

# Statistical Consistency of Finite-dimensional Unregularized Linear Classification

Matus Telgarsky\*

## Abstract

This manuscript studies statistical properties of linear classifiers obtained through minimization of an unregularized convex risk over a finite sample. Although the results are explicitly finite-dimensional, inputs may be passed through feature maps; in this way, in addition to treating the consistency of logistic regression, this analysis also handles boosting over a finite weak learning class with, for instance, the exponential, logistic, and hinge losses. In this finite-dimensional setting, it is still possible to fit arbitrary decision boundaries: scaling the complexity of the weak learning class with the sample size leads to the optimal classification risk almost surely.

## 1 Introduction

Binary linear classification operates as follows: obtain a new instance, determine a set of real-valued features, form their weighted combination, and output a label which is positive iff this combination is nonnegative. The interpretability, empirical performance, and theoretical depth of this scheme have all contributed to its continued popularity (Freund and Schapire, 1997, Friedman et al., 2000, Caruana and Niculescu-Mizil, 2006).

In order to obtain the coefficients in the above weighting, convex optimization is typically employed. Specifically, rather than just trying to pick the weighting which makes the fewest mistakes over a finite sample — which is computationally intractable — consider instead paying attention to the *amount* by which these combinations clear the zero threshold, a quantity called the *margin*. Applying a convex penalty to these margins yields a convex optimization procedure, specifically one which can be specialized into both logistic regression and AdaBoost.

Statistical analyses of this scheme predominately follow two paths. The first path is a parameter estimation approach; positive and negative instances are interpreted as drawn from a family of distributions, indexed by the combination weights above, and the convex scheme is performing a maximum likelihood search for these parameters (Friedman et al., 2000). This provides one way to analyze logistic regression, specifically the ability of the above convex optimization to recover these parameters; these analyses of course require such parameters to exist, and usually for the full problem to obey certain regularity conditions (Lebanon, 2008, Gouriéroux and Monfort, 1981).

The second approach is focused on the case of binary classification, with an interpretation of the data generation process taking a background role. Indeed, in this setting, optimal parameters may simply fail to exist (Schapire, 2010), and the convex optimization procedure can produce unboundedly large weightings. Analyses first focused on the separable case, showing that AdaBoost approximately maximizes *normalized* margins, and that this leads to good generalization (Schapire and Freund, in preparation, Chapter 5 and the references therein). It is historically

---

\*Department of Computer Science and Engineering, University of California, San Diego. Email: <mtelgars@cs.ucsd.edu>.

interesting that this setting, which entails the non-existence of the best parameters, is diametrically opposed to the parameter estimation setting above.

In order to produce a more general analysis, it was necessary to control the unbounded iterates. This has been achieved either implicitly through regularization (Blanchard et al., 2003), or explicitly with an early stopping rule (Bartlett and Traskin, 2007, Zhang and Yu, 2005, Schapire and Freund, in preparation). Those analyses which handle the case of AdaBoost (cf. the work of Bartlett and Traskin (2007) and Schapire and Freund (in preparation, Chapter 12)), are sensitive both to the choice of exponential loss, to the choice of minimization scheme, and to the choice of stopping condition.

The goal of this manuscript is to analyze the setting of minimizing an unregularized convex loss applied to a finite sample (i.e., just like logistic regression and AdaBoost), but for a large class of loss functions, and without any demands on the optimization algorithm beyond an ability to attain arbitrarily small error.

## 1.1 Contribution

In more detail, the primary characteristics of the presented analysis are as follows.

**Any minimization scheme.** The oracle producing approximate solutions to the convex problem can output iterates which have any norm; they must simply be close in objective value to the optimum. The intent of this choice is twofold: for practitioners, it means that focusing on minimizing this objective value suffices; for theorists, it means that the wild deviations caused by these unbounded norms are not actually an issue.

**Many convex losses.** The analysis applies to any convex loss which is positive at the origin, and zero in the limit. (Some results also require differentiability at the origin.) In particular, the analysis handles the popular choice of using the logistic loss, but also applies to the exponential and hinge losses. (For a discussion on the difficulties of generalizing from the exponential loss, please see the work of Bartlett and Traskin (2007, Section 4).)

The main limitation of the presented analysis is that the set of features, or *weak learners*, must be finite. This weakness can be circumvented in the setting of boosting, where the complexity of the feature set can increase with the availability of data; it will be shown that the popular choice of decision trees fit this regime nicely.

## 1.2 Outline

A summary of the manuscript, and its organization, are as follows. Briefly, primary notation and technical background appear in Section 2.

Section 3 presents an impossibility result, which forces the structure of subsequent content. Specifically, with no bound on the iterates, it is in general impossible to control the deviations between the *empirical convex risk* (the convex surrogate risk over the observed finite sample), and the *true convex risk* (the convex surrogate risk over the source distribution).

The solution is to break the input space into two pieces: a *hard core*, where there exists an imperfect yet optimal parameter vector, 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 4.

