

UC San Diego

UC San Diego Electronic Theses and Dissertations

Title

Duality and Data Dependence in Boosting /

Permalink

<https://escholarship.org/uc/item/04p733x7>

Author

Telgarsky, Matus

Publication Date

2013

Peer reviewed|Thesis/dissertation

UNIVERSITY OF CALIFORNIA, SAN DIEGO

Duality and Data Dependence in Boosting

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

in

Computer Science

by

Matus Telgarsky

Committee in charge:

Professor Sanjoy Dasgupta, Chair
Professor Kamalika Chaudhuri
Professor Patrick Fitzsimmons
Professor Yoav Freund
Professor Philip Gill
Professor Robert Schapire

2013

Copyright

Matus Telgarsky, 2013

All rights reserved.

The Dissertation of Matus Telgarsky is approved and is acceptable in quality and form for publication on micro Im and electronically:

Chair

University of California, San Diego

2013

EPIGRAPH

In mathematics you don't understand things. You just get used to them.

John von Neumann

TABLE OF CONTENTS

Signature Page	iii
Epigraph	iv
Table of Contents	v
List of Figures	viii
Acknowledgements	ix
Vita	x
Abstract of the Dissertation	xi
Chapter 1 Introduction	1
1.1 Boosting	1
1.2 The Hypothesis Class \mathcal{H}	3
1.3 Primal Optimization Problem	4
1.4 Dual Optimization Problem	7
1.5 Algorithm	11
1.6 Other Approaches	14
1.7 Basic Structures	15
1.8 Analysis Overview	18
Appendix 1.A Banach Spaces and Linear Operators	20
Appendix 1.B Convexity Properties of \mathbb{R} and \mathbb{B}	25
Appendix 1.C Line Search Properties	34
1.C.1 Line Search Options 1 and 2	35
1.C.2 Line Search Option 3: the Wolfe Search	36
Appendix 1.D Existence of Hard Cores	41
Appendix 1.E Bibliographic Notes	42
I Just a Finite Sample	45
Chapter 2 Problem Structure	46
2.1 Overview	46
2.2 Strong Duality	48
2.3 Generalized Weak Learning Rate	51
2.4 The Hard Core	53
2.4.1 Weak Learnability	54
2.4.2 Attainability	55
2.4.3 General Setting	58
Appendix 2.A Generalizing the Weak Learning Rate	62
2.A.1 Choosing a Generalization to (b)	62
2.A.2 Proof of Theorem 2.3.6	65
Appendix 2.B Proof of 0-coercivity Characterization	70
Appendix 2.C Bibliographic Notes	70
Chapter 3 Optimization Guarantees	71

3.1	Overview	71
3.2	Weak Learnability	73
3.3	Attainability	75
3.4	General Setting	78
Appendix 3.A Miscellaneous Technical Material		82
3.A.1	The Logistic Loss is within $L_{f_{so}}$	82
3.A.2	$L_{f_{so}}$ implies \mathcal{L} has Lipschitz Gradients	83
3.A.3	Proof of Lemma 3.3.2	84
3.A.4	Splitting Distances along $A_0; A_+$	85
3.A.5	Proof of Theorem 3.4.6	85
Appendix 3.B Bibliographic Notes		87

II Statistical Behavior 89

Chapter 4 The Primal and the Dual		90
4.1	Strong Duality	90
Appendix 4.A Deferred Proofs		92
Appendix 4.B Bibliographic Remarks		98
Chapter 5 Finite Hypothesis Classes		99
5.1	Overview	100
5.2	An Impossibility Result	101
5.3	Hard Cores	102
5.4	Hard Cores and Convex Risk	103
5.5	Deviation Inequalities	104
Appendix 5.A Technical Preliminaries		107
Appendix 5.B Structure of L over $S_D(H; \cdot)$		110
Appendix 5.C Deferred Material from Section 5.2		115
Appendix 5.D Deferred Material from Section 5.3		116
5.D.1	Primal Hard Cores	116
5.D.2	Proof of Theorem 5.3.1	121
Appendix 5.E Deferred Material from Section 5.4		122
Appendix 5.F Deferred Material from Section 5.5		123
Appendix 5.G Bibliographic Notes		127
Chapter 6 Infinite Hypothesis Classes		129
6.1	Overview	129
6.2	Consistency Statement and Analysis Sketch	131
6.3	The Separable Case $L(\cdot) = 0$	132
6.3.1	The Quantity $\mathcal{L}(\cdot)$	133
6.3.2	Proof Sketch of Theorem 6.3.2	135
6.4	The Nonseparable Case $L(\cdot) > 0$	136
6.4.1	Curvature	138
6.4.2	Proof of Theorem 6.4.2	139
Appendix 6.A The Family of Dense Classes $\mathcal{F}_{ds}(\cdot)$		140
Appendix 6.B Loss Function Classes L_{lg} and L_{2d}		143
Appendix 6.C Duality Properties of		145
Appendix 6.D Reweighted Margin Deviations (with p Fixed)		150
Appendix 6.E Deferred Material from Section 6.3		152
6.E.1	Deviations of $\mathcal{L}(\cdot)$	152

6.E.2 Other Results	158
6.E.3 Optimization Guarantees	161
6.E.4 Statistical Guarantees	162
Appendix 6.F Deferred Material from Section 6.4	163
6.F.1 Proof of Proposition 6.4.1	163
6.F.2 Proof of Lemma 6.4.5	164
6.F.3 Optimization Guarantees	171
6.F.4 Statistical Guarantees	172
Appendix 6.G Proof of Consistency	174
Appendix 6.H Bibliographic Notes	175
 Bibliography	 178
 Index	 181

LIST OF FIGURES

Figure 1.1.	The exponential and logistic losses, along with their conjugates	9
Figure 1.2.	Stylized primal and dual problems.	10
Figure 1.3.	The boosting algorithm: coordinate descent applied to ℓ^b A. Note that Ad-aBoost (Freund and Schapire, 1997) is recovered by choosing the exponential loss $\ell = \ell_{\text{exp}}$, exact unconstrained step sizes (option 1), and exact coordinate selection ($\alpha = 1$).	13
Figure 1.4.	Bracketing and bisection search for step size satisfying Wolfe conditions.	39
Figure 1.5.	The Wolfe conditions.	40
Figure 2.1.	Geometric view of the primal and dual problem, under weak learnability.	55
Figure 2.2.	Geometric view of the primal and dual problem, under attainability.	57
Figure 2.3.	Geometric view of the primal and dual problem in the general case.	60
Figure 5.1.	Statistical difficulties with unconstrained minimization	101

ACKNOWLEDGEMENTS

Chapters 1, 2, and 3 contain material from the Journal of Machine Learning Research, volume 13, pages 561-606, 2012. The dissertation author was the sole author of this paper.

Chapters 1, 4, and 6 contain material from the Conference on Learning Theory, 2013. The dissertation author was the sole author of this paper.

Chapters 1 and 5 contains material being prepared for submission. The dissertation author is the sole author.

Much of this work was supported by the NSF under grants IIS-0713540 and IIS-0812598.

The author thanks the following people for valuable discussions throughout the project: Jake Abernethy, Akshay Balsubramani, Sanjoy Dasgupta, Daniel Hsu, Brian McFee, Aditya Menon, Indraneel Mukherjee, Alexander Rakhlin, Jeroen Rombouts, Robert Schapire, Karthik Sridharan, Manfred Warmuth, and numerous reviewers. The author wishes to give special mention to: his advisor Sanjoy Dasgupta for support and mathematical nutrition; Robert Schapire for insight, support, discussions, and for sharing open problems, early book drafts, and enthusiasm; Daniel Hsu for insight, support, discussions, for sparking the author's interest in boosting by initiating a reading project, and for setting the research path of this thesis by suggesting the author solve the problem in Chapter 3; his mother, Anna Telgarsky, for a work ethic.

VITA

- 2003 Diploma in Violin Performance, The Juilliard School
- 2007 B.S. in Computer Science and Discrete Math, Carnegie Mellon University
- 2013 Ph.D. in Computer Science, University of California, San Diego

ABSTRACT OF THE DISSERTATION

Duality and Data Dependence in Boosting

by

Matus Telgarsky

Doctor of Philosophy in Computer Science

University of California, San Diego, 2013

Professor Sanjoy Dasgupta, Chair

Boosting algorithms produce accurate predictors for complex phenomena by welding together collections of simple predictors. In the classical method AdaBoost, as well as its immediate variants, the welding points are determined by convex optimization; unlike typical applications of convex optimization in machine learning, however, the AdaBoost scheme eschews the usual regularization and constraints used to control numerical and statistical properties.

On the other hand, the data and simple predictors impose rigid structure on the behavior of AdaBoost variants, and moreover convex duality provides a lens to resolve this rigidity. This structure is fundamental to the properties of these methods, and in particular leads to numerical and statistical convergence rates.

Chapter 1

Introduction

1.1 Boosting

Let a room with heaps of books be given; since the finiteness of human life precludes reading all of them, what is the best way to identify a high quality subset? It is unclear how to define a good book. On the other hand, it is easy to identify many characteristics which are correlated with the quality of a book: for example, it should have spaceships, and it should contain information, but reading it should not lead to sadness for no reason.

This is precisely the idea behind boosting: rather than carefully reasoning about an appropriate model for a prediction problem, simply identify a set of characteristics correlated with the intended values, and leave it to an algorithm to combine these scraps of information into an accurate model.

The key question here is of course what such a scheme can hope to achieve. An early concrete question along these lines was provided by Kearns and Valiant (1989), where the task was to identify individual books as good or bad (for an abstract notion of book). In more detail, suppose there exists a positive constant $\epsilon > 0$ so that for any distribution over the set of books, some characteristic has correlation at least ϵ with the quality of the books according to the selected distribution; is it possible to construct a circuit which uses correlation values for a set of books to separate the good ones from the bad?

The definition of ϵ suggests an iterative scheme, since an initial circuit gives rise to a distribution emphasizing those books where the circuit errs, which in turn by the stated assumption leads to a characteristic that is ϵ -correlated with book quality. If this characteristic can be baked into the initial circuit in such a way that accuracy improves, then perhaps the goal

may be achieved by recursing.

The landmark result here, provided by Schapire (1990), is that this can indeed be done. The circuit provided by Schapire's boosting method was a ternary tree, where leaves measured correlation values for a provided book, and internal nodes took a majority vote over their three children.

The next major result, by Freund (1995), was that not only did a single majority over a set of correlation values suffice, but moreover only $O(\ln(1/\epsilon))$ correlates were used, where $\epsilon > 0$ is the target accuracy (e.g., choosing $\epsilon = 1/(\text{number of books})$ suffices to classify the books perfectly). This scheme unfortunately required ϵ as an input parameter; fortunately, the two authors of the two preceding methods teamed up to produce an Adaptive Boosting method, AdaBoost, which does not require knowledge of ϵ , and preserved both the majority structure and the $O(\ln(1/\epsilon))$ size bound (Freund and Schapire, 1997).

This thesis studies AdaBoost and its immediate variants; the two main themes are as follows.

Duality. The coefficient values in the weighted majority function produced by AdaBoost are the solution to a convex optimization problem; the dual optimization is a maximum entropy problem over a space of weightings on individual books. The primal player seeks to identify good books; the dual player seeks a distribution which decorrelates all the characteristics from the quality of the books (thus forcing $\epsilon = 0$). How much does duality dictate these basic behavior of the methods?

Data-dependence. The scalar ϵ is a function of data; specifically, it is a function of the set of books, and the set of chosen characteristics. The condition $\epsilon > 0$ implies the size bound $O(\ln(1/\epsilon))$ stated above. What other quantities and structures dictate the behavior of the algorithm, for example when $\epsilon = 0$? Can these objects be derived from duality structure?

The purpose of this thesis is to study the duality structure and data-dependent structure of the optimization problems corresponding to AdaBoost variants, and thereafter to apply this lens to the algorithms themselves, deriving both numerical and statistical guarantees.

This introductory chapter is primarily devoted to establishing a language with which to discuss AdaBoost variants (Section 1.2 through Section 1.7, with additional technical development

in Appendix 1.A through Appendix 1.D). A summary of the main results of the thesis, as well as a more detailed overview of the remaining chapters, appears in Section 1.8.

1.2 The Hypothesis Class \mathcal{H}

Let X denote an abstract input space (abstract room of abstract books). The aforementioned characteristics or correlates are modeled as a set of functions \mathcal{H} , where each $h \in \mathcal{H}$ has domain X and bounded range $[-1; +1]$, meaning $h : X \rightarrow [-1; +1]$. At times, it will be necessary to constrain every $h \in \mathcal{H}$ to be binary, meaning $h : X \rightarrow \{-1; +1\}$.

Throughout this thesis, μ will denote a weighting over \mathcal{H} , which leads to a map of the form

$$x \mapsto \int_{h \in \mathcal{H}} h(x) \mu(h) dh$$

By the boundedness of $h \in \mathcal{H}$, this map also has bounded range, specifically $[-k_1; +k_1]$.

In order to more conveniently talk about this map, some further notation is beneficial. Firstly, let μ denote counting measure over \mathcal{H} , whereby the space $\ell^1(\mathcal{H})$ is simply all weightings of \mathcal{H} which are absolutely summable; that is to say, $\mu \in \ell^1(\mathcal{H})$ implies

$$\sum_{h \in \mathcal{H}} |\mu(h)| < 1 \quad \text{and} \quad \sup_{x \in X} \int_{h \in \mathcal{H}} h(x) \mu(h) dh \leq \sup_{x \in X} \int_{h \in \mathcal{H}} |h(x)| \mu(h) dh < 1;$$

whereby the second property provides that the map $x \mapsto \int_{h \in \mathcal{H}} h(x) \mu(h) dh$ is well-defined whenever $\mu \in \ell^1(\mathcal{H})$.

The algorithms considered here will only work with choices $\mu \in \ell^1(\mathcal{H})$ which have finitely many nonzero coordinates. The reason for working with the larger space $\ell^1(\mathcal{H})$ is that it is a complete normed vector space, which is essential to many continuity properties, which moreover are essential to many duality arguments. To this end, define the simpler notation $\mathcal{L}^1(\mathcal{H}) := \ell^1(\mathcal{H})$.

As an additional convenience, let \mathcal{H} be a linear operator over $\mathcal{L}^1(\mathcal{H})$ defined as

$$(\mathcal{H})_x := (\mathcal{H})(x) := \int_{h \in \mathcal{H}} h(x) \mu(h) dh = \int_{h \in \mathcal{H}} h(x) \mu(h) dh;$$

where the preceding development established this sum is well-defined whenever $\mu \in \mathcal{L}^1(\mathcal{H})$. Note that this expression could have been defined directly over $\mathcal{L}^1(\mathcal{H})$ via some theory of integration over Banach spaces, for instance the Bochner integral, but pulling in that machinery here is unnecessary. For

further properties of H , please see Section 1.A.

In some of the analysis, namely in Chapters 2 and 3, both H and X will be finite, whereby H is simply a matrix. But either way, H is a linear operator, and H is well-defined for \mathcal{X} .

1.3 Primal Optimization Problem

Identifying good and bad books is a classification problem; the boosting methods here will choose a set of coefficients α , and rather than outputting $(H^\alpha)_x$ for a book x , this value will be thresholded, meaning the final output is

$$x \mapsto \mathbb{1}[(H^\alpha)_x \geq 0];$$

which is 1 when $(H^\alpha)_x \geq 0$, and -1 otherwise.

It is desirable for this map to output the correct value for every $x \in \mathcal{X}$. In order to discuss correctness, the relevant space is not \mathcal{X} but $\mathcal{X} \times \mathcal{Y}$, the set $\mathcal{X} \times \mathcal{Y}$ distinguishing good and bad books. Correspondingly, let \mathbb{P} denote a probability measure over $\mathcal{X} \times \mathcal{Y}$; at times, an i.i.d. sample $(x_i; y_i)_{i=1}^m$ from \mathbb{P} of size m will be given, and the corresponding empirical measure will be \mathbb{P}_m ; many results hold for either measure, and thus the measure variable will often be employed. The σ -algebra is always the Borel σ -algebra; this is largely irrelevant, with the exception of material discussing how to approximate arbitrary (Borel!) measurable functions. Additionally, let \mathbb{P}^X denote the marginal distribution over \mathcal{X} , with corresponding conditional probability $\mathbb{P}(y = +1 | x)$ (analogously define \mathbb{P}^Y , $\mathbb{P}^{X|Y}$, etc.).

With this measure notation in place, the goal of the algorithm is to have as few errors as possible over \mathbb{P} ; that is, define the risk R of a measurable map $f : \mathcal{X} \rightarrow \mathbb{R}$ as

$$R(f) := \int_{\mathcal{X} \times \mathcal{Y}} \mathbb{1}[(f(x) < 0 \wedge y = +1) \vee (f(x) \geq 0 \wedge y = -1)] d\mathbb{P}(x, y)$$

(with \mathbb{P} defined analogously with b replacing \mathbb{P}), and thus the goal is to approximately solve the optimization problem

$$\inf_{f \in \mathcal{F}} R(f) \tag{1.3.1}$$

There are two difficulties here. The first is that the algorithm will not have access to \mathbb{P} , but rather just an i.i.d. draw $(x_i; y_i)_{i=1}^m$ from \mathbb{P} . To overcome this, the empirical measure \mathbb{P}_m corresponding

to this draw is used instead, meaning R is minimized rather than R . Of course, replacing R with b is a dangerous proposition, and occupies much of this thesis. The main issue, as the problem is currently stated, is to constrain the complexity of H , or rather its convex hull; namely, it is essential to establish the convergence of these scores R to those of b uniformly over the class of functions $x \in \mathcal{X}$ (H).

The second issue is that this optimization problem is not tractable in the general case; for example, if the in mal error rate is 1% (because, say, H is too small, or simply the conditional probability $\Pr[Y = +1 | X]$ is not 0 or 1 almost everywhere), it is NP-Hard to solve this problem with fewer than 49% errors (Guruswami and Raghavendra, 2006). This problem will be overcome by replacing the in mand with a convex surrogate, developed as follows.

First, for added convenience, define another linear operator A , whereby A is now a function of both $x \in \mathcal{X}$ and $y \in \mathcal{Y}$:

$$(A)(x; y) := (A)_{x,y} := y(H)_x = \sum_{h \in \mathcal{H}} y h(x) (h);$$

note that this quantity is well-defined for $\mathcal{X} \times \mathcal{Y}$ exactly for the same reasons that H is well-defined. Now let ψ be a nondecreasing convex function with $\psi(0) = 0$ and in particular $\psi'(0) > 0$; for every $(x; y) \in \mathcal{X} \times \mathcal{Y}$,

$$((H)_x < 0 \wedge y = +1) \vee ((H)_x \geq 0 \wedge y = -1) \implies (A)_{x,y} \geq 0;$$

whereby

$$1 [((H)_x < 0 \wedge y = +1) \vee ((H)_x \geq 0 \wedge y = -1)] \geq \frac{\psi((A)_{x,y})}{\psi'(0)}; \quad (1.3.2)$$

This in turn motivates approximate minimization of the convex surrogate

$$\inf_{x,y} \sum_{d} \psi((A)_{x,y}) d(x; y); \quad (1.3.3)$$

where d is either a or b .

Of course, the inequality in eq. (1.3.2) is one-sided; ideally, approximate solutions to eq. (1.3.3) will also be approximate optima to the original (potentially nonconvex) problem in

eq. (1.3.1) Establishing a stronger relationship of this form is indeed possible, although it depends on a number of properties of ϕ and H (Zhang, 2004, Bartlett et al., 2006), and will be a key aspect of Chapter 6. For a hint as to why such a result is possible, note that if ϕ is differentiable and increasing at the origin (which is in fact a necessary condition (Bartlett et al., 2006, Theorem 2)), then it is nearly linear in some tiny open ball around the origin; as such, minimizing it can be expected to push incorrect predictions in the right direction, without hurting correct predictions too much.

The approach of AdaBoost (followed in this thesis) is to minimize the convex surrogate in eq. (1.3.3). The statistical issues with this optimization are more severe than with the non-convex formulation in eq. (1.3.1), since the function $\phi: \mathbb{R} \rightarrow \mathbb{R}$, unlike R , is sensitive to the norm $\|\cdot\|_1$; in particular, in order to achieve good statistical properties, it will sometimes be desirable to ensure that $\|\cdot\|_1$ is small.

Before adjourning this section, a few more remarks on these convex losses are appropriate, which leads to additional notation. First, the main losses considered throughout are

$\ell_{\text{exp}}(x) := \exp(x)$	exponential loss
$\ell_{\text{log}}(x) := \ln(1 + \exp(x))$	logistic loss;
$\ell_{\text{hinge}}(x) := \max\{0, x + 1\}$	hinge loss
$\ell_{\text{russ}}(x) := 0.5(x + 1)^2 \mathbb{1}_{[-1 < x < 0]} + (x + 0.5) \mathbb{1}_{[0 \leq x]}$	smoothed hinge loss

The exponential loss ℓ_{exp} was the original choice of AdaBoost. The logistic loss is popular in practice but less well understood theoretically, and therefore will receive preference in this thesis (in particular, some results in Chapter 6 hold for Lipschitz losses like ℓ_{log} , but not for the original choice ℓ_{exp}). The hinge loss ℓ_{hinge} is popular in machine learning, and an interesting counterpart to the preceding two because it is neither strictly convex nor differentiable; however, the only chapter to allow the hinge loss throughout is Chapter 5. Lastly, the smoothed hinge loss ℓ_{russ} is differentiable, but still not strictly convex, which will allow it to be used in a few places in Chapter 6 where ℓ_{hinge} is disallowed; this interesting loss is a rescaling of one implicit in a proof by Impagliazzo (1995, Proof of Lemma 1), and a translated version was also discussed by Zhang (2004, Section 3.6).

Finally, a few more notational conveniences. For any $f \in L^1(\mathcal{X})$ and $f \in L^1(\mathcal{Y})$, define the

shorthands $L(f) := \int_{\mathcal{X} \times \mathcal{Y}} (f(x; y)) d(x; y) = \int_{\mathcal{X} \times \mathcal{Y}} (f) d$ and $\mathbb{L}(f) := \int_{\mathcal{X} \times \mathcal{Y}} (f(x; y)) db(x; y) = \int_{\mathcal{X} \times \mathcal{Y}} (f) db$, where these definitions also demonstrate the practice of dropping integration variables. As such, a succinct way to write the central optimization problem from eq. (1.3.3), now specialized for b , is

$$\inf_{\mathcal{L}} \mathbb{L}(A) : \mathbb{L} \in \mathcal{L} \quad (1.3.4)$$

Whenever \mathcal{L} or \mathbb{L} are passed a set \mathcal{S} , let $L(\mathcal{S})$ and $\mathbb{L}(\mathcal{S})$ denote minimization over that set, meaning for instance

$$L(A) := \inf_{\mathcal{L}} \int_{\mathcal{X} \times \mathcal{Y}} L(A) : \mathcal{L} \in \mathcal{L} \quad \text{and} \quad \mathbb{L}(A) := \inf_{\mathbb{L}} \int_{\mathcal{X} \times \mathcal{Y}} \mathbb{L}(A) : \mathbb{L} \in \mathbb{L};$$

with corresponding shorthands $R(\mathcal{H})$ and $\mathbb{R}(\mathcal{H})$ appearing frequently. Lastly, let brackets provide shorthand definitions of sets, meaning for instance

$$[A \in \mathcal{L}] = \{ \mathcal{L} : \int_{\mathcal{X} \times \mathcal{Y}} \mathcal{L}(A) \leq \epsilon \} = \{ \mathcal{L} : \int_{\mathcal{X} \times \mathcal{Y}} \mathcal{L}(A) \leq \epsilon \}$$

1.4 Dual Optimization Problem

The primal problem, presented in eq. (1.3.3) and eq. (1.3.4), is concerned with coefficients over the hypothesis space \mathcal{H} . On the other hand, the dual is concerned with weightings over the instance space. To start, note the following weak duality result, which is proved in Section 1.B.

Proposition 1.4.1. Let finite measure μ , hypothesis class \mathcal{H} , and convex loss $\ell : \mathcal{R} \rightarrow \mathcal{R}_+$ with $\lim_{z \rightarrow 1} \ell(z) = 0$ be given. Then

$$\inf_{\mathcal{L}} \int_{\mathcal{X} \times \mathcal{Y}} \mathcal{L}(A) d \leq \sup_{\mathcal{P}} \int_{\mathcal{X} \times \mathcal{Y}} \ell(p) d : \mathcal{P} \in \mathcal{L}^1(\mu); \mathcal{P} \geq 0 \text{ -a.e.}; \int_{\mathcal{X} \times \mathcal{Y}} \mathcal{P}(A) d = 0;$$

where ℓ° denotes the Fenchel conjugate of ℓ , discussed below.

It would be nice to replace the inequality with an equality; this in turn is provided by strong duality statements, presented in Chapter 2 and Chapter 4, which will additionally provide optimality guarantees. Even so, this weak duality statement already identifies a few key properties of the dual.

First note that the Fenchel conjugate of a function $\psi : \mathbb{R} \rightarrow \mathbb{R}$ is defined as

$$\psi^*(\lambda) = \sup_{x \in \mathbb{R}} \lambda x - \psi(x).$$

The conjugation operation is a cornerstone of convex analysis, and studying the conjugacy properties of ψ and L is essential to the development here. As discussed below, the choices of ψ most relevant here have a conjugate which resembles entropy. General information on Fenchel conjugates can be found in any treatment of convex analysis, for instance in the finite dimensional setting by Rockafellar (1970), and Hiriart-Urruty and Lemaréchal (2001) and Borwein and Lewis (2000), and in the infinite dimensional setting by Zalinescu (2002) and Borwein and Zhu (2005).

When μ is a discrete measure over a finite subset of $\mathbb{X} = \{1, \dots, n\}$ (e.g., when μ is δ_b), strong duality follows with a fair bit of generality, as is proved in Chapter 2. On the other hand, in the general case of an arbitrary (Borel) measure over $\mathbb{X} = \{1, \dots, n\}$, strong duality is only proved (in Chapter 4) when ψ is additionally Lipschitz, which rules out the exponential loss used by AdaBoost. There are indeed a number of steps that break down for the exponential loss, and in particular the dual should be over a larger space than $L^1(\mu)$, though this thesis only provides hints to the correct choice; please see Chapter 4 for further discussion. Note that the restriction of ψ Lipschitz was reasonable since this allows for analysis of the logistic loss, which as discussed previously has been somewhat neglected despite its empirical success.

In the case of the conjugacy theory of L (upon which the duality builds), it is the combination of H finite and μ infinite which leads to problems; if either μ is finite (cf. Proposition 1.B.2) or H is finite (cf. lemma 5.A.2), then L is much more well behaved.

This dual captures the duality structure stated in Section 1.1; namely, the dual space $L^1(\mu)$ is over reweightings of the instance space the instance space $\mathbb{X} = \{1, \dots, n\}$, and moreover the dual feasible set is all reweightings which decorrelate all primal weightings, written symbolically here as $\sum_{i \in \mathbb{X}} \lambda_i A_i = 0$ for all $\lambda \in \mathbb{R}^{\mathbb{X}}$.

The machine learning literature follows both the convention here of having losses nonde-

¹You should call it entropy, [...] no one really knows what entropy really is" | John von Neumann.

Figure 1.1. The exponential loss $\ell_{\text{exp}}(z) = \exp(z)$, the logistic loss $\ell_{\text{log}}(z) = \ln(1 + \exp(z))$, and their conjugates.

creasing (see also (Boucheron et al., 2005b)), and the mirrored convention of nonincreasing losses. The choice here was made so that gradients and gradient-like functions of primal quantities will directly lead to nonnegative weightings over the source distribution; in this way, the gradients and dual domain are more readily interpreted in the context presented in Section 1.1.

As mentioned above, a few more remarks on the structure of ℓ are pertinent. First, the exponential loss ℓ_{exp} is conjugate to the Boltzmann-Shannon entropy (see Borwein and Lewis, 2000, closing commentary, Section 3.3), defined as

$$\ell_{\text{exp}}(\lambda) = \begin{cases} \ln(\lambda) & \text{when } \lambda > 0; \\ 0 & \text{when } \lambda = 0; \\ 1 & \text{otherwise;} \end{cases}$$

whereas the logistic loss ℓ_{log} is conjugate to the Fermi-Dirac entropy (see Borwein and Lewis, 2000, closing commentary, Section 3.3), defined as

$$\ell_{\text{log}}(\lambda) = \begin{cases} (1-\lambda)\ln(1-\lambda) + \lambda\ln(\lambda) & \text{when } \lambda \in (0; 1); \\ 0 & \text{when } \lambda \in \{0; 1\}; \\ 1 & \text{otherwise;} \end{cases}$$

For a plot of these losses and their conjugates, please see Figure 1.1.

(a) Stylized primal problem.

(b) Stylized dual problem.

Figure 1.2. Stylized primal and dual problems; see for instance the weak duality relation in Proposition 1.4.1. On the primal problem, the arrow signifies that the objective value may fail to ever curve back upward; correspondingly, not only may minimizers fail to exist, but moreover it may be necessary for a sequence of iterates approaching the optimal value to have unbounded norms. On the other hand, the dual problem possesses curvature (cf. Proposition 1.4.2). The green area signifies that the dual is a constrained optimization problem: in particular, it is maximizing this curved shape, subject to the constraint that all weak predictors are decorrelated from the target prediction value $f = 1; +1g$.

More generally, while the primary losses considered here do not have minima (they tend to zero), their conjugates have a bowl-like structure, attaining a minimum and increasing to either side of it.

Proposition 1.4.2. Suppose $\phi : \mathbb{R} \rightarrow \mathbb{R}_+$ is differentiable, strictly convex, and $\lim_{z \rightarrow \pm\infty} \phi(z) = 0$.

Then

$$\phi^*(\lambda) = \begin{cases} 1 & \text{when } \lambda < 0; \\ 0 & \text{when } \lambda = 0; \\ (\phi^*(0); 0) & \text{when } \lambda \in (0; \phi'(0)); \\ \phi^*(0) & \text{when } \lambda = \phi'(0); \\ (\phi^*(0); 1] & \text{when } \lambda > \phi'(0); \end{cases}$$

Correspondingly, note the stylized versions of the primal and dual optimization problems as depicted in Figure 1.2.

1.5 Algorithm

What is the best way to minimize the primal optimization problem in eq. (1.3.4)? Ideally, the method somehow fits with the computational model suggested in Section 1.1, namely that a good correlation variable $h \in H$ can be found, given a reweighting of the provided examples $f(x_i; y_i) g_{i=1}^m$.

The choice followed by AdaBoost, and continued here, is to use coordinate descent: each iteration updates a coefficient vector $\beta \in \mathbb{R}^H$ by computing a gradient, and then updating the coordinate of h (identified by some $h \in H$) corresponding to a large entry in the gradient.

Concretely, coordinate descent is applied to the function $\mathcal{L}(\beta; A)$; as is developed in Proposition 1.B.2, the gradient of $\mathcal{L}(\beta; A)$ satisfies the familiar relation $\nabla_{\beta} \mathcal{L}(\beta; A) = A^T \mathcal{L}'(A \beta)$, but more importantly, it satisfies the pairing

$$\nabla_{\beta} \mathcal{L}(\beta; A) \cdot e_h = \frac{1}{m} \sum_{i=1}^m y_i \cdot \mathcal{L}'(A \beta)_{x_i, y_i} h(x_i); \quad (1.5.1)$$

where $e_h \in \mathbb{R}^H$ is the indicator vector for a single $h \in H$. Notice in particular that this precisely matches the computational model of interest here: this inner product is simply the correlations between h and the labels $f y_i g_{i=1}^m$, reweighted by $\nabla_{\beta} \mathcal{L}(A \beta)$. For each $h \in H$, supposing h has finite support (which it will within the algorithm), this expression is expressly computable (given, of course, some computable approximation of the reals and \mathcal{L}'). As a technical note, the gradient used here is the Gateaux derivative, which can handle an infinite coordinate space H , but crucially does not disagree with the usual notion when everything is finite; again, these details are presented together with basic convexity properties of \mathcal{L} in Proposition 1.B.2.

Coordinate descent methods typically choose the best coordinate. Since H may be infinite, a best coordinate may not exist; but even if it did, it is useful for computational reasons to relax this requirement. As such, the algorithm here takes in a parameter $\eta \in (0, 1)$, and each coordinate update chooses some $h_t \in H$ so that

$$\nabla_{\beta} \mathcal{L}(A \beta_t) \cdot A e_{h_t, j} = \sup_{h \in H} \nabla_{\beta} \mathcal{L}(A \beta_t) \cdot A e_{h, j} = \eta \nabla_{\beta} \mathcal{L}(A \beta_t) \cdot A k_t;$$

where the middle expression matches embeds eq. (1.5.1) from above (and thus fits the computational model), and the correspondence with $\eta \nabla_{\beta} \mathcal{L}(A \beta_t) \cdot A k_t$ in the last expression is provided by

Lemma 1.A.5, which is again consistent with the finite case. Note that the absolute value means that a descent step will move either in direction $+h_t$ or $-h_t$; many treatments of AdaBoost omit this choice, but instead stipulate that H is closed under negation, meaning $\|h\|_2 = \|-h\|_2$.

The full algorithm appears in Figure 1.3. It contains two choices which have not yet been pointed out.

1. There are two possible stopping criteria. One is based on the duality gap; unfortunately, as discussed Chapter 2, this choice is not computationally feasible, since it requires determining the kernel of A^T , whereas the computational model stressed here does not even require A and H to be ever be stored in memory. The second choice is to simply stop after m iterations, where $m \in (0; 1)$ is provided in advance. A popular choice in practice, which is not studied here, is to use cross-validation.
2. Three line search choices are allowed. When the unconstrained step exists, and η is binary and ℓ is the exponential loss, then this step recovers the original choice used in AdaBoost. The unconstrained step is the most difficult to analyze, since it may elect to take extremely large steps, potentially damaging the aforementioned suggestion that solutions with small norm are preferable. This difficulty will be acute in Chapter 6.

The second choice, or rather range of choices, is based upon computing a quadratic upper bound to the line search problem, and choosing its minimizer. In some cases, this step also has a simple expression, and does not lead to the difficulties of the unconstrained step.

Lastly, the Wolfe search is somewhere between these two; in a sense, it attempts to minimize the univariate line search problem, but subject to a constraint that disallows steps from being too large. The implementation of this method is similar to binary search, and does not need to know any properties of the functions it is minimizing (in the convex case). It takes two parameters, which are set to $1/2$ and $1/3$ in this thesis for simplicity.

Details on all three line searches may be found in Section 1.C. Each provides nearly the same optimization guarantees, however, as discussed above, they vary in terms of the norms of coefficient vectors they lead to.

Routine. Boost .

Input. Loss ℓ and empirical measure \mathbb{P}_n (granting m, ℓ, r, \mathbb{P}), hypothesis class H (granting $A, A^>$), coordinate search parameter $\alpha \in (0; 1]$.

Output. Coefficient vectors $f_i, g_{i=1}^t$.

1. Set $\alpha_0 := 0$.

2. For $t = 1; 2; \dots$:

(a) Decide whether to exit loop over t as follows.

Option 1. Exit if $\alpha_t \geq m^\alpha$, for some provided $\alpha \in (0; 1)$.

Option 2. Exit if duality gap discussed in Remark 2.2.5 is upper bounded by some provided $\epsilon > 0$.

(b) Choose approximate best coordinate (weak learner) $h_t \in H$ satisfying

$$\mathbb{P}_n \ell(A_{t-1})^> A_{e_{h_t}} \leq \inf_{h \in H} \mathbb{P}_n \ell(A_{t-1})^> A_{e_h} = \alpha_t \mathbb{P}_n \ell(A_{t-1})^> A_{k_t} :$$

(c) Set descent direction $v_t \in \mathbb{R}^d$ where

$$\alpha_t \mathbb{P}_n \ell(A_{t-1})^> A_{k_t} \leq \mathbb{P}_n \ell(A_{t-1})^> A_{v_t} \leq \alpha_t \mathbb{P}_n \ell(A_{t-1})^> A_{k_t} :$$

(d) Set $\alpha_t := \arg \min_{\alpha > 0} \mathbb{P}_n \ell(A_{t-1} + \alpha v_t)$, and choose a step α_t as follows.

Option 1. If $\alpha_t < 1$, set $\alpha_t = \alpha_t$ (i.e., make an optimal unconstrained step).

Option 2. If ℓ has Lipschitz gradients with parameter B_t as discussed in Section 1.C, choose any $\alpha_t \in \mathbb{R}^+$ such that $\mathbb{P}_n \ell(A_{t-1})^> A_{v_t} = B_t \alpha_t$.

Option 3. Choose any α_t satisfying the Wolfe conditions (please see Section 1.C.2.)

(e) Update $\alpha_t := \alpha_{t-1} + \alpha_t v_t$.

3. Return $f_i, g_{i=1}^t$.

Figure 1.3. The boosting algorithm: coordinate descent applied to ℓ on A . Note that AdaBoost (Freund and Schapire, 1997) is recovered by choosing the exponential loss $\ell = \exp$, exact unconstrained step sizes (option 1), and exact coordinate selection ($\alpha = 1$).

1.6 Other Approaches

Before descending into a labyrinth of analysis of this algorithm and optimization problem, it is a good time to consider other approaches to the same problem.

It should be stressed that tools from convex optimization typically work with objective functions which are strongly convex (lower bounded by a quadratic), or constrained. It is also common to assume the existence of minimizers, which is characterized in Chapter 2 as imposing quite a bit of curvature structure on the objective function. On the other hand, the original AdaBoost analysis of Freund and Schapire (1997), under the assumption $\gamma > 0$, is characterized in Chapter 2 as being, in a strong sense, diametrically opposed to the typical setting of convex optimization. These two distinct settings, and their interplay, is a recurring element in this thesis.

Back to task, here are some other ways to solve this problem.

More coordinates. Leave the objective function alone, but take more than one coordinate at a time. Without any further assumptions or algorithmic modifications, it could be the case that all the chosen coordinates are similar, and therefore not much is gained over choosing a single coordinate (Long and Servedio, 2011). On the other hand, possibly with additional assumptions or an explicit choice of H as possessing some sort of independence between elements, there could be a way to still make good progress (Dasgupta and Long, 2003).

Subgradient Descent. Leave the objective function alone, suppose ℓ is Lipschitz, assume H is finite, and simply execute subgradient descent. The analysis of subgradient descent typically assumes a minimizer, but an inspection of the standard proof reveals that no optimality properties of this minimizer are invoked, and the proof may be simply restated with respect to an arbitrary reference point. In particular, reworking the clean presentation by Vandenberghe (2013), it follows, for any $C > 0$, that

$$\ell(A_{\text{best-of-}t}) \leq \inf_{k \leq k_2} \ell(A_k) + \frac{CL}{t};$$

where $A_{\text{best-of-}t}$ is the element encountered with lowest value ℓ in the first t iterations, and L is a Lipschitz parameter. This rate is slower than those presented in Chapter 3, and of course the computational model has been departed (the operator A , now simply a matrix, needs to be explicitly computed), however the method is simple and the comparison to a

solution with small norm is statistically pleasing.

Naive Regularization. Another choice is to choose some regularization weight $\lambda > 0$, and replace eq. (1.3.4) with

$$\min_{A} \frac{1}{n} \mathbb{E}(\ell(A)) + \lambda \|A\|_1; \quad (1.3.5)$$

where the use of "min" is justified, meaning a minimizer always exists. From here, a variety of techniques may be employed to solve this problem; but of course, there is a question of how to choose the parameter λ .

Statistical properties of penalized boosting estimators were studied by Blanchard et al. (2003), Lugosi and Vayatis (2004). Since it is not entirely clear how to implement these algorithms, one simple choice achieving a similar effect is to explicitly curtail the step size in Boost, which was shown by Zhang and Yu (2005) to lead to good numerical and statistical rates. While this thesis does not close the question of how eq. (1.3.4) should best be penalized, it is interesting to point out that there are cases, as studied in Chapter 6, where excellent rates are achieved despite producing solutions with large norms.

Different Objective. Another choice is to simply rewrite the objective function entirely; for instance, Ratsch and Warmuth (2005), Shalev-Shwartz and Singer (2008) add both regularization and constraints. Some basic guarantees of this scheme have been provided, but neither this thesis, nor other existing literature, comprehensively identify the respective benefits or disadvantages of the scheme considered here versus these constrained schemes.

1.7 Basic Structures

To give a taste of future material and to clarify a few ambiguous claims, this section develops a few basic quantities relevant throughout the analysis.

The quantity ϵ , discussed in Section 1.1, is typically called the weak learning rate and the condition $\epsilon > 0$ is called the weak learning assumption. In the original discussion, it was not clear whether this value is measured over n or b , thus consider the following definition.

Definition 1.7.1 (Schapire and Freund (2012, Chapter 2)) A class H is weakly PAC-learnable with rate ϵ if for any measure μ over $X \times Y$, there exists $h \in H$ with $\int \mu(h(x) - y) dx < \epsilon$.

Additionally, class H and empirical measure μ are empirically weakly learnable with rate ϵ if there exists $h \in H$ so that $\int y h(x) p(x; y) d(x; y) \geq \epsilon$ for every reweighting p of measure μ .

As discussed previously, the assumption $\gamma > 0$ (which is implied by $\epsilon > 0$) is greatly beneficial to optimization, with AdaBoost needing only $O(\ln(1/\epsilon) / \gamma^2)$ iterations to achieve accuracy $\epsilon > 0$.

What are some ways for the weak learning assumption $\gamma > 0$ to fail? An easy case is that the sample $\{(x_i; y_i)\}_{i=1}^m$ has $x_i = x_j$ for some $i \neq j$, but $y_i \neq y_j$: the dual player can simply place positive probability on these two examples, and zero probability elsewhere. Additionally, it is possible that the sample is separable but each element only breaks even by a constant number of examples, meaning $\gamma = O(1/m)$, which will be momentarily identified as inadequate for statistical guarantees based on γ .

Remark 1.7.2. The largest achievable empirical weak learning rate γ is equivalent to a measure of the separability of an instance, called the maximum l_1 margin (Shalev-Shwartz and Singer, 2008). Moreover, AdaBoost, with a number of step size and loss function choices, is able to achieve nearly maximum l_1 margins in a small number of iterations (Telgarsky, 2013). When γ is positive and not too small, for instance when $\frac{1}{\sqrt{V(H)}} = \Omega(\gamma)$ as $m \rightarrow \infty$ (where $V(H)$, the VC dimension of H , measures the complexity of H), AdaBoost exhibits good generalization performance (Schapire et al., 1997).

As the analysis in this thesis is primarily concerned with worst case guarantees, this route of margin assumptions and analysis is not explicitly considered, though ideas from it appear in various places. Since large or even positive margins can not be guaranteed, thus guarantees based on γ and L are more generally applicable.

It is worth mentioning, however, when there is prior knowledge on the structure of the prediction problem encoded in \mathcal{X} , it may be possible to engineer a hypothesis class which ensures large margins, and moreover, the engineering effort involved in implementing AdaBoost variants is typically largely spent in this process.

As a consequence, it is arguably AdaBoost's greatest strength that it possesses both an excellent set of margin-based bounds under favorable values of γ , and a strong set of fallback guarantees based on loss minimization properties.

In the previous section, the question of \mathcal{H} possessing minimizers was raised; how does

this relate to the weak learning assumption? As discussed previously, and as will be proved in Chapter 3 for choices of η which include η_{exp} and η_{log} , this assumption implies that $L^{\eta}(A_{\eta})$ decreases exponentially quickly, which in particular means $L^{\eta}(A) = 0$. On the other hand, since $\eta_{\text{exp}} > 0$ and $\eta_{\text{log}} > 0$, it follows that $L^{\eta}(A) > 0$ for either loss and any $\eta > 0$: in particular, the empirical weak learning assumption $(b) > 0$ and the existence of minimizers never coincide. But what is a more complete characterization of the objective function?

To this end, another basic structure related to L and L^{η} is the hard core, defined as follows.

Definition 1.7.3. Given finite hypothesis class H and probability measure μ , let $D(H; \mu)$ denote reweightings of μ which decorrelate every regressor $h \in H$; that is,

$$D(H; \mu) := \{p \in L^1(\mu) : p \geq 0 \text{ a.e.}, \int_{\mathcal{X}} h(x) p(x) d\mu(x) = 0 \text{ for all } h \in H\}$$

Correspondingly, $S_D(H; \mu)$ tracks the supports of these weightings:

$$S_D(H; \mu) := \{X \subseteq \mathcal{X} : \exists p \in D(H; \mu) \text{ with } p > 0 \text{ on } X\}$$

A hard core $C \subseteq \mathcal{X}$ for $(H; \mu)$ is a maximal element of $S_D(H; \mu)$; that is,

$$C \in S_D(H; \mu) \text{ and } \mu(C \setminus C') = 0 \text{ for all } C' \in S_D(H; \mu)$$

(“Maximal”, in the presence of measures, will always mean up to sets of measure zero.)

In words, $S_D(H; \mu)$ is a collection of sets for which there exists reweightings $p \in L^1(\mu)$ which decorrelate every function $h \in H$ from the prediction problem encoded in μ , meaning $\int_{\mathcal{X}} h(x) p(x) d\mu(x) = \int_{\mathcal{X}} h(x) \mu(x, y) p(x, y) d\mu(x, y) = 0$. And amongst the elements of $S_D(H; \mu)$, a hard core C is maximal in a measure-theoretic sense, meaning no element $C' \in S_D(H; \mu)$ exceeds C by a set of positive measure.

Crucially, hard cores always exist.

Theorem 1.7.4. Every finite hypothesis class H and probability measure μ admits a hard core, and any two hard cores disagree only on a set of measure zero.

To prove this, note first that $S_D(H; \mu)$ is never empty—it will at least contain the empty set (with corresponding weighting $p = 0$). Furthermore, $S_D(H; \mu)$ is closed under countable

unions; but simply unioning together all of $S_D(H; \cdot)$ will not produce a hard core, since the corresponding object need not be measurable; the full proof is in Section 1.D.

The hard core provides a handle on the relationship between weak learnability and minimizers. Specifically, in the finite sample optimization guarantees of Chapters 2 and 3, any finite sample can be split uniquely into two pieces, where one piece satisfies the weak learning assumption, and the other has minimizers, and moreover considering these two pieces separately leads to a general analysis. Thereafter, Chapter 5 uses hard cores to similarly split the statistical deviations into two pieces, and thereafter control each piece separately.

1.8 Analysis Overview

The main byproducts of this thesis are a number of worst case and high probability guarantees for AdaBoost variants as defined by Boost in Figure 1.3; in particular, although it is a fruitful and important topic, this thesis does not focus on scenarios where $(b) > 0$ is favorably non-small as discussed in Remark 1.7.2. Concretely, the first set of guarantees are of numerical convergence, namely showing the ability of Boost to minimize L or in some cases Φ (Chapter 3 and Chapter 6). The second set of guarantees are of statistical convergence, proving that the output sequence of iterates also minimizes L or in some cases Φ (Chapter 5 and Chapter 6), which in turn can lead to statements of statistical consistency (Chapter 6).

In order to systematically develop these byproducts, this thesis devotes significant effort to identifying the structure of the primal and dual optimization problems identified in Section 1.3 and Section 1.4, and moreover their impact on the behavior of the algorithm. This investigation, initiated in this introductory chapter (particularly in the immediately preceding Section 1.7), is the focus of Chapter 2 and Chapter 4, but can also be found throughout the other chapters as well.

In more detail, the complete layout of this thesis is as follows. Note that most proofs are relegated to appendices closing each chapter.

Chapter 2. This chapter and the following focus on finite samples, meaning the empirical objective L ; additionally, most results assume that H has some form of bounded complexity, whereby the operators A and H , restricted to the sample, may be written as finite dimensional matrices.

This chapter first proves a strong duality result, which is used to characterize hard cores over finite samples. Additionally, this chapter develops an analog to β which is never zero, and thus is a data-dependent quantity which characterizes the behavior of the algorithm without any assumptions.

Chapter 3. Leveraging the structures in Chapter 2, optimization guarantees follow. In particular, when either $\beta > 0$ (the hard core is empty) or minimizers exist (the hard core is the entire sample), the number of iterations to reach accuracy $\epsilon > 0$ is $O(\ln(1/\epsilon))$. For general instances, the problem may be broken into the hard core and its complement, analyzed using the preceding tools on both pieces, and then stitched together to achieve a rate $O(\epsilon)$, with a matching lower bound for the logistic loss ℓ_{\log} .

Chapter 4. This next chapter starts the second part of the thesis, which focuses on statistical questions. Chapter 4 is a gentle introduction, presenting a strong duality analog to Proposition 1.4.1, and identifying where difficulties begin to arise when the measure is not assumed to correspond to a finite sample.

Chapter 5. The first set of statistical guarantees are against an arbitrary measure μ (with potentially infinite support), but still assuming H is finite. While this is not quite the desired level of generality, this family of results have two points of value. First, the results are stronger than those presented later in Chapter 6 for H infinite: more loss functions are handled, the deviation bounds are tighter, and less demands are made on the minimization algorithm. Second, note that this guarantee may be invoked by starting with an infinite H , and then feeding an appropriate initialization to an algorithm; in a sense, this indicates that an effective penalization scheme for boosting is to determine how much data is available, and then appropriately curtail the complexity of H before invoking the algorithm.

The method of proof here is to develop the structure of hard cores more carefully in the case of a general measure (and not simply an empirical measure corresponding to a finite sample as in Chapter 2).

Chapter 6. Lastly, losses which are Lipschitz (plus a few other regularity conditions) are handled in the case that H is infinite; amongst the four losses presented in Section 1.3, only the logistic loss ℓ_{\log} survives, though some of the analysis also holds for ℓ_{russ} . The analysis here

considers two cases: either the minimal surrogate risk over the distribution is zero, or it is positive. When it is zero ($L(A) = 0$), the convergence rate is fast, and moreover based on an analog to which is statistically stable and in a sense better adapted to Lipschitz losses. On the other hand, when the minimal surrogate risk is positive ($L(A) > 0$), it is possible to exhibit distribution-dependent controls on the norms of the output predictors, which again leads to statistical controls.

Appendix 1.A Banach Spaces and Linear Operators

The purpose of this section is to identify basic properties of the Banach spaces and linear operators relevant in this thesis, which are fundamental to the behavior of convex duality theory (since the dual is nothing but a set of linear functions and a topology that renders them continuous); in particular, this section will provide the definition of the dual pairing $\langle \cdot, \cdot \rangle$, which is used in the definition of Fenchel conjugate.

H and A are mappings which produce bounded functions; the bulk of the analysis, however, considers them as producing functions over $L^1(X)$ and $L^1(\mu)$ as follows (where μ is a probability distribution over $X \in \mathbb{R}^d$).

Lemma 1.A.1. Let μ be a probability measure over $X \in \mathbb{R}^d$, and let μ^X denote the marginal distribution over X .

1. The definition of H and A is valid for arbitrary weightings μ ; in particular, $\text{supp}(\mu)$ is countable, and

$$\begin{aligned} (H)_x &= \int_Z h(x) \mu(h) d\mu(h) = \int_{h \in \text{supp}(\mu)} (H)_x(h); \\ (A)_{x,y} &= \int_Z y \cdot (H)_x(h) d\mu(h) = \int_{h \in \text{supp}(\mu)} y \cdot (H)_x(h); \end{aligned}$$

2. H and A are linear operators.

3. $H : L^1(X) \rightarrow \mathbb{R}$ and $A : L^1(\mu) \rightarrow \mathbb{R}$ are continuous linear operators (with unit norm).

Proof. If μ is a counting measure, then $\sum_j \mu(h_j) < \infty$, and since μ is a counting measure, it follows that $\text{supp}(\mu)$ is countable. Furthermore, for any $x; y; h$, $|y \cdot (H)_x(h)| \leq |y| \cdot |h(x)| \leq |y| \cdot \sum_j |h_j(x)| \leq |y| \cdot \sum_j \mu(h_j) < \infty$, and thus the rescalings

$h(x)$ and $\gamma h(x)$ are both in \mathcal{H} , and in particular

$$(H \gamma)_x = \int_Z h(x) \gamma(h) d\mu(h) = \int_{h^2 \text{supp}(\gamma)} (H \gamma)_x \gamma(h);$$

and similarly for A .

It follows by definition (and another check for integrability) that $(H(a \gamma_1 + b \gamma_2))_x = a(H \gamma_1)_x + b(H \gamma_2)_x$, and thus H is a linear operator; the proof for A is the same.

Lastly, H is continuous with unit norm, since boundedness of each γ combined with μ being a probability measure gives

$$\begin{aligned} \sup_{\gamma \in \mathcal{H}} \|H \gamma\|_2 &= \sup_{\gamma \in \mathcal{H}} \left(\int_Z (H \gamma)_x^2 d\mu(x) \right)^{1/2} \\ &\leq \sup_{\gamma \in \mathcal{H}} \left(\int_Z \int_{h^2 \text{supp}(\gamma)} h(x) \gamma(h) d\mu(h) d\mu(x) \right)^{1/2} \\ &\leq \sup_{x,h} |h(x)| \left(\int_{h^2 \text{supp}(\gamma)} \gamma(h) d\mu(h) \right)^{1/2} \\ &= 1. \end{aligned}$$

(The proof for A is the same, since $\|f\|_1 \leq \|f\|_2$ implies $\|j\gamma h(x)\|_2 = \|j\gamma h(x)\|_1$.) □

Note, H and A may also be defined as Bochner (or similar) integrals.

Next, to develop the adjoint of $A^>$, relevant dual spaces need to be established (the adjoint of $H^>$ does not appear, but is similar).

Lemma 1.A.2. If μ is a probability measure over X $f \in L^1(\mu)$, then $L^1(\mu)$ is isometrically isomorphic to $L^1(\mu)$, and in particular for every $Q \in L^1(\mu)$ there exists $q \in L^1(\mu)$ so that $Q(f) = \int_{\mathbb{R}} qf d\mu$ for every $f \in L^1(\mu)$. Similarly, recalling $\mu = L^1(\mu)$ where μ is counting measure over some class \mathcal{H} , the dual $L^1(\mu)$ is isometrically isomorphic to $L^1(\mu)$, and once again elements of $L^1(\mu)$ can be written as integrals over \mathcal{H} with an element of $L^1(\mu)$.

Proof. The first relationship follows since μ is a probability measure and thus μ -finite (Folland, 1999, Theorem 6.15), and the second is a general property of counting measures (even though the cardinality of \mathcal{H} may preclude μ from being σ -finite) (Folland, 1999, Exercises 3.15 and 6.25). □

This thesis always identifies the above dual spaces by the provided isometric isomorphism,

a fact which will be crucial in the convex duality theory of L^1 (cf. Lemma 4.1.1). The following definition clarifies this impact.

Definition 1.A.3. Let X and X^* be a Banach space and a dual space. Given $x \in X$ and $x^* \in X^*$, define the dual pairing $\langle x, x^* \rangle = x^*(x)$ (Rudin, 1973, Chapter 4). When X is identified via isometric isomorphism, this dual pairing provides the means to compute these linear functions. In particular, in light of Lemma 1.A.2:

$$\text{Given } \mu \in \mathcal{P}(X) \text{ and } q \in L^1(\mu), \int q d\mu = \int \langle x, q \rangle d\mu(x).$$

$$\text{Given measure } \mu \text{ over } X, f, g \in L^1(\mu), \text{ as well as } p \in L^1(\mu), \int \langle f, p \rangle d\mu = \int f(x) p(x) d\mu(x).$$

$$\text{Given empirical measure } \mu_n \text{ corresponding to a finite sample of size } n, \text{ as well as } f \in L^1(\mu_n) \text{ and } p \in L^1(\mu_n), \int \langle f, p \rangle d\mu_n = \frac{1}{n} \sum_{i=1}^n f(x_i) p(x_i) = \int f p d\mu_n.$$

Remark 1.A.4. Note that the dual pairings with respect to $L^1(\mu)$ and $L^1(\mu_n)$ are not consistent, which for instance causes the artifact of an $\frac{1}{n}$ appearing in the dual form over \mathbb{D} as compared with L^1 (concretely compare Theorem 2.2.1 to Proposition 1.4.1 and lemma 4.1.1). By making the dual pairing over $L^1(\mu_n)$ agree with the standard scalar product, however, duality and in particular gradient manipulations of $L^1(\mu_n)$ are consistent with the standard forms over finite dimensional Euclidean spaces, which hopefully improves readability of any such manipulations, particularly those in Chapters 2 and 3, and moreover allows the application of results from the literature which assumed the dual pairing is simply the scalar product.

To avoid confusion, dual pairing notation will often be avoided, but notice that it is at the heart of the convex duality relations, since the Fenchel conjugate to a function $g: X \rightarrow \mathbb{R}$ is defined as

$$g^*(x^*) = \sup_{x \in X} \langle x, x^* \rangle - g(x).$$

Turning back to H and A , the adjoint A^* has the following structure.

Lemma 1.A.5. Let probability measure μ over X , $f, g \in L^1(\mu)$ and any H be given.

1. Considering A as a linear operator from $L^1(\mu)$ to $L^1(\mu)$, its adjoint $A^*: L^1(\mu) \rightarrow L^1(\mu)$ is the unique continuous linear operator satisfying $\langle A^* p, f \rangle = \langle p, A f \rangle$, where $p \in L^1(\mu)$ and $f \in L^1(\mu)$ (and dual spaces have been identified via isomorphism as in Lemma 1.A.2).

2. Again identifying $A^\triangleright p$ for $p \in L^1(X, \mu)$ with an element of $L^1(X, \mu)$,

$$\begin{aligned} \|A^\triangleright p\|_1 &= \sup_{\|h\|_2=1} \int_Z |y h(x) p(x; y)| d(x; y) \\ &= \sup_{\|h\|_2=1} \int_Z |A^\triangleright p|_{x; y} p(x; y) d(x; y) \end{aligned}$$

3. The map $p \mapsto \|A^\triangleright p\|_1$ is a convex function over $L^1(X, \mu)$, and is lower semi-continuous in the weak* topology (i.e., the weak topology induced on $L^1(X, \mu)$ by $L^\infty(X, \mu)$).

Proof. 1. Recall by Lemma 1.A.1 that A is a continuous linear operator; the basic properties of A^\triangleright follow by properties of adjoints of continuous linear operators (Rudin, 1973, Theorem 4.10) combined with the isometric isomorphism of the relevant dual spaces as provided by Lemma 1.A.2.

2. Let $p \in L^1(X, \mu)$ be given. Since the isometric isomorphism provided by Lemma 1.A.2 allows $A^\triangleright p$ to be identified with an element of $L^1(X, \mu)$, the operator norm of $A^\triangleright p$ is simply the $L^1(X, \mu)$ norm of the element it has been identified with by the isomorphism. Since μ is a counting measure, letting $e_h \in L^1(X, \mu)$ be an indicator function for a single $h \in H$ (it is 1 on h and 0 elsewhere), using the above adjoint relation $\|A^\triangleright p\|_1 = \sup_{\|h\|_2=1} \int_Z |y h(x) p(x; y)| d(x; y)$, and using the definition of norms on $L^1(X, \mu)$,

$$\begin{aligned} \|A^\triangleright p\|_1 &= \inf_{a \geq 0} \{ \int_Z |y h(x) p(x; y)| d(x; y) > a \} \\ &= \inf_{a \geq 0} \{ \int_Z |y h(x) p(x; y)| d(x; y) > a \} \\ &= \inf_{a \geq 0} \{ \int_Z |y h(x) p(x; y)| d(x; y) > a \} \\ &= \sup_{\|h\|_2=1} \int_Z |y h(x) p(x; y)| d(x; y) \end{aligned}$$

where the last equality can be established by noting the domain of the infimum includes all $a \geq 0$ satisfying $a < \int_Z |y h(x) p(x; y)| d(x; y)$, but no values satisfying $a < \int_Z |y h(x) p(x; y)| d(x; y)$.

Next, to show

$$\begin{aligned} \sup_{h \in H} \int_{\mathcal{X}} y h(x) p(x; y) d(x; y) & \\ & = \sup_{\substack{A \subseteq \mathcal{X}; y \\ k_1 \leq 1}} \int_{\mathcal{X}} y p(x; y) d(x; y) \end{aligned}$$

one direction is immediate, since positive and negative copies e_h of the indicator elements satisfy $e_h \in H$ and $k_1 e_h k_1 = 1$. For the other direction, let $\epsilon > 0$ be arbitrary, and choose any $\delta > 0$ which is within ϵ of the supremum on the right side of the display. Then, since $\text{supp}(p(\cdot; y))$ is countable (via Lemma 1.A.1), and since $k_1 \leq 1$ implies $k_1 A \subseteq A$, the dominated convergence theorem (Folland, 1999, Theorem 2.25 (summation form)) may be applied (with dominating function 1), and

$$\begin{aligned} \sup_{\substack{A \subseteq \mathcal{X}; y \\ k_1 \leq 1}} \int_{\mathcal{X}} y p(x; y) d(x; y) & \\ & + \int_{\mathcal{X}} y p(x; y) d(x; y) \\ & = \int_{\mathcal{X}} y h(x) p(x; y) d(x; y) \\ & + \int_{\mathcal{X}} y p(x; y) d(x; y) \\ & = \int_{\mathcal{X}} y p(x; y) d(x; y) \\ & + k_1 \sup_{h \in H} \int_{\mathcal{X}} y h(x) p(x; y) d(x; y) \end{aligned}$$

since $k_1 \leq 1$, and since $\epsilon > 0$ was arbitrary, the result follows.

3. For the last part, define a convex indicator over $L^1(\mathcal{X})$ as

$$(f) = \begin{cases} 0 & \text{when } f \in A : \int_{\mathcal{X}} y p(x; y) d(x; y) \leq k_1; \\ 1 & \text{otherwise} \end{cases}$$

(Note that (f) is not necessarily lower semi-continuous over $L^1(\mathcal{X})$, since as discussed shortly in Lemma 1.A.6, the subspace A might not be closed.) The conjugate of (f) is, for any

$p \in L^1(\mathbb{R})$,

$$\begin{aligned} (p) &= \sup_{\substack{Z \\ f: \mathbb{R} \rightarrow \mathbb{R}; k \leq 1, \int f = A}} \\ &= \sup_{\substack{Z \\ (A) p: \mathbb{R} \rightarrow \mathbb{R}; k \leq 1}} \\ &= kA \geq pk_1; \end{aligned}$$

where the last step used the earlier equalities for $A \geq 0$. Since $p \in L^1(\mathbb{R})$ is the conjugate of a convex function, it is lower semi-continuous in the weak* topology (Zalinescu, 2002, Theorem 2.3.1(i)).

□

Lastly, note the following properties of the sets H and A .

Lemma 1.A.6. Let any H and any probability measure μ over X be given. Then H and A are subspaces, but it is possible that neither is closed in its respective $L^1(X)$ and $L^1(\mathbb{R})$ topology (indeed, Example 6.3.4 provides the counterexample).

Proof. Since \mathbb{R} is a Banach space and H and A are linear operators, it follows that H and A are subspaces.

For the lack of closure, consider the setting of Example 6.3.4, and in particular building a sequence of functions f_k, g_k which are a combination of k thresholds, and predict correctly on the last k intervals. This sequence has a limit point in $L^1(X)$ (in particular, it is a countable sum of indicators over intervals), but no such function is in H , which is therefore not closed in $L^1(X)$. To obtain a similar result for A , define $g_k(x; y) := f_k(x)$. □

Appendix 1.B Convexity Properties of ψ and \mathbb{B}

First, the following basic properties of ψ are used throughout the thesis.

Lemma 1.B.1. Suppose $\psi: \mathbb{R} \rightarrow \mathbb{R}_+$ is convex with $\lim_{z \rightarrow 1} \psi(z) = 0$.

1. ψ is lower semi-continuous, whereby ψ is convex lower semi-continuous, and $\psi = \psi^*$.
2. $\psi(z) = 1$ for $z < 0$, and $\psi(0) = 0$.

3. Let $L := \sup_{x \neq y} \frac{|f(x) - f(y)|}{|x - y|}$ denote the tightest Lipschitz constant for f . If $L < 1$, then $f'(x) = 1$ when $x > 0$, and $f'(x) < 1$ when $x \in [0, 1]$. If $L = 1$, then f' is finite along $[0, 1)$.
4. If $g \in \partial f(0)$ is any subgradient of f at the origin and $f'(0) > 0$, then $f'(x) < 0$ for $x \in (0, g)$, and f' attains its minimum value $f'(0)$ at g .
5. If f is differentiable and strictly convex, then f' is strictly convex and continuously differentiable over the interior of its domain $\text{dom}(f) = f^{-1}(\mathbb{R}) : f'(x) < 1/g$.

Proof. 1. Since f' is finite everywhere, it is continuous (thus lower semi-continuous), and thus $f'' = f'$ and f' is convex lower semi-continuous (Rockafellar, 1970, Theorem 12.2).

2. For any $x \in \mathbb{R}$ and subgradient $g_x \in \partial f(x)$, $f'(0) = f'(x) + g_x(0 - x)$. Since $\lim_{z \downarrow 1} f'(z) = 0$ and f' is convex, it follows that f' is nondecreasing, meaning $g_x \geq 0$, and thus, for any $x < 0$,

$$f'(x) = \sup_{x \in \mathbb{R}} \frac{f(x) - f(0)}{x - 0} = \sup_{x \in \mathbb{R}} \frac{f(x) - f(0) - g_x x}{x - 0} = \sup_{x < 0} \frac{f(x) - f(0) - g_x x}{x - 0} = 1 :$$

Additionally, since $\inf_x f'(x) = 0$,

$$f'(0) = \sup_x \frac{f(x) - f(0)}{x - 0} = \inf_x f'(x) = 0 :$$

3. Next, suppose f' has tightest Lipschitz parameter $L \in [0, 1]$, whereby the any subgradient g_x at a point x satisfies $|g_x| \leq L$. Consequently, proceeding just as in the study of the case $L < 1$, for any $x > 0$ (which is vacuous when $L = 1$),

$$f'(x) = \sup_{x \in \mathbb{R}} \frac{f(x) - f(0) - g_x x}{x - 0} = \sup_{x > 0} \frac{f(x) - f(0) - g_x x}{x - 0} = 1 :$$

On the other hand, let $h \in (0, 1)$ be arbitrary, whereby there must exist $x > y$ with

$$h < \frac{f(x) - f(y)}{x - y}$$

(where the absolute values were dropped since $x > y$ and f' is nondecreasing). Taking any $h \in (0, 1)$, note

$$h < \frac{f(x) - f(y)}{x - y} = \frac{f(x) - (f(x) + h(y - x))}{x - y} = h :$$

Consequently, by the Fenchel-Young inequality,

$$\langle h, x \rangle - f^*(x) < 1 :$$

Since f^* is convex, it is finite over a convex set. Since h was arbitrary, it follows that f^* is finite over $[0; \infty)$. If $\alpha = 1$, the proof is done; otherwise, when $\alpha < 1$, since f^* is lower semi-continuous and finite over $[0; \infty)$, it must also hold that $f^*(\alpha) < 1$.

4. Let $g \in f^*(0)$ be given; by the Fenchel-Young inequality and $f^*(0) > 0$,

$$\langle g, 0 \rangle - f^*(0) < 0:$$

Since $f^*(0) = 0$ and f^* is closed and convex, the first part follows. For the second part, since f^* is closed and convex $g \in f^*(0)$ implies $0 \in \langle g, \cdot \rangle$ (Rockafellar, 1970, Theorem 23.5), which is precisely the first order optimality condition (Borwein and Lewis, 2000, Proposition 3.1.5).

5. Since f^* is strictly convex, then f^* is continuously differentiable on $\text{int}(\text{dom}(f^*))$ (Hiriart-Urruty and Lemaréchal, 2001, Theorem E.4.1.1). Since f^* is differentiable, then f^* is strictly convex along every convex subset of \mathbb{R} (Hiriart-Urruty and Lemaréchal, 2001, Theorem E.4.1.2). Now consider any $\alpha \in \text{int}(\text{dom}(f^*))$; since f^* is lower semi-continuous (as established above), $f^* = r^*$ ($r^*(\cdot) = r^*(r^*(\cdot))$) (Hiriart-Urruty and Lemaréchal, 2001, Corollary E.1.3.6, Proposition E.1.4.3), which means $r^*(\mathbb{R}) \subseteq \text{int}(\text{dom}(f^*))$. Since f^* is convex, $\text{int}(\text{dom}(f^*))$ is a convex subset of $r^*(\mathbb{R})$, which means, by the above, that f^* is strictly convex on $\text{int}(\text{dom}(f^*))$.

It remains to be shown that f^* is strictly convex on $\text{dom}(f^*)$, note just the interior. As such, let $\alpha_1 < \alpha_2$ and $\alpha_2 \in (0; 1)$ be given with $\alpha_1, \alpha_2 \in \text{int}(\text{dom}(f^*))$. Set $\alpha_3 := \alpha_1 + (1 - \alpha_2) \alpha_2 \in \text{int}(\text{dom}(f^*))$, whereby

$$\frac{\alpha_1 + \alpha_3}{2} \in \text{int}(\text{dom}(f^*)) \quad \text{and} \quad \frac{\alpha_2 + \alpha_3}{2} \in \text{int}(\text{dom}(f^*));$$

thus

$$\begin{aligned}
 \psi(\lambda_1 + (1 - \lambda_1)\lambda_2) &= \lambda_1 \frac{1 + \lambda_1^3}{2} + (1 - \lambda_1) \frac{1 + \lambda_2^3}{2} \\
 &< \lambda_1 \frac{1 + \lambda_1^3}{2} + (1 - \lambda_1) \frac{1 + \lambda_2^3}{2} \\
 &= \lambda_1 \frac{(1 + \lambda_1) + (1 + \lambda_1^3)}{2} + (1 - \lambda_1) \frac{(1 + \lambda_2) + (1 + \lambda_2^3)}{2} \\
 &\quad + \frac{(1 + \lambda_1) + (1 + \lambda_1^3)}{2} \lambda_1 + \frac{(1 + \lambda_2) + (1 + \lambda_2^3)}{2} (1 - \lambda_1) \\
 &= \lambda_1 \psi(\lambda_1) + (1 - \lambda_1) \psi(\lambda_2):
 \end{aligned}$$

Therefore, to show ψ is strictly convex over $\text{dom}(\psi)$, first choose $\lambda_1 = 0$ in the above derivation, whereby ψ is strictly convex over $\text{int}(\text{dom}(\psi)) \cap [0, \infty)$. If $\text{dom}(\psi) = \mathbb{R}_+$, the proof is done; otherwise, as established above, $\text{dom}(\psi) = [0; \lambda_2]$, and the result follows by choosing $\lambda_2 = \lambda_1$.

□

Note that loss structure presented in Section 1.4 provides a subset of the guarantees from above.

Proof of Proposition 1.4.2. This result is a combination of various statements in Lemma 1.B.1; in particular, the uniqueness of the optimum $\psi^*(0)$ is from strict convexity of ψ . □

Next, note that \mathcal{L} also possesses reasonable convexity structure; the relevant part to this chapter is the existence and form of Gâteaux derivatives.

Proposition 1.B.2. Let convex $\psi: \mathbb{R} \rightarrow \mathbb{R}_+$ be given, along with empirical measure \mathbb{P}_n over $X \times \mathbb{R}$ (corresponding to some finite sample $(x_i; y_i)_{i=1}^n$), and corresponding map \mathcal{L} defined as $\mathcal{L}(f) = \int_{\mathbb{R}} \psi(f(z)) d\mathbb{P}_n(z)$, and hypothesis class \mathcal{H} with corresponding map $A: \mathcal{H} \rightarrow L^1(\mathbb{P}_n)$ be given.

1. \mathcal{L} and $\mathcal{L} \circ A$ are convex.
2. \mathcal{L} is continuous over $L^1(\mathbb{P}_n)$, and $\mathcal{L} \circ A$ is continuous over \mathcal{H} .

3. $\mathbb{L}(p) = \int_{\mathbb{R}^n} (mp) db$ and $\mathbb{L}(0) = 0$. Furthermore, if ψ has tightest Lipschitz constant $L := \sup_{x \neq y} |\psi(x) - \psi(y)| / |x - y|$, then $\mathbb{L}(p) < 1$ if $p \in L^1(b)$ satisfies $p \geq 0$ and $\int p db = m$ everywhere.
4. Given $f \in L^1(b)$, then $\mathbb{L}(f) \in L^1(b)$ and $\int \mathbb{L}(f) db = \int f db$ everywhere. Additionally, given $A \in \mathcal{A}$, $\mathbb{L}(A) \in L^1(b)$ and $\int \mathbb{L}(A) db = \int A db$.
5. If ψ is differentiable, then for any $f \in L^1(b)$ the Gâteaux derivative $\mathbb{L}(f)$ exists (and is an element of $L^1(b)$), and moreover $(\mathbb{L}(f))_{x_i, y_i} = \psi'(f(x_i; y_i))$ almost everywhere. Similarly, again assuming differentiability of ψ , then for any $A \in \mathcal{A}$ the Gâteaux derivative $(\mathbb{L}(A))_{x_i, y_i}$ exists (now an element of $L^1(b)$), and moreover satisfies $(\mathbb{L}(A))_{x_i, y_i} = \psi'(A(x_i; y_i))$ and for any $h \in \mathbb{R}^n$

$$\begin{aligned} D_{r(\mathbb{L}(A))_{x_i, y_i}} \psi &= \frac{1}{m} \sum_{i=1}^n \int_{\mathbb{R}^n} \psi'(A(x_i; y_i)) h(x_i) db \\ &= \frac{1}{m} \sum_{i=1}^n \int_{\mathbb{R}^n} \psi'(A(x_i; y_i)) h(x_i) db \end{aligned}$$

Evaluated at a coordinate $e_h \in \mathbb{R}^n$ (for $h \in \mathbb{R}^n$), this simplifies to

$$D_{r(\mathbb{L}(A))_{x_i, y_i}} \psi(e_h) = \frac{1}{m} \sum_{i=1}^n \psi'(A(x_i; y_i)) h(x_i)$$

6. If ψ is differentiable and strictly convex, then \mathbb{L} is strictly convex, and continuously differentiable along the interior of its domain $\text{dom}(\mathbb{L}) = \{p \in L^1(b) : \mathbb{L}(p) < 1\}$.

Proof. 1. Let $f, g \in L^1(b)$ and $\lambda \in [0, 1]$ be given; since ψ is convex,

$$\begin{aligned} \int \psi(\lambda f + (1-\lambda)g) db &\leq \int (\lambda \psi(f) + (1-\lambda)\psi(g)) db \\ &= \lambda \int \psi(f) db + (1-\lambda) \int \psi(g) db \end{aligned}$$

For $A \in \mathcal{A}$, let $\lambda \in [0, 1]$ and $\mu \in [0, 1]$ be given; since $A : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is linear (cf. Lemma 1.A.1),

$$A(\lambda x + (1-\lambda)y) = \lambda A(x) + (1-\lambda)A(y),$$

whereby convexity now follows by convexity of ψ ; alternatively, $A \in L^1(b)$ and $A \in L^1(b)$

may be plugged into the above derivation.

2. For \mathbb{D} , let $f \in L^1(b)$ and $\epsilon > 0$ be given; it suffices to exhibit $\delta > 0$ such that $\|f - f^0\|_{k_1} < \delta$ implies $|\mathbb{D}(f) - \mathbb{D}(f^0)| < \epsilon$. To this end, define real $M := \max_{i \in [m]} \int f(x_i; y_i)$ and compact set $S := [-M - 1; M + 1]$; since ψ is convex over \mathbb{R} , it has a Lipschitz constant L over S (Hiriart-Urruty and Lemaréchal, 2001, Theorem B.3.1.2). With these definitions in place, choose $\delta := \min\{\epsilon/m; \epsilon/L\}$. Now let $f^0 \in L^1(b)$ with $\|f - f^0\|_{k_1} < \delta$ be arbitrary. It follows by $\delta \leq \epsilon/m$ and choice of S that f^0 , evaluated over the sample, lies within S . Secondly, since $\delta \leq \epsilon/L$,

$$\begin{aligned} |\mathbb{D}(f) - \mathbb{D}(f^0)| &= \left| \frac{1}{m} \sum_{i=1}^m \psi(f(x_i; y_i)) - \psi(f^0(x_i; y_i)) \right| \\ &\leq \frac{L}{m} \sum_{i=1}^m |f(x_i; y_i) - f^0(x_i; y_i)| = L \|f - f^0\|_{k_1} < \epsilon. \end{aligned}$$

It follows that \mathbb{D} is continuous. This also grants continuity of \mathbb{D}^*A , since A is a continuous linear operator (cf. lemma 1.A.1).

3. For any $p \in L^1(b)$,

$$\begin{aligned} \mathbb{D}^*(p) &= \sup_{f \in L^1(b)} \langle f; p \rangle - \mathbb{D}(f) \\ &= \frac{1}{m} \sum_{i=1}^m \sup_{f_i \in \mathbb{R}} \langle p(x_i; y_i) f_i - \psi(f_i) \rangle \\ &= \int (mp) db \end{aligned}$$

(The equivalence classes of functions comprising $L^1(b)$ each contain simple functions; for the general case $L^1(b)$, more work is necessary to develop the conjugacy relation (cf. Lemma 4.1.2).) From here it directly follows that $\mathbb{D}^*(0) = 0$ since Lemma 1.B.1 granted $\psi(0) = 0$.

For the second part, if $p(x_i; y_i) \in [0; \epsilon/m]$ for some $i \in [m]$, then Lemma 1.B.1 provides $\psi(\epsilon/m) = \epsilon$, whereby the above conjugacy rule provides $\mathbb{D}^*(p) = \epsilon$. On the other hand, if $p \in [0; \epsilon/m]$ everywhere, then Lemma 1.B.1 provides $\psi(\epsilon/m)$ is finite everywhere, and the above conjugacy relation gives $\mathbb{D}^*(p) < \epsilon$.

4. Given $f \in L^1(b)$, by the Fenchel-Young inequality, a function $p \in L^1(b)$ satisfies $p \in \text{subd}(f)$ if $\int p + \int f = \int pf$ (Zalinescu, 2002, Theorem 2.4.2.iii). Writing this out, $p \in \text{subd}(f)$ if

$$\int p(x_i; y_i) f(x_i; y_i) = \frac{1}{m} \sum_i (f(x_i; y_i)) + \frac{1}{m} \sum_i (mp(x_i; y_i)):$$

For every i with $p(x_i; y_i) \in \text{subd}(f(x_i; y_i))=m$, Fenchel-Young again grants $p(x_i; y_i) f(x_i; y_i) = f(x_i; y_i) + mp(x_i; y_i)$ (Zalinescu, 2002, Theorem 2.4.2.iii). On the other hand, if some i has $p(x_i; y_i) \notin \text{subd}(f(x_i; y_i))=m$, then Fenchel-Young holds with strict inequality (Zalinescu, 2002, Theorem 2.4.2.iii); consequently $p \in \text{subd}(f)=m$ everywhere.

For $\text{subd}(A)$, since A is a continuous linear operator (cf. Lemma 1.A.1), and since it was established above that subd is continuous everywhere on $L^1(b)$, the subdifferential rule $\text{subd}(A)(\cdot) = A^* \text{subd}(A)$ holds (Borwein and Zhu, 2005, Theorem 4.3.3).

5. When \cdot is differentiable, subd now holds a single equivalence class of a.e. equivalent functions, and similarly for $\text{subd}(A)$, which means a Gâteaux derivative exists in both cases, and each time is equal to the corresponding equivalence class identified in the subdifferential rule (Borwein and Zhu, 2005, Theorem 4.2.9).

In order to simplify the expression for $r(\text{subd}(A))(\cdot)$ for $\cdot \in \mathbb{R}^n$, which is an element of $L^1(\cdot)$, by the earlier dual pairing and the above expansion of subd , for any $\cdot \in \mathbb{R}^n$,

$$\begin{aligned} D_{r(\text{subd}(A))}(\cdot); \cdot &= D_{A^* r(\text{subd}(A))}; \cdot \\ &= D_{r(\text{subd}(A))}; A \cdot \\ &= \sum_{i=1}^n (\cdot)_i ((A)_{x_i, y_i}) = m \sum_{h \in H} y_i \sum_{h \in H} h(x_i) \cdot(h) \\ &= \frac{1}{m} \sum_{i=1}^n \sum_{h \in H} y_i (\cdot)_i ((A)_{x_i, y_i}) h(x_i) \cdot(h) \\ &= \frac{1}{m} \sum_{h \in H} \sum_{i=1}^n y_i (\cdot)_i ((A)_{x_i, y_i}) h(x_i) \cdot(h) \end{aligned}$$

where Fubini's theorem justifies swapping the final two summations, which is easily valid since $\sum_{h \in H} |h(x)| < 1$ and $\sup_{h \in H} |h(x)| < 1$. Specializing this expression for $\cdot \in \mathbb{R}^n$ is direct: at

any coordinate e_h , by the above rule for $r^{\mathbb{L}}$ and properties of the adjoint $A^>$,

$$D_{r^{\mathbb{L}}(A)}(h); e_h = A^> r^{\mathbb{L}}(A); e_h = r^{\mathbb{L}}(A); h = \frac{1}{m} \sum_{i=1}^n y_i^{-1} ((A)_{x_i, y_i}) h(x_i):$$

6. By Lemma 1.B.1, ψ is strictly convex and continuously differentiable over the interior of its domain. Thus, given $\lambda \in (0, 1)$ and $p_1, p_2 \in \text{dom}(\mathbb{L})$ with $p_1 \neq p_2$, there must be at least one example $f(x_i; y_i) \in p_2(x_i; y_i)$, whereby the above conjugacy relation grants

$$\begin{aligned} \mathbb{L}(\lambda p_1 + (1 - \lambda)p_2) &= \frac{1}{m} \sum_{i=1}^n (\lambda (mp_1(x_i; y_i)) + (1 - \lambda)mp_2(x_i; y_i)) \\ &< \frac{1}{m} \sum_{i=1}^n (\lambda (mp_1(x_i; y_i)) + (1 - \lambda)(mp_2(x_i; y_i))) \\ &= \mathbb{L}(p_1) + (1 - \lambda)\mathbb{L}(p_2): \end{aligned}$$

Next, given some $f \in L^1(\mu)$, the same derivation as for $r^{\mathbb{L}}$ provides that $r^{\mathbb{L}}(f)_{x_i, y_i} = r^{\mathbb{L}}(f(x_i; y_i)) = m$, whereby $r^{\mathbb{L}}$ is continuously differentiable since Lemma 1.B.1 granted that ψ is continuously differentiable.

□

Notice that this appendix makes no assertions regarding \mathbb{L} . While the above results for ψ and \mathbb{L} may seem remedial, the analogous properties for \mathbb{L} are not simply more difficult to prove due to technicalities (cf. Lemma 4.1.2), but even basic properties such as finiteness of \mathbb{L} are false in the generality stated above for \mathbb{L} (cf. Proposition 4.1.3). Note that another way to obtain fairly general guarantees is to relax the finiteness condition on the measure, but instead require μ to be finite (cf. Lemma 5.A.2). Consequently, things only break down (for the choice of Banach spaces here) in the case that both the measure over $X \times Y$ and the measure over H have infinite support.

Lastly, the weak duality result may be proved as follows.

Proof of Proposition 1.4.1. To start, since ψ is convex and finite, it is continuous and thus measurable. Moreover, given any $p \in L^1(\mu)$ (whereby ψ is measurable), since ψ is bounded below by Lemma 1.B.1 and since μ is a finite measure, $\int_{\{\psi < 0\}} \psi d\mu < \infty$. Consequently, only the positive content of ψ can lead to an infinite integral, thus $\int \psi d\mu$ is well-defined.

Next, for any $f \in L^1(\Omega)$ and any $p \in L^1(\Omega)$, the definition of Fenchel conjugate provides

$$p(z)f(z) = \int_{\Omega} (f(z) - p(z)) dz:$$

Consequently, for any $p \in L^1(\Omega)$,

$$\int_{\Omega} p(z) f(z) dz = \int_{\Omega} p(z) f(z) dz - \int_{\Omega} p(z) f(z) dz + \int_{\Omega} p(z) f(z) dz = \int_{\Omega} p(z) f(z) dz - \int_{\Omega} p(z) f(z) dz + \int_{\Omega} p(z) f(z) dz; \quad (1.B.3)$$

where $\int_{\Omega} p$ is measurable and the corresponding integral is well-defined as verified above.

Now let χ_A denote the indicator function over A , meaning it is 0 when its argument resides in A^c , otherwise it is infinite. By direct inspection, its conjugate for any $q \in L^1(\Omega)$ is

$$(\chi_A)^*(q) = \sup_{h \in L^1(\Omega)} \int_{\Omega} h(z) q(z) dz = \int_{\Omega} q(z) dz \chi_A(q): \quad (1.B.4)$$

For any $\chi_A \in L^1(\Omega)$ and $p \in L^1(\Omega)$, the definition of Fenchel conjugate provides

$$\begin{aligned} \int_{\Omega} p(z) \chi_A(z) dz &= \int_{\Omega} p(z) \chi_A(z) dz - \int_{\Omega} p(z) \chi_A(z) dz + \int_{\Omega} p(z) \chi_A(z) dz; \\ \int_{\Omega} p(z) \chi_A(z) dz &= \int_{\Omega} p(z) \chi_A(z) dz - \int_{\Omega} p(z) \chi_A(z) dz + \int_{\Omega} p(z) \chi_A(z) dz; \end{aligned}$$

Combining this with eq. (1.B.3) and eq. (1.B.4),

$$\begin{aligned} \int_{\Omega} p(z) \chi_A(z) dz + \int_{\Omega} p(z) \chi_A(z) dz - \int_{\Omega} p(z) \chi_A(z) dz &= \int_{\Omega} p(z) \chi_A(z) dz + \int_{\Omega} p(z) \chi_A(z) dz - \int_{\Omega} p(z) \chi_A(z) dz \\ &= \int_{\Omega} p(z) \chi_A(z) dz + \int_{\Omega} p(z) \chi_A(z) dz - \int_{\Omega} p(z) \chi_A(z) dz \\ &= 0: \end{aligned}$$

Rearranging,

$$\int_{\Omega} p(z) \chi_A(z) dz + \int_{\Omega} p(z) \chi_A(z) dz = \int_{\Omega} p(z) \chi_A(z) dz:$$

Since $\chi_A \in L^1(\Omega)$ and $p \in L^1(\Omega)$ were arbitrary, $\sup_{p \in L^1(\Omega)}$ may be applied to the right hand side

and $\inf_{\lambda \geq 0} \mathcal{L}(\lambda)$ may be applied to the left hand side while preserving the inequality, meaning

$$\inf_{\lambda \geq 0} \mathcal{L}(\lambda) = \sup_{p \in \mathcal{P}} \mathcal{L}^*(p) \quad \text{where } \mathcal{P} = \{p \in \mathbb{R}^d : p \geq 0\}$$

In order to obtain the desired statement from here, note first that $\mathcal{L}^*(0) = 0$ for $\lambda \geq 0$, and thus the expression may be dropped. Second, to adjust the dual, first note that $\mathcal{L}^*(0) = 0$ by Lemma 1.B.1, thus 0 is always a dual feasible point (with value 0), meaning constraints can be added to the dual formulation so long as the value is not forced below zero. To start, note by Lemma 1.B.1 that $\mathcal{L}^*(r) = 1$ when $r < 0$, thus the condition $p \geq 0$ -a.e. may be safely added, since $\{[p < 0]\} \cap \{p \geq 0\} = \emptyset$ entails a dual objective value of 1. Lastly, the statement $\mathcal{L}^*(A) = 0$ for all A is just a fancy way of writing $\inf_{\lambda \geq 0} \mathcal{L}(\lambda) < 1$. \square

Appendix 1.C Line Search Properties

This section provides basic properties of the three line search choices, most importantly a lower bound on the amount they decreased the (primal) objective function in each iteration.

These line searches depend on Lipschitz properties of gradients; concretely, since Boost is coordinate descent applied to \mathcal{L}^b , the gradient mapping $r(\mathcal{L}^b)$ will need to be Lipschitz with respect to the $L^1(b)$ norm over a set S , meaning there exists some B_2 so that for any $\lambda, \lambda' \geq 0$,

$$\|r(\mathcal{L}^b)(\lambda) - r(\mathcal{L}^b)(\lambda')\|_1 \leq B_2 |\lambda - \lambda'| \quad (1.C.1)$$

An inequality of this type will be proved separately in Chapter 3 and in Chapter 6, but each time will make use of the following lemma which takes care of the first few steps.

Lemma 1.C.2. Let convex differentiable $\mathcal{L} : \mathbb{R} \rightarrow \mathbb{R}_+$ be given, along with empirical measure \mathbb{P} over $X = \{(x_i, y_i)\}_{i=1}^n$. Then for any $\lambda, \lambda' \geq 0$,

$$\mathcal{L}(\lambda) - \mathcal{L}(\lambda') \leq \frac{1}{m} \sum_{i=1}^n [y_i \lambda - y_i \lambda' + \mathcal{L}(\lambda) - \mathcal{L}(\lambda')] = \mathcal{L}(\lambda) - \mathcal{L}(\lambda') + \frac{1}{m} \sum_{i=1}^n (y_i \lambda - y_i \lambda')$$

Proof. By the form of $r(\mathbb{L} A)$ from Proposition 1.B.2 and the form of $kA^> k_1$ from Lemma 1.A.5,

$$\begin{aligned}
 & r(\mathbb{L} A)(\cdot) - r(\mathbb{L} A)(\theta_1) \\
 = & \sup_{\substack{X \\ \frac{1}{m} \sum_{i=1}^n X_i}} \left(\sum_{h \in \mathcal{H}} y_i(\cdot^0((A)_{x_i; y_i}) - \cdot^0((A^0)_{x_i; y_i})) h(x_i) v(h) : v \in \mathcal{V}^2; kvk_1 \leq 1 \right) \\
 & \sup_{\substack{X \\ \frac{1}{m} \sum_{i=1}^n X_i}} \left(\sum_{h \in \mathcal{H}} j(\cdot^0((A)_{x_i; y_i}) - \cdot^0((A^0)_{x_i; y_i})) j(y_i, h(x_i)) v(h) : v \in \mathcal{V}^2; kvk_1 \leq 1 \right) \\
 & \frac{1}{m} \sum_{i=1}^n j(\cdot^0((A)_{x_i; y_i}) - \cdot^0((A^0)_{x_i; y_i}))
 \end{aligned}$$

as desired. □

1.C.1 Line Search Options 1 and 2

The following guarantee holds for the unconstrained line search (option 1), and the quadratic upper bound line search (option 2), and moreover provides a derivation of the latter choice.

Lemma 1.C.3. Let hypothesis class \mathcal{H} , empirical measure b , integer t denoting an iteration, and convex differentiable loss $\cdot : \mathbb{R} \rightarrow \mathbb{R}_+$ be given. Suppose $\mathbb{L} A$ has Lipschitz gradients with parameter B_t with respect to norm $L^1(b)$ over the t^{th} level set

$$S_t := \left\{ \theta : \mathbb{L}(A) - \mathbb{L}(A_{t-1}) \leq \theta \right\}$$

(Please see eq(1.C.1) for more on this definition.) Suppose step size η_t is chosen according to one of the first two step choices in Figure 1.3, meaning either $\eta_t = \eta_t$ or $\eta_t = \frac{2}{r} \frac{\mathbb{L}(A_{t-1}) - Av_t}{B_t}$; $\eta_t \leq \eta_t$. Then

$$\begin{aligned}
 & \eta_t \frac{kA^> r \mathbb{L}(A_{t-1}) k_1}{B_t}; \\
 \mathbb{L}(A_{t-1}) - \mathbb{L}(A_{t-1}) & \leq \frac{2kA^> r \mathbb{L}(A_{t-1}) k_1^2}{2B_t}.
 \end{aligned}$$

Proof. Consider any $\theta \in S_t$ so that $\theta := \eta_t + v_t \in S_t$. By the form of $r(\mathbb{L} A)$ in Proposition 1.B.2,

$$\int_0^{\eta_t} r(\mathbb{L}(A_{t-1} + r(A - A_{t-1}))) dr = \int_0^{\eta_t} \mathbb{L}(A_{t-1} + r(A - A_{t-1})) dr = \mathbb{L}(A) - \mathbb{L}(A_{t-1});$$

As such, by the Lipschitz gradient property (which holds for every element between t_{i-1} and t_i),

$$\begin{aligned}
 \mathbb{E} \mathbb{L}(A(t_{i-1} + v_t)) &= \mathbb{E} \mathbb{L}(A(t_{i-1})) + \mathbb{E} \int_{t_{i-1}}^{t_i} \nabla \mathbb{L}(A(t_{i-1} + r(A - A(t_{i-1}))) \nabla \mathbb{L}(A(t_{i-1}); v_t) dr \\
 &\quad + \mathbb{E} \int_{t_{i-1}}^{t_i} \frac{1}{2} \nabla^2 \mathbb{L}(A(t_{i-1} + r(A - A(t_{i-1}))) \nabla^2 \mathbb{L}(A(t_{i-1}); v_t) k v_t k_1^2 dr \\
 &= \mathbb{E} \mathbb{L}(A(t_{i-1})) + \mathbb{E} \int_{t_{i-1}}^{t_i} \nabla \mathbb{L}(A(t_{i-1} + r(A - A(t_{i-1}))) \nabla \mathbb{L}(A(t_{i-1}); v_t) dr \\
 &\quad + \mathbb{E} \int_{t_{i-1}}^{t_i} \frac{1}{2} \nabla^2 \mathbb{L}(A(t_{i-1} + r(A - A(t_{i-1}))) \nabla^2 \mathbb{L}(A(t_{i-1}); v_t) k v_t k_1^2 dr \\
 &= \mathbb{E} \mathbb{L}(A(t_{i-1})) + \mathbb{E} \int_{t_{i-1}}^{t_i} \nabla \mathbb{L}(A(t_{i-1} + r(A - A(t_{i-1}))) \nabla \mathbb{L}(A(t_{i-1}); v_t) dr \\
 &\quad + \frac{1}{2} \mathbb{E} \int_{t_{i-1}}^{t_i} \nabla^2 \mathbb{L}(A(t_{i-1} + r(A - A(t_{i-1}))) \nabla^2 \mathbb{L}(A(t_{i-1}); v_t) k v_t k_1^2 dr
 \end{aligned}$$

This final expression defines a univariate quadratic with minimum $\mathbb{L}(A(t_{i-1} + v_t)) = \mathbb{L}(A(t_{i-1})) + \frac{1}{2} \mathbb{E} \int_{t_{i-1}}^{t_i} \nabla^2 \mathbb{L}(A(t_{i-1} + r(A - A(t_{i-1}))) \nabla^2 \mathbb{L}(A(t_{i-1}); v_t) k v_t k_1^2 dr$ (where necessarily $t_{i-1} + v_t \in S_t$). This function has slopes everywhere exceeding $\frac{1}{2} \mathbb{E} \int_{t_{i-1}}^{t_i} \nabla^2 \mathbb{L}(A(t_{i-1} + r(A - A(t_{i-1}))) \nabla^2 \mathbb{L}(A(t_{i-1}); v_t) k v_t k_1^2 dr$ along $[t_{i-1}; t_i]$ (for either choice of step size), and so $t_{i-1} \leq t_i$. (Indeed, these bounds give a derivation for the second step size choices.) To get the second guarantee, note that plugging into the above quadratic and simplifying via

$$\mathbb{E} \int_{t_{i-1}}^{t_i} \nabla^2 \mathbb{L}(A(t_{i-1} + r(A - A(t_{i-1}))) \nabla^2 \mathbb{L}(A(t_{i-1}); v_t) k v_t k_1^2 dr$$

gives the desired minimum quadratic upper bound. □

1.C.2 Line Search Option 3: the Wolfe Search

The Wolfe line search is required to output a step size satisfying two conditions; after providing this basic definition and an analog to the single step guarantee established in Lemma 1.C.3 for the other two line searches, this subsection then closes by providing an efficient means for computing a step size satisfying the Wolfe conditions.

Definition 1.C.4 (Wolfe line search). The Wolfe line search chooses any t_i which satisfies the

following conditions (where this thesis makes the simple choice $c_1 = 1/3$ and $c_2 = 1/2$):

$$\begin{aligned} \mathbb{E} \ell(A_{t-1} + v_t) &\leq \mathbb{E} \ell(A_{t-1}) + c_1 \mathbb{E} \langle \nabla \ell(A_{t-1}), v_t \rangle; \\ \mathbb{E} \ell(A_{t-1}) &\leq \frac{1}{3} k_A \mathbb{E} \|v_t\|; \end{aligned} \tag{1.C.5}$$

$$\begin{aligned} \mathbb{E} \ell(A_{t-1} + v_t) &\leq \mathbb{E} \ell(A_{t-1}) + c_2 \mathbb{E} \langle \nabla \ell(A_{t-1}), v_t \rangle; \\ \mathbb{E} \langle \nabla \ell(A_{t-1}), v_t \rangle &\leq \frac{1}{2} k_A \mathbb{E} \|v_t\|. \end{aligned} \tag{1.C.6}$$

This method admits an efficient implementation (cf. Figure 1.4).

The basic guarantee of the Wolfe search is as follows; due to its similarity Lemma 1.C.3, it will be possible to invoke these single-step guarantees, and to then analyze the different step sizes simultaneously.

Lemma 1.C.7. Let hypothesis class \mathcal{H} , empirical measure b , integer t denoting an iteration, and convex differentiable loss $\ell : \mathbb{R}^d \rightarrow \mathbb{R}_+$ be given. Suppose ℓ has Lipschitz gradients with parameter B_t with respect to norm $L^1(b)$ over the t^{th} level set

$$S_t := \{A \in \mathbb{R}^d : \ell(A) \leq \ell(A_{t-1})\}.$$

(Please see eq(1.C.1) for more on this definition.) Suppose step size η_t satisfies the Wolfe conditions for some $0 < c_1 < c_2 < 1$. Then

$$\begin{aligned} \eta_t &\leq \frac{(1 - c_2) k_A \mathbb{E} \langle \nabla \ell(A_{t-1}), v_t \rangle}{B_t}; \\ \ell(A_t) - \ell(A_{t-1}) &\leq \frac{c_1 (1 - c_2) k_A \mathbb{E} \langle \nabla \ell(A_{t-1}), v_t \rangle^2}{B_t}. \end{aligned}$$

Proof. By the first Wolfe condition (cf. eq. (1.C.5)), the step size η_t is a descent step, meaning $A_t = A_{t-1} + \eta_t v_t \in S_t$, and thus $\langle \nabla \ell(A_{t-1}), v_t \rangle \leq 0$. By the definition of B_t ,

$$\begin{aligned} \mathbb{E} \langle \nabla \ell(A_{t-1}), v_t \rangle &\leq \mathbb{E} \langle \nabla \ell(A_{t-1}), v_t \rangle; \\ \mathbb{E} \langle \nabla \ell(A_{t-1}), v_t \rangle &\leq \mathbb{E} \langle \nabla \ell(A_{t-1}), v_t \rangle; \\ \mathbb{E} \langle \nabla \ell(A_{t-1}), v_t \rangle &\leq \mathbb{E} \langle \nabla \ell(A_{t-1}), v_t \rangle; \\ \mathbb{E} \langle \nabla \ell(A_{t-1}), v_t \rangle &\leq \mathbb{E} \langle \nabla \ell(A_{t-1}), v_t \rangle; \\ \mathbb{E} \langle \nabla \ell(A_{t-1}), v_t \rangle &\leq \mathbb{E} \langle \nabla \ell(A_{t-1}), v_t \rangle; \end{aligned}$$

The rest of the proof is just as for standard Wolfe search guarantees and direct from the Wolfe conditions (Nocedal and Wright, 2006, Theorem 3.2). First, subtracting $\frac{D}{A^T r} \mathbb{L}(A^T r; v_t)$ from both sides of eq. (1.C.6) gives

$$\frac{D}{A^T r} \mathbb{L}(A^T r; v_t) - \frac{E}{A^T r} \mathbb{L}(A^T r; v_t) = (c_2 - 1) \frac{D}{A^T r} \mathbb{L}(A^T r; v_t) - \frac{E}{A^T r} \mathbb{L}(A^T r; v_t);$$

which can be combined with the above derivation to yield

$$\frac{(c_2 - 1) \frac{D}{A^T r} \mathbb{L}(A^T r; v_t) - \frac{E}{A^T r} \mathbb{L}(A^T r; v_t)}{B_t} = \frac{(1 - c_2) k A^T r \mathbb{L}(A^T r; v_t) k_1}{B_t}$$

Plugging this into eq. (1.C.5) gives

$$\mathbb{L}(A^T r; v_t) - \mathbb{L}(A^T r; v_{t-1}) = \frac{2c_1(1 - c_2) k A^T r \mathbb{L}(A^T r; v_t) k_1^2}{B_t};$$

□

For whatever reason, the Wolfe search is less popular than the backtracking line search, even though the latter requires a guess for the upper bound on the Lipschitz parameter; by contrast, the Wolfe search needs no parameters. Consequently, this thesis elects to explain and investigate the method in a little more generality than is necessary.

Consider any convex differentiable function h over some finite dimensional Euclidean space, a current iterate x , and a descent direction v (that is, $r h(x)^T v < 0$). By convexity, the linearization of h at x in direction v , symbolically $h(x) + r h(x)^T v$, will lie below the function. But, by continuity, it must be the case that, for any $c_1 \in (0; 1)$, the ray $h(x) + c_1 r h(x)^T v$, depicted in Figure 1.5, must lie above h for some small region around x ; this gives the first Wolfe condition (presented earlier in eq. (1.C.5)), also known as the Armijo condition (cf. Nocedal and Wright (2006, Equation 3.4) and Bertsekas (1999, Exercise 1.2.16)):

$$h(x + v) \geq h(x) + c_1 r h(x)^T v; \tag{1.C.8}$$

Unfortunately, this rule may grant only very limited decrease in objective value, since $\epsilon > 0$ can be chosen arbitrarily small and still satisfy the rule; thus, the second Wolfe condition (presented earlier in eq. (1.C.6)), also called a curvature condition, which depends on $c_2 \in (c_1; 1)$, forces the

Routine. Wolfe .
 Input. Convex function h , iterate x , descent direction v .
 Output. Step size $\alpha > 0$ satisfying (1.C.8) and (1.C.9).

1. Bracketing step.
 - (a) Set $\alpha_{\max} := 1$.
 - (b) While α_{\max} satisfies (1.C.8):
 - Set $\alpha_{\max} := 2 \alpha_{\max}$.
2. Bisection step.
 - (a) Set $\alpha_{\min} := 0$ and $\alpha := \alpha_{\max} = 2$.
 - (b) While α does not satisfy (1.C.8) and (1.C.9):
 - i. If α violates (1.C.8):
 - Set $\alpha_{\max} := \alpha$.
 - ii. Else, α must violate (1.C.9):
 - Set $\alpha_{\min} := \alpha$.
 - iii. Set $\alpha := (\alpha_{\min} + \alpha_{\max})/2$.
 - (c) Return α .

Figure 1.4. Bracketing and bisection search for step size satisfying Wolfe conditions.

step to be farther away:

$$\alpha \nabla h(x + \alpha v) \cdot v \geq c_2 \alpha \nabla h(x) \cdot v \quad (1.C.9)$$

This requires the new gradient (in direction v) to be closer to 0, mimicking first order optimality conditions for the exact line search. Note that the new gradient (in direction v) may in fact be positive; this does not affect the analysis.

An algorithm to find a point satisfying these conditions, presented in Figure 1.4, is simple enough: grow α as quickly as possible, and then bisect backwards for a satisfactory point. As compared with the presentation in Nocedal and Wright (2006, Algorithm 3.5), α_{\max} is searched for rather than provided, and convexity removes the need for interpolation.

Proposition 1.C.10. Given a continuously differentiable convex bounded below function h , iterate x , and direction v , Wolfe terminates with an $\alpha > 0$ satisfying (1.C.8) and (1.C.9).

Proof. The bracketing search must terminate: v is a descent direction, so the linearization at x with slope $\nabla h(x) \cdot v$ will eventually intersect h (since h is bounded below).

The remainder of this proof is illustrated in Figure 1.5. Let α_1 be the greatest positive

Appendix 1.D Existence of Hard Cores

In order to prove Theorem 1.7.4, this section first proves the countable additivity of $S_D(H; \cdot)$, and thereafter proves existence via an optimization on the size of elements in $S_D(H; \cdot)$ (a similar technique can be used as the first step of the proof of the Hahn Decomposition theorem (Folland, 1999, Theorem 3.3)).

Lemma 1.D.1. Given any $f, c_i g_{i=1}^1$ with $c_i \geq 0$ and $f, p_i g_{i=1}^1$ with $p_i \in D(H; \cdot)$ and $\sum_{i=1}^{\infty} c_i k_{p_i} < 1$, the limit object $p_1 := \sum_{i=1}^{\infty} c_i p_i$ exists, and satisfies $p_1 \in D(H; \cdot)$.

Proof. Let $f, c_i g_{i=1}^1$ and $f, p_i g_{i=1}^1$ be given as specified. First, by the monotone convergence theorem and since $\sum_{i=1}^{\infty} c_i k_{p_i} < 1$, the function $p_1 := \sum_{i=1}^{\infty} c_i p_i$ exists (i.e., all limits converge pointwise), is measurable, and satisfies $\sum_{i=1}^{\infty} c_i p_i < 1$, meaning $p_1 \in L^1(\cdot)$ (Folland, 1999, Theorem 2.15). Moreover, $\sum_{i=1}^{\infty} c_i k_{p_i} < 1$, thus $p_1 \in L^1(\cdot)$ as well. Now let any $\mu \in \mathbb{R}^n$ be given; note that $\sum_{i=1}^{\infty} c_i \int p_i(H) d\mu < 1$. Thanks to this, by the dominated convergence theorem (Folland, 1999, Theorem 2.25),

$$\begin{aligned} \int p_1(x; y) y(H) d\mu(x; y) &= \int \sum_{i=1}^{\infty} c_i p_i(x; y) y(H) d\mu(x; y) \\ &= \sum_{i=1}^{\infty} \int c_i p_i(x; y) y(H) d\mu(x; y) \\ &= \sum_{i=1}^{\infty} c_i \int p_i(x; y) y(H) d\mu(x; y) \\ &= 0. \end{aligned} \quad \square$$

Lemma 1.D.2. $S_D(H; \cdot)$ is closed under countable unions.

Proof. Let any collection $f, C_i g_{i=1}^1$ with $C_i \in S_D(H; \cdot)$ and corresponding weighting $p_i \in D(H; \cdot)$ be given. Define

$$C := \bigcup_{i=1}^{\infty} C_i \quad \text{and} \quad p := \sum_{i=1}^{\infty} \frac{p_i}{2^i \max\{1, k_{p_i}\}}.$$

By Lemma 1.D.1, p is well-defined and satisfies $p \in D(H; \cdot)$. Note further that $C = \{p > 0\}$, and thus $C \in S_D(H; \cdot)$. □

Proof of Theorem 1.7.4. Consider the optimization problem

$$d := \sup \{ \int (C) : C \in S_D(H; \gamma) \}$$

Since S_D is nonempty (always contains; corresponding top = $0 \leq D(H; \gamma)$) and $(X f_{i=1}^1 + 1 g) < 1$, the supremum is finite. Let $\{C_i\}_{i=1}^{\infty}$ be a maximizing sequence, and define $D_j := \int_{i=1}^j C_i$ and $D := \int_{i=1}^{\infty} D_j = \int_{i=1}^{\infty} C_i$. By Lemma 1.D.2, $D_j \in S_D(H; \gamma)$ for every j , and since $(D_j) \leq (C_j)$, it follows that $\int D_j$ must also be a maximizing sequence to the above supremum. Finally, since Lemma 1.D.2 also grants $D \in S_D(H; \gamma)$, then by continuity of measures from below (Folland, 1999, Theorem 1.8(c)),

$$\int (D) = \lim_{j \rightarrow \infty} \int (D_j) = d$$

Since $D \in S_D(H; \gamma)$ attains the supremum, it is a hard core; in particular, hard cores exist.

If there are two hard cores D_1 and D_2 , then the definition provides that each satisfies $D_i \in S_D(H; \gamma)$, and $\int (D_2 \wedge D_1) = 0$ and $\int (D_1 \wedge D_2) = 0$, meaning the symmetric difference has measure zero, which gives the second property. \square

Appendix 1.E Bibliographic Notes

As discussed in Section 1.1, the original boosting algorithm is due to Schapire (1990), who also gave the basic guarantees closing the question posed by Kearns and Vazirani (1994). The complexity of the circuit produced by this algorithm was reduced to a single majority by Freund (1995), and thereafter both authors together developed AdaBoost, which is the topic of study in this thesis (Freund and Schapire, 1997). Most recently, Schapire and Freund have also put out an excellent, comprehensive book on many aspects of AdaBoost (Schapire and Freund, 2012). It is noteworthy that while boosting variants abound, practical boosting implementations are still very close to the original (Caruana and Niculescu-Mizil, 2006).

The algorithm was not originally phrased as minimization of a convex loss; this view appears to have been first presented by (Breiman, 1999). Subsequently, the algorithm was studied for many alternate losses, primarily the logistic loss, by many authors Friedman et al. (2000), Friedman (2000), Lafferty (1999), Mason et al. (2000), Collins et al. (2002), Du y and Helmbold (2000).

Occasionally, authors express the view that the loss minimization view of boosting is somehow incomplete or inaccurate. One reason for this is that the primal optimization problem given in Equation (1.3.4) does not place any restrictions on the complexity of the learned function; since H is typically chosen to be some large VC class, its linear span $\text{span}(H)$ typically has infinite VC dimension, which means it is statistically unstable (Devroye et al., 1996, Theorem 14.3). Consequently, if the objective function is not modified in some way, statistical stability must be achieved via some aspect of the minimization algorithm. Boosting in turn accomplishes this in two ways: first via coordinate descent, which means that the t^{th} iterated is supported on at most t coordinates, and secondly via some sort of stopping rule.

Another reason why focusing on the primal objective function is perhaps distasteful is that it omits any mention of the weights over examples and the basic computational model of boosting. This thesis aims to point out that the convex dual highlights that the weighting structure was present in the convex formulation all along.

It was arguably a great boon that the original authors did not formulate the algorithm with loss minimization in mind. Had that been the case, it would have been more expedient to analyze the scheme under the settings more common for convex optimization (strong convexity, existence of minima, etc). As discussed throughout this chapter, and highlighted in the problem structure in Chapter 2, the typical convex case, and the original case of boosting ($\gamma > 0$), are very different. As such, thanks to the original algorithm not being phrased as loss minimization, the world now has these two very different analyses, which can be understood via hard cores to rather fundamentally partition the space of possibilities.

Regarding hard cores, the term itself is from cryptography (Goldreich and Levin, 1989). Closer to the spirit here, Impagliazzo (1995) showed via a "hard core lemma" that the existence of functions which are somewhat hard for circuits of a given size implies the existence of functions which are very hard for circuits of a smaller size. Similarly, the hard core C of an instance here is a subset of the inputs which is very difficult, and captures the basic difficulty of the problem (in Chapter 3 and Chapter 6, when the hard core is empty, the instance is easy). Interestingly, the only existing proofs of fast numerical convergence of AdaBoost to $\epsilon(A)$ use hard cores; one example is Chapter 3, and the other is the proof due to Mukherjee et al. (2011, See "zero loss set", the complement of the hard core). The hard core may also be used to develop a theory of margins when separability does not hold (Telgarsky, 2013).

Acknowledgements

This chapter is based on work by the dissertation author which appeared in the Journal of Machine Learning Research, Volume 13, pages 561-606, 2012, as well as the Conference on Learning Theory, 2013, and lastly work under preparation.

Part I

Just a Finite Sample

Chapter 2

Problem Structure

This chapter develops the structure of the primal optimization problem over finite samples (e.g., eq. (1.3.3) with $\mu = b$), which is then used in Chapter 3 to produce numerical convergence rates. Concretely, these two chapters operate under the following two assumptions.

1. Only the empirical measure $\hat{\mu}$ corresponding to a finite sample $\{(x_i; y_i)\}_{i=1}^m$ is considered. The material in Chapter 5 may be used to produce statistical guarantees (i.e., to infer behavior over \mathcal{X} from the behavior over b).
2. H must satisfy the following notion of bounded complexity: the set of images $\{h(x_i)\}_{i=1}^m : h \in H$ has some finite size $< \infty$. This condition is satisfied if for instance the set of possible images $\{h(x) : x \in \mathcal{X}; h \in H\}$ has some finite size $< \infty$ (whereby the above image set has size at most k^m), or for instance when H has finite VC dimension $V(H) < \infty$ (whereby the Sauer-Shelah lemma grants that the size is eventually $\leq (m+1)^{V(H)}$).

These two assumptions together entail that the linear operators H and A may be represented as matrices in $\mathbb{R}^{m \times n}$ without affecting any of the mathematical development; indeed, to succinctly state these two assumptions, the results will require $A \in \mathbb{R}^{m \times n}$. In this setting, $L^1(b)$ is isometrically isomorphic to \mathbb{R}^m with the l_1 norm; thus for any vector $\alpha \in \mathbb{R}^m$, the object $A \in \mathbb{R}^{m \times n}$ is simultaneously considered as an element of \mathbb{R}^m , and also an element of $L^1(b)$ via isometric isomorphism. Similarly, note the simplified notation $(A)_i := (A)_{x_i; y_i}$.

2.1 Overview

This chapter will study the basic structure of the optimization problem for \mathcal{D} under the condition $A \in \mathbb{R}^{m \times n}$. All results will hold for the following family of loss functions.

Definition 2.1.1. Let L_{fs} contain all functions $\psi : \mathbb{R} \rightarrow \mathbb{R}_{++}$ which are strictly convex, continuously differentiable, and satisfy $\lim_{z \rightarrow 1} \psi(z) = 0$.

The first step, in Section 2.2, is to provide a strong duality analog to Proposition 1.4.1, but now under the assumption of a finite sample. This result will be proved under slightly more generality than stated above; in particular, the measure will be b , but μ may be infinite. Even so, focusing on b means that only \mathcal{D} and \mathcal{D} must be considered, which were shown to be well behaved in Proposition 1.B.2. Moreover, the dual space, rather than being $L^1(\mu)$ (which is a bad choice for the exponential loss, as discussed in Chapter 4), is simply \mathbb{R}^m (where all metric topologies agree since the space has finite dimension).

Section 2.2 will also flesh out the duality-based stopping condition in Figure 1.3, which is simplified via a few more definitions. First note the following kernel and image definitions for A and $A^>$, which are stated in full generality, but will only be used for $A \in \mathbb{R}^{m \times n}$.

Definition 2.1.2. Given a linear operator $B : X \rightarrow \mathbb{R}$, define its image and kernel as

$$\text{im}(B) := \{ Bx : x \in X \} \quad \text{and} \quad \ker(B) := \{ x \in X : Bx = 0 \}$$

Notice for instance that signed decorrelating weightings are simply $\ker(A^>)$; the following definition in turn ensures nonnegativity and simplifies many future expressions.

Definition 2.1.3. Given operator $A : L^1(\mu) \rightarrow \mathbb{R}$, and identifying $\ker(A)$ with a subset of $L^1(\mu)$ via isometric isomorphism as per Lemma 1.A.5), define $D_A := \ker(A) \cap \{ p \in L^1(\mu) : p \geq 0 \}$ -a.e.g.

Notice that this expression contains most of the domain of the dual in Proposition 1.4.1; when $A \in \mathbb{R}^{m \times n}$, the earlier tacit isometric isomorphisms between $L^1(\mu)$ and \mathbb{R}^m allow D_A to be identified with $\ker(A^>) \cap \mathbb{R}_+^m$. Note however that \mathcal{D} is not finite everywhere over D_A ; when $\psi \in L_{fs}$ is Lipschitz with parameter L , then \mathcal{D} is only finite over the cube $[0, L^{-1}]^m$ (cf. Proposition 1.B.2). On the other hand, given $p \in D_A \cap \{0\}$, the rescaling $\tilde{p} := p / (\|p\|_1)$ always satisfies $\tilde{p} \in \ker(A^>)$ and $\mathcal{D}(\tilde{p}) < 1$.

The following notation for projections and distances will also be used in the duality-based stopping condition.

Definition 2.1.4. Let $C \subseteq \mathbb{R}^m$ be a closed convex set, and define the following projection and distance shorthands, where $p \in [1; \infty]$, and projections in the case of non-uniqueness (for $p = 1$ and $p = \infty$ norms) make some arbitrary consistent choice:

$$P_C^p(x) = \arg \min_{y \in C} \|x - y\|_p; \quad \text{and} \quad D_C^p(x) = \|x - P_C^p(x)\|_p;$$

Excluding appendices, the chapter has two more sections. First, Section 2.4 pins down the structure of hard cores under the assumption $A \in \mathbb{R}^{m \times n}$. Second, Section 2.3 provides a generalization of the quantity β which is positive under conditions which will hold for the optimization problem in general, and thus the modified quantity will appear in all convergence rates of Chapter 3, and not simply those corresponding to the case $\beta > 0$.

2.2 Strong Duality

Theorem 2.2.1. Let empirical probability measure b , hypothesis class \mathcal{H} (finite or infinite), and loss $\ell : \mathbb{R} \rightarrow \mathbb{R}_+$ be given with $\lim_{z \rightarrow 1} \ell(z) = 0$. Then

$$\inf_{\lambda} \int \ell(A \lambda) db : \lambda \in \mathcal{D} = \max_{p \in \mathcal{D}_A} \int \ell(p) db : \lambda \in \mathcal{D}_A; \quad (2.2.2)$$

and a pair $(\lambda^* ; p^*) \in \mathcal{L}^1(b)$ is optimal if $p^* \in \mathcal{D}_A$ and $p^* \in \mathcal{D}_A$ everywhere. Furthermore, if ℓ is differentiable and strictly convex, then the dual optimum $\lambda^* \in \mathcal{D}_A$ is unique, and λ^* is optimal if $\lambda^* \in \mathcal{D}_A = \mathcal{D}_A$ everywhere.

Since the structure of \mathcal{D} and \mathcal{D}_A was already handled by Proposition 1.B.2, the proof is little more than an invocation of an appropriate Fenchel duality rule and some remedial manipulations. Additionally, recall, as discussed in Remark 1.A.4, that the appearance of λ in the dual (which was absent in the weak duality formula from Proposition 1.4.1) is an artifact of using inconsistent dual pairings for \mathcal{D} and \mathcal{L} .

Special notation is used for the dual optimum λ^*_A in this chapter and Chapter 3, since $A \in \mathbb{R}^{m \times n}$ will be shown to entail very strict properties. As such, use the notation as follows.

Definition 2.2.3. When $A \in \mathbb{R}^{m \times n}$ and $\ell \in \mathcal{L}_{fs}$, let λ^*_A denote the dual optimum as in Theorem 2.2.1.

Proof of Theorem 2.2.1. Consider the two Fenchel problems

$$V_p := \inf_{\mathbb{Z}} \int (A^T p) db + \psi(p);$$

$$V_d := \sup_p \int (mp) db - \psi^*(p);$$

where ψ denotes a convex indicator function (0 on S , $+\infty$ on S^c), and the conjugacy relation $(\int (A^T p) db)^* = \int (mp) db - \psi^*(p)$ was provided by Proposition 1.B.2.

In order to show $V_p = V_d$, since the dual space $L^1(b)$ is finite dimensional, it suffices to apply a variant of Slater's condition making use of quasi relative interiors (Zalinescu, 2002, Corollary 2.8.5 using condition (x)); namely, the relative interior of $\text{dom}(\int (A^T p) db)$ is simply $L^1(b)$ as provided by Proposition 1.B.2 and the quasi relative interior of ψ includes the origin, so the constraint qualification is satisfied, whereby $V_p = V_d$, there is attainment in the dual, and p^* is optimal if there exists a dual optimum p with $\int (A^T p) db = f_0^*$ and $p \in \text{supp}(\int (A^T p) db)$. By the form of $\int (A^T p) db$ given in Proposition 1.B.2, $\int (A^T p) db = \int p^T A^T db = m$, and the first condition can be simplified to say $\int (A^T p) db = 0$ for all A^T , which in turn is equivalent to $p \in \ker(A^T)$. Since $\int (A^T p) db \in \mathbb{R}_+$, it follows that $p \in D_A$. Furthermore, if ψ is strictly convex and differentiable, by Proposition 1.B.2 it follows that $\int (A^T p) db$ is also strictly convex whereby p must be unique, and moreover the simplified optimality condition follows by the derivative rules given in Proposition 1.B.2.

In order to translate the dual expression into the desired form, first note that

$$f_0^*(q) = \sup_{\mathbb{R}^m} \int h; q \quad (q) = f_0(q):$$

As such, since $\int (A^T p) db = 0$ by Proposition 1.B.2 and $f_0^*(A^T p) = 0$ by the above, the dual value is always at least 0, thus values below this may be safely removed from the dual domain. First, $\int (A^T p) db = 1$ when p is not nonnegative everywhere by proposition 1.B.2, thus $p \geq 0$ may be stipulated. This combined with $f_0^*(A^T p) = 0$ if $\int (A^T p) db = 0$ for all p means the dual domain may be replaced with D_A .

Lastly, the simplified form of the optimality criteria are from Proposition 1.B.2, where the dual optimum is unique by strict convexity of $\int (A^T p) db$ again by Proposition 1.B.2. \square

Remark 2.2.4. Continuing the connection to maximum entropy, for any $\psi \in L_{fs}$, by Lemma 1.B.1, $\int (A^T p) db$ attains its maximal value at the vector $p \in \mathbb{R}^m$ with $p_i = \psi'(0) = m$ for all

coordinates $i \in [m]$, which is a rescaling of the uniform distribution. But recalling that the algorithm is initialized with $\beta_0 = 0$, then Proposition 1.B.2 grants $r \mathcal{L}(A_0) = r \mathcal{L}(0) = p$: that is, the initial dual iterate $p = r \mathcal{L}(A_0)$ is the unconstrained optimum! Let $\beta_t := r \mathcal{L}(A_t)$ denote the t^{th} dual iterate; since $r \mathcal{L}(r \mathcal{L}(z)) = z$ for any z (Zalinescu, 2002, Theorem 2.4.2.ii), then for any $\beta \in D_A \cap \ker(A^>)$,

$$r \mathcal{L}(\beta); \quad \sum_{i=1}^m h(A_{i,t}; \beta) = 0:$$

This allows the dual optimum to be rewritten as

$$\begin{aligned} \hat{\beta}_A &= \arg \min_{\beta \in D_A} \int_{\mathcal{Z}} \ell(\beta) d\mathbf{b} \\ &= \arg \min_{\beta \in D_A} \int_{\mathcal{Z}} \ell(\beta) d\mathbf{b} + \int_{\mathcal{Z}} \ell(\beta_t) d\mathbf{b} - \int_{\mathcal{Z}} \ell(\beta_t) d\mathbf{b} + \int_{\mathcal{Z}} \ell(\beta_t) d\mathbf{b} \\ &= \arg \min_{\beta \in D_A} \int_{\mathcal{Z}} \ell(\beta) d\mathbf{b} + \int_{\mathcal{Z}} \ell(\beta_t) d\mathbf{b} - \int_{\mathcal{Z}} \ell(\beta_t) d\mathbf{b} + \int_{\mathcal{Z}} \ell(\beta_t) d\mathbf{b} \end{aligned}$$

that is, the dual optimum $\hat{\beta}_A$ is the Bregman projection (according to ℓ) onto D_A of any dual iterate $\beta_t = r \mathcal{L}(A_t)$. In particular, $\hat{\beta}_A$ is the Bregman projection onto the feasible set of the unconstrained optimum $\beta_0 = r \mathcal{L}(A_0)$!

The connection to Bregman divergences runs deep; in fact, just as the specification of Boost in Figure 1.3 is purely in terms of primal variables (whereas the original formulation used both primal and dual variables), it is also possible to specify the algorithm purely in terms of dual variables, where the line search performs Bregman projections (Kivinen and Warmuth, 1999, Collins et al., 2002).

Before adjourning this section, the above strong duality rule and projection notation together lead to the duality gap stopping criterion (option 2 in Figure 1.3).

Remark 2.2.5. Given any β_t , the corresponding dual variable $\beta_t := r \mathcal{L}(A_t)$ is not necessarily dual feasible, which can be achieved by instead working with the projection $\hat{\beta}_t := P_S^1(\beta_t)$ where $S := D_A \cap \{p \in \mathbb{R}^m : \mathcal{L}(p) < 1/g\}$, the latter being a closed convex set. In order to ensure accuracy $\epsilon > 0$, it suffices by Theorem 2.2.1 to check $\mathcal{L}(A_t) + \mathcal{L}(\hat{\beta}_t) < 1/g + \epsilon$.

Unfortunately, this duality gap computation has two significant weaknesses. First, it is painfully expensive to compute; said another way, this projection does not fit the computational model of boosting, where A is never intended to be explicitly constructed or stored. (What is $\ker(A^>)$ when H consists of size-bounded decision trees?) Concretely, the primal-dual relationship

in Theorem 2.2.1 can be written

$$\inf_{\mathcal{L}} \mathbb{L}(\cdot) : \mathbb{L}(\cdot) \in \text{im}(A) = \max_{\mathcal{L}} \mathbb{L}(\cdot) : \mathbb{L}(\cdot) \in \ker(A^T) = \text{im}(A)^\circ$$

(since $\text{dom}(\mathbb{L}) \subseteq \mathbb{R}_+^m$ encodes the orthant constraint), where the computational black box gives elements of $\text{im}(A)$, but what is needed in the dual is a black box for $\ker(A^T) = \text{im}(A)^\circ$.

The second weakness is that the optimization guarantees of Chapter 3 do not indicate how fast $\mathbb{L}(A_t) + \mathbb{L}(\cdot_t)$ decreases. The guarantees themselves only indicate how fast the error $\mathbb{L}(A_t) - \mathbb{L}(A)$ decreases, and while the proofs do go through a duality gap obtained via projections, these projections are onto a set which does not seem easily computed by the algorithm.

2.3 Generalized Weak Learning Rate

The weak learning rate β was critical to the original convergence analysis of AdaBoost, providing a handle on the progress of the algorithm. But to be useful, this value must be positive, which was precisely the condition granted by the (empirical) weak learning assumption. This section will generalize the weak learning rate β into a quantity which can be made positive for any boosting instance.

Now let β denote the largest possible value for β ; this corresponds to the optimization problem

$$\beta = \inf_{\substack{\mathcal{L} \\ k \geq 1}} \max_{j \in [n]} \sum_{i=1}^n (x_i)_j y_i h_j(x_i) = \inf_{\substack{\mathcal{L} \\ k \geq 1}} \frac{\sum_{i=1}^n (x_i)_j y_i h_j(x_i)}{k} = \inf_{\substack{\mathcal{L} \\ k \geq 1}} \frac{\sum_{i=1}^n (x_i)_j y_i h_j(x_i)}{k} \quad (2.3.1)$$

where the dual form, with extraneous zero terms, will lead the path to the eventual generalization. For now, this simplified form leads to the following alternate characterization.

Proposition 2.3.2. A boosting instance is empirically weakly learnable (i.e., it is possible to choose $\beta > 0$) if and only if $D_A \neq \{0\}$.

Proof. Suppose $D_A \neq \{0\}$; since the first term in eq. (2.3.1) is of a continuous function over a compact set, it has some minimizer \mathcal{L}^* . But $\sum_{i=1}^n (x_i)_j y_i h_j(x_i) = 1$, meaning $\mathcal{L}^* \in D_A$, and so $\beta > 0$.

and (b) may be chosen positive. On the other hand, if $D_A \notin \{0\}$, take any $0 < \epsilon < \min\{D_A, \epsilon_0\}$; then

$$0 < (b) = \inf_{z \in \mathbb{R}_+^m \setminus \{0\}} \frac{kA^> z - k_1}{\|z\|} = \frac{kA^> \epsilon_0 - k_1}{\epsilon_0} = 0:$$

□

As such, to weaken the assumption $(b) > 0$ is the same as moving away from $D_A = \{0\}$, and a way to brute force this is to take various zeroes in the primal expression of eq. (2.3.1), and replace them with something else. This generalization depends on a generic set S , which was alluded to in the duality gap discussions in Remark 2.2.5, and will only be set in Chapter 3.

Definition 2.3.3. Given a matrix $A \in \mathbb{R}^{m \times n}$ and a set $S \subseteq \mathbb{R}^m$, define

$$(A; S) := \inf_{z \in S \setminus \ker(A^>)} \frac{kA^> z - k_1}{\|z\|} : z \in S \setminus \ker(A^>) :$$

First note that in the scenario of weak learnability (i.e., $D_A = \{0\}$ by Proposition 2.3.2), the choice $S = \mathbb{R}_+^m$ allows the new notion to exactly cover the old one: $(A; \mathbb{R}_+^m) = (b)$ (b).

To get a better handle on the meaning of S , recall the projection and distance notation P and D from Section 2.1. Suppose, for some t , that $r \in \mathbb{B}(A^> t) \cap S \setminus \ker(A^>)$; then the minimum within $(A; S)$ may be instantiated with $r \in \mathbb{B}(A^> t)$, yielding

$$(A; S) = \inf_{z \in S \setminus \ker(A^>)} \frac{kA^> z - k_1}{\|z\|} = \frac{kA^> r - k_1}{\|r\|} = \frac{kA^> r - k_1}{\|r\|} \frac{\|r\|}{\|r\|} = \frac{kA^> r - k_1}{\|r\|} \frac{\|r\|}{\|P_{S \setminus \ker(A^>)}^1(r)\|} : (2.3.4)$$

Rearranging this,

$$(A; S) = \frac{kA^> r - k_1}{\|r\|} \frac{\|r\|}{\|P_{S \setminus \ker(A^>)}^1(r)\|} = \frac{kA^> r - k_1}{\|P_{S \setminus \ker(A^>)}^1(r)\|} : (2.3.5)$$

This is helpful because the right hand side appears in standard guarantees for single-step progress in descent methods. Meanwhile, the left hand side has reduced the influence of r to a single number, and the normed expression is the distance to a restriction of dual feasible set, which will converge to zero if the minimum is to be approached, so long as this restriction contains the dual optimum.

This will be exactly the approach taken in Chapter 3; indeed, the first step towards convergence rates, Proposition 3.1.2, will use exactly the upper bound in eq. (2.3.5). The detailed work that remains is then dealing with the distance to the dual feasible set. The choice of δ will be made to facilitate the production of these bounds, and will depend on the optimization structure revealed in Section 2.4.

In order for these expressions to mean anything, $(A; S)$ must be positive.

Theorem 2.3.6. Let matrix $A \in \mathbb{R}^{m \times n}$ and polyhedron $S \subseteq \mathbb{R}^m$ be given with $S \cap \ker(A^T) \neq \emptyset$; and $S \setminus \ker(A^T) \neq \emptyset$. Then $(A; S) > 0$.

The proof, material on other generalizations of (b), and discussion on the polyhedrality of S can all be found in Section 2.A.

As a final connection, since $A^T P_{S \setminus \ker(A^T)}^1(x) = 0$ for any x , note that

$$(A; S) = \inf_{x \in S \cap \ker(A^T)} \frac{kA^T x_1}{P_{S \setminus \ker(A^T)}^1(x)k_1} = \inf_{x \in S \cap \ker(A^T)} \frac{kA^T (P_{S \setminus \ker(A^T)}^1(x))k_1}{P_{S \setminus \ker(A^T)}^1(x)k_1};$$

which now makes use of both extraneous zeros in eq. (2.3.1). To connect this to standard notions from linear algebra, if A were square and full rank, then the relationship between $kA^T x_2$ and kx_2 could be bounded by the smallest and largest singular values of A . The above expression has two modifications: it is a more general operator norm (with l_1 over the domain and l_1 over the image), and explicitly avoids problems caused by A not being full rank.

2.4 The Hard Core

When $A \in \mathbb{R}^{m \times n}$, the hard core may be specified in great detail. By Theorem 1.7.4, hard cores always exist, and since $\mathcal{L}^1(b)$ is identified here with \mathbb{R}^m , a hard core C is not only unique, but moreover it can be identified simply as a subset of the index set $[n]$.

It will be useful to view the rows $f_{a_i} g_{i=1}^m$ of A as a collection of m points in \mathbb{R}^n . Due to the form of \mathcal{L}^1 for $\lambda \in L_{fs}$, Boost is therefore searching for a halfspace, parameterized by a vector α , which contains all of these points. Sometimes such a halfspace may not exist, and applies a smoothly increasing penalty to points that are farther and farther outside it. It will be shown that the best choices of α place the entire hard core on the boundary of the halfspace defined by α (where the degenerate halfspace corresponding to $\alpha = 0$ is understood to have all of \mathbb{R}^m as its boundary), whereas the complement of the hard core falls interior to the halfspace.

2.4.1 Weak Learnability

Recall Proposition 2.3.2, which provided that an instance is empirically weakly learnable if $D_A = \{0\}$. With this in mind, the following theorem establishes four equivalent formulations of empirical weak learnability.

Theorem 2.4.1. For any $A \in \mathbb{R}^{m \times n}$, the following conditions are equivalent.

1. $\exists \gamma \in \mathbb{R}^n, A \gamma \geq \mathbf{1}$,
2. $\exists \gamma \in L_{fs}, \mathcal{L}(A) = 0$,
3. $\exists \gamma \in L_{fs}, \hat{\gamma}_A = 0$,
4. $D_A = \{0\}$.

First note that item 4 indicates (via Proposition 2.3.2) that this indeed captures the setting of weak learnability, which is equivalent to $|C| = 0$ where C is the hard core.

Recall the earlier discussion of boosting as searching for a halfspace containing the points $\{a_i\}_{i=1}^m = \{e_i^T A\}_{i=1}^m$; item 1 encodes precisely this statement, and moreover that there exists such a halfspace with these points interior to it. Note that this statement also encodes the margin separability equivalence of weak learnability due to Shalev-Shwartz and Singer (2008); specifically, if labels are bounded away from 0 and each point a_i (row of A) is replaced with $y_i a_i$, the definition of A grants that positive examples will land on one side of the hyperplane, and negative examples on the other.

Item 4 and item 1 can be interpreted geometrically, as depicted in Figure 2.1: the dual feasibility statement is that no convex combination of $\{a_i\}_{i=1}^m$ will contain the origin.

From this, the fact $\mathcal{L}(A) = 0$ (Item 2) is immediate. Lastly, when $D_A = \{0\}$, the dual optimum satisfies $\hat{\gamma}_A = 0$ (Item 3).

Proof of Theorem 2.4.1. (Item 1 \Rightarrow item 2.) Let $\gamma \in \mathbb{R}^n$ with $A \gamma \geq \mathbf{1}$, any increasing sequence $\{c_i\}_{i=1}^m$, and any $\gamma \in L_{fs}$ be given. Then, since $\mathcal{L}(A) > 0$ and $\lim_{x \rightarrow 1} \gamma(x) = 0$,

$$\inf \mathcal{L}(A) = \lim_{i \rightarrow \infty} \mathcal{L}(c_i A) = 0 = \inf \mathcal{L}(A);$$

(Item 2 \Rightarrow item 3.) Let any $\gamma \in L_{fs}$ be given. The point 0 is always dual feasible, and

Figure 2.1. Geometric view of the primal and dual problem, under weak learnability. The vertices of the pentagon denote the points $f_{a_i} g_{i=1}^m$. The arrow, denoting λ in item 1 of Theorem 2.4.1, defines a homogeneous halfspace containing these points; on the other hand, their convex hull does not contain the origin. Please see Theorem 2.4.1 and its discussion.

$\lambda^*(0) = 0$ (cf. Lemma 1.B.1) and the form of λ^* (cf. Proposition 1.B.2) grants

$$\inf \lambda^*(A) = 0 = \lambda^*(0):$$

Since the dual optimum is unique (cf. Theorem 2.2.1), $\lambda^*_A = 0$.

(Item 3 \Rightarrow item 4.) Choose any $\lambda \in L_{fs}$, whereby the assumption grants a unique optimum $\lambda^*_A = 0$. Suppose contradictorily that there exists another dual feasible point $\lambda \in \mathbb{R}_+^m$. By Lemma 1.B.1, λ^* is negative along $[0 \setminus \{0\}]^m \cap f_0 g$, which has nonzero intersection with $[\lambda^*_A; \lambda]$; consequently, there must exist $\lambda^0 \in D_A \cap f_0 g$ with $\lambda^*(\lambda^0) > 0 = \lambda^*(\lambda^*_A)$, which contradicts the optimality of λ^*_A .

(Item 4 \Rightarrow item 1.) This case is a one direction of Gordan's theorem (Borwein and Lewis, 2000, Theorem 2.2.6). \square

2.4.2 Attainability

For strictly convex functions, there is a nice characterization of attainability, which will require the following definition.

Definition 2.4.2 (Hiriart-Urruty and Lemaëchal (2001, Section B.3.2)). A closed convex function h is called θ -coercive when all level sets are compact. (That is, for any $\gamma \in \mathbb{R}$, the set $\{x : f(x) \leq \gamma\}$ is compact.)

Proposition 2.4.3. Suppose h is strictly convex, closed, and proper. Then $\inf_x h(x)$ is attainable

is 0-coercive.

Armed with this notion, it is now possible to build an attainability theory for $\mathbb{L}^b A$. Some care must be taken with the above concepts, however; note that while \mathbb{L}^b is strictly convex, $\mathbb{L}^b A$ need not be (for instance, if there exist nonzero elements $\ker(A)$, then moving along these directions does not change the objective value). Therefore, 0-coercivity statements will refer to the function

$$(\mathbb{L}^b + \text{im}(A))(x) = \begin{cases} \mathbb{L}^b(x) & \text{when } x \in \text{im}(A); \\ 1 & \text{otherwise.} \end{cases}$$

This function is effectively taking the epigraph of \mathbb{L}^b , and intersecting it with a slice representing $\text{im}(A) = \{Ax : x \in \mathbb{R}^n\}$, the set of points considered by the algorithm. As such, it is merely a convenient way of dealing with $\ker(A)$ as discussed above.

Theorem 2.4.4. For any $A \in \mathbb{R}^{m \times n}$, the following conditions are equivalent.

1. $\exists \delta \in \mathbb{R}^n, A \in \mathbb{R}^{m \times n} \setminus \{0\}$,
2. $\exists \delta \in L_{fs}, \mathbb{L}^b + \text{im}(A)$ is 0-coercive,
3. $\exists \delta \in L_{fs}, A \in \mathbb{R}_{++}^m$,
4. $D_A \setminus \mathbb{R}_{++}^m \neq \emptyset$;

Following the discussion above, item 2 is the desired attainability statement.

Next, note that item 4 is equivalent to the expression $|C| = m$ (where C is the hard core), i.e. there exists a reweighting of b with $p > 0$ which decorrelates every $h \in H$ from the target.

For a geometric interpretation, consider item 1 and item 4. Item 1 says that any halfspace containing some a_i within its interior must also fail to contain some a_j (with $i \neq j$). (Item 1 also allows for the scenario that no valid enclosing halfspace exists, i.e. $= 0$.) Item 4 states that the origin 0 is contained within a positive convex combination of a_i $g_{i=1}^m$ (alternatively, the origin is within the relative interior of these points). These two scenarios appear in Figure 2.2.

Finally, note item 3: it is not only the case that there are dual feasible points fully interior to \mathbb{R}_{++}^m , but furthermore the dual optimum is also interior. This will be crucial in the convergence rate analysis, since it will allow the dual iterates to never be too small.

Figure 2.2. Geometric view of the primal and dual problem, under attainability. Once again, the $f_i, g_{i=1}^m$ are the vertices of the pentagon. This time, no (closed) homogeneous halfspace containing all the points will contain one strictly, and the relative interior of the pentagon contains the origin. Please see Theorem 2.4.4 and its discussion.

Proof of Theorem 2.4.4. (Item 1 \Rightarrow item 2.) Let $d \in \mathbb{R}^m \setminus \{0\}$, $\lambda \in \mathbb{R}^n$, and $\psi \in L_{fs}$ be arbitrary. To show 0-coercivity, it suffices (Hiriart-Urruty and Lemaréchal, 2001, Proposition B.3.2.4.iii) to show

$$\lim_{t \downarrow 0} \frac{\psi(A + td) + \inf_{x \in \text{im}(A)} \psi(A + td) - \psi(A)}{t} > 0. \quad (2.4.5)$$

If $d \notin \text{im}(A)$ and $t > 0$, then $\inf_{x \in \text{im}(A)} \psi(A + td) = \psi(A)$. Now suppose $d \in \text{im}(A)$; by item 1, since $d \neq 0$, then $d \in \mathbb{R}^m$, meaning there is at least one positive coordinate. But then, since $\psi > 0$ and ψ is convex,

$$\begin{aligned} (2.4.5) \quad & \lim_{t \downarrow 0} \frac{\sum_{j=1}^m d_j \psi(e_j^T (A + td)) - \psi(A)}{t} \\ & \lim_{t \downarrow 0} \frac{\sum_{j=1}^m d_j \psi(e_j^T A) + t \sum_{j=1}^m d_j \psi(e_j^T A) - \psi(A)}{t} \\ & = \sum_{j=1}^m d_j \psi(e_j^T A); \end{aligned}$$

which is positive by the selection of d_j and since $\psi > 0$.

(Item 2 \Rightarrow item 3.) Choose any $\psi \in L_{fs}$. Since the infimum is attainable, designate any x^* satisfying $\inf \psi(A) = \psi(A)$ (note, although ψ is strictly convex, $\psi|_A$ need not be, thus uniqueness is not guaranteed!). By the optimality conditions on the duality rule in Theorem 2.2.1, $\lambda_A = \lambda^0(A) = \psi > 0$ everywhere.

(Item 3 \Rightarrow item 4.) This holds since $D_A \psi = \lambda_A^0$ and $\lambda_A^0 \in \mathbb{R}_{++}^m$ for every $\psi \in L_{fs}$.

(Item 4 =) item 1.) This case is directly handled by Stiemke's theorem (Borwein and Lewis, 2000, Exercise 2.2.8). □

2.4.3 General Setting

So far, the scenarios of weak learnability and attainability corresponded to the extremal hard core cases of $C_j \neq \emptyset$; m_j . The situation in the general setting $1 \leq j \leq m-1$ is basically as good as one could hope for: it interpolates between the two extremal cases.

As a first step, partition A into two submatrices according to C .

Definition 2.4.6. Let $A \in \mathbb{R}^{m \times n}$ and corresponding hard core C be given. Partition A by rows into two matrices $A_0 \in \mathbb{R}^{m_0 \times n}$ and $A_+ \in \mathbb{R}^{m_+ \times n}$, where A_+ has rows corresponding to C , and $m_+ = |C|$. For convenience, permute the examples so that

$$A = \begin{pmatrix} h \\ A_0 \\ A_+ \end{pmatrix} :$$

(This merely relabels the coordinate axes, and does not change the optimization problem.) Note that this decomposition is unique, since C is uniquely specified.

As a first consequence, this partition cleanly decomposes the dual feasible set D_A into D_{A_0} and D_{A_+} .

Proposition 2.4.7. For any $A \in \mathbb{R}^{m \times n}$, $D_{A_0} = \{0\}$, $D_{A_+} \setminus \mathbb{R}_{++}^{m_+} \neq \emptyset$, and

$$D_A = D_{A_0} \cup D_{A_+} :$$

Furthermore, no other partition of A into $B_0 \in \mathbb{R}^{z \times n}$ and $B_+ \in \mathbb{R}^{p \times n}$ satisfies these properties.

Proof. It must hold that $D_{A_0} = \{0\}$, since otherwise there would exist $\lambda \in \ker(A_0^>) \setminus \mathbb{R}_+^{m_0}$ with $\lambda \neq 0$, which could be extended to $\lambda = \begin{pmatrix} 0 \\ \lambda \end{pmatrix} \in D_A$ and the positive coordinate of λ could be added to C , contradicting the maximality of C .

The property $D_{A_+} \setminus \mathbb{R}_{++}^{m_+} \neq \emptyset$ is a consequence of Theorem 2.4.4 and the definition of C .

For the decomposition, note first that certainly every $\lambda \in D_{A_0} \cup D_{A_+}$ satisfies $\lambda \in D_A$. Now suppose contradictorily that there exists $\lambda \in D_A \setminus (D_{A_0} \cup D_{A_+})$. There must exist $j \in [m] \setminus C$ with $(\lambda)_j > 0$, since otherwise $\lambda \in \{0\} \cup D_{A_+}$; but that means j should have been included in C , a contradiction.

Lastly, the uniqueness is a consequence of the uniqueness of hard cores of \mathbb{R}^n (cf. Theorem 1.7.4, where uniqueness follows from the structure of \mathbb{R}^n). \square

The main result of this section will have the same two main ingredients as Proposition 2.4.7:

The full boosting instance may be uniquely decomposed into two pieces A_0 and A_+ , each of which individually behave like the weak learnability and attainability scenarios.

The subinstances have a somewhat independent effect on the full instance.

Theorem 2.4.8. Let $A \in \mathbb{R}^{m \times n}$ be given. Let $B_0 \in \mathbb{R}^{z \times n}$, $B_+ \in \mathbb{R}^{p \times n}$ be any partition of A by rows. The following conditions are equivalent.

1. $\exists \gamma \in \mathbb{R}^n$, $B_0 \gamma \wedge B_+ \gamma = 0$ and $\exists \delta \in \mathbb{R}^n$, $B_+ \delta \in \mathbb{R}^p \setminus \{0\}$,
2. $\exists \gamma \in L_{fs}$, $\mathcal{D}(A) = \mathcal{D}(B_+)$, and $\mathcal{D}(B_0) = 0$, and $\mathcal{D}_{+ \text{im}(B_+)}$ is 0-coercive,
3. $\exists \gamma \in L_{fs}$, $A = \begin{pmatrix} B_0 \\ B_+ \end{pmatrix}$ with $\gamma_{B_0} = 0$ and $\gamma_{B_+} \in \mathbb{R}_{++}^p$,
4. $D_{B_0} = 0$, and $D_{B_+} \setminus \mathbb{R}_{++}^p \neq \emptyset$; and $D_A = D_{B_0} \cup D_{B_+}$.

Stepping through these properties, notice that item 4 mirrors the expression in Proposition 2.4.7. But that proposition also granted that this representation was unique, thus only one partition of A satisfies the above properties, namely $A_0; A_+$. Since this theorem is stated as a series of equivalences, any one of these properties can in turn be used to identify the hard core set C .

To continue with geometric interpretations, notice that item 1 states that there exists a halfspace strictly containing those points in $[m] \times \mathbb{C}$, with all points of C on its boundary; furthermore, trying to adjust this halfspace to contain elements of C will place others outside it. With regards to the geometry of the dual feasible set as provided by item 4, the origin is within the relative interior of the points corresponding to C , however the convex hull of the other $[m] \setminus C$ points can not contain the origin. Furthermore, if the origin is written as a convex combination of all points, this combination must place zero weight on the points with indices $[m] \setminus C$. This scenario is depicted in Figure 2.3.

In items 2 and 3, B_0 mirrors the behavior of weakly learnable instances in Theorem 2.4.1, and analogously B_+ follows instances with minimizers from Theorem 2.4.4. The interesting addition, as discussed above, is the independence of these components: item 2 provides that the

Figure 2.3. Geometric view of the primal and dual problem in the general case. There is a closed homogeneous halfspace containing the points $a_i g_{i=1}^m$, where the hard core lies on the halfspace boundary, and the other points are within its interior; moreover, there does not exist a closed homogeneous halfspace containing all points but with strict containment on a point in the hard core. Finally, although the origin is in the convex hull of $f a_i g_{i=1}^m$, any such convex combination places zero weight on points outside the hard core. Please see Theorem 2.4.8 and its discussion.

in mum of the combined problem is the sum of the in ma of the subproblems, while item 3 provides that the full dual optimum may be obtained by concatenating the subproblems' dual optima.

Proof of Theorem 2.4.8. (Item 1 \Rightarrow item 2.) Let $\lambda \in L_{fs}$ and $B_0 \in \mathbb{R}^z$ and $B_+ = 0$, and let $f c_i g_{i=1}^1 \rightarrow 1$ be an arbitrary sequence increasing without bound. Lastly, let $f_i g_i^1$ be a minimizing sequence for $\inf \mathcal{L}(B_+)$. Then

$$\begin{aligned} \inf m_+ \mathcal{L}(B_+) &= \lim_{i \uparrow} (m_+ \mathcal{L}(B_+ - i) + m_0 \mathcal{L}(c_i B_0)) = \inf m \mathcal{L}(A) \\ &= \inf (m_+ \mathcal{L}(B_+) + m_0 \mathcal{L}(B_0)) = \inf m_+ \mathcal{L}(B_+); \end{aligned}$$

which used the fact that $\mathcal{L}(B_0) = 0$ since $\mathcal{L} \geq 0$. Since the chain of inequalities starts and ends the same, it must be a chain of equalities, which means $\inf \mathcal{L}(B_0) = 0$. To show 0-coercivity of $\mathcal{L} + \inf_{im(B_+)}$, note the second part of item 1 is one of the conditions of Theorem 2.4.4.

(Item 2 \Rightarrow item 3.) First, by Theorem 2.4.1, $\inf \mathcal{L}(B_0) = 0$ means $\lambda_{B_0} = 0$ and

$D_{B_0} = \{0\}$. Thus

$$\begin{aligned} m\mathbb{L}(\hat{A}) &= \sup_{z \in D_A} m\mathbb{L}(z) \\ &= \sup_{z \in D_{B_0}} m_0\mathbb{L}(z) + \sup_{p \in D_{B_+}} m_+\mathbb{L}(p) \\ &= \inf_{z \in \mathbb{R}^n} m_+\mathbb{L}(B_+ z) = \inf_{z \in \mathbb{R}^n} m\mathbb{L}(A z) = m\mathbb{L}(\hat{A}): \end{aligned}$$

Combining this with $\mathbb{L} = \mathbb{R}^n$ and $\mathbb{L}(0) = 0$ (cf. Proposition 1.B.2 Lemma 1.B.1), $m\mathbb{L}(\hat{A}) = m_+\mathbb{L}(\hat{B}_+) = m\mathbb{L}(\hat{B}_+)$. But Theorem 2.2.1 shows \hat{A} was unique, which gives the result. To obtain $\hat{B}_+ \in \mathbb{R}^p_{++}$, use Theorem 2.4.4 with the 0-coercivity of $\mathbb{L} + \text{im}(B_+)$.

(Item 3 =) item 4.) Since $\hat{B}_0 = 0$, it follows by Theorem 2.4.1 that $D_{B_0} = \{0\}$. Furthermore, since $\hat{B}_+ \in \mathbb{R}^p_{++}$, it follows that $D_{B_+} \setminus \mathbb{R}^p_{++} = \emptyset$. Now suppose contradictorily that $D_A \not\subseteq D_{B_0} \cup D_{B_+}$; since it always holds that $D_A \subseteq D_{B_0} \cup D_{B_+}$, this supposition grants the existence of $(z, p) \in D_A$ where $z \in \mathbb{R}^n \setminus \{0\}$.

Consider the element $q := z + \hat{A} p$, which has more nonzero entries than \hat{A} , but still $q \in D_A$ since D_A is a convex cone. Let I_q index the nonzero entries of q , and let A_q be the restriction of A to the rows I_q . Since $q \in D_A$, meaning q is nonnegative and $q \in \ker(A^>)$, it follows that the restriction of q to its positive entries is within $\ker(A_q^>)$ (because only zeros of q and matching rows of A are removed, dot products between q with rows of $A^>$ are the same as dot products between the restriction of q and rows of $A_q^>$), and so $q \in D_{A_q}$, meaning $D_{A_q} \setminus \mathbb{R}^p_{++}$ is nonempty. Correspondingly, by Theorem 2.4.4, the dual optimum \hat{A}_q of this restricted problem will have only positive entries. But by the same reasoning granting that q restricted to I_q is within D_{A_q} , it follows that the full optimum \hat{A} , restricted to I_q , must also be within D_{A_q} (since, by q 's construction, \hat{A} 's zero entries are a superset of the zero entries of q). Therefore this restriction \hat{A}_q of \hat{A} to I_q will have at least one zero entry, meaning it can not be equal to \hat{A}_q ; but Theorem 2.2.1 provided that the dual optimum is unique, thus $\mathbb{L}(\hat{A}_q) > \mathbb{L}(\hat{A}_q)$. Finally, produce \hat{A}_q from \hat{A}_q by inserting a zero for each entry of I_q ; the same reasoning that allows feasibility to be maintained while removing zeros allows them to be added, and thus $\hat{A}_q \in D_A$. But this is a contradiction: since $\mathbb{L}(0) = 0$ (cf. Lemma 1.B.1), both \hat{A}_q and the optimum \hat{A}

have zero contribution to the objective along the entries outside of q , and thus

$$m\mathbb{L}(\hat{A}_q) = \sum_j I_{qj} \mathbb{L}(\hat{A}_q) > \sum_j I_{qj} \mathbb{L}(\hat{A}) = m\mathbb{L}(\hat{A});$$

meaning \hat{A}_q is feasible and has strictly greater objective value than the optimum \hat{A} .

(Item 4 =) item 1.) Unwrapping the definition of D_A , the assumed statements imply

$$(8 \ 0 \ 2 \ \mathbb{R}_+^z \ \text{nf} \ 0g; \ + \ 2 \ \mathbb{R}_+^p \ B_0^z \ + \ B_+^z \ + \ \notin \ 0) \wedge (9 \ + \ 2 \ \mathbb{R}_+^p \ B_+^z \ + \ = \ 0):$$

Applying Motzkin's transposition theorem (cf. (Ben-Israel, 2002) or (Dantzig and Thapa, 2003, Theorem 2.16)) to the left statement and Stiemke's theorem (cf. (Borwein and Lewis, 2000, Exercise 2.2.8), which is implied by Motzkin's theorem) to the right yields

$$(9 \ 2 \ \mathbb{R}^n \ B_0 \ 2 \ \mathbb{R}^z \ \wedge \ B_+ \ 2 \ \mathbb{R}^p) \wedge (8 \ 2 \ \mathbb{R}^n \ B_+ \ 62\mathbb{R}^p \ \text{nf} \ 0g);$$

which implies the desired statement. \square

Remark 2.4.9. Notice the dominant role A and corresponding hard core C play in the structure of the optimization solutions. For every $i \in [m] \setminus C$, the corresponding dual weights go to zero (i.e., $(r \ \mathbb{L}(A \ t))_i \neq 0$), and the corresponding primal margins grow unboundedly (i.e., $e_i^T A \ t \rightarrow 1$, since otherwise $\inf \mathbb{L}(A_0) > 0$). This is completely unaffected by the choice of $\lambda \in \mathbb{L}_{fs}$. Furthermore, whether this instance is weak learnable, attainable, or neither is dictated purely by A (respectively $|C| = 0$, $|C| = m$, or $|C| \in [1; m - 1]$).

Where different loss functions disagree is how they assign dual weight to the points in C . In particular, each $\lambda \in \mathbb{L}_{fs}$ (and corresponding \mathbb{L}) defines a notion of entropy via \mathbb{L} . The dual optimization in Theorem 2.2.1 can then be interpreted as selecting the max entropy choice (per \mathbb{L}) amongst those convex combinations of C equal to the origin.

Appendix 2.A Generalizing the Weak Learning Rate

2.A.1 Choosing a Generalization to (b)

Any generalization $\hat{0}$ of (b) should satisfy the following properties.

When empirical weak learnability holds (i.e., when $(b) > 0$), then $\hat{0} = (b)$.

For any boosting instance, $\gamma \in (0; 1)$.

γ provides an expression similar to eq. (2.3.5), which allows the full gradient to be converted into a notion of suboptimality in the dual.

Taking the form of the classical weak learning rate from eq. (2.3.1) as a model, the template generalized weak learning rate is

$$\gamma(A; S; C; D) := \inf_{S \in C} \frac{kA^> k_1}{\inf_{S \in D} k_1};$$

for some sets S , C , and D (for instance, the classical weak learning rate uses $S = R_+^m$ and $C = D = \{0\}$). In order to provide an expression similar to eq. (2.3.5), the domain of the minimum must include every suboptimal dual iterate $r \in \mathcal{B}(A, t)$.

Any choice C which does not include all of $\ker(A^>)$ is immediately problematic: this allows $S \in C \setminus \ker(A^>)$ to be selected, whereby $kA^> k_1 = 0$ and $\gamma = 0$. But note that without being careful about D , it is still possible to force the value 0.

Remark 2.A.1. Another generalization is to define

$$\gamma(A) := \gamma(A; R_+^m; \ker(A^>); \{0\}) = \inf_{S \in R_+^m \setminus \ker(A^>)} \frac{kA^> k_1}{k_1}. \quad (2.A.2)$$

This form agrees with the original γ when weak learnability holds, and leads to a very convenient analog to eq. (2.3.5).

Unfortunately, $\gamma(A)$ may be zero. Specifically, take A to be the matrix S defined in Section 3.4, due to Schapire (2010), where

$$S = \frac{\gamma(0)}{3} \begin{pmatrix} h_1 & i \\ 1 & 0 \end{pmatrix}.$$

Furthermore, for any $\gamma \in (0; 1)$, define

$$a := \begin{pmatrix} h_1 & i \\ 0 & 1 \end{pmatrix} \in \text{im}(S); \quad \gamma := (1 \quad \gamma) \begin{pmatrix} 1=2 \\ 1=2 \\ 0 \end{pmatrix} + \begin{pmatrix} h_1 & i \\ 1 & 0 \end{pmatrix} \in \ker(S^>).$$

Then

$$\inf_{S \in R_+^m \setminus \ker(S^>)} \frac{kS^> k_1}{k_1} = \inf_{\gamma \in (0; 1)} \frac{kS^> (\begin{pmatrix} h_1 & i \\ 1 & 0 \end{pmatrix} + \gamma \begin{pmatrix} 1=2 \\ 1=2 \\ 0 \end{pmatrix}) k_1}{k_1} = \inf_{\gamma \in (0; 1)} \frac{1}{1} = 0:$$

The natural correction to these worries is to set $C = D = \ker(A^>)$. But there is still sensitivity due to S .

Remark 2.A.3. Set $A := (1; 1)^>$, meaning $\ker(A^>) = \{z \in \mathbb{R}^2 : z \perp (1; 1)^>\}$, and $S = B((1; 1)^>; \sqrt{2})$, the ball of radius $\sqrt{2}$ around $(1; 1)^>$; note that $S \setminus \ker(A^>) \neq \emptyset$. Consider $\rho(A; S; \ker(A^>); \ker(A^>))$, and the sequence $\{g_i\}_{i=1}^\infty$ where

$$g_i = (1; 1)^> + \frac{1}{\sqrt{i^2 + 1}} \begin{pmatrix} 2i \\ 3i + 1 \end{pmatrix}$$

Note that $\|g_i - (1; 1)^>\|_2 = \sqrt{2}$, thus $g_i \in S$. Furthermore, $A^> g_i \neq 0$, so $g_i \in S \setminus \ker(A^>)$. As such,

$$\rho(A; S; \ker(A^>); \ker(A^>)) = \inf_i \frac{A^> g_i}{\|g_i\|_2} \tag{2.A.4}$$

$$= \frac{k(1; 1)^> (1; 1)^> + \frac{1}{\sqrt{i^2 + 1}} \begin{pmatrix} 2i \\ 3i + 1 \end{pmatrix}}{\sqrt{2} \sqrt{i^2 + 1}} \tag{2.A.5}$$

Using $\sqrt{1 + y} \leq 1 + \frac{y}{2}$, the numerator has upper bound

$$\begin{aligned} k(1; 1)^> (1; 1)^> + \frac{1}{\sqrt{i^2 + 1}} \begin{pmatrix} 2i \\ 3i + 1 \end{pmatrix} &= \sqrt{2} \sqrt{i^2 + 1} + \frac{2i}{\sqrt{i^2 + 1}} \\ &= 2i \left(\sqrt{1 + \frac{1}{i^2}} + 1 \right) \\ &= 2i \left((2 + \frac{1}{i^2})^{1/2} + 1 \right) = 2 + \frac{1}{i} \end{aligned}$$

The denominator is

$$\begin{aligned} \sqrt{2} \sqrt{i^2 + 1} &= \sqrt{2} \sqrt{i^2 + 1} \\ &= \sqrt{2} \left(\sqrt{i^2 + 1} + \sqrt{i^2 + 1} \right) \\ &= 2\sqrt{2} \sqrt{i^2 + 1} \end{aligned}$$

Thus eq. (2.A.5) is bounded above by $\frac{2 + \frac{1}{i}}{2\sqrt{2} \sqrt{i^2 + 1}} \rightarrow 0$.

The difficulty here was the curvature of S , which allowed elements arbitrarily close to $\ker(A^>)$ without actually being inside this subspace. This possibility is averted here by requiring

polyhedrality of S . This choice is sufficiently rich to allow the various dual-distance upper bounds of Chapter 3.

2.A.2 Proof of Theorem 2.3.6

The proof of Theorem 2.3.6 requires a few steps, but the strategy is straightforward. First note that $(A; S)$ can be rewritten as

$$\begin{aligned}
 (A; S) &= \inf_{x \in S \setminus \ker(A^>)} \frac{kA^> x k_1}{k P_{S \setminus \ker(A^>)}^1(x) k_1} \\
 &= \inf_{x \in S \setminus \ker(A^>)} \frac{kA^> (P_{S \setminus \ker(A^>)}^1(x)) k_1}{k P_{S \setminus \ker(A^>)}^1(x) k_1} \\
 &= \inf_{v \in \mathbb{R}^m} \frac{kA^> v k_1}{k v k_1} : v \in \mathbb{R}^m, v \in S, v \in P_{S \setminus \ker(A^>)}^1(0) \\
 &= \inf_{v \in \mathbb{R}^m} kA^> v k_1 : k v k_1 = 1, v \in S, v \in P_{S \setminus \ker(A^>)}^1(0) ; \quad (2.A.6)
 \end{aligned}$$

where the second equivalence uses $P_{S \setminus \ker(A^>)}^1(0) = 0$.

In the final form, $v \in \ker(A^>)$, and so $A^> v = 0$; that is to say, the numerator is positive for every element of its domain. The difficulty is that the domain of the minimum, written in this way, is not obviously closed; thus one can not simply assert the minimum is attainable and positive.

The goal then will be to reparameterize the minimum to have a compact domain. For technical convenience, the result will be mainly proved for the ℓ_2 norm (where projections behave nicely), and norm equivalence will provide the final result.

Lemma 2.A.7. Given $A \in \mathbb{R}^{m \times n}$ and a polyhedron $S \subset \mathbb{R}^m$ with $S \setminus \ker(A^>) \neq \emptyset$; and $S \cap \ker(A^>) \neq \emptyset$;

$$\inf_{x \in S \setminus \ker(A^>)} \frac{kA^> (P_{S \setminus \ker(A^>)}^2(x)) k_2}{k P_{S \setminus \ker(A^>)}^2(x) k_2} : x \in S \setminus \ker(A^>) \neq \emptyset \quad (2.A.8)$$

To produce the desired reparameterization of this minimum, the following characterization of polyhedral sets will be used.

Definition 2.A.9. For any nonempty polyhedral set $S \subset \mathbb{R}^m$, let H_S index a finite (but possibly empty) collection of affine functions $g : \mathbb{R}^m \rightarrow \mathbb{R}$ so that $S = \{ x \in \mathbb{R}^m : g(x) \leq 0 \}$ (with the convention that $S = \mathbb{R}^m$ when $H_S = \emptyset$). For any $x \in S$, let $I_S(x)$ denote the active set for x :

$I_S(x) = \{ i : g_i(x) = 0 \}$. Lastly, define a relation \sim_S over points in S : given $x, y \in S$, $x \sim_S y$ if

$I_S(x) = I_S(y)$. Observe that \sim_S is an equivalence relation over points within S , and let C_S be the set of equivalence classes.

The equivalence relation \sim_S thus partitions S into the members of C_S , each of which has a very convenient structure.

Lemma 2.A.10. Let a polyhedral set $S \subseteq \mathbb{R}^m$ be given, and x a nonempty $F \in C_S$. Then F is convex, and F is equal to its relative interior (i.e., $F = \text{ri}(F)$). Finally, fixing an arbitrary $z_0 \in F$, the normal cone at any point $z \in F$ is orthogonal to the vector space parallel to the affine hull of F (i.e., $N_F(z) = (a \perp F) \cap \{z\}^\circ = (a \perp F) \cap \{z_0\}^\circ$).

Throughout the remainder of this section, normal and tangent cones will be considered at points within a set $F \in C_S$. As Lemma 2.A.10 establishes, any $se \in F \in C_S$ is relatively open ($F = \text{ri}(F)$), however, the required properties of normal and tangent cones, as developed by Hiriart-Urruty and Lemaréchal (2001, Sections A.5.2 and A.5.3), suppose closed convex sets. But it is always the case that $\text{ri}(F) = \text{ri}(\text{cl}(F))$ (Hiriart-Urruty and Lemaréchal, 2001, Proposition A.2.1.8); as such, the normal and tangent cones at the desired relative interior points may just as well be constructed against $\text{cl}(F)$, and thus the aforementioned properties safely hold.

Proof. If $S = \mathbb{R}^m$ (meaning H_S is empty) or $\dim(F) = 0$ (F is a single point), everything follows directly, thus suppose $S \subsetneq \mathbb{R}^m$, and x a nonempty $F \in C_S$ with $\dim(F) > 0$.

Let any $x_0, x_1 \in F$ and $\lambda \in [0, 1]$ be given, and define $x := (1 - \lambda)x_0 + \lambda x_1$. Since each $g \in \text{defining } S$ is a ne,

$$g(x) = (1 - \lambda)g(x_0) + \lambda g(x_1) \quad (2.A.11)$$

By construction of C_S , $g(x_0) = 0$ if $g(x_1) = 0$ and otherwise both are negative, thus $g(x) = 0$ if $g(x_0) = g(x_1) = 0$, meaning $I_S(x) = I_S(x_0) = I_S(x_1)$, so $x \in F$ and F is convex.

Now let any $y_0 \in F$ be given; $y_0 \in \text{ri}(F)$ when there exists a $\delta > 0$ so that

$$B(y_0; \delta) \setminus a(F) \cap F \quad (2.A.12)$$

(Hiriart-Urruty and Lemaréchal, 2001, Definition A.2.1.1). To this end, first define δ to be half the distance to the closest hyperplane defining S which is not active for y_0 :

$$\delta := \frac{1}{2} \min_{H \in \text{nl}_S(y_0)} \min_{k \perp y_0} \|y_0 - k\|_2 : y_0 \in \mathbb{R}^m; g(y_0) = 0$$

Since there are only finitely many such hyperplanes, and the distance to each is nonzero, $\delta > 0$. Let any $y \in B(y_0; \delta) \setminus a(F)$ be given; by definition of $a(F)$, there must exist $\lambda \in \mathbb{R}$ and $y_1 \in F$ so that $y = (1 - \lambda)y_0 + \lambda y_1$. By eq. (2.A.11), for any $\lambda \in I_S(y_0) = I_S(y_1)$,

$$g(y) = (1 - \lambda)g(y_0) + \lambda g(y_1) = 0:$$

On the other hand, for any $\lambda \in H_S \cap I_S(y_0)$, it must be the case that $g(y) < 0$, since $y \in B(y_0; \delta)$, and due to the choice of δ . Returning to the definition of relative interior in eq. (2.A.12), it follows that $y_0 \in \text{ri}(F)$, and $\text{ri}(F) = F$ since $y_0 \in F$ was arbitrary.

For the final property, for any $z_0; z \in \text{ri}(F) = F$, the tangent cone $T_F(z)$ has form $a(F) + f(z)$ (see Hiriart-Urruty and Lemaréchal, 2001, Proposition A.5.2.1 and discussion within Section A.5.3), and note $a(F) + f(z) = a(F) + f(z_0) + z - z_0 + f(z_0) = a(F) + f(z_0)$. Lastly, $N_F(z) = T_F(z)^\circ$ (Hiriart-Urruty and Lemaréchal, 2001, Proposition A.5.2.4). \square

The relevance to eq. (2.A.8) and eq. (2.A.6) is that projections from polyhedron S onto $S \setminus \ker(A^>)$ (itself a polyhedron, as is verified in the proof of Lemma 2.A.7) must land on some equivalence class of $\mathcal{C}_{S \setminus \ker(A^>)}$, and these projections are easily characterized.

Lemma 2.A.13. Let any nonempty polyhedra $S \subseteq \mathbb{R}^m$ and $K \subseteq \mathbb{R}^m$ be given, and x any nonempty $F \subseteq \mathcal{C}_{S \setminus K}$ and $x_F \in F$. Define

$$P_F := \{c \mid \exists \lambda \in \mathbb{R}^2_{S \setminus K}(\lambda) : c > 0; \lambda \in S; \lambda \in P_{S \setminus K}^2(\lambda) \in F\};$$

$$D_F := N_F(x_F) \setminus \{y \mid x_F : y \in \mathbb{R}^m; \exists \lambda \in I_S(x_F) : g(y) = 0\};$$

where $N_F(x_F)$ is the normal cone of F at x_F . Then $P_F = D_F$.

Note that the final active set $I_S(x_F)$ is with respect to S , not $S \setminus K$.

Proof. () Let any $\lambda \in S$ with $\lambda := P_{S \setminus K}^2(\lambda) \in F$ be given, where the latter is well-defined since F and hence $S \setminus K$ are nonempty. By Lemma 2.A.10, $\lambda \in \text{ri}(F)$, and $N_F(\lambda) = N_F(x_F)$, meaning $\lambda \in N_F(x_F)$ (Hiriart-Urruty and Lemaréchal, 2001, Proposition A.5.3.3). Since $\lambda \in S$, for

any $z \in I_S(x_F) \cap H_S, g(z) = 0$, so

$$\begin{aligned} \inf_{y \in I_S(x_F)} g(y) - 0 &= \inf_{y \in I_S(x_F)} g(y) - g(x_F) \\ &= \inf_{y \in I_S(x_F)} (g(y) - g(x_F)) \\ &= \inf_{y \in I_S(x_F)} g(y - x_F) \end{aligned}$$

the final equality following since $g(x_F) = g(x_F) = 0$ and g defines an affine hyperplane, meaning the corresponding affine halfspace is closed under translations by x_F . This holds for all $z \in I_S(x_F)$, thus $z \in D_F$, and since D_F is a convex cone, for any $\alpha > 0, \alpha z \in D_F$.

() Define

$$\delta := \min_{\{z \in H_S \cap I_S(x_F); z \in D_F; g(z) = 0\}} \|z - x_F\|_2$$

For any fixed x_F , this minimum is positive since $g(x_F) < 0$, while polyhedrality of S grants that $I_S(x_F)$ ranges over a finite set, together meaning $\delta > 0$. Now let any $v \in D_F$ be given, and set

$z := x_F + v = (2\|v\|_2 k)$. The form of D_F immediately grants $g(z) = 0$ for $z \in I_S(x_F)$, but notice for $z \in H_S \cap I_S(x_F)$, it still holds that $g(z) = 0$, since $g(x_F) < 0$ and $\|x_F - z\|_2 < \delta$. So $v = (2\|v\|_2 k) \in P_{S \setminus K}^2(x_F)$ where $z \in S$ and $P_{S \setminus K}^2(x_F) = x_F + 2F$, meaning $v \in P_F$. \square

The result now follows by considering all elements of $C_{S \setminus \ker(A^>)}$.

Proof of Lemma 2.A.7. For convenience, set $K := \ker(A^>)$. Note that K (and hence $S \setminus K$) is a polyhedron; indeed, it has the form

$$\begin{aligned} K = \ker(A^>) &= \{f \in \mathbb{R}^m : A^> f = 0\} \\ &= \bigcap_{i=1}^n \{f \in \mathbb{R}^m : e_i^> A^> f = 0\} \end{aligned}$$

Next, note $C_{S \setminus K}$ has at least one nonempty equivalence class, since $S \setminus K$ is nonempty by assumption. Rewriting eq. (2.A.8) as in eq. (2.A.6), and fixing an x_F within each nonempty

For $C_{S \setminus K}$, Lemma 2.A.13 grants

$$\begin{aligned}
 (2.A.8) &= \inf_{\substack{F \in \mathbb{R}^n \\ F \neq 0}} \inf_{\substack{v \in \mathbb{R}^n \\ \|v\|_2 = 1}} \frac{A^T v}{\|v\|_2} : \|v\|_2 = 1; \exists c > 0; \exists S \subseteq \mathbb{R}^n \quad P_{S \setminus K}^2(\cdot) = cv \\
 &= \min_{\substack{F \in \mathbb{R}^n \\ F \neq 0}} \inf_{\substack{v \in \mathbb{R}^n \\ \|v\|_2 = 1}} \frac{A^T v}{\|v\|_2} : \|v\|_2 = 1; \exists c > 0; \exists S \subseteq \mathbb{R}^n \quad P_{S \setminus K}^2(\cdot) = cv; P_{S \setminus K}^2(\cdot) \geq F \\
 &= \min_{\substack{F \in \mathbb{R}^n \\ F \neq 0}} \inf_{\substack{v \in \mathbb{R}^n \\ \|v\|_2 = 1}} \frac{A^T v}{\|v\|_2} : \|v\|_2 = 1; v \in N_F(x_F); \exists S \subseteq \mathbb{R}^n \quad g(x_F + v) \geq 0 :
 \end{aligned}$$

Since $S \cap \ker(A^T) \neq \emptyset$; and $S \setminus \ker(A^T)$, at least one in them has a nonempty domain (for the others, take the convention that their value is $+\infty$). Each in them with a nonempty domain in this final expression is of a continuous function over a compact set (in fact, a polyhedral cone intersected with the boundary of the unit l^2 ball), and thus it has a minimizer v , which corresponds to some $c \in P_{S \setminus K}^2(\cdot) \cap \ker(A^T)$, where $c > 0$. It follows that

$$A^T v = c A^T (\cdot) \in \ker(A^T);$$

meaning each of these in them is positive. But since S is polyhedral, C_S has finitely many equivalence classes $\{C_S^j \mid j \in \{1, \dots, H\}\}$, meaning the outer minimum is attained and positive. \square

Finally, as mentioned above, the desired result follows by norm equivalence.

Proof of Theorem 2.3.6. For the upper bound, note as in the proof of Lemma 2.A.7 that $S \setminus \ker(A^T) \neq \emptyset$; and the in hand is positive for every element of the domain, so the in them is finite. For the lower bound, by Lemma 2.A.7 and norm equivalence,

$$\begin{aligned}
 (A; S) &= \frac{\inf_{S \cap \ker(A^T)} \|A^T\|_{k_1}}{\inf_{S \setminus \ker(A^T)} \|A^T\|_{k_1}} \\
 &\geq \frac{1}{\inf_{S \cap \ker(A^T)} \|A^T\|_{k_2}} > 0: \quad \square
 \end{aligned}$$

Appendix 2.B Proof of 0-coercivity Characterization

Proof of Proposition 2.4.3. Suppose a minimizer x exists, and define

$$S_b := \{x + v : \|v\|_2 = b\};$$

$$r_b := \inf_{x \in S_b} f(x);$$

For any $b > 0$, it cannot hold that $r_b < f(x)$, since then there exists x with $f(x) < f(x)$; furthermore, it also can not hold that $r_b = f(x)$, since the closure property only guarantees there exists x_b with $x_b \in S_b$ and $f(x_b) = r_b$, and by strict convexity, the point $x^0 := (x_b + x)/2$ has $f(x^0) < f(x)$. Next, convexity grants $2(r_b - f(x)) \leq f(x_2) - f(x)$, since for any x_2 with $\|x_2 - x\|_2 = 2b$,

$$r_b - f(x) \leq \frac{1}{2}(f(x_2) - f(x));$$

since x_2 was arbitrary, it holds for the infimum over all such x_2 , which gives the inequality. It follows that all sublevel sets are compact.

On the other hand, if all sublevel sets are compact, since f is closed and proper, a minimizer exists. \square

Appendix 2.C Bibliographic Notes

Bibliographic notes on hard cores appear in Section 1.E. Hints of the dual problem may be found in many works, most notably those of Kivinen and Warmuth (1999), Collins et al. (2002), which demonstrated that boosting is seeking a difficult distribution over training examples via iterated Bregman projections.

Acknowledgements

This chapter is based on work by the dissertation author which appeared in the Journal of Machine Learning Research, Volume 13, pages 561–606, 2012.

Chapter 3

Optimization Guarantees

This chapter develops rates of numerical convergence for boost under the setting of Chapter 2; namely, only a finite sample is considered (meaning the measure is \mathbb{P} and the surrogate risk is \mathbb{L}), and the image set $\{h(x_i)\}_{i=1}^m : h \in H$ is assumed finite, whereby A may simply be written as a matrix $A \in \mathbb{R}^{m \times n}$.

Looking forward, if H is simply finite, then the numerical convergence rates here can be combined with the statistical guarantees of Chapter 5. Developing a statistical analysis from the rates here for infinite H seems challenging, however, as the bounds here grow quickly with the sample size m . Consequently, Chapter 6, which provides guarantees for infinite H , will provide a separate numerical analysis; note however that the numerical rates there are much slower than the ones here.

Chapter 6 and the present chapter not only disagree on the size of n which they can handle; the later chapter elects to handle only certain Lipschitz losses, for instance the logistic loss ℓ_{\log} . Interestingly, that analysis fails in a number of places for the exponential loss, and meanwhile the analysis here has terrible constants growing quickly with the sample size in the case of the logistic loss. As such, the mechanisms in this thesis are still not adequately sensitive to these losses.

3.1 Overview

Convergence rates will be proved for the following family of loss functions, which we denote \mathcal{L}_{fs} by requiring some derivative bounds over any level set.

Definition 3.1.1. \mathcal{L}_{fs} contains all functions ℓ satisfying the following properties. First, $\ell \in \mathcal{L}_{fs}$. Second, for any $x \in \mathbb{R}^m$ satisfying $\mathbb{L}(x) = \mathbb{L}(A_0) = \ell(0)$, and for any coordinate $(x)_i$, there exist

constants $\gamma > 0$ and $\beta > 0$ such that $\phi^0(x)_i = \gamma \phi(x)_i$ and $\phi^1(x)_i = \beta \phi(x)_i$.

The exponential loss ϕ_{exp} is in this family with $\gamma = \beta = 1$ since $\exp(\cdot)$ is a fixed point with respect to the differentiation operator. Furthermore, as is verified in Remark 3.A.1, the logistic loss ϕ_{log} is also in this family, with $\gamma = 2^m = (m \ln(2))$ and $\beta = 1 + 2^m$ (which may be loose). In a sense, γ and β encode how similar some $2 L_{\text{fso}}$ is to the exponential loss, and thus these parameters can degrade radically. However, outside the weak learnability case, the other terms in the bounds here can also incur a large penalty with the exponential loss, and there is some evidence that this is unavoidable (see the lower bounds in Mukherjee et al. (2011) or the upper bounds in Ratsch et al. (2001)).

The rates are presented in three separate sections, mirroring the optimization structure in Chapter 2. First, Sections 3.2 and 3.3 respectively consider weak learnable and attainable instances (i.e., instances where the hard core is respectively empty or contains every example), which entails that the rate is $O(\ln(1/\epsilon))$. The general case is presented in Section 3.4, where the rate degrades to $O(1/\epsilon)$ (with a matching lower bound for ϕ_{log}); this degradation is due to the preceding cases being fast for very different reasons, and as will be discussed after the proof, these different reasons come into conflict in the general case. In each section, there will be some discussion on how to slightly weaken the conditions on L_{fso} , however in each case this will lead to degraded rates. For example, the relationship between function values and derivative values in the definition of L_{fso} may seem unusual, however it will be the key to the $O(\ln(1/\epsilon))$ rate under the empirical weak learning assumption.

Each proof will use the following basic guarantee on the progress granted by a single step of the algorithm. This fact mostly follows from standard guarantees of line search methods, however it also leverages basic properties of $(A; S)$, exactly as presented in eq. (2.3.5); namely, $(A; S)$ allows $A^>$ to be removed from a normed expression.

Proposition 3.1.2. For any $t, \epsilon \in [0, 1]$, $A \in \mathbb{R}^{m \times n}$, and S for $\phi(A_t)$ with $(A; S) > 0$,

$$\phi(A_{t+1}) - \phi(A) \leq \phi(A_t) - \phi(A) - \frac{\epsilon^2 (A; S)^2 D_{S \setminus \ker(A^>)}^1 (r \phi(A_t))^2}{6 \phi(A_t)};$$

Proof of Proposition 3.1.2. If ϵ_t is optimal, then $r \phi(A_t) = 0$, and moreover since every line search option performs descent, $\epsilon_{t+1} = \epsilon_t$, whereby the desired guarantee holds with equality (potentially using the convention $0/0=1$ in the second term). As such, suppose in the remainder

of the proof that ℓ_t is not optimal.

Proceeding as in eq. (2.3.4), since $\ell_t \in \text{ker}(A_t)$,

$$(A; S) = \inf_{\substack{2 \in \text{ker}(A_t) \\ \ell_t \in \text{ker}(A_t)}} \frac{k_A > k_1}{D_{S \setminus \text{ker}(A_t)}^1(\ell_t)} = \frac{k_A > r \ell_t k_1}{D_{S \setminus \text{ker}(A_t)}^1(r \ell_t)}:$$

By Lemma 3.A.2, $\ell_t \in A$ has Lipschitz gradients with parameter ℓ_t with respect to $L^1(b)$ over the $(t+1)^{\text{st}}$ level set, whereby the preconditions for the guarantees for all three line searches in Lemmas 1.C.3 and 1.C.7 are met, yielding

$$\ell_t \in \ell_{t+1} \leq \frac{2k_A > r \ell_t k_1^2}{6 \ell_t} + \frac{(A; S)^2 D_{S \setminus \text{ker}(A_t)}^1(r \ell_t)^2}{6 \ell_t}:$$

Subtracting ℓ_t from both sides and rearranging yields the statement. \square

From here, the three different cases will manage the dual distance $D_{S \setminus \text{ker}(A_t)}^1(r \ell_t)$, specifically producing a relation to $\ell_t \in \ell_t$, the total suboptimality in the preceding iteration; from there, standard tools in convex optimization will yield convergence rates. Matching the problem structure revealed in Section 2.4, first the extremal cases of weak learnability and attainability will be handled, and only then the general case.

3.2 Weak Learnability

First consider the case that the empirical weak learning assumption holds, which implies the dual set D_A is exactly $\{0\}$ by Proposition 2.3.2, which in turn implies the hard core C is empty.

Theorem 3.2.1. Suppose $A \in \mathbb{R}^{m \times n}$, and the corresponding hard core C satisfies $|C| = 0$. Let any $\ell \in L_{\text{fso}}$ be given. Then $(A; \mathbb{R}_+^m) > 0$, and for any $t \geq 0$,

$$\ell_t \in \ell_0 \leq 1 + \frac{(A; \mathbb{R}_+^m)^2}{6 \ell^2} t:$$

Proof. By Theorem 2.4.1, $D_A = \{0\}$, meaning

$$D_{D_A}^1(r \ell_t) = \inf_{\substack{2 \in D_A \\ \ell_t \in D_A}} \text{kr} \ell_t \quad k_1 = \text{kr} \ell_t k_1 \quad \ell_t = :$$

Next, R_+^m is polyhedral, and Theorem 2.4.1 grants $R_+^m \setminus \ker(A^>) \neq \emptyset$; and $R_+^m \cap \ker(A^>) \neq \emptyset$; so Theorem 2.3.6 provides $(A; R_+^m) > 0$. Since $\mathbb{L}(A_t) \in R_+^m$, all conditions of Proposition 3.1.2 are met, and using $\mathbb{L}(A) = 0$ (again by Theorem 2.4.1),

$$\mathbb{L}(A_{t+1}) - \mathbb{L}(A_t) \leq \frac{2(A; R_+^m)^2 f(A_t)^2}{6^2 \mathbb{L}(A_t)} = \mathbb{L}(A_t) \left(1 - \frac{2(A; R_+^m)^2}{6^2} \right); \quad (3.2.2)$$

and recursively applying this inequality yields the result. \square

Specializing this analysis to the exponential loss (where $\sigma = 1$), assuming that a best choice $h_t \in H$ exists (whereby $\sigma = 1$), and recalling from Section 2.3 that $(A; R_+^m) = (b)$ (the empirical weak learning rate), the bound becomes $(1 - (b)^2/6)^t$. This bound in turn recovers the bound due to Schapire and Singer (1999), albeit with the denominator degraded from 2 to 6 (which can be recovered by choosing only step size options 1 & 2); this is despite the analysis here going through a very different path.

In general, still making use of $(A; R_+^m) = (b)$, solving for t in the expression

$$= \frac{\mathbb{L}(A_t) - \mathbb{L}(A)}{\mathbb{L}(A_0) - \mathbb{L}(A)} \leq \left(1 - \frac{2(b)^2}{6^2} \right)^t \exp \left(\frac{t^2 (b)^2}{6^2} \right)$$

reveals that $t \leq \frac{6^2}{2(b)^2} \ln(1/\epsilon)$ iterations suffice to reach suboptimality ϵ . Recall that σ and γ , in the case of the logistic loss, have only been bounded by quantities like $\epsilon^m/2$. While it is unclear if this analysis of σ and γ was tight, note that it is plausible that the logistic loss is slower than the exponential loss in this scenario, as it works less in initial phases to correct minor margin violations. As stated previously, however, it may simply be the case that the mechanism used here to simultaneously analyze the logistic and exponential losses is simply insensitive to the logistic loss.

Remark 3.2.3. The rate $O(\ln(1/\epsilon))$ depended crucially on both $\sigma \leq \epsilon^0$ and $\gamma \leq \epsilon^{00}$. If for instance the second inequality were replaced with $\gamma \leq C$, then eq. (3.2.2) would instead have form $\mathbb{L}(A_{t+1}) - \mathbb{L}(A_t) \leq \mathbb{L}(A_t)^2 O(1)$, which by an application of standard convex optimization tools would grant a rate $O(1/\epsilon)$ (Shalev-Shwartz and Singer, 2008, Lemma 20). For functions which asymptote to zero (i.e., everything in L_{fs}), satisfying this milder second order condition is quite easy. The real mechanism behind producing a fast rate is $\sigma \leq \epsilon^0$, which guarantees that the flattening of the objective function is concomitant with low objective values.

3.3 Attainability

Consider now the case of attainability. Recall from Theorem 2.4.4 and Proposition 2.4.3 that attainability occurred along with a stronger property: the function $\mathcal{L} + \text{im}(A)$ is 0-coercive, meaning it has compact level sets. Recall further that it was necessary to restrict \mathcal{L} to the slice $\text{im}(A)$ via $\mathcal{L} + \text{im}(A)$, since A may have linearly dependent columns, which immediately destroys any hope of compact level sets.

This level set structure has an immediate consequence to the task of relating $\mathcal{L}(A - t)$ to the dual distance $D_{S \setminus \ker(A)}^1(r, \mathcal{L}(A - t))$. \mathcal{L} is a strictly convex function, which means it is strongly convex over any compact set. Strong convexity in the primal corresponds to upper bounds on second derivatives in the dual (occasionally termed strong smoothness), which in turn can be used to relate distance and objective values. Making these statements rigorous will also provide the choice of polyhedron $S \in (A; S)$: unlike the case of weak learnability which chose $S = \mathbb{R}_+^m$, a compact set containing the initial level set will be used.

Theorem 3.3.1. Suppose $A \in \mathbb{R}^{m \times n}$, and the corresponding hard core C satisfies $|C| = m$. Let any $\epsilon \in \mathbb{R}_{fso}$ be given. Then there exists a (compact) tightest axis-aligned rectangle \mathcal{C} containing the initial level set $f \times \mathbb{R}^m : (\mathcal{L} + \text{im}(A))(x) \leq \mathcal{L}(A_0) + \epsilon$, and \mathcal{L} is strongly convex with modulus $c > 0$ over C . Finally, $(A; r, \mathcal{L}(C)) > 0$, and for all $\epsilon > 0$ and t ,

$$\mathcal{L}(A - t) - \mathcal{L}(A) \leq (\mathcal{L}(A_0) - \mathcal{L}(A)) + 1 - \frac{c^2 (A; r, \mathcal{L}(C))^{2t}}{3 \mathcal{L}(A_0)}.$$

As in Section 3.2, when ϵ_0 is suboptimal, this bound may be rearranged to say that $t \geq \frac{3 \mathcal{L}(A_0)}{c (A; r, \mathcal{L}(C))^2} \ln(1/\epsilon)$ iterations suffice to reach suboptimality.

To make sense of this bound and its proof, the essential object is \mathcal{C} , whose properties are captured in the following lemma, which is stated with some slight generality in order to allow reuse in Section 3.4.

Lemma 3.3.2. Let $A \in \mathbb{R}^{m \times n}$, corresponding hard core C with $|C| = m$, any $\epsilon \in \mathbb{R}_{fso}$, and any $d = \inf \mathcal{L}(A)$ be given. Then there exists a (compact nonempty) tightest axis-aligned rectangle $C = f \times \mathbb{R}^m : (\mathcal{L} + \text{im}(A))(x) \leq d + \epsilon$. Furthermore, the dual imager $\mathcal{L}(C) \subseteq \mathbb{R}^m$ is also a (compact nonempty) axis-aligned rectangle, and moreover $\mathcal{L}(C) \subseteq \text{int}(\text{dom}(\mathcal{L})) \subseteq \mathbb{R}_+^m$. Finally, $(r, \mathcal{L}(C))$ contains dual feasible points (i.e., $(r, \mathcal{L}(C)) \in D_A(\mathcal{L}; \epsilon)$).

A full proof may be found in Section 3.A.3; the principle is that $\|C\| = m$ provides 0-coercivity of $\mathbb{D}^+_{\text{im}(A)}$, and thus the initial level set is compact. To later show $\langle A; S \rangle > 0$ via Theorem 2.3.6, S must be polyhedral, and to apply Proposition 3.1.2, it must contain the dual iterates from $\mathbb{D}^+(A)g_{t=1}^1$; the easiest choice then is to take the bounding box \mathcal{C} of the initial level set, and use its dual map $\mathbb{D}(\mathcal{C})$. To exhibit dual feasible points within $r \mathbb{D}(\mathcal{C})$, note that \mathcal{C} will contain a primal minimizer, and optimality conditions grant that $r \mathbb{D}(\mathcal{C})$ contains the dual optimum.

With the polyhedron in place, Proposition 3.1.2 may be applied, so what remains is to control the dual distance. Again, this result will be stated with some extra generality in order to allow reuse in Section 3.4.

Lemma 3.3.3. Let $A \in \mathbb{R}^{m \times n}$, any $\gamma \in L_{\text{fso}}$, and any compact set S with $r \mathbb{D}(S) \setminus \ker(A^\top) \neq \emptyset$; be given. Then \mathbb{D} is strongly convex over S , and taking $c > 0$ to be the modulus of strong convexity, for any $x \in S \setminus \text{im}(A)$,

$$\mathbb{D}(x) - \mathbb{D}(A) \geq \frac{1}{2c} \inf_{r \mathbb{D}(S) \setminus \ker(A^\top)} \|r \mathbb{D}(x) - k_1^2\|$$

Before presenting the proof, it can be sketched quite easily. Using the Fenchel-Young inequality and the dual optimization problem (cf. Theorem 2.2.1), primal suboptimality can be converted into a Bregman divergence in the dual. If there is strong convexity in the primal, it allows this Bregman divergence to be converted into a distance via standard tools in convex optimization (Shalev-Shwartz and Singer, 2008, Lemma 18). Although \mathbb{D} lacks strong convexity in general, it is strongly convex over any compact set.

Proof of Lemma 3.3.3. Consider the optimization problem

$$\inf_{x \in S} \inf_{\substack{D \\ k_2=1}} \|r \mathbb{D}(x)\|; \quad E = m^{-1} \inf_{x \in S} \inf_{\substack{D \\ k_2=1}} \sum_{i=1}^m \chi^m(\cdot)^2; \quad \chi^m(\cdot) \geq 0;$$

since S is compact and $\chi^m(\cdot)^2$ and $\|\cdot\|^2$ are continuous, the in mum is attained. But $\chi^m > 0$ and $\neq 0$, meaning the in mum c is nonzero, and moreover it is the modulus of strong convexity of \mathbb{D} over S (Hiriart-Urruty and Lemaréchal, 2001, Theorem B.4.3.1.iii).

Now let any $x \in S \setminus \text{im}(A)$ be given, define $D = r \mathbb{D}(S) \cap \mathbb{R}_+^m$, and for convenience set $K := \ker(A^\top)$. Consider the dual element $P_{D \setminus K}^2(r \mathbb{D}(x))$ (which exists since $D \setminus K \neq \emptyset$); due to

the projection, it is dual feasible, and thus it must follow from Theorem 2.2.1 and the form of \mathfrak{L} in Proposition 1.B.2 that

$$\mathfrak{L}(A) = \sup_{x \in \mathbb{R}^m} \mathfrak{L}(x) : \mathfrak{L}(x) \in \mathcal{P}_{D \setminus K}^2(r \mathfrak{L}(x)) :$$

Furthermore, since $x \in \text{im}(A)$,

$$\mathfrak{L}(x) - \mathfrak{L}(A) = 0 :$$

Combined with the Fenchel-Young inequality and the general property $\mathfrak{L}(x) = r \mathfrak{L}(r \mathfrak{L}(x))$ (Rockafellar, 1970, Theorem 23.5),

$$\begin{aligned} \mathfrak{L}(x) - \mathfrak{L}(A) &= \mathfrak{L}(x) - \mathfrak{L}(A) + \mathfrak{L}(P_{D \setminus K}^2(r \mathfrak{L}(x))) \\ &= \mathfrak{L}(P_{D \setminus K}^2(r \mathfrak{L}(x))) - \mathfrak{L}(A) + \mathfrak{L}(r \mathfrak{L}(x)) - \mathfrak{L}(A) \\ &= \mathfrak{L}(P_{D \setminus K}^2(r \mathfrak{L}(x))) - \mathfrak{L}(r \mathfrak{L}(x)) + \mathfrak{L}(r \mathfrak{L}(x)) - \mathfrak{L}(A) \\ &= \mathfrak{L}(P_{D \setminus K}^2(r \mathfrak{L}(x))) - \mathfrak{L}(r \mathfrak{L}(x)) + \mathfrak{L}(r \mathfrak{L}(x)) - \mathfrak{L}(A) \end{aligned} \tag{3.3.4}$$

$$\frac{1}{2c} \|r \mathfrak{L}(x) - P_{D \setminus K}^2(r \mathfrak{L}(x))\|_2^2 ; \tag{3.3.5}$$

where the last step follows by a property of Bregman divergences in the dual (Shalev-Shwartz and Singer, 2008, Lemma 18), noting that both $r \mathfrak{L}(x)$ and $P_{D \setminus K}^2(r \mathfrak{L}(x))$ are in $r \mathfrak{L}(S) = D$, and \mathfrak{L} is strongly convex with modulus c over S . To finish, rewrite P as an infimum and use $\|k\|_2 \leq \|k\|_1$. \square

The desired result now follows readily.

Proof of Theorem 3.3.1. Invoking Lemma 3.3.2 with $d = \mathfrak{L}(A_0)$ immediately provides a compact tightest axis-aligned rectangle C containing the initial level set $S := \{x \in \mathbb{R}^m : (\mathfrak{L} + \text{im}(A))(x) \leq \mathfrak{L}(A_0)\}$. Crucially, since the objective values never increase S and C contain every iterate $f(A_t)_{t=1}^1$.

Applying Lemma 3.3.3 to the set C (by Lemma 3.3.2, $r \mathfrak{L}(C) \setminus \ker(A^T) \in \mathcal{C}$), then for any t ,

$$\mathfrak{L}(A_t) - \mathfrak{L}(A) \leq \frac{1}{2c} \|r \mathfrak{L}(A_t) - P_{r \mathfrak{L}(C) \setminus \ker(A^T)}^1(r \mathfrak{L}(A_t))\|_2^2 ;$$

where $c > 0$ is the modulus of strong convexity of \mathfrak{L} over C .

Finally, if there are suboptimal iterates, then $r \mathcal{L}(C) \cap r \mathcal{L}(S)$ contains points that are not dual feasible, meaning $r \mathcal{L}(C) \cap \ker(A^T) \neq \emptyset$; since Lemma 3.3.2 also provided $r \mathcal{L}(C) \cap D_A \neq \emptyset$; and $r \mathcal{L}(C)$ is a hypercube, it follows by Theorem 2.3.6 that $\langle A; r \mathcal{L}(C) \rangle > 0$. Plugging this into Proposition 3.1.2 and using $\mathcal{L}(A_t) \subseteq \mathcal{L}(A_0)$ gives

$$\mathcal{L}(A_{t+1}) \subseteq \mathcal{L}(A) \cap \mathcal{L}(A_t) \subseteq \mathcal{L}(A) \cap \frac{2 \langle A; r \mathcal{L}(C) \rangle^2 D_{r \mathcal{L}(C) \setminus \ker(A^T)}^1 (r \mathcal{L}(A_t))^2}{6 \mathcal{L}(A_t)} \\ (\mathcal{L}(A_t) \cap \mathcal{L}(A)) \subseteq \frac{c^2 \langle A; r \mathcal{L}(C) \rangle^2}{3 \mathcal{L}(A_0)};$$

and the result again follows by recursively applying this inequality. \square

Remark 3.3.6. The key conditions on \mathcal{L}_{fso} , namely the existence of constants granting $\mathcal{L} \cap \mathcal{L}^0$ and $\mathcal{L}^{\infty} \subseteq \mathcal{L}$ within the initial level set, are much more than are needed in this setting. Inspecting the presented proofs, it entirely suffices that on any compact set in \mathbb{R}^m , \mathcal{L} has quadratic upper and lower bounds (equivalently, bounds on the smallest and largest eigenvalues of the Hessian), which are precisely the weaker conditions used in previous treatments (Bickel et al., 2006, Ratsch et al., 2001).

These quantities are therefore necessary for controlling convergence under weak learnability. To see how the proofs of this section break down in that setting, consider the central Bregman divergence expression in eq. (3.3.4). What is really granted by attainability is that every iterate lies well within the interior of $\text{dom}(\mathcal{L})$, and therefore these Bregman divergences, which depend on $r \mathcal{L}$, can not become too wild. On the other hand, with weak learnability, all dual weights go to zero (cf. Theorem 2.4.1), which means that $\mathcal{L} \rightarrow \mathcal{L}^{\infty}$, and thus the upper bound in eq. (3.3.5) ceases to be valid. As such, another mechanism is required to control this scenario, which is precisely the role of $\mathcal{L} \cap \mathcal{L}^0$ and $\mathcal{L}^{\infty} \subseteq \mathcal{L}$.

3.4 General Setting

The key development of Section 2.4.3 was that general instances may be decomposed uniquely into two smaller pieces, one satisfying attainability and the other satisfying weak learnability, and that these smaller problems behave somewhat independently. This independence is leveraged here to produce convergence rates relying upon the existing rate analysis for the attainable and weak learnable cases. The mechanism of the proof is as straightforward as one

could hope for: decompose the dual distance into the two pieces, handle them separately using preceding results, and then stitch them back together.

Theorem 3.4.1. Let $A \in \mathbb{R}^{m \times n}$, corresponding hard core C with $1 \leq j \leq m$, and any loss $\ell \in L_{fso}$ be given. Recall from Section 2.4.3 the partition of the rows of A into $A_0 \in \mathbb{R}^{m_0 \times n}$ and $A_+ \in \mathbb{R}^{m_+ \times n}$, and suppose the axes of \mathbb{R}^m are ordered so that $A = \begin{pmatrix} A_0 \\ A_+ \end{pmatrix}$. Set C_+ to be the tightest axis-aligned rectangle $C_+ \subseteq \mathbb{R}^{m_+} : (\ell + \text{im}(A_+))(x) \leq \ell(A_0)g$, and $w := \sup_{x \in C_+} \text{kr} \ell(A_+ \cdot x) - \text{P}_{\text{r} \ell(C_+) \setminus \ker(A_+)}^1(\text{r} \ell(A_+ \cdot x))k_1$. Then C_+ is compact, $w < 1$, ℓ has modulus of strong convexity $c > 0$ over C_+ , and $(A; \mathbb{R}^{m_0} \text{r} \ell(C_+)) > 0$. Using these terms, for all t ,

$$\ell(A \cdot t) - \ell(A) \leq \frac{n \cdot 2\ell(A_0)}{(t+1) \min\{1; 2(A; \mathbb{R}^{m_0} \text{r} \ell(C_+))\}^2} + w = (3(1+w=2c))^2$$

The new term, w , appears when stitching together the two subproblems. For choices of $\ell \in L_{fso}$ where $\text{dom}(\ell)$ is a compact set, this value is easy to bound; for instance, the logistic loss, where $\text{dom}(\ell) = [0; 1]^m$, has $w = \sup_{x \in \text{dom}(\ell)} \text{kr} \ell(x) - 0k_1 = 1$ (since $0 \in \text{dom}(\ell)$). And with the exponential loss, taking $S := \{f \in \mathbb{R}^n : \ell(A \cdot f) \leq \ell(A_0)g\}$ to denote the initial level set, since 0 is always dual feasible,

$$w = \sup_{x \in S} \text{kr} \ell(A \cdot x)k_1 = \sup_{x \in S} \ell(A \cdot x) - \ell(A_0) = 1 :$$

Note that rearranging the rate from Theorem 3.4.1 will provide that $O(1/w)$ iterations suffice to reach suboptimality, whereas the earlier scenarios needed on $O(\ln(1/w))$ iterations. The exact location of the degradation will be pinpointed after the proof, and is related to the introduction of w .

Proof of Theorem 3.4.1. By Theorem 2.4.8, $m_+ \ell(A_+) = m \ell(A)$, and the form of ℓ gives $m \ell(A \cdot t) = m_0 \ell(A_0 \cdot t) + m_+ \ell(A_+ \cdot t)$, thus

$$m \ell(A \cdot t) - m \ell(A) = m_0 \ell(A_0 \cdot t) + m_+ \ell(A_+ \cdot t) - m_+ \ell(A_+) : \tag{3.4.2}$$

For the left term, since $\ell(x) = \sum_{j=1}^n \ell_j(x_j)$,

$$\ell(A_0 \cdot t) = \text{kr} \ell(A_0 \cdot t)k_1 = \text{kr} \ell(A_0 \cdot t) - \text{P}_{A_0}^1(\text{r} \ell(A_0 \cdot t))k_1; \tag{3.4.3}$$

which used the fact (from Theorem 2.4.8) that $\mathbb{A}_0 = f \mathbb{0}g$.

For the right term of eq. (3.4.2), recall from Theorem 2.4.8 that $\mathbb{L}^+_{\text{im}(A_+)}$ is 0-coercive, thus the level set $S_+ := \{x \in \mathbb{R}^{m_+} : (\mathbb{L}^+_{\text{im}(A_+)}) (x) \leq \mathbb{L}(A_0)g\}$ is compact. For all t , since $\mathbb{L} \geq 0$ and the objective values never increase,

$$m\mathbb{L}(A_0) - m\mathbb{L}(A_t) = m_0\mathbb{L}(A_0) + m_+\mathbb{L}(A_+ - t) - m_+\mathbb{L}(A_+ - t);$$

in particular, $A_+ - t \in S_+$. It is crucial that the level set compares against $\mathbb{L}(A_0)$ and not $\mathbb{L}(A_+)$.

Continuing, Lemma 3.3.2 may be applied to A_+ with value $d = \mathbb{L}(A_0)$, which grants a tightest axis-aligned rectangle $C_+ \subseteq \mathbb{R}^{m_+}$ containing S_+ , and moreover $\mathbb{L}(C_+) \setminus \ker(A_+^>) \neq \emptyset$. Applying Lemma 3.3.3 to A_+ and C_+ , \mathbb{L} is strongly convex with modulus $c > 0$ over C_+ , and for any t ,

$$\mathbb{L}(A_+ - t) - \mathbb{L}(A_+) \leq \frac{1}{2c} \text{kr} \mathbb{L}(A_+ - t) - P_{\mathbb{L}(C_+) \setminus \ker(A_+^>)}^1(\text{r} \mathbb{L}(A_+ - t))k_1^2; \quad (3.4.4)$$

Next, set $w := \sup_t \text{kr} \mathbb{L}(A_+ - t) - P_{\mathbb{L}(C_+) \setminus \ker(A_+^>)}^1(\text{r} \mathbb{L}(A_+ - t))k_1$; $w < 1$ since S_+ is compact and $\mathbb{L}(C_+) \setminus \ker(A_+^>)$ is nonempty. By the definition of w ,

$$D_{\mathbb{L}(C_+) \setminus \ker(A_+^>)}^1(\text{r} \mathbb{L}(A_+ - t))^2 \leq w D_{\mathbb{L}(C_+) \setminus \ker(A_+^>)}^1(\text{r} \mathbb{L}(A_+ - t));$$

which combined with eq. (3.4.4) yields

$$\mathbb{L}(A_+ - t) - \mathbb{L}(A_+) \leq \frac{w}{2c} D_{\mathbb{L}(C_+) \setminus \ker(A_+^>)}^1(\text{r} \mathbb{L}(A_+ - t)); \quad (3.4.5)$$

To merge the subproblem dual distance upper bounds eq. (3.4.3) and eq. (3.4.5) via Lemma 3.A.3, it must be shown that $(\mathbb{R}_+^{m_0} \times \mathbb{L}(C_+)) \setminus D_A \neq \emptyset$. But this follows by construction and Theorem 2.4.8, since $f \mathbb{0}g = \mathbb{A}_0 \in \mathbb{R}_+^m$, $\text{r} \mathbb{L}(C_+) \setminus \mathbb{A}_+ \neq \emptyset$; by Lemma 3.3.2, and the decomposition $D_A = \mathbb{A}_0 \times \mathbb{A}_+$. Returning to the total suboptimality expression eq. (3.4.2),

these dual distance bounds yield

$$m\mathbb{L}(A_t) - m\mathbb{L}(A) = m_0 D_{A_0}^1(r\mathbb{L}(A_0 - t)) + m_+ w=(2c) D_{r\mathbb{L}(C_+) \setminus \ker(A_+)}^1(r\mathbb{L}(A_+ - t)) \\ - m_+ w=(2c) D_{(R_+^{m_0} \cap r\mathbb{L}(C_+) \setminus \ker(A_+))}^1(r\mathbb{L}(A_t));$$

the second step using Lemma 3.A.3.

To finish, note $R_+^{m_0} \cap r\mathbb{L}(C_+)$ is polyhedral, and

$$(R_+^{m_0} \cap r\mathbb{L}(C_+)) \cap \ker(A_+) \neq \emptyset \quad \text{for } \mathbb{L}(A_t) \geq_{t=1} \ker(A_+) \neq \emptyset;$$

since no primal iterate is optimal and thus $r\mathbb{L}(A_t)$ is not dual feasible by optimality conditions; combined with the above derivation $(R_+^{m_0} \cap r\mathbb{L}(C_+)) \setminus D_A \neq \emptyset$; Theorem 2.3.6 may be applied, meaning $(A; R_+^{m_0} \cap r\mathbb{L}(C_+)) > 0$. As such, all conditions of Proposition 3.1.2 are met, and making use of $\mathbb{L}(A_t) - \mathbb{L}(A_0)$,

$$\mathbb{L}(A_{t+1}) - \mathbb{L}(A) = \mathbb{L}(A_t) - \mathbb{L}(A) \\ + \frac{2(A; R_+^{m_0} \cap r\mathbb{L}(C_+))^2 D_{(R_+^{m_0} \cap r\mathbb{L}(C_+) \setminus \ker(A_+))}^1(r\mathbb{L}(A_t))^2}{6\mathbb{L}(A_t)} \\ - \frac{2(A; R_+^{m_0} \cap r\mathbb{L}(C_+))^2 (\mathbb{L}(A_t) - \mathbb{L}(A))^2}{6\mathbb{L}(A_0) (w=(2c))^2}.$$

Applying a standard conversion from convex optimization, for instance as stated by Shalev-Shwartz and Singer (2008, Lemma 20) and invoking it with parameters

$$t := \frac{\mathbb{L}(A_t) - \mathbb{L}(A)}{\mathbb{L}(A_0)} \quad \text{and} \quad r := \frac{1}{2} \min \left(1; \frac{2(A; R_+^{m_0} \cap r\mathbb{L}(C_+))^2}{3(w=(2c))^2} \right);$$

the result follows. \square

In order to produce a rate $O(\ln(1/\epsilon))$ under attainability, strong convexity related the suboptimality to a squared dual distance k_1^2 (cf. eq. (3.3.5)). On the other hand, the rate $O(\ln(1/\epsilon))$ under weak learnability came from a fortuitous cancellation with the denominator $\mathbb{L}(A_t)$ (cf. eq. (3.2.2)), which is equal to the total suboptimality since Theorem 2.4.1 provides $\mathbb{L}(A) = 0$. But in order to merge the subproblem dual distances via Lemma 3.A.3, the di ering

properties granting fast rates must be ignored. (In the case of attainability, this process introduces w.)

This incompatibility is not merely an artifact of the analysis. Intuitively, the finite and infinite margins sought by the two pieces $A_0; A_+$ are in conflict. For a beautifully simple, concrete case of this, consider the following matrix, due to Schapire (2010):

$$S := \begin{pmatrix} 2 & 3 \\ 6 & 1 \\ 8 & +1 \\ 4 & 1 \\ 1 & 1 \end{pmatrix} \begin{matrix} z \\ z \\ z \\ z \\ z \end{matrix}$$

The optimal solution here is to push both coordinates of x unboundedly positive, with margins approaching $(0; 1)$. But pushing any coordinate $(x)_i$ too quickly will increase the objective value, rather than decreasing it. In fact, this instance will provide a lower bound, and the mechanism of the proof shows that the primal weights grow extremely slowly, as $\mathcal{O}(\ln(t))$.

Theorem 3.4.6. Consider $\ell_{\log} = \ln(1 + \exp(\cdot))$ L_{fso} , the logistic loss, and suppose the line search is exact (cf. option 1 in Figure 1.3). Then for any $t \geq 1$, $\mathcal{L}(S_t) - \mathcal{L}(S) \leq 1/(24t)$.

Finally, note that this third setting does not always entail slow convergence. Again taking the view of the rows of S being points $f_{s_i} g_{i=1}^3$, consider the effect of rotating the entire instance around the origin by $\pi/4$. The optimization scenario is unchanged, however coordinate descent can now be arbitrarily close to the optimum in one iteration by pushing a single primal weight extremely high.

Appendix 3.A Miscellaneous Technical Material

3.A.1 The Logistic Loss is within L_{fso}

Remark 3.A.1. This remark develops bounds on the quantities \mathcal{L} for the logistic loss $\ell_{\log} = \ln(1 + \exp(\cdot))$. First note that the initial level set $S_0 := \{x \in \mathbb{R}^m : \mathcal{L}(x) \leq \mathcal{L}(A_0)\}$ is contained within a cube $(-1/b; b]^m$, where $b = m \ln(2)$; this follows since $\mathcal{L}(A_0) = \mathcal{L}(0) = \ln(2)$, whereas $\ell_{\log}(m \ln(2)) = \ln(1 + \exp(m \ln(2))) > m \ln(2)$.

For convenience, the analysis will be mainly written with respect to $b = m \ln(2)$. Let any $x \in (-1/b; b]$ be given, and note $\ell_{\log}^0 = \exp(\cdot) = (1 + \exp(\cdot))$, and $\ell_{\log}^{00} = \exp(\cdot) = (1 + \exp(\cdot))^2$.

To determine \mathcal{L} , note $1 - 1/(1 + \exp(x)) = 1/(1 + \exp(b))$. Since \ln is concave, it follows for all

$z \in [1; 1 + \exp(b)]$ that the secant line through $(1; 0)$ and $(1 + \exp(b); \ln(1 + \exp(b)))$ is a lower bound:

$$\ln(z) \geq \frac{\ln(1 + \exp(b)) - 0}{1 + \exp(b) - 1} (z - 1) = \frac{\ln(1 + \exp(b))}{1 + \exp(b)} (z - 1)$$

As such, for $x \in (1; b]$, $\ln(1 + \exp(x)) \geq \frac{\exp(x)}{1 + \exp(x)} \ln(1 + \exp(b))$, so

$$\frac{\ln(1 + \exp(x))}{\ln(1 + \exp(b))} \geq \frac{\exp(x)}{(1 + \exp(x))^2 \ln(1 + \exp(x))} \geq \frac{\exp(b)}{(1 + \exp(x))^2 \ln(1 + \exp(b))} \geq \frac{\exp(b)}{\ln(1 + \exp(b))}$$

Consequently, a sufficient choice is $\beta := \exp(b) = \ln(1 + \exp(b)) - 2^m = (m \ln(2))$.

For $\ln(x) \geq \frac{\beta}{x}$, using $\ln(x) \geq x - 1$,

$$\frac{\ln(x)}{\ln(1 + \exp(x))} \geq \frac{\ln(1 + \exp(x))}{\frac{\exp(x)}{1 + \exp(x)}} \geq \frac{\exp(x)}{1 + \exp(x)}$$

That is, it suffices to set $\beta := 1 + \exp(b) = 1 + 2^m$.

3.A.2 $L_{f_{\text{so}}}$ implies \mathcal{L}^b has Lipschitz Gradients

Lemma 3.A.2. Let $\mathcal{H} \in L_{f_{\text{so}}}$, any H (infinite is fine), and any empirical measure \mathbb{b} over $X \in \mathbb{R}^d; +1g$ be given. Then \mathcal{L}^b has Lipschitz gradients with parameter $\mathcal{L}(A_{t-1})$ with respect to norm $L^1(\mathbb{b})$ over the t^{th} level set

$$S_t := \left\{ \mathbb{b} \in \mathcal{L}(A_{t-1}) : \mathcal{L}(A_{t-1}) \leq \beta \right\}$$

Proof. Starting from Lemma 1.C.2, for any $\mathbb{b} \in S_t$, since $\mathcal{L}(A_{t-1}) \leq \beta$ (whereby $\mathcal{L}(A_{t-1}) \leq \beta$), and also

using the definition of the level set S_t ,

$$\begin{aligned} & r \inf_{x \in \mathbb{R}^n} (A_0 + r(A_1 - A_0))x \\ &= \frac{1}{m} \sum_{i=1}^m \int_0^1 \int_{\mathbb{R}^n} \lambda^i \left((A_0 + r(A_1 - A_0))x \right) dx \\ &= \frac{1}{m} \sum_{i=1}^m \int_0^1 \int_{\mathbb{R}^n} \lambda^i \left((A_0 + r(A_1 - A_0))x \right) dx \\ &= \frac{1}{m} \sum_{i=1}^m \int_0^1 \int_{\mathbb{R}^n} \lambda^i \left((A_0 + r(A_1 - A_0))x \right) dx \\ &= \frac{1}{m} \sum_{i=1}^m \int_0^1 \int_{\mathbb{R}^n} \lambda^i \left((A_0 + r(A_1 - A_0))x \right) dx \\ &= \frac{1}{m} \sum_{i=1}^m \int_0^1 \int_{\mathbb{R}^n} \lambda^i \left((A_0 + r(A_1 - A_0))x \right) dx \end{aligned}$$

as desired. □

3.A.3 Proof of Lemma 3.3.2

Proof of Lemma 3.3.2. Since $\inf \inf_{x \in \mathbb{R}^m} (L + \text{im}(A))(x)$ is nonempty. Since $|C| = m$, Theorem 2.4.4 provides $L + \text{im}(A)$ is 0-coercive, meaning S_d is compact.

Now consider the rectangle C defined as a product of intervals $C = \prod_{i=1}^m [a_i; b_i]$, where

$$a_i := \inf \{x_i : x \in S_d\}; \quad b_i := \sup \{x_i : x \in S_d\}.$$

By construction, $C \subseteq S_d$, and furthermore any smaller axis-aligned rectangle must violate some minimum or supremum above, and so must fail to include a piece of S_d . In particular, the tightest rectangle exists, and it is C .

Next, note that $r(L(x)) = (\lambda^1(x_1)=m; \lambda^2(x_2)=m; \dots; \lambda^m(x_m)=m)$, thus $D = \prod_{i=1}^m \lambda^i([a_i; b_i])=m$, an axis-aligned rectangle in the dual. Since L is strictly convex and $\text{dom}(L) = \mathbb{R}$, both $\lambda^i(a_i)$ and $\lambda^i(b_i)$ are within $\text{int}(\text{dom}(\lambda^i))$ (for all i), and so $r(L(C)) \subseteq \text{int}(\text{dom}(L))$.

Finally, Proposition 2.4.3 grants that $L + \text{im}(A)$ has a minimizer; thus choose any $x \in \mathbb{R}^n$ so that $L(A) = \inf L(A)$. By the optimality conditions in Theorem 2.2.1, $\lambda_A = r(L(A))$. But the dual optimum is dual feasible, and $A \in S_d$, so

$$r(L(C)) \setminus D_A \subseteq r(L(A)) \setminus D_A = \emptyset \quad \square$$

3.A.4 Splitting Distances along $A_0; A_+$

Lemma 3.A.3. Let $A = \begin{pmatrix} h & i \\ A_0 & A_+ \end{pmatrix}$ be given as in Theorem 3.4.1, and let a set $S = S_0 \cup S_+$ be given with $S_0 \subset \mathbb{R}^{m_0}$ and $S_+ \subset \mathbb{R}^{m_+}$ and $S \setminus D_A \in \mathcal{E}$; . Then, for any $\begin{pmatrix} 0 \\ + \end{pmatrix} \in \mathbb{R}^m$ with $0 \in \mathbb{R}^{m_0}$ and $+ \in \mathbb{R}^{m_+}$,

$$D_{S \setminus D_A}^1(\begin{pmatrix} 0 \\ + \end{pmatrix}) = D_{S_0 \setminus D_{A_0}}^1(0) + D_{S_+ \setminus D_{A_+}}^1(+):$$

Proof. Recall from Theorem 2.4.8 that $D_A = D_{A_0} \cup D_{A_+}$, thus

$$S \setminus D_A = (S_0 \setminus D_{A_0}) \cup (S_+ \setminus D_{A_+});$$

and $S \setminus D_A \in \mathcal{E}$; grants that $S_0 \setminus D_{A_0} \in \mathcal{E}$; and $S_+ \setminus D_{A_+} \in \mathcal{E}$; . Define now the notation $[]_0 : \mathbb{R}^m \rightarrow \mathbb{R}^{m_0}$ and $[]_+ : \mathbb{R}^m \rightarrow \mathbb{R}^{m_+}$, which respectively select the coordinates corresponding to the rows of A_0 , and the rows of A_+ .

Let $\begin{pmatrix} 0 \\ + \end{pmatrix} \in \mathbb{R}^m$ be given; in the above notation, $0 = []_0$ and $+ = []_+$. By the above Cartesian product and intersection properties,

$$\begin{pmatrix} 2 \\ 4 \end{pmatrix} P_{S_0 \setminus D_{A_0}}^1(0) \cap \begin{pmatrix} 3 \\ 5 \end{pmatrix} P_{S_+ \setminus D_{A_+}}^1(+)} \subset S \setminus D_A;$$

and so

$$D_{S \setminus D_A}^1(\begin{pmatrix} 0 \\ + \end{pmatrix}) = D_{S_0 \setminus D_{A_0}}^1(0) + D_{S_+ \setminus D_{A_+}}^1(+):$$

On the other hand, since $P_{S \setminus D_A}^1(\begin{pmatrix} 0 \\ + \end{pmatrix}) \subset (S_0 \setminus D_{A_0}) \cup (S_+ \setminus D_{A_+})$,

$$D_{S_0 \setminus D_{A_0}}^1(0) + D_{S_+ \setminus D_{A_+}}^1(+)} \supseteq 0 \cdot [P_{S \setminus D_A}^1(\begin{pmatrix} 0 \\ + \end{pmatrix})]_0 + + \cdot [P_{S \setminus D_A}^1(\begin{pmatrix} 0 \\ + \end{pmatrix})]_+ = D_{S \setminus D_A}^1(\begin{pmatrix} 0 \\ + \end{pmatrix}):$$

□

3.A.5 Proof of Theorem 3.4.6

Proof of Theorem 3.4.6. This proof proceeds in two stages: first the gap between any solution with l_1 norm B is shown to be large, and then it is shown that the l_1 norm of the Boost solution (under logistic loss) grows slowly.

To start, $\ker(S^>) = \{z(1; 1; 0) : z \in \mathbb{R}^g\}$, and \hat{z}_{\log} is maximized at $\hat{z}_{\log}^0(0)$ with value $\hat{z}_{\log}^0(0)$ (cf. Lemma 1.B.1). Thus $\hat{z}_S = (\hat{z}_{\log}^0(0); \hat{z}_{\log}^0(0); 0) = 3$, and $\mathcal{L}(S) = \mathcal{L}(\hat{z}_S) = 2\hat{z}_{\log}^0(0) = 3 = 2 \ln(2) = 3$.

Next, by calculus, given any B ,

$$\begin{aligned} \inf_{k, k_1, B} \mathcal{L}(S) &= \mathcal{L}(S) = \mathcal{L}(S) \Big|_{\substack{B=2 \\ B=2}} = 2 \ln(2) = 3 \\ &= (2 \ln(2) + \ln(1 + \exp(-B))) = 3 - 2 \ln(2) = 3 \\ &= \ln(1 + \exp(-B)) = 3: \end{aligned}$$

Now to bound the l_1 norm of the iterates. By the nature of exact line search, the coordinates of \hat{z} are updated in alternation (with arbitrary initial choice); thus let u_t denote the value of the coordinate updated in iteration t , and v_t be the one which is held fixed. (In particular, $v_t = u_{t-1}$.)

The objective function, written in terms of $(u_t; v_t)$, and ignoring the term $1 = m = 1 = 3$ which will be deleted after differentiation, is

$$\begin{aligned} & \ln(1 + \exp(v_t - u_t)) + \ln(1 + \exp(u_t - v_t)) + \ln(1 + \exp(-u_t - v_t)) \\ &= \ln(2 + \exp(v_t - u_t) + \exp(u_t - v_t) + 2 \exp(-u_t - v_t) + \exp(-2u_t) + \exp(-2v_t)) : \end{aligned}$$

Due to the use of exact line search, and the fact that u_t is the new value of the updated variable, the derivative with respect to u_t of the above expression must equal zero. In particular, producing this equality and multiplying both sides by the (nonzero) denominator yields

$$\exp(v_t - u_t) + \exp(u_t - v_t) - 2 \exp(-u_t - v_t) - 2 \exp(-2u_t) = 0 :$$

Multiplying by $\exp(u_t + v_t)$ and rearranging, it follows that, after line search, u_t and v_t must satisfy

$$\exp(2u_t) = \exp(2v_t) + 2 \exp(-u_t - v_t) + 2 : \quad (3.A.4)$$

First it will be shown for $t = 1$, by induction, that $u_t = v_t$. The base case follows by inspection (since $u_0 = v_0 = 0$ and so $u_1 = \ln(2)$). Now the inductive hypothesis grants $u_t = v_t$; the case $u_t = v_t$ can be directly handled by eq. (3.A.4), thus suppose $u_t > v_t$. But previously, it

was shown that the optimal l_1 bounded choice has both coordinates equal; as such, the current iterate, with coordinates $(u_t; v_t)$, is worse than the iterate $(u_t; u_t)$, and thus the line search will move in a positive direction, giving $u_{t+1} = v_{t+1}$.

It will now be shown by induction that, for $t \geq 1$, $u_t \leq \frac{1}{2} \ln(4t)$. The base case follows by the direct inspection above. Applying the inductive hypothesis to the update rule above, and recalling $v_{t+1} = u_t$ and that the weights increase (i.e., $u_{t+1} = v_{t+1} = u_t$),

$$\exp(2u_{t+1}) = \exp(2u_t) + 2 \exp(u_t - u_{t+1}) + 2 \exp(2u_t) + 2 \exp(u_t - u_t) + 2 \exp(4t + 4) - 4(t + 1):$$

To finish, recall by Taylor expansion that $\ln(1 + q) \geq q - \frac{q^2}{2}$; consequently for $t \geq 1$

$$\mathbb{E}(\log(S_{t+1})) - \mathbb{E}(\log(S_t)) \leq \inf_{k \in \mathcal{K}_1} \mathbb{E}(\log(S_{t+1})) - \mathbb{E}(\log(S_t)) \leq \frac{1}{3} \ln \left(1 + \frac{1}{4t} \right) \leq \frac{1}{12t} - \frac{1}{6} \frac{1}{4t^2} \leq \frac{1}{24t}:$$

□

Appendix 3.B Bibliographic Notes

The study of general numerical convergence properties of AdaBoost variants was initiated by Collins et al. (2002), which provided convergence for a variety of step sizes and losses. The first rate was $O(\exp(-\epsilon^2))$ to achieve accuracy > 0 , proved in a very general setting by Bickel et al. (2006); this work later served as the basis for the first consistency proof of AdaBoost, due to Bartlett and Traskin (2007).

In terms of rates under various assumptions, there is of course the original $O(\ln(1/\epsilon) \epsilon^{-2})$ rate due to Freund and Schapire (1997) when $(b) > 0$; a reasonable fast convergence rate was also shown with the logistic loss by Duchi and Helmbold (2000). On the other end of the spectrum, when minimizers exist, Ratsch et al. (2001) established a fast (Q-linear) convergence rate using techniques due to Luo and Tseng (1993). As discussed in Chapter 2, hard cores identify these two cases as extremal (hard core is empty or everything), and the general case is a combination of the two.

Parallel to the original work behind this thesis chapter, Mukherjee et al. (2011) established general convergence under the exponential loss, with a rate of (ϵ) . Just as with this chapter, their analysis split the problem into two pieces according to hard cores. Notably, they also proved

a $O(\epsilon^{-5})$ iteration bound to beat a predictor of some fixed norm. This second bound ends up being the essential one when controlling statistical properties in the general case, and a similar bound will be produced for the logistic loss (and similar losses) in Chapter 6, and moreover this bound was used to produce a streamlined consistency analysis of AdaBoost (Schapire and Freund, 2012, Theorem 12.2).

There are boosting variants with good convergence rates, for instance by Ratsch and Warmuth (2002), Warmuth et al. (2007), Shalev-Shwartz and Singer (2008). These variants, however, adjust the objective function, in particular adding regularization and constraints, and thus are incomparable. These works mainly focus on producing good margin guarantees; for similar guarantees for AdaBoost and related discussion, please see the work of Telgarsky (2013).

Acknowledgements

This chapter is based on work by the dissertation author which appeared in the Journal of Machine Learning Research, Volume 13, pages 561-606, 2012.

Part II

Statistical Behavior

Chapter 4

The Primal and the Dual

Part II of this thesis considers statistical issues, namely controlling the difference between \mathbb{E} and L as a function of the sample size and various properties of the hypothesis class H , loss ℓ , and source distribution μ .

To warm up to this problem, this chapter merely provides a strong duality analog to the weak duality result Proposition 1.4.1. Unlike the strong duality result for finite samples (cf. Theorem 2.2.1), this duality result requires the loss ℓ to be Lipschitz, and consequently this section discusses some of the difficulties presented by the exponential loss.

4.1 Strong Duality

Without further ado, the strong duality result follows.

Lemma 4.1.1. Let $\ell : \mathbb{R} \rightarrow \mathbb{R}_+$ be convex with $\lim_{z \rightarrow 1} \ell(z) = 0$ and finite tightest Lipschitz constant $L := \sup_{x \neq y} |\ell(x) - \ell(y)| / |x - y| < \infty$ be given. Additionally, let finite nonnegative measure μ over $X \times Y$ and hypothesis class H be given. Then

$$\inf_{h \in H} \int \ell(A_h) d\mu = \max_{\{p : p \in L^1(\mu); p \geq 0; \int p(x,y) d\mu = 0\}} \int \ell(p) d\mu \quad \mu\text{-a.e.}; \int p(x,y) d\mu = 0$$

Additionally, $\int \ell(p) d\mu$ is optimal if there exists $p \in L^1(\mu)$ such that $p(x,y) \leq \ell(A_h(x,y))$ for $\mu\text{-a.e.}$ (x,y) and p is feasible ($p \geq 0; \int p(x,y) d\mu = 0$).

Comparing Lemma 4.1.1 to Theorem 2.2.1, the latter, in addition to handling non-Lipschitz losses, has an artifact term. This in turn is a consequence of using the "wrong" dual pairing $\int p f_i = \sum_i p_i f_i$ rather than $\int p f_i = \sum_i p_i f_i$, a choice which was discussed in Remark 1.A.4.

In order to prove this duality result, the first step is to establish the convexity structure of L , which is analogous to the structure of \mathbb{D} presented in Proposition 1.B.2, although once again \mathbb{D} exhibits the artifact scaling term m .

Lemma 4.1.2. Let $\rho : \mathbb{R} \rightarrow \mathbb{R}_+$ be convex with $\lim_{z \rightarrow 1} \rho(z) = 0$ and finite tightest Lipschitz constant $L := \sup_{x \neq y} \frac{|\rho(x) - \rho(y)|}{|x - y|} < \infty$, and let μ be a probability measure over $Z := \{x \in \mathbb{R}^d : |x| \leq 1\}$.

1. If $f \in L^1(\mu)$, then $\int_{\mathbb{R}^d} (f(z)) d \mu(z)$ is well-defined and finite.
2. ρ^* is convex lower semi-continuous over $L^1(\mu)$.
3. Its conjugate $(\rho^*)^*$ is also convex lower semi-continuous as a function over $L^1(\mu)$.
4. If $p \in L^1(\mu)$, then $(\rho^*)^*(p) = \int_{\mathbb{R}^d} \rho^*(p) d \mu$, which is finite $(\rho^*)^*(p) < \infty$ if $p \in [0, \infty]$ -a.e..
5. If $f \in L^1(\mu)$, then $\int_{\mathbb{R}^d} \rho^*(f) d \mu < \infty$; and moreover $p \in \text{supp}(\rho^*(f))$ if $p(x) \in \text{supp}(f(z))$ for -a.e. z .

The proof of Lemma 4.1.2 can be expected to be somewhat similar to the analogous proofs for \mathbb{D} in Proposition 1.B.2, however measurability issues prevent the integrals being treated identically to the finite sums in \mathbb{D} . To circumvent measurability issues, the proof as usual must go through simple functions and limits thereof, which predictably adds some length.

From here, the proof of the strong duality rule in Lemma 4.1.1, proceeds by invoking an appropriately strengthened Fenchel duality rule, namely one strong enough to handle the Banach spaces in question (Zalinescu, 2002, Corollary 2.8.5 using condition (vii)). Note that this particular rule can not be applied if $\mu \in L^1(\mu)$ is replaced with, say, all weightings over H with finite support, which is not a complete normed vector space under the $L^1(\mu)$ norm.

Before turning to the complete proofs, the difficulties with the exponential loss stem from the fact that it is not in general finite over $L^1(\mu)$.

Proposition 4.1.3. Let μ denote the standard Gaussian measure over \mathbb{R}^d , and define $f(x) := x^2$ and $f_i(x) := x^2 \mathbb{1}[|x| \leq i]$. Then $f_i \in L^1(\mu)$, $f \in L^1(\mu)$, and $\|f_i - f\|_1 \rightarrow 0$, but

$$\int_{\mathbb{R}^d} \exp(f_i(x)) d \mu(x) < \infty \quad \text{and} \quad \int_{\mathbb{R}^d} \exp(f(x)) d \mu(x) = \infty :$$

Proof. To start, $\int_{\mathbb{R}} f d\mu = 1$ (variance of a standard Gaussian), and thus $\int_{\mathbb{R}} f_i d\mu \rightarrow \int_{\mathbb{R}} f d\mu$ by the monotone convergence theorem (and $\sum_{i=1}^{\infty} \int_{\mathbb{R}} f_i d\mu < \infty$). But

$$\begin{aligned} \int_{\mathbb{R}} \exp(f_i(x)) d\mu(x) &= \int_{\mathbb{R}} e^{i^2 x^2} d\mu(x) < 1; \\ \int_{\mathbb{R}} \exp(f(x)) d\mu(x) &= \int_{\mathbb{R}} \frac{1}{2} e^{x^2-2} dx = 1; \end{aligned}$$

□

In particular, a step in the proof of Lemma 4.1.1 breaks down when this is combined with the possibility of A not being a closed subspace (cf. Lemma 1.A.6). This does not mean the situation with μ_{\exp} is hopeless; rather, it is simply the case that $L^1(\mu)$ (and dual space $L^1(\mu)$) is the wrong lens with which to inspect \mathbb{R}_{\exp}^* .

Appendix 4.A Deferred Proofs

First, a proof of the convexity properties of \mathbb{R}_{\exp}^* .

Proof of Lemma 4.1.2. Let $f \in L^1(\mu)$ be arbitrary. Since μ is convex and finite, it is continuous, so $\mu \circ f$ is measurable, and moreover it is nonnegative thus $\int_{\mathbb{R}} \mu \circ f d\mu$ is well-defined. Additionally,

$$\begin{aligned} \int_{\mathbb{R}} \mu(f(z)) d\mu(z) &= \int_{\mathbb{R}} \mu(0) d\mu(z) + \int_{\mathbb{R}} (\mu(f(z)) - \mu(0)) d\mu(z) \\ &= \mu(0)(Z) + \int_{\mathbb{R}} |f(z) - 0| d\mu(z) = \mu(0)(Z) + \|f\|_1 < \infty; \end{aligned}$$

Next, for any $f_1, f_2 \in L^1(\mu)$ and $\lambda \in [0, 1]$,

$$\begin{aligned} \int_{\mathbb{R}} \mu(\lambda f_1(z) + (1-\lambda)f_2(z)) d\mu(z) &= \int_{\mathbb{R}} (\mu(\lambda f_1(z) + (1-\lambda)f_2(z))) d\mu(z) \\ &= \lambda \int_{\mathbb{R}} \mu(f_1(z)) d\mu(z) + (1-\lambda) \int_{\mathbb{R}} \mu(f_2(z)) d\mu(z); \end{aligned}$$

whereby \mathbb{R}_{\exp}^* is convex. Since it is finite over $L^1(\mu)$ (as above), it is necessarily lower semi-continuous.

Since \mathbb{R}_{\exp}^* is convex lower semi-continuous, so is its conjugate (Zalinescu, 2002, Theorem 2.3.3), where the dual space $L^1(\mu)$ is identified with $L^1(\mu)$ as per the isomorphism statements in Lemma 1.A.2.

The remainder of this proof will reason about the conjugate to \mathbb{R}^d . First let $p \in L^1(\mathbb{R}^d)$ be given with $\int_{\mathbb{R}^d} p(x) dx > 0$; it will follow that $\int_{\mathbb{R}^d} p(x) dx = 1$. Define the sets

$$S := p^{-1}((-\infty; 0)), \quad S_0 := p^{-1}([0; \infty]); \quad S_+ := p^{-1}((-\infty; 1]);$$

as well as, for every $c \in \mathbb{R}$, the reals

$$g_c \in \mathcal{C}(c); \quad g_0 \in \mathcal{C}(0); \quad g_+ \in \mathcal{C}(+c);$$

and lastly the simple functions

$$f_c(z) := c1_{[z \in S]} + 01_{[z \in S_0]} + c1_{[z \in S_+]};$$

$$g_c(z) := g_c^{-1}1_{[z \in S]} + g_0^{-1}1_{[z \in S_0]} + g_+^{-1}1_{[z \in S_+]};$$

By these choices f_c and g_c are measurable and within $L^1(\mathbb{R}^d)$, and moreover $g_c \in \mathcal{C}(f_c)$ everywhere.

As such,

$$\begin{aligned} \int_{\mathbb{R}^d} p(x) dx &= \sup_{f \in L^1(\mathbb{R}^d)} \int_{\mathbb{R}^d} (f(x) p(x)) dx \\ &= \sup_{c \in \mathbb{R}} \int_{\mathbb{R}^d} (f_c(x) p(x)) dx \\ &= \sup_{c \in \mathbb{R}} \left(\int_{S} (f_c(x) p(x)) dx + \int_{S_0} (f_c(x) p(x)) dx + \int_{S_+} (f_c(x) p(x)) dx \right) \\ &= \sup_{c \in \mathbb{R}} \left(c \int_S p(x) dx + 0 \int_{S_0} p(x) dx + c \int_{S_+} p(x) dx \right) \\ &= 1; \end{aligned}$$

the last step following since $g_c \in \mathcal{C}(f_c)$ everywhere and $\int_{S \cup S_+} p(x) dx > 0$. As such, $\int_{\mathbb{R}^d} p(x) dx = 0$, and since $\int_{\mathbb{R}^d} p(x) dx = 1$ by properties of \int (cf. Lemma 1.B.1), it follows that $\int_{\mathbb{R}^d} p(x) dx = \int_{\mathbb{R}^d} p(x) dx = 1$.

In the remainder of the proof of this conjugacy property, suppose $p \in \mathcal{C}(0; \infty)$ -a.e..

Now consider the case that p is a simple function with $p \in \mathcal{C}(0; \infty)$ everywhere. Since \int is finite over $[0; \infty]$ (cf. Lemma 1.B.1), p is within the relative interior of the domain of \int everywhere, and thus $\mathcal{C}(p(z))$ is a nonempty set for every $z \in \mathbb{R}^d$ (Rockafellar, 1970, Theorem 23.4).

Consequently, construct q so that $q(z) \in \partial \varphi(p(z))$ everywhere, and moreover q is also a simple function (i.e., pick the same subgradient along each of the finitely many regions composing Ω); these choices will ensure that there are no measurability issues with q (otherwise, the arguments pass through for arbitrary $p \in L^1(\Omega)$); additionally, $q \in L^1(\Omega)$ since μ is a finite measure. Since φ is lower semi-continuous, $q(z) \in \partial \varphi(p(z))$ implies $p(z) \in \partial \varphi^*(q(z)) = \partial \varphi^*(q)$, and the Fenchel-Young inequality implies

$$\varphi^*(p(z)) = p(z)q(z) - \varphi(q(z)) = p(z)q(z) - \varphi^*(q(z)):$$

As such,

$$\begin{aligned} \int_{\Omega} \varphi^*(p) d\mu &= \sup_{f \in L^1(\Omega)} \int_{\Omega} (fp - \varphi(f)) d\mu \\ &= \int_{\Omega} (qp - \varphi(q)) d\mu \\ &= \int_{\Omega} \varphi^*(p) d\mu \end{aligned}$$

Now using the fact that $p(z) \in \partial \varphi(q(z))$,

$$\begin{aligned} \int_{\Omega} \varphi^*(p) d\mu &= \sup_{f \in L^1(\Omega)} \int_{\Omega} (fp - \varphi(f)) d\mu \\ &= \sup_{f \in L^1(\Omega)} \int_{\Omega} (fp - \varphi(q) - p(f - q)) d\mu \\ &= \sup_{f \in L^1(\Omega)} \int_{\Omega} (pq - \varphi(q)) d\mu \\ &= \int_{\Omega} \varphi^*(p) d\mu; \end{aligned}$$

combining these two inequalities, $\int_{\Omega} \varphi^*(p) d\mu = \int_{\Omega} \varphi^*(p) d\mu$.

Now consider the case that $p \in L^1(\Omega)$ is just measurable. Since the simple functions are dense in $L^1(\Omega)$ (Folland, 1999, Theorem 6.8), there exists a simple function $p_i \in L^1(\Omega)$ with $\|p - p_i\|_1 \leq 1/i$, and moreover p_i may be clamped to the range $[-i, i]$ (with i sufficiently large to make this interval nonempty), whereby this clamped simple function p_i satisfies $\|p - p_i\|_1 \leq 2/i$. Since $\int_{\Omega} \varphi^*$ is lower semi-continuous,

$$\int_{\Omega} \varphi^*(p) d\mu = \liminf_i \int_{\Omega} \varphi^*(p_i) d\mu = \liminf_i \int_{\Omega} \varphi^*(p_i) d\mu = \int_{\Omega} \varphi^*(p) d\mu;$$

where the last step used the dominated convergence theorem applied with dominating constant map $z \mapsto \sup_{q \in [0, 1]} |q|$, which is finite since ψ is continuous over the compact set $[0, 1]$ (cf. Lemma 1.B.1).

Next consider the case that measurable $p \in L^1(\Omega; \mathbb{R})$ -a.e.; then $\varphi(z) := p(z)1_{[p(z) \in (0, 2]} + (2-p(z))1_{[p(z) \in (2, \infty)]}$ satisfies $(\varphi^*, d^*)(p) = (\varphi^*, d^*)(\varphi)$ by definition of the conjugate (the integrals ignore measure zero sets), whereby $(\varphi^*, d^*)(p) = (\varphi^*, d^*)(\varphi) = \int \varphi^*(p) d\mu = \int \varphi^*(\varphi) d\mu$.

Lastly, suppose measurable $p \in L^1(\Omega; \mathbb{R})$ -a.e.. For each i , define $p_i = (1 - \frac{1}{i})p + \frac{1}{i}$. Then $p_i \in L^1(\Omega; \mathbb{R})$ -a.e., and $\|p_i - p\|_1 \rightarrow 0$, whereby the lower semi-continuity of (φ^*, d^*) and dominated convergence theorem cover this case in the same way as the move away from simple functions.

Note lastly that these last choices provide a finite integral, since $\sup_{z \in [0, 1]} |z| < 1$ as above, and μ is a finite measure.

Lastly, for the subdifferential rule, let f be given, and first suppose $p \in L^1(\Omega; \mathbb{R})$ -a.e. z . Then, by the Fenchel-Young equality (applied to φ at each $f(z)$) and the conjugacy relation $(\varphi^*)^* = \varphi$, established above,

$$\int \varphi^*(f(z)) d\mu(z) = \int (\varphi^*(p(z)) + p(z)f(z)) d\mu(z) = \int \varphi^*(p) d\mu + \int p f d\mu;$$

by the Fenchel-Young equality (now applied a single time to φ^* (Zalinescu, 2002, Theorem 2.4.2.iii)), this implies that $p \in \partial \varphi^*(f)$. This establishes one direction of the desired i , but also that the subdifferential is nonempty, since a valid choice of p always exists (for instance, take $p(z) \in \partial \varphi^*(f(z))$ when $|f(z)| < 1$, and $p(z) = 0$ when $|f(z)| = 1$, where the latter adjustment is made only over a null set since $f \in L^1(\Omega)$).

The other direction is established via contrapositive; namely, let $p \in L^1(\Omega; \mathbb{R})$ be given, and suppose the set

$$S := \{z \in \Omega : p(z) \notin \partial \varphi^*(f(z))\}$$

satisfies $\mu(S) > 0$. Since the Fenchel-Young inequality is an equality only in the case that a subdifferential relationship is satisfied, it follows that

$$\varphi^*(f(z)) + \varphi^*(p(z)) > p(z)f(z) \quad \forall z \in S;$$

Consequently, since S has positive measure, and again using the relation $(\cdot)_d = \cdot$ as established above,

$$\begin{aligned} \int_{S^c} f + \int_S f &= \int_{S^c} f + \int_S f \\ &> \int_{S^c} f + \int_S f \\ &= \int f \end{aligned}$$

thus, again invoking the fact that equality is only achieved in the Fenchel-Young inequality when a subdifferential relation is satisfied (Zalinescu (2002, Theorem 2.4.2.iii)), it follows that $\int f = \int (\cdot)_d(f)$. \square

Lastly, the proof of the duality relation is as follows.

Proof of Lemma 4.1.1. Consider the following two Fenchel problems:

$$\begin{aligned} V_p &:= \inf_{\lambda} \int (A - \lambda)_d + \int 0_d : \lambda \in \mathbb{R}^2; \\ V_d &:= \sup_{p \in L^1(\cdot)} \int (p)_{f \circ g} \end{aligned}$$

where $f \circ g$ is the indicator for the set $f \circ g$,

$$f \circ g(\cdot) = \begin{cases} 0 & \text{when } \cdot = 0; \\ 1 & \text{otherwise,} \end{cases}$$

and is the conjugate to $\int 0_d$. In order to show $V_p = V_d$, also that attainment occurs in the dual, and lastly the optimality condition, an appropriate Fenchel duality rule will be applied (Zalinescu, 2002, Corollary 2.8.5 using condition (vii)), which requires the verification of the following properties.

First note that $\int (\cdot)_d$ and $\int (\cdot)$ are both convex lower semi-continuous, and moreover mutually conjugate (cf. Lemma 4.1.2). The function $\int 0_d = \int 0_d$ is immediately convex lower semi-continuous (over \mathbb{R}), and thus its conjugate $f \circ g$ is similarly convex lower semi-continuous, and the two are mutually conjugate (Zalinescu, 2002, Theorem 2.3.3).

Both $L^1(\Omega)$ and $\mathcal{M} = L^1(\Omega)$ are Banach and therefore Fréchet spaces. (The present proof is one of the reasons \mathcal{M} was taken to be a Banach space and not merely, say, weightings with finite support as used by the algorithm).

Let $\text{dom}(f) = \{x : f(x) < \infty\}$ denote the effective domain of a convex function, meaning those values where it is finite. As provided by Lemma 4.1.2, $\text{dom}(f) \subseteq \mathbb{R}^d$, and thus, since $A : \mathcal{M} \rightarrow L^1(\Omega)$,

$$A(\text{dom}(f)) \subseteq \text{dom}(f) \subseteq \mathbb{R}^d = A + L^1(\Omega) = L^1(\Omega);$$

which settles the constraint qualification. (Recall that A is not necessarily a closed subspace (cf. Lemma 1.A.6); thus problems would occur here if this proof were attempted for $\mathcal{M} = \text{exp}$, as $\text{dom}(f) \subseteq \mathbb{R}^d$ would not swallow the closure issues of A .)

This completes the conditions necessary for the Fenchel duality result.

To adjust the statement into the desired form, Lemma 4.1.2 provided that \mathbb{R}^d is finite if its input lies within $[0; \infty]$ -a.e. (thus other values may safely be discarded from the optimization problem, which always has feasible point $0 \in L^1(\Omega)$), and secondly $\langle A^\top p, \cdot \rangle < \infty$ if $\langle A^\top p, \cdot \rangle = 0$ (recall the form of $\langle A^\top p, \cdot \rangle$ in Lemma 1.A.5).

Lastly, for the optimality conditions, the invoked Fenchel duality rule provided by Zalinescu (2002, Corollary 2.8.5 using condition (vii)) also grants that λ^* is optimal if there exists $p \in L^1(\Omega)$ which satisfies $p \in \mathcal{M}^*(A)$ and $\langle A^\top p, \cdot \rangle \in \mathcal{M}^*(0)$. To translate this into the desired form, first note that Lemma 4.1.2 provided $p \in \mathcal{M}^*(A)$ if $p \in \mathcal{M}^*(A)$ -a.e.. Secondly, $\mathcal{M}^*(0) = \{f \geq 0\}$ everywhere (indeed, "everywhere" is the same as -a.e., where μ is counting measure over \mathbb{H} and $\mathcal{M} = L^1(\Omega)$), thus $\langle A^\top p, \cdot \rangle = 0$ everywhere, which is equivalent to $\langle A^\top p, \cdot \rangle = 0$ for all $\lambda \in \mathcal{M}$ by the form of A^\top developed in Lemma 1.A.5, which in turn is equivalent to $\langle A^\top p, \cdot \rangle = 0$ again by Lemma 1.A.5. What remains is to necessitate $p \in [0; \infty]$ -a.e.. Since $\langle \cdot, \cdot \rangle$ is Lipschitz with tightest constant $\|A^\top\|$, and $\langle \cdot, \cdot \rangle$ is increasing, it follows that $\langle \cdot, \cdot \rangle \in [0; \infty]$ everywhere, thus the above property $p \in \mathcal{M}^*(A)$ -a.e. grants one direction. On the other hand, if $p \notin [0; \infty]$ for a set of positive measure, then $p \notin \mathcal{M}^*(A)$ for a set of positive measure by Lemma 4.1.2, and thus the optimality condition $p \in \mathcal{M}^*(A)$ can not hold. \square

Appendix 4.B Bibliographic Remarks

While duality statements abound in machine learning, they appear to be rare when considering the optimization over the distribution. The results here were based on general Fenchel rules presented by Zalinescu (2002) and Borwein and Zhu (2005), however production of such general rules is the topic of great research in the convex analysis literature.

Acknowledgements

This chapter is based on work by the dissertation author which appeared in the Conference on Learning Theory, 2013.

Chapter 5

Finite Hypothesis Classes

This chapter provides the first set of statistical guarantees. In a sense, these results are strong, since they hold for a very relaxed class of convex losses | including not just ℓ_{exp} and ℓ_{log} , but even ℓ_{russ} , which is not strictly convex, and ℓ_{hinge} , which is neither strictly convex nor differentiable | and moreover the algorithmic component is completely unspecified: the guarantees hold for solutions which are merely ϵ -optimal, for a provided value of ϵ .

What's the catch? The hypothesis class H must be finite. To put this in context, here are some ways to produce a boosting algorithm with statistical guarantees.

1. Add some sort of penalization or regularization to the (primal) objective function, as investigated by, for instance, Blanchard et al. (2003), Lugosi and Vayatis (2004).
2. Adjust the algorithm to take smaller steps, thus explicitly producing iterates with small norm as discussed in Chapter 1; this choice is investigated by Zhang and Yu (2005).
3. Apply an early stopping criterion, for instance the second option in Figure 1.3 (enforce $\|w\| \leq m^a$); this option was investigated for the exponential loss by Bartlett and Traskin (2007), Schapire and Freund (2012), and is investigated in Chapter 6 for losses similar to the logistic loss.
4. This chapter says that it is enough to ensure H is finite, and that some black box minimization scheme is provided for the corresponding choice of ϵ (and thus L), meaning the results of Chapter 3 suffice.

There is a productive way to work around this catch: start with an infinite class H_0 , and based on the availability of data, cut it down to some finite subset $H \subset H_0$. For many classes, it can be shown that as the sample size increases, the degradation incurred by this finitization goes to zero.

5.1 Overview

Throughout this chapter, \mathcal{Z} will be identified with some \mathbb{R}^n , where $n := |\mathcal{H}| < \infty$. Additionally, the relevant set of losses is as follows.

Definition 5.1.1. \mathcal{L}_c contains all convex losses $\ell : \mathbb{R} \rightarrow \mathbb{R}_+$ which are positive at the origin, and satisfy $\lim_{z \rightarrow 1} \ell'(z) = 0$.

This class contains all losses discussed in this thesis: ℓ_{exp} ; ℓ_{log} ; ℓ_{russ} ; ℓ_{hinge} . Note that the algorithm in Figure 1.3 is only suited to the minimization of continuously differentiable choices, so for ℓ_{hinge} some other choice is needed, for instance subgradient descent, or some faster alternative (Nesterov, 2003).

As discussed previously, the guarantees will be roughly of the form: if a near-optimal solution over \mathcal{B} is given, then the quality over \mathcal{L} will not be too bad.

Asking for a relationship between \mathcal{B} and \mathcal{L} , however, is in general too much to ask without further conditions on the provided solutions. In particular, Section 5.2 presents a negative result, namely that there exist nearly optimal choices over \mathcal{B} which are moreover very stable (have good margin properties), however their behavior over \mathcal{L} is arbitrary bad. This result determines the nature of the subsequent analysis and resulting bounds.

The solution is to once again rely on hard cores: the input space is broken into two pieces, one being the hard core, where there exist optimal choices that make a few mistakes, and the hard core's complement, where it is possible to have zero mistakes, albeit giving up on the existence of a minimizer to the true convex risk. This material appears in Section 5.3, and must make a few concessions as compared with the development in Chapter 2.

The hard core has direct entailments on the structure of the convex risk. Specifically, Section 5.4 establishes first that the true risk has quantifiable curvature over the hard core, and effectively zero error over the rest of the space; this mirrors the basic properties of hard cores for finite samples provided by Theorem 2.4.8. This section then goes on to show how, with high probability, this structure carries over to any sampled instance.

The significance of first proving properties of the true risk, and then carrying them over to the sample, is that quantities dictating the structure of the empirical convex risk are sample independent. Consequently, finite sample guarantees, which appear in Section 5.5, display a number of terms which are properties of the true convex risk, and not simply opaque random

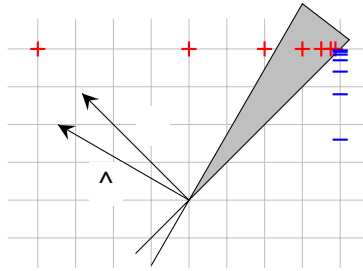


Figure 5.1. Unconstrained minimization can have some bad statistical properties; please see Proposition 5.2.1 and the surrounding discussion.

variables derived from the sample.

5.2 An Impossibility Result

The stated goal of allowing iterates to have unbounded norms is at odds with the task of bounding the convex risk L .

Proposition 5.2.1. There exists a probability measure \mathbb{P} , input space X , and finite set of hypotheses H with the following characteristics.

1. X is the square $[-1; +1]^2$, and H consists of the two projection maps.
2. \mathbb{P} has countable support.
3. There exists a perfect separator, albeit with zero margin.
4. For any $\lambda \in L_c$, $L(A^\lambda) = 0$.
5. Let any finite sample $f(x_i; y_i)_{i=1}^m$, any $b > 0$, and any $\lambda \in L_c$ be given. Then there exists a maximum margin solution $\hat{\lambda}$, i.e., a solution satisfying

$$\arg \min_{\lambda \in L_c} \frac{\sum_{i=1}^m y_i (H^\lambda)(x_i)}{\|\lambda\|_1} = \sup_{\lambda \in L_c} \left(\arg \min_{\lambda \in L_c} \sum_{i=1}^m y_i (H^\lambda)(x_i) : \|\lambda\|_1 = 1 \right)$$

which has $\mathbb{P}(H^{\hat{\lambda}}) = 0$ and $\mathbb{E}(A^{\hat{\lambda}}) = 1/b$ but $L(A^{\hat{\lambda}}) = 0$.

A full proof is provided in Section 5.C, but the mechanism is simple enough to appear as a picture. Consider the linear classification problem in Figure 5.1, which has positive ("+") and negative ("-") examples along two lines. Optimal solutions to R are of the form $c \lambda$, where $\lambda = (-1; +1)$ and $c > 0$ (note $\lim_{c \rightarrow \infty} L(c\lambda) = L(A^\lambda) = 0$). Unfortunately, the positive and

negative examples are staggered; as a result, for any sample, every max margin predictor, which is determined solely by the rightmost "+" and uppermost "-", will fail to agree with the optimal predictor on some small region. A positive probability mass of points fall within this region, and so, by considering scalings c^k as c^{-1} , the convex risk L may be made arbitrary large.

The statement of Proposition 5.2.1 is encumbered with details in order to convey the message that not only do such examples exist, they are fairly benign; indeed, the example depends on the additional regularity of large margin solutions. The only difficulty is the lack of any norm constraint on permissible iterates.

On the other hand, notice that the classification risk R is not only small, but its empirical counterpart \hat{R} provides a reasonable estimate as n increases. Furthermore, if the distribution were adjusted slightly so that every 2^{-n} made some mistake, then these unbounded iterates would fail to exist: the huge penalty for predictions very far from correct would constrain the norms of all good predictors.

The preceding paragraph describes the exact strategy of the remainder of this chapter: linear classification problems are split into two pieces, one where optimization may produce unboundedly large iterates with small classification risk, and another piece where iterates are bounded thanks to the presence of difficult examples.

5.3 Hard Cores

The choice here is to split the problem along a hard core. Recall, from Section 1.7, that the hard core was the support $\{p > 0\}$ of a reweighting $p \in L^1(\cdot)$ which decorrelated every predictor, meaning $\sum_{i \in R} p_i = 0$. The choice of $L^1(\cdot)$ matches the strong duality result for Lipschitz losses in Lemma 4.1.1; while this may seem to conflict with the generality of L_c , which contains non-Lipschitz losses, strong duality will only be invoked in proofs of hard core properties, and need only apply for a single loss. Thereafter, when applying these properties, the finiteness of H will be essential. Recall additionally from Section 1.7 that hard cores are guaranteed to exist (cf. Theorem 1.7.4).

The following result shows that hard cores achieve the goal laid out at the closing of Section 5.2; this statement can be seen as a distillation of the components of Theorem 2.4.8 which are needed for the statistical proofs here. The proof of this statement, which is somewhat involved, appears in Section 5.D.

Theorem 5.3.1. Let finite hypothesis class H , probability measure μ , and hard core C be given. The following statements hold.

1. There exists a sequence ϵ_i, g_i^1 with $y(H_{i-1})(x) = 0$ for μ -a.e. $(x; y) \in C$, and $y^0(H_{i-1})(x^0) \geq 1$ for μ -a.e. $(x^0; y^0) \in C^c$.
2. Every $\mu \in \mathcal{R}^n$ satisfies either $\mu(C \setminus [y(H)(x) = 0]) = \mu(C)$ or $\mu(C \setminus [y(H)(x) < 0]) > 0$.

The first property provides the existence of a sequence which is not only very good μ -a.e. over C^c , but furthermore does not impact the value of H over C ; that is to say, this sequence can grow unboundedly, and have unboundedly positive margins over C^c , while optimization over C can effectively proceed independently. On the other hand, C is difficult: every predictor is either abstaining μ -a.e., or makes errors on a set of positive measure.

Finally, corresponding to the hard core, it will be convenient in future sections to use adopt the following notation for risk specialized to some region.

Definition 5.3.2. Given a set C (typically C or C^c), loss ℓ , function class $F \subseteq L^1(\mu)$, and any $f \in F$, define

$$L(f; C) := \int_C \ell(f(x; y)) d\mu(x; y) \quad \text{and} \quad L(F; C) := \inf_{f \in F} L(f; C);$$

with analogous definitions for \mathbb{D} , and additionally for \mathcal{R} and \mathcal{R}^n over functions in $L^1(\mu^{\otimes n})$.

5.4 Hard Cores and Convex Risk

The hard core imposes the following structure on \mathcal{L} . As provided by Theorem 5.3.1, there is a sequence which does arbitrarily well over C^c , without impacting predictions over C . On the other hand, since mistakes must occur over C , convex losses within \mathcal{L}_c will be forced to avoid large predictors.

Theorem 5.4.1. Let finite hypothesis class H , probability measure μ , hard core C , and loss $\ell \in \mathcal{L}_c$ be given.

1. There exists a sequence ϵ_i, g_i^1 with $y(H_{i-1})(x) = 0$ for μ -a.e. $(x; y) \in C$, and moreover $\lim_{i \rightarrow \infty} \mu(\{(A_i)_{x^0; y^0} = 0\}) = 0$ for μ -a.e. $(x^0; y^0) \in C^c$.

- Let any $\epsilon > 0$ be given. Then there exists $\delta \in \mathbb{R}$ and a set N with $(N) = 0$ so that for every $\mathcal{H} \in \mathbb{R}^n$ with $L(\mathcal{H}; C) \leq L(\mathcal{H}^0; C) + \delta$, there exists a representation $\mathcal{H}^0 \in \mathbb{R}^n$ with $\mathcal{H} = \mathcal{H}^0$ over $C \cap N$, and $k \leq k_1 + \epsilon$.

The structural properties of the true convex risk transfer over, with high probability, to any sampled problem. Crucially, the various bounds are quantified outside the probability; that is to say, they do not depend on the sample.

Theorem 5.4.2. Let finite hypothesis class H , probability measure \mathbb{P} , hard core C , and loss $\ell \in L_c$ be given.

- With probability 1 over the draw of a finite sample, there exists $\mathcal{H} \in \mathbb{R}^n$ so that every $(x_i; y_i) \in C^c$ satisfies $y_i(\mathcal{H})(x_i) > 0$, and every $(x_i^0; y_i^0) \in C$ satisfies $y_i^0(\mathcal{H})(x_i^0) = 0$.
- Given any empirical suboptimality $\epsilon > 0$, there exist $c > 0$ and $b > 0$ so that for any $\delta > 0$, with probability at least $1 - \delta$ over a draw of m points where m_C , the number of points landing in C , has bound

$$m_C \leq \frac{c^2}{\delta} (\ln(n) + \ln(1/\delta));$$

then every ϵ -suboptimal $\mathcal{H} \in \mathbb{R}^n$ over the sample restricted to C , meaning

$$L(\mathcal{H}; C) \leq L(\mathcal{H}^0; C) + \epsilon;$$

has a representation \mathcal{H}^0 with $k \leq k_1 + b$ which has $\mathcal{H} = \mathcal{H}^0$ over the sample restricted to C , and in general \mathbb{P} -a.e. over C .

5.5 Deviation Inequalities

With the structure of the convex risk in place, the stage is set to establish deviation inequalities. These will be stated in terms of both a convex risk L , but also the classification risk R . In order to make this correspondence, this chapter relies on the following standard techniques due to Zhang (2004) and Bartlett, Jordan, and McAuliffe (2006).

Definition 5.5.1. Let F denote the set of measurable functions over X , and let F_Y denote the family of maps $f(x; y) \mapsto yf(x) : f \in F$. (Note that H is to A as F is to F_Y : this definitions exist for convenient use of L .)

Proposition 5.5.2. (See also Bartlett et al. (2006).) Let any $\ell \in L_c$ be given with ℓ differentiable at 0. There exists an associated function $\rho : [0, 1] \rightarrow [0, 1)$ with the following properties. First, for any probability measure μ and any $f : X \rightarrow \mathbb{R}$, $(R(f) - R(F)) \leq L(\int (x; y) \mu - \int y f(x)) \leq L(F_Y)$. Second, the inverse ρ^{-1} exists over $[0, 1)$, and satisfies $\rho^{-1}(r) \neq 0$ as $r \neq 0$.

Definition 5.5.3. Given $\ell \in L_c$ with ℓ differentiable at 0, let ρ , called the ρ -transform, be as in Proposition 5.5.2.

The general use of ρ is through its inverse, which provides, for any ℓ satisfying the conditions to Proposition 5.5.2,

$$\begin{aligned} R(H) - R(F) &= \rho^{-1}(L(A) - L(F_Y)) \\ &= \rho^{-1}(L(A) - L(A) + L(A) - L(F_Y)) : \end{aligned}$$

Note first that ρ^{-1} may be unwieldy, it is frequently easy to provide a useful upper bound. For instance, the exponential loss has $\rho^{-1}(r) \leq 2^p \bar{r}$, the logistic loss has $\rho^{-1}(r) \leq 4^p \bar{r}$, and the hinge loss has $\rho^{-1}(r) = r$ (Zhang, 2004, Bartlett et al., 2006). Secondly, the value of the second decomposition above is that by choosing \bar{r} large enough, $L(A) - L(F_Y)$ may be forced arbitrarily small.

Theorem 5.5.4. Let finite hypothesis class H , probability measure μ , hard core C , and loss $\ell \in L_c$ be given. Let a suboptimality tolerance $\epsilon > 0$ be given; results will depend on reals $\alpha > 0$ and $b > 0$ determined by the preceding terms. The following statements simultaneously hold with any probability $1 - \epsilon$ over the draw of m samples (with $\epsilon := \epsilon/8$ for convenience), and any weighting $w \in \mathbb{R}^n$ which is ϵ -suboptimal (with ϵ) for the corresponding surrogate empirical risk problem, meaning $L(A) \leq L(A) + \epsilon$.

1. Let m_C and m_+ respectively denote the number of samples falling into C and C^c . Then

$$\begin{aligned} m_C &\leq m \cdot (C) \cdot \frac{p}{\ln(1 - \epsilon/2m)} ; \\ m_+ &\leq m \cdot (C^c) \cdot \frac{p}{\ln(1 - \epsilon/2m)} ; \end{aligned}$$

2. The true classification risk over the unbounded portion, C^c , has bound

$$\frac{R(H; C^c)}{(C^c)} < \frac{1}{(0)} + 2 \frac{\frac{1}{2} (n \ln(2m_+ + 1) + \ln(4 = \eta))}{(0)m_+} + \frac{4(n \ln(2m_+ + 1) + \ln(4 = \eta))}{m_+}. \quad (5.5.5)$$

If moreover $\hat{c}(0) = m$, then

$$\frac{R(H; C^c)}{(C^c)} < \frac{4(n \ln(2m_+ + 1) + \ln(4 = \eta))}{m_+}. \quad (5.5.6)$$

3. Suppose

$$m_C = c^2(\ln(n) + \ln(6 = \eta)):$$

The true surrogate risk over the bounded portion, C , has bound

$$L(A; C) < L(A; C) + \frac{c \frac{1}{\ln(n)} + 4 \frac{1}{\ln(2 = \eta)}}{p \frac{1}{m_C}}; \quad (5.5.7)$$

Additionally, if \hat{c} is differentiable at 0, the classification risk has bound

$$R(H; C) < R(F_Y; C) + \frac{c \frac{1}{\ln(n)} + 4 \frac{1}{\ln(2 = \eta)}}{p \frac{1}{m_C}} + L(A; C) < L(F_Y; C); \quad (5.5.8)$$

4. Suppose, for simplicity, that

$$m = \max \{2 \ln(1 = \eta)\} = \min \{ (C)^2; (C^c)^2; 2c^2(\ln(n) + \ln(1 = \eta)) \} = (C)$$

(where bounds are interpreted to hold trivially when denominators contain 0) and additionally that $\hat{c}(0) = m$ and \hat{c} is differentiable at 0. Then the true classification risk of the full

problem has bound

$$R(H) - R(F) \leq 1 + \frac{c^{\frac{p}{2}} \overline{\ln(n)} + 4^{\frac{p}{2}} \overline{\ln(2\epsilon)}}{m^{\frac{p}{2}} (C^{\circ})} + L(\text{span}(H); C) - L(F; C) + \frac{8(n \ln(m^{\frac{p}{2}} (C^{\circ}) + 1) + \ln(4\epsilon))}{m^{\frac{p}{2}} (C^{\circ})}.$$

As a consequence of these bounds, it is reasonable to choose n so that $\ln(|H|) = o(m)$; for instance, $|H| = \exp(\frac{p}{m})$ gives roughly a rate of $m^{-1/4}$.

Appendix 5.A Technical Preliminaries

Lemma 5.A.1. Let any $\ell \in L_c$ be given. Then ℓ is continuous, measurable, and nondecreasing. Subgradients exist everywhere, and satisfy $\ell(0) \in \mathbb{R}_{++}$.

Proof. Since ℓ is finite everywhere, it is continuous (Rockafellar, 1970, Corollary 10.1.1), and thus measurable (Folland, 1999, Corollary 2.2). Since convex functions are subdifferentiable everywhere along the relative interior of their domains (which in this case is just \mathbb{R}), it follows that ℓ has subgradients everywhere (Rockafellar, 1970, Theorem 23.4).

If ℓ were not nondecreasing, there would exist $x < y$ with $\ell(x) > \ell(y)$; but that means every subgradient $g \in \partial \ell(x)$ satisfies

$$\ell(y) \leq \ell(x) + g(y - x);$$

and thus $g < 0$. But then, for any $z < x$, $\ell(z) \leq \ell(x) + g(z - x)$, which in particular contradicts $\lim_{z \uparrow 1} \ell(z) = 0$ (indeed, it implies $\lim_{z \uparrow 1} \ell(z) = 1$), thus ℓ is nondecreasing.

Next, since ℓ is nondecreasing, $\ell \in \mathbb{R}_+$. However, since $\ell(0) > 0$, it follows that $\ell(0) \in \mathbb{R}_{++}$, since otherwise $\lim_{z \uparrow 1} \ell(z) = 0$ would be contradicted. \square

The following grants continuity of the map $\ell: \mathbb{R}^d \rightarrow \mathbb{R}$ for all of L_c , and in this way is stronger than the analogous result for general choices of ℓ (cf. Lemma 4.1.2). Unsurprisingly,

the proof crucially hinges upon \mathbb{R}^n being finite dimensional (in particular, the compactness of l_1 balls is used, which is of course in general false in the l_1 topology in general).

Lemma 5.A.2. Let finite hypothesis class H , finite measure μ , and loss $\ell \in L_c$ be given. Then the map $\mu \mapsto \int \ell(y, H(x)) d\mu(x)$ is convex and continuous.

Proof. Convexity follows from convexity of ℓ ; namely, let μ, ν and $\lambda \in [0, 1]$ be given; then

$$\begin{aligned} \int \ell(y, H(\lambda\mu + (1-\lambda)\nu)) d(x; y) &= \int \ell(y, H(x)) \lambda d\mu(x) + (1-\lambda) \int \ell(y, H(x)) d\nu(x) \\ &= \lambda \int \ell(y, H(x)) d\mu(x) + (1-\lambda) \int \ell(y, H(x)) d\nu(x) \\ &= \lambda \int \ell(y, H(x)) d\mu(x) + (1-\lambda) \int \ell(y, H(x)) d\nu(x) \end{aligned}$$

For continuity, it suffices to show that the map commutes with limits. In particular, let a convergent sequence μ_i be given, and without loss of generality suppose there exists finite $B > 0$ so that $\|\mu_i - \mu\|_1 \leq B$ (whereby $\|\mu\|_1 \leq B$ since $\mu_i \geq 0$). Since $\sup_{h, x, y} |j_h(x)| \leq 1$ and ℓ is nondecreasing, it follows for every i and $(x; y) \in X \times Y$ that

$$\ell(y, H(x)) \leq \|\mu_i - \mu\|_1 + \ell(y, H(x)) \leq B + \ell(y, H(x))$$

Furthermore, note that the map $f(x; y) = \ell(y, H(x))$ satisfies $f \in L^1(\mu)$ since μ is a finite measure, i.e., $\int f d\mu \leq \|\mu\|_1 \leq B < \infty$. Consequently, by the dominated convergence theorem with dominating map f and the continuity of ℓ (cf. Lemma 5.A.1),

$$\lim_{i \rightarrow \infty} \int \ell(y, H(x)) d\mu_i(x) = \int \ell(y, H(x)) d\mu(x) = \int \ell(y, H(x)) d\mu(x)$$

as desired. \square

Proposition 5.A.3. Let finite hypothesis class H , probability measure μ , and loss $\ell \in L_c$ be given. Then given a bound b on the l_1 norm of considered predictors, there exists $\epsilon(b)$ so that, for any $\delta > 0$, with probability at least $1 - \delta$ over the draw of m points from μ , every $\hat{\mu} \in \mathcal{R}^n$ with $\|\hat{\mu} - \mu\|_1 \leq b$ satisfies

$$L(\hat{\mu}) - L(\mu) \leq \frac{c \sqrt{p \ln(n)} + \sqrt{p \ln(2/\delta)}}{\sqrt{m}}$$

Proof. Let bound b and loss $\ell \in L_c$ be given. Define a truncation

$$\hat{\ell}(z) := \begin{cases} \ell(z) & \text{when } z \leq b, \\ \ell(b) & \text{otherwise.} \end{cases}$$

Since ℓ is nondecreasing (cf. Lemma 5.A.1), $\hat{\ell}(z) \leq \ell(b)$, and furthermore $\hat{\ell}$ is Lipschitz with a constant that may be measured at b ; indeed, since ℓ is finite everywhere, it has bounded subdifferential sets (Rockafellar, 1970, Theorem 23.4), and thus, taking any $z_1, z_2 \in \mathbb{R}$ and supposing without loss of generality that $z_1 \leq z_2$,

$$\begin{aligned} | \hat{\ell}(z_2) - \hat{\ell}(z_1) | &= | \ell(z_2) - \ell(z_1) | \\ &= \sup \{ \ell(z_2) - (\ell(z_2) + h(z_2 - z_1)) : g \in \partial \ell(z_2) \} \\ &= |z_2 - z_1| \sup \{ |g| : g \in \partial \ell(z_2) \} \\ &< 1; \end{aligned}$$

correspondingly, set a Lipschitz constant $L := \sup \{ |g| : g \in \partial \ell(b) \}$.

Note that for every $x \in \mathbb{R}^n$ with $|x| \leq b$, $\sup_{x \in \mathbb{R}^n} | \ell(x) - \ell(b) | \leq bL$, and thus $L(A) = \mathbb{E} \ell(A)$ and $\mathbb{E} \hat{\ell}(A) = \mathbb{E} \ell(A) - bL$. Lastly, the desired constant c , which does not depend on n , or m , will be $c := \max \{ 2L, b \} \sqrt{\frac{2 \ln(2)}{m}}$.

Now let a sample of size m be given, and let $R_m(\text{span}(H); b)$ denote the Rademacher complexity of $\{ H : |x| \leq b \}$. By properties of Rademacher complexity and a few appeals to McDiarmid's inequality, (Boucheron et al., 2005a, Theorem 3.1, and the proof of Theorem 4.1), with probability at least $1 - \delta$ over the draw of this sample,

$$\sup_{|x| \leq b} L(A) - \mathbb{E} \hat{\ell}(A) = \sup_{|x| \leq b} \ell(A) - \mathbb{E} \ell(A) \leq 2L R_m(\text{span}(H); b) + \sqrt{\frac{2 \ln(2/\delta)}{m}}. \quad (5.A.4)$$

Next, taking $R_m(S)$ denote the Rademacher complexity of a set of functions S , by $R_m(\text{span}(H); b) = b R_m(\text{span}(H); 1) = b R_m(H)$ and an appeal to Massart's Finite Lemma (Boucheron et al., 2005a, Theorem 3.3)

$$R_m(\text{span}(H); b) \leq b \sqrt{\frac{2 \ln(n)}{m}}.$$

Plugging this into eq. (5.A.4) and recalling the choice $c = \max_{z \in S} f(z)$, the result follows. \square

Lemma 5.A.5. Let $S \subseteq \mathbb{R}$ and convex $f : S \rightarrow \mathbb{R}$ be given. If $x, y \in S$ are given with $x < y$ and $f(x) < f(y)$, then for every $S \ni z > y$, $f(y) < f(z)$.

Proof. Write y as a combination of x and z :

$$y = x \frac{z - y}{z - x} + z \frac{y - x}{z - x} :$$

By convexity and $f(y) > f(x)$,

$$\begin{aligned} f(y) &= f\left(x \frac{z - y}{z - x} + z \frac{y - x}{z - x}\right) > f(x) \frac{z - y}{z - x} + f(z) \frac{y - x}{z - x} \\ &= f(x) \frac{z - y}{z - x} + z \frac{y - x}{z - x} \\ &= f(y) : \end{aligned}$$

Rearranging and using $x < y$, it follows that $f(y) < f(z)$. \square

Appendix 5.B Structure of L over $S_D(H; \cdot)$

The following theorem leads to a number of properties presented in Sections 5.3 and 5.4; it is easiest to prove them at once, as a ring of implications. Note the resemblance to Theorem 2.4.4 from Chapter 2.

Theorem 5.B.1. Let finite hypothesis class H , probability measure μ , and a set D be given. The following statements are equivalent.

1. For every $\lambda \in \mathbb{R}^n$, either $\mu(D \setminus [y(H) \cdot \lambda = 0]) = \mu(D)$ or $\mu(D \setminus [y(H) \cdot \lambda < 0]) > 0$.
2. Given any λ and any $\eta \in L_c$, there exists a bound b and a null set $N \subseteq X$ with $\mu(N) = 0$ (i.e., $\mu(N) = 0$) so that for every λ -suboptimal weighting \hat{h} over D , meaning any weighting satisfying

$$L(\hat{h}; D) \leq L(h; D) + b ;$$

there exists η^0 with $\|\eta^0\| \leq b$ and $\hat{h} = \eta^0$ over $D \setminus N$.

3. $D \subset S_D(H; \mu)$.

The following structural lemma is crucial. The finiteness of H is used in two places. Less essentially, it ensures the existence of orthogonal projections (since subspaces of \mathbb{R}^n are always closed), but this could be circumvented with points that are almost closest. The more crucial reliance is in exhibiting a positive constant $c > 0$ on the difference of norms between H and H^\perp (over K^\perp); this in turn depends on the compactness of the L^1 ball, which is false for general vector spaces (with compactness measured in the L^1 topology).

Lemma 5.B.2. Let finite hypothesis class H , probability measure μ , and a set D be given. Define the set

$$K := \{f \in \mathbb{R}^n : y(H) \cdot x = 0 \text{ for } \mu\text{-a.e. } (x; y) \in D\}$$

The following statements hold.

1. K is a closed subspace, as is its orthogonal complement K^\perp , and orthogonal projection onto K^\perp is well-defined.
2. There exists a set N with $\mu(N) = 0$ so that, for any $f \in \mathbb{R}^n$, the orthogonal projection $P_{K^\perp}(f) \in K^\perp$ satisfies $H \cdot P_{K^\perp}(f) = H \cdot f$ everywhere over $D \cap N^c$.
3. There exists a constant $c > 0$ so that, for any $f \in \mathbb{R}^n$ with $\mu(D \setminus [A \leq 0]) > 0$, $\|A\|_{L^1(D)} = \|f\|_{L^1(D)} > c$, where $L^1(D)$ is the L^1 metric with respect to the measure defined by $\mu_D(S) = \mu(D \cap S)$ for any measurable set S .

Proof. (Item 1) Direct from its construction, K is a subspace. Since \mathbb{R}^n is finite dimensional, it follows that both K and K^\perp are closed subspaces, whereby the orthogonal projection P_{K^\perp} exists.

(Item 2) Given the subspace pair K and K^\perp , for any $f \in \mathbb{R}^n$, there exists the decomposition $f = P_K(f) + P_{K^\perp}(f)$, where $P_{K^\perp}(f) \in K^\perp$. By definition, $H \cdot P_K(f) = 0$ μ -a.e. over D , and thus $H \cdot f = H \cdot P_{K^\perp}(f)$ μ -a.e. over D .

Now let Q be any countable dense subset of \mathbb{R}^n . For each $q_i \in Q$, define $N_i := [H \cdot q_i \leq 0]$, where the above provides $\mu(N_i) = 0$. Set $N := \bigcup_i N_i$, which is measurable since it is a countable union, and moreover $\mu(N) = 0$ by μ -additivity. It will now be argued that the projections onto K^\perp give equivalences over $D \cap N^c$.

To this end, let any $\alpha \in \mathbb{R}^n$, any $(x; y) \in D \times \mathbb{N}$, and any $\epsilon > 0$ be given. Since Q is a countable dense subset of \mathbb{R}^n , there exists $i \in Q$ with $\|k_i - k_2\| = \|k_i - k_1\| = 2$. Now let P^ϵ denote the orthogonal projection operator onto K^ϵ ; by the triangle inequality, Cauchy-Schwarz, and since orthogonal projection is a contraction,

$$\begin{aligned} 0 &\leq \|H(x) - (H^\epsilon)^\epsilon(x)\| = \|H(x) - (HP^\epsilon)^\epsilon(x)\| \\ &= \|H(\alpha + i)(x) - (HP^\epsilon(\alpha + i))^\epsilon(x)\| \\ &\leq \|H(\alpha)(x) - (H^\epsilon)^\epsilon(x)\| + \|H(i)(x)\| + \|HP^\epsilon(\alpha)(x)\| \\ &\leq \epsilon + \sup_{h;x} \|h(x)\| \|k_i\| + \sup_{h;x} \|h(x)\| \|k_1\| P^\epsilon k_2 \|k_i\| \\ &= 0 + \epsilon + \epsilon = 2\epsilon \end{aligned}$$

Taking $\epsilon \rightarrow 0$, it follows that $H = H^\epsilon$ over $D \times \mathbb{N}$.

(Item 3) For the final part, if every $\alpha \in \mathbb{R}^n$ has $(D \setminus [A \leq 0]) = \emptyset$, there is nothing to show, so suppose there exists $\alpha \in \mathbb{R}^n$ with $(D \setminus [A \leq 0]) > \emptyset$. Consider the optimization problem

$$\inf_{\alpha \in \mathbb{R}^n; (D \setminus [A \leq 0]) > \emptyset} \frac{\|A - k_{L^1(D)}\|}{\|k_1\|} = \inf_{\alpha \in K^\epsilon; \|k_1\| = 1} \|A - k_{L^1(D)}\|$$

The latter is a minimization of a continuous function over a nonempty compact set, and thus attains a minimizer. But $\alpha \in K^\epsilon$ and $\|k_1\| = 1$, thus $\|A - k_{L^1(D)}\| > 0$. The result follows with $c := \|A - k_{L^1(D)}\| > 0$. □

Proof of Theorem 5.B.1. (Item 1 \Rightarrow Item 2.) The proof will proceed by contradiction. Specifically, suppose Item 1, but contradictorily that Item 2 fails to hold; in particular, there exists $\alpha \in \mathbb{R}^n$ and $\epsilon \in L_c$ so that for any $b > 0$ and any null set $N \subseteq \mathbb{R}^n$, there exists $\hat{\alpha} \in \mathbb{R}^n$ with $L(\hat{\alpha}) \subseteq L(\alpha; D) + \epsilon$, but there is no representation $\alpha \in \mathbb{R}^n$ satisfying $\|k_1\| = b$ and $H^\epsilon = H^\alpha$ over $D \times \mathbb{N}$. As such, let N be the null set, as provided by Lemma 5.B.2, so that every $\alpha \in \mathbb{R}^n$ has $H = H^\epsilon$ everywhere on $D \times \mathbb{N}$. Additionally, let ϵ and δ be as given by the contradiction; therefore, for every $i \in \mathbb{Z}_{++}$, there exists α_i with $L(\alpha_i; D) \subseteq L(\alpha; D) + \delta$, but there is no representation α_i with $H_{\alpha_i} = H_{\alpha_i}^\epsilon$ over $D \times \mathbb{N}$ and $\|k_1\| = \alpha_i$.

For every i , by Lemma 5.B.2 and the choice of N , the projection P_i^ϵ of α_i onto K^ϵ satisfies $H_{\alpha_i} = H_{P_i^\epsilon}^\epsilon$ over $D \times \mathbb{N}$ (whereby $L(\alpha_i; D) \subseteq L(\alpha; D) + \delta$ still holds), and thus the

above assumed contradiction provides $k_i \geq k_{i-1}$. Note that $f_i = k_i g_{i=1}^1$ lies in a compact set (the unit L^1 ball), and thus let $f_i^{(2)} g_{i=1}^1$ be a subsequence of $f_i = g_{i=1}^1$ with $k_i^{(2)} = k_i \geq k_{i-1} \geq 2 \in \mathbb{R}^n$. Since $K^?$ is a closed subspace, it follows that $2 \in K^?$ as well, and furthermore $k_i = 1$. Since $k_i \geq k_{i-1}$ and the subsequence $f_i^{(2)} g_{i=1}^1$ can be formed simply via deletions, it follows that $k_i^{(2)} = k_i \geq 1$, and in particular $\lim_{i \rightarrow \infty} k_i^{(2)} = 1$, and lastly $L(A_i^{(2)}; D) \leq L(A; D) + \epsilon$ still holds. Furthermore, since Lemma 5.B.2 also provides $\epsilon > 0$ with $k_A \leq k_{L^1(D)} = k_1 \leq c$ for $2 \in K^?$, it follows that $k_A \leq k_{L^1(D)} \leq c$, which is only possible if $(D \setminus [y(H) \leq 0]) > 0$.

By assumption (i.e., by Item 1), since $(D \setminus [y(H) \leq 0]) > 0$, then $(D \setminus [y(H) < 0]) > 0$; for convenience, define the set $P := [y(H) < 0]$. Thus, for any $2 \in \mathbb{R}^n$, taking any $g \in \mathbb{R}^n$ (note $g > 0$ via Lemma 5.A.1),

$$\begin{aligned}
 & \lim_{t \downarrow 0} \frac{\int_D (y(H(\cdot + t)))(x) d(x; y) - \int_D (y(H)(x)) d(x; y)}{t} \\
 & \lim_{t \downarrow 0} \frac{\int_{D \setminus P} (y(H(\cdot + t)))(x) d(x; y) - \int_D (y(H)(x)) d(x; y)}{t} \\
 & \lim_{t \downarrow 0} \frac{\int_{D \setminus P} (0 + g(y(H(\cdot + t)))(x)) d(x; y) - \int_D (y(H)(x)) d(x; y)}{t} \\
 & = g \int_{D \setminus P} y(H)(x) d(x; y) \\
 & \quad + \lim_{t \downarrow 0} \frac{\int_{D \setminus P} (0 + g(y(H)(x))) d(x; y) - \int_D (y(H)(x)) d(x; y)}{t} \\
 & > 0:
 \end{aligned} \tag{5.B.3}$$

The above derivation shows that $\int_D (y(H(\cdot + t))) d$ grows arbitrarily large with t , and in particular the ray $f + t : t \geq 0$ must exit (and never return to) the desired ϵ -sublevel set

$$C := \{f + t : t \geq 0\} : L(A; D) \leq L(A; D) + \epsilon$$

To develop the contradiction, it will be shown that the construction of C indicates that this ray should in fact not be leaving C ; the proof will be similar to one due to Hiriart-Urruty and Lemaréchal (2001, Proposition A.2.2.3).

Since $\int_D \cdot d$ is convex and continuous (via Lemma 5.A.2, using finite measure μ_D defined as $\mu_D(S) = \mu(D \setminus S)$ for all measurable S), the sublevel sets of $\int_D \cdot d$ are closed convex sets, and in particular C is closed and convex. By construction of C and continuity of H , recalling $\lim_{i \rightarrow \infty} k_i^{(2)} = 1$ and $f_i^{(2)} \in C$, and using the closed convexity of C , every $0 \in C$ and

$t > 0$ satisfies

$$H^{0+} + tH = \lim_{i \rightarrow \infty} \left(1 - \frac{t}{k_i^{(2)} k_1} \right) H^{0+} + \frac{t}{k_i^{(2)} k_1} H_i^{(2)} \in C :$$

This holds for all $t > 0$, which contradicts the fact, derived in eq. (5.B.3), that the ray $f^{0+} + t : t \rightarrow 0^+$ must exit and subsequently never return to C .

(Item 2 \Rightarrow) Item 3.) Instantiate Item 2 with $\lambda = \lambda_{\log}$ and $\beta = 1$, and let the corresponding b and N be given. Pick any minimizing sequence $\{g_{i=1}^{(1)}\}$ for $L(\cdot; D)$, meaning $L(A_i^{(1)}; D) \rightarrow L(A; D)$. Next produce a subsequence $\{g_{i=1}^{(2)}\}$ of $\{g_{i=1}^{(1)}\}$ by removing all elements with $L(A_i^{(1)}; D) > L(A; D) + 1$ (this procedure must be possible, since otherwise $\{g_{i=1}^{(1)}\}$ is not a minimizing sequence). By the assumed properties of Item 2, a sequence $\{g_{i=1}^{(3)}\}$ may be produced from $\{g_{i=1}^{(2)}\}$ where $k_i^{(3)} k_1 \leq b$, and $H_i^{(2)} = H_i^{(3)}$ over $D \cap N$, which in particular means $\{g_{i=1}^{(3)}\}$ is also a minimizing sequence for $L(\cdot; D)$. But this is now a minimizing sequence lying within the l_1 ball in \mathbb{R}^n of radius b , a compact set, thus some subsequence $\{g_{i=1}^{(4)}\}$, has a limit $g \in \mathbb{R}^n$. Since $\int_{\mathbb{R}^n} (y(H) \cdot x) d_{\mu}$ is lower semi-continuous by Lemma 4.1.2 and finite everywhere on \mathbb{R}^n (and thus continuous), it follows that g attains the desired minimal value, meaning $L(A; D) = L(A; D)$.

Applying the duality relation in Chapter 4 to $L(\cdot; D)$ (i.e., applying Lemma 4.1.1 to $\mathbb{R}^n \setminus d_D$), the existence of a primal minimum grants the existence of a dual maximum p satisfying $p \in D(H; D)$, and moreover

$$p(x; y) \in D(\int (y(H) \cdot x)) = \text{fr}_{\log}(\int (y(H) \cdot x)) g = \frac{1}{1 + \exp(y(H) \cdot x)} \quad \text{for } \mu\text{-a.e. } (x; y):$$

As such, the choice $p^0(x; y) := 1/(1 + \exp(y(H) \cdot x))$ satisfies $p^0 \in D$ -a.e., and thus $p^0 \in D(H; D)$; moreover $p^0 > 0$ everywhere, since $1/(1 + \exp(z)) > 0$ for all $z \in \mathbb{R}$.

This reweighting p^0 was with respect to μ_D , so to finish, define $p(x; y) := p^0(x; y) 1_{[(x; y) \in D]}$. By construction, $[p > 0] = D$. Finally, given any $z \in \mathbb{R}^n$,

$$\begin{aligned} \int_{\mathbb{R}^n} y(H) \cdot x p(x; y) d_{\mu}(x; y) &= \int_{\mathbb{R}^n} y(H) \cdot x p^0(x; y) 1_{[(x; y) \in D]} d_{\mu}(x; y) \\ &= \int_{\mathbb{R}^n} y(H) \cdot x p^0(x; y) d_{\mu}(x; y) \\ &= 0: \end{aligned}$$

It follows that $p \in D(H; \cdot)$, and that $D \in S_D(H; \cdot)$.

(Item 3 \Rightarrow Item 1.) Let $p \in D(H; \cdot)$ with $D = \{p > 0\}$ be given, and take any $\mu \in \mathbb{R}^n$ satisfying $\mu(D \setminus \{y(H(\cdot))x > 0\}) > 0$. But notice then, since p decorrelates H ,

$$\begin{aligned} 0 &= \int_{\mathbb{Z}} p(x; y) y(H(\cdot))(x) d(x; y) \\ &= \int_{D; y(H(\cdot))(x) > 0} p(x; y) y(H(\cdot))(x) d(x; y) + \int_{D; y(H(\cdot))(x) < 0} p(x; y) y(H(\cdot))(x) d(x; y); \end{aligned}$$

From this it follows that

$$\int_{D; y(H(\cdot))(x) < 0} p(x; y) y(H(\cdot))(x) d(x; y) = \int_{D; y(H(\cdot))(x) > 0} p(x; y) y(H(\cdot))(x) d(x; y) > 0;$$

where the inequality follows from $\mu(D \setminus \{y(H(\cdot))x > 0\}) > 0$ (Folland, 1999, Proposition 2.23(b)).

The result follows. □

Appendix 5.C Deferred Material from Section 5.2

Proof of Proposition 5.2.1. As stated in the proposition, set $X = [-1; +1]^2$, and H to be the two projection maps $h_1(x) = x_1$ and $h_2(x) = x_2$. Next define a set of positive instances $\{p_i\}_{i=1}^{\infty}$, and their corresponding probability mass:

$$p_i = \frac{1}{2} \cdot \frac{1}{4^{2^i}}; \quad (p_i) = 2^{-i-1};$$

Here are the negative instances:

$$n_i = \frac{1}{2} \cdot \frac{1}{3 \cdot 4^{2^i}}; \quad (n_i) = 2^{-i-1};$$

Notice that μ has countable support, and $\mu(X) = 1$. Furthermore, the vector $\mu = (\frac{1}{2}; \frac{1}{2})$ is a perfect separator: given any positive example p_i , $\mu(H(\cdot))(p_i) > 0$, and given negative example n_i , $\mu(H(\cdot))(n_i) < 0$. Note however that, as required by the proposition statement, the margins go to zero. However, given any $\epsilon \in \mathbb{R}_c$, since $\lim_{z \rightarrow 1} \mu(z) = 0$,

$$0 = \inf_{c \in \mathbb{R}_c} L(A) = \lim_{c \rightarrow 1} \int_{\mathbb{Z}} \mu(y_i(H(c))_{z_i}) d(z_i; y_i) = 0;$$

The key property of this construction is that the positive and negative examples are staggered; this will cause max margin solutions to avoid . As such, let any finite sample of size m be given. If all drawn examples have the same class, then $\hat{A} = (1 - y; 1 + y)$ (which is a maximum margin solution) has either n_1 or p_1 on the wrong side of the separator, and by choosing $c > 0$ large enough, $L(c\hat{A}) > b$ and $\mathbb{E}L(c\hat{A}) < 1 - b$.

As such, henceforth suppose there is at least one positive example, and at least one negative example. Suppose p_j and n_k respectively denote a sampled positive point and sampled negative point n_k having highest index among positive and negative examples; these maxima exist since m is finite.

Every max margin solution is determined solely by p_j and n_k . To obtain one of them, define

$$h := \frac{(1 + (n_k)_2) = (2 + (p_j)_1 + (n_k)_2)}{(1 + (p_j)_1) = (2 + (p_j)_1 + (n_k)_2)} i$$

To verify that this is a max margin solution, note that for any sampled (positive or negative) point z_i with label $y_i \in \{-1, +1\}$,

$$y_i(h \cdot z_i - (h \cdot p_j)) = (h \cdot n_k) = h \cdot n_k = \frac{(p_j)_1 (n_k)_2}{2 + (p_j)_1 + (n_k)_2} > 0.$$

By construction, however, $(p_j)_1 \neq (n_k)_2$, meaning h is not a rescaling of \hat{A} . As such, h is wrong for either all large p_i or n_i , and taking $\hat{A} = qh$ with q large, it follows that $L(\hat{A}) > b$ and $\mathbb{E}L(\hat{A}) < 1 - b$. □

Appendix 5.D Deferred Material from Section 5.3

Throughout this section, the following notation for measures will be employed

Definition 5.D.1. Given a finite measure μ and a set P , let μ_P be the restriction of μ to P : for any measurable set S , $\mu_P(S) = \mu(S \cap P)$. Note also that $d_{\mu_P}(x; y) = \mathbb{1}((x; y) \in P) d_{\mu}(x; y)$.

5.D.1 Primal Hard Cores

In light of the duality relationship for L (cf. Lemma 4.1.1), the definition for hard cores, provided in Section 5.3, is tied to the convex dual to L . Analogously, it is possible to define a primal form of hard cores, which will lead to a proof of Theorem 5.3.1.

Definition 5.D.2. Define $S_P(H; \cdot)$ to contain all sets C for which there exists a sequence $g_{i=1}^1$ satisfying the following properties.

1. Every g_i and $(x; y) \in C$ satisfies $y(H_{g_i})x = 0$.
2. For ϵ -almost-every $(x; y)$ in C^c , $y(H_{g_i})x > \epsilon$.

A primal hard core P is a (ϵ -a.e.) minimal set within $S_P(H; \cdot)$:

$$P \in S_P(H; \cdot) \quad \text{and} \quad \forall C \in S_P(H; \cdot) \quad (P \cap C) = \emptyset$$

Lemma 5.D.3. $S_P(H; \cdot)$ is closed under countable intersections.

Proof. To start, note that $S_P(H; \cdot)$ is closed under finite intersections as follows. Let $C_i, g_{i=1}^p$ be given with $C_i \in S_P(H; \cdot)$, and let corresponding sequences $g_{j=1}^{(i)}$ be given as granted by the definition of $S_P(H; \cdot)$. Define $C := \bigcap C_i$ and $g_j := \bigwedge_i g_j^{(i)}$. By construction, for every $(x; y) \in C$ and pair $(i; j)$, $y(H_{g_j^{(i)}})x = 0$, and thus $y(H_{g_j})x = 0$. Next, for each C_i , define $C_i^0 \subset C_i^c$ with $(C_i^0)^c = (C_i^c)^c$ so that, for every $(x; y) \in C_i^0$, $y(H_{g_j^{(i)}})x > \epsilon$. Correspondingly, define $C^0 := \bigcup C_i^0$, where $(C^0)^c = (C^c)^c$. Now let any $(x; y) \in C^0$ and any $B > 0$ be given. For each i , there are two cases: either this is an area where $y(H_{g_j^{(i)}})x > \epsilon$, or $y(H_{g_j^{(i)}})x = 0$. In the first case, let T_i denote an integer, as granted by $y(H_{g_j^{(i)}})x > \epsilon$, so that for all $j \leq T_i$, $y(H_{g_j^{(i)}})x > B$. For those i where $(x; y) \notin C_i^0$ (but still $(x; y) \in C^0$), due to the ruled out nullsets, $y(H_{g_j^{(i)}})x = 0$, safely set $T_i = 0$. To finish, taking $T := \max_i T_i$, it follows that for every $j > T$, $y(H_{g_j})x > B$, whereby it follows that $y(H_{g_j})x > \epsilon$ over C^0 , and thus over C^c ϵ -a.e.

Now let a countable family $\{D_i, g_{i=1}^1\}$ be given with $D_i \in S_P(H; \cdot)$, and define $D = \bigcap D_i$. Consider the optimization problem

$$p := \inf \int \exp(-y(H_{g_i})x) d_{D^c}(x; y) : \int \exp(-y(H_{g_i})x) d_{D^c}(x; y) < \infty; \forall (x; y) \in D, y(H_{g_i})x = 0$$

Define $E_j := \bigcap_{i \leq j} D_i$, whereby $D := \bigcap E_j$. Since $\int \exp(-y(H_{g_i})x) d_{D^c}(x; y) < \infty$, by continuity of measures from above (Folland, 1999, Theorem 1.8(d)), for any $\epsilon > 0$ there exists E_k with $\int \exp(-y(H_{g_i})x) d_{D^c}(x; y) > \int \exp(-y(H_{g_i})x) d_{E_k^c}(x; y) - \epsilon$. Since it was shown above that $S_P(H; \cdot)$ is closed under finite intersections, $E_k = \bigcap_{i \leq k} D_i \in S_P(H; \cdot)$; consequently, let $g_{i=1}^1$ to be a sequence of predictors certifying that $E_k \in S_P(H; \cdot)$, as according to the definition; since $E_k \supset D$, it follows that $y(H_{g_i})x = 0$ over

E_k and thus over D as well. It follows by the dominated convergence theorem (with dominating function $(x; y) \mapsto \exp(0)$) that

$$p = \lim_{i \rightarrow \infty} \int_{D^c} \exp(-y(H_i)x) d_{D^c}(x; y) = \int_{D^c} \exp(0) d_{D^c}(x; y) = \int_{D^c} 1 d_{D^c}(x; y) < \infty$$

Since ϵ was arbitrary, it follows that $p = 0$.

As such, for any $n \in \mathbb{Z}_{++}$, choose $\epsilon_n \in \mathbb{R}^n$ with $y(H_n)x = \epsilon_n$ over D satisfying

$$\int_D \exp(-y(H_n)x) d_{D^c}(x; y) < \epsilon_n^2$$

By Markov's inequality, it follows that

$$\int_{D^c} (\exp(-y(H_n)x) - \epsilon_n) d_{D^c}(x; y) < \epsilon_n$$

Since $\exp(-y(H_n)x) = \exp(0) = 1$ over D , it follows that $f_n(x; y) = \exp(-y(H_n)x) \mathbf{1}_{D^c}(x; y)$ is Cauchy in measure, and moreover converges in measure to the function $(x; y) \mapsto \mathbf{1}_{D^c}(x; y)$. Consequently, there exists a subsequence n_j with $\exp(-y(H_{n_j})x) \rightarrow \mathbf{1}_{D^c}(x; y)$ μ -a.e. (Folland, 1999, Theorem 2.30). This is only possible if $y(H_{n_j})x \rightarrow 0$ μ -a.e. $(x; y) \in D^c$, and the result follows, with $f_{n_j} \mathbf{1}_{D^c}$ as the certifying sequence for D , since every $y(H_{n_j})x = 0$ for $(x; y) \in D$ by construction. \square

Theorem 5.D.4. Every finite hypothesis class H and probability measure μ admits a primal hard core.

Proof. Consider the optimization problem

$$p := \inf \{ \int_C f(x; y) d_{D^c}(x; y) : C \in \mathcal{S}_P(H; \mu) \}$$

Since \mathcal{S}_P is nonempty (it always contains $X \times \{1\} + \mathbf{1}_g$ with certifying sequence $\epsilon_i = 0$ for every i) and μ is a finite nonnegative measure, the infimum is finite. Let $f_{n_j} \mathbf{1}_{D_j}$ be a minimizing sequence, and define $D_j := \bigcup_{i=1}^j C_i$ and $D := \bigcup_{j=1}^{\infty} D_j = \bigcup_{i=1}^{\infty} C_i$. By Lemma 5.D.3, $D_j \in \mathcal{S}_P(H; \mu)$ for every j , and since $D_j \subset D_{j+1}$, it follows that $f_{n_j} \mathbf{1}_{D_j}$ must also be a minimizing sequence to the above infimum. Finally, since μ is finite and Lemma 5.D.3 also grants $D \in \mathcal{S}_P(H; \mu)$, then

by continuity of measures from above (Folland, 1999, Theorem 1.8(d)),

$$(D) = \lim_{j \uparrow \infty} (D_j) = p:$$

Since $D \in S_P(H; \cdot)$ attains the infimum, it is a primal hard core. \square

With existence of primal hard cores out of the way, the next key is the equivalence to (dual) hard cores.

Theorem 5.D.5. Let finite hypothesis class H , probability measure μ , a (dual) hard core C , as well as a primal hard core P be given. Then C and P agree on all but a null set.

The proof uses the following technical lemma.

Lemma 5.D.6. Let finite hypothesis class H , probability measure μ , $C_1 \in S_P(H; \cdot)$, as well as a $g \in \mathbb{R}^n$ be given, with $y(H_2)x \geq 0$ for $(x; y) \in C_1$ (but potentially $y(H_2)x < 0$ elsewhere). Then $C_1 \cap [y(H_2)x > 0] \in S_P(H; \cdot)$.

Proof. Let $C_1; g$ be given as specified. Let $\{g_i^1\}_{i=1}^\infty$ be a certifying sequence for C_1 . Define $P := [y(H_2)x > 0]$ and $C_3 := C_1 \cap P = C_1 \cap [y(H_2)x > 0]$; the goal is to show $C_3 \in S_P(H; \cdot)$.

Let $i \in \mathbb{Z}_{++}$ be arbitrary; the following steps will construct $\{g_i^4\}$, a certifying sequence for C_3 , meaning $C_3 \in S_P(H; \cdot)$.

First, let c be sufficiently large so that $\{g_i^2\} := c \cdot g$ satisfies

$$\sum \exp(-y(H_2)x) \mu(x; y) < 1/i^2:$$

By Markov's inequality, it follows that

$$\mu([\exp(-y(H_2)x) \geq 1/i]) \leq \sum \exp(-y(H_2)x) \mu(x; y) < 1/i. \quad (5.D.7)$$

Consequently define $P_i := [y(H_2)x > \ln(i)]$, where the above statements show $\mu(P_i) > \mu(P) - 1/i$.

Next, since $\exp(-y(H_2)x) \leq 1$ $[(x; y) \in C_1]$ -a.e. and $\sum_{i=1}^\infty \exp(-y(H_2)x) < \infty$, by Egorov's Theorem (Folland, 1999, Theorem 2.33), this convergence is uniform over a subset $S_i \subset C_1$ with $\mu(S_i) > \mu(C_1) - 1/i$. In particular, there exists an integer

T_i so that, for any $(x; y) \in S_i \setminus C_1$,

$$y(H_{T_i}^{(1)})x > k_i^{(2)} k_1 + \ln(i):$$

As such, define $H_i^{(3)} := H_{T_i}^{(1)} + H_i^{(2)}$. First, for any $(x; y) \in C_3$ and any i ,

$$y(H_i^{(3)})x = 0 = y(H_i^{(1)})x = y(H_{T_i}^{(1)})x:$$

On the other hand, for any $(x; y) \in S_i \setminus C_1$, since $\sup_{x,y} |y(x)| = 1$,

$$\begin{aligned} y(H_i^{(3)})x &= y(H_{T_i}^{(1)})x + y(H_i^{(2)})x \\ &> k_i^{(2)} k_1 + \ln(i) - k_i^{(2)} k_1 = \ln(i): \end{aligned}$$

Lastly, as shown above, for any $(x; y) \in P_i$,

$$y(H_i^{(3)})x = 0 + y(H_i^{(2)})x = \ln(i):$$

Combining the above facts,

$$(\int \exp(-y(H_i^{(3)})x) \mathbb{1}_{[(x; y) \in C_3]} \mathbb{1}_{[i=1]} < (C_1^c \cap S_i) + (P \cap P_i) \mathbb{1}_{[i=1]}:$$

It follows that $\exp(-y(H_i^{(3)})x) \mathbb{1}_{[(x; y) \in C_3]} \mathbb{1}_{[i=1]}$ is Cauchy in measure (and the sequence of functions is Cauchy in measure), and thus there is a subsequence $\{i_{j=1}^{(4)}\}$ for which $\exp(-y(H_{i_j}^{(4)})x)$ converges to $\mathbb{1}_{[(x; y) \in C_3]} \mathbb{1}_{[i=1]}$ -a.e. (Folland, 1999, Theorem 2.30). It follows that $\{i_{j=1}^{(4)}\}$ is the desired sequence certifying that $C_3 \in \mathcal{S}_P(H; \cdot)$. \square

Proof of Theorem 5.D.5. If $(P \cap C) > 0$, then by the maximality of C , P is a set of positive measure away from any element of $\mathcal{S}_D(H; \cdot)$, and in particular $P \in \mathcal{S}(H; \cdot)$, and thus Theorem 5.B.1 provides the existence of $\mathbb{2} \in \mathbb{R}^n$ with $(P \setminus [y(H_{\mathbb{2}})x = 0]) = (P)$ and $(P \setminus [y(H_{\mathbb{2}})x > 0]) > 0$. But then, by Lemma 5.D.6, $P^0 := P \cap [y(H_{\mathbb{2}})x > 0]$ has $(P^0) < (P)$ but $P^0 \in \mathcal{S}_P(H; \cdot)$, which contradicts the -a.e. minimality of hard core P .

Now suppose $(C \cap P) > 0$, and let \cdot_C denote the restriction of \cdot to C : for any C ,

$c(C) := (C \setminus C)$. Consider the optimization problem

$$\inf_{z \in \mathbb{R}^n} \int_C \exp(-y(H)(x)) d_c(x; y)$$

Consider the sublevel set of 1-suboptimal points for this problem. By Theorem 5.B.1, there exists B so that each z in this sublevel set has $\|z\| \leq B$ with $H = H^0$ -a.e. and $k_1 \leq B$. However, by the definition of P , there exists a sequence $\{g_i\}_{i=1}^\infty$ which is zero over P and approaches 1 -a.e. over P^c , and in particular over the positive measure set $C \cap P^c$. Thus, taking any z in the 1-suboptimal set, notice that

$$\lim_{i \rightarrow \infty} \int_C \exp(-y(H + g_i)(x)) d_c(x; y) = \int_C \exp(-y(H)(x)) \mathbb{1}_{\{x; y \in P^c\}} d_c(x; y) =: p$$

Since z has a bounded representation $\exp(-y(H)(x)) \leq 0$, which combined with $(C \cap P^c) > 0$ grants $p < \int_C \exp(-y(H)(x)) d_c(x; y)$ (Folland, 1999, Theorem 2.23(b)). But since $\int_C \exp(-y(H)(x)) d_c(x; y)$ is continuous in z for any finite measure (cf. Lemma 5.A.2), there must exist a large j so that

$$\int_C \exp(-y(H + g_j)(x)) d_c(x; y) < \int_C \exp(-y(H)(x)) d_c(x; y);$$

and moreover $y(H + g_j)(x) > B$ for a subset of C with positive measure. But that means $z + g_j$ is in the 1-sublevel set (since it beats z , which was chosen to reside in the 1-sublevel set), but can not have a representation with norm at most B (since H is a bounded linear operator with operator norm 1), contradicting Theorem 5.B.1. \square

5.D.2 Proof of Theorem 5.3.1

This is now just a consequence of the equivalence to primal hard cores, and the structure over C developed in Theorem 5.B.1 (which was used to prove the equivalence to primal hard cores as well).

Proof of Theorem 5.3.1. The second property is direct from Theorem 5.B.1. For the first property, since primal hard cores exist and are -a.e. equivalent to (dual) hard cores (cf. Theorem 5.D.5), the statement thus follows by taking the sequence provided by the definition of any primal hard core. \square

Appendix 5.E Deferred Material from Section 5.4

Proof of Theorem 5.4.1. (Item 1) Let $f_{i,g_{i=1}^1$ be given as per Theorem 5.3.1. Automatically, $y(H_i)x = 0$ for $(x; y) \in C$. And since $y^0(H_i)x^{0^m} = 1$ for ϵ -a.e. $(x^0, y^0) \in C^c$, it follows from the definition of L_C that $\lim_{i \rightarrow \infty} \int (y^0(H_i)x) = 0$.

(Item 2) This is a consequence of Theorem 5.B.1. □

Proof of Theorem 5.4.2. (Item 1) Let a sequence $f_{i,g_{i=1}^1$ be given as provided by Theorem 5.3.1. In particular, $\exp(y(H_i)x) \leq 1$ for $(x; y) \in C$ ϵ -a.e.. Now choose a finite sample size m ; by Egorov's Theorem (Folland, 1999, Theorem 2.33), for any $\delta > 0$, there exists S with $\mu(S) > \mu(X) - \delta$ over which this convergence is uniform. As such, choose δ so that $\exp(y(H_i)x) < 1 - \delta$ over $S \setminus C^c$, meaning in particular $y(H_i)x > 0$ for every $(x; y) \in S \setminus C^c$. The probability over a draw of m points that some within C^c are misclassified by \hat{y} has upper bound

$$\Pr[\exists (x_i; y_i) \in C^c : y(H_i)x \leq 0] \leq \mu(C^c \setminus [y(H_i)x \leq 0]) < \delta$$

Since δ can be made arbitrarily small, the probability of failure is zero. Furthermore, since \hat{y} satisfies $y(H_i)x = 0$ ϵ -a.e. over C (cf. Theorem 5.3.1), it also follows that, with probability 1, \hat{y} abstains on every example falling within C .

(Item 2) Let $\epsilon > 0$ and $\delta \in L_C$ be given. Choose $\delta > 0$, as provided by Theorem 5.4.1, so that every $\hat{y} \in R^n$ with $L(\hat{y}; C) \leq L(y; C) + \delta$ has a representation $\hat{y} = \sum_{k=1}^m \alpha_k H_k$, where $H_k = H_k^0$ everywhere along $C \cap N$, where $\mu(N) = 0$; henceforth, discard the (measure zero) event that any example falls within N . Additionally, choose $c > 0$ as provided by Proposition 5.A.3 so that, given m_C i.i.d. points within $C \cap N$, every $\hat{y} \in R^n$ with $\|\hat{y}\| \leq b$ satisfies, with probability at least $1 - \delta$,

$$|L(\hat{y}; C) - L(y; C)| \leq c^0(C) \frac{\sum_{k=1}^m \frac{1}{\ln(n)} + \sum_{k=1}^m \frac{1}{\ln(2^m)}}{m_C}; \tag{5.E.1}$$

where the extra $c^0(C)$ comes from applying Proposition 5.A.3 to \hat{y} conditioned on points falling within C , which can be related to the above left hand side via this normalization by $c^0(C)$. Henceforth, for simplicity, set $c := c^0(C)$.

Now consider any $\hat{y} \in R^n$ with no representation $\hat{y} = \sum_{k=1}^m \alpha_k H_k$ so that $\hat{y} \neq H^0$ over $C \cap N$, which directly entails, by Theorem 5.4.1, that $L(\hat{y}; C) \leq L(y; C) + \delta > L(y; C) + 4\delta$. Additionally choose

and any $\beta \in \mathbb{R}^n$ with $L(\hat{A}; C) - L(A; C) < 1$, whereby the choice of $\beta > 0$ indicates that, without loss of generality, $k_1 \leq \beta$. Since $\mathbb{R}^n \ni (A, d_C)$ is continuous (cf. Lemma 5.A.2), considering the line segment $\beta + (1 - \beta) \cdot \mathbb{1}_{[0; 1]}$, there must exist \hat{A} along this line segment with with

$$\beta + 3 L(\hat{A}; C) - L(A; C) > \beta + 4;$$

let \hat{A}^0 be a representation with $k_1 \leq \beta$ and $H^{\hat{A}} = H^{\hat{A}^0}$ over $C \times \mathbb{N}$ (and thus it holds for every example). Applying the deviation inequality in eq. (5.E.1) twice,

$$\begin{aligned} \mathbb{P}(A^{\hat{A}}; C) - \mathbb{P}(A; C) &\leq L(A^{\hat{A}^0}; C) - L(A; C) + 2c \frac{p \overline{\ln(n)} + p \overline{\ln(2=)}}{m_C} \\ &= L(A^{\hat{A}^0}; C) - L(A; C) + (L(A; C) - L(A; C)) \\ &\quad + 2c \frac{p \overline{\ln(n)} + p \overline{\ln(2=)}}{m_C} \\ &> (\beta + 3) - (1) + 2c \frac{p \overline{\ln(n)} + p \overline{\ln(2=)}}{m_C} \\ &\quad ; \end{aligned}$$

where the last step used the lower bound βm_C . Returning to \mathbb{R}^n as specified above, convexity, in the form of Lemma 5.A.5, grants that $\mathbb{P}(A; C) < \mathbb{P}(A^{\hat{A}}; C)$ implies $\mathbb{P}(A^{\hat{A}}; C) - \mathbb{P}(A; C) > 0$, and thus

$$\mathbb{P}(A; C) - \mathbb{P}(A; C) + \mathbb{P}(A^{\hat{A}}; C) - \mathbb{P}(A; C) > 0;$$

Since β was arbitrary, it follows that every β with no representation $k_1 > \beta$ that has agreement of $H^{\hat{A}}$ and $H^{\hat{A}^0}$ -a.e. over C does not lie in the empirical β -sublevel set. Since $\mathbb{P}(\cdot; C)$ is convex and continuous (cf. Lemma 5.A.2), the β -sublevel set is nonempty, and thus every β within it has a representation $k_1 \leq \beta$. □

Appendix 5.F Deferred Material from Section 5.5

Proof of Proposition 5.5.2. This proof is essentially a repackaging of various results and comments due to Bartlett et al. (2006). Fix any $\beta \in \mathbb{R}_c$; β is convex, increasing at 0, and differentiable at 0, which grants that the corresponding β -transform is classification calibrated (Bartlett et al., 2006, Theorem 6, although note losses in this thesis are increasing rather than decreasing). It follows

that $(R(f) - R(F)) - L(f) - L(F_Y)$; (Bartlett et al., 2006, Theorem 3, part 3(c)).

Next, $\phi(0) = 0$ (Bartlett et al., 2006, Lemma 5, part 8), $\phi(r) > 0$ when $r > 0$ (Bartlett et al., 2006, Lemma 5, part 9(b)), and since ϕ is convex by construction (Bartlett et al., 2006, Definition 2), it follows by Lemma 5.A.5 that ϕ is increasing. Since ϕ is continuous as well, (Bartlett et al., 2006, Lemma 5, part 6), it follows that ϕ has a well-defined inverse along the image $\phi([0; 1])$. Finally, the fact that $\phi^{-1}(r) \neq 0$ as $r \neq 0$ is due to Bartlett et al. (2006, Theorem 3, part 3(b)). □

Proof of Theorem 5.5.4. Throughout this proof, $\delta := \epsilon/8$ will be the failure probability of various crucial events; the final statement is obtained by unioning them together, and subsequently throwing them all out. Note also that some of the statements vacuously hold if $\phi(C) = 0$ or $\phi(C) = (X - f - 1; +1)g$ (i.e., when terms depending on either appear in denominators); interpret these expressions as simply being ∞ , whereby the bounds hold automatically.

(Item 1) Let S_C and S_{C^c} respectively denote the set of samples landing in C and C^c , where the notation proposed in the theorem statement provides $m_C = |S_C|$ and $m_{+} = |S_{+}|$. By a Chernoff bound (Kearns and Vazirani, 1994, Theorem 9.2), basic deviations for these quantities are

$$\begin{aligned} \Pr[|S_C| < (\phi(C) - \delta)m] & \leq \exp(-m^2 \delta^2 / 2) = \delta^q, \\ \Pr[|S_{+}| < (\phi(C^c) - \delta)m] & \leq \exp(-m^2 \delta^2 / 2) = \delta^q, \end{aligned}$$

where $\delta = \frac{q}{2m} \ln \frac{1}{\delta}$, and \Pr denotes the product measure corresponding to m copies of \mathcal{X} . Label these failure events F_1 and F_2 , and henceforth rule them out.

(Item 2) As provided by Theorem 5.4.2, there exists $\gamma \in \mathbb{R}^n$ with $\gamma_i(H) x_i > 0$ for all $(x_i; y_i)$ falling in C^c , and $\gamma_i(H) x_i = 0$ for those landing in C . Consequently,

$$\begin{aligned} \mathbb{E}(A) &= \inf_{c > 0} \mathbb{E}(A(\gamma + c); C) + \mathbb{E}(A(\gamma + c); C^c) \\ &= \inf_{c > 0} \mathbb{E}(A(\gamma); C) \\ &= \mathbb{E}(A) : \end{aligned}$$

Combining this with

$$\mathbb{P}(A; C^c) + \mathbb{P}(A; C) = \mathbb{P}(A) + \mathbb{P}(A);$$

it follows that

$$\mathbb{P}(A; C^c) - \mathbb{P}(A) = \mathbb{P}(A; C) - \mathbb{P}(A);$$

Next, since $\hat{\rho}(0) > 0$ and $\hat{\rho}$ is nondecreasing (cf. Lemma 5.A.1),

$$\hat{\rho}(H; C^c) = \frac{\mathbb{P}(A; C^c)}{\hat{\rho}(0)} = \frac{\mathbb{P}(A; C) - \mathbb{P}(A)}{\hat{\rho}(0)};$$

To obtain eq. (5.5.5) from here, first notice that S_+ , the portion of the sample falling within C^c , can be interpreted as an i.i.d. sample from the probability measure $(\cdot \mid C^c) = (\cdot) - (C)$. Next, the VC dimension of $\text{span}(H)$ is the VC dimension of linear separators over the transformed space

$$f((h_1(x); h_2(x); \dots; h_n(x)); y) : (x; y) \in X \times \{-1, +1\};$$

namely, it is n . As such, eq. (5.5.5) follows by an application of a relative deviation version of the VC Theorem (Boucheron et al., 2005a, discussion preceding Corollary 5.2). albeit with an adjustment to apply this uniform deviation bound with $\hat{\rho}$ conditioned on C^c , and then rescaling to remove this conditioning (just as in the proof of Theorem 5.4.2).

To obtain eq. (5.5.6), note that $\hat{\rho}(0) = m$ means there are no mistakes over C^c :

$$\hat{\rho}(0) > m \implies \max_{i \in [m_+]} \hat{\rho}(y_i(H) x_i) > 0;$$

that is to say, since $\hat{\rho}$ is nondecreasing, for every $x_i; y_i \in S_+$, $0 < y_i(H) x_i$. Plugging $\hat{\rho}(H) = 0$ into the same relative deviation bound as before (Boucheron et al., 2005a, discussion preceding Corollary 5.2), the second bound follows.

(Item 3) By Theorem 5.4.2, there exist constants $b > 0$ and $c(b)$, depending on

$H; ; C$, so that with probability at least $1 - \epsilon$, if $m_C \geq c^2(\ln(n) + \ln(1/\epsilon))$, then every ϵ -suboptimal predictor over C , and in particular \hat{A} , has a representation A^ϵ which is equivalent to \hat{A} -a.e. over C , and satisfies $\|A^\epsilon - \hat{A}\|_1 \leq b$. As such, since

$$\mathbb{E}L(A; C) = \mathbb{E}L(A^\epsilon; C) \quad \text{and} \quad \mathbb{E}L(A; C) = \mathbb{E}L(A^\epsilon; C);$$

an application of Proposition 5.A.3 grants

$$\begin{aligned} L(A; C) &= L(A^\epsilon; C) \\ &= \mathbb{E}L(A^\epsilon; C) + \frac{c \sqrt{\frac{2 \ln(n)}{m_C}} + 4 \sqrt{\frac{2 \ln(1/\epsilon)}{m_C}}}{\sqrt{m_C}} \\ &= \mathbb{E}L(A; C) + \frac{c \sqrt{\frac{2 \ln(n)}{m_C}} + 4 \sqrt{\frac{2 \ln(1/\epsilon)}{m_C}}}{\sqrt{m_C}} \\ &= \mathbb{E}L(A; C) + \frac{c \sqrt{\frac{2 \ln(n)}{m_C}} + 4 \sqrt{\frac{2 \ln(1/\epsilon)}{m_C}}}{\sqrt{m_C}}; \end{aligned}$$

where as in the proof of Theorem 5.4.2, the constant also captures a correction since Proposition 5.A.3 is applied to condition on points falling in C . Next, noting that Theorem 5.4.1 provides that a minimizing sequence to $L(A; C)$ can be taken without loss of generality to lie within a compact set (e.g., points with l_1 norm at most b); combined with the continuity of $\int \mathbb{R}^d(A) d_C$ (cf. Lemma 5.A.2), it follows that a minimizer exists. By an application of McDiarmid's inequality, with probability at least $1 - \epsilon$,

$$\mathbb{E}L(A; C) \leq \mathbb{E}L(A; C) + L(A; C) + c \frac{\sqrt{2 \ln(1/\epsilon)}}{\sqrt{m_C}};$$

(Note, $L(A; C)$ is independent of the sample, thus McDiarmid succeeds, with constant $c \sqrt{2 \ln(1/\epsilon)}$ (b) since $L(A; C)$ is in this initial sublevel set.) Combining these two pieces, it follows that

$$L(A; C) \leq \mathbb{E}L(A; C) + \frac{c \sqrt{\frac{2 \ln(n)}{m_C}} + 4 \sqrt{\frac{2 \ln(1/\epsilon)}{m_C}}}{\sqrt{m_C}};$$

which is precisely eq. (5.5.7).

To produce eq. (5.5.8), the definition of the ϵ -transform (cf. Proposition 5.5.2), combined

with Equation (5.5.7), provides

$$\begin{aligned}
 R(H; C) - R(F; C) &= L(A; C) - L(F_Y; C) \\
 &= L(A; C) - L(A; C) + L(A; C) - L(F_Y; C) \\
 &= \frac{c^p \ln(n) + 4^p \ln(2^q)}{m_C} + L(A; C) - L(F_Y; C)
 \end{aligned}$$

(Item 4) Combining the lower bound on m with Item 1,

$$\begin{aligned}
 m_+ &= m(C^c) = 2; \\
 m_C &= m(C) = 2 \cdot c^2 (\ln(n) + \ln(1 - \epsilon));
 \end{aligned}$$

the first two bounds will allow expressions to be simplified, whereas the last bound will allow an invocation of item 3.

As such, combining all preceding bounds (and making use of the refinement over C^c when $\epsilon = 0$),

$$\begin{aligned}
 R(H) - R(F) &= (R(H; C) - R(F; C)) + (R(H; C^c) - R(F; C^c)) \\
 &= \frac{c^p \ln(n) + 4^p \ln(2^q)}{m_C} + L(A; C) - L(F_Y; C) \\
 &\quad + (C^c) \frac{4(n \ln(2m_+ + 1) + \ln(4 - \epsilon))}{m_+} \\
 &= (C^c) \frac{c^p \frac{m_+}{2} \ln(n) + 4^p \ln(2^q)}{m(C)} + L(A; C) - L(F_Y; C) \\
 &\quad + \frac{8(n \ln(m(C^c) + 1) + \ln(4 - \epsilon))}{m};
 \end{aligned}$$

□

Appendix 5.G Bibliographic Notes

As discussed at the beginning of the chapter, there are a variety of ways to produce general statistical guarantees for boosting methods; results of this type will be discussed in more

detail in the bibliographic remarks in Chapter 6.

Focusing then on the finite-dimensional case with a goal of handling general losses, there are very few results, since the problem is not regularized or constrained in any way. Similarly, there are statistical guarantees in the presences of margins bounded below by a positive constant (Schapire et al., 1997), but this of course does not hold in general. If the goal is relaxed to consistency and just handling the logistic loss, the only results appear to be those related to the consistency of maximum likelihood (Lebanon, 2008, Gourieroux and Monfort, 1981), which require a number of conditions that may fail in general (in particular, the minimum of the optimization problem may fail to exist, which corresponds to a correct model failing to exist).

Note that discussion on hard cores and some of their history may be found in Section 1.E.

Acknowledgements

This chapter is based on work by the dissertation author which is currently under preparation.

Chapter 6

Infinite Hypothesis Classes

This chapter provides numerical and statistical guarantees while allowing both the support of μ and the size of H to be infinite; however, the guarantees only hold for losses which, like the logistic loss, are Lipschitz.

As has been the case throughout this thesis, the strategy is to investigate the structure imposed on the optimization problem by μ and H , and thereafter provide the numerical and statistical guarantees.

Unlike the preceding chapters, hard cores are not used, but the structure is similar. The first case is when the optimal convex risk over the distribution, $L(A)$, is zero, which is analogous to the case that all hard cores have measure zero (cf. Proposition 6.3.6). In this situation, the analysis exhibits a quantity similar to the weak learning rate which can once again establish quick decrease of the primal objective function. On the other hand, when $L(A) > 0$, the analysis shows that every sample with high probability has some tricky examples, which in turn introduce curvature and constrain the norms of the chosen predictors; from here, a convergence analysis follows readily.

6.1 Overview

This chapter will establish statistical properties of Boost for any of the three line searches, together with the second stopping condition which ensures roughly m^a iterations occur. The basic consistency guarantee, as well as an overview of the analysis, appear in Section 6.2. The class of loss functions to which the results apply is as follows.

Definition 6.1.1. Let L_{2d} denote twice continuously differentiable convex losses. Additionally, let L_{lg} contain all differentiable convex Lipschitz losses $\ell : \mathbb{R} \rightarrow \mathbb{R}_+$ with tightest Lipschitz

constant $B_2 := \sup_{x \neq y} |\eta(x) - \eta(y)| = |x - y|$ as follows. First, every η has $\eta \in [0; 2]$ everywhere and $\eta \in [1; 2]$ over \mathbb{R}_+ for some $0 < \eta_1 < \eta_2$. Second, every η has Lipschitz gradients with (tightest) parameter B_2 , meaning $B_2 := \sup_{x,y \in \mathbb{R}} (\eta'(y) - \eta'(x)) = (y - x) < 1$.

Although the most general guarantees require $\eta \in L_{1g} \setminus L_{2d}$, the separable case needs only $\eta \in L_{1g}$, which allows consideration of the interesting piecewise quadratic $\text{loss}_{\text{russ}}(x)$, which was used by Impagliazzo (1995, Proof of Lemma 1) in the study of hard cores.

Section 6.3 considers the separable case $\eta(A) = 0$. Similarly to the analogous case in Chapter 5, the convex risk L is not controlled in this scenario, but rather only the classification risk R . The analysis has a minor handicap; rather than applying to the scalar coefficients η_t , it applies to the coefficient vector $\hat{\eta}$, defined as follows.

Definition 6.1.2. Given the sequence of coefficient vectors $\{g_t\}_{t=1}^T$, let $\hat{\eta}$ denote the vector achieving the best empirical classification risk, meaning $\hat{\eta} := \arg \min_{\eta \in \mathcal{H}} \sum_{t=1}^T \eta(g_t)$.

While this additional indirection $|\hat{\eta}|$ which can be computed easily as the algorithm progresses $|\hat{\eta}|$ is perhaps unnecessary, it leads to a simple analysis in the case $\eta(A) = 0$. In light of the fact that similar structure arose in Chapter 5, it is interesting to wonder if, in this regime, boosting methods are in fact better adapted to R than to L .

In the separable case, the hypothesis class must have finite VC dimension, denoted as follows.

Definition 6.1.3. Let F_{vc} contain all classes \mathcal{H} of finite VC dimension, denoted $V(\mathcal{H}) < \infty$.

The nonseparable follows in Section 6.4, and holds for losses with $\eta \in L_{2d} \setminus L_{1g}$. Additionally, the hypothesis class must not only have finite VC dimension, indicating it is not too large, it must also satisfy the following condition which indicates it is not too small.

Definition 6.1.4. Let $F_{\text{ds}}(\cdot)$ contain every class \mathcal{H} whose linear span $\text{span}(\mathcal{H})$ is dense (in the $L^1(\cdot)$ topology) in the collection of all bounded measurable functions over X .

Conditions similar to those defining $F_{\text{ds}}(\cdot)$ are usually called dense class assumptions (Bartlett and Traskin, 2007, Condition 1, Denseness), or completeness assumptions (Breiman, 2000, Definition 1); for a more extensive discussion of these conditions, please see Section 6.A; for the time being, the important point is that reasonable elements of $F_{\text{vc}} \setminus F_{\text{ds}}(\cdot)$ exist; in particular, the following result provides that if $X = \mathbb{R}^d$, then decision lists and decision trees with d axis-aligned splits suffice.

Proposition 6.1.5. Suppose $X = \mathbb{R}^d$, and let μ be a Borel probability measure over X . If $\text{span}(H)$ contains all indicators of products of half-open intervals of the form $\prod_{i=1}^d [a_i; b_i)$, where $a_i < b_i$, then $H \in \mathcal{F}_{\text{ds}}(\mu)$.

As a final bit of notation, define the shorthands $L := L(A)$ and $R := R(H)$.

Lastly, note that this section will be forced to develop its own convergence analysis of Boost, and in particular can not use those present in Chapter 3. The reason is simply that the many "constants" in the bounds there are dependent on the sample, and more importantly do not concentrate with high probability in a way which allows those bounds to be used here.

6.2 Consistency Statement and Analysis Sketch

The analysis considers two cases: either $L = 0$ (separable) or $L > 0$ (nonseparable). (By Proposition 6.B.3, $L = 0$ implies finite samples have a separating choice ≥ 2 almost surely.)

When the instance is separable, the improvement in objective value $L(A)$ in early iterations may be lower bounded by a margin-based quantity related to the classical weak learning rate; while this quantity is a random variable, with high probability it can be lower bounded by the analogous quantity over the distribution (which will be shown positive if the instance is separable). The bulk of the analysis is in constructing and controlling this quantity; the optimization and generalization analysis thereafter is straightforward, yielding a rate of roughly $O(1/m^{1-\beta})$ when $\beta = 2/3$.

When the instance is not separable, every weak learner makes a fair number of mistakes, and thus the algorithm makes more hesitant progress. Concretely, with high probability, the norms of the iterates are bounded, and moreover the quantity $\frac{1}{m} \sum_{i=1}^m \langle \nabla_{x_i, y_i} \ell(A), \nabla_{x_i, y_i} \ell(A) \rangle$, which is roughly the Hessian in axis-aligned directions (and relevant to coordinate descent), is also lower bounded. This in turn allows adaptation of the optimization analysis due to Zhang and Yu (2005). While the rate in this case is still roughly $O(1/m^{1-c})$, unfortunately the exponent $c > 1$ depends both on μ and on H (but is of course finite).

As a final point of interest, each case, in order to respectively establish either fast decrease or the norm constraints, considers the behavior of the reweighted average margins

$$\int (A) p d = \int (A)_{x,y} p(x;y) d(x;y); \quad (6.2.1)$$

where $\epsilon \in (0, 1]$ and $p \in L^1(\mathcal{X})$. In the separable case, this quantity is studied for a single good as p varies, whereas the nonseparable case studies a single bad as p varies.

Combining these finite sample results with the Borel-Cantelli lemma gives the following.

Theorem 6.2.2. Let $\ell \in L_{lg} \setminus L_{2d}$, probability measure μ over $\mathcal{X} = \{1, +1\}^g$, binary class $H \in \mathcal{F}_{vc} \setminus \mathcal{F}_{ds}(\mathcal{X})$, and any stopping parameter $\alpha \in (0, 1]$ be given. Let A_m denote the output of Figure 1.3 when run on m examples, and let R_α to denote the Bayes error rate. Then $R(A_m) \leq R_\alpha + \alpha^m$ almost surely.

6.3 The Separable Case ($L = 0$)

The rates in the case $L = 0$ will depend on the following quantity $\mathcal{Z}(\mu, H, \epsilon, \alpha)$, which directly embeds the reweighted margin expression in eq. (6.2.1). This quantity in turn depends on a set of reweightings $D(\mu, H, \epsilon, \alpha)$ which is a slightly constrained version of the set $D(H; \epsilon)$ used to define hard cores in Section 1.7.

Definition 6.3.1. Let μ be any probability measure over $\mathcal{X} = \{1, +1\}^g$ (relevant choices are μ and b), and let $\alpha \in [0, 1]$ be given. Define a permissible set of densities (with respect to μ)

$$D(\mu, H, \epsilon, \alpha) := \{p \in L^1(\mathcal{X}) : p \geq 0 \text{ -a.e.; } \int p = 1; \int p \ell \geq 1 - \alpha\}$$

with the convention $1/0 = 1$ in the case $\alpha = 0$. Additionally define

$$\mathcal{Z}(\mu, H, \epsilon, \alpha) := \inf_{p \in D(\mu, H, \epsilon, \alpha)} \sup_{h \in H} \int p(h - \ell)$$

(When μ is a discrete measure, $\mathcal{Z}(\mu, H, \epsilon, \alpha)$ is almost equivalent to AvgMin_k as developed by Shalev-Shwartz and Singer (2008, Section 4.1).)

This quantity will play a role analogous to the weak learning rate γ in AdaBoost, which guarantees the algorithm makes speedy progress in certain separable cases. The correspondence between these two quantities will occupy much of this section; but first, note the primary guarantee in the separable case.

Theorem 6.3.2. Let $\ell \in L_{lg}$ (with parameters $\epsilon_1; \epsilon_2; B_2$) and any $H \in \mathcal{F}_{vc}$ be given, and suppose $L = 0$. Let any error tolerance $\epsilon \in (0, 1]$ and any confidence parameter $\alpha \in (0, 1]$ be given, and

for convenience set $\delta := \epsilon/2$; by these choices, $\delta > 0$. Suppose Figure 1.3 is run with stopping parameter $\epsilon \in (0, 1)$, and the sample size m satisfies

$$m \geq \max \left(\frac{2}{(\delta/2)^2} \ln \frac{4}{\delta} ; \frac{24B_2^*(0)}{(\delta/2)^2} \right)^{1/a} :$$

Then, with probability at least $1 - \delta$, the algorithm's output \hat{h} satisfies

$$R(\hat{h}) \leq \frac{r}{m^{1-a}} \frac{(V(H) + 1) \ln(2em) + \ln(8/\delta)}{m^{1-a}} + 4 \frac{(V(H) + 1) \ln(2em) + \ln(8/\delta)}{m^{1-a}} :$$

To simplify this bound, first note that, for the logistic loss, $B_2 = 1/4$, $\epsilon = 1/2$, and $\delta = 1/2$ (cf. Lemma 6.B.1). Ignoring these terms, as well as $\delta(0) = 1$, (which can be set to $1/2$), and δ ; the choices $a = 1/2$ and $r := O(m^{-a/2}) = O(m^{-1/4})$ grant that $O(m^{1/2})$ iterations suffice to achieve classification risk $O(m^{-1/2})$, whereas the choices $a = 2/3$ and $r := O(m^{-a/2}) = O(m^{-1/3})$ provide that $O(m^{1/3})$ suffice to achieve error $O(m^{-1/3})$.

6.3.1 The Quantity $\epsilon(\delta)$

To develop the meaning and necessity of $\epsilon(\delta)$, first recall the classical notion of weak learnability, as defined in Section 1.7. Both the weak learning rate $\epsilon(\delta)$ and its empirical analogue $\epsilon(b)$ are close in definition of $\epsilon(\delta)$ and $\epsilon(b)$; indeed, the latter subscripted quantities replace the inequalities and quantifiers in the former quantities with maxima and suprema, which leads to the following correspondence.

Proposition 6.3.3. Let class H , probability measure μ (over X of size n), and empirical counterpart b be given. Then the weak PAC-learning rate $\epsilon(\delta)$ satisfies $\epsilon(\delta) \leq \epsilon(b)$, and the empirical weak learning rate $\epsilon(b)$ satisfies $\epsilon(b) \leq \epsilon(\delta)$.

The following example highlights why $\epsilon(\delta)$ can be problematic, even when $L = 0$.

Example 6.3.4. (Nightmare scenario #1.) Suppose $X = (0; 1]$, and

$$\Pr[y = +1 | x] = \begin{cases} 1 & \text{when } x \in (1/(i+1); 1/i] \text{ for some odd integer } i; \\ 0 & \text{when } x \in (1/(i+1); 1/i] \text{ for some even (positive) integer } i. \end{cases}$$

Let H consist of threshold functions (decision stumps). Given any integer k , a combination of thresholds may be constructed which is correct on k intervals, and thus $L = 0$ by considering

$k \geq 1$. Unfortunately, the norm of these solutions also grows unboundedly, suggesting and () are tiny. Indeed, consider a distribution over X which is uniform on k of the intervals, and zero elsewhere. Any threshold is incorrect on nearly half of these intervals, and by considering $k \geq 1$, it follows that $\phi(\cdot) = 0$.

In precise terms, this nightmare, and suggested sequence of distributions, provide the following property.

Proposition 6.3.5. There exist choices for H and γ so that $L = 0$, but $\phi(\cdot) = 0$ and, with any probability $1 - \epsilon$ and sample size m large enough that $k := \lfloor \frac{m^{1-\epsilon}}{3 \ln(4-\epsilon)} \rfloor$ satisfies $k \geq 2$, then (b) $\phi(b) = O((\ln(4-\epsilon) + \ln(m))^{1-\epsilon})$, where $\epsilon = O(1/k) = O(m^{-1-\epsilon})$ and the $O(\cdot)$ only suppresses terms independent of m and ϵ .

But something is wrong here | Example 6.3.4 seems quite easy! The reason otherwise is that it simply tries too hard: Example 6.3.4 is easy if giving up on an ϵ -fraction of the data is acceptable. This reasoning leads to the relaxation (), which, in contrast to Proposition 6.3.5, carries the following guarantee.

Proposition 6.3.6. Let probability μ over X be given and class H be given.

1. Let loss $\ell \geq L_{\mu}$ be given. Then $L = 0$ if $\mu(C) > 0$ for all $C \in H$ with $\ell(C) = 0$ for every hard core C .
2. Let any $\epsilon > 0$, confidence parameter $\delta \in (0, 1]$, and empirical measure $\hat{\mu}$ be given. Then with probability at least $1 - \delta$,

$$(b) \quad \ell(\hat{\mu}) \leq \frac{1}{2m} \ln \frac{2}{\epsilon} :$$

Note that the first item above also indicates that the two cases considered in this chapter could equivalently be stated in terms of hard cores: the separable case $L = 0$ is equivalent to all hard cores having zero measure, and the nonseparable case $L > 0$ is equivalent to all hard cores having positive measure.

In order to prove this result, and also a few other components in the proof of Theorem 6.3.2, the following dual representation of () is used. A similar result was proved by Shalev-Shwartz and Singer (2008, see the quantity AvgMin_k) in the case of measures with finite support and

nite cardinality hypothesis classes; the proof here invokes Sion's Minimax Theorem (Komiya, 1988), which operates in fairly general topological vector spaces.

Lemma 6.3.7. Let probability measure μ over X , $f \in \mathcal{L}_g$, any H , and any $\gamma \in (0, 1]$ be given. Then

$$\begin{aligned} \mathcal{R}_\gamma(H) &= \min_{p \in \mathcal{D}(X)} \sup_{Z} \int (A^\gamma) p d\mu : \int Z d\mu \leq \gamma; k_1 \leq 1 \\ &= \sup_{p \in \mathcal{D}(X)} \min_{Z} \int (A^\gamma) p d\mu : \int Z d\mu \leq \gamma; k_1 \leq 1 \\ &= \min_{p \in \mathcal{D}(X)} k A^\gamma p k_1; \end{aligned}$$

where A^γ is the unique adjoint operator to A (cf. Lemma 1.A.5), and

$$k A^\gamma p k_1 = \sup_{h \in H} \int h(x) p(x; y) d\mu(x; y)$$

In order to use this to prove the first part of Theorem 6.3.2, first note that whenever $\mathcal{R}_\gamma(H) = 0$, there exists a dual element $p \in \mathcal{D}(X)$ certifying this property, which in turn can be related to the duality structure of L (as presented in Lemma 4.1.1), and gives the result. For the second part of Theorem 6.3.2, similarly the infimum in the definition of $\mathcal{R}_\gamma(H)$ can be removed by considering a single bad certificate $p \in \mathcal{D}(X)$, and the supremum can be removed by considering a single good Z . The certificate p can be shown to have a simple structure (it emphasizes margin violations for the fixed good Z), and in turn the deviations are easy to control.

6.3.2 Proof Sketch of Theorem 6.3.2

The pieces are in place to establish the finite sample guarantees in Theorem 6.3.2. First, note the following empirical risk guarantee.

Lemma 6.3.8. Let any $\gamma \in \mathcal{L}_g$, empirical measure b , and H be given. Suppose Figure 1.3 is run with any of the three step size choices for T iterations, let $\mathcal{R}_t = \mathcal{R}(H_t)$ denote the classification error of H_t , and set $\mathcal{R}_t^0 := \mathcal{R}_{t-1}$ for convenience. Then

$$\mathcal{R}(A_T) \leq \mathcal{R}_T^0 + \frac{\sum_{t=1}^T (\mathcal{R}_{t-1} - \mathcal{R}_t^0)(b)^2}{6B_2}$$

Notice that this result indicates that the convex risk decreases quickly in the presence of

classification errors. The proof, sketched as follows, is fairly straightforward. First, standard properties of the line search choices show that $\|b\|$ drops in round t proportionally to $kA^T r \|b(A_{t-1})\|_{k_1}$. Considering $r \|b(A_{t-1})\|$ as a reweighting of b , this expression appears in the dual form of dual form of (b) as presented in Lemma 6.3.7. In order to make the correspondence precise $\|b(A_{t-1})\|$ must be rescaled to unit norm; but, by the Lipschitz property, the rescaling is by at most $\|r\|_{k_1}$. After some algebra, and summing across all iterations, the result follows.

From here, there is little to do. By Lemma 6.3.8, until some iteration has low error, progress is quick. The selection rule (returning \hat{w} with minimal classification risk) ensures there are no problems if the classification risk happens to go back up, and Proposition 6.3.6 allows (c) to replace (b). As this reasoning provides a direct guarantee on the empirical classification risk, standard uniform convergence techniques give the result.

6.4 The Nonseparable Case ($L > 0$)

When $L > 0$, the essential object will be an optimum to the convex dual of the central optimization problem $\inf L(A)$, specified as follows. The following slight refinement of the duality result from Lemma 4.1.1 additionally specifies a few properties of the dual optimum.

Proposition 6.4.1. Let loss $\ell \in L_{lg}$ (with tightest Lipschitz parameter ℓ_2), class H , and probability measure μ over $X \times \{-1, +1\}$ be given. Then

$$\inf_{A \in \mathcal{A}} \ell(A) = \max_{\{p \in \mathbb{R}^d : \|p\|_2 \leq L^{-1}(\cdot); p \in [0, \ell_2] \text{ -a.e.}; kA^T p\|_{k_1} = 0\}} \ell^*(p)$$

where ℓ^* is the Fenchel conjugate to ℓ , and the adjoint A^T is as in Lemma 6.3.7 and Lemma 1.A.5. Additionally, the dual optimum p satisfies $\langle p, y_i \rangle \geq \ell_2$, where $\ell_2 > 0$ whenever the optimal value L is positive, and moreover p has the explicit form $p := \ell^{-1}(\ell_2)$, where ℓ^{-1} is the (well-defined) inverse of ℓ along $[0, \ell_2]$.

The strategy in the nonseparable case is to exhibit curvature in the objective function (i.e., a lower bound on the second-order expression $\sum_i p_i \ell((A_{x_i, y_i}))$), and the dual optimum p will provide the mechanism; making these statements precise is the topic of this section.

Theorem 6.4.2. Let loss $\ell \in L_{lg} \setminus L_{2d}$, binary class $H \in \mathcal{F}_{vc}$, probability measure μ over $X \times \{-1, +1\}$ with empirical counterpart b corresponding to a sample of size n , time horizon

Let $\alpha \in (0, 1)$, and any confidence level $1 - \alpha \in (0, 1]$ be given. Suppose $\epsilon > 0$, and let $p \in L^1(\mathcal{Z})$ denote the dual optimum as in Proposition 6.4.1, with corresponding real number β so that $(\beta, \gamma) \in \mathcal{S}$. Define the quantities

$$c := \frac{16\beta(0)}{\beta(0)} \max_{k \in \mathcal{K}} \{1; \frac{1}{\beta} \max_{k \in \mathcal{K}} \{k p_{k_1} g\}\}; \quad B_1 := \frac{1}{8} \inf_{z \in \mathcal{Z}[\beta; \gamma]} \beta(z);$$

$$R_i := \frac{p_i}{\beta} \frac{\beta(0) \max_{k \in \mathcal{K}} \{5; 2B_1 = B_2 g\}}{2B_1};$$

and suppose the sample size is large enough to satisfy $m \geq \frac{2}{\alpha} \ln(\frac{1}{\alpha})$ and

$$\frac{2(R_t + 2c)k p_{k_1}}{m^{1-\alpha}} \leq \frac{1}{2} \frac{p}{2V(H) \ln(m+1) + \beta} + \frac{p}{2 \ln(\frac{1}{\alpha})} \leq \frac{c}{8}$$

(which happens for all large m since $\lim_{m \rightarrow \infty} \frac{p}{\ln(m)} = \lim_{m \rightarrow \infty} \frac{p}{\ln(m) m^{1-\alpha}} = 0$). Then it holds that the above values c, B_1 , and R_i (for $0 < i \leq t$) are all positive, and moreover the following statements hold simultaneously with probability at least $1 - \epsilon$.

1. The natural coefficient vector $\beta \in \mathcal{S}$ satisfies

$$L(A_t) \leq \inf_{k \in \mathcal{K}} L(A_t) + m^{-\alpha} + R_t \frac{1}{m} \ln \frac{6}{\beta}$$

$$+ \frac{2 R_t}{m^{1-\alpha}} \frac{1}{2} \frac{p}{2V(H) \ln(m+1) + \beta} + \beta(R_t) \frac{p}{2 \ln(\frac{1}{\alpha})}$$

$$+ \beta(0) \frac{k p_{k_1}}{k p_{k_1} + m^{-\alpha}} \stackrel{3B_1 = 9B_2}{\leq} \frac{1}{4B_2 R_1} :$$

2. If $H \in \mathcal{F}_{ds}(X)$ (where X is the marginal of \mathcal{X} over X), and letting R_β denote the Bayes error rate, there exists $\beta: \mathcal{R} \rightarrow \mathcal{R}$ satisfying $\beta(A_t) \leq R_\beta \leq L(A_t) + L$ and $\beta(z) \leq 0$ as $z \leq 0$. (For instance, when $\beta = \log$, then $\beta(z) = z^2/2$.)

3. The returned coefficients $\hat{\beta}$ satisfy

$$R(\hat{\beta}) \leq R(A_t) + 4 \frac{1}{\beta(A_t)} \frac{(V(H) + 1) \ln(2em) + \ln(24\beta)}{m^{1-\alpha}}$$

$$+ 8 \frac{(V(H) + 1) \ln(2em) + \ln(24\beta)}{m^{1-\alpha}} :$$

This bound is inferior to the guarantee in the separable case; while it is still of the form

$1 = m^{1-c}$, the exponent $1-c$ is distribution-dependent. The source of weakness is the optimization guarantee (cf. Lemma 6.4.6), which is brute-forced and should be improvable.

6.4.1 Curvature

Recall that the dual optimum p satisfies $kA^T p \leq k_1$, which implies $\sum_{A_i < 0} p_i = 0$ for every ϵ (cf. Lemma 1.A.5). To see how this helps locate bad examples and produce curvature, note the rearrangement

$$\sum_{A_i < 0} p_i = \sum_{A_i > 0} p_i;$$

meaning p has been reweighted by p so that negative and positive margins are equal (in a sense, p renders every ϵ equivalent to random guessing). Since p is fairly well-behaved (it is within $[0; \epsilon]$ -a.e. (where ϵ is the Lipschitz constant for σ), and is fairly flat since $\|p\|_1 \leq \epsilon$), then some algebra allows the removal of p from the above display, which yields the statement: if A has many good margins, it also has many bad margins. This constrains the norms of solutions found by the algorithm, and generates curvature in the sense that progress in any direction quickly leads to L increasing.

Of course, p could have instead been directly constructed from the presence of noise, but then the results would not be applicable to cases where H itself is noiseless, but H is simply very weak. The following example emphasizes this role of noise, but also shows that the above development overlooked the effect of sampling.

Example 6.4.3. (Nightmare scenario #2.) Pick any X , (marginal) distribution μ over X , hypothesis class $H \subseteq \mathcal{F}_{ds}(X)$, and any $\epsilon > 0$. Define the conditional density $\Pr[y = +1 | x]$ to be 0.9 when $(H)_x \geq 0$, and 0.1 otherwise when $(H)_x < 0$. By this construction, H attains the Bayes error rate (which is 0.1), and every other H does at best this well. Any weighting with favorable convex risk $L(A)$ will necessarily have a small norm in consequence of the guaranteed 10% classification error.

Unfortunately, finite samples look slightly different. Suppose $X = \mathbb{R}^d$ and μ is absolutely continuous with respect to Lebesgue measure. With probability 1, a random sample of any size will contain no noise, and $\text{span}(H)$ has a perfect predictor H^* (over the sample); in particular, nothing inhibits the norms of solutions over \mathcal{D} .

In this example, the good predictor H^* is potentially very complex, as it is fitting noise.

The solution here will be to only control those predictors with small norms; note that this deviation inequality embeds the reweighted average margin expression from eq. (6.2.1).

Lemma 6.4.4. Let probability measure \mathbb{P} over $X = \{1, \dots, m\}$ with empirical counterpart \mathbb{P}_n , any hypothesis class $\mathcal{H} \subset \mathbb{R}^m$, reweighting \mathbb{P}_n with $k_A > 0$, and norm bound C be given. Then, with probability at least $1 - \epsilon$, \mathbb{P}_n -a.e., and

$$\sup_{\substack{h \in \mathcal{H} \\ \|h\| \leq C}} \sum_{i=1}^n (A_i) \mathbb{P}_n(h) \leq \frac{2Ck_A}{m^{1-\epsilon}} \mathbb{P} \left(\frac{2V(H) \ln(m+1)}{2\epsilon} + \frac{\mathbb{P} \left(\frac{1}{2\epsilon} \right)}{2\epsilon} \right) :$$

Armed with these tools, the structure of the nonseparable problem is as follows. Note that the term B_1 is the aforementioned curvature lower bound, and furthermore the facts $\mathbb{P}_{i=1}^m \alpha_i = 1$ and $\mathbb{P}_{i=1}^m \alpha_i^2 < 1$ mean that the step sizes exactly hit the constrained step size regime studied by Zhang and Yu (2005, Equation (4)).

Lemma 6.4.5. Suppose the setting and quantities in the preamble of Theorem 6.4.2; the following statements hold simultaneously with probability at least $1 - \epsilon$.

1. Every \mathbb{P}_n with $k_A > R_t + 4c$ and $\mathbb{P}_n(A_i) < 2^{-t}$ has $\mathbb{P}_{i=1}^m \alpha_i \log(\mathbb{P}_n(A_i)) \geq B_1$.
2. For every choice of step size $\alpha_i \leq R_i$ and

$$\alpha_i \leq \min \left(\frac{9k_A - r \mathbb{P}_n(A_{i-1}) k_A^2}{4 B_1^2}, \frac{\max\{5, 2B_2 = B_1 g(\mathbb{P}_n(A_{i-1}), \mathbb{P}_n(A_i))\}}{2B_1} \right) :$$

3. Let \mathbb{P}_n with $k_A > R_1 \frac{\mathbb{P}_n(A_{i-1})}{\mathbb{P}_n(A_i)} = R_{t-1}$ and $\alpha_i := \min_{i \in [t-1]} \mathbb{P}_n(A_i) - \mathbb{P}_n(A_t) = 0$ be arbitrary. For every choice of step size,

$$\alpha_i \leq \frac{(\mathbb{P}_n(A_{i-1}) - \mathbb{P}_n(A_t))}{2B_2(k_A + R_1 \frac{\mathbb{P}_n(A_{i-1})}{\mathbb{P}_n(A_i)})} \quad \text{and} \quad \sum_{i=1}^t \alpha_i \leq \frac{\mathbb{P}_n(A_{t-1})}{4B_2 R_1} :$$

6.4.2 Proof of Theorem 6.4.2

The convergence analysis due to Zhang and Yu (2005) can be adjusted to the present setting (where step and coordinate selection are decoupled), yielding the following guarantee. Note that Lemma 6.4.5 also allows the application of the analysis due to Bartlett and Traskin (2007) (again with decoupling modifications), however this leads to a rate of roughly $O(\frac{1}{\sqrt{\ln(m)}})$.

Lemma 6.4.6. Let $\mathcal{F} \in L_{lg} \setminus L_{2d}$ with Lipschitz gradient parameter B_2 , binary class H , time horizon t , and empirical probability measure b be given. Let $\epsilon > 0$ be arbitrary, and suppose there exists $c_3 > 0$ with $c_3 \leq k A^r \mathbb{E}(\|A_{i-1}\|_{k_1})$ for all $0 \leq i \leq t$. Then

$$\mathbb{E}(\|A_t\|) - \mathbb{E}(\|A_0\|) \leq \mathbb{E}(\|A_0\|) + \frac{k k_1}{k k_1 + \dots} \epsilon^{c_3/2} = (6B_2) \epsilon^{c_3/2} :$$

From here, there is little to do: the conditions for this rate are met with high probability thanks to Lemma 6.4.5, and the rest is standard uniform convergence.

Appendix 6.A The Family of Dense Classes $F_{ds}(\mu)$

As the goal of a consistency analysis is to show that the Bayes predictor is approximated arbitrarily well, necessarily the function class considered by a purportedly consistent algorithm must be very large.

As discussed in Section 6.1, one choice is the class $F_{ds}(\mu)$ of functions dense according to $L^1(\mu)$ in the family of bounded measurable functions. A partial survey of density assumptions in other work is as follows.

Breiman (2000, Definition 1) works with a similar definition: the relevant metric is $L^2(\mu)$, and the closure must contain $L^2(\mu)$, where μ is constrained to be continuous with respect to Lebesgue measure. By contrast, the metric for $F_{ds}(\mu)$ is $L^1(\mu)$, where μ is an arbitrary measure over the Borel σ -algebra, and the closure of the class must contain bounded measurable functions, which are a subspace of $L^1(\mu)$, which is contained within $L^2(\mu)$.

Proposition 6.1.5, which will be proved shortly, states that it suffices for $\text{span}(H)$ to contain boxes formed by half-open intervals. This result was stated by Breiman (1999, Proposition 1) with an abbreviated proof for his setting of Lebesgue-continuous measures, thus the present result can be taken as merely proving that result with slightly more generality and verbosity.

The closest assumption and family of results to those here were provided by Zhang (2004, Section 4); while an analog to Proposition 6.1.5 is not shown there, the proofs rely on a form of Lusin's Theorem, which is used in Proposition 6.1.5 as well; indeed, the proofs here owe their existence to those earlier ones by Zhang (2004, Section 4).

Another approach, suggested by Lugosi and Vayatis (2004, Theorem 1 and subsequent remarks), and later used by Bartlett and Traskin (2007, Condition 1) and Schapire and Freund (2012, eq. (12.11)), is to require the weaker condition that

$$\inf_{g \in H} \int \ell(g(x)) d\mu(x) = \inf_{f \text{ measurable from } X \text{ to } \mathbb{R}} \int \ell(f(x)) d\mu(x)$$

for a verification that this property is indeed weaker, see Lemma 6.A.1. Lugosi and Vayatis (2004, Lemma 1) show that this assumption is satisfied by classes whose convex hull contains indicators of all Borel sets, and thus Lemma 6.A.1 can be considered a simplification which succeeds to grant consistency with more computationally tractable classes (like decision lists and trees).

As discussed above, the essential property of $\mathcal{F}_{ds}(\cdot)$ is that it implies the weaker condition used by Lugosi and Vayatis (2004, Theorem 1 and subsequent remarks), which in turn is directly needed for the classification calibration methods in the consistency proof (cf. Theorem 6.2.2). The Lipschitz condition here is not crucial, and for instance can be removed by adjusting $\mathcal{F}_{ds}(\cdot)$ to require approximants to a function to carry nearly the same uniform bound.

Lemma 6.A.1. Let distribution μ over X , ℓ a Lipschitz convex loss, class $H \subseteq \mathcal{F}_{ds}(X)$, and nonnegative Lipschitz convex loss ℓ be given (with Lipschitz constant L). Then

$$\inf_{g \in H} \int \ell(g(x)) d\mu(x) = \inf_{f \text{ measurable from } X \text{ to } \mathbb{R}} \int \ell(f(x)) d\mu(x)$$

Proof. One direction is immediate, since H defines a family of measurable functions.

Going the other direction, first define, for any measurable f , a clamping

$$[f]_r(z) := \begin{cases} f(z) & \text{when } |f(z)| \leq r, \\ r & \text{otherwise} \end{cases}$$

For any $r > 0$, based on four cases for the structure of f , a clamping value r is defined as follows in order to satisfy, for any f and z , $|([f]_r(z) - f(z))| \leq r$.

If $\lim_{z \rightarrow 1} f(z) = \lim_{z \rightarrow -1} f(z) < 1$, then f is a constant function, and $r = 0$ suffices.

If $\lim_{z \rightarrow 1} f(z) = \lim_{z \rightarrow -1} f(z) = 1$, then f has compact level sets, and in particular an

r exists so that

$$f(z) := \inf_q \int (q + g|f(z) - q|)^r dz$$

It follows that $\int ([f]_r)^r dz = \int f^r dz$.

If $\lim_{z \rightarrow 1} \int f^r dz < 1$ and $\lim_{z \rightarrow 1} \int f^r dz = 1$, then set

$$r := \inf \int f^r dz : \int f^r dz = \int (q + g|f - q|)^r dz$$

Unlike the preceding two cases, clamping here can increase the value, but not by more than

If $\lim_{z \rightarrow 1} \int f^r dz = 1$ and $\lim_{z \rightarrow 1} \int f^r dz < 1$, then this case is handled by the preceding one by considering the reflection $z \rightarrow 1 - z$.

Consequently, let f_i, g_i be a minimizing sequence for the target in the above so that

$$\int (f_i - g_i)^2 dx + \int (f_i - g_i)^2 dx \leq 2^{-i} + \inf \int (f - g)^2 dx : f \text{ measurable from } X \text{ to } \mathbb{R}$$

Each f_i might not be bounded, so define $[f_i]_{r_i}$ where $r_i := 2^i$; by this choice,

$$\begin{aligned} \int (f_i - g_i)^2 dx &= \int ([f_i]_{r_i} - g_i)^2 dx \\ &\leq 2^{-i} + \int (f_i - [f_i]_{r_i})^2 dx \\ &\leq 2^{-i+1} + \inf \int (f - g)^2 dx : f \text{ measurable from } X \text{ to } \mathbb{R} \end{aligned}$$

Lastly, since $\text{span}(H)$ is dense in the $L^1(X)$ metric, let $h_i \in \text{span}(H)$ satisfy $\|h_i - g_i\|_1 \leq 2^{-i}$; since \int is Lipschitz with constant 2 , then

$$\begin{aligned} \int (h_i - g_i)^2 dx &\leq \int (h_i - [f_i]_{r_i})^2 dx + \int ([f_i]_{r_i} - g_i)^2 dx \\ &\leq \int (h_i - [f_i]_{r_i})^2 dx + 2 \int (f_i - [f_i]_{r_i})^2 dx \\ &\leq \int (h_i - [f_i]_{r_i})^2 dx + 2 \cdot 2^{-i+1} \\ &\leq (2 + 2)2^{-i} + \inf \int (f - g)^2 dx : f \text{ measurable from } X \text{ to } \mathbb{R} \end{aligned}$$

and the result follows. □

To close, the proof of Proposition 6.1.5, which avoids strong structural assumptions on the measure (for instance, a relationship to Lebesgue measure) via an invocation of Lusin's Theorem.

Proof of Proposition 6.1.5. Let $\epsilon > 0$ and bounded measurable g with $\|g\|_{\infty} := \sup_x |g(x)| \leq 2$ (0; 1) (when $\|g\|_{\infty} \leq 2$, then $g \in \text{span}(H)$ and the proof is complete). By Lusin's Theorem, there exists compactly-support continuous $h \in L^1(\mu)$ which satisfies $\|h - g\|_{\infty} \leq \epsilon$ (Folland, 1999, Theorem 7.10). Let C denote the compact support of h ; continuity over a compact subset of \mathbb{R}^d means uniform continuity, and therefore let $\delta > 0$ be sufficiently small that the bounding box of C may be partitioned into finitely many cubes of side length δ (products of half-open intervals of length δ) so that, for any x_1 and x_2 within a single cube, $|h(x_1) - h(x_2)| \leq \epsilon$. Now let f be a sum of indicators of these cubes, where each indicator is weighted by $h(x)$ with x being an arbitrary point in the corresponding cube. By construction and since $\text{span}(H)$ contains such cubes, $f \in \text{span}(H)$, and moreover $\|f - h\|_1 \leq \epsilon$ since μ is a probability measure, which provides

$$\|f - g\|_1 \leq \|f - h\|_1 + \|h - g\|_1 \leq \epsilon + \int_{[h \neq g]} |h - g| d\mu \leq \epsilon + 2\epsilon \mu([h \neq g]) \leq 3\epsilon$$

□

Appendix 6.B Loss Function Classes L_{lg} and L_{2d}

First, note that L_{lg} and L_{2d} contain a few useful things.

Lemma 6.B.1. $\sigma_{\log} \in L_{lg}$ with parameters $B_2 = 1/4$, $\alpha_1 = 1/2$, $\alpha_2 = 1$. $\sigma_{\text{russ}} \in L_{lg}$ with parameters $B_2 = \alpha_1 = \alpha_2 = 1$. Lastly, $\exp \in L_{2d}$ and $\sigma_{\log} \in L_{2d}$.

Proof. For the logistic loss σ_{\log} , note $0 \leq \sup_x \sigma_{\log}''(x) \leq 1/4$, thus the mean value theorem grants Lipschitz gradients with parameter $B_2 = 1/4$. σ_{\log} 's Lipschitz parameters are $\alpha_1 = 1/2$ and $\alpha_2 = 1$.

Since σ_{russ} is not twice differentiable, gradient slopes must be checked manually. To start,

note

$$\psi_{\text{russ}}(x) = \begin{cases} 0 & \text{when } x \leq -1; \\ x + 1 & \text{when } x \in (-1; 0); \\ 1 & \text{when } x \geq 0; \end{cases}$$

whereby $\psi_1 = \psi_2 = 1$. Within each line segment, the gradient slopes are 0, 1, and 0. By manually checking pairs $x < y$ in the first and second, first and third, and second and third intervals, the tightest Lipschitz constant on the gradients is 1.

The containments within L_{2d} are direct. □

Next, note that the Lipschitz gradient properties on ψ also cover \mathbb{B} .

Lemma 6.B.2. Let convex differentiable $\psi : \mathbb{R} \rightarrow \mathbb{R}$ with Lipschitz gradient parameter B_2 , any H , and any empirical measure μ over $X = \mathbb{R}^d$ be given. Then $\mathbb{B}(A)$ has Lipschitz gradients with parameter B_2 with respect to norm $L^1(\mu)$ over the t^{th} level set

$$S_t := \left\{ \mu : \mathbb{B}(A) \subseteq \mathbb{B}(A_{t-1}) \right\}$$

Proof. Starting from Lemma 1.C.2, for any $\mu \in S_t$, since $\psi'' \geq 0$ (whereby $\psi' \circ \mu = \psi'' \circ \mu$), and also using the definition of the level set S_t ,

$$\begin{aligned} r(\mathbb{B}(A))(\mu) - r(\mathbb{B}(A))(\mu^0) &\leq \frac{1}{m} \sum_{i=1}^n \left(\psi'(A_{x_i, y_i}) - \psi'(A_{x_i^0, y_i^0}) \right) \\ &\leq \frac{1}{m} \sum_{i=1}^n B_2 |j(A_{x_i, y_i}) - j(A_{x_i^0, y_i^0})| \\ &\leq \frac{1}{m} \sum_{i=1}^n B_2 \sum_{h \in H} |h(x_i) - h(x_i^0)| \\ &\leq B_2 k \sum_{i=1}^n \mu^0(x_i^0) \end{aligned}$$

as desired. □

As a final basic result about ψ , note that the terminology "separable" is at least somewhat justified.

Proposition 6.B.3 (See also Theorem 5.4.2) Suppose $\psi : \mathbb{R} \rightarrow \mathbb{R}_{++}$ is convex with $\lim_{z \rightarrow 1} \psi'(z) = 0$, and let any H and any probability measure μ over $X = \mathbb{R}^d$ be given. Suppose $\inf_{\mathbb{R}^d} \mathbb{B}(A) d =$

0.

1. For any $\epsilon > 0$, there exists n so that $(\|A_n - A\|) < \epsilon$.
2. With probability 1 over the draw of a sample $(x_i; y_i)_{i=1}^m$ (for any $m < \infty$), there exists n so that $(A_n)_{x_i; y_i} > 1$ for every i .
3. In general, there does not exist n so that $(\|A_n - A\|) = 0$ (indeed, Example 6.3.4 provides a counterexample).

Proof. Let $\epsilon > 0$ be given, and choose $\delta_{i=1}^m$ so that $\sum_{i=1}^m \delta_i = \epsilon$. Since $(A_i) \geq 0$ -a.e., by Egorov's theorem there exists S with $(S) > 1 - \delta$ so that $(A_i) \geq 0$ uniformly on S (Folland, 1999, Theorem 2.33). But since $\epsilon > 0$ everywhere and $\lim_{z \rightarrow 1} \epsilon(z) = 0$ and ϵ is convex, it must be the case that $(A_i) \geq 1 - \delta$ uniformly on S , and so there exists n with $(A_n)_{x_i; y_i} > 1 - \delta$ on S , which gives the first result.

For the second result, take any $\epsilon > 0$, and choose δ as granted by the first part. Let b denote the empirical measure over the provided sample; then

$$\Pr[\exists i (A_n)_{x_i; y_i} < 1 - \delta] = (\|A_n - A\|)^m (1 - \delta)^m$$

Since $\epsilon > 0$ was arbitrary, the second result follows.

For the third result, recall that Example 6.3.4 (whose properties are provided in Proposition 6.3.5) gave an instance where every element of $\text{span}\{h\}$ makes some mistakes. □

Appendix 6.C Duality Properties of

In order to develop (\cdot) , the set $D(\cdot)$ must first be studied.

Proposition 6.C.1. (Basic properties of $D(\cdot)$.) Let μ be an arbitrary probability measure over $X = \{f \in \mathbb{R}^d; \|f\|_1 \leq 1\}$, and let $\alpha \in [0; 1]$ be arbitrary. The set $D(\cdot)$ has the following properties.

1. $D(\cdot)$ is convex.
2. $D(\cdot)$ is closed in the $L^1(\cdot)$ topology.
3. If $\alpha > 0$, then $D(\cdot)$ is closed in the $L^1(\cdot)$ topology, and also closed in the weak* topology (i.e., the weak topology induced upon $L^1(\cdot)$ by $L^1(\cdot)$).

4. $D(\cdot)$ is compact in the weak* topology or $L^1(\cdot)$ (as discussed in the preceding point).
5. $D(\cdot)$ is not guaranteed to be compact in the $L^1(\cdot)$ or $L^1(\cdot)$ topologies; indeed, it is not compact when $\Omega = \{1, 2\}$, $X = [0; 1]$, the marginal distribution μ^X is uniform on $[0; 1]$, and the conditional distribution $\Pr(Y = 1 | X = x)$ is arbitrary.

Proof. 1. For convexity, let any $\mu \in D(\cdot)$ and $p_1, p_2 \in D(\cdot)$ be given, and define sets $N_j := \mu_j^{-1}(\{1\})$ for $j \in \{1, 2\}$, where necessarily $\mu(N_j) = 0$. The goal is to show $\mu := (1-\alpha)\mu_1 + \alpha\mu_2 \in D(\cdot)$.

Define $N := N_1 \cup N_2$ (where again $\mu(N) = 0$). First, for any $(x, y) \in N^c$,

$$\mu(x, y) = (1-\alpha)\mu_1(x, y) + \alpha\mu_2(x, y) = 0 + \alpha \cdot 0 = 0;$$

whereby it follows that $\mu \geq 0$ -a.e.. Second,

$$\int \mu(x, y) dx dy = (1-\alpha) \int \mu_1(x, y) dx dy + \alpha \int \mu_2(x, y) dx dy = 1;$$

again using the convention $\int \mu = 0$, whereby $\int \mu(x, y) dx dy = 1$ as desired. Lastly,

$$\int \mu(x, y) dx dy = \int_{N^c} ((1-\alpha)\mu_1 + \alpha\mu_2) dx dy = (1-\alpha) \int \mu_1 + \alpha \int \mu_2 = 1;$$

meaning all conditions are met, and $\mu \in D(\cdot)$. Since μ_1, μ_2 were arbitrary, it follows that $D(\cdot)$ is convex.

2. For closure within $L^1(\cdot)$, since $L^1(\cdot)$ is a metric space, it is first countable, and thus it suffices to check that any sequence $\mu_j \in D(\cdot)$ with $\mu_j \rightarrow \mu$ in $L^1(\cdot)$ satisfies $\mu \in D(\cdot)$ (Folland, 1999, Proposition 4.6). Given any such sequence $\mu_j \in D(\cdot)$, choose a subsequence μ_{j_k} so that $\mu_{j_k} \rightarrow \mu$ -a.e. (Folland, 1999, Corollary 2.32).

Let N_p be the (null) set of points for which convergence fails, and additionally, for each i , define $N_i := \mu_{j_k}^{-1}(\{1\})$; lastly, set $N := N_p \cup (\cup_i N_i)$, where again $\mu(N) = 0$. Thus for

any $(x; y) \in N^c$,

$$\begin{aligned} p(x; y) &= q(x; y) + (p(x; y) - q(x; y)) \\ &= q(x; y) + j |p(x; y) - q(x; y)| \\ \liminf_{i \rightarrow \infty} q(x; y) + j |p(x; y) - q(x; y)| \\ &= 0; \end{aligned}$$

thus $p \geq 0$ a.e. Additionally,

$$k p k_1 = \int_N |p| = \int_N p = \int_N p_i + \int_N (p_i - p) = k p_i k_1 + \int_N (p_i - p);$$

whereby

$$k p k_1 - k p_i k_1 = \int_N (p_i - p) \leq k p_i - p k_1 \leq 0;$$

and $k p k_1 = 1$ as desired.

For the last property, if $\epsilon = 0$, there is nothing to show, thus suppose $\epsilon \in (0; 1]$, set $P_i := q^{-1}((1-\epsilon; 1])$, and $Z := N_p \setminus (\cup_i P_i)$, whereby it follows that

$$(Z) = 0; \quad q \leq p \text{ over } Z^c; \quad q \geq 1-\epsilon \text{ over } Z^c;$$

Then, for any $(x; y) \in Z^c$,

$$p(x; y) = q(x; y) + (p(x; y) - q(x; y)) \leq \limsup_{i \rightarrow \infty} q(x; y) + j |p(x; y) - q(x; y)| \leq 1-\epsilon;$$

which establishes $k p k_1 = 1-\epsilon$, and thus $p \in D(\epsilon)$.

- Note firstly that if $\epsilon = 0$, then $D(\epsilon)$ can contain members which are not elements of $L^1(\mu)$, and thus discussing this set in the $L^1(\mu)$ topology does not make sense. For the remainder of this case, suppose $\epsilon > 0$.

Just as in the case of $L^1(\mu)$, for $L^1(\mu)$ it suffices to let a sequence $f_i, g_{i=1}^1 \in D(\epsilon)$ be given with $p_i \rightarrow p$ in the $L^1(\mu)$ topology, and to show that $p \in D(\epsilon)$. Notice however, since μ is

a probability measure, that

$$\|p_i - p_j\|_1 = \sum_k |p_{ik} - p_{jk}| = \sum_k |p_{ik} - p_{ik} + p_{ik} - p_{jk}| \leq \sum_k |p_{ik} - p_{jk}| = \|p_i - p_j\|_1;$$

meaning $p_i \in L^1(\Omega)$ as well, which by the preceding case provides that $D(\epsilon)$ as desired.

Lastly, since $D(\epsilon)$ is convex and additionally closed according to $L^1(\Omega)$, then it is also weak* closed (Rudin, 1973, Theorem 3.12).

4. Again suppose $\epsilon > 0$, and define

$$B := \{p \in L^1(\Omega) : \|p - g\|_1 \leq \epsilon\}$$

By Alaoglu's Theorem (Folland, 1999, Theorem 5.18), B_1 is compact in the weak* topology, thus $B = \epsilon B_1$ is weak*-compact as well. The result follows since $D(\epsilon)$ is a weak*-closed subset of B , and closed subsets of compact sets are compact (Folland, 1999, Theorem 4.22).

5. Noncompactness can be understood from the fact that norm balls are in general not compact, but an explicit construction is provided for completeness. Since both $L^1(\Omega)$ and $L^\infty(\Omega)$ are metric spaces, to prove non-compactness, it suffices to prove $D_{1=2}(\epsilon)$ is not totally bounded. In particular, a countably infinite subset of $C = D_{1=2}(\epsilon)$ will be constructed satisfying the property $(f, g) \in C \times C$ with $f \neq g$ implies $\|f - g\|_1 = 1 = 2\epsilon$ and $\|f - g\|_\infty = 2\epsilon$, which suffices to show that C (and thus $D_{1=2}(\epsilon)$) is not totally bounded (in either metric) for the following reason. Let S be any finite subset of $L^1(\Omega)$ or $L^\infty(\Omega)$. Since C and S have respectively infinite and finite cardinalities, there must exist $h \in S$ which is a closest element in S to two distinct functions $f \neq g$ in C . Let $\|\cdot\|$ denote either norm under consideration, and note that

$$1 = 2\epsilon \leq \|f - g\| \leq \|f - h\| + \|h - g\| \leq 2 \max\{\|f - h\|, \|h - g\|\};$$

which means that one of these two distances is at least ϵ . Since S was an arbitrary finite set, it follows that there is no finite set of balls of radius $\epsilon/2$ which covers C , and thus C and $D_{1=2}(\epsilon)$ are not totally bounded according to either norm.

The construction is as follows. For every positive integer $i \in \mathbb{Z}_{++}$, define the function

$$f_i(x; y) := 2^{-i} \sum_{j=0}^{2^i - 1} x^{2j} [(2j)2^{-i-1}; (2j+1)2^{-i-1}] :$$

Define $C := \{f_i : i \in \mathbb{Z}_{++}\}$. By construction, $C \subseteq D_{1=2}(\cdot)$ (i.e., $\|f_i\|_1 = 1$ and $\|f_i\|_2 = 2$), and moreover $i \neq j$ implies f_i and f_j disagree on exactly half of their support, which yields $\|f_i - f_j\|_1 = 1 = 2^{-1}$ and $\|f_i - f_j\|_2 = \|f_i\|_2 = 2$.

□

With the structure of $D(\cdot)$ established, the basic duality structure of (\cdot) follows. Note that the value of establishing the weak*-compactness of $D(\cdot)$ is to grant an application of Sion's minimax Theorem without making any topological assumptions on H (or rather, on the subspace H). Additionally, Lemma 6.3.7 in Section 6.3 is a combination of this result and part of Lemma 1.A.5.

Lemma 6.C.2. Let probability measure μ over $X = \{-1; +1\}$, any H , and any $\alpha \in [0; 1]$ be given. Then

$$\begin{aligned} \alpha &= \min_{p \in D(\cdot)} \sup_{g \in Z} \int (A) p d\mu : 2 \\ &= \sup_{p \in D(\cdot)} \min_{g \in Z} \int (A) p d\mu : 2 \\ &= \min_{p \in D(\cdot)} \|kA\|_1 ; \end{aligned}$$

where $\|kA\|_1$ is discussed in Lemma 1.A.5.

Proof of Lemma 6.C.2. Before applying the duality result, it must be established that the various infima are attained. To start, consider the dual expression $\inf_{p \in D(\cdot)} \|kA\|_1$, and let $\{p_i\}_{i=1}^\infty$ with $p_i \in D(\cdot)$ be a minimizing sequence to the infimum. Since Proposition 6.C.1 establishes that $D(\cdot)$ is weak*-compact, there is a subsequence $\{p_{i_j}\}_{j=1}^\infty$ which weak*-converges to some $p \in D(\cdot)$ (Folland, 1999, Theorem 4.29). But Lemma 1.A.5 established that $\|kA\|_1$ is weak* lower semi-continuous, and since it is finite over $L^1(\cdot)$, it is therefore weak* continuous, and therefore the limit point p attains the infimum. Furthermore, Lemma 1.A.5 provides that $\|kA\|_1$ is the same as the first infimum, whereby both expressions attain their minimizers and are equal.

The middle expression is the easiest; once again constructing a weak*-convergent sequence $q_n \rightarrow q$ with $q_n \in D(\cdot)$, the definition of weak*-convergence explicitly grants $\int \phi f d q_n \rightarrow \int \phi f d q$ for every $f \in L^1(\cdot)$, and since $A \in L^1(\cdot)$ is held fixed within this inner expression, it follows that q attains the infimum.

What remains is to swap minimization and maximization. This in turn follows by Sion's minimax theorem (Komiya, 1988); to verify this application, note that $(p; \cdot) \rightarrow \int \phi d p$ is linear and continuous in both parameters (indeed, this is by construction, since $L^1(\cdot)$ is isometrically isomorphic to the topological dual to $L^\infty(\cdot)$, and the weak* topology over $L^\infty(\cdot)$ ensures that this integral relation is continuous for every $\phi \in L^1(\cdot)$), also that $D(\cdot)$ is a topological vector space, and lastly that $D(\cdot)$ is a convex compact subset of a topological vector space (namely, the weak* topology, and not the $L^1(\cdot)$ topology, where $D(\cdot)$ is not necessarily compact as per Proposition 6.C.1). \square

Appendix 6.D Reweighted Margin Deviations (with k Fixed)

Lemma 6.D.1. Let probability measure μ over X with μ and empirical counterpart b , any hypothesis class $H \subseteq F_{vc}$, reweighting $p \in L^1(\cdot)$, and norm bound C be given. Then, with probability at least $1 - \epsilon$, $p \in [0; k p_{k_1}]$ b-a.e., and

$$\sup_{k \leq k_1} \frac{1}{m} \sum_{i=1}^m (A_i) p d b - \frac{1}{m} \sum_{i=1}^m (A_i) p d \mu \leq \frac{2Ck p_{k_1}}{m^{1/2}} \sqrt{2V(H) \ln(m+1)} + \frac{p}{2 \ln(1/\epsilon)} :$$

Proof of Lemma 6.D.1. First, define a simplified reweighting $p^0(x; y) := p(x; y) \mathbb{1}[|p(x; y)| \leq k p_{k_1}]$; by the definition of $k p_{k_1}$, then $p^0 = p$ -a.e., and thus, with probability 1, any finite sample of any size has p^0 and p agreeing. The proof will work with p^0 , which satisfies $\sup_{x,y} |p^0(x; y)| \leq k p_{k_1}$, and then close by discarding a measure zero set and thus relating p to p^0 .

The main part of the proof is an almost standard application of Rademacher complexity techniques for voted classifiers (Boucheron et al., 2005a, Theorem 4.1 and its proof, which controls for a surrogate loss and not just the classification loss); the only modification will be to work with a loss function which is sensitive to each example in the sample $\mathcal{L} = \sum_{i=1}^m f(x_i; y_i) g_{i-1}^m$, which will require a slightly refined Lipschitz contraction principle for Rademacher complexities (Shalev-Shwartz, 2009, Section 22.2, Lemma 15).

Specifically, define the loss

$$((A)_{x,y}) := \begin{cases} Cp^0(x; y) & \text{when } |j(A)_{x,y}| \leq C; \\ (A)_{x,y} p^0(x; y) & \text{when } |j(A)_{x,y}| > C; \\ + Cp^0(x; y) & \text{when } |j(A)_{x,y}| > C; \end{cases}$$

Since $\sup_{h,x} |jh(x)| \leq 1$ and $k \leq k_1 \leq C$, it follows that $|j(A)_{x,y}| \leq C$, and thus the extremal cases are never encountered, meaning

$$((A)_{x,y}) = (A)_{x,y} p^0(x; y);$$

and by construction ϕ is Lipschitz with parameter kpk_1 (as a function of (A)) and (A) has uniform bound Cpk_1 .

As such, letting R denote Rademacher complexity, by the Lipschitz contraction principle for per-coordinate losses (Shalev-Shwartz, 2009, Section 22.2, Lemma 15), behavior of Rademacher complexity on convex hulls (Boucheron et al., 2005a, Theorem 3.3), and relationship between Rademacher complexity and VC dimension (Boucheron et al., 2005a, See the display after eq. (7)),

$$R((A)) \leq kpk_1 R((A)) \leq kpk_1 CR(H) \leq kpk_1 C \sqrt{\frac{2V(H) \ln(m+1)}{m}};$$

This handling of a per-coordinate Lipschitz loss may be inserted into a standard deviation bound for uniformly bounded Lipschitz losses (Boucheron et al., 2005a, Theorem 4.1 and its proof) | albeit with an extra factor two to control deviations in both directions | and it follows, with probability at least $1 - \epsilon$, that

$$\sup_{k \leq k_1 \leq C} \int (A) p^0 db - \int (A) p^0 d \leq \frac{2Ckpk_1}{m^{1/2}} \sqrt{2 \int \frac{p}{2V(H) \ln(m+1)} + \int \frac{p}{2 \ln(1/\epsilon)}};$$

To complete the proof, recall that $p^0 = p$ -a.e., and a measure zero event was discarded, whereby $p^0 = p$ b-a.e. as well. □

Appendix 6.E Deferred Material from Section 6.3

6.E.1 Deviations of \hat{p}_m ()

This subsection establishes the following one-sided deviation bound on \hat{p}_m ().

Lemma 6.E.1. Let any H , any $\alpha \in (0; 1]$, any confidence parameter $\beta \in (0; 1]$, and any probability measure \mathbb{P} with empirical counterpart \hat{p}_m be given. Then with probability at least $1 - \beta$,

$$(b) \quad \hat{p}_m \leq \frac{1}{\alpha} + \frac{1}{2m} \ln \frac{2}{\beta} :$$

The difficulty in the analysis is that the definition of \hat{p}_m () involves an infimum over $\mathcal{P}^2(D)$ () and a supremum over \mathcal{P}^2 with $k_1 = 1$. The proof strategy employed here is to consider a single good choice for \mathcal{P}^2 , and to consider the effect on deviations as α varies. These deviations do not appear to be amenable to the usual approach, as \mathcal{P}^2 () is massive: it is in general not compact in the relevant metric topologies (cf. Proposition 6.C.1), and does not obviously possess other structure granting a uniform convergence result. The approach here is to instead identify that the dual optimum has very simple structure, and moreover this structure is robust to sampling.

Considering again the definition of \hat{p}_m (), while it is true that $\mathcal{P}^2(D)$ is defined over a potentially massive space, when placed in the expression $\sup_{\mathcal{P}^2} \int p d\mu$, all that matters is the behavior of p for each value of A , which ranges over $[-1; +1]$. That is to say, p is really reweighting the univariate margin distribution of A , and the best it can do is emphasize bad margins. In particular, the following lemma proves basic properties of an idealized univariate distillation of this scenario.

Lemma 6.E.2. Let a probability measure μ supported on $[-1; +1]$ and some $\alpha \in (0; 1]$ be given. Correspondingly define

$$\begin{aligned} S &:= \{c \in [-1; +1] : \mu((-\infty; c]) \geq \alpha\} \\ c &:= \sup S ; \\ l &:= (-\infty; c) ; \\ p(r) &:= \frac{1}{\alpha} 1_{[r \in l]} + \frac{(\mu(l))}{(\alpha - \mu(l))} 1_{[r = c]} ; \end{aligned}$$

with the convention $0=0 = 0$ in the definition of p . These objects have the following properties.

1. S is the closed interval $[1; c]$.
2. (I) , and $(I [f c g])$.
3. $k p k_1 = 1$ and $k p k_1 = 1 =$.
4. The optimization problem

$$\inf_{\substack{Z \\ p \in L^1(\cdot); k p k_1 = 1; p \in [0; 1] \text{ -a.e.}}} \int p(r) d(r)$$

is minimized at p .

Proof. First note that $1 \in S$, since p is supported on $[1; +1]$ and thus $([1; 1]) = 0$.

Next, S is an interval, since if $1 < c_1 < c_2$ and $c_2 \in S$, then $([1; c_1]) < ([1; c_2])$ and thus $c_1 \in S$.

To show that S is indeed a closed interval, consider any increasing sequence c_i with $c_i \in S$, thus $c_i \leq c$ for some $c \in [1; +1]$ since $[1; +1]$ is compact and the sequence is increasing. Then

$$([1; c]) = \bigcup_{i=1}^{\infty} ([1; c_i])$$

and thus, by continuity of measures Folland (1999, Theorem 1.8),

$$([1; c]) = \left(\bigcup_{i=1}^{\infty} ([1; c_i]) \right) = \liminf_{i \rightarrow \infty} ([1; c_i]) = \limsup_{i \rightarrow \infty} ([1; c_i]) ;$$

meaning $c \in S$ and S is closed.

Since S is a closed interval, then $c = \sup S \in S$, and it follows by the preceding properties that $S = [1; c]$.

By definition, for every $c \in S$, it holds that $([1; c]) > 0$, thus $c \in S$ implies that $(I) > 0$.

Next, for every positive integer $i \in \mathbb{Z}_{++}$, it holds by definition of c that $c + 1/i \notin S$, and thus, again by continuity of measures Folland (1999, Theorem 1.8),

$$([1; c]) = \left(\bigcap_{i=1}^{\infty} ([1; c + 1/i]) \right) = \liminf_{i \rightarrow \infty} ([1; c + 1/i]) ;$$

For the norms of p (which is a simple function over the Borel σ -algebra), notice that

$$\|p\|_1 = \int (|p|) + \frac{\int (|p|)}{\int (f \circ g)} \int (f \circ g) = 1:$$

Moreover, $p = 1$ on $(-\infty; c)$, and $p = 0$ on $(c; \infty)$; to show $\|p\|_1 = 1$, the behavior of p on c is all that needs to be checked. Since $(-\infty; c] \cup \{c\} = \mathbb{R}$, then

$$\int (|p|) = \int_{(-\infty; c]} (|p|) + \int_{\{c\}} (|p|) = \int (f \circ g);$$

so $\int (|p|) = 1$. Additionally $\int (|p|) = 1$ implies $\int p = 0$, and thus $\|p\|_1 = 1$ as desired.

Lastly, for the minimization problem, consider any feasible p (meaning $\|p\|_1 = 1$ and $\int p = 0$) with $\int p > \int p^0$. But since p is as large as possible along \mathbb{R} , it follows that $p < p^0$ for a positive measure subset of \mathbb{R} , and $p > p^0$ for a positive measure subset of $(c; \infty)$. Consequently $\int p < \int p^0$. Since p^0 was arbitrary, it follows that p is a minimal choice. □

The task now is to map the optimization over $D(\mathbb{R})$ down to this idealized univariate search problem. Temporarily adopting notation from probability theory, a first step in this direction would be to write

$$\int_{\mathbb{R}} (A) p d = E((A) p) = E(E((A) p | (A) = r));$$

where the latter notation signifies a conditional expectation with respect to the σ -algebra generated by events such that (A) falls in some Borel subset of \mathbb{R} (recall that all σ -algebras here are Borel). In some circumstances, the function $E((A) p | (A) = r)$ can be converted into integration over a function that takes r as input, which would directly allow conversion to the above univariate idealization; these techniques generally require assumptions on f and g which would rather be avoided here (Durrett, 2010, Section 5.1.3, regular conditional probabilities). As such, the following result exhibits the desired correspondence manually, albeit keeping the above idea in mind.

Lemma 6.E.3. Let any $\alpha \in (0; 1]$, any probability measure μ over $X = [-1; +1]$, any H , and any $\nu \in \mathcal{P}(\mathbb{R})$ with $\int \nu = 1$ be given. Define a probability measure ρ over \mathbb{R} as the pushforward of

through A , meaning, for any Borel subset S of R ,

$$(S) := ((A)^{-1}(S)):$$

Then μ is supported on $[-1; +1]$, and moreover the function $p(x; y) := p((A)^{-1}_{x,y})$, where p is as defined in Lemma 6.E.2, is a (feasible) minimizer to the optimization problem

$$\inf_{\mu \in \mathcal{D}(\mu)} \int (A)^{-1} p d\mu : p \in \mathcal{D}(\mu) :$$

Proof. Since $\|k\|_1 = 1$ and A is a continuous linear operator with unit norm (cf. Lemma 1.A.1, or recall the definition of A and the property $\sup_{x,h} |h(x)| \leq 1$), then $|j(A)^{-1}_{x,y}| \leq 1$, and thus $(x; y) \in \mathcal{D}((A)^{-1}_{x,y})$ maps $X \in [-1; +1]$ to $[-1; +1]$, and so the corresponding pushforward measure μ is supported on $[-1; +1]$. Therefore Lemma 6.E.2 provides the structure of $p \in \mathcal{D}(\mu)$ attaining the minimum in

$$\inf_{\mu \in \mathcal{D}(\mu)} \int p(r) d\mu(r) : p \in \mathcal{D}(\mu); \|k\|_1 = 1; p \in [0; 1] \text{ -a.e. } \mu :$$

Setting $p := p((A)^{-1})$ as in the statement, by the above optimality guarantee and by properties of pushforward measures (Resnick, 1999, Theorem 5.5.1),

$$\begin{aligned} \inf_{\mu \in \mathcal{D}(\mu)} \int p(r) d\mu(r) &: p \in \mathcal{D}(\mu); \|k\|_1 = 1; p \in [0; 1] \text{ -a.e. } \mu & (6.E.4) \\ &= \int p(r) d\mu(r) \\ &= \int (A)^{-1} p d\mu \\ &= \inf_{\mu \in \mathcal{D}(\mu)} \int (A)^{-1} p d\mu : \end{aligned}$$

Now let $\epsilon > 0$ and $\mu \in \mathcal{D}(\mu)$ be arbitrary. A corresponding element q with $\|k\|_1 = \|k\|_1$ and $\|k\|_1 \leq \|k\|_1 + \epsilon$ will be constructed as follows in order to satisfy

$$\int (A)^{-1} p d\mu - \int q(r) d\mu(r) \leq \epsilon :$$

Cover $[-1; +1]$ with at most $1/\epsilon$ disjoint half-open intervals I_i of the form $[-1 +$

$i \geq 0$; $1 + (i + 1) \leq \dots$ where i is a nonnegative integer and $0 := (1 + d_1 = e)$. De ne

$$q(r) := \prod_{i=1}^k (l_i)^{-1} [r \geq l_i] \quad \text{pd};$$

with the convention $0=0 = 0$ (i.e., $q(l_i) = 0$ when $(l_i) = 0$). By this choice, $kqk_1 = kpk_1$, and

$$kqk_1 = \int q(r) d(r) = \prod_{i=1}^k \int [r \geq l_i] \text{pd} = kpk_1:$$

More importantly, using Fubini's Theorem to interchange the integrals over \dots and \dots ,

$$\begin{aligned} \int (A) \text{pd} \int r q(r) d(r) &= \prod_{i=1}^k \int [r \geq l_i] \text{pd} \int r (l_i)^{-1} \text{pd} d(r) \\ &= \prod_{i=1}^k \int [r \geq l_i] \text{pd} \int r (l_i)^{-1} d(r) d \\ &= \prod_{i=1}^k \int [r \geq l_i] \text{pd} \int r (l_i)^{-1} d(r) d : \end{aligned}$$

Since \dots and p were arbitrary,

$$\begin{aligned} \inf \int (A) \text{pd} : p \in D(\dots) \\ \inf \int r p(r) d(r) : p \in L^1(\dots); kpk_1 = 1; p \in [0; 1] \text{ -a.e.}; \end{aligned}$$

which combined with the inequalities starting with eq. (6.E.4) provides that p is indeed a minimizer. □

With these tools in place, the proof of Lemma 6.E.1 follows.

Proof of Lemma 6.E.1. Consider the form of (\dots) provided by Lemma 6.C.2, whereby the supremum over \dots is on the outside. Let $\epsilon > 0$ be arbitrary, choose \dots which is within $\epsilon > 0$ of achieving the supremum, and let p be an optimal dual element as provided by Lemma 6.E.3, together meaning

$$(\dots) + \inf_{p \in D(\dots)} \int (A) \text{pd} = \int (A) p d : \tag{6.E.5}$$

Now consider the behavior of p over b . By construction, $kpk_1(b) = 1$, however

$k p_{k_{L^1(b)}}$ is a random variable; but by Hoeffding's inequality, with probability at least $1 - \frac{1}{2m}$,

$$k p_{k_{L^1(b)}} - k p_{k_{L^1(\cdot)}} \leq \frac{1}{2m} \ln \frac{2}{\epsilon};$$

henceforth discard this failure event.

Next instantiate another dual optimum p^b via Lemma 6.E.3, but now over the empirical measure; since λ is primal feasible in the definition of (b), and again using the form from Lemma 6.C.2 with the supremum on the outside, it follows that

$$(b) \quad \int (A) p^b db \tag{6.E.6}$$

Now recall the exact form of p and p^b as provided by Lemma 6.E.3 (and more specifically Lemma 6.E.2), which are both exactly ϵ up to some point, within $[0; 1-\epsilon]$ at that point (potentially distinct for p^b and p), and zero thereafter; if $k p_{k_{L^1(b)}} \geq 1$, then

$$\int p^b db = \int p db - \int p db \leq 1;$$

whereas $k p_{k_{L^1(b)}} < 1$ implies

$$\int p^b db = \int p db = 1$$

In either case, using as usual the fact $\sup_{x,y} |j(A)_{x,y} - j(A)_{x,y}| \leq k_1$, and additionally the controls on $k p_{k_{L^1(b)}}$ from above,

$$\begin{aligned} \int (A) p^b db - \int (A) p db &= \int p^b db - \int p db \\ &\leq \frac{1}{2m} \ln \frac{2}{\epsilon}; \end{aligned}$$

Combining this with eqs. (6.E.5) and (6.E.6),

$$(b) \int (A) p^b db \leq \int (A) p db - \frac{1}{2m} \ln \frac{2}{s} :$$

Since $s > 0$ was arbitrary, the result follows. □

6.E.2 Other Results

Lemma 6.E.7. Let μ be a probability measure on $X = \{1, +1\}^g$. If $0 \leq \mu_1 \leq \mu_2 \leq 1$, then $0 \leq \mu_1(\cdot) \leq \mu_2(\cdot) \leq 1$.

Proof. Let $0 \leq \mu_1 \leq \mu_2 \leq 1$ be given; then $D_1(\cdot) \leq D_2(\cdot)$ by definition, and thus $\mu_1(\cdot) \leq \mu_2(\cdot)$. Next, $\mu_1(\cdot) \leq 0$ follows by Lemma 6.C.2 since $kA^g \geq pk_1 = 0$, or by considering the effect of the primal player choosing $\mu = 0 \leq 2$. For the upper bound, since $\sup_{h,x} |h(x)| \leq 1$, then $\mu_2(\cdot) \leq \sup_{k_1} \{k_1 : \mu_2 \leq k_1\} \leq 1$: □

Proof of Proposition 6.3.3. Since every $p \in L^1(\cdot)$ with $\int p k_1 = 1$ and $p \geq 0$ -a.e. defines a probability measure $p d\mu$ (ignoring a μ -null set which does not affect that value of integration with respect to μ), and since $\int e_h k_1 = 1$ for every $h \in H$,

$$\inf_{\mu} \sup_{h \in H} \int h(x) d\mu(x; y) : \mu \text{ is a Borel probability measure over } X = \{1, +1\}^g$$

$$= \inf_{\mu} \sup_{k_1 \geq 1} \int (A)_{x,y} p(x; y) d\mu(x; y) : p \in L^1(\cdot); \int p k_1 = 1; p \geq 0 \text{ -a.e.}$$

$$= \mu_0(\cdot):$$

For $\mu_0(\cdot) = \mu_0(\cdot)$ with μ_0 a discrete measure over a finite set, the proof is as above (indeed with a tiny refinement, since in this case both $\mu_0(\cdot)$ and $\mu_0(\cdot)$ consider the same set of weightings over X). □

Proof of Proposition 6.3.5. For convenience, define $i := (1, \dots, i+1, i]$. This proof will proceed by establishing, for every k , a bounded weighting p_k , which will establish an upper bound on $\mu_0(\cdot)$

for some $\epsilon > 0$ which is a function of k . The result will then follow for $\phi_0(\cdot)$ by the monotonicity of $\phi(\cdot)$ as a function of ϵ (cf. Lemma 6.E.7), and result for $\phi_0(b)$ will use deviation bounds on $\phi(\cdot)$ and again the monotonicity property.

Define p_k to be positive over intervals I_i with $1 \leq i \leq 2k$, and zero elsewhere as follows. For any $x \in I_i$, $p_k(x; y) = i(i+1)/(2k)$. By this choice,

$$\int_{I_i} p(x; y) d(x; y) = \frac{i(i+1)}{2k} \frac{1}{i} \frac{1}{i+1} = \frac{1}{2k};$$

It follows that $\int p_k = 1$, $\int p_k^2 = 2k+1$, and p_k makes ϕ look like the uniform distribution over k consecutive intervals.

Now consider any hypothesis H , with some threshold r . If r lies outside this set of intervals, then H is equally correct and incorrect, thus $\int y h(x) p(x; y) d(x; y) = 0$. Otherwise, suppose there are a intervals before the threshold, and b intervals after it; H must be incorrect on at least $a/2 - 1$ of the left intervals, and $b/2 - 1$ of the right intervals; since $2k - 1 = a + b - 2k$ (this proof is charitable), thus at least $k - 2$ intervals are predicted completely incorrectly. Consequently,

$$\int y h(x) p(x; y) d(x; y) \leq \frac{k+2}{2k} \frac{k-2}{2k} = \frac{2}{k};$$

Thus the form of $\phi(\cdot)$ from Lemma 6.C.2 provides $\phi_{1=2k}(\cdot) \leq 2/k$, and so $\phi_0(\cdot) = 0$ by monotonicity (Lemma 6.E.7), and ϕ_0 by Proposition 6.3.3.

The remainder of the proof considers finite sample effects. Let $\epsilon > 0$, and a sample of size $m \geq 2$ be given. Choose integer $k := \lceil \frac{m^{1/4}}{3} \rceil \lceil \frac{1}{2 \ln(4/\epsilon)} \rceil$ (where the lower bound on m provides $k \geq 1$), and consider the behavior of density p_k , defined as above. Note first that $\int p_k^2 \leq 2k+1 \leq 3k \leq \frac{m^{1/4}}{3} \lceil \frac{1}{2 \ln(4/\epsilon)} \rceil$. Next, with probability at least $1 - \epsilon/2$, Hoeffding's bound grants

$$\int p_k d\mathbb{P} - \int p_k d\mathbb{P}_k \leq \frac{1}{2m} \ln \frac{4}{\epsilon} \leq \frac{1}{2m^{1/4}} \leq \frac{1}{2};$$

Now define $p_k^0 := p_k - \int p_k d\mathbb{P}$, which means $\int p_k^0 d\mathbb{P} = 0$, and furthermore $\int p_k^0 \leq 2 \int p_k \leq 2k+1$.

Since H has VC dimension $V(H) = 2$, Lemma 6.D.1 grants, with probability at least $1 - \epsilon$,

$$\sup_{h \in H} \int y h(x) p_k^0(x; y) d b(x; y) - \int y h(x) p_k^0(x; y) d (x; y) \leq \frac{2k p_k^0 k_1}{m} \left(4 \frac{1}{\ln(m+1)} + \frac{1}{2 \ln(2)} \right)$$

Now set $\tilde{p} := 1 - k p_k^0 k_1 \frac{1}{2 \ln(4)} = (2m^{1-4})$. Then $\tilde{p} \in D$ (b), and the above computations provide

$$(b) \quad \frac{16 \left(\frac{1}{\ln(m+1)} + 1 \right)}{m^{1-4}} + \frac{2}{k} \frac{16 \left(\frac{1}{\ln(m+1)} + 1 \right)}{m^{1-4}} + \frac{2}{m^{1-4} \left(3 \frac{1}{2 \ln(4)} \right) - 1};$$

and lastly Lemma 6.E.7 and Proposition 6.3.3 grant (b) \Rightarrow (b). □

Proof of Proposition 6.3.6. 1. This proof will start by proving $L > 0 \Leftrightarrow (C) > 0$ for some hard core $C \in \mathcal{C}$ ($\epsilon > 0$ for some $\epsilon > 0$). The result then follows by negating each of the three statements, and then simplifying them with the facts $L = 0$, $(C) = 0$, and $(\epsilon) = 0$ for all $\epsilon > 0$. To ease notation, let $S_h(H; \epsilon)$ denote the set of hard cores.

$(L > 0 \Rightarrow \exists C \in S_h(H; \epsilon) (C) > 0)$. If $L > 0$, then by Lemma 4.1.1 there exists a dual element p which is dual feasible, meaning $\int p d L^1(\epsilon) > 0$ and $k A^> p k_1 = 0$ and $([p < 0]) = 0$, and moreover p satisfies $\int p d R^1(\epsilon) < 0$, which means $([p > 0]) > 0$ since $\int p d (0) = 0$ by Lemma 1.B.1. As such, simplifying $k A^> p k_1 = 0$ via Lemma 1.A.5, $p \in D(H; \epsilon)$, and thus $[p > 0] \in S_D(H; \epsilon)$, and any hard core C satisfies $(C) = ([p > 0]) > 0$ as desired.

$((\exists C \in S_h(H; \epsilon) (C) > 0) \Rightarrow \exists \epsilon > 0 (\epsilon) = 0)$. Let $p \in S_D(H; \epsilon)$ be given so that $[p > 0] = C$ for any hard core C (they all have the same measure). Note that $k p k_1 > 0$, and define $\tilde{p} := p - k p k_1$, which satisfies $k \tilde{p} k_1 = 1$, but still $\tilde{p} \in L^1(\epsilon)$ and $k A^> \tilde{p} k_1 = 0$. Since \tilde{p} is a probability measure and $k \tilde{p} k_1 = 1$, it must hold that $k \tilde{p} k_1 = 1$, and thus the choice $\tilde{p} := 1 - k p k_1$ satisfies $\tilde{p} \in (0; 1]$. More importantly, since $\tilde{p} \in D(\epsilon)$, it follows that $(\tilde{p}) = 0$ as desired.

$((\exists \epsilon > 0 (\epsilon) = 0) \Rightarrow L > 0)$. Let $\epsilon > 0$ with $(\epsilon) = 0$ be given; then attainment in the duality formula in Lemma 6.C.2 provides the existence of $p \in L^1(\epsilon)$ with $k p k_1 = 1$ and $k A^> p k_1 = 0$. By Lemma 1.B.1, there exists $c > 0$ so that \tilde{p} is strictly negative along $(0; c)$. Consequently, $\tilde{p} := c p - k p k_1$ satisfies $k \tilde{p} k_1 > 0$, and $\tilde{p} \in (0; c)$ -a.e., thus $\tilde{p} < 0$ -a.e.,

and also $kA^>pk_1 = 0$; together, it follows that $L_{R, (p)d} > 0$ as desired.

2. This result is the same as Lemma 6.E.1. □

Proof of Lemma 6.3.7. This result is the combination of Lemma 6.C.2 and Lemma 1.A.5. □

6.E.3 Optimization Guarantees

Note that the following proof does not overtly use convexity; convexity however is used both algorithmically by the line searches (otherwise they are not efficient), and for their guarantees (cf. Lemmas 1.C.3 and 1.C.7).

Proof of Lemma 6.3.8. Consider any $0 \leq t \leq T - 1$. Since $\ell \in L_{ig}$,

$$kr \ell(A_t)k_1 = \frac{1}{m} \sum_{\substack{i \in [m] \\ (A_t)_{x_i, y_i} = 0}} \ell^0((A_t)_{x_i, y_i}) \quad t \geq 1:$$

Combining this with the fact that $\ell^0 \in [0, \ell_2]$, the vector $p_t := r \ell(A_t) = kr \ell(A_t)k_1$ satisfies $kp_tk_1 = 1$ and

$$kp_tk_1 = \frac{\ell_2}{kr \ell(A_t)k_1} = \frac{\ell_2}{t+1} = \frac{1}{t}.$$

where ℓ_1 is as provided in the statement (and $\ell_1 \geq 1$ since $t \geq 1$ and $\ell_1 \geq \ell_2$). Recalling the dual form $\ell_1(b) = \min \{kA^>pk_1 : p \in D_{\ell_1}(b)\}$ from Lemma 6.C.2, and noting that $p_t \in D_{\ell_1}(b)$,

$$kA^>r \ell(A_t)k_1 = kr \ell(A_t)k_1 kA^>p_tk_1 = \frac{\ell_1(b)}{t+1}.$$

Plugging this into the single-step guarantees from the three line search choices (cf. Lemmas 1.C.3 and 1.C.7 and the Lipschitz gradient guarantee from Lemma 6.B.2),

$$\begin{aligned} \ell(A_t) - \ell(A_{t-1}) &\leq \frac{2kA^>r \ell(A_{t-1})k_1^2}{6B_2} \\ \ell(A_t) - \ell(A_{t-1}) &\leq \frac{(\ell_1 - t \ell_1 \ell_1^{-1}(b))^2}{6B_2}. \end{aligned}$$

The desired result comes by summing across all iterations and noting $\ell(A_0) = \ell(0) = \ell^0(0)$. □

6.E.4 Statistical Guarantees

Proof of Theorem 6.3.2. The first step of the proof is to show $\mathbb{P}(H^\wedge)$. Thus consider the case that every iteration has $\mathbb{P}(H_t) > \epsilon$; by the monotonicity of $\mathbb{P}(H_t)$ (cf. Lemma 6.E.7), positivity of $\mathbb{P}(H_t)$ (cf. Proposition 6.3.6), together with the bound on m (and the deviations on $\mathbb{P}(H_t)$ in Proposition 6.3.6), with probability at least $1 - \epsilon = 2$,

$$\mathbb{P}(b) - \mathbb{P}(b) - \mathbb{P}(b) \leq \frac{1}{2m} \ln \frac{4}{\epsilon} - \frac{\mathbb{P}(b)}{2} > 0;$$

where the last equality is also by Proposition 6.3.6. Thus, by Lemma 6.3.8, and the monotonicity of $\mathbb{P}(H_t)$ in t , and the second lower bound on m (and thus on m^a),

$$\begin{aligned} \mathbb{P}(A_T) &\leq \sum_{t=1}^T \frac{\mathbb{P}(H_t) \mathbb{P}(b)^2}{6B_2} \\ &\leq \sum_{t=1}^T \frac{m^a (\mathbb{P}(H_t))^2}{24B_2} \\ &\leq 0; \end{aligned}$$

a contradiction since \mathbb{P} is nonnegative, and moreover positive on regions where it makes mistakes, therefore the above indicates $\mathbb{P}(A_T) = 0 < \epsilon$. As such, thanks to the final step of Figure 1.3 picking out the iterate with lowest classification error, $\mathbb{P}(H^\wedge)$.

What remains is to establish a deviation inequality. Let H_t denote the hypothesis class used by predictor \hat{f}_t (i.e., $H_t = \{f_{i=1}^t c_i h_i : c_i \in \mathbb{R}; h_i \in \mathcal{H}_g\}$), and let $S_{H_t}(m)$ denote the corresponding shatter coefficient when H_t is applied to the sample of size m (Boucheron et al., 2005a, Section 3). It follows (Schapire and Freund, 2012, Lemma 4.5) that

$$S_{H_t}(m) \leq \frac{em}{t} + \frac{em}{V(H)^t} :$$

Plugging this and $t = m^a$ into an appropriate VC theorem and simplifying (Boucheron et al., 2005a, Theorem 5.1 and subsequent discussion), with probability at least $1 - \epsilon = 2$,

$$\begin{aligned} R(A^\wedge) &\leq \frac{1}{m} \ln(S_{H_t}(2m)) + \ln(8/\epsilon) + 4 \frac{\ln(S_{H_t}(2m)) + \ln(8/\epsilon)}{m} \\ &\leq \frac{1}{m^{1-a}} \frac{(V(H) + 1) \ln(2em) + \ln(8/\epsilon)}{m^{1-a}} + 4 \frac{(V(H) + 1) \ln(2em) + \ln(8/\epsilon)}{m^{1-a}}. \end{aligned}$$

□

Appendix 6.F Deferred Material from Section 6.4

6.F.1 Proof of Proposition 6.4.1

Lemma 6.F.1. Let $\psi \in L^1_{[0,g]}$ and any $g \in \mathbb{R}$ be given.

1. The restriction of ψ to $[0; g]$, denoted $\psi_{[0;g]}$, is a (decreasing) bijection between $[0; g]$ and $[\psi(0); 0]$.
2. Let $f(z) := (\psi_{[0;g]})^{-1}(z)$ denote the inverse of $\psi_{[0;g]}$. If μ is a probability measure, and $p \in L^1(\mu)$, and $r := \int_{\mathbb{R}} p \, d\mu \in [0; \psi(0)]$, then $(\{p \leq f(r)\}) = f(r)$, with $f(r) > 0$ if $r > 0$.

Proof. Choose any $g \in \mathbb{R}$; by Lemma 1.B.1, ψ is 0 at 0, negative along $(0; g)$, and attains its minimum at g . Since ψ is finite for every $z \in (0; g)$, then $\psi^{-1}(z)$ exists (Rockafellar, 1970, Theorem 23.4), and every $x \in \psi^{-1}(z)$ satisfies $x \leq \psi^{-1}(z)$ (Rockafellar, 1970, Theorem 23.5), and so ψ is strictly convex along $(0; g)$ (Hiriart-Urruty and Lemaréchal, 2001, Theorem E.4.1.2), meaning ψ is injective along $[0; g]$. Since $\psi(0) = 0$, and $\psi(g) = \psi(0)$ (by the Fenchel-Young inequality), and since ψ is lower semi-continuous (Rockafellar, 1970, Theorem 12.2), then ψ is also surjective from $[0; g]$ to $[\psi(0); 0]$.

Now let f denote the inverse map (from $[\psi(0); 0]$ to $[0; g]$), and let probability measure μ , function $p \in L^1(\mu)$, and scalar $r := \int_{\mathbb{R}} p \, d\mu \in [0; \psi(0)]$ be given (where the containment provides $f(r)$ is valid). By Jensen's inequality,

$$r = \int_{\mathbb{R}} \psi(p) \, d\mu;$$

since f is a decreasing map, this implies $\int_{\mathbb{R}} p \, d\mu \leq f(r)$. Furthermore,

$$f(r) = \int_{\{p \leq f(r)\}} p \, d\mu + \frac{f(r)}{2} + \int_{\{p > f(r)\}} p \, d\mu = \int_{\{p \leq f(r)\}} p \, d\mu + \frac{f(r)}{2} + (\int_{\{p > f(r)\}} p \, d\mu - \int_{\{p > f(r)\}} f(r) \, d\mu);$$

meaning $(\int_{\{p \leq f(r)\}} p \, d\mu = f(r)) = f(r)$ as desired.

Lastly, the statement $f(r) > 0$ if $r > 0$ follows from the bijectivity of $\psi_{[0;g]}$. □

Proof of Proposition 6.4.1. The basic duality relation is provided by Lemma 4.1.1. Since the optimal value satisfies $L^*(0) \leq L^*(0)$ (since μ is a probability measure, μ is primal feasible, and $\mu \in \mathcal{P}(X)$), then Lemma 6.F.1 may be applied with parameter p being the dual optimum and parameter r being the corresponding objective value $L^*(0)$. \square

6.F.2 Proof of Lemma 6.4.5

The first step is to use p to show that if μ has low error and norm, then H will have very small margins over some positive measure set.

Lemma 6.F.2. Let convex differentiable $\phi : \mathbb{R} \rightarrow \mathbb{R}_+$ with $\phi'(0) > 0$, any class H , and any probability measure μ over $X = \{1, +1\}$ with empirical counterpart $\hat{\mu}$ be given. Suppose the following quantities and constants exist.

1. Suppose there exists $p \in L^1(\mu)$ with $p \in [0, k_1]$ μ -a.e., and that there exists $\epsilon > 0$ with $b(\mu - \hat{\mu}) \leq \epsilon$.

2. Set

$$c := \frac{\phi'(0)}{\phi''(0)} \max\left\{1, \frac{1}{\epsilon} \max_{f \in H} \int f d\mu - \int f d\hat{\mu}\right\};$$

and suppose there exist $C > 0$ and $D = ck_1^2 \leq 8$ so that every μ with $k_1 \leq C$ satisfies $\int \mu^R(A) d\mu \leq D$.

If μ satisfies $k_1 \leq C$ and $b(\mu - \hat{\mu}) \leq \epsilon$, then $L^*(A) \geq 2L^*(0)$.

Proof of Lemma 6.F.2. First consider the case that $b(\mu - \hat{\mu}) \leq \epsilon$. By subgradient rules for convex functions, since $\phi''(0) > 0$ and $\phi'(0) > 0$, and using the definition of c ,

$$\begin{aligned} \int \mu^R(A) d\mu &= \int \mu^R(A) d\mu \\ &\geq \int_{[y(H)]_x \leq c} \mu^R(A) d\mu \\ &\geq \int_{[y(H)]_x \leq c} (\phi'(0) + \phi''(0)(A - 0)) d\mu \\ &\geq \frac{\phi'(0)c}{8} \\ &\geq 2L^*(0); \end{aligned}$$

meaning $L^*(A) \geq 2L^*(0) = 2L^*(0)$.

Now consider the remaining possibility that $b(\mu - \hat{\mu}) > \epsilon$. Since by assumption $b(\mu - \hat{\mu}) \leq 1$, it follows that $b(\mu - \hat{\mu}) \leq 1$. In turn, it also holds

that

$$b([y(H)_x \leq c] \setminus [p \leq 1]) = 4:$$

Next, the definition of D provides that

$$D = \int_{[y(H)_x \leq 0]} A \, p \, d b + \int_{[y(H)_x > 0]} y(H)_x \, p(x; y) \, d b(x; y);$$

meaning

$$\int_{[y(H)_x \leq 0]} y(H)_x \, p(x; y) \, d b(x; y) + \int_{[y(H)_x > 0]} y(H)_x \, p(x; y) \, d b(x; y) = D;$$

As such, since $\mathbb{2} [0; kpk_1]$ over the sample,

$$\begin{aligned} kpk_1 \int_{[y(H)_x \leq 0]} y(H)_x \, d b(x; y) &= \int_{[y(H)_x \leq 0]} y(H)_x \, p(x; y) \, d b(x; y) \\ &+ \int_{[y(H)_x > 0]} y(H)_x \, p(x; y) \, d b(x; y) = D \\ &+ \int_{[y(H)_x \leq c] \setminus [p \leq 1]} y(H)_x \, p(x; y) \, d b(x; y) = D \\ &+ \frac{c^2}{4} = D; \end{aligned}$$

Turning back to \mathbb{L}^b and proceeding similarly to the earlier case,

$$\begin{aligned} \int_{[y(H)_x \leq 0]} \hat{y}(A) \, d b &= \int_{[y(H)_x \leq 0]} \hat{y}(A) \, d b \\ &+ \int_{[y(H)_x > 0]} (\hat{y}(0) + \hat{y}^0(0)(A)) \, d b \\ &= \frac{\hat{y}^0(0)}{kpk_1} + \frac{c^2}{4} = D \\ &= 2 \hat{y}^0(0); \end{aligned}$$

which again yields $\mathbb{L}^b(A) = 2 \hat{y}^0(0) = 2 \mathbb{L}^b(0) = 2 \mathbb{L}^b(A_0)$. □

Next, these small margins in turn cause the line search to not look too far, meaning the next iterate will also have some small margins. Note that H is assumed binary; this is in order to

changes in H to changes in \cdot .

Lemma 6.F.3. Let convex $\ell : \mathbb{R} \rightarrow \mathbb{R}_+$, binary H , and probability measure \mathbb{P} with empirical counterpart $\hat{\mathbb{P}}$ be given. Let positive reals C, c be given so that ℓ with $k_1 \leq C + 2c$ and $b(\hat{\mathbb{P}}[H \leq c]) \leq 1 - \epsilon$ implies $\ell(A) \leq \ell^*(0)$. Then for any ℓ with $k_1 \leq C$ and $\ell(A) \leq \ell^*(0)$, the set of line search candidates

$$S := \{ \ell^0 \in \mathbb{R} : \ell(A^0) < \ell^*(0); \exists h \in H, \ell^0 := \ell + e_h \}$$

satisfies $b(\hat{\mathbb{P}}[H \leq c]) < 1 - \epsilon$ for every $\ell^0 \in S$.

Proof. Let ℓ with $k_1 \leq C$ be given, and consider any ℓ^0 of the form $\ell^0 := \ell + e_h$ for some $\ell \in \mathbb{R}$ and $h \in H$. The desired statement will be shown by contrapositive; namely, $b(\hat{\mathbb{P}}[H \leq c]) \leq 1 - \epsilon$ implies $\ell^0 \in S$.

For any example (x, y) , since H has binary predictors, the map $\ell \mapsto \ell(y(H(\ell + e_h)))_x$ is constant for $\ell > 2c$ and for $\ell < -2c$; consequently, since $\hat{\mathbb{P}}$ is a discrete measure over a finite set, the map $\ell \mapsto b(\hat{\mathbb{P}}[y(H(\ell + e_h))_x \leq c])$ is also constant for $\ell > 2c$ and $\ell < -2c$. As such, the existence of ℓ^0 as above implies the existence of $\ell^{00} := \ell + e_h$ where $h \in H$ is as before, ℓ^{00} and ℓ^0 have the same sign, and $|\ell^0 - \ell^{00}| \leq 2c$; in other words, ℓ^{00} is along the path from ℓ to ℓ^0 , but moreover satisfies $k_1 \leq 2c$. But this means $k_1 \leq k_1 + k_1 \leq C + 2c$, whereby the stated assumptions combined with $b(\hat{\mathbb{P}}[H \leq c]) \leq 1 - \epsilon$ provide $\ell(A^{00}) \leq \ell^*(0)$, thus $\ell^{00} \in S$. Furthermore, since ℓ^{00} is along the path from ℓ to ℓ^0 , and $\ell(A) \leq \ell^*(0)$ and $\ell(A^{00}) \leq \ell^*(0) \leq \ell(A)$, it follows by convexity that $\ell(A^0) \leq \ell^*(0)$, and thus $\ell^0 \in S$ as well. \square

These small margin controls directly give a bound on step sizes.

Lemma 6.F.4. (See also Bartlett and Traskin (2007, eq. (28)).) Let $\ell \in L_{1g} \setminus L_{2d}$ be given with Lipschitz gradient parameter B_2 , binary class H , time horizon t , and empirical probability measure $\hat{\mathbb{P}}$ corresponding to a sample of size n be given. Let positive real $B_1 > 0$ be given so that for any ℓ with

$$k_1 \leq \frac{s}{t \max\{5, 2B_2\} B_1 g(\ell(A_0) - \ell(A_t))},$$

and any line search candidate $\ell^0 := \ell + e_h$ for $h \in H$, $\ell \in \mathbb{R}$, and satisfying $\ell(A^0) \leq \ell(A)$, then $B_1 \frac{1}{m} \sum_{i=1}^m \ell^0(A^0_i)$. The following properties hold.

1. For every integer $0 \leq i < t$, an optimal step α_i exists, and every step size choice satisfies

$$\alpha_i^2 \leq \min \left(\frac{9k_A^2 r \ell(A_{i-1})k_1^2}{4^2 B_1^2}; \frac{\max\{5; 2B_2=B_1g(\ell(A_{i-1}) - \ell(A_i))\}}{2B_1} \right);$$

2. For every integer $0 \leq i < t$ and any sequence of step size choices,

$$k_i \leq k_1 \prod_{j=1}^i \frac{p_j}{q_j} \leq \frac{p_i}{q_i} \frac{\max\{5; 2B_2=B_1g(\ell(A_0) - \ell(A_i))\}}{2B_1} \leq \frac{p_i}{q_i} \frac{\max\{5; 2B_2=B_1g\}}{2B_1};$$

Proof. This proof establishes both properties simultaneously by induction on i . In the base case $i = 0 = 0$ and there is nothing to show, thus suppose $i \geq 1$.

Define the interval $I_i := \{x \in [0; \ell(A_{i-1} + v_i)) \mid \ell(A_{i-1}) \leq x \leq \ell(A_i)\}$. Combining the inductive hypothesis (controlling $k_{i-1}k_1$) with the assumptions on line search candidates means the second-order lower bound B_1 is active along I_i . Now consider a Taylor expansion of ℓ , but in the direction reverse to Lemma 1.C.3, and using the fact that H is binary; then for any $x \in I_i$ and some $\xi \in [x; x + v_i]$,

$$\begin{aligned} \ell(A_{i-1} + v_i) &= \ell(A_{i-1}) + \frac{D}{A} \ell(A_{i-1}; v_i) + \frac{E}{2m} \sum_{i=1}^m \ell^{(i)}(Az_i)(Av_i)^2 \\ &\geq \ell(A_{i-1}) - k_A^2 \ell(A_{i-1})k_1 + \frac{2}{2} B_1 v_i^2 \end{aligned}$$

This last expression defines a univariate quadratic which lies below $\ell(A_{i-1} + v_i)$ along I_i (outside of I_i , the constraints granting the lower bound B_1 may be violated). Consequently, I_i is bounded, and the optimal step α_i must exist, and moreover satisfies

$$\alpha_i := \frac{k_A^2 \ell(A_{i-1})k_1}{B_1}; \tag{6.F.5}$$

where α_i is the minimizer to the above quadratic. Plugging α_i back into the quadratic, for any $x \in I_i$,

$$\ell(A_{i-1} + v_i) - \ell(A_{i-1}) \geq \frac{k_A^2 \ell(A_{i-1})k_1^2}{2B_1}; \tag{6.F.6}$$

Now consider the first two step size options in Figure 1.3; combining eq. (6.F.6) with

Lemma 1.C.3 and the Lipschitz gradient guarantee in Lemma 6.B.2,

$$\frac{2}{i} \frac{2}{i} \frac{kA > \mathbb{L}(A_{i-1})k_1^2}{B_1^2} \frac{2B_2(\mathbb{L}(A_{i-1}) - \mathbb{L}(A_i))}{2B_1^2}$$

as desired.

For option 3 (the Wolfe search), combining eq. (1.C.5) with eq. (6.F.6) grants

$$\frac{2}{i} \frac{9(\mathbb{L}(A_{i-1}) - \mathbb{L}(A_i))^2}{2kA > r \mathbb{L}(A_{i-1})k_1^2} \frac{9(\mathbb{L}(A_{i-1}) - \mathbb{L}(A_i))}{2^2 B_1} \frac{9kA > r \mathbb{L}(A_{i-1})k_1^2}{4^2 B_1^2},$$

which establishes the first inductive property for all step sizes.

The second statement is just Cauchy-Schwarz combined with the bound on $\frac{2}{i}$:

$$\begin{aligned} k_i k_1 & \sum_{j=1}^i \langle v_j, k_i \rangle \\ & \leq \sqrt{\sum_{j=1}^i \langle v_j, v_j \rangle} \sqrt{\sum_{j=1}^i \langle k_i, k_i \rangle} \\ & \leq \sqrt{\sum_{j=1}^i \frac{1}{\rho_i} \langle X_j, X_j \rangle} \sqrt{\sum_{j=1}^i \frac{1}{\rho_i} \langle X_j, X_j \rangle} \\ & \leq \sqrt{\sum_{j=1}^i \frac{1}{\rho_i} \frac{1}{\max\{5, 2B_2=B_1\}g} \frac{(\mathbb{L}(A_{j-1}) - \mathbb{L}(A_j))}{2B_1}} \\ & = \sqrt{\sum_{j=1}^i \frac{1}{\rho_i} \frac{1}{\max\{5, 2B_2=B_1\}g} \frac{(\mathbb{L}(A_0) - \mathbb{L}(A_i))}{2B_1}}; \end{aligned}$$

and nonnegativity of ρ_i and the fact that all steps perform descent grants $\mathbb{L}(A_0) - \mathbb{L}(A_i) \geq 0$. □

After some algebra, the upper bound also grants a lower bound; due to this indirection, it should be possible to improve this bound. Note that the beginning of this derivation, when initially lower bounding ρ_i , uses derivations similar to those used by Zhang and Yu (2005) and Bartlett and Traskin (2007).

Lemma 6.F.7. Let $\mathcal{L} \in L_{lg} \setminus L_{2d}$ be given with Lipschitz gradient parameter B_2 , binary class H , time horizon t , and empirical probability measure ρ corresponding to a sample of size n be given. Suppose there exists $c_2 > 0$ with $k_i k_1 \leq c_2 \rho_i$ for all $0 \leq i \leq t$. Additionally, let $C_2 > 0$ and

2 be given with $k_1 \leq C_2$ and $\epsilon := \min_{i \in [t-1]} \mathbb{E}[\mathbb{L}(A_i)] - \mathbb{E}[\mathbb{L}(A)] \geq 0$. Then

$$i \leq \frac{(\mathbb{E}[\mathbb{L}(A_{i-1})] - \mathbb{E}[\mathbb{L}(A)])}{2B_2(k_1 + c_2 \frac{1}{i})}$$

and

$$\sum_{i=1}^t \frac{1}{2B_2 c_2} \left(2^{\frac{p}{t-1}} + \frac{c_2}{C_2} + \frac{2C_2}{c_2} \ln \frac{C_2}{C_2 + c_2 \frac{1}{t-1}} \right);$$

or more simply $\sum_{i=1}^t \frac{1}{2B_2 c_2} \left(2^{\frac{p}{t-1}} - 1 \right) = (4B_2 c_2)^{-1}$ if $C_2 \geq c_2 \frac{1}{t-1}$.

Proof. For any $1 \leq i \leq t$, note by Lemma 1.A.5 that

$$\begin{aligned} k_{i-1} \mathbb{E}[\mathbb{L}(A_{i-1})] &= \sup_{k \leq k_{i-1}} \mathbb{E}[\mathbb{L}(A_{i-1})]; \\ &= \frac{1}{k_{i-1}} \mathbb{E}[\mathbb{L}(A_{i-1})]; \\ &= \frac{1}{k_{i-1}} (\mathbb{E}[\mathbb{L}(A_{i-1})] - \mathbb{E}[\mathbb{L}(A)]); \end{aligned}$$

Thus, by the lower bounds in Lemmas 1.C.3 and 1.C.7 for every step size (combined with the Lipschitz gradient guarantee in Lemma 6.B.2),

$$i \leq \frac{k_{i-1} \mathbb{E}[\mathbb{L}(A_{i-1})] - \mathbb{E}[\mathbb{L}(A)]}{2B_2} \frac{1}{k_{i-1}} \frac{1}{k_1} \frac{1}{2B_2(k_{i-1} + k_1)};$$

Combining this with the provided upper bound on k_{i-1} ,

$$i \leq \frac{1}{2B_2(k_1 + k_{i-1}k_1)} \frac{1}{2B_2(k_1 + c_2 \frac{1}{i-1})}$$

As a consequence of this, and recalling the simplification $k_1 \leq C_2$,

$$\begin{aligned} \sum_{i=1}^t \frac{1}{2B_2 c_2} &= \frac{1}{2B_2 c_2} \left(\frac{c_2}{C_2} + \sum_{i=1}^{t-1} \frac{c_2}{C_2 + c_2 \frac{1}{i}} \right) \\ &= \frac{1}{2B_2 c_2} \left(\frac{c_2}{C_2} + \int_0^{t-1} \frac{c_2 dx}{C_2 + c_2 \frac{1}{x}} \right) \\ &= \frac{1}{2B_2 c_2} \left(\frac{c_2}{C_2} + 2^{\frac{p}{t-1}} \frac{2C_2 \ln(C_2 + c_2 \frac{1}{t-1})}{c_2} \right) \\ &= \frac{1}{2B_2 c_2} \left(\frac{c_2}{C_2} + 2^{\frac{p}{t-1}} + \frac{2C_2}{c_2} \ln \frac{C_2}{C_2 + c_2 \frac{1}{t-1}} \right); \end{aligned}$$

When $C_2 = c_2^p t^{-1}$, it suffices to instantiate the above bound with $C_2^0 := c_2^p t^{-1}$ (whereby it still holds that $k \leq k_1 = C_2 = C_2^0$), and then rearrange, noting $1 - \ln(2) \geq 1/4$ and deleting the nonnegative standalone term $c_2 = C_2$. \square

By combining the above chain of results, the proof of Lemma 6.4.5 follows.

Proof of Lemma 6.4.5. Recalling the structure from Proposition 6.4.1, let dual optimum p and real $\epsilon > 0$ be given so that $([p - \epsilon])$. By Hoeffding's inequality and the first lower bound on m , with probability at least $1 - \epsilon/4$

$$b([p - \epsilon]) \leq ([p - \epsilon]) \frac{s}{2m} \ln \frac{4}{\epsilon} \leq \frac{\epsilon}{2}$$

Henceforth disregard the corresponding failure event.

Next define

$$D := \frac{2(R_t + 2c)kpk_1}{m^{1-2}} \leq 2^p \frac{2V(H) \ln(m+1) + p}{2 \ln(4/\epsilon)}$$

The second condition on m grants $D \leq \epsilon^2/8$, whereby Lemma 6.4.4 grants, with probability at least $1 - \epsilon/4$, that every \mathcal{L} with $k \leq k_1 = R_t + 4c$ satisfies

$$\mathbb{P}(A) \leq D \leq \frac{\epsilon^2}{8}$$

Discard the corresponding failure event as well; unioning this with the earlier failure event, the remaining steps hold with probability at least $1 - \epsilon/2$.

It follows from Lemma 6.F.2 that every such \mathcal{L} with $k \leq k_1 = R_t + 2c$ either satisfies $b([jH - j - c]) < 1 - \epsilon/8$, or else $\mathbb{P}(A) \leq \epsilon^2/8$. This in turn means the preconditions to Lemma 6.F.3 are met (with $C := R_t$), and in particular, for any \mathcal{L} with $k \leq k_1 = R_t$ and $\mathbb{P}(A) \leq \epsilon^2/8$, the set of line search candidates

$$S := \{ \ell \in \mathcal{L} : \mathbb{P}(A_\ell) < \epsilon^2/8 \} \cup \{ \ell \in \mathcal{L} : R_\ell \leq R_t + 2c \}$$

satisfies $b([jH^c] \setminus c) < 1 - \epsilon$ for every $0 \leq S$. But then, for every $0 \leq S$ (and note $2 \leq S$),

$$\frac{1}{m} \sum_{i \in [m]} \mathbb{E} \left[\ell(A_{x_i, y_i}) \right] \leq \frac{1}{m} \sum_{i \in [m]} \mathbb{E} \left[\ell(A_{x_i, y_i}) \right] + b([jH^c] \setminus c) \inf_{z \in [c, +c]} \mathbb{E} \left[\ell(z) \right] + B_1$$

This establishes the first desired statement.

For the second statement (upper bounds on ϵ_i and $k_i \leq k_1$), note that the above properties satisfy the preconditions to Lemma 6.F.4, whereby the desired upper bounds follow.

Similarly, for the third statement (lower bounds on ϵ_i and $k_i \geq k_1$), the preconditions for Lemma 6.F.7 are now met. □

6.F.3 Optimization Guarantees

As stated previously, the following proof is a reworking of a proof due to Zhang and Yu (2005), albeit with the present decoupling of line search and coordinate selection.

Proof of Lemma 6.4.6. Let ϵ_i be arbitrary. The first step of this proof is to develop two lower bounds on $k_i \geq r \mathbb{E}[\ell(A_{i-1})] k_1$. First, just as in the proof of Lemma 6.F.7,

$$\begin{aligned} k_i \geq r \mathbb{E}[\ell(A_{i-1})] k_1 &= \sup_{k \geq k_1} \mathbb{E} \left[\ell(A_{i-1}) \right]; \\ &= \frac{1}{k_{i-1} k_1} \mathbb{E} \left[\ell(A_{i-1}) \right]; \\ &= \frac{1}{k_{i-1} k_1} (\mathbb{E}[\ell(A_{i-1})] - \mathbb{E}[\ell(A)]); \end{aligned}$$

The second lower bound is provided by assumption, and thus, by the guarantee on any line search as in Lemmas 1.C.3 and 1.C.7 (together with the Lipschitz gradient property in Lemma 6.B.2),

$$\begin{aligned} \mathbb{E}[\ell(A_i)] - \mathbb{E}[\ell(A)] &\leq \mathbb{E}[\ell(A_{i-1})] - \mathbb{E}[\ell(A)] + \frac{2k_i \geq r \mathbb{E}[\ell(A_{i-1})] k_1^2}{6B_2} \\ &\leq \mathbb{E}[\ell(A_{i-1})] - \mathbb{E}[\ell(A)] + \frac{c_3^2 \epsilon_i (\mathbb{E}[\ell(A_{i-1})] - \mathbb{E}[\ell(A)])}{6B_2 k_{i-1} k_1} \\ &= \mathbb{E}[\ell(A_{i-1})] - \mathbb{E}[\ell(A)] \left(1 - \frac{c_3^2 \epsilon_i}{6B_2 k_{i-1} k_1} \right) \\ &\leq \mathbb{E}[\ell(A_{i-1})] - \mathbb{E}[\ell(A)] \exp \left(-\frac{c_4 \epsilon_i}{k_{i-1} k_1} \right); \end{aligned}$$

where the last step took $c_4 := c_3^2 = (6B_2)^2$ for convenience. Iterating this bound,

$$\mathbb{E}(A_t) - \mathbb{E}(A) \leq \mathbb{E}(A_0) - \mathbb{E}(A) + \exp(-c_4) \sum_{i=1}^t \frac{1}{k + k_{i-1}k_1}.$$

Focusing on the summation, define $S_i := \sum_{j=1}^i k_j$ with $S_0 := 0$, whereby $k_j k_1 \geq S_{j-1}$. Using this (see also the similar derivation by Zhang and Yu (2005, Proof of Lemma 4.2)),

$$\begin{aligned} \sum_{i=1}^t \frac{1}{k + k_{i-1}k_1} &= \sum_{i=1}^t \frac{1}{k + S_{i-1}} \\ &= \sum_{i=1}^t \left(\frac{k + S_i}{k + S_{i-1}} - \frac{k + S_{i-1}}{k + S_{i-1}} \right) \\ &= \sum_{i=1}^t \ln \frac{k + S_i}{k + S_{i-1}} = \ln \frac{k + S_t}{k} : \end{aligned}$$

Plugging this into the preceding display and collecting terms, the result follows. □

Note that the substitution $k_j k_1 \geq S_{j-1}$ at the end of the proof of Lemma 6.4.6 works around the fact that $\sum_{j=1}^i k_j$ could be much larger than $k_j k_1$; this issue is frequently avoided in the literature by assuming that H is closed under negation, whereby $v_i = e_{h_i}$ in each round (i.e., rather than $v_i = 2^i e_{h_i}$).

6.F.4 Statistical Guarantees

Proof of Theorem 6.4.2. Let $\epsilon > 0$ be arbitrary, and choose δ with $k + k_1 \geq R_{t-1}$ so that

$$L(A) \leq \inf_{k + k_1 \geq R_{t-1}} L(A) :$$

By McDiarmid's inequality and the fact that $\sup_{x,y} |j(A)_{x,y} - j| \leq R_{t-1}$, with probability at least $1 - \epsilon/6$,

$$\mathbb{E}(A) \leq R_{t-1} \frac{2}{m} \ln \frac{6}{\epsilon} + L(A) + R_{t-1} \frac{2}{m} \ln \frac{6}{\epsilon} + \inf_{k + k_1 \geq R_{t-1}} L(A) :$$

Now let $\beta \in (0, 1/2)$ be arbitrary, and consider two cases for the difference $\Delta := \min_{i \in [t-1]} \mathbb{E}(A_i) - \mathbb{E}(A)$.

If $(t - 1)^b$, then

$$\mathbb{P}(A_t) \leq \mathbb{P}(A_{t-1}) + (t - 1)^b + \mathbb{P}(A) :$$

Otherwise $> (t - 1)^b$. Thus Lemma 6.4.5 grants, with probability at least $1 - \epsilon = 2$,

$$\sum_{i=1}^t \frac{9k_A > r \mathbb{P}(A_{i-1})k_1^2}{4^2 B_1^2} \quad \text{and} \quad \sum_{i=1}^t \frac{P_{t-1}}{4B_2 R_1} ;$$

which can then be plugged into Lemma 6.4.6 (with $c_3 := 2 B_1 = 3$), together with the lower bound on $\mathbb{P}(A)$, to yield

$$\begin{aligned} \mathbb{P}(A_t) - \mathbb{P}(A) &\leq \mathbb{P}(A_0) - \mathbb{P}(A) + \frac{k k_1}{k k_1 + \sum_{i=1}^t k_1} \quad ! c_3^2 = (6B_2) \\ &\leq \mathbb{P}(A_0) - \mathbb{P}(A) + \frac{k k_1}{k k_1 + (t - 1)^{1=2} b = (4B_2 R_1)} \quad {}^3 B_1 = (9B_2) : \end{aligned}$$

Summing these two bounds gives a relation which holds in general; consequently, with probability at least $1 - \epsilon = 2$ (due to the invocation of Lemma 6.4.5),

$$\mathbb{P}(A_t) - \mathbb{P}(A) \leq (t - 1)^b + \epsilon(0) \frac{k k_1}{k k_1 + (t - 1)^{1=2} b = (4B_2 R_1)} \quad {}^3 B_1 = (9B_2) :$$

The first desired claim follows (with probability $1 - \epsilon = 6$) by recalling the earlier application of McDiarmid's inequality to $\mathbb{P}(A)$, combined with $\epsilon \neq 0$, the choice $b := 1 = 4$ (which is not optimal, but neither is the exponent $B_1^3 = (9B_2)$), and standard Rademacher bounds for voting classifiers applied to Lipschitz losses (Boucheron et al., 2005a, Theorem 4.1, eq. (8), and their proofs, which control for L), which makes use of the bound $k_t k_1 \leq R_t$ (granted by the earlier instantiation of Lemma 6.4.5), and simplifying via $t - 1 \leq t$,

$$\begin{aligned} L(A_t) &\leq \inf_{k, k_1, R_{t-1}} L(A) + t^{-1=4} + R_{t-1} \frac{s}{m} \ln \frac{6}{\epsilon} \\ &\quad + \frac{2 \cdot 2 R_t}{m^{1=2}} 2^p \frac{2V(H) \ln(m+1)}{2V(H) \ln(m+1)} + \epsilon(R_t)^p \frac{2 \ln(6)}{2 \ln(6)} \\ &\quad + \epsilon(0) \frac{k k_1}{k k_1 + t^{1=4} = (4B_2 R_1)} \quad {}^3 B_1 = (9B_2) : \end{aligned}$$

Plugging in $t = m^a$ gives the first result.

For the second guarantee, since $H \in \mathcal{F}_{ds}(X)$, then Lemma 6.A.1 grants

$$\inf_{g \in H} L(g) = \inf_{f \in \mathcal{F}_b} \int_Z \ell(yf(x)) d(x; y)$$

where \mathcal{F}_b is the family of Borel measurable functions from X to \mathbb{R} . From here, since ℓ is classification calibrated (Bartlett et al., 2006, Theorem 2, noting that the present thesis instead takes losses to be nondecreasing rather than nonincreasing), there exists a function $R: \mathbb{R} \rightarrow \mathbb{R}$ satisfying

$$R(A_t) - R(L(A_t)) = \inf_{f \in \mathcal{F}_b} \int_Z \ell(yf(x)) d(x; y) - \inf_{g \in H} L(g)$$

(where the last step used the previous display), and moreover $\lim_{z \rightarrow 0} R(z) = 0$ (Bartlett et al., 2006, Theorem 1). The specialization for ℓ_{\log} is due to Zhang (2004, Subsection 3.5 and Corollary 3.1).

The final guarantee follows by applying a version of the VC theorem to the predictors, and follows the exact strategy as in Theorem 6.3.2 (but using failure probability $= \epsilon$), and making use of the equality

$$R(\hat{H}_t) - R(H_t) = (R(\hat{H}_t) - R(H_t)) + (R(H_t) - R(H_t)) + (R(H_t) - R(H_t));$$

and the fact that $R(\hat{H}_t) \leq R(H_t)$. □

Appendix 6.G Proof of Consistency

Proof of Theorem 6.2.2. This proof is a standard application of the Borel-Cantelli Lemma; for an exposition on such applications, please see the proof of consistency of AdaBoost due to Schapire and Freund (2012, Proof of Corollary 12.3). In particular, let any $\epsilon > 0$ be given, and let E_m be the event that the output \hat{H}_m , trained on m examples, has classification risk $R(\hat{H}_m)$ exceeding the Bayes risk R_γ by more than ϵ ; to prove consistency, it suffices (thanks to Borel-Cantelli) to show that $\sum_m \Pr(E_m) < 1$.

There are two cases to consider: either $L = 0$, or $L > 0$. In the case that $L = 0$,

instantiate the finite sample guarantee in Theorem 6.3.2 for each $m \geq 1$ with $\epsilon = 2$ and $\delta := 1/m^2$; there is a real $M < \infty$ where $m > M$ provides the preconditions on the bound are met and the bound is at most ϵ , and thus

$$\sum_{m=1}^{\infty} \Pr(E_{m; \epsilon}) \leq \sum_{m=1}^{\infty} \left(1 + \sum_{m=M+1}^{\infty} \frac{1}{m^2} \right) \epsilon < 1 :$$

When $L > 0$, once again instantiate a relevant finite sample guarantee, this time from Theorem 6.4.2, with $\delta := 1/m^2$. It will be necessary to use all three guarantees; first, let M_1 be sufficiently large so that the third guarantee provides $R(H^m) \leq R(H_m^0) + \epsilon$ with failure probability m^{-2} for all $m > M_1$, where H_m^0 is the last iterate considered by the algorithm when run on m examples (thus H_m^0 is basically H_m^a , modulo rounding issues). Next, by the second guarantee, there exists ϵ_0 small enough so that

$$\epsilon_0 \leq (R(H_m^0) - R^*) \implies R(H_m^0) \leq R^* + \epsilon_0 :$$

As such, now consider the first guarantee, where the goal will be to establish that $L(H_m^0) \leq L^* + \epsilon_0$ for all large m . But note first that the quantity $R_{t-1} \leq 1$ as $m \geq 1$, which combined with

$$\inf_{C > 0} fL(A) : \epsilon = 2 \leq \inf_{C > 0} fL(A) : \epsilon = 2 ; k \leq k_1 \leq Cg ;$$

guarantees that, if $m > M_2$ (for some M_2), then $\inf_{k \leq k_{R_t}} L(A) \leq L^* + \epsilon_0/2$. Finally, the rest of the terms in the first guarantee are at most $\epsilon_0/2$ for $m > M_3$ (for some M_3) with the same m^{-2} failure probability. As such, similarly to before, $\sum_{m=1}^{\infty} \Pr(E_{m; \epsilon}) \leq \max\{M_1; M_2; M_3\}g + \epsilon_0/6 < 1$. \square

Appendix 6.H Bibliographic Notes

The consistency of AdaBoost was first analyzed under various regularization strategies. Most notably, the work of Blanchard et al. (2003) and Lugosi and Vayatis (2004) studied the solutions of penalized estimators; the former work in particular achieving excellent finite sample guarantees, with convex risk decaying roughly as $\mathcal{O}(m^{-1/2})$ (where m is the sample size), with improvements under various noise conditions. This work, however, did not demonstrate tractable algorithms to produce these estimators, which was a goal of the work by Zhang and Yu (2005);

namely, there it is shown that merely constraining the step size taken by an AdaBoost-style scheme (with a variety of losses) suffices to achieve a convex risk rate of roughly $O(m^{-1/4})$ (in the case of the logistic loss), which includes the effect of approximate solutions produced by the algorithm. As will be discussed later, the present chapter, in the nonseparable case, fits well with the development by Zhang and Yu (2005).

Two works give a consistency analysis of AdaBoost without any algorithmic modifications, under the condition that the algorithm is stopped after m^a iterations (with arbitrary $a \in (0; 1)$). The first such analysis, due to Bartlett and Traskin (2007), was focused on establishing consistency, and established the convex risk decays roughly as $O(m^{-a} \sqrt{\ln(m)})$; the analysis depends on a curvature lower bound, which follows from a lower bound on the convex risk since the exponential loss is equal to its derivative. This derivative structure is of course not present with the logistic loss, and the analysis in this chapter must find another way. A streamlined consistency analysis of AdaBoost appears in the textbook of Schapire and Freund (2012, Theorem 12.2), with a rate of roughly $O(m^{-1/9})$ (by choosing $a := 5/9$); the analysis is short and clean, but it is not clear how to decouple the exponential loss.

In the separable case, the analysis in this chapter relies upon ideas from weak learnability, just as with the original analysis of AdaBoost (under margin assumptions) (Freund and Schapire, 1997). The relaxed notion of margin here is very close to the quantity AvgMin_k as developed by Shalev-Shwartz and Singer (2008, Section 4.1); the main contrasting point is that this thesis chapter is concerned with statistical properties, and in particular how these relaxed margin properties behave under sampling. The optimization analysis in the separable case here shares ideas both with the original AdaBoost analysis (Freund and Schapire, 1997), but also with the literature on hard cores (Impagliazzo, 1995, Barak et al., 2009); one distinction is that the latter methods take the target weak learning rate as input, whereas here (and in general with adaptive boosting), it must be found by the algorithm. Interestingly, the loss function implicit in the boosting algorithm due to Impagliazzo (1995, Proof of Lemma 1), ℓ_{russ} as discussed in Section 1.3, achieves superior constants to the logistic loss in Theorem 6.3.2; nearly the same loss was presented and praised by Zhang (2004, see the definition at the end of Section 4.6).

As stated previously, the nonseparable case fits well with the scheme laid down by Zhang and Yu (2005), where the algorithm is modified to constrain step sizes. Indeed, the analysis here first establishes that the iterates are well-behaved with exactly the sorts of norm bounds needed

by the analysis of Zhang and Yu (2005) (compare for instance the summability conditions (Zhang and Yu, 2005, Equation (4)) with Lemma 6.4.5). In order to produce these results, the present work uses a dual optimum as a witness to the difficulty of the convex risk problem over the source distribution; this technique follows structural properties extrapolating from the finite setting in Chapter 2. The corresponding statistical analysis in Chapter 3 appears statistically unstable, and the subsequent analysis here follows a similar path to the one by Zhang and Yu (2005), with additional help from Bartlett and Traskin (2007). One interesting distinction between the present work and those by Bartlett and Traskin (2007) and Zhang and Yu (2005) is that the latter two require a more strenuous algorithm: the weak learner and step size selection must be performed simultaneously. Decoupling these does not appear to impact the rates, however, this distinction prevents those results from being directly invoked here, meaning they must instead be reworked.

Lastly, note that the translation between convex and classification risks follows standard results on classification calibration as first developed by Zhang (2004), and later extended by Bartlett et al. (2006).

Acknowledgements

This chapter is based on work by the dissertation author which appeared in the Conference on Learning Theory, 2013.

Bibliography

- Boaz Barak, Moritz Hardt, and Satyen Kale. The uniform hardcore lemma via approximate bregman projections. In *SODA*, pages 1193{1200, 2009.
- Peter L. Bartlett and Mikhail Traskin. AdaBoost is consistent. *Journal of Machine Learning Research* 8:2347{2368, 2007.
- Peter L. Bartlett, Michael I. Jordan, and Jon D. McAuliffe. Convexity, classification, and risk bounds. *Journal of the American Statistical Association*, 101(473):138{156, 2006.
- Adi Ben-Israel. Motzkin's transposition theorem, and the related theorems of Farkas, Gordan and Stiemke. In M. Hazewinkel, editor, *Encyclopaedia of Mathematics*, Supplement III 2002.
- Dimitri P. Bertsekas. *Nonlinear Programming*. Athena Scientific, 2 edition, 1999.
- Peter J. Bickel, Yaacov Ritov, and Alon Zakai. Some theory for generalized boosting algorithms. *Journal of Machine Learning Research* 7:705{732, 2006.
- Gilles Blanchard, Gabor Lugosi, and Nicolas Vayatis. On the rate of convergence of regularized boosting classifiers. *Journal of Machine Learning Research* 4:861{894, 2003.
- Jonathan Borwein and Adrian Lewis. *Convex Analysis and Nonlinear Optimization*. Springer Publishing Company, Incorporated, 2000.
- Jonathan M. Borwein and Qiji J. Zhu. *Techniques of Variational Analysis*. Springer, 1 edition, 2005.
- Stéphane Boucheron, Olivier Bousquet, and Gabor Lugosi. Theory of classification: a survey of recent advances. *ESAIM: Probability and Statistics*, 9:323{375, 2005a.
- Stéphane Boucheron, Olivier Bousquet, and Gabor Lugosi. Theory of classification: A survey of some recent advances. *ESAIM: Probability and Statistics*, 9:323{375, 2005b.
- Leo Breiman. Prediction games and arcing algorithms. *Neural Computation*, 11:1493{1517, October 1999.
- Leo Breiman. Some infinity theory for predictor ensembles, 2000. Berkeley statistics technical report 577.
- Rich Caruana and Alexandru Niculescu-Mizil. An empirical comparison of supervised learning algorithms. pages 161{168, 2006.

- Michael Collins, Robert E. Schapire, and Yoram Singer. Logistic regression, AdaBoost and Bregman distances. *Machine Learning*, 48(1-3):253{285, 2002.
- George B. Dantzig and Mukund N. Thapa. *Linear Programming 2: Theory and Extensions* Springer, 2003.
- Sanjoy Dasgupta and Phil Long. Boosting with diverse base classifiers. *IrCOLT*, 2003.
- L. Devroye, L. Györfi, and G. Lugosi. *A probabilistic theory of pattern recognition*. Springer, 1996.
- Nigel Druy and David Helmbold. Potential boosters? In *NIPS*, pages 258{264. MIT Press, 2000.
- Rick Durrett. *Probability: Theory and Examples* Cambridge University Press, 4 edition, 2010.
- Gerald B. Folland. *Real analysis: modern techniques and their applications* Wiley Interscience, 2 edition, 1999.
- Yoav Freund. Boosting a weak learning algorithm by majority. *Information and Computation*, 121(2):256{285, 1995.
- Yoav Freund and Robert E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *J. Comput. Syst. Sci.*, 55(1):119{139, 1997.
- Jerome Friedman, Trevor Hastie, and Robert Tibshirani. Additive logistic regression: a statistical view of boosting. *Annals of Statistics*, 28(2):337{407, 2000.
- Jerome H. Friedman. Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, 29:1189{1232, 2000.
- Oded Goldreich and Leonid Levin. A hard-core predicate for all one-way functions. *STOC*, pages 25{32, 1989.
- Christian Gourieroux and Alain Monfort. Asymptotic properties of the maximum likelihood estimator in dichotomous logit models. *Journal of Econometrics*, 17(1):83{97, 1981.
- Venkatesan Guruswami and Prasad Raghavendra. Hardness of learning halfspaces with noise. In *FOCS*, 2006.
- Jean-Baptiste Hiriart-Urruty and Claude Lemaréchal. *Fundamentals of Convex Analysis* Springer Publishing Company, Incorporated, 2001.
- Russell Impagliazzo. Hard-core distributions for somewhat hard problems. In *FOCS*, pages 538{545, 1995.
- Michael Kearns and Leslie Valiant. Cryptographic limitations on learning finite automata and boolean formulae. *STOC*, pages 433{444, 1989.
- Michael Kearns and Umesh Vazirani. *An introduction to computational learning theory*. MIT Press, 1994.
- Jyrki Kivinen and Manfred K. Warmuth. Boosting as entropy projection. In *COLT*, pages

- 134{144, 1999.
- Hidetoshi Komiya. Elementary proof for sion's minimax theorem. *Kodai Mathematical Journal*, 11(1):5{7, 1988.
- John L. Elert. Additive models, boosting, and inference for generalized divergences. *In Proc. 12th Annu. Conf. on Comput. Learning Theory*, pages 125{133. ACM Press, 1999.
- Guy Lebanon. Consistency of the maximum likelihood estimator, 2008. URL <http://www.cc.gatech.edu/~lebanon/notes/mleConsistency.pdf>.
- Phil Long and Rocco Servedio. Algorithms and hardness results for parallel large margin learning. *In NIPS*, 2011.
- Gabor Lugosi and Nicolas Vayatis. On the bayes-risk consistency of regularized boosting methods. *Annals of Statistics*, 32(1):30{55, 2004.
- Zhi-Quan Luo and Paul Tseng. Error bounds and convergence analysis of feasible descent methods: a general approach. *Annals of Operations Research* 46(1):157{178, 1993.
- Llew Mason, Jonathan Baxter, Peter L. Bartlett, and Marcus R. Frean. Functional gradient techniques for combining hypotheses. *In A.J. Smola, P.L. Bartlett, B. Schölkopf, and D. Schuurmans, editors, Advances in Large Margin Classifiers*, pages 221{246, Cambridge, MA, 2000. MIT Press.
- Indraneel Mukherjee, Cynthia Rudin, and Robert Schapire. The convergence rate of AdaBoost. *In COLT*, 2011.
- Yurii Nesterov. *Introductory Lectures on Convex Optimization: A Basic Course*. Springer, 1 edition, 2003.
- Jorge Nocedal and Stephen J. Wright. *Numerical optimization*. Springer, 2 edition, 2006.
- Gunnar Ratsch and Manfred Warmuth. Efficient margin maximizing with boosting. *Journal of Machine Learning Research* 6:2153{2175, 2005.
- Gunnar Ratsch and Manfred K. Warmuth. Maximizing the margin with boosting. *In COLT*, pages 334{350, 2002.
- Gunnar Ratsch, Sebastian Mika, and Manfred K. Warmuth. On the convergence of leveraging. *In NIPS*, pages 487{494, 2001.
- Sidney I. Resnick. *A Probability Path*. Birkhäuser, 5 edition, 1999.
- R. Tyrrell Rockafellar. *Convex Analysis*. Princeton University Press, 1970.
- Walter Rudin. *Functional Analysis*. McGraw-Hill Book Company, 1973.
- Robert E. Schapire. The strength of weak learnability. *Machine Learning*, 5:197{227, July 1990.
- Robert E. Schapire. The convergence rate of AdaBoost. *In COLT*, 2010.

- Robert E. Schapire and Yoav Freund. *Boosting: Foundations and Algorithms* MIT Press, 2012.
- Robert E. Schapire and Yoram Singer. Improved boosting algorithms using confidence-rated predictions. *Machine Learning*, 37(3):297{336, 1999.
- Robert E. Schapire, Yoav Freund, Peter Bartlett, and Wee Sun Lee. Boosting the margin: A new explanation for the effectiveness of voting methods. In *ICML* , pages 322{330, 1997.
- Shai Shalev-Shwartz. *Introduction to Machine Learning, Course Notes*. 2009.
- Shai Shalev-Shwartz and Yoram Singer. On the equivalence of weak learnability and linear separability: New relaxations and efficient boosting algorithms. In *COLT* , pages 311{322, 2008.
- Matus Telgarsky. Margins, shrinkage, and boosting. In *ICML* , 2013.
- Lieven Vandenberghe. *Ee 236c: Subgradient method*. 2013.
- Manfred K. Warmuth, Karen A. Glocer, and Gunnar Ratsch. Boosting algorithms for maximizing the soft margin. In *NIPS*, 2007.
- Constantin Zalinescu. *Convex analysis in general vector spaces* World scientific, 2002.
- Tong Zhang. Statistical behavior and consistency of classification methods based on convex risk minimization. *The Annals of Statistics*, 32:56{85, 2004.
- Tong Zhang and Bin Yu. Boosting with early stopping: Convergence and consistency *The Annals of Statistics*, 33:1538{1579, 2005.

Index

A , 5
 A_0 , 58
 $(A\lambda)_i$, 46
 $(A\lambda)_{x;y}$, 5
 A_+ , 58
 C , 17
 $D(H, \nu)$, 17
 D , 47
 D_A , 47
dual pairing, 22
 e_h , 11
 F , 104
Fenchel conjugate, 8, 22
 γ , 1, 15
 $\gamma(A, S)$, 52
 $\gamma(\nu)$, 16
 H , 3
 H , 3
hard core, 17
im, 47
 h , i , 22
ker, 47
 \cdot , 3
 L , 131
 L , 6
 \hat{L} , 6
 $L(\cdot; C)$, 103
 ℓ_{exp} , 6
 L_{fs} , 47
 L_{fso} , 71
 L_C , 100
 ℓ_{hinge} , 6
Lipschitz gradient, 34
 ℓ_{\log} , 6
 ℓ_{russ} , 6
 $\hat{\mu}$, 4
 μ , 4
 m , 4
 μ_P , 116
 P , 47
 ψ_A , 48
 R , 131
 R , 4
 ρ , 3, 11
 $S_D(H, \nu)$, 17

strong duality, 7

weak duality, 7

weak learning rate, 15

WOLFE, 39

Wolfe conditions, 36

Wolfe search, 12, 36

X , 3