer to third layer from top to bottom or right to left on the vertical and horizontal axes, respectively.

5 and 6 and the keywords associated with each subtopic that do not overlap. Looking at words such as “detection” and “assay” associated with Topic 5 and “surveillance”, “case”, “season”, and “year” associated with Topic 5-3, it appears that Topic 5-3 is more associated with detecting and monitoring the prevalence of cases of influenza in the general populace in a given flu season. On the other hand, the presence of keywords “patient”, “hospital”, “clinical”, and “study” associated with the parent topic, Topic 6, as well as “patient”, “child”, and “respiratory” associated with Topic 6-3, it seems that Topic 6-3, while also related to influenza studies, deals more specifically with cases in a hospital setting, perhaps specifically related to children, and examining the relationship with respiratory illness in general.

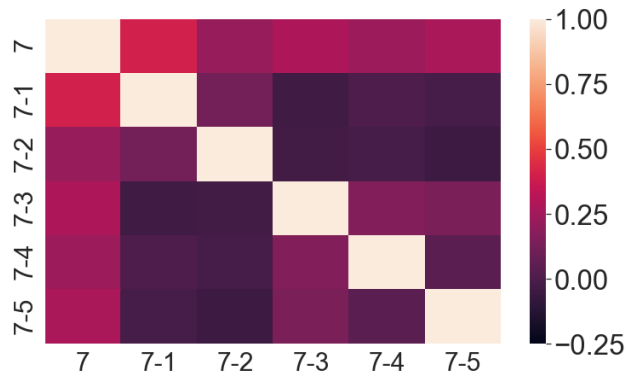


Figure 4.4: Topic similarity between Topic 7 and its subtopics measured by WDM: Topic 7 has high topic similarity with its five subtopics (7-1, 7-2, 7-3, 7-4, 7-5) and the five topics have low similarity between themselves.

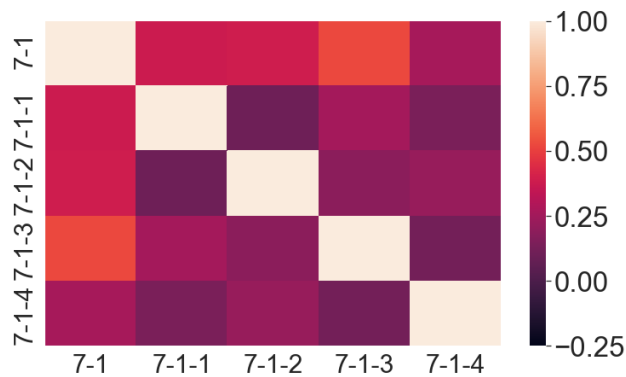


Figure 4.5: Topic similarity between Topic 7-1 and its subtopics measured by WDM: Topic 7-1 has high topic similarity with its four subtopics (7-1-1, 7-1-2, 7-1-3, 7-1-4) and the four topics have low similarity between themselves.

A study of similar subtopics such as these show the effectiveness of the tree in separating related topics into more dissimilar supertopics to make navigation to articles of interest clear. However, Algorithm 4 allows for an article to be assigned to more than one subtopic, acknowledging that a single article may of equal interest to researchers investigating different, but related topics.

4.5 Implementation

In this section, we discuss the details of the implementation of HNMF and the construction of the hierarchical structure.

Algorithm 5 Determine optimal number of topics

- 1: **Input:** integer q , corpus matrix X
 - 2: Determine a range for the potential topic number $[k_1, k_2]$ by plotting increment in variance explained by adding one more cluster to X
 - 3: Randomly select $q + 1$ seeds for initialization
 - 4: **for** integer k in $[k_1, k_2]$ **do**
 - 5: Generate topic sets $\{T_j\}_{j=1}^{q+1}$ from NMF initialized by random seed j
 - 6: Generate S_{kj} for $j = 1, 2, \dots, q$ where S_{kj} is the cosine similarity matrix between topics in T_j, T_{j+1}
 - 7: **for** $S_{kj}, j = 1, 2, \dots, q$ **do**
 - 8: $LSS_k = \emptyset$
 - 9: Add $lss = \min(\max(s_{a.}), \max(s_{.b}))$ to LSS_k , where s_{ab} is the (a, b) th entry of the matrix S_{kj}
 - 10: **end for**
 - 11: **end for**
 - 12: **return** $k^* = \arg \max_k(\text{median}(LSS_k))$
-

4.5.1 Determining number of topics in each layer

As previously discussed, the latent topics discovered by NMF are sensitive to the initial state of the algorithm, leading to different dictionaries for each topic. In order to reduce this sensitivity, we seek to find an appropriate number of topics, k^* , in each layer such that if a k^* -topic NMF is initialized using any two random seeds, the content in the topics discovered should be similar, as measured by cosine similarity. We define this as a *consistent* number of topics. Algorithm 5 summarizes the process to find the “best” number of topics, as defined in this manner, for a corpus matrix X .

In Algorithm 5, first the increment in proportion of variance explained by adding one more cluster to split the corpus matrix X is plotted. This is calculated by looking at the singular values of X . By examining this plot (Figure 4.6), a range $[k_1, k_2] = [7, 11]$ in which a potential optimal number of topics, k^* can be found is obtained by noting where the proportion of variance explained starts to level off.

To determine the value of k^* in this range, first, $q + 1$ random seeds are randomly selected, where q is a sufficiently large number. In this case, $q = 30$ was used. For each number of topics $k \in [k_1, k_2]$, topic sets are generated $\{T_j\}_{j=1}^{q+1}$ using each of the $q + 1$ random seeds for initializing NMF.

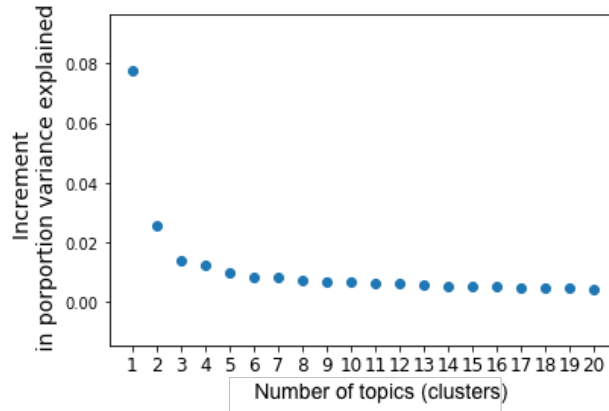


Figure 4.6: Plot of marginal increment in proportion of variance explained by adding another cluster to split X . It is determined that the ideal number of clusters/topics likely lies in the range $[7, 11]$, as this is where the plot starts to level off.

Then, the cosine similarity is calculated between each of the k topics for every consecutive pair of T_j 's. The similarity scores between the topics for each pair (T_j, T_{j+1}) are stored in a matrix $S_{kj} \in \mathbb{R}^{k \times k}$. Therefore, q of such matrices are generated for each $k \in [k_1, k_2]$. For a fixed k , the minimum of all maximum entries from each column and row of each similarity matrix S_{kj} is defined to be least seed similarity (lss) score for that k . The set containing the q , lss scores for a given number of topics k is denoted LSS_k . A consistent number of topics should have an overall high similarity between the topics generated for each seed. Therefore, we choose k^* in $[k_1, k_2]$ to be the “best” number of topics if the median of all its lss scores is the highest.

The boxplot in Figure 4.7 shows the distribution of the lss scores for k in $[7, 11]$. In this case, 8 is chosen as the “best” number of topics since it results in the highest median lss score.

4.5.2 Implementation of hierarchical NMF

A hierarchical NMF (see Algorithm 4) is applied to cluster the articles, where the number of topics in each layer is determined by Algorithm 5. The hierarchical tree structure is established from top to bottom and consists of three layers on this data set (see Figure 4.1).

To generate topics in the first layer, NMF is applied to the matrix X containing all the vectorized

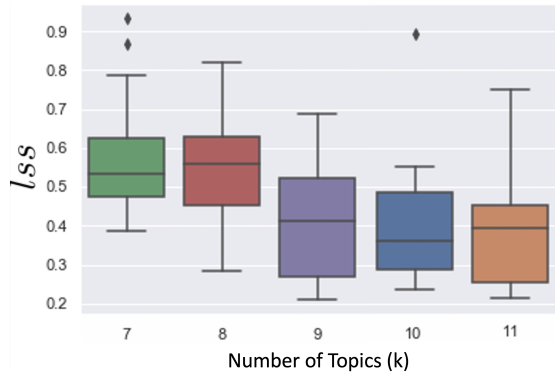


Figure 4.7: Box plot of LSS_k : Topic number 8 is the “best” as it has the highest median lss (least seed similarity) score and should be expected to yield consistent results with random seeds.

articles, resulting in a factorization with 8 topics, as determined by Algorithm 5. Next, a threshold α (in this case, $\alpha = 0.05$) is chosen, and the articles in X are assigned into a topic class X_1, \dots, X_8 if their corresponding document-topic correlation in the H matrix is greater than α . Note that by this definition, one article could be assigned to one or more topic class. After this, any articles not classified to one of the 8 topics are assigned to the “Extra Document” corpus, X_e . Now, the second layer of the tree consists of text corpora X_1, \dots, X_8 .

For each X_i , $i = 1, 2, \dots, 8$ in the second layer, the topic is further subdivided into a third layer if the number of articles assigned to a topic class i is more than some m (in this analysis, we chose $m = 1400$). If it is determined that text corpus X_i needs to be divided further using NMF, the number of subtopics is chosen by Algorithm 5 and again, articles from X_i are assigned to each subtopic based on the threshold α . As before, any articles that do not receive a classification are assigned to X_e . This process is continued for each level in the tree until each leaf contains no more than m articles.

Finally, the cosine similarity between each article in X_e and the dictionary associated to each leaf (topic in the lowest layer in a given branch) is calculated. Note that the dictionary of a leaf is a column of the W matrix of its parent topic. Then the articles in X_e are assigned to the leaf with the highest cosine similarity. After this reassignment, the number of articles associated with each leaf is calculated again, and any leaves containing more than m articles are further subdivided.

We note that in this framework, newly published papers could be added to the tree by first assigning them to X_e and then distributing them as described above. However, since the addition of new papers may also necessitate the introduction of new topics, future work includes extending the tool to an online version that would allow for new topics to be added as new papers appear.

4.6 Conclusion and future work

HNMF is used to organize existing literature on coronaviruses and pandemics, and early literature on COVID-19 into an interactive structure easily searchable by researchers and available to use through a corresponding website. The topics discovered by HNMF reveal that early research of interest to the COVID-19 research community divides into diverse areas such as research related to other coronaviruses, research related to other respiratory diseases, virology and genetic research, as well as research relating to the public health response. A topic coherence metric reveals that the topics discovered are consistent and semantically meaningful, while a topic similarity metric reveals that the topics differ sufficiently from one another to allow for a diversity of choice and areas of interest on the part of the user.

In the future, we hope to regularly update the hierarchical structure as well as the associated website as new research papers are added, both by adding new papers and by adding and deleting classifications as new research topics emerge. We hope to do this using an online version of the HNMF algorithm such as the one in [[TCL18](#)].

4.7 Appendix

4.7.1 Keywords removed

Here is the list of top 30 keywords from the topic which is identified as not the content of the publishing information : without, also, biorxiv, perpetuity copyright, ccbyncnd international, ccbyncnd, peerreviewed copyright, perpetuity peerreviewed, medrxiv preprint, made available, preprint peer-

reviewed, available authorfunder, license made, international license, granted,perpetuity, display preprint, preprint perpetuity, license display, authorfunder granted, granted medrxiv, medrxivlicense, peerreviewed, holder, authorfunder, copyright, copyright holder, holder preprint, covid³, license, medrxiv, preprint.

4.7.2 Topic keywords

Following tables list ten most probable keywords associated with each topic and subtopic in the tree generated by HNMF and the C_V coherence score associated with each. These keywords are visible to website users to enable them to make choices to navigate through the tree. Note that for the first layer we gave suggested topic titles. Not being experts in the field, these are only suggestions to give an idea of the types of research someone may be looking for within that topic.

³We note that we performed a semantic comparison of topics generated using our algorithm both including and removing the word “covid”. No significant differences in topic interpretation were found between the two results indicating that including the word did not add additional information to the topic modeling in this case.

Topic Number	Key Words	Possible Topic
1	merscov, bat, sarscov, camel, mers, merscov infection, ace, human, virus, rbd	Coronaviruses affecting Humans
2	cell, mouse, expression, infection, virus, gene, response, viral, cytokine, immune	Cellular immune response to viral infection
3	sequence, gene, genome, strain, rna, primer, ibv, nucleotide, sample, using	Genetic characteristics of the virus
4	protein, binding, residue, peptide, domain, structure, compound, membrane, cell, activity	Cell-protein interaction
5	virus, influenza, child, rsv, respiratory, infection, viral, sample, assay, detection	Detection and biological study of respiratory viruses
6	patient, hospital, study, pneumonia, clinical, day, infection, treatment, case, symptom	Clinical and hospital studies (esp. of respiratory illnesses)
7	health, model, disease, case, epidemic, outbreak, public, country, population, transmission	Infection models and experiments related to public health
8	vaccine, antibody, pedv, mouse, serum, antigen, pig, strain, protein, response	Vaccine development (esp. of the coronavirus PEDV)

Table 4.2: The top 10 keywords associated with each of the 8 topics in the first layer of the tree

Topic Number	Key Words	C_V
1-1	camel, merscov, dromedary, human, sample, dromedary camel, animal, herd, sequence, study	0.56
1-2	sarscov, ace, protein, ncov, rbd, binding, sars, residue, sequence, virus	0.61
1-3	patient, case, merscov, infection, mers, hospital, outbreak, disease, respiratory, day	0.63
1-4	bat, specie, virus, sequence, bat specie, human, sample, host, covs, study	0.57
1-5	cell, merscov, mouse, protein, antibody, vaccine, virus, response, infection, serum	0.55

Table 4.3: The top 10 keywords associated with each of the subtopics of Topic 1 in the 2nd layer of the tree and the C_V coherence score

Topic Number	Key Words	C_V
2-1	infection, response, lung, immune, cytokine, mouse, tlr, ifn, virus, macrophage	0.59
2-2	virus, cell, viral, infection, infected, replication, vero, culture, antiviral, vero cell	0.61
2-3	cell, antibody, tumor, antigen, culture, human, line, cell line, surface, patient	0.41
2-4	protein, expression, gene, cell, figure, pathway, using, sirna, activity, level	0.55
2-5	mouse, cns, brain, day, demyelination, cell, astrocyte, mhv, day pi, spinal	0.74

Table 4.4: The top 10 keywords associated with each of the subtopics of Topic 2 in the 2nd layer of the tree and the C_V coherence score

Topic Number	Key Words	C_V
3-1	bat, specie, covs, cov, bat specie, virus, sequence, human, sample, coronaviruses	0.59
3-2	gene, rna, protein, cell, expression, mrna, sequence, codon, virus, orf	0.61
3-3	sequence, virus, genome, read, viral, analysis, human, tree, using, specie	0.47
3-4	sample, primer, assay, pcr, probe, detection, dna, virus, reaction, amplification	0.70
3-5	strain, pedv, sequence, pedv strain, aa, vp, nt, diarrhea, gene, china	0.55
3-6	ibv, strain, chicken, vaccine, ibv strain, isolates, virus, bird, gene, flock	0.72

Table 4.5: The top 10 keywords associated with each of the subtopics of Topic 3 in the 2nd layer of the tree and the C_V coherence score

Topic Number	Key Words	C_V
4-1	compound, activity, inhibitor, drug, derivative, mmol, docking, pro, protease, ic	0.62
4-2	nsp, rna, sarscov, nsp nsp, mm, replication, protein, activity, domain, rdrp	0.61
4-3	protein, sequence, interaction, gene, analysis, also, study, function, method, used	0.40
4-4	protein, cell, antibody, mm, sarscov, using, min, expression, serum, recombinant	0.61
4-5	virus, cell, viral, rna, replication, infection, membrane, host, hcv, er	0.64
4-6	peptide, residue, structure, sarscov, binding, fusion, domain, figure, sequence, hr	0.63

Table 4.6: The top 10 keywords associated with each of the subtopics of Topic 4 in the 2^{nd} layer of the tree and the C_V coherence score

Topic Number	Key Words	C_V
5-1	assay, pcr, sample, detection, primer, sensitivity, specimen, method, amplification, probe	0.74
5-2	child, rsv, hmpv, infection, study, hbov, asthma, infant, respiratory, age	0.79
5-3	influenza, ili, virus, surveillance, sari, case, influenza virus, year, study, season	0.66
5-4	patient, respiratory, infection, study, pneumonia, viral, virus, bacterial, pneumoniae, pathogen	0.60
5-5	virus, cell, influenza, infection, influenza virus, protein, viral, antibody, ha, mouse	0.56
5-6	virus, sample, sequence, human, read, hbov, genome, viral, sequencing, study	0.52

Table 4.7: The top 10 keywords associated with each of the subtopics of Topic 5 in the 2^{nd} layer of the tree and the C_V coherence score

Topic Number	Key Words	C_V
6-1	patient, sarscov, display, ct, case, reserved reuse, allowed permission, reuse allowed, permission display, wuhan	0.61
6-2	cap, patient, antibiotic, pneumonia, study, pneumoniae, bacterial, pathogen, infection, culture	0.64
6-3	virus, infection, respiratory, patient, viral, child, rsv, influenza, study, respiratory virus	0.59
6-4	sars, patient, hospital, contact, case, transmission, sars patient, outbreak, staff, care	0.67
6-5	patient, study, treatment, cell, disease, group, level, lung, therapy, day	0.43

Table 4.8: The top 10 keywords associated with each of the subtopics of Topic 6 in the 2^{nd} layer of the tree and the C_V coherence score

Topic Number	Key Words	C_V
7-1	health, public, public health, care, patient, disease, emergency, hospital, system, response	0.67
7-2	disease, animal, human, virus, pathogen, specie, host, infection, vaccine, zoonotic	0.61
7-3	model, individual, epidemic, network, parameter, infected, node, contact, number, rate	0.57
7-4	data, model, study, used, analysis, case, variable, using, method, time	0.44
7-5	case, available display, international made, display, made available, day, wuhan, number, china, international	0.60

Table 4.9: The top 10 keywords associated with each of the subtopics of Topic 7 in the 2^{nd} layer of the tree and the C_V coherence score

Topic Number	Key Words	C_V
8-1	pig, serum, antibody, virus, piglet, group, sample, day, prrsv, tgev	0.58
8-2	mouse, cell, response, group, merscov, immunized, immunization, dna, protein, antibody	0.55
8-3	vaccine, virus, response, influenza, disease, vaccination, immune, human, development, antigen	0.51
8-4	pedv, strain, piglet, pedv strain, ped, cell, gene, diarrhea, sequence, pig	0.56
8-5	protein, antibody, epitope, mabs, peptide, serum, sarscov, mab, elisa, binding	0.61
8-6	ibv, chicken, strain, bird, ibv strain, group, virus, vaccine, ib, egg	0.72

Table 4.10: The top 10 keywords associated with each of the subtopics of Topic 8 in the 2nd layer of the tree and the C_V coherence score

Topic Number	Key Words	C_V
7-1-1	china, research, chinese, disaster, government, social, also, development, policy, people	0.62
7-1-2	patient, hospital, care, infection, staff, medical, health care, nurse, physician, healthcare	0.66
7-1-3	health, public health, public, disease, surveillance, country, system, global, laboratory, outbreak	0.65
7-1-4	pandemic, influenza, participant, respondent, sars, study, outbreak, risk, public, information	0.55

Table 4.11: The top 10 keywords associated with each of the subtopics of Topic 7-1 in the 3rd layer of the tree and the C_V coherence score

Topic Number	Key Words	C_V
7-4-1	study, risk, participant, age, influenza, respondent, factor, country, health, population	0.52
7-4-2	sample, rat, cell, group, animal, cat, used, using, study, protein	0.35
7-4-3	model, data, case, outbreak, surveillance, disease, epidemic, transmission, influenza, time	0.54
7-4-4	air, particle, concentration, wind, velocity, ventilation, flow, airflow, temperature, room	0.73
7-4-5	calf, diarrhea, farm, colostrum, milk, fecal, cow, dairy, herd, day	0.72

Table 4.12: The top 10 keywords associated with each of the subtopics of Topic 7-4 in the 3rd layer of the tree and the C_V coherence score

CHAPTER 5

Adversarial Learning in Distributed Systems

This work is in collaboration with Longxiu Huang and Deanna Needell and is public on [LHN22]. The problem and approach were suggested by the collaborators. I developed the numerical methods and all the simulations as well as fine-tuning the details of the model. Longxiu Huang and I developed the theoretical convergence analysis.

5.1 Background

As machine-learning algorithms gain popularity in industrial applications, it is critical to make them and their optimization subroutines to be robust and adversary-tolerant. Some types of adversaries include evasion [GMP18], data poisoning [GFH20] and model extraction [WXG21, KTP19]. Large-scale machine-learning problems are typically run on distributed systems and an attack in this setting is the Byzantine attack [LSP82] where the individual computing units (also known as ‘workers machines’ or simply ‘workers’) may return adversarial results. A common approach to address this is to utilize *redundancy*; that is, to request the same computation from multiple workers. The main challenge with such an approach is to leverage the outputs from these workers efficiently, and in such a way that even seemingly catastrophic adversarial outputs can be identified and tolerated.

Let’s consider the optimization problem of the following form:

$$\min_{x \in \mathbb{R}^{d_2}} \sum_{i=1}^{d_1} f_i(x) \tag{5.1}$$

where d_1 is a finite integer. To solve the problem in an iterative approach, we have the updating rule:

$$x_{j+1} = x_j + \gamma_j \sum_{i=1}^{d_1} \nabla f_i(x_j) \quad (5.2)$$

with some step-size γ_j . Such objective functions lend themselves naturally to distributed algorithms. In the distributed setting, the central server distributes f_i among the workers. Each worker returns the corresponding gradient $\nabla f_i(x_j)$ and the central server aggregates those returns to compute or approximate the updating step (5.2). In particular, we consider to solve the over-determined linear system

$$Ax = b. \quad (5.3)$$

This problem can be modeled as a least squares problem $\min_x \|Ax - b\|_2^2$ and the least squares problem can be rewritten in the form of (5.1) with $f_i(x_j) = \frac{1}{2}(A_i x_j - b_i)^2$, where $A \in \mathbb{R}^{d_1 \times d_2}$, $b \in \mathbb{R}^{d_1}$, A_i is the i -th row of A , and b_i is the i -th component of b . In this chapter, we mainly focus on developing algorithms to solve (5.3). However, the algorithms can be easily generalized for (5.1). The central server partitions the data matrix A into rows A_i which are distributed among the workers. In the linear setting, each worker only needs to return the scalar $A_i x_j - b_i$ instead of the gradient $(A_i x_j - b_i)A_i^\top$. Then the central server aggregates those returns and approximate the updates in (5.2).

In this work, we consider the setting that some of the workers are adversarial, i.e., the workers return noisy results or enormously large results. Our goal is to develop a variant of the randomized Kaczmarz (RK) method [SV09] for adversarial workers to solve the linear system $Ax = b$. For readers' convenience, we restate the RK method in Algorithm 6. We assume that there is one

Algorithm 6 Randomized Kaczmarz Algorithm

- 1: Select a row index $i_j \in [d_1]$ with probability $p_{i_j} = \frac{\|A_{i_j}\|_2^2}{\|A\|_F^2}$
 - 2: Update $x_{j+1} = \arg \min_{x \in \mathbb{R}^{d_2}} \|x - x_j\|$ s.t. $A_{i_j} x_{j+1} = b_{i_j}$
 - 3: Repeat until convergence
-

central server w_c and N workers in total, among which p fraction of the unknown workers are

adversarial and there are k error categories in total. During the initial data distribution, each row A_r is distributed to N_r workers. Among those N_r workers, workers in ℓ -th category C_ℓ take up a ratio of $p_{r,\ell}$ fraction of all workers, and the total adversarial rate for row r is $p_r = \sum_{\ell=1}^k p_{r,\ell}$. We assume $p_{r,\ell} < 1 - p_r$, for all r, ℓ . Our approach utilizes simple statistics to identify and ignore adversarial results, and thus the setting in which the adversarial workers communicate and select among k types of errors to output is the most challenging for our approach. An adversarial worker w_s^r in category C_ℓ returns the residual $c_s^r = b_r + e_{\ell,r} - \langle x_j, A_r \rangle$, $e_{\ell,r} \in \mathbb{R}$, and a reliable worker returns $c_s^r = b_r - \langle x_j, A_r \rangle$.

5.1.1 Contribution

Our main contributions are threefold: (i) develop efficient methods and algorithms to guarantee accurate estimates for the true solution when adversaries are present, (ii) identify the adversarial workers efficiently, (iii) provide theoretical convergence analysis with part of workers being adversarial for solving large-scale linear systems.

5.1.2 Related work

Kaczmarz method. The Kaczmarz method is an iterative method for solving linear systems that was first proposed by [Kac37]. The method is also known under the name *Algebraic Reconstruction Technique* (ART) in computer tomography [GHJ75, HM93, Nat01] and has found various applications ranging from computer tomography to digital signal processing. Later Strohmer et al. [SV09] proposed a randomized version of Kaczmarz method, where the probability of each row being selected is set to be proportional to the Euclidean norm of the row and prove the exponential bound on the expected rate of convergence. While we consider consistent linear systems, others have analyzed variants of the Kaczmarz methods to handle inconsistent linear systems ([PP16, Pop99, BW21, MNR15]). For example, in [Nee09], the author proved that RK converges for inconsistent linear systems to a horizon that depends upon the size of the largest

entry of the noise. An adaptive maximum-residual sampling strategy has also been analyzed for the inconsistent extension [PP16]. The randomized Kaczmarz method has also been studied in the context of solving systems of linear inequalities [LL10, Agm54, BW19].

Robust optimization. A practical challenge in optimization problem is that there are almost always adversaries present due to mistakes in data collection and data transmission, adversarial or non-responsive workers (also known as ‘stragglers’), or modern storage systems that can introduce corruptions. There have been ongoing researches on mitigating the issues with straggling workers. For example, in [GBH70, KSD17], the authors introduced several encoding schemes that embed the redundancy directly in the data itself to mitigate the effect of straggling. Later, Bitar et al. [BWR19] proposed an approximate gradient-coding scheme for straggler mitigation when the stragglers are uniformly random. An important branch of advances in the analysis of SGD-type methods deal with robustness to adversaries from the data. In [CLZ19] and [HNR22], quantile-based methods were designed to solve corrupted linear equations. To deal with adversarial workers, Yang et al. [YB19] proposed a variant of the gradient descent method based on the geometric median in the setting where the workers split all the data. Alistarh et al. [AAL18] discussed the problem of stochastic optimization in an adversarial setting where the workers sample data from a distribution and an α fraction of them may adversarially return any vector. These methods work only when the adversary rate is less than $\frac{1}{2}$, whereas our algorithm is able to converge to the exact solution even with an adversary rate higher than $\frac{1}{2}$ by utilizing redundancy.

5.2 Method

In this section, we introduce a simple but efficient **mode**-based method to effectively solve linear systems in the presence of the adversarial workers and identify the potential adversarial workers (which are put in a block-list). The method detects the mode category based on the returned category size. More specifically, for each row, the central worker groups the same results and find results from the group with the largest size (i.e., **mode**). Among those modes for all rows, the central server

then updates the guess with the mode with the largest size. If there is only one row, the central server updates the guess with the mode.

Given n workers to deal with one specific row r , the expected number of workers from C_ℓ is $np_{r,\ell}$ ¹ and number of non-adversarial workers is $n(1 - p_r)$. In practice, we uniformly randomly choose a group with the maximum group size, the result will be used to update the guess as long as the group size is greater than $n(1 - p_r)$ (see Alg. 7 Line 7). If a user chooses to implement the algorithm with the block-list, the block-list is updated through a frequency-based approach throughout the iterations: each row has a counter and the counter records if a worker is chosen but fails to be the mode during each iteration. After certain number of iterations, the worker with the largest count in each counter is identified as the potential adversarial worker (see Alg. 7 Line 16–18) and is put in the block-list. Once a worker is in the block-list, it will not be revisited. More details are referred to Alg. 7. In the next section, the theoretical results are based on Alg. 7.

5.3 Theoretical results

In this section, we study the mode distributions and convergence of our method from a theoretical perspective. For the readers' convenience, we first summarize some important notation in Table 5.1.

5.3.1 Mode distribution

Algorithm 7 utilizes the mode to identify adversaries and achieve convergence. In this section, we compute the probability that a category ℓ is the mode for each row during each iteration. For simplicity, let C_0 denote the category of “good” workers. For each row r , the fraction of good workers holding row r is $1 - p_r$. We use d_0 rows for the computation per iteration. Recall that each row r is held by N_r workers (fixed). Among those N_r workers, workers in the category ℓ take up a fraction of $p_{r,\ell}$. At each iteration, the central worker uniformly randomly chooses a set of row

¹The central worker w_c uses first m iterations to determine the number of different groups of results during each iteration and take the maximum number.

Algorithm 7 DISTRIBUTED RANDOMIZED KACZMARZ WITH/WITHOUT BLOCK-LIST

- 1: **Input:** Worker sets D_r , the number of used worker for each row n_r , a counter vector $E_r = 0 \in \mathbb{R}^{N_r}$ for each row r , the number of used rows d_0 , MaxIter, Tol, $c_s = 2 \times \text{Tol}$, checking period T , Blocklist_flag
 - 2: **if** Blocklist_flag **then**
 - 3: initialize **block-list** B
 - 4: **end if**
 - 5: **while** $j < \text{MaxIter}$ and $|c_s| > \text{Tol}$, **do**
 - 6: The central worker w_c uniformly randomly select a subset of rows $\tau \subset [d_1]$
 - 7: Sample $w_1^r, \dots, w_{n_r}^r$ for each row $r \in \tau$, uniformly from D_r
 - 8: Broadcast A_r to $w_1^r, \dots, w_{n_r}^r$
 - 9: w_s^r returns $c_s^r = \frac{\langle A_r, x_i \rangle - b_r + e_l}{\|A_r\|^2}$, if $w_s \in C_l$
 - 10: **for** $r \in \tau$ **do**
 - 11: w_c splits $\{c_s^r\}_{s=1}^{n_r}$ into groups G_1, \dots, G_{k_r}
 - 12: and choose from groups $G_{s^*}^r$ where $G_{s^*}^r = \max_{s'} |G_{s'}^r|$ and $G_{s^*}^r$ satisfies $|G_{s^*}^r| \geq n_r(1 - \sum_{l=1}^{k_r} p_{r,l})$
 - 13: **end for**
 - 14: $c_{s^*} = \max_{r \in \tau} |c_{s^*}^r|$, where $c_{s^*}^r \in G_{s^*}^r$
 - 15: Update $x^{j+1} = x^j + c_{s^*} A_{i_j}^\top$
 - 16: Update $E(s) = E(s) + 1$, if $c_s^r \notin G_{s^*}^r$ ²
 - 17: **if** Blocklist_flag & $\text{mod}(j, T) = 0$ **then**
 - 18: Update B by checking the value of entries in E
 - 19: $D = D \setminus B$
 - 20: **end if**
 - 21: Update $j = j + 1$
 - 22: **end while**
 - 23: **if** Blocklist_flag **then**
 - 24: **Output:** x^j and B
 - 25: **else**
 - 26: **Output:** x^j
 - 27: **end if**
-

indices of size d_0 and requests the corresponding workers to return their results. More specifically, given a set of row indices τ_i at i -th iteration, the central server first finds the modes among the results from each row $r \in \tau_i$ and among those modes, chooses the mode with the largest group size (“the majority vote”).

Table 5.1: Notation table

A	The data matrix A , $A \in \mathbb{R}^{d_1 \times d_2}$
\tilde{A}	The row normalized version of matrix A
N	Number of workers in total
N_r	Number of workers holding row r
n_r	Number of workers chosen for row r
C_ℓ	ℓ -th error category
k	Number of error categories in total
$e_{r,\ell}$	The error of the ℓ -th error category for a row r
e_r	The vector form of errors in all error categories of row r , $e_r := (e_{r,1}, \dots, e_{r,k})$
e	The maxtrix form of errors in all error categories of all rows, $e := \{e_{r,\ell}\}_{r,\ell}$
d_0	Number of rows chosen
$p_{r,\ell}$	The adversarial rate of workers holding row r in error category ℓ
$\hat{q}_{\text{mode}}^{\ell,r}$	Probability that there is a mode among the outputs of chosen workers of row r and the mode is in the category ℓ
q^r	Probability that there is a mode among the outputs of chosen workers for row r
$[d_1]$	The set of the integers from 1 to d_1 , $[d_1] := \{1, \dots, d_1\}$
τ_i	The index set of chosen rows at i -th iteration, $ \tau_i = d_0$
τ_i'	The index set of chosen rows that have a mode, $\tau_i' \subset \tau_i$
$t_i = t(x_{i-1}, \tau_i)$	The index of the row that has the largest mode number

For any row r , let $a_{g,\ell}^r$ be the coefficient of the term $x^{n_r - g}$ of the polynomials

$$\prod_{\ell'=0, \ell' \neq \ell}^k \sum_{j=0}^{g-1} \binom{N_r p_{r,\ell'}}{j} x^j.$$

Let b_i^r be the coefficient of the term x^{n_r} of the polynomial

$$\prod_{\ell=0}^k \left(\sum_{j=0}^{g-1} \binom{N_r p_{r,\ell}}{j} x^j \right)$$

Lemma 5.3.1. *For a row r , the probability that the mode is in the category ℓ with mode number g is $\mathbb{P}(r \text{ mode}, g, \ell) = \frac{\binom{N_r p_{r,\ell}}{g} a_{g,\ell}^r}{\binom{N_r}{n_r}}$.*

Using Lemma 5.3.1, we obtain the following conclusions by going over all possible mode

numbers or all error categories:

Lemma 5.3.2. *For row r , the probability that the category ℓ is the mode is*

$$\hat{q}_{mode}^{\ell,r} = \sum_{g=g_0(r)}^{n_r} \frac{\binom{N_r p_{r,\ell}}{g} a_{g,\ell}^r}{\binom{N_r}{n_r}}.$$

where $g_0(r) = \max(\lceil \frac{n_r}{k+1} \rceil, \lceil n_r(1 - p_{r,0}) \rceil)$. Thus, the probability that there is a mode with mode number g for the calculation of row r is

$$q_g^r = \sum_{\ell=0}^k \frac{\binom{N_r p_{r,\ell}}{g} a_{g,\ell}^r}{\binom{N_r}{n_r}}.$$

Additionally, the probability that there is a mode for the calculation of row r is

$$q^r = \sum_{\ell=0}^k \hat{q}_{mode}^{\ell,r} = \sum_{g=g_0(r)}^{n_r} \sum_{\ell=0}^k \frac{\binom{N_r p_{r,\ell}}{g} a_{g,\ell}^r}{\binom{N_r}{n_r}}, \quad (5.4)$$

where $\binom{n_r}{g} = 0$ when $n_r < g$.

In the following lemma, we also calculate the probability $\mathbb{P}(t, \ell, g | \tau_i, x_{i-1})$ that a mode produced from row t in the category C_ℓ with a mode number g when rows τ_i are used in the computation and the previous estimate x_{i-1} is given. For simplicity, we omit the condition of τ_i, x_{i-1} in the notation and denote $\mathbb{P}(t, \ell, g | \tau_i, x_{i-1})$ by $\mathbb{P}(t, \ell, g)$.

Lemma 5.3.3. *Given the previous estimate x_{i-1} and row indices τ_i , we have*

$$\mathbb{P}(t_i, \ell, g) = \frac{\binom{N_{t_i} p_{t_i,\ell}}{g} a_{g,\ell}^{t_i}}{\binom{N_{t_i}}{n_{t_i}}} \prod_{s \in \tau_i \setminus t_i} \frac{b_g^s}{\binom{N_s}{n_s}}.$$

Proof. The probability that the mode of row t_i is produced by category ℓ with group size g

$$\mathbb{P}(t_i, l, g) = \mathbb{P}(t \text{ mode}, g, \ell) \times \mathbb{P}(t \text{ is the mode with the largest mode number} | t \text{ mode}, l, g) \quad (5.5)$$

$$= \frac{\binom{N_{t_i} p_{t_i, \ell}}{g} a_{g, \ell}^{t_i}}{\binom{N_{t_i}}{n_{t_i}}} \times \mathbb{P}(t_i \text{ is the mode with the largest mode number} | t_i \text{ mode}, l, g) \quad (5.6)$$

$$= \frac{\binom{N_{t_i} p_{t_i, \ell}}{g} a_{g, \ell}^{t_i}}{\binom{N_{t_i}}{n_{t_i}}} \prod_{s \in \tau_i \setminus t_i} \frac{b_g^s}{\binom{N_s}{n_s}}. \quad (5.7)$$

□

Corollary 5.3.4. *Iterating over all error categories, we get the probability that row t_i produces the mode with mode number g ($g \leq n_{t_i}$):*

$$\mathbb{P}(t_i, g) = \sum_{\ell=0}^k \mathbb{P}(t_i, \ell, g) = \sum_{\ell=0}^k \frac{\binom{N_{t_i} p_{t_i, \ell}}{g} a_{g, \ell}^{t_i}}{\binom{N_{t_i}}{n_{t_i}}} \prod_{s \in \tau_i \setminus \{t_i\}} \frac{b_g^s}{\binom{N_s}{n_s}} = q_g^t \prod_{s \in \tau_i \setminus t_i} \frac{b_g^s}{\binom{N_s}{n_s}}. \quad (5.8)$$

5.3.2 Convergence without block-list

Let t_i be the row chosen at i -th iteration to update the guess x . For the convergence analysis, we consider solving

$$\begin{aligned} A_{t_i} x &= b_{t_i}, \\ A_{t_i} x &= b_{t_i} + e_{t_i, 1}, \\ &\vdots \\ A_{t_i} x &= b_{t_i} + e_{t_i, k}, \end{aligned}$$

with probability q_0, q_1, \dots, q_k respectively. For simplicity, we assume $\tau_i \sim \text{unif}\left(\binom{[d_1]}{d_0}\right)$. We would either have the iteration

$$x_i = x_{i-1} - \frac{\langle A_{t_i}, x_{i-1} \rangle - b_{t_i}}{\|A_{t_i}\|^2} A_{t_i}^\top,$$

or

$$x_i = x_{i-1} - \frac{\langle A_{t_i}, x_{i-1} \rangle - (b_{t_i} + e_{t_i, \ell})}{\|A_{t_i}\|^2} A_{t_i}^\top,$$

for $\ell = 1, \dots, k$, and A_{t_i} is the t_i -th row of matrix A .

In the following analysis, let \mathbb{E}_{τ_i} denote expectation with respect to the uniformly random

sample τ_i conditioned upon the sampled τ_j for $j < i$, and \mathbb{E} denote expectation with respect to all random samples τ_j for $1 \leq j \leq i$ where i is understood to be the last iteration in the context in which \mathbb{E} is applied.

Lemma 5.3.5. *The conditional expectation of the squared convergence error at i -th iteration can be decomposed into three parts: the squared convergence error at $(i - 1)$ -th iteration $\|x_{i-1} - x^*\|_2^2$, the conditional expectation of the squared errors (normalized) from the adversarial workers $\mathbb{E}_{\tau_i} \mathbb{E}_{t_i} \mathbb{E}_{\ell_{t_i}} \frac{e_{t_i, \ell_{t_i}}^2}{\|A_{t_i}\|_2^2}$, and the conditional expectation of the squared residual (normalized) $\mathbb{E}_{\tau_i} \mathbb{E}_{t_i} \frac{\langle A_{t_i}^T, x_{i-1} - x^* \rangle^2}{\|A_{t_i}\|_2^2}$:*

$$\mathbb{E}\|x_i - x^*\|_2^2 = \|x_{i-1} - x^*\|_2^2 + \mathbb{E}_{\tau_i} \mathbb{E}_{t_i} \mathbb{E}_{\ell_{t_i}} \frac{e_{t_i, \ell_{t_i}}^2}{\|A_{t_i}\|_2^2} - \mathbb{E}_{\tau_i} \mathbb{E}_{t_i} \frac{\langle A_{t_i}^T, x_{i-1} - x^* \rangle^2}{\|A_{t_i}\|_2^2}. \quad (5.9)$$

Proof. We start by decomposing the error:

$$\|x_i - x^*\|_2^2 = \|x_{i-1} - \frac{\langle A_{t_i}^T, x_{i-1} \rangle - (b_{t_i} + e_{t_i, \ell_{t_i}})}{\|A_{t_i}\|_2^2} A_{t_i}^\top - x^*\|_2^2 \quad (5.10)$$

$$= \|x_{i-1} - x^*\|_2^2 + \frac{(\langle A_{t_i}^T, x_{i-1} - x^* \rangle - e_{t_i, \ell_{t_i}})^2}{\|A_{t_i}\|_2^2} \quad (5.11)$$

$$- \frac{2}{\|A_{t_i}\|_2^2} \langle x_{i-1} - x^*, A_{t_i}^T \rangle (\langle A_{t_i}^T, x_{i-1} - x^* \rangle - e_{t_i, \ell_{t_i}}) \quad (5.12)$$

$$= \|x_{i-1} - x^*\|_2^2 - \frac{\langle A_{t_i}^T, x_{i-1} - x^* \rangle^2}{\|A_{t_i}\|_2^2} + \frac{e_{t_i, \ell_{t_i}}^2}{\|A_{t_i}\|_2^2} \quad (5.13)$$

Take the expectation, we get

$$\mathbb{E}\|x_i - x^*\|_2^2 = \|x_{i-1} - x^*\|_2^2 + \mathbb{E}_{\tau_i} \mathbb{E}_{t_i} \mathbb{E}_{\ell_{t_i}} \frac{e_{t_i, \ell_{t_i}}^2}{\|A_{t_i}\|_2^2} - \mathbb{E}_{\tau_i} \mathbb{E}_{t_i} \frac{\langle A_{t_i}^T, x_{i-1} - x^* \rangle^2}{\|A_{t_i}\|_2^2} \quad (5.14)$$

□

Next we compute the conditional expectation of the error part from the adversarial workers $\mathbb{E}_{\tau_i} \mathbb{E}_{t_i} \mathbb{E}_{\ell_{t_i}} \frac{e_{t_i, \ell_{t_i}}^2}{\|A_{t_i}\|_2^2}$ and the residual part $\mathbb{E}_{\tau_i} \mathbb{E}_{t_i} \frac{\langle A_{t_i}^T, x_{i-1} - x^* \rangle^2}{\|A_{t_i}\|_2^2}$ separately in the following lemmas.

Lemma 5.3.6. *The conditional expectation of squared residual can be bounded by:*

$$\mathbb{E}_{\tau_i} \mathbb{E}_{t_i} \frac{\langle A_{t_i}^T, x_{i-1} - x^* \rangle^2}{\|A_{t_i}\|_2^2} \geq Q_{\min} \frac{d_0}{d_1} \sigma_{\min}^2(\tilde{A}) \|x_{i-1} - x^*\|^2. \quad (5.15)$$

Moreover, we have

$$\|x_{i-1} - x^*\|^2 - \mathbb{E}_{\tau_i} \mathbb{E}_{t_i} \frac{\langle A_{t_i}^T, x_{i-1} - x^* \rangle^2}{\|A_{t_i}\|_2^2} \leq (1 - Q_{\min} \frac{d_0}{d_1} \sigma_{\min}^2(\tilde{A})) \|x_{i-1} - x^*\|^2, \quad (5.16)$$

where

$$Q_{\min} = \min_{g, t_i, \tau_i} \sum_{g=g_0(t_i)}^{n_t} q_g^{t_i} \prod_{s \in \tau_i \setminus \{t_i\}} \frac{b_g^s}{\binom{N_s}{n_s}}.$$

Proof. The expectation of the squared residual can be calculated:

$$\mathbb{E}_{\tau_i} \mathbb{E}_{t_i} \frac{\langle A_{t_i}^T, x_{i-1} - x^* \rangle^2}{\|A_{t_i}\|_2^2} = \mathbb{E}_{\tau_i} \sum_{t_i \in \tau_i} \sum_{g=g_0(t_i)}^{n_{t_i}} \sum_{\ell=0}^k \mathbb{P}(t_i, \ell, g) \left\| \frac{A_{t_i}^T (x_{i-1} - x^*)}{\|A_{t_i}\|} \right\|^2 \quad (5.17)$$

$$= \sum_{\tau_i \in \binom{[d_1]}{d_0}} p_{x_{i-1}(\tau_i)} \sum_{t_i \in \tau_i} \sum_{g=g_0(t_i)}^{n_{t_i}} \mathbb{P}(t_i, g) \left\| \frac{A_{t_i}^T (x_{i-1} - x^*)}{\|A_{t_i}\|} \right\|^2 \quad (5.18)$$

$$= \sum_{\tau_i \in \binom{[d_1]}{d_0}} p_{x_{i-1}(\tau_i)} \sum_{t_i \in \tau_i} \sum_{g=g_0(t_i)}^{n_{t_i}} q_g^{t_i} \prod_{s \in \tau_i \setminus \{t_i\}} \frac{b_g^s}{\binom{N_s}{n_s}} \left\| \frac{A_{t_i}^T (x_{i-1} - x^*)}{\|A_{t_i}\|} \right\|^2. \quad (5.19)$$

Let

$$Q_{\min} = \min_{t, \tau_i} \sum_{g=g_0(t)}^{n_t} q_g^t \prod_{s \in \tau_i \setminus \{t\}} \frac{b_g^s}{\binom{N_s}{n_s}}.$$

Then

$$\mathbb{E}_{\tau_i} \mathbb{E}_{t_i} \frac{\langle A_{t_i}^T, x_{i-1} - x^* \rangle^2}{\|A_{t_i}\|_2^2} \geq Q_{\min} \sum_{\tau_i \in \binom{[d_1]}{d_0}} p_{x_{i-1}(\tau_i)} \sum_{t \in \tau_i} \left\| \frac{A_{t_i}^T (x_{i-1} - x^*)}{\|A_{t_i}\|} \right\|^2 \quad (5.20)$$

$$\geq p_{x_{i-1}}(\tau_i) Q_{\min} \sum_{\tau_i \in \binom{[d_1]}{d_0}} \sum_{t_i \in \tau_i} \left\| \frac{A_{t_i}^T (x_{i-1} - x^*)}{\|A_{t_i}\|} \right\|^2 \quad (5.21)$$

$$\geq p_{x_{i-1}}(\tau_i) Q_{\min} \binom{d_1 - 1}{d_0 - 1} \|\tilde{A}^T (x_{i-1} - x^*)\|_2^2 \quad (5.22)$$

$$\geq p_{x_{i-1}}(\tau_i) Q_{\min} \binom{d_1 - 1}{d_0 - 1} \sigma_{\min}^2(\tilde{A}) \|x_{i-1} - x^*\|_2^2, \quad (5.23)$$

where $\tilde{A} \in \mathbb{R}^{d_1 \times d_2}$ is obtained by normalizing each row of A . In fact, $p_{x_{i-1}}(\tau_i)$ is independent of x_{i-1} since

$$\tau_i \sim \text{unif}\left(\binom{[d_1]}{d_0}\right).$$

We have

$$p_{x_{i-1}}(\tau_i) = 1 / \binom{d_1}{d_0} = \frac{d_0!(d_1 - d_0)!}{d_1!}.$$

Note that $p_{x_{i-1}}(\tau_i) \binom{d_1 - 1}{d_0 - 1} = \frac{d_0}{d_1}$, $Q_{\min} \in (0, 1)$. Therefore, we have

$$\begin{aligned} \mathbb{E}_{\tau_i} \mathbb{E}_{t_i} \frac{\langle A_{t_i}^T, x_{i-1} - x^* \rangle^2}{\|A_{t_i}\|_2^2} &\geq Q_{\min} \frac{d_0}{d_1} \sigma_{\min}^2(\tilde{A}) \|x_{i-1} - x^*\|^2, \\ \|x_{i-1} - x^*\|^2 - \mathbb{E}_{\tau_i} \mathbb{E}_{t_i} \frac{\langle A_{t_i}^T, x_{i-1} - x^* \rangle^2}{\|A_{t_i}\|_2^2} &\leq (1 - Q_{\min} \frac{d_0}{d_1} \sigma_{\min}^2(\tilde{A})) \|x_{i-1} - x^*\|^2 \end{aligned} \quad (5.24)$$

□

Lemma 5.3.7. *The expectation of the squared error from the adversarial workers can be bounded by:*

$$\begin{aligned} \mathbb{E}_{\tau_i} \mathbb{E}_t \mathbb{E}_\ell \frac{e_{t,\ell}^2}{\|A_t\|_2^2} &\leq \sum_{t \in [d_1]} q_t \|\tilde{e}_t\|_2^2, \\ \text{where } \tilde{e}_{\max}^2 &= \max_{t_i, \ell} \frac{e_{t_i, \ell}^2}{\|A_{t_i}\|_2^2}, \\ Q_{\max}(t, g, \tau_i \setminus \{t\}) &= \max_\ell \frac{\binom{N_t p_{t,\ell}}{g} a_{g,\ell}^t}{\binom{N_t}{n_t}} \prod_{s \in \tau_i \setminus \{t\}} \frac{b_g^s}{\binom{N_s}{n_s}}, \\ q_t &= \sum_{\tilde{\tau}_{t,i} \in \binom{[d_1 - 1]}{d_0 - 1}} \sum_{g=g_0(t)}^{n_{t_i}} \frac{1}{\binom{d_1}{d_0}} Q_{\max}(t, g, \tilde{\tau}_{t,i}). \end{aligned} \quad (5.25)$$

Proof. The expectation of the squared error from the adversarial workers:

$$\mathbb{E}_{\tau_i} \mathbb{E}_t \mathbb{E}_\ell \frac{e_{t,\ell}^2}{\|A_t\|_2^2} = \mathbb{E}_{\tau_i} \sum_{t \in \tau_i} \sum_{g=g_0(t)}^{n_t} \sum_{\ell=0}^k \frac{\binom{N_t p_{t,\ell}}{g} a_{g,\ell}^t}{\binom{N_t}{n_t}} \prod_{s \in \tau_i \setminus \{t\}} \frac{b_g^s}{\binom{N_s}{n_s}} \frac{e_{t,\ell}^2}{\|A_t\|_2^2} \quad (5.26)$$

Let $\tilde{e}_{t,\ell}^2 = \frac{e_{t,\ell}^2}{\|A_t\|_2^2}$, $\tilde{e}_t = (e_{t,1}, \dots, e_{t,k})$, and $Q_{\max}(t, g, \tau_i \setminus \{t\}) = \max_\ell \frac{\binom{N_t p_{t,\ell}}{g} a_{g,\ell}^t}{\binom{N_t}{n_t}} \prod_{s \in \tau_i \setminus \{t\}} \frac{b_g^s}{\binom{N_s}{n_s}}$.

Then

$$\begin{aligned} (5.25) &= \mathbb{E}_{\tau_i} \sum_{t \in \tau_i} \sum_{g=g_0(t)}^{n_t} \sum_{\ell=0}^k \frac{\binom{N_t p_{t,\ell}}{g} a_{g,\ell}^t}{\binom{N_t}{n_t}} \prod_{s \in \tau_i \setminus \{t\}} \frac{b_g^s}{\binom{N_s}{n_s}} \tilde{e}_{t,\ell}^2 \\ &\leq \mathbb{E}_{\tau_i} \sum_{t \in \tau_i} \sum_{g=g_0(t)}^{n_t} Q_{\max}(t, g, \tau_i \setminus \{t\}) \sum_{\ell=0}^k \tilde{e}_{t,\ell}^2 \\ &= \mathbb{E}_{\tau_i} \sum_{t \in \tau_i} \sum_{g=g_0(t)}^{n_t} Q_{\max}(t, g, \tau_i \setminus \{t\}) \|\tilde{e}_t\|_2^2 \\ &= \sum_{\tau_i \in \binom{[d_1]}{d_0}} \frac{1}{\binom{d_1}{d_0}} \sum_{t \in \tau_i} \sum_{g=g_0(t)}^{n_t} Q_{\max}(t, g, \tau_i \setminus \{t\}) \|\tilde{e}_t\|_2^2 \\ &= \sum_{t \in [d_1]} \sum_{\tilde{\tau}_{t,i} \in \binom{[d_1] \setminus \{t\}}{d_0-1}} \sum_{g=g_0(t)}^{n_t} \frac{1}{\binom{d_1}{d_0}} Q_{\max}(t, g, \tilde{\tau}_{t,i}) \|\tilde{e}_t\|_2^2 = \sum_{t \in [d_1]} q_t \|\tilde{e}_t\|_2^2 \end{aligned} \quad (5.27)$$

with $q_t = \sum_{\tilde{\tau}_{t,i} \in \binom{[d_1] \setminus \{t\}}{d_0-1}} \sum_{g=g_0(t)}^{n_t} \frac{1}{\binom{d_1}{d_0}} Q_{\max}(t, g, \tilde{\tau}_{t,i})$. Notice that

$$\sum_{g=g_0(t)}^{n_t} \sum_{\ell=0}^k \frac{\binom{N_t p_{t,\ell}}{g} a_{g,\ell}^t}{\binom{N_t}{n_t}} \prod_{s \in \tau_i \setminus \{t\}} \frac{b_g^s}{\binom{N_s}{n_s}} \leq 1,$$

we thus have $\sum_{g=i_0(t)}^{n_t} Q_{\max}(t, g, \tau_i \setminus \{t\}) \leq 1$. Thus,

$$(5.25) \leq \sum_{\tau_i \in \binom{[d_1]}{d_0}} \frac{1}{\binom{d_1}{d_0}} \sum_{t \in \tau_i} \sum_{g=g_0(t)}^{n_t} Q_{\max}(t, g, \tau_i \setminus \{t\}) \|\tilde{e}_t\|_2^2 \leq \frac{d_0}{d_1} \sum_{t=1}^{d_1} \|\tilde{e}_t\|_2^2. \quad (5.28)$$

□

Combine lemma 5.3.6 and lemma 5.3.7, we have the following theorem.

Theorem 5.3.8. Let $A \in \mathbb{R}^{d_1 \times d_2}$ with $d_1 \geq d_2$ and $b, e_1, \dots, e_k \in \mathbb{R}^{d_1}$. Assume that we solve $Ax^* = b$ via Algorithm 8; then

$$\mathbb{E}\|x_i - x^*\|_2^2 \leq \alpha^i \|x_0 - x^*\|_2^2 + \frac{1 - \alpha^{i+1}}{1 - \alpha} \sum_{t \in [d_1]} q_t \|\tilde{e}_t\|_2^2, \quad (5.29)$$

where

$$\begin{aligned} \alpha &= 1 - Q_{\min} \frac{d_0}{d_1} \sigma_{\min}^2(\tilde{A}), \\ Q_{\min} &= \min_{g, t_i, \tau_i} \sum_{g=g_0(t_i)}^{n_t} q_g^{t_i} \prod_{s \in \tau_i \setminus \{t_i\}} \frac{b_g^s}{\binom{N_s}{n_s}}, \\ Q_{\max} &= \max_{g, t_i, \tau_i} \sum_{g=g_0(t_i)}^{n_{t_i}} q_g^{t_i} \prod_{s \in \tau_i \setminus \{t_i\}} \frac{b_g^s}{\binom{N_s}{n_s}}, \\ \|\tilde{e}_t\|_2^2 &= \sum_{\ell=0}^k \tilde{e}_{t,\ell}^2, \quad \tilde{e}_{t,\ell}^2 = \frac{e_{t,\ell}^2}{\|A_t\|_2^2}, \\ q_t &= \sum_{\tilde{\tau}_{t,i} \in \binom{[d_1-1]}{d_0-1}} \sum_{g=g_0(t)}^{n_{t_i}} \frac{1}{\binom{d_1}{d_0}} Q_{\max}(t, g, \tilde{\tau}_{t,i}), \end{aligned} \quad (5.30)$$

and \tilde{A} is the row normalized version of matrix A and $\sigma_{\min}^2(\tilde{A})$ is the smallest eigenvalue of \tilde{A} .

Lemma 5.3.9. Let $\tilde{A} \in \mathbb{R}^{d_1 \times d_2}$ with $d_1 \geq d_2$ and each row is normalized. Then $\sigma_{\min}(\tilde{A}) \leq \sqrt{d_1/d_2}$.

Proof. Notice that $\sum_{i=1}^{d_2} \sigma_i^2 = d_1$, we have $\sigma_{\min}^2 \leq d_1/d_2$. \square

Remark 5.3.10. When $d_0 \leq d_2$ we have $0 < 1 - \frac{d_0}{d_1} Q_{\min} \sigma_{\min}^2(\tilde{A}) < 1$

From (5.30), q_t is not a simple linear function of d_0 and increasing d_0 may not necessarily decrease q_t . An example in Table 5.2 shows that increasing d_0 , to some extent, can decrease q_t and therefore, improves the speed of convergence. For more details about finding the optimal d_0 , one can refer to remark 5.6.3 and remark 5.6.4 in the appendix 5.6.2. Meanwhile, one should be aware of the increasing d_0 leads to more communication cost. Thus, in practice, finding an optimal d_0 is not just minimizing q_t but also reducing the communication cost.

When adversaries $\{\tilde{e}_t\}_t$ are relatively small, the method without the block-list can guarantee a convergence error with the same or smaller magnitude as $\{\tilde{e}_t\}_t$ using the right parameters. However, from (5.30), when adversaries $\{\tilde{e}_t\}_t$ are relatively larger, the convergence is not guaranteed. Therefore, it is crucial to introduce the block-list method to exclude the adversarial workers and with block-list the convergence is also guaranteed in any hostile environment. In the following sections, we first provide theoretical error bound for the method without block-list. Next we reason why the block-list is necessary in certain cases and show the effectiveness of the block-list method.

Table 5.2: Total number of workers $N = 10$, number of error categories $k = 3$.

p	n	d_0	Q	q_t
0.6	5	2	4.25×10^{-3}	8.5×10^{-4}
	5	3	3.3×10^{-4}	9.9×10^{-5}
	5	5	3.63×10^{-6}	1.82×10^{-6}

5.3.3 Block-list method

According to Alg. 7 after S iterations, the worker w^* with the largest non-mode (being chosen but not the mode) counts $c_{w^*}^+$ will be set in the block-list. To evaluate the effectiveness of the method, it's necessary to calculate the probability that a bad/good worker is put in the block-list. This problem can be reformulated mathematically as the following:

Problem 5.3.11. Let $c_w^+(S) := c_w^+$, $c_w^0(S) := c_w^0$ be the counters of the worker w is non-mode, and mode or in no mode case respectively, among S iterations. Then we have $0 \leq c_w^+, c_w^0 \leq S$ and $\sum_{w=1}^N (c_w^+ + c_w^0) = nS$. The probability w^* is in the block-list after S iterations can be calculated as follows:

$$\mathbb{P}_{bl}(w^*) = \mathbb{P}(c_{w^*}^+ > c_w^+, \forall w \neq w^*) \text{ s.t. } \sum_{w=1}^N c_w^+ + c_w^0 = nS \quad (5.31)$$

Note that this probability can be calculated by using integer dynamic programming or estimated by Monte Carlo simulations.

Lemma 5.3.12. *Run Alg. 7 with S iteration. Then the conditional probability that a good worker w_0 is in the block-list is*

$$\frac{p_0 \mathbb{P}_{bl}(w_0)}{\sum_{\ell=0}^k p_\ell \mathbb{P}_{bl}(w_\ell)}.$$

Similarly, the conditional probability that a bad worker $w_{\ell'}$ in category ℓ' is in the block-list is

$$\frac{p_{\ell'} \mathbb{P}_{bl}(w_{\ell'})}{\sum_{\ell=0}^k p_\ell \mathbb{P}_{bl}(w_\ell)}.$$

To illustrate how the quantities changes with respect to S in Lemma 5.3.12, we consider the following example.

Remark 5.3.13. *Assume that there are two categories of workers i.e. $k = 0, 1$, and 5 workers $w_i, i = 1, \dots, 5$ in total with $w_1, w_2 \in C_1, w_3, w_4, w_5 \in C_0$. Let $n = 3$. Note that $\mathbb{P}_{bl}(w_1) = \mathbb{P}_{bl}(w_2) := \mathbb{P}_{bl}^1$ and $\mathbb{P}_{bl}(w_3) = \mathbb{P}_{bl}(w_4) = \mathbb{P}_{bl}(w_5) := \mathbb{P}_{bl}^0$. The probability is estimated by Monte Carlo simulations. We simulated the experiment 100 times and count the numbers of experiments where each worker is listed in the block-list. Those numbers are used to calculate the frequency and estimate the probability. The estimated results are summarized in Table 5.3. Table 5.3 shows*

Table 5.3: Conditional probability of being in block-list.

S	5	10	50	100
\mathbb{P}_{bl}^1	0.403	0.452	0.5	0.5
\mathbb{P}_{bl}^0	0.065	0.032	~ 0	~ 0

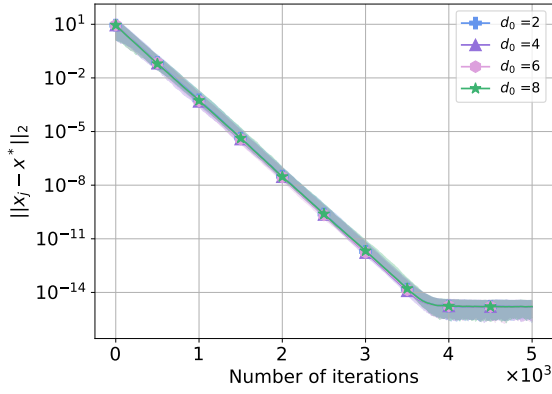
that the probability of an adversarial worker in the block-list increases as the number of iterations S increases. Meanwhile, the probability of a good worker in the block-list decreases. Using the method with the block-list, we are able to avoid choosing the results from the adversarial workers. As a results, the probability of using the adversarial workers decreases, i.e., q_t decreases.

5.4 Simulations

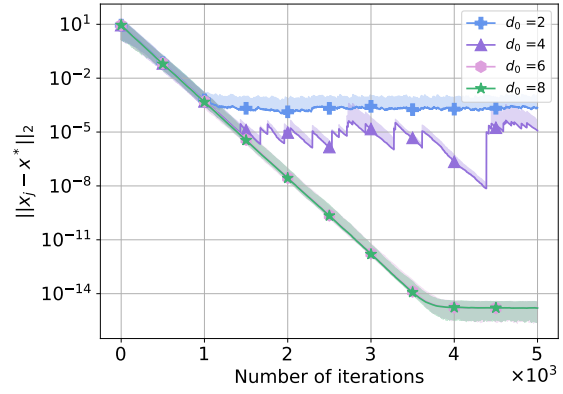
In this section, we test the performance of our approaches for solving consistent linear systems. In the simulations, we randomly generate a row-normalized matrix $A \in \mathbb{R}^{1200 \times 50}$, $x \in \mathbb{R}^{50}$ from normalized Gaussian distribution and set $b = Ax$. For each row, there are the same number of error categories k and the adversarial rate p_r is same, i.e., $p_r = p/k$, where p is the total adversarial rate. The linear system $Ax = b$ is solved via Alg. 7 with and without block-list. At each iteration, d_0 rows of A are uniformly randomly chosen. For each row r in the selected d_0 rows, n_r workers are randomly selected from N_r workers to participate the calculation. The simulation shows how the number of used rows d_0 , the number of used workers n_r , the total adversary rate p and the number of the error categories k affect the performance.

Fig. 5.1 and 5.2 present the effects of the number of the used rows d_0 and the convergence results for our distributed Randomized Kaczmarz method with and without the block-list. The maximum of adversary e is 10^{-3} and 500, respectively. In this example, increasing the number of used rows d_0 from 2 to 4 improves the convergence in both with and without block-list cases, regardless of the adversaries' magnitude. Fig. 5.1 shows the convergence results when $\|e\|_\infty = 10^{-3}$. Using the block-list, the convergence are fast over all choices of d_0 when the adversarial rate $p = 0.2$ (Fig. 5.1a); the larger the number of used rows d_0 , the faster the convergence when the adversarial rate $p = 0.6$ (Fig. 5.1c). In Fig. 5.1b, when $d_0 = 2, 4$, the central server possibly uses a corrupted step-size to update and oscillate around the solution and thus, the error converges to a range of magnitude from 10^{-3} to 10^{-5} ; when $d_0 = 6, 8$, the convergence error goes to 0 after 4000 iterations. In Fig. 5.1d, the convergence errors of all choices of d_0 are in the range of $(10^{-3}, 10^{-4})$. When the magnitude of $\|e\|$ is small, the method without the block-list can converge when increasing the number of used rows d_0 . However, when the magnitude of the adversaries and the adversarial rate are large (5.2d), the convergence is no longer guaranteed without a block-list.

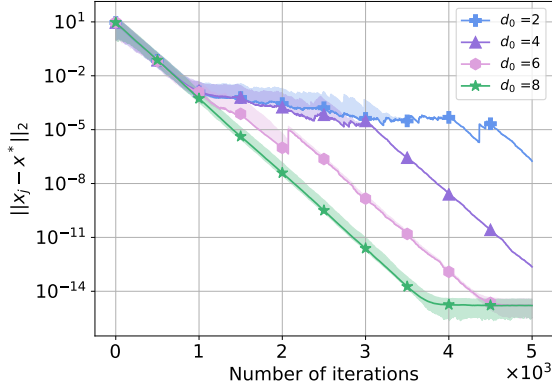
Fig. 5.2 shows the convergence results when $\|e\|_\infty = 5 \times 10^2$. In Fig. 5.2a, with the block-list, the method reaches an accuracy of 10^{-14} regardless of the value of d_0 . Similar oscillations exist in



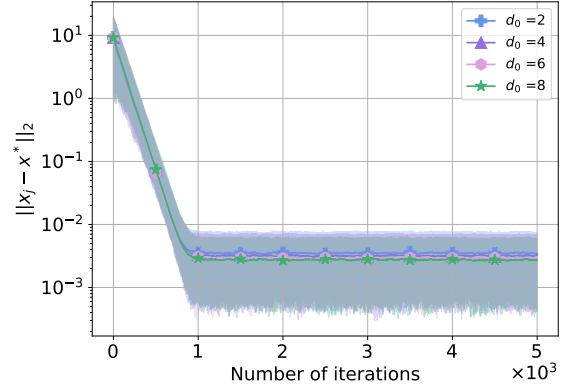
(a) Error vs. d_0 : $p = 0.2$; using Alg. 7 with block-list.



(b) Error vs. d_0 : $p = 0.2$; using Alg. 7 without block-list.



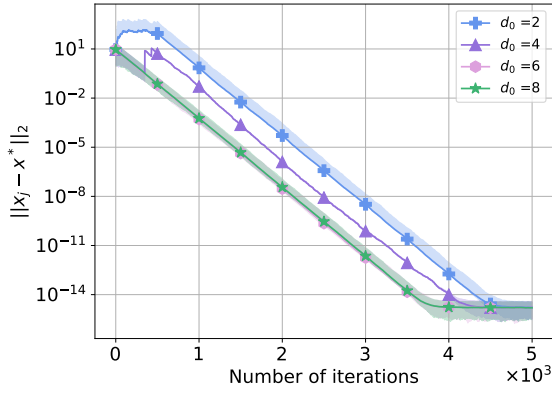
(c) Error vs. d_0 : $p = 0.6$; using Alg. 7 with block-list.



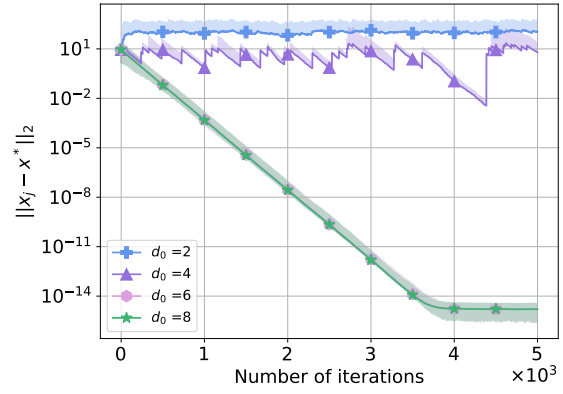
(d) Error vs. d_0 : $p = 0.6$; using Alg. 7 without block-list.

Figure 5.1: Effects of the number of used rows d_0 on convergence: $N_r = 20, n_r = 4, k = 3, \|e\|_\infty = 10^{-3}$. The error norms were averaged over 50 trials (the solid lines) with 90% percentiles (the shaded areas).

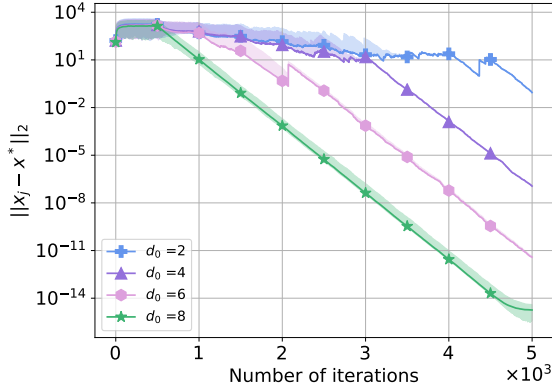
Fig. 5.2b and Fig. 5.2c without the block-list. In particular, without the block-list, the convergence is not guaranteed as we increase to d_0 when $p = 0.6$. This suggests that in an environment with larger outliers, it is efficient to use the block-list method. Fig. 5.3 presents the effect of the number of the chosen workers n_r when the adversary rate p is 0.2 and 0.6. As the number of chosen workers n increases from 3 to 7, the convergence is faster for both without and with the block-list. The method with the block-list, in general, guarantees better convergence compared to the one without. The trade-off is the extra storage for the block-list. Without the block-list, when $p = 0.2$, the convergence oscillates for $n_r = 3$ and when the adversarial rate is increased to $p = 0.6$, the



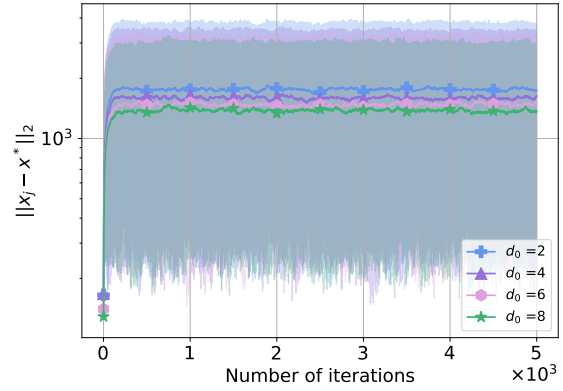
(a) Error vs. d_0 : $p = 0.2$; using Alg. 7 with block-list.



(b) Error vs. d_0 : $p = 0.2$; using Alg. 7 without block-list.



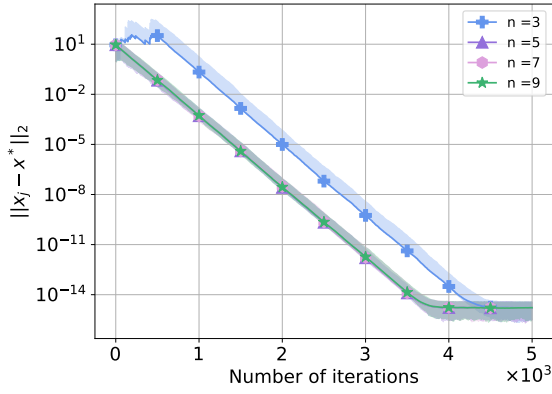
(c) Error vs. d_0 : $p = 0.6$; using Alg. 7 with block-list.



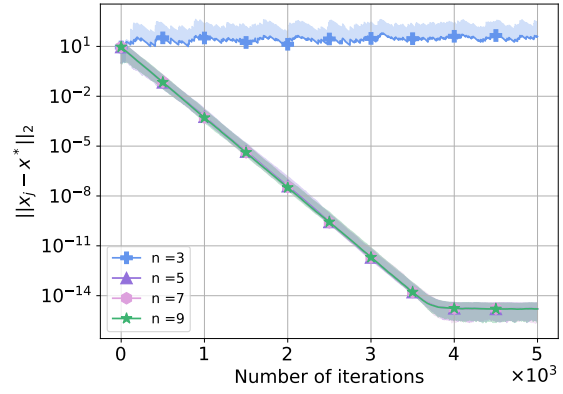
(d) Error vs. d_0 : $p = 0.6$; using Alg. 7 without block-list.

Figure 5.2: Effects of different data sizes d_0 on convergence: $N_r = 20$, $n_r = 4$, $k = 3$, and $\|e\|_\infty = 500$. The error norms were averaged over 50 trials (the solid lines) with 90% percentiles (the shaded areas).

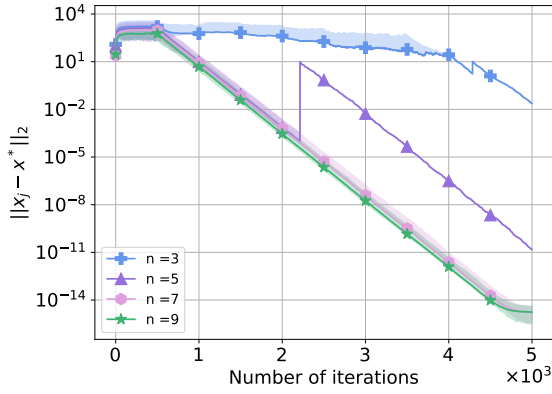
convergence oscillates for $n_r = 3$ and $n_r = 5$. Fig. 5.5 presents the effect of the adversary rate. As the adversary rate p increases, the accuracy decreases. Even though the adversary rate is large, the final results using the block-list method are still satisfying. Without the block-list, when the adversarial rate $p > 0.5$, the central server fails to approach the true solution due to the adversarial workers. This again shows the importance and effectiveness of using the block-list, especially in a highly hostile environment with a higher adversarial rate and a higher magnitude of the adversary. In addition, Fig. 5.4 shows the effect of the number of category types k with the block-list. The method converges as $k \rightarrow \infty$, i.e., with random noises when the adversarial rate is 0.2 and 0.6. In



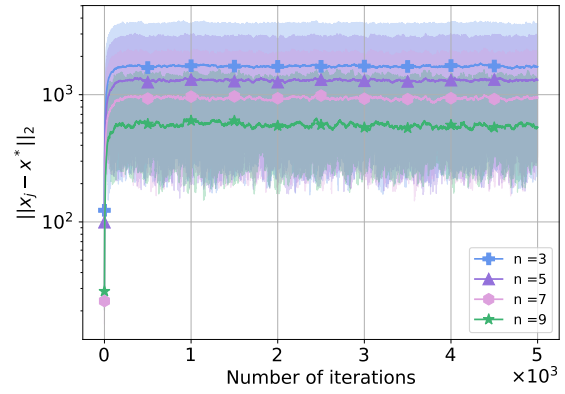
(a) Error vs. n_r : $p = 0.2$; using Alg. 7 with block-list.



(b) Error vs. n_r : $p = 0.2$; using Alg. 7 without block-list.



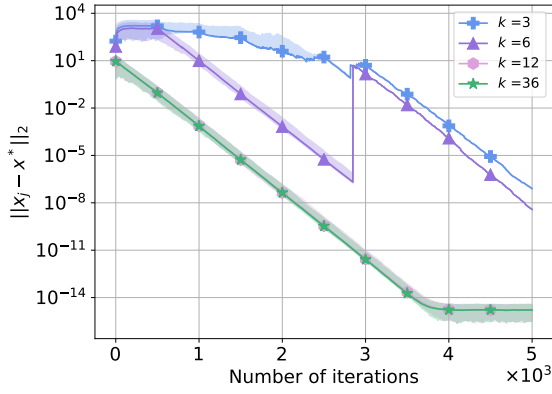
(c) Error vs. n_r : $p = 0.6$; using Alg. 7 with block-list.



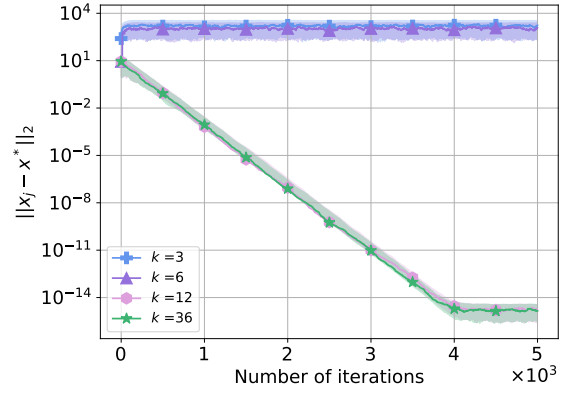
(d) Error vs. n_r : $p = 0.6$; using Alg. 7 without block-list.

Figure 5.3: Effects of the number of used workers n_r : $k = 3$, $d_0 = 6$, $N_r = 20$, $\|e\|_\infty = 5 \times 10^2$.

Fig. 5.6, we use the Wisconsin (Diagnostic) Breast Cancer data set, which includes data points whose features are computed from a digitized image of a fine needle aspirate (FNA) of a breast mass and describe characteristics of the cell nuclei present in the image ([DG17]). We set up the experiment similar to [HNR22]: the collection of data points forms our matrix $A \in \mathbb{R}^{569 \times 10}$. We then normalize A and construct x and b using a Gaussian distribution to form a consistent system. The convergence results in Fig. 5.6 show the effectiveness of our method solving this linear systems in a relatively safer environment with an adversarial rate $p = 0.3$ (Fig. 5.6a) and a more hostile environment with an adversarial rate $p = 0.6$ (Fig. 5.6b). When $p = 0.3$, the method converges within 1000 iterations, and as d_0 increases, the convergence speed becomes faster. Meanwhile,

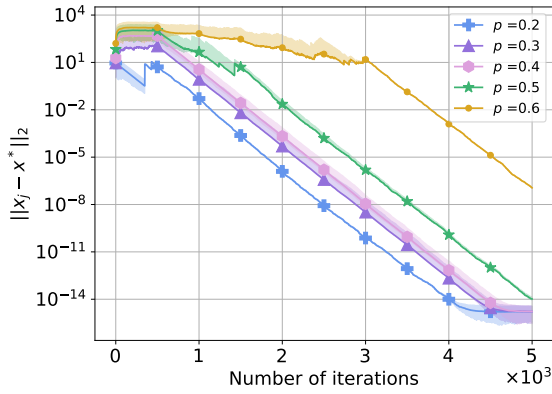


(a) Error vs. k : $p = 0.6$; using Alg. 7 with block-list.

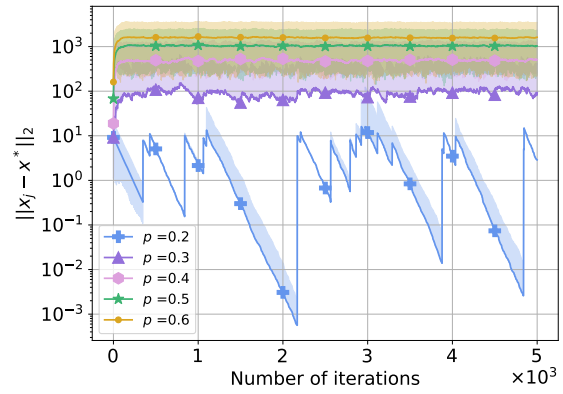


(b) Error vs. k : $p = 0.6$; using Alg. 7 without block-list.

Figure 5.4: Effects of the number of categories k . $N_r = 20, n_r = 4, d_0 = 4, \|e\|_\infty = 5 \times 10^2$



(a) Error vs. p : with block-list, $\|e\|_\infty = 10^{-3}$



(b) Error vs. p : without block-list, $\|e\|_\infty = 500$

Figure 5.5: Effects of the adversarial rate p , $d_0 = 3, N_r = 20, n_r = 4$, and $k = 3$. Squared error norms were averaged over 50 trials with the 90% percentiles

when $p = 0.3$, the method converges within 1500 iterations, and $d_0 = 8$ has the fastest convergence speed among all choices of d_0 .

Lastly, we study the effects of the number of iterations taken to update the block-list S . In Table 5.4, we calculate the accuracy of the block-list method when $S = 200, 500, 1000, 2000$. The two examples in the table show that as S increases, the accuracy is higher.

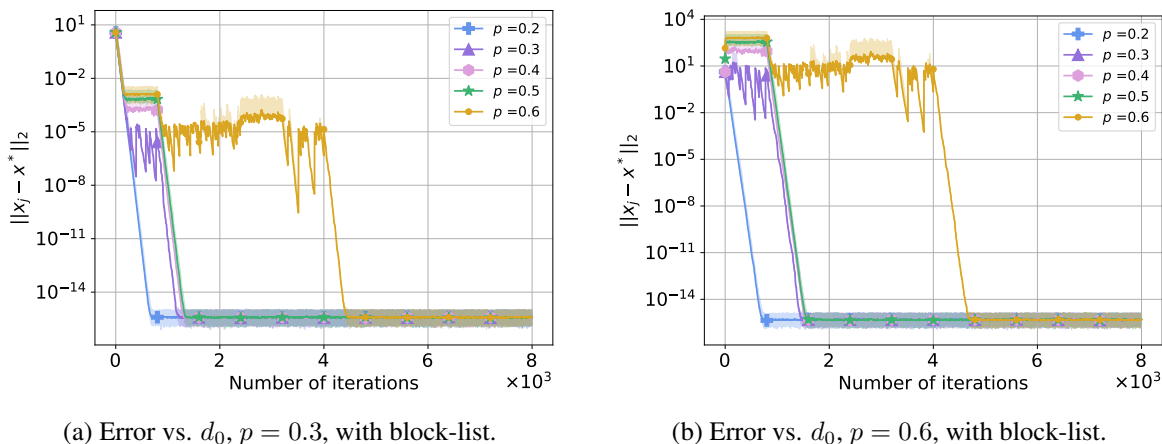


Figure 5.6: Effects of number of used row d_0 using the Breast Cancer Wisconsin data set, $N_r = 10$, $n_r = 4$, $k = 3$, $\|e\|_\infty = 500$.

Table 5.4: Accuracy of the method with the block-list when number of iterations to update the block-list $S = 200, 500, 1000, 2000$, $k = 3$, $N_r = 20$, and $n_r = 4$.

S	200	500	1000	2000
$p = 0.6, d_0 = 8$	0.75	0.792	0.875	0.875
$p = 0.4, d_0 = 6$	0.75	0.9375	1	1

5.5 Conclusion and future work

It is of great significant for optimization algorithms to be robust and resistant to adversaries. In this work, we propose efficient algorithms based on **mode** for solving large-scale linear systems with the presence of the adversarial workers. This kind of adversary has plenty of applications in the real life, e.g. IoT (Internet of Things). We provide theoretical convergence guarantee and our experiments support these theoretical results, as well as illustrate that the methods converge in many

scenarios. The methods are able to deal with different adversarial rates. In particular, the method with block-list is able to handle the adversarial rate $p > 0.5$, and at the same time, identify the adversarial workers. In the numerical simulations, we also present the effects of several important parameters of the adversaries and of anti-adversary strategies, namely, the number of used rows d_0 , the number of error categories k , the adversary rate p , and the number of chosen workers n at each iteration.

Currently, our methods assume that the good workers are more than each category of adversarial workers. It is straightforward to adjust the method so that the case where the good workers are the minority can be dealt with. We implemented the algorithms in a sequential manner to mimic distributed computation. However, to consider the storage overhead, one should deploy the algorithm in distributed systems. For future work, one can generalize this method to non-linear convex problems and perhaps non-convex problems.

5.6 Appendix

5.6.1 Single row convergence without block-list

We present the algorithm and the theory for a special case when the number of used rows d_0 is 1.

5.6.1.1 Algorithm

5.6.1.2 Convergence

Theorem 5.6.1. *Let $A \in \mathbb{R}^{d_1 \times d_2}$ with $d_1 \geq d_2$ and $b, e_1, \dots, e_k \in \mathbb{R}^{d_1}$. Assume that we solve $Ax^* = b$ via Algorithm 8, then*

$$\mathbb{E}\|x_i - x^*\|_2^2 \leq \alpha^{i+1}\|x_0 - x^*\|_2^2 + \frac{1 - \alpha^{i+1}}{1 - \alpha} \frac{1}{\|A\|_F^2} \sum_{\ell=1}^k q_\ell \|e_\ell\|^2, \quad (5.32)$$

Algorithm 8 DISTRIBUTED RANDOMIZED KACZMARZ WITH BLOCK-LIST

- 1: **Input:** Initialize **block-list** B , good worker set $D = [N]$, a counter vector $E = 0 \in \mathbb{R}^n$, MaxIter , Tol , $c_s = 2 \text{Tol}$, checking period T .
 - 2: **while** $j < \text{MaxIter}$ and $|c_s| > \text{Tol}$, **do**
 - 3: The central worker w_c selects a row index $i_j \in [m]$ with probability $p_{i_j} = \frac{\|A_{i_j}\|_2^2}{\|A\|_F^2}$
 - 4: Sample w_1, \dots, w_n uniformly from D
 - 5: Broadcast A_{i_j} to w_1, \dots, w_n
 - 6: w_s returns $c_s = \frac{\langle A_{i_j}, x_i \rangle - b_{i_j}}{\|A_{i_j}\|^2} + e_l$, if $w_s \in C_l$
 - 7: w_c splits $\{c_s\}_{s=1}^n$ into groups G_1, \dots, G_k and randomly choose from groups G_s that satisfy $|G_s| \geq n(1-p)$
 - 8: Update $x^{j+1} = x^j + c_{s_0} A_{i_j}^\top$
 - 9: Update $E(s) = E(s) + 1$, if $c_s \notin G_{s_0}$
 - 10: **if** $\text{mod}(j, T) = 0$ **then**
 - 11: Update B by checking the value of entries in E
 - 12: $D = D \setminus B$
 - 13: **end if**
 - 14: Update $j = j + 1$
 - 15: **end while**
 - 16: **Output:** x^j and B
-

where $\sigma_{\min}^2(A)$ is the smallest singular value of A , $\alpha = 1 - \frac{\sigma_{\min}^2(A)}{\|A\|_F^2}$ and $q_\ell = \frac{\hat{q}_{\text{mode}}^\ell}{q}$.

Additionally, if $\|e_\ell\| \leq C$, we have

$$\mathbb{E}\|x_i - x^*\|_2^2 \leq \alpha^{i+1} \|x_0 - x^*\|_2^2 + \frac{1 - \alpha^{i+1}}{1 - \alpha} \frac{Cq_0}{\|A\|_F^2}. \quad (5.33)$$

In (5.33) we use the fact that $\sum_{\ell=0}^k q_\ell = 1$. Furthermore, we assume that $\mathbb{E}\|e_\ell\|^2 = d\sigma_\ell^2$ at each iteration.

To provide a quantitative understanding of Theorem 5.6.1, we present several examples in Tables 5.5 and 5.6. For simplicity, assume that each error category has the same fraction $p_\ell = p/k$. Thus, all $\hat{q}_{\text{mode}}^\ell$ are equal. Here q_0 is the probability that the algorithm chooses the right mode and q is the probability that there is a mode. In these two tables, we present the values for $\hat{q}_{\text{mode}}^\ell$, \hat{q}_{mode}^0 , q and q_0 by varying the number of error categories k , the number of chosen workers n and the adversarial rate p . These two tables are generated by solving a linear system with a row-normalized matrix

$A \in \mathbb{R}^{1000 \times 100}$. As k increases, q_ℓ decreases and q_0 increases. Therefore, the error bound in equation (5.33) decreases with respect to k and thus reaches better convergence. When k is large enough, $q_\ell \approx 0$. Therefore, when the noise is uniformly random error and there is a mode for the step-size, the mode will be the correct mode. As n increases, there is a similar decrease effect and therefore a better convergence.

Table 5.5: Total number of workers $N = 100$, number of chosen workers $n = 5$.

p	k	\hat{q}_{mode}^ℓ	\hat{q}_{mode}^0	q	q_0
0.8	5	0.1	0.16	0.67	0.15
	10	0.04	0.21	0.57	0.36
	15	0.02	0.23	0.48	0.46
0.2	3	0.002	0.63	0.64	0.98
	5	8×10^{-4}	0.65	0.65	0.99
	10	2×10^{-4}	0.66	0.67	0.99
	15	2×10^{-4}	0.685	0.689	0.99

Table 5.6: Total number of workers $N = 100$, number of error categories $k = 5$.

p	n	\hat{q}_{mode}^ℓ	\hat{q}_{mode}^0	q	q_0
0.8	10	0.099	0.18	0.67	0.26
	15	0.099	0.2	0.7	0.29
	20	0.097	0.23	0.71	0.31
0.2	10	7×10^{-6}	0.904	0.90	$1 - 5 \times 10^{-6}$
	15	5×10^{-7}	0.97	0.97	$1 - 3 \times 10^{-6}$
	20	1×10^{-7}	0.99	0.99	$1 - 6 \times 10^{-7}$

Theorem 5.6.2. Let n , k , and each entry in $\{\hat{n}_i\} < n$ be non-negative integers. The number of solutions to the equations:

$$\begin{cases} n_1 + n_2 + \cdots + n_k = n \\ \text{s.t. } 0 \leq n_i \leq \hat{n}_i \end{cases}$$

is the coefficient of term x^n in the polynomial

$$f(x) = \left(\sum_{j=0}^n \binom{\hat{n}_i}{j} x^j \right)^k$$

Proof. Consider we have k bins. The constrained problem is equivalent to choosing n apples (different from each other) from k baskets; in each basket n_i has \hat{n}_i apples, in total.

We first consider the problem without the constraint and every apple is the same. The number of ways to choose the apples equals the coefficient of the term x^n of the polynomial $f(x) = \left(\sum_{j=0}^{\hat{n}_i} x^j\right)^k$. With the constraint and apples being different, there are $\binom{\hat{n}_i}{j}$ ways to choose j apples from basket i . Let j vary and since we have k baskets, in total the number of ways to choose equals the coefficient of the term x^n in the polynomial $f(x)$. \square

Here is the proof of Theorem 5.6.1:

Proof. To prove Equation (5.32), at each iteration, we consider solving $Ax = b$, $Ax = b + e_1, \dots$, $Ax = b + e_k$ with probability q_0, q_1, \dots, q_k , respectively. Therefore, for the $(i+1)$ -th step, we have the iteration

$$x_{i+1} = x_i - \frac{\langle A_j, x_i \rangle - b_j}{\|A_j\|^2} (A_j)^\top,$$

or

$$x_{i+1} = x_i - \frac{\langle A_j, x_i \rangle - (b_j + e_\ell(j))}{\|A_j\|^2} (A_j)^\top.$$

for $\ell = 1, \dots, k$, A_j is the j -th row of matrix A .

Notice that when $x_{i+1} = x_i - \frac{\langle A_j, x_i \rangle - b_j}{\|A_j\|^2} (A_j)^\top$, we have

$$\begin{aligned} & \mathbb{E}_j \|x_{i+1} - x^*\|_2^2 \\ &= \mathbb{E}_j \left\| x_i - \frac{\langle A_j, x_i \rangle - b_j}{\|A_j\|^2} (A_j)^\top - x^* \right\|_2^2 \\ &\leq \left(1 - \frac{\sigma_{\min}^2(A)}{\|A\|_F^2} \right) \|x_i - x^*\|_2^2. \end{aligned} \tag{5.34}$$

When $x_{i+1} = x_i - \frac{\langle A_j, x_i \rangle - (b_j + e_\ell(j))}{\|A_j\|_2^2} (A_j)^\top$, we have

$$\begin{aligned}
& \mathbb{E}_j \|x_{i+1} - x^*\|_2^2 \\
&= \mathbb{E}_j \left\| x_i - \frac{\langle A_j, x_i \rangle - (b_j + e_\ell(j))}{\|A_j\|_2^2} A_j^\top(j, :) - x^* \right\|_2^2 \\
&\leq \left(1 - \frac{\sigma_{\min}^2(A)}{\|A\|_F^2} \right) \|x_i - x^*\|_2^2 + \mathbb{E}_j \frac{e_\ell^2(j)}{\|A_j\|_2^2} \\
&= \left(1 - \frac{\sigma_{\min}^2(A)}{\|A\|_F^2} \right) \|x_i - x^*\|_2^2 + \frac{\|e_\ell\|_2^2}{\|A\|_F^2}.
\end{aligned} \tag{5.35}$$

Combining (5.34) and (5.35), we have

$$\begin{aligned}
& \mathbb{E}_j \|x_{i+1} - x^*\|_2^2 \\
&\leq \left(1 - \frac{\sigma_{\min}^2(A)}{\|A\|_F^2} \right) \|x_i - x^*\|_2^2 + \frac{1}{\|A\|_F^2} \sum_{\ell=1}^k q_\ell \|e_\ell\|_2^2.
\end{aligned} \tag{5.36}$$

Set $\alpha = 1 - \frac{\sigma_{\min}^2(A)}{\|A\|_F^2}$. Therefore,

$$\begin{aligned}
& \mathbb{E} \|x_{i+1} - x^*\|_2^2 \\
&\leq \alpha^{i+1} \|x_0 - x^*\|_2^2 + \frac{1 - \alpha^{i+1}}{1 - \alpha} \frac{1}{\|A\|_F^2} \sum_{\ell=1}^k q_\ell \|e_\ell\|_2^2.
\end{aligned} \tag{5.37}$$

□

5.6.2 Finding the optimal d_0

Remark 5.6.3. *If we assume each row is held by the same number of workers and used workers, i.e. $n_r \equiv n$, $N_r \equiv N$ and the probability of each error category ℓ ($\ell \neq 0$) for each row is the same, i.e.*

$p_{r,\ell} \equiv p/k$. Then $g_0(r) = \max(\lceil \frac{n_r}{k+1} \rceil, \lceil n_r(1-p_r) \rceil)$. Moreover, we have

$$\begin{aligned}
b_g^r &\equiv b_g, \text{ the coefficient of term } x^n \text{ of } \left(\sum_{j=0}^{g-1} \binom{Np/k}{j} x^j \right)^{k-1} \left(\sum_{j=0}^{g-1} \binom{N(1-p)}{j} x^j \right), \\
a_{g,\ell}^r &\equiv a_g, \text{ the coefficient of term } x^{n-g} \text{ of } \left(\sum_{j=0}^{g-1} \binom{Np/k}{j} x^j \right)^{k-1}, \\
q_g^t &\equiv q_g = \frac{\binom{Np/k}{g} a_g}{\binom{N}{n}}, g_0(t) = \max(\lceil \frac{n}{k+1} \rceil, \lceil n(1-p) \rceil), \\
Q(t, \tau_i) &\equiv \sum_{g=g_0}^n q_g \left(\frac{b_g}{\binom{N}{n}} \right)^{d_0-1} = Q_{\max} = Q_{\min} := Q, \\
q_t &= \frac{d_0}{d_1} Q_{\max}.
\end{aligned} \tag{5.38}$$

Let $\alpha(d_0) = \alpha = 1 - Q \frac{d_0}{d_1} \sigma_{\min}^2(\tilde{A})$. Then the convergence error is bounded:

$$\begin{aligned}
\mathbb{E} \|x_i - x^*\|_2^2 &\leq \alpha \|x_{i-1} - x^*\|_2^2 + \sum_{t \in [d_1]} q_t \|\tilde{e}_t\|_2^2 \\
&= \alpha^i \|x_0 - x^*\|_2^2 + \frac{1 - \alpha^{i+1}}{1 - \alpha} \sum_{t \in [d_1]} q_t \|\tilde{e}_t\|_2^2.
\end{aligned} \tag{5.39}$$

To study the relation between d_0 and the convergence, consider

$$\frac{\partial \alpha(d_0)}{\partial d_0} \propto - \sum_{g=g_0}^n q_g \left(1 + d_0 \log\left(\frac{b_g}{\binom{N}{n}}\right) \right) \left(\frac{b_g}{\binom{N}{n}} \right)^{d_0-1}. \tag{5.40}$$

If $d_0 \geq -\frac{1}{\log(\frac{b_g}{\binom{N}{n}})}$ for all g , then $\frac{\partial \alpha(d_0)}{\partial d_0} \geq 0$. This implies that as d_0 increases, $\alpha(d_0)$ increases.

Remark 5.6.4. When $g_0 = n$, we have

$$\frac{\partial \alpha(d_0)}{\partial d_0} \propto - \left(1 + d_0 \log\left(\frac{b_n}{\binom{N}{n}}\right) \right) \left(\frac{b_n}{\binom{N}{n}} \right)^{d_0-1},$$

and to reach the fastest convergence rate, $d_0 = -\frac{1}{\log(\frac{b_g}{\binom{N}{n}})}$. One can explore the minimizers for α in more general cases, where multiple local minimizers could present in the landscape.

5.6.3 Other proofs

Here is the proof of lemma 5.3.12:

Proof. The probability that a good worker is put in the block-list is

$$p_0 N \mathbb{P}_{bl}(w_0).$$

The probability that a bad worker $w_{\ell'}$ in category ℓ' is put in the block-list is:

$$p_{\ell'} N \mathbb{P}_{bl}(w_{\ell'}).$$

The probability that a worker, either good or bad, is put in the block-list is

$$\sum_{\ell=0}^k p_{\ell} N \mathbb{P}_{bl}(w_{\ell}).$$

□

CHAPTER 6

Conclusion

In this thesis, we studied modeling methods for different aspects of COVID-19 and algorithms for least-squares problems in distributed systems with adversarial workers.

In chapter 2, we analyzed theoretically and quantitatively the *finite-size effects* arising in stochastic compartmental models. We applied a martingale approach can be applied to the stochastic model and show that the fluid limit is indeed the deterministic SIR model by providing a bound of the variances of the martingale. We found a theoretical explanation for the *finite-size effects* by observing that the stochastic component of the martingale formulation scales as the inverse of the square root of the population size. A larger variance both in the outbreak size and its temporal behavior arises as population size decreases. Our work provides a good guide for authorities of smaller populations to estimate risk over time in order to prepare for the outbreak. It is important to bear in mind that the broader variations in the pandemic caused by the smaller population would lead to a wide outcome when it comes to estimating risk.

In chapter 3, we proposed a policy-making model coupling the SIR model for one region and multiple regions. We introduced an existing approach of optimal control in the literature and reproduce the results using our method. We also discussed the different policies and pandemic dynamics resulting from different minimal policy time intervals and different parameters. In chapter 4, we proposed a HNMF to organize existing literature on coronaviruses and pandemics, and early literature on COVID-19 into an interactive structure easily searchable by researchers and available to use through a corresponding website. The topics discovered by HNMF reveal that early research of interest to the COVID-19 research community divides into diverse areas such as research related

to other coronaviruses, research related to other respiratory diseases, virology and genetic research, as well as research relating to the public health response.

Lastly, we proposed efficient algorithms and provide theoretical convergence guarantee for solving the least squares problem with the presence of the adversarial workers in distributed systems (chapter 5). Our methods are able to deal with different adversarial rates as large as $p > 0.5$. Additionally, the algorithm identifies the adversarial workers. We also present the effect of several important parameters of the adversaries and of the anti-adversaries strategy, namely, the number of error categories k , the adversary rate p , and the number of chosen workers n at each iteration.

The models we proposed are inspired by the real-world applications and yet they still have limitations. For example, the SIR model and the excitation matrix used in our analysis assume a perfect mixing in the population. In the policy function, the action parameter α is a heuristic representation of the lockdown, social distancing and mask policies. It remains to be discussed how other policies, for example, vaccination policies, affect the spreading in the different stages of a pandemic. In our study of the adversarial learning, we assumed that the good workers are the majority compared to workers from other error categories. However, the method can be easily generalized to the case where the good workers are minority. Our models are by no means the most practical. However, we believe that they improve upon previous works, serve as landmark for real situations and provide insights for future research.

REFERENCES

- [AA03] Linda J.S. Allen and Edward J. Allen. “A Comparison of Three Different Stochastic Population Models with Regard to Persistence Time.” *Theoretical Population Biology*, **64**(4):439–449, 2003.
- [AAL18] Dan Alistarh, Zeyuan Allen-Zh, and Jerry Li. “Byzantine Stochastic Gradient Descent.” 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.
- [AG18] Elena Almaraz and Antonio Gómez-Corral. “On SIR-models with Markov-modulated Events: Length of an Outbreak, Total Size of the Epidemic and Number of Secondary Infections.” *Discrete and Continuous Dynamical Systems. Series B. A Journal Bridging Mathematics and Sciences*, **23**(6):2153–2176, 2018.
- [Agm54] Shmuel Agmon. “The Relaxation Method for Linear Inequalities.” *Canadian Journal of Mathematics*, **6**:382–392, 1954.
- [AKL21] Parham Azimi, Zahra Keshavarz, Jose Guillermo Cedeno Laurent, Brent Stephens, and Joseph G. Allen. “Mechanistic Transmission Modeling of COVID-19 on the Diamond Princess Cruise Ship Demonstrates the Importance of Aerosol Transmission.” *Proceedings of the National Academy of Sciences*, **118**(8), 2021.
- [All08] Linda J. S. Allen. “An Introduction to Stochastic Epidemic Models.” In *Mathematical epidemiology*, volume 1945 of *Lecture Notes in Math.*, pp. 81–130. Springer, Berlin, 2008.
- [All11] Linda J. S. Allen. *An Introduction to Stochastic Processes with Applications to Biology*. CRC Press, Boca Raton, FL, second edition, 2011.
- [App09] David Applebaum. *Lévy Processes and Stochastic Calculus*, volume 116 of *Cambridge Studies in Advanced Mathematics*. Cambridge University Press, Cambridge, second edition, 2009.
- [ASM21] Alaa Abd-Alrazaq, Jens Schneider, Borbala Mifsud, Tanvir Alam, Mowafa Househ, Mounir Hamdi, and Zubair Shah. “A Comprehensive Overview of the COVID-19 Literature: Machine Learning–Based Bibliometric Analysis.” *Journal of Medical Internet Research*, **23**(3):e23703, Mar 2021.
- [ASN17] Melissa Ailem, Aghiles Salah, and Mohamed Nadif. “Non-negative Matrix Factorization Meets Word Embedding.” In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 1081–1084, 2017.

- [BAC20] Tahar Zamene Boulmezaoud, Luis Álvarez, Miguel Colom, and Jean-Michel Morel. “A Daily Measure of the SARS-CoV-2 Effective Reproduction Number for all Countries.” *Image Processing On Line*, **10**:191–210, 2020.
- [Bai75] Norman T. J. Bailey. *The Mathematical Theory of Infectious Diseases and Its Applications*. Hafner Press [Macmillan Publishing Co., Inc.] New York, second edition, 1975.
- [Bar60] M. S. Bartlett. “The Critical Community Size for Measles in the United States.” *Journal of the Royal Statistical Society: Series A (General)*, **123**(1):37–44, 1960.
- [BBC21] Tom Belin, Andrea Bertozzi, Nishchal Chaudhary, Todd Graves, Jeffrey Guterman, M. Claire Jarashow, Roger J. Lewis, Joe Marion, Frederic Schoenberg, Megha Shah, Juliana Tolles, Elizabeth Traub, Kert Viele, and Fei Wu. “Projections of Hospital-based Healthcare Demand due to COVID-19 in Los Angeles County May 24, 2021.”, 2021.
- [BD21] Pierre-Alexandre Bliman and Michel Duprez. “How Best Can Finite-time Social Distancing Reduce Epidemic Final Size?” *Journal of Theoretical Biology*, **511**:110557, 2021.
- [BDP21] Pierre-Alexandre Bliman, Michel Duprez, Yannick Privat, and Nicolas Vauchelet. “Optimal Immunity Control and Final Size Minimization by Social Distancing for the SIR Epidemic Model.” *Journal of Optimization Theory and Applications*, 03 2021.
- [BEC06] Donald S. Burke, Joshua M. Epstein, Derek A.T. Cummings, Jon I. Parker, Kenneth C. Cline, Ramesh M. Singa, and Shubha Chakravarty. “Individual-based Computational Modeling of Smallpox Epidemic Control Strategies.” *Academic Emergency Medicine*, **13**(11):1142–1149, 2006.
- [BF07] Martin C. J. Bootsma and Neil M. Ferguson. “The Effect of Public Health Measures on the 1918 Influenza Pandemic in U.S. Cities.” *Proceedings of the National Academy of Sciences*, **104**(18):7588–7593, 2007.
- [BFM20] Andrea L. Bertozzi, Elisa Franco, George Mohler, Martin B. Short, and Daniel Sledge. “The Challenges of Modeling and Forecasting the Spread of COVID-19.” *Proceedings of the National Academy of Sciences*, **117**(29):16732–16738, 2020.
- [Bic02] Klaus Bichteler. *Stochastic Integration with Jumps*, volume 89 of *Encyclopedia of Mathematics and its Applications*. Cambridge University Press, Cambridge, 2002.
- [BIH18] Elizabeth Buckingham-Jeffery, Valerie Isham, and Thomas House. “Gaussian Process Approximations for Fast Inference from Infectious Disease Data.” *Mathematical Biosciences*, **301**:111–120, 2018.
- [bio20a] bioRxiv. “COVID-19 SARS-CoV-2 Preprints from medRxiv and bioRxiv.” <https://connect.biorxiv.org/relate/content/181>, 2020. Accessed: 2020-04-20.

- [Bio20b] National Institute of Health National Center for Biotechnology Information. “COVID-19 Article Collection from LitCovid.” [https://www.ncbi.nlm.nih.gov/research/pubtator/?view=docsum&query=\\$LitCovid](https://www.ncbi.nlm.nih.gov/research/pubtator/?view=docsum&query=$LitCovid), 2020. Accessed: 2020-04-20.
- [BKL09] Steven Bird, Ewan Klein, and Edward Loper. *Natural Language Processing with Python: Analyzing Text with the Natural Language Toolkit.* ” O’Reilly Media, Inc.”, 2009.
- [BKM02] SM Blower, Katia Koelle, and John Mills. “Health Policy Modeling: Epidemic Control, HIV Vaccines, and Risky Behavior.” *Quantitative evaluation of HIV prevention programs*, pp. 260–289, 2002.
- [BL07] David M. Blei and John D. Lafferty. “A correlated topic model of Science.” *The Annals of Applied Statistics*, **1**(1):17 – 35, 2007.
- [Br] Bo Martin Bibby and Michael Sørensen. “Martingale Estimation Functions for Discretely Observed Diffusion Processes.” *Official Journal of the Bernoulli Society for Mathematical Statistics and Probability*.
- [BRQ20] Farrah Kristel Batista, Angel Martín del Rey, and Araceli Queiruga-Dios. “A Review of SEIR-D Agent-Based Model.” In Enrique Herrera-Viedma, Zita Vale, Peter Nielsen, Angel Martin Del Rey, and Roberto Casado Vara, editors, *Distributed Computing and Artificial Intelligence, 16th International Conference, Special Sessions*, pp. 133–140, Cham, 2020. Springer International Publishing.
- [Buc08] Ioan Buciu. “Non-negative Matrix Factorization, A New Tool for Feature Extraction: Theory and Applications.”, 2008.
- [BW19] Zhong-Zhi Bai and Wen-Ting Wu. “On Partially Randomized Extended Kaczmarz Method for Solving Large Sparse Overdetermined inconsistent linear systems.” *Linear Algebra and Its Applications*, **578**:225–250, 2019.
- [BW21] Zhong-Zhi Bai and Wen-Ting Wu. “On Greedy Randomized Augmented Kaczmarz Method for Solving Large Sparse Inconsistent Linear Systems.” *SIAM Journal on Scientific Computing*, **43**(6):A3892–A3911, 2021.
- [BWR19] Rawad Bitar, Mary Wootters, and Salim El Rouayheb. pp. 1–5, 2019.
- [Cena] Centers for Disease Control and Prevention. “CDC Museum COVID-19 Timeline.” <https://www.cdc.gov/museum/timeline/covid19.html>. Accessed: 2022-06-20.
- [Cenb] Centers for Disease Control and Prevention. “U.S. State and Territorial Stay-At-Home Orders: March 15, 2020 – August 15, 2021 by County by Day.” <https://data.cdc.gov/Policy-Surveillance/U-S-State-and-Territorial-Stay-At-Home-Orders-Marc/y2iy-8irm>. Accessed: 2022-06-20.

- [CGS13] Fabienne Comte, Valentine Genon-Catalot, and Adeline L. Samson. “Nonparametric Estimation for Stochastic Differential Equations with Random Effects.” *Stochastic Processes and their Applications*, **123**(7):2522–2551, 2013.
- [CLZ19] Yuejie Chi, Yuanxin Li, Huishuai Zhang, and Yingbin Liang. “Median-Truncated Gradient Descent: A Robust and Scalable Nonconvex Approach for Signal Estimation.” 2019.
- [Cru21] Princess Cruises. “Princess Cruise Lines (2020) Diamond Princess Updates.”, May 2021.
- [CTW19] Yifan Chang, Justin C Tzou, Michael Jeffrey Ward, and Jun cheng Wei. “Refined Stability Thresholds for Localized Spot Patterns for the Brusselator Model in \mathbb{R}^2 .” *European Journal of Applied Mathematics*, **30**(4):791–828, 2019.
- [CW14] Kai lai Chung. and Ruth J. Williams. *Introduction to Stochastic Integration*. Modern Birkhäuser Classics. Birkhäuser/Springer, New York, second edition, 2014.
- [DCL19] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding.” pp. 4171–4186, 2019.
- [DDG20] Ensheng Dong, Hongru Du, and Lauren Gardner. “An Interactive Web-based Dashboard to Track COVID-19 in Real Time.” *The Lancet Infectious Diseases*, 2020.
- [Deb54] Gerard Debreu. “Valuation Equilibrium and Pareto Optimum.” *Proceedings of the National Academy of Sciences*, **40**(7):588–592, 1954.
- [Dep22] Department of Public Health, Los Angeles. “Department of Public Health, Los Angeles.” <http://publichealth.lacounty.gov/>, 3 2022. Accessed: 2022-06-20.
- [DG17] Dheeru Dua and Casey Graff. “UCI Machine Learning Repository.”, 2017.
- [DP] Centers for Disease Control and Prevention. “COVID-19 research articles downloadable database.” <https://www.cdc.gov/library/researchguides/2019novelcoronavirus/researcharticles.html>. Accessed: 2020-04-20.
- [Dur96] Richard Durrett. *Stochastic Calculus*. Probability and Stochastics Series. CRC Press, Boca Raton, FL, 1996. A practical introduction.
- [Dur99] Rick Durrett. *Essentials of Stochastic Processes*. Springer Texts in Statistics. Springer-Verlag, New York, 1999.
- [Dur02] Rick Durrett. *Probability Models for DNA Sequence Evolution*. Probability and its Applications (New York). Springer-Verlag, New York, 2002.

- [EFG20] El-Mahdi El-Mhamdi, Sadegh Farhadkhani, Rachid Guerraoui, Arsany Guirguis, Lê Nguyễn Hoàng, and Sébastien Rouault. “Collaborative Learning as an Agreement Problem.” *Clinical Orthopaedics and Related Research*, **abs/2008.00742**, 2020.
- [EK86] Stewart N. Ethier and Thomas G. Kurtz. *Markov Processes*. Wiley Series in Probability and Mathematical Statistics: Probability and Mathematical Statistics. John Wiley & Sons, Inc., New York, 1986. Characterization and convergence.
- [FL98] Drew Fudenberg and David K. Levine. *The Theory of Learning in Games*, volume 1 of *MIT Press Books*. The MIT Press, 1998.
- [FLN20] Neil M Ferguson, Daniel Laydon, Gemma Nedjati-Gilani, Natsuko Imai, Kylie Ainslie, Marc Baguelin, Sangeeta Bhatia, Adhiratha Boonyasiri, Zulma Cucunubá, Gina Cuomo-Dannenburg, Amy Dighe, Ilaria Dorigatti, Han Fu, Katy Gaythorpe, Will Green, Arran Hamlet, Wes Hinsley, Lucy C Okell, Sabine Van Elsland, Hayley Thompson, Robert Verity, Erik Volz, Haowei Wang, Yuanrong Wang, Patrick GT Walker, Caroline Walters, Peter Winskill, Charles Whittaker, Christl A Donnelly, Steven Riley, and Azra C Ghani. “Impact of Non-pharmaceutical Interventions (NPIs) to Reduce COVID-19 Mortality and Healthcare Demand.” *Report 9, Imperial College COVID-19 Response Team, Imperial College London, London, United Kingdom*, 2020.
- [FXM14] Jiashi Feng, Huan Xu, and Shie Mannor. “Distributed Robust Learning.”, 2014.
- [GBH70] Richard Gordon, Robert Bender, and Gabor T. Herman. “Algebraic Reconstruction Techniques (ART) for Three-dimensional Electron Microscopy and X-ray Photography.” *Journal of Theoretical Biology*, **29**(3):471–481, 1970.
- [GEG03] Rebecca F Grais, J Hugh Ellis, and Gregory E Glass. “Assessing the Impact of Air-line Travel on the Geographic Spread of Pandemic Influenza.” *European journal of epidemiology*, **18**(11):1065–1072, 2003.
- [GFH20] Jonas Geiping, Liam Fowl, W. Ronny Huang, Wojciech Czaja, Gavin Taylor, Michael Moeller, and Tom Goldstein. “Witches’ Brew: Industrial Scale Data Poisoning via Gradient Matching.” *Clinical Orthopaedics and Related Research*, **abs/2009.02276**, 2020.
- [GHH20] Rachel Grotheer, Longxiu Huang, Yihuan Huang, Alona Kryshchenko, Oleksandr Kryshchenko, Pengyu Li, Xia Li, Elizaveta Rebrova, Kyung Ha, and Deanna Needell. “COVID-19 Literature Topic-Based Search via Hierarchical NMF.” In *Workshop on NLP for COVID-19 (Part 2) at EMNLP 2020*, 2020.
- [GHJ75] Richard Gordon, Gabor T Herman, and Steven A Johnson. “Image Reconstruction from Projections.” *Scientific American*, **233**(4):56–71, 1975.

- [GHM19] Mengdi Gao, Jamie Haddock, Denali Molitor, Deanna Needell, Eli Sadvnik, Tyler Will, and Runyu Zhang. “Neural Nonnegative Matrix Factorization for Hierarchical Multilayer Topic Modeling.” In *2019 IEEE 8th International Workshop on Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP)*, pp. 6–10. IEEE, 2019.
- [Gil76] Daniel T Gillespie. “A General Method for Numerically Simulating the Stochastic Time Evolution of Coupled Chemical Reactions.” *Journal of Computational Physics*, **22**(4):403–434, 1976.
- [GKG21] Xiaolong Geng, Gabriel G. Katul, Firas Gerges, Elie Bou-Zeid, Hani Nassif, and Michel C. Boufadel. “A Kernel-modulated SIR model for COVID-19 Contagious Spread from County to Continent.” *Proceedings of the National Academy of Sciences*, **118**(21):e2023321118, 2021.
- [GMP18] Ian Goodfellow, Patrick McDaniel, and Nicolas Papernot. “Making Machine Learning Robust against Adversarial Inputs.” *Communications of the ACM*, **61**(7):56–66, jun 2018.
- [GPV88] M. Z. Guo, George C. Papanicolaou, and S. R. Srinivasa Varadhan. “Nonlinear Diffusion Limit for a System with Nearest Neighbor Interactions.” *Communications in Mathematical Physics*, **118**(1):31–59, 1988.
- [GWW19] Daniel Gomez, Michael J. Ward, and Juncheng Wei. “The Linear Stability of Symmetric Spike Patterns for a Bulk-membrane Coupled Gierer-Meinhardt Model.” *SIAM Journal on Applied Dynamical Systems*, **18**(2):729–768, 2019.
- [GZ17] Volker Gerdtts and Alexander Zakhartchouk. “Vaccines for Porcine Epidemic Diarrhea Virus and Other Swine Coronaviruses.” *Veterinary microbiology*, **206**:45–51, 2017.
- [Ham18] Fujihiro Hamba. “Turbulent Energy Density in Scale Space for Inhomogeneous Turbulence.” *Journal of Fluid Mechanics*, **842**:532–553, 2018.
- [Het00] Herbert W. Hethcote. “The Mathematics of Infectious Diseases.” *SIAM Review*, **42**(4):599–653, 2000.
- [HM93] Gabor T Herman and Lorraine B Meyer. “Algebraic Reconstruction Techniques Can be Made Computationally Efficient (Positron Emission Tomography Application).” *IEEE transactions on medical imaging*, **12**(3):600–609, 1993.
- [HM21] Jamie Haddock and Anna Ma. “Greed Works: An Improved Analysis of Sampling Kaczmarz–Motzkin.” *SIAM Journal on Mathematics of Data Science*, **3**(1):342–368, 2021.

- [HMZ18] Ismail Hameduddin, Charles Meneveau, Tamer A. Zaki, and Dennice F. Gayme. “Geometric decomposition of the conformation tensor in viscoelastic turbulence.” *Journal of Fluid Mechanics*, **842**:395–427, 2018.
- [HNR22] Jamie Haddock, Deanna Needell, Elizaveta Rebrova, and William Swartworth. “Quantile-Based Iterative Methods for Corrupted Systems of Linear Equations.” *SIAM Journal on Matrix Analysis and Applications*, **43**(2):605–637, 2022.
- [HPS72] Paul G. Hoel, Sidney C. Port, and Charles J. Stone. *Introduction to Stochastic Processes*. Houghton Mifflin Co., Boston, Mass., 1972. The Houghton Mifflin Series in Statistics.
- [HW20] Labour Ministry of Health and Japan Welfare. “Ministry of Health, Labor and Welfare, Japan (2020) About Coronavirus Disease 2019 (COVID-19).”, March 2020.
- [HWY92] Sheng Wu He, Jia Gang Wang, and Jia An Yan. *Semimartingale Theory and Stochastic Calculus*. Kexue Chubanshe (Science Press), Beijing; CRC Press, Boca Raton, FL, 1992.
- [IMM18] Simon J. Illingworth, Jason P. Monty, and Ivan Marusic. “Estimating Large-scale Structures in Wall Turbulence Using Linear Models.” *Journal of Fluid Mechanics*, **842**:146–162, 2018.
- [Ish91] Valerie Isham. “Assessing the Variability of Stochastic Epidemics.” *Mathematical Biosciences*, **107**(2):209–224, 1991.
- [Ish93] Valerie Isham. “Stochastic Models for Epidemics with Special Reference to AIDS.” *The Annals of Applied Probability*, **3**(1):1–27, 1993.
- [Ish05] Valerie Isham. “Stochastic Models for Epidemics.” In *Celebrating statistics*, volume 33 of *Oxford Statistical Science Series*, pp. 27–54. Oxford Univ. Press, Oxford, 2005.
- [Jim18] Javier Jiménez. “Coherent Structures in Wall-bounded Turbulence.” *Journal of Fluid Mechanics*, **842**:P1, 100, 2018.
- [JML21] Feiran Jia, Aditya Mate, Zun Li, Shahin Jabbari, Mithun Chakraborty, Milind Tambe, Michael Wellman, and Yevgeniy Vorobeychik. “A Game-Theoretic Approach for Hierarchical Policy-Making.” *arXiv preprint arXiv:2102.10646*, 2021.
- [Kac37] Stefan Kaczmarz. “Angenäherte Auflösung Von Systemen Llinearer Gleichungen.” *Bulletin International de l’ Académie Polonaise des Sciences et des Lettres*, pp. 355–357, 1937.
- [Kag20] Kaggle. “COVID-19 Open Research Dataset Challenge (CORD-19).”, 2020. Accessed: 2020-04-20.

- [KCM05] Isthinayagy Krishnarajah, Alex Cook, Glenn Marion, and Gavin Gibson. “Novel moment closure approximations in stochastic epidemics.” *Bulletin of Mathematical Biology*, **67**(4):855–873, 2005.
- [KCP15] Da Kuang, Jaegul Choo, and Haesun Park. “Nonnegative Matrix Factorization for Interactive Topic Modeling and Document Clustering.” In *Partitional Clustering Algorithms*, pp. 215–243. Springer, 2015.
- [KL99] Claude Kipnis and Claudio Landim. *Scaling Limits of Interacting Particle Systems*, volume 320 of *Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences]*. Springer-Verlag, Berlin, 1999.
- [KMS17] Istvan Z Kiss, Joel C Miller, and Peter Simon. (Book) *Mathematics of Epidemics on Networks: from Exact to Approximate Models*. Springer, 2017.
- [KMW27] William Ogilvy Kermack, A. G. McKendrick, and Gilbert Thomas Walker. “A Contribution to the Mathematical Theory of Epidemics.” *Proceedings of the Royal Society of London. Series A, Containing Papers of a Mathematical and Physical Character*, **115**(772):700–721, 1927.
- [KOV89] Claude Kipnis, Stefano Olla, and S. R. Srinivasa Varadhan. “Hydrodynamics and Large Deviation for Simple Exclusion Processes.” *Communications on Pure and Applied Mathematics*, **42**(2):115–137, 1989.
- [KP13] Da Kuang and Haesun Park. “Fast Rank-2 Nonnegative Matrix Factorization for Hierarchical Document Clustering.” In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 739–747, 2013.
- [KR02] Matt J Keeling and Pejman Rohani. “Estimating Spatial Coupling in Epidemiological Systems: a Mechanistic Approach.” *Ecology Letters*, **5**(1):20–29, 2002.
- [KSD17] Can Karakus, Yifan Sun, Suhas Diggavi, and Wotao Yin. “Straggler Mitigation in Distributed Optimization through Data Encoding.” In *Advances in Neural Information Processing Systems*, NIPS’17, p. 5440–5448, 2017.
- [KSK15] Matt Kusner, Yu Sun, Nicholas Kolkin, and Kilian Weinberger. “From Word Embeddings to Document Distances.” In *International conference on machine learning*, pp. 957–966, 2015.
- [KT81] Samuel Karlin and Howard M. Taylor. *A Second Course in Stochastic Processes*. Academic Press, Inc. [Harcourt Brace Jovanovich, Publishers], New York-London, 1981.
- [KTP19] Kalpesh Krishna, Gaurav Singh Tomar, Ankur P. Parikh, Nicolas Papernot, and Mohit Iyyer. “Thieves on Sesame Street! Model Extraction of BERT-based APIs.” *Clinical Orthopaedics and Related Research*, **abs/1910.12366**, 2019.

- [Kun86] Hiroshi Kunita. *Lectures on Stochastic Flows and Applications*, volume 78 of *Tata Institute of Fundamental Research Lectures on Mathematics and Physics*. Published for the Tata Institute of Fundamental Research, Bombay; by Springer-Verlag, Berlin, 1986.
- [KW18] Theodore Kolokolnikov and Juncheng Wei. “Pattern Formation in a Reaction-diffusion System with Space-dependent Feed Rate.” *SIAM Review*, **60**(3):626–645, 2018.
- [KWT18] Theodore Kolokolnikov, Michael Ward, Justin Tzou, and Juncheng Wei. “Stabilizing a Homoclinic Stripe.” *Philosophical Transactions of the Royal Society A. Mathematical, Physical and Engineering Sciences*, **376**(2135):20180110, 13, 2018.
- [KZR21] Nitin Kamra, Yizhou Zhang, Sirisha Rambhatla, Chuizheng Meng, and Yan Liu. “Pol-sird: Modeling Epidemic Spread under Intervention Policies.” *Journal of Healthcare Informatics Research*, **5**(3):231–248, 2021.
- [LHN22] Xia Li, Longxiu Huang, and Deanna Needell. “Distributed Randomized Kaczmarz for the Adversarial Workers.”, 2022.
- [Lig80] Thomas M. Liggett. “Interacting Markov Processes.” In *Biological growth and spread (Proceedings of a Conference Held at Heidelberg, 1979)*, volume 38 of *Lecture Notes in Biomathematics*, pp. 145–156. Springer, Berlin-New York, 1980.
- [Lig85] Thomas M. Liggett. *Interacting Particle Systems*, volume 276 of *Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences]*. Springer-Verlag, New York, 1985.
- [Lig10] Thomas M. Liggett. *Continuous Time Markov Processes*, volume 113 of *Graduate Studies in Mathematics*. American Mathematical Society, Providence, RI, 2010. An introduction.
- [LL10] Dennis Leventhal and Adrian S. Lewis. “Randomized Methods for Linear Constraints: Convergence Rates and Conditioning.” *Mathematics of Operations Research*, **35**(3):641–654, 2010.
- [Llo01] Alun L. Lloyd. “Realistic Distributions of Infectious Periods in Epidemic Models: Changing Patterns of Persistence and Dynamics.” *Theoretical Population Biology*, **60**(1):59 – 71, 2001.
- [LM11] Erik Lewis and George Mohler. “A Nonparametric EM Algorithm for Multiscale Hawkes Processes.” *Journal of Nonparametric Statistics*, **1**(1):1–20, 2011.
- [LMR19] Xin Liu, Anuj Mubayi, Dominik Reinhold, and Liu Zhu. “Approximation Methods for Analyzing Multiscale Stochastic Vector-borne Epidemic Models.” *Mathematical Biosciences*, **309**:42–65, 2019.

- [Los20] Office of Los Angeles Mayor Eric Garcetti. “Covid-19: Keeping Los Angeles safe.” <https://coronavirus.lacity.org/english/mayor-garcetti-issues-safer-home-emergency-order-stopping-non-essential-activities-outside>, 2020. Accessed: 2022-09-06.
- [LS99] Daniel D. Lee and H. Sebastian Seung. “Learning the Parts of Objects by Non-negative Matrix Factorization.” *Nature*, **401**:788–791, 1999.
- [LSP82] Leslie Lamport, Robert Shostak, and Marshall Pease. “The Byzantine Generals Problem.” *ACM Transactions on Programming Languages and Systems*, **4**(3):382–401, jul 1982.
- [LTH18] Zhiming Li, Zhidong Teng, Dujun Hong, and Xiaoping Shi. “Comparison of Three SIS Epidemic Models: Deterministic, Stochastic and Uncertain.” *Journal of Intelligent & Fuzzy Systems*, **35**:5785–5796, 2018.
- [LWL22] Xia Li, Chuntian Wang, Hao Li, and Andrea L. Bertozzi. “A Martingale Formulation for Stochastic Compartmental Susceptible-infected-recovered (SIR) Models to Analyze Finite Size Effects in COVID-19 Case Studies.” *Networks and Heterogeneous Media*, **17**(3):311–331, 2022.
- [M82] Michel Métivier. *Semimartingales*, volume 2 of *de Gruyter Studies in Mathematics*. Walter de Gruyter & Co., Berlin-New York, 1982. A course on stochastic processes.
- [MBX21] Mélodie Monod, Alexandra Blenkinsop, Xiaoyue Xi, Daniel Hebert, Sivan Bershan, Simon Tietze, Marc Baguelin, Valerie C. Bradley, Yu Chen, Helen Coupland, Sarah Filippi, Jonathan Ish-Horowicz, Martin McManus, Thomas Mellan, Axel Gandy, Michael Hutchinson, H. Juliette T. Unwin, Sabine L. van Elsland, Michaela A. C. Vollmer, Sebastian Weber, Harrison Zhu, Anne Bezancon, Neil M. Ferguson, Swapnil Mishra, Seth Flaxman, Samir Bhatt, and Oliver Ratmann. “Age Groups that Sustain Resurging COVID-19 Epidemics in the United States.” *Science*, **371**(6536), 2021.
- [MCC13] Tomas Mikolov, Kai Chen, Greg S. Corrado, and Jeffrey Dean. “Efficient Estimation of Word Representations in Vector Space.”, 2013.
- [MNR15] Anna Ma, Deanna Needell, and Aaditya Ramdas. “Convergence Properties of the Randomized Extended Gauss–Seidel and Kaczmarz Methods.” *SIAM Journal on Matrix Analysis and Applications*, **36**(4):1590–1604, 2015.
- [MP80] Michel Métivier and Jean Pellaumail. *Stochastic Integration*. Academic Press [Harcourt Brace Jovanovich, Publishers], New York-London-Toronto, Ont., 1980. Probability and Mathematical Statistics.
- [MSC13] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. “Distributed Representations of Words and Phrases and Their Compositionality.” In *Advances in neural information processing systems*, pp. 3111–3119, 2013.

- [MSH12] Brian Macdonald, Paulo Shakarian, Nicholas Howard, and Geoffrey Moores. “Spread-ers in the Network SIR Model: an Empirical Study.” *arXiv preprint arXiv:1208.4269*, 2012.
- [MSS20] George Mohler, Frederic Schoenberg, Martin B Short, and Daniel Sledge. “Analyzing the World-Wide Impact of Public Health Interventions on the Transmission Dynamics of COVID-19.” *arXiv preprint arXiv:2004.01714*, 2020.
- [MWT11] David Mimno, Hanna Wallach, Edmund Talley, Miriam Leenders, and Andrew McCallum. “Optimizing Semantic Coherence in Topic Models.” In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing*, pp. 262–272, 2011.
- [Nat01] Frank Natterer. *The Mathematics of Computerized Tomography*. SIAM, 2001.
- [Nee09] Deanna Needell. “Randomized Kaczmarz Solver for Noisy Linear Systems.” *BIT Numerical Mathematics*, **50**:395–403, 2009.
- [New] Official Site of the State of New Jersey. “Has the State Used Data to Make Decisions and Slow the Spread of COVID-19?: FAQ.” <https://covid19.nj.gov/faqs/nj-information/the-latest-data/how-is-the-state-using-data-to-make-decisions-and-slow-the-spread-of-covid-19>.
- [PHB18] Allison L. Pitt, Keith Humphreys, and Margaret L. Brandeau. “Modeling Health Benefits and Harms of Public Policy Responses to the US Opioid Epidemic.” *American Journal of Public Health*, **108**(10):1394–1400, 2018. PMID: 30138057.
- [PNI18] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. “Deep Contextualized Word Representations.” pp. 2227–2237, June 2018.
- [PNW19] Frédéric Paquin-Lefebvre, Wayne Nagata, and Michael J. Ward. “Pattern Formation and Oscillatory Dynamics in a Two-dimensional Coupled Bulk-surface Reaction-diffusion System.” *SIAM Journal on Applied Dynamical Systems*, **18**(3):1334–1390, 2019.
- [Pop99] Constantin Popa. “Characterization of the Solutions Set of Inconsistent Least-squares Problems by an Extended Kaczmarz Algorithm.” *Korean Journal of Computational and Applied Mathematics*, **6**(1):51–64, 1999.
- [PP16] Stefania Petra and Constantin Popa. “Single Projection Kaczmarz Extended Algorithms.” *Numerical Algorithms*, **73**(3):791–806, 2016.
- [Pro05] Philip E. Protter. *Stochastic Integration and Differential Equations*, volume 21 of *Stochastic Modelling and Applied Probability*. Springer-Verlag, Berlin, 2005. Second edition. Version 2.1, Corrected third printing.

- [PZ07] Szymon Peszat and Jerzy Zabczyk. *Stochastic Partial Differential Equations with Lévy Noise*, volume 113 of *Encyclopedia of Mathematics and its Applications*. Cambridge University Press, Cambridge, 2007. An Evolution Equation Approach.
- [Rao99] B. L. S. Prakasa Rao. *Statistical Inference for Diffusion Type Processes*, volume 8 of *Kendall's Library of Statistics*. Edward Arnold, London; Oxford University Press, New York, 1999.
- [RBH15] Michael Röder, Andreas Both, and Alexander Hinneburg. “Exploring the Space of Topic Coherence Measures.” In *Proceedings of the eighth ACM international conference on Web search and data mining*, pp. 399–408, 2015.
- [RHW15] Jonathan Le Roux, John R Hershey, and Felix Weninger. “Deep NMF for Speech Separation.” In *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 66–70. IEEE, 2015.
- [RMK18] Marian-Andrei Rizoiu, Swapnil Mishra, Quyu Kong, Mark Carman, and Lexing Xie. *SIR-Hawkes: Linking Epidemic Models and Hawkes Processes to Model Diffusions in Finite Populations*, pp. 419–428. International World Wide Web Conferences Steering Committee, Republic and Canton of Geneva, CHE, 2018.
- [RSW20] Joacim Rocklöv, Henrik Sjödin, and Annelies Wilder-Smith. “COVID-19 Outbreak on the Diamond Princess Cruise Ship: Estimating the Epidemic Potential and Effectiveness of Public Health Countermeasures.” *Journal of Travel Medicine*, **27**(3):taaa030, 2020.
- [RWL06] Shigui Ruan, Wendi Wang, and Simon A. Levin. “The Effect of Global Travel on the Spread of SARS.” *Mathematical Biosciences and Engineering*, **3**(1):205–218, 2006.
- [SB88] Gerard Salton and Christopher Buckley. “Term-weighting Approaches in Automatic Text Retrieval.” *Information processing & management*, **24**(5):513–523, 1988.
- [Sch20] Semantic Scholar. “CORD-19: COVID-19 Open Research Dataset.” <https://pages.semanticscholar.org/coronavirus-research>, 2020. Accessed: 2020-04-20.
- [SNT17] Xiaoxia Sun, Nasser M. Nasrabadi, and Trac D. Tran. “Supervised Multilayer Sparse Coding Networks for Image Classification.” *CoRR*, **abs/1701.08349**, 2017.
- [SP] The White House Office Of Science and Technology Policy. “Call to Action to the Tech Community on New Machine Readable COVID-19 Dataset.” <https://www.whitehouse.gov/briefings-statements/call-action-tech-community-new-machine-readable-covid-19-dataset/>. Accessed: 2020-03-20.
- [SS17] Shaheen Syed and Marco Spruit. “Full-text or Abstract? Examining Topic Coherence Scores Using Latent Dirichlet Allocation.” In *2017 IEEE International conference on data science and advanced analytics (DSAA)*, pp. 165–174. IEEE, 2017.

- [SV06] Daniel W. Stroock and S. R. Srinivasa Varadhan. *Multidimensional Diffusion Processes*. Classics in Mathematics. Springer-Verlag, Berlin, 2006. Reprint of the 1997 edition.
- [SV09] Thomas Strohmmer and Roman Vershynin. “A Randomized Kaczmarz Algorithm with Exponential Convergence.” *Journal of Fourier Analysis and Applications*, **15**(2):262, 2009.
- [Tar71] Robert Tarjan. “Depth-first search and linear graph algorithms.” In *12th Annual Symposium on Switching and Automata Theory (swat 1971)*, pp. 114–121, 1971.
- [TBZ16] George Trigeorgis, Konstantinos Bousmalis, Stefanos Zafeiriou, and Björn W Schuller. “A Deep Matrix Factorization Method for Learning Attribute Representations.” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **39**(3):417–429, 2016.
- [TCL18] Ding Tu, Ling Chen, Mingqi Lv, Hongyu Shi, and Gencai Chen. “Hierarchical Online NMF for Detecting and Tracking Topic Hierarchies in A Text Stream.” *Pattern Recognition*, **76**:203–214, 2018.
- [TKH20] Stefan Thurner, Peter Klimek, and Rudolf Hanel. “A network-based Explanation of Why Most COVID-19 Infection Curves are Linear.” *Proceedings of the National Academy of Sciences*, **117**(37):22684–22689, 2020.
- [TT94] Waiyuan Tan and Sichin Tang. “A General Markov Model of the HIV Epidemic in Populations Involving Both Sexual Contact and IV Drug Use.” *Mathematical and Computer Modelling*, **19**(10):83–132, 1994.
- [TV03] Bálint Tóth and Benedek Valkó. “Onsager Relations and Eulerian Hydrodynamic Limit for Systems with Several Conservation Laws.” *Journal of Statistical Physics*, **112**(3-4):497–521, 2003.
- [TW18] Wang Hung Tse and Michael J. Ward. “Asynchronous Instabilities of Crime Hotspots for a 1-D Reaction-diffusion Model of Urban Crime with Focused Police Patrol.” *SIAM Journal on Applied Dynamical Systems*, **17**(3):2018–2075, 2018.
- [TWW18] Justin C. Tzou, Michael Jeffrey Ward, and Jun cheng Wei. “Anomalous Scaling of Hopf Bifurcation Thresholds for the Stability of Localized Spot Patterns for Reaction-diffusion Systems in Two Dimensions.” *SIAM Journal on Applied Dynamical Systems*, **17**(1):982–1022, 2018.
- [Uch08] Masayuki Uchida. “Approximate Martingale Estimating Functions for Stochastic Differential Equations with Small Noises.” *Stochastic Processes and their Applications*, **118**(9):1706–1721, 2008.
- [Var95] S. R. Srinivasa Varadhan. “Entropy Methods in Hydrodynamic Scaling.” In *Proceedings of the International Congress of Mathematicians, Vol. 1, 2 (Zürich, 1994)*, pp. 196–208. Birkhäuser, Basel, 1995.

- [Var00] S. R. Srinivasa Varadhan. “Lectures on Hydrodynamic Scaling.” In *Hydrodynamic Limits and Related Topics (Toronto, ON, 1998)*, volume 27 of *Fields Institute Communications*, pp. 3–40. Amer. Math. Soc., Providence, RI, 2000.
- [War18] Michael Jeffrey Ward. “Spots, Traps, and Patches: Asymptotic Analysis of Localized Solutions to Some Linear and Nonlinear Diffusive Systems.” *Nonlinearity*, **31**(8):R189–R239, 2018.
- [Wika] Wikipedia. “California COVID-19 Timeline.” https://en.wikipedia.org/wiki/Timeline_of_the_COVID-19_pandemic_in_California. Accessed: 2022-06-20.
- [Wikb] Wikipedia. “Covid-19 pandemic in France.” https://en.wikipedia.org/wiki/COVID-19_pandemic_in_France#Timeline_of_measures. Accessed: 2022-09-06.
- [WTG11] Gerhard-Wilhelm Weber, Pakize Taylan, Zafer-Korcan Görgülü, Houssam Abdul-Rahman, and Acu Bahar. “Parameter Estimation in Stochastic Differential Equations.” In *Dynamics, Games and Science. II*, volume 2 of *Springer Proceedings in Mathematics*, pp. 703–733. Springer, Heidelberg, 2011.
- [WXG21] Xinran Wang, Yu Xiang, Jun Gao, and Jie Ding. “Information Laundering for Model Privacy.” In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.
- [XLG03] Wei Xu, Xin Liu, and Yihong Gong. “Document Clustering Based on Non-negative Matrix Factorization.” In *Proceedings of the 26th annual international ACM SIGIR conference on Research and development in informaion retrieval*, pp. 267–273, 2003.
- [Yan08] Ping Yan. “Distribution Theory, Stochastic Processes and Infectious Disease Modelling.” In *Mathematical Epidemiology*, volume 1945 of *Lecture Notes in Mathematics*, pp. 229–293. Springer, Berlin, 2008.
- [YB19] Zhixiong Yang and Waheed U Bajwa. “ByRDIE: Byzantine-resilient Distributed Coordinate Descent for Decentralized Learning.” *IEEE Transactions on Signal and Information Processing over Networks*, **5**(4):611–627, 2019.
- [YC19] Ping Yan and Gerardo Chowell. *Beyond the Initial Phase: Compartment Models for Disease Transmission*, pp. 135–182. Springer International Publishing, Cham, 2019.
- [YLB19] Baichuan Yuan, Hao Li, Andrea L. Bertozzi, P. Jeffrey Brantingham, and Mason A. Porter. “Multivariate Spatiotemporal Hawkes Processes and Network Reconstruction.” *SIAM Journal on Mathematics of Data Science*, **1**(2):356–382, 2019.