The hard core has direct entailments on the structure of the convex risk. Specifically, Section 5 establishes first that the true risk has quantifiable curvature over the hard core, and effectively zero error over the rest of the space. Additionally, 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 6, display a number of terms

which are properties of the true convex risk, and not simply opaque random variables derived from the sample. It is thus possible to control many such bounds together; the eventual consistency results, appearing in Section 7, simply combine the finite sample guarantees, which all share the same primary structural quantities, together with standard probability techniques. As discussed previously, in order to fit arbitrary decision boundaries, structural risk minimization is employed, and it is furthermore established that decision trees with a constraint on the location of splits meet the requisite structural risk minimization condition.

Note that all proofs, as well as some supporting technical material, appear in a variety of appendices.

## 2 Notation

**Definition 2.1.** Instances  $x \in \mathcal{X}$  will have associated labels  $y \in \mathcal{Y} = \{-1, +1\}$ .  $\mu$  will always denote a probability measure over  $\mathcal{X} \times \mathcal{Y}$ , with only occasional mention of the related  $\sigma$ -algebra.  $\diamond$

To achieve generality sufficient to treat boosting, instances will not be worked with directly, but instead through a family of feature maps, or weak learners.

**Definition 2.2.** Let  $\mathcal{H} = \{h_i\}_{i=1}^n$  denote a finite set of (measurable) functions  $\mathcal{H} \ni h : \mathcal{X} \rightarrow [-1, +1]$ . Call a pair  $(\mathcal{H}, \mu)$  a *linear classification problem*. For convenience, let  $H$  denote a (bounded) linear operator with elements of  $\mathcal{H}$  as abstract columns: given any weighting  $\lambda \in \mathbb{R}^n$ ,

$$H\lambda = \sum_{i=1}^n \lambda_i h_i.$$

For convenience, define related classes of functions

$$\begin{aligned} \text{span}(\mathcal{H}, b) &:= \{H\lambda : \lambda \in \mathbb{R}^n, \|\lambda\|_1 \leq b\}, \\ \text{span}(\mathcal{H}) &:= \bigcup_{b=1}^{\infty} \text{span}(\mathcal{H}, b) = \{H\lambda : \lambda \in \mathbb{R}^n\}. \end{aligned} \quad \diamond$$

The class  $\text{span}(\mathcal{H})$  will be the search space for linear classification; if for instance  $\mathcal{H}$  consists of projection maps, then this is the standard setting of linear regression, however in general it can be viewed as a boosting problem. That the range of the function family is fixed specifically to  $[-1, +1]$  is irrelevant, however compactness of this output space is used throughout.

**Definition 2.3.**  $\Phi$  contains all convex losses  $\phi$  which are positive at the origin, and satisfy  $\lim_{z \rightarrow -\infty} \phi(z) = 0$ .  $\diamond$

This manuscript makes the choice of writing losses as nondecreasing functions; in this notation, three examples are the exponential loss  $\exp(z)$ , logistic loss  $\ln(1 + \exp(z))$ , and hinge loss  $\max\{0, 1 + z\}$ . Some of the consistency results will also require the loss to be differentiable at the origin; this requirement, which is satisfied by the three preceding examples, will be explicitly stated.

**Definition 2.4.** Given a probability measure  $\mu$ , a loss  $\phi \in \Phi$ , a function class  $\mathcal{F}$ , and arbitrary element  $f \in \mathcal{F}$ , the corresponding risk functional, and optimal risk, are

$$\mathcal{R}_\phi(f) := \int \phi(-yf(x)) d\mu(x, y), \quad \mathcal{R}_\phi(\mathcal{F}) = \inf_{f \in \mathcal{F}} \mathcal{R}_\phi(f).$$

When a sample  $\mathcal{S} := \{(x_i, y_i)\}_{i=1}^m$  is provided, let  $\mathcal{R}_\phi^m$  denote the corresponding empirical risk, meaning the convex risk corresponding to the empirical measure  $\mu_m(C) := m^{-1} \sum_{i=1}^m \mathbb{1}((x_i, y_i) \in C)$ , thus  $\mathcal{R}_\phi^m(f) = m^{-1} \sum_i \phi(-y_i f(x_i))$ . Lastly, let  $\mathcal{L}$  denote the classification risk  $\mathcal{L}(y, y') := \mathbb{1}(y \neq y')$ , and overload the notation for risks so that

$$\mathcal{R}_\mathcal{L}(f) := \int \mathcal{L}(y, 2 \cdot \mathbb{1}(f(x) \geq 0) - 1) d\mu(x, y), \quad \mathcal{R}_\mathcal{L}(\mathcal{F}) = \inf_{f \in \mathcal{F}} \mathcal{R}_\mathcal{L}(f). \quad \diamond$$

Typically, some function class  $\mathcal{H}$ , a particular weighting  $\lambda \in \mathbb{R}^n$ , and perhaps a sample of size  $m$  will be available, and example relevant risks are  $\mathcal{R}_\phi(H\lambda)$ ,  $\mathcal{R}_L^m(H\lambda)$ ,  $\mathcal{R}_\phi(\text{span}(\mathcal{H}))$ .

**Definition 2.5.** The requirement placed on the minimization oracle is that, for any  $\mathcal{H}$ ,  $\phi \in \Phi$ , finite sample of size  $m$ , and *suboptimality*  $\rho > 0$ , the oracle can produce  $\lambda \in \mathbb{R}^n$  with  $\mathcal{R}_\phi^m(H\lambda) \leq \mathcal{R}_\phi^m(\text{span}(\mathcal{H})) + \rho$ .  $\diamond$

The theorems themselves will avoid any reliance on this oracle, and their guarantees will hold with any  $\rho$ -suboptimal  $\lambda$  as input; this manuscript is concerned with statistical properties of these predictors. However, note briefly that for many losses of interest, in particular the hinge, logistic, and exponential losses, oracles satisfying the above guarantee exist.

**Proposition 2.6** ((Nesterov, 2003, Telgarsky, 2012)). *Let a linear classification problem  $(\mathcal{H}, \mu)$ , finite sample of size  $m$ , and suboptimality  $\rho > 0$  be given. Suppose:*

1. *Either  $\phi$  is Lipschitz continuous, attains its infimum, and subgradient descent is employed;*
2. *Or  $\phi$  is in the convex cone generated by the logistic and exponential losses, and coordinate descent is employed (as in AdaBoost);*

*then  $\text{poly}(1/\rho)$  iterations suffice to produce a  $\rho$ -suboptimal iterate  $\lambda \in \mathbb{R}^n$ .*

(The proof, in Appendix D, is mostly a reduction to known results regarding subgradient and coordinate descent.)

Lastly, this manuscript adopts a form of event-defining notation common in probability theory.

**Definition 2.7.** Given a function  $f : A \rightarrow B$  and binary relation  $\sim$ , define  $[f \sim b] := \{a \in A : f(a) \sim b\}$ ; for example  $[f > 0] := \{a \in A : f(a) > 0\} = f^{-1}((0, \infty))$ . At times, the variables will also be provided, for instance  $[bf(a) > 0] = \{(a, b) \subseteq A \times B : bf(a) > 0\}$ .  $\diamond$

### 3 An impossibility result

The stated goal of allowing iterates to have unbounded norms is at odds with the task of bounding the convex risk  $\mathcal{R}_\phi$ .

**Proposition 3.1.** *There exists a linear classification problem  $(\mathcal{H}, \mu)$  with the following characteristics.*

1.  *$\mathcal{X}$  is the square  $[-1, +1]^2$ , and  $\mathcal{H}$  consists of the two projection maps.*
2.  *$\mu$  has countable support.*
3. *There exists a perfect separator, albeit with zero margin.*
4. *For any  $\phi \in \Phi$ ,  $\mathcal{R}_\phi(\text{span}(\mathcal{H})) = 0$ .*
5. *Let any finite sample  $\{(x_i, y_i)\}_{i=1}^m$ , any  $b > 0$ , and any  $\phi \in \Phi$  be given. Then there exists a maximum margin solution  $\hat{\lambda}$ , i.e., a solution satisfying*

$$\arg \min_{i \in [m]} \frac{y_i(H\hat{\lambda})(x_i)}{\|\hat{\lambda}\|_1} = \sup \left\{ \arg \min_{i \in [m]} y_i(H\lambda)(x_i) : \lambda \in \mathbb{R}^n, \|\lambda\|_1 = 1 \right\},$$

*which has  $\mathcal{R}_\phi(H\hat{\lambda}) \geq b$ .*

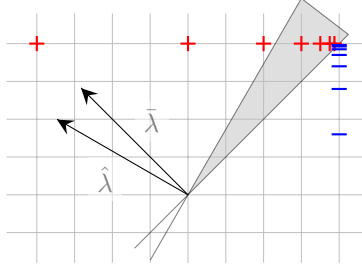


Figure 1: A bad example for unconstrained linear classification; please see Proposition 3.1.

A full proof is provided in Appendix E, but the mechanism is simple enough to appear as a picture. Consider the linear classification problem in Figure 1, which has positive (“+”) and negative (“-”) examples along two lines. Optimal solutions to  $\mathcal{R}_{\mathcal{L}}$  are of the form  $c\bar{\lambda}$ , where  $\bar{\lambda} = (-1, +1)$  and  $c > 0$  (note  $\lim_{c \uparrow \infty} \mathcal{R}_{\phi}(c\bar{\lambda}) = \mathcal{R}_{\phi}(\text{span}(\mathcal{H})) = 0$ ). Unfortunately, the positive and negative examples are staggered; as a result, for any sample, every max margin predictor  $\hat{\lambda}$ , 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  $c\bar{\lambda}$  as  $c \uparrow \infty$ , the convex risk  $\mathcal{R}_{\phi}$  may be made arbitrary large.

The statement of Proposition 3.1 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  $\mathcal{R}_{\mathcal{L}}$  is not only small, but its empirical counterpart  $\mathcal{R}_{\mathcal{L}}^m$  provides a reasonable estimate as  $m$  increases. Furthermore, if the distribution were adjusted slightly so that every  $\lambda \in \mathbb{R}^n$  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 the manuscript: linear classification problems are split into two pieces, one where optimization may produced unboundedly large iterates with small classification risk, and another piece where iterates are bounded thanks to the presence of difficult examples.

## 4 Hard cores

One way to split a linear classification problem into two pieces, one bounded and one unbounded, is to identify a *hard core* of very difficult instances. (Note, forms of the hard core have been previously used to study linear classification (Impagliazzo, 1995, Mukherjee et al., 2011, Telgarsky, 2012).)

**Definition 4.1.** Given a linear classification problem  $(\mathcal{H}, \mu)$ , let  $\mathcal{D}(\mathcal{H}, \mu)$  denote reweightings of  $\mu$  which decorrelate every regressor  $H\lambda$ ; that is,

$$\mathcal{D}(\mathcal{H}, \mu) := \left\{ p \in L^1(\mu) : p \geq 0, \forall \lambda \in \mathbb{R}^n \cdot \int y(H\lambda)(x)p(x, y)d\mu(x, y) = 0 \right\}.$$

Correspondingly,  $\mathcal{S}_{\mathcal{D}}(\mathcal{H}, \mu)$  tracks the supports of these weightings:

$$\mathcal{S}_{\mathcal{D}}(\mathcal{H}, \mu) := \{[p > 0] : p \in \mathcal{D}(\mathcal{H}, \mu)\}.$$

A *hard core*  $\mathcal{C} \subseteq \mathcal{X} \times \mathcal{Y}$  for  $(\mathcal{H}, \mu)$  is a maximal element of  $\mathcal{S}_{\mathcal{D}}(\mathcal{H}, \mu)$ ; that is,

$$\mathcal{C} \in \mathcal{S}_{\mathcal{D}}(\mathcal{H}, \mu) \quad \text{and} \quad \forall C \in \mathcal{S}_{\mathcal{D}}(\mathcal{H}, \mu) \cdot \mu(\mathcal{C} \setminus C) \geq 0 \text{ and } \mu(C \setminus \mathcal{C}) = 0.$$

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

Momentarily it will be established that hard cores split problems in the desired way; but first, note that hard cores actually exist.

**Theorem 4.2.** *Every linear classification problem  $(\mathcal{H}, \mu)$  has a hard core.*

To prove this, first observe that  $\mathcal{S}_{\mathcal{D}}(\mathcal{H}, \mu)$  is nonempty: it always contains  $\emptyset$ , with corresponding reweighting  $p(x, y) = 0$ . In order to produce a hard core, it does not suffice to simply union the contents of  $\mathcal{S}_{\mathcal{D}}(\mathcal{H}, \mu)$ , since the resulting set may fail to be measurable, and it is entirely unclear if a corresponding  $p \in \mathcal{D}(\mathcal{H}, \mu)$  can be found. Instead, the full proof in Appendix F constructs the hard core via an optimization, and the observation that  $\mathcal{S}_{\mathcal{D}}(\mathcal{H}, \mu)$  is closed under countable unions.

With the basic sanity check of existence out of the way, notice that hard cores achieve the goal laid out at the closing of Section 3. The proof, which is somewhat involved, appears in Appendix F.

**Theorem 4.3.** *Let problem  $(\mathcal{H}, \mu)$  and hard core  $\mathcal{C}$  be given. The following statements hold.*

1. *There exists a sequence  $\{\lambda_i\}_{i=1}^{\infty}$  with  $y(H\lambda_i)(x) = 0$  for  $\mu$ -a.e.  $(x, y) \in \mathcal{C}$ , and  $y'(H\lambda_i)(x') \uparrow \infty$  for  $\mu$ -a.e.  $(x', y') \in \mathcal{C}^c$ .*
2. *Every  $\lambda \in \mathbb{R}^n$  satisfies either  $\mu(\mathcal{C} \cap [y(H\lambda)(x) = 0]) = \mu(\mathcal{C})$  or  $\mu(\mathcal{C} \cap [y(H\lambda)(x) < 0]) > 0$ .*

The first property provides the existence of a sequence which is not only very good  $\mu$ -a.e. over  $\mathcal{C}^c$ , but furthermore does not impact the value of  $H\lambda$  over  $\mathcal{C}$ ; that is to say, this sequence can grow unboundedly, and have unboundedly positive margins over  $\mathcal{C}^c$ , while optimization over  $\mathcal{C}$  can effectively proceed independently. On the other hand,  $\mathcal{C}$  is difficult: every predictor is either abstaining  $\mu$ -a.e., or makes errors on a set of positive measure.

Finally, corresponding to the hard core, it is useful to specialize the definition of risk to consider regions.

**Definition 4.4.** Given a set  $C$  (typically  $\mathcal{C}$  or  $\mathcal{C}^c$ ), loss  $\phi$ , function class  $\mathcal{F}$ , and any  $f \in \mathcal{F}$ , define

$$\mathcal{R}_{\phi;C}(f) := \int \phi(-yf(x))\mathbb{1}((x, y) \in C)d\mu(x, y), \quad \mathcal{R}_{\phi;C}(\mathcal{F}) := \inf_{f \in \mathcal{F}} \mathcal{R}_{\phi;C}(f),$$

with analogous definitions for  $\mathcal{R}_{\phi;C}^m$ ,  $\mathcal{R}_{\mathcal{L};C}^m$ , etc.  $\diamond$

## 5 Hard cores and convex risk

The hard core imposes the following structure on  $\mathcal{R}_{\phi}$ . As provided by Theorem 4.3, there is a sequence which does arbitrarily well over  $\mathcal{C}^c$ , without impacting predictions over  $\mathcal{C}$ . On the other hand, since mistakes must occur over  $\mathcal{C}$ , convex losses within  $\Phi$  will be forced to avoid large predictors.

**Theorem 5.1.** *Let problem  $(\mathcal{H}, \mu)$ , hard core  $\mathcal{C}$ , and loss  $\phi \in \Phi$  be given.*

1. *There exists a sequence  $\{\lambda_i\}_{i=1}^{\infty}$  with  $y(H\lambda_i)(x) = 0$  for  $\mu$ -a.e.  $(x, y) \in \mathcal{C}$ , and  $\lim_{i \rightarrow \infty} \phi(-y'(H\lambda_i)(x')) = 0$  for  $\mu$ -a.e.  $(x', y') \in \mathcal{C}^c$ .*
2. *Let any  $\rho > 0$  be given. Then there exists  $c_{\rho} \in \mathbb{R}$  and a set  $N_{\rho}$  with  $\mu(N_{\rho}) = 0$  so that for every  $\lambda \in \mathbb{R}^n$  with  $\mathcal{R}_{\phi;\mathcal{C}}(H\lambda) \leq \mathcal{R}_{\phi;\mathcal{C}}(\text{span}(\mathcal{H})) + \rho$ , there exists a representation  $\lambda' \in \mathbb{R}^n$  with  $H\lambda = H\lambda'$  over  $\mathcal{C} \setminus N_{\rho}$ , and  $\|\lambda'\|_1 \leq c_{\rho}$ .*

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 5.2.** *Let problem  $(\mathcal{H}, \mu)$ , hard core  $\mathcal{C}$ , and loss  $\phi \in \Phi$  be given.*

1. *With probability 1 over the draw of a finite sample, there exists  $\lambda \in \mathbb{R}^n$  so that every  $(x_i, y_i) \in \mathcal{C}^c$  satisfies  $y_i(H\lambda)(x_i) > 0$ , and every  $(x'_i, y'_i) \in \mathcal{C}$  satisfies  $y'_i(H\lambda)(x'_i) = 0$ .*
2. *Given any empirical suboptimality  $\rho > 0$ , 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  $m_{\mathcal{C}}$ , the number of points landing in  $\mathcal{C}$ , has bound*

$$m_{\mathcal{C}} \geq c^2(\ln(n) + \ln(1/\delta)),$$

*then every  $\rho$ -suboptimal  $\lambda \in \mathbb{R}^n$  over the sample restricted to  $\mathcal{C}$ , meaning*

$$\mathcal{R}_{\phi; \mathcal{C}}^m(H\lambda) \leq \mathcal{R}_{\phi; \mathcal{C}}^m(\text{span}(\mathcal{H})) + \rho,$$

*has a representation  $\lambda'$  with  $\|\lambda'\|_1 \leq b$  which has  $H\lambda = H\lambda'$  over the sample restricted to  $\mathcal{C}$ , and in general  $\mu$ -a.e. over  $\mathcal{C}$ .*

## 6 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  $\mathcal{R}_{\phi}$ , but also the classification risk  $\mathcal{R}_{\mathcal{L}}$ . In order to make this correspondence, this manuscript relies on standard techniques due to Zhang (2004) and Bartlett, Jordan, and McAuliffe (2006).

**Definition 6.1.** Let  $\mathfrak{F}$  denote the set of measurable functions over  $\mathcal{X}$ . ◇

**Proposition 6.2** (Bartlett et al. (2006)). *Let any  $\phi \in \Phi$  be given with  $\phi$  differentiable at 0. There exists an associated function  $\psi : [0, 1] \rightarrow [0, \infty)$  with the following properties. First, for any probability measure  $\mu$  and any  $f : \mathcal{X} \rightarrow \mathbb{R}$ ,  $\psi(\mathcal{R}_{\mathcal{L}}(f) - \mathcal{R}_{\mathcal{L}}(\mathfrak{F})) \leq \mathcal{R}_{\phi}(f) - \mathcal{R}_{\phi}(\mathfrak{F})$ . Second, the inverse  $\psi^{-1}$  exists over  $[0, \infty)$ , and satisfies  $\psi^{-1}(r) \downarrow 0$  as  $r \downarrow 0$ .*

**Definition 6.3.** Given  $\phi \in \Phi$ , let  $\psi$ , called the  $\psi$ -transform, be as in Proposition 6.2. ◇

The general use of  $\psi$  is through its inverse, which provides

$$\begin{aligned} \mathcal{R}_{\mathcal{L}}(H\lambda) - \mathcal{R}_{\mathcal{L}}(\mathfrak{F}) &\leq \psi^{-1}(\mathcal{R}_{\phi}(H\lambda) - \mathcal{R}_{\phi}(\mathfrak{F})) \\ &= \psi^{-1}(\mathcal{R}_{\phi}(H\lambda) - \mathcal{R}_{\phi}(\text{span}(\mathcal{H})) + \mathcal{R}_{\phi}(\text{span}(\mathcal{H})) - \mathcal{R}_{\phi}(\mathfrak{F})). \end{aligned}$$

Although  $\psi^{-1}$  may be unwieldy, it is frequently easy to provide a useful upper bound. For instance, the exponential loss has  $\psi^{-1}(r) \leq 2\sqrt{r}$ , the logistic loss has  $\psi^{-1}(r) \leq 4\sqrt{r}$ , and the hinge loss has  $\psi^{-1}(r) = r$  (Zhang, 2004, Bartlett et al., 2006).

**Theorem 6.4.** *Let  $(\mathcal{H}, \mu)$ ,  $\mathcal{C}$ , and  $\phi \in \Phi$  be given. Let a suboptimality tolerance  $\rho > 0$  be given; results will depend on reals  $c > 0$  and  $b > 0$  determined by the preceding terms. The following statements simultaneously hold with any probability  $1 - \delta$  over the draw of  $m$  samples (with  $\delta' := \delta/8$  for convenience), and any weighting  $\lambda \in \mathbb{R}^n$  which is  $\epsilon$ -suboptimal (with  $\epsilon \leq \rho$ ) for the corresponding surrogate empirical risk problem, meaning  $\mathcal{R}_{\phi}^m(H\lambda) \leq \mathcal{R}_{\phi}^m(\text{span}(\mathcal{H})) + \epsilon$ .*

1. *Let  $m_{\mathcal{C}}$  and  $m_+$  respectively denote the number of samples falling into  $\mathcal{C}$  and  $\mathcal{C}^c$ . Then*

$$\begin{aligned} m_{\mathcal{C}} &\geq m \left( \mu(\mathcal{C}) - \sqrt{\ln(1/\delta')/(2m)} \right), \\ m_+ &\geq m \left( \mu(\mathcal{C}^c) - \sqrt{\ln(1/\delta')/(2m)} \right). \end{aligned}$$

2. The true classification risk over the unbounded portion,  $\mathcal{C}^c$ , has bound

$$\mathcal{R}_{\mathcal{L};\mathcal{C}^c}(H\lambda) \leq \frac{\epsilon}{\phi(0)} + 2\sqrt{\frac{2\epsilon(n \ln(2m_+ + 1) + \ln(4/\delta'))}{\phi(0)m_+}} + \frac{4(n \ln(2m_+ + 1) + \ln(4/\delta'))}{m_+}. \quad (6.5)$$

If moreover  $\epsilon < \phi(0)/m$ , then

$$\mathcal{R}_{\mathcal{L};\mathcal{C}^c}(H\lambda) \leq \frac{4(n \ln(2m_+ + 1) + \ln(4/\delta'))}{m_+}. \quad (6.6)$$

3. Suppose

$$m_{\mathcal{C}} \geq c^2(\ln(n) + \ln(6/\delta')).$$

The true surrogate risk over the unbounded portion has bound

$$\mathcal{R}_{\phi;\mathcal{C}}(H\lambda) - \mathcal{R}_{\phi;\mathcal{C}}(\text{span}(\mathcal{H})) \leq \epsilon + \frac{c\left(\sqrt{\ln(n)} + 4\sqrt{\ln(2/\delta')}\right)}{\sqrt{m_{\mathcal{C}}}}, \quad (6.7)$$

Additionally, if  $\phi$  is differentiable at 0, the classification risk has bound

$$\begin{aligned} \mathcal{R}_{\mathcal{L};\mathcal{C}}(H\lambda) - \mathcal{R}_{\mathcal{L};\mathcal{C}}(\mathfrak{F}) &\leq \psi^{-1}\left(\epsilon + \frac{c\left(\sqrt{\ln(n)} + 4\sqrt{\ln(2/\delta')}\right)}{\sqrt{m_{\mathcal{C}}}}\right) \\ &\quad + \mathcal{R}_{\phi;\mathcal{C}}(\text{span}(\mathcal{H})) - \mathcal{R}_{\phi;\mathcal{C}}(\mathfrak{F}). \end{aligned} \quad (6.8)$$

4. Suppose, for simplicity, that

$$m \geq \max\left\{2\ln(1/\delta')/\min\{\mu(\mathcal{C})^2, \mu(\mathcal{C}^c)^2\}, 2c^2(\ln(n) + \ln(1/\delta'))/\mu(\mathcal{C})\right\}$$

(where bounds are interpreted to hold trivially when denominators contain 0) and additionally that  $\epsilon < \phi(0)/m$  and  $\phi$  is differentiable at 0. Then the true classification risk of the full problem has bound

$$\begin{aligned} \mathcal{R}_{\mathcal{L}}(H\lambda) - \mathcal{R}_{\mathcal{L}}(\mathfrak{F}) &\leq \psi^{-1}\left(\epsilon + \frac{c\sqrt{2}\left(\sqrt{\ln(n)} + 4\sqrt{\ln(2/\delta')}\right)}{\sqrt{m\mu(\mathcal{C})}}\right) \\ &\quad + \mathcal{R}_{\phi;\mathcal{C}}(\text{span}(\mathcal{H})) - \mathcal{R}_{\phi;\mathcal{C}}(\mathfrak{F}) \\ &\quad + \frac{8(n \ln(m\mu(\mathcal{C}^c) + 1) + \ln(4/\delta'))}{m\mu(\mathcal{C}^c)}. \end{aligned}$$

## 7 Consistency

In order for the predictors to converge to the best choice, near-optimal choices must be available. Correspondingly, the first consistency result makes a strong assumption about the function class, albeit one which may be found in many treatments of the consistency of boosting (cf. the work of Bartlett and Traskin (2007) and Schapire and Freund (in preparation, Chapter 12)).



**Theorem 7.1.** Let  $(\mathcal{H}, \mu)$  and  $\phi \in \Phi$  be given with  $\phi$  differentiable at 0. Suppose  $\mathcal{R}_\phi(\text{span}(\mathcal{H})) = \mathcal{R}_\phi(\mathfrak{F})$ . Then there exists a sequence of sample sizes  $\{m_i\}_{i=1}^\infty \uparrow \infty$ , and empirical suboptimality tolerances  $\{\epsilon_i\}_{i=1}^\infty \downarrow 0$ , so that every sequence of  $\epsilon_i$ -suboptimal weightings  $\{\lambda_i\}_{i=1}^\infty$  (i.e.,  $\mathcal{R}_\phi^{m_i}(H\lambda_i) \leq \epsilon_i + \mathcal{R}_\phi^{m_i}(\text{span}(\mathcal{H}))$ ) satisfies  $\mathcal{R}_\mathcal{L}(H\lambda_i) \rightarrow \mathcal{R}_\mathcal{L}(\mathfrak{F})$  almost surely.

This additional assumption is hard to justify in the presence of only finitely many hypotheses. To mitigate this, this manuscript follows an approach remarked upon by Schapire and Freund (in preparation, Chapter 12): to consider an increasing sequence of classes which asymptotically grant the desired expressiveness property.

**Definition 7.2.** Let a probability measure  $\mu$  be given. A family of finite hypothesis classes  $\{\mathcal{H}_i\}_{i=1}^\infty$  is called a *linear structural risk minimization family* for  $\mu$ , or simply L-SRM family, if for any  $\phi \in \Phi$  and tolerance  $\epsilon > 0$ , there exists  $j$  so that  $\mathcal{R}_\phi(\text{span}(\mathcal{H}_j)) < \mathcal{R}_\phi(\mathfrak{F}) + \epsilon$ .  $\diamond$

The significance of this definition will be clear momentarily, as it grants a stronger consistency result. But first notice that straightforward classes satisfy the L-SRM condition.

**Proposition 7.3.** Suppose  $\mathcal{X} = \mathbb{R}^d$ , and let a probability measure  $\mu$  be given where  $\mu_\mathcal{X}$ , the marginal over  $\mathcal{X}$ , is a Borel probability measure. Let  $\mathcal{H}_i$  denote the collection of decision trees with axis aligned splits with thresholds taken from  $\{-i, -i + 1/i, \dots, i - 1/i, i\}$ . Then  $\{\mathcal{H}_i\}_{i=1}^\infty$  is an L-SRM family.

Proving this fact, as with many classical universal approximation theorems (Kolmogorov, 1957, Cybenko, 1989), relies on basic properties of continuous functions over compact sets. In order to reduce to this scenario from the general scenario of measurable functions  $\mathfrak{F}$ , Lusin's Theorem is employed, just as with similar results due to Zhang (2004, Section 4).

Now that the existence of reasonable L-SRM families is established, note the corresponding consistency result.

**Theorem 7.4.** Let probability measure  $\mu$  and loss  $\phi \in \Phi$  be given with  $\phi$  differentiable at 0, as well as an L-SRM  $\{\mathcal{H}_i\}_{i=1}^\infty$  for  $\mu$ . Then there exists a sequence of sample sizes  $\{m_i\}_{i=1}^\infty$ , a subsequence of classes  $\{\mathcal{H}_{j_i}\}_{i=1}^\infty$ , and suboptimality tolerances  $\{\epsilon_i\}_{i=1}^\infty$ , so that the every sequence of regressors  $\{H_{j_i}\lambda_i\}_{i=1}^\infty$   $\epsilon_i$ -suboptimal for the corresponding empirical problem satisfies  $\mathcal{R}_\mathcal{L}(H_{j_i}\lambda_i) \rightarrow \mathcal{R}_\mathcal{L}(\mathfrak{F})$  almost surely.

This manuscript is basically saying that constraining learning at the level of the weak learning oracle is sufficient for consistency. Of course, it could be argued that it is more elegant to instead apply a regularizer to the objective function (with data-dependent parameter choice), and permit a powerful weak learning class of infinite size. But such a discussion is beyond the scope of this manuscript.

## References

- 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.
- Gilles Blanchard, Gábor Lugosi, and Nicolas Vayatis. On the rate of convergence of regularized boosting classifiers. *Journal of Machine Learning Research*, 4:861–894, 2003.
- Stéphane Boucheron, Olivier Bousquet, and Gábor Lugosi. Theory of classification: a survey of recent advances. *ESAIM: Probability and Statistics*, 9:323–375, 2005.
- Rich Caruana and Alexandru Niculescu-Mizil. An empirical comparison of supervised learning algorithms. pages 161–168, 2006.

- George Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals, and Systems*, 2:303–314, 1989.
- Gerald B. Folland. *Real analysis: modern techniques and their applications*. Wiley Interscience, 2 edition, 1999.
- 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.
- Christian Gourieroux and Alain Monfort. Asymptotic properties of the maximum likelihood estimator in dichotomous logit models. *Journal of Econometrics*, 17(1):83–97, 1981.
- 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 Umesh Vazirani. *An introduction to computational learning theory*. MIT Press, 1994.
- Andrei N. Kolmogorov. On the representation of continuous functions of several variables as superpositions of continuous functions of one variable and addition. *Dokl. Acad. Nauk SSSR*, 114(5): 953–956, 1957. Translation to English: V. M. Volosov.
- Guy Lebanon. Consistency of the maximum likelihood estimator, 2008. URL <http://www.cc.gatech.edu/~lebanon/notes/mleConsistency.pdf>.
- 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.
- R. Tyrrell Rockafellar. *Convex Analysis*. Princeton University Press, 1970.
- R. Tyrrell Rockafellar. Integrals which are convex functionals. 39:439–469, 1971.
- Walter Rudin. *Functional Analysis*. McGraw-Hill Book Company, 1973.
- Robert E. Schapire. The convergence rate of AdaBoost. In *COLT*, 2010.
- Robert E. Schapire and Yoav Freund. *Boosting: Foundations and Algorithms*. MIT Press, in preparation.
- Matus Telgarsky. A primal-dual convergence analysis of boosting. *JMLR*, 13:561–606, 2012.
- Constantin Zălinescu. *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.

## A Technical Preliminaries

**Lemma A.1.** *Let any  $\phi \in \Phi$  be given. Then  $\phi$  is continuous, measurable, and nondecreasing. Subgradients exist everywhere, and satisfy  $\partial\phi(0) \subseteq \mathbb{R}_{++}$ . Lastly, the conjugate  $\phi^*$  satisfies  $\text{dom}(\phi^*) \subseteq \mathbb{R}_+$  and  $\phi^*(0) = 0$ .*

*Proof.* Since  $\phi$  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  $\phi$  has subgradients everywhere (Rockafellar, 1970, Theorem 23.4).

If  $\phi$  were not nondecreasing, there would exist  $x < y$  with  $\phi(x) > \phi(y)$ ; but that means every subgradient  $g \in \partial\phi(x)$  satisfies

$$\phi(y) \geq \phi(x) + g(y - x),$$

and thus  $g < 0$ . But then, for any  $z < x$ ,  $\phi(z) \geq \phi(x) + g(z - x)$ , which in particular contradicts  $\lim_{z \rightarrow -\infty} \phi(z) = 0$  (indeed, it implies  $\lim_{z \rightarrow -\infty} \phi(z) = \infty$ ), thus  $\phi$  is nondecreasing.

Next, since  $\phi$  is nondecreasing,  $\partial\phi \subseteq \mathbb{R}_+$ . However, since  $\phi(0) > 0$ , it follows that  $\partial\phi(0) \subset \mathbb{R}_{++}$ , since otherwise  $\lim_{z \rightarrow -\infty} \phi(z) = 0$  would be contradicted.

Turning to  $\phi^*$ , first note

$$\phi^*(0) = \sup_z 0 \cdot z - \phi(z) = 0.$$

Lastly, since  $\phi$  is nondecreasing, then for any  $g < 0$ ,

$$\phi^*(g) = \sup_z gz - \phi(z) \geq \sup_{z < 0} gz - \phi(z) = \infty.$$

That is to say,  $\text{dom}(\phi^*) \subseteq \mathbb{R}_+$ . □

**Proposition A.2.** *Let a linear classification problem  $(\mathcal{H}, \nu)$  and loss  $\phi \in \Phi$  be given. Then given a bound  $b$  on the  $l^1$  norm of considered predictors, there exists  $c \geq \phi(b)$  so that, for any  $\delta > 0$ , with probability at least  $1 - \delta$  over the draw of  $m$  points from  $\nu$ , every  $\lambda \in \mathbb{R}^n$  with  $\|\lambda\|_1 \leq b$  satisfies*

$$|\mathcal{R}_\phi(H\lambda) - \mathcal{R}_\phi^m(H\lambda)| \leq \frac{c \left( \sqrt{\ln(n)} + \sqrt{\ln(2/\delta)} \right)}{\sqrt{m}}.$$

*Proof.* Let bound  $b$  and loss  $\phi \in \Phi$  be given. Define a truncation

$$\hat{\phi}(z) := \begin{cases} \phi(z) & \text{when } z \leq b, \\ \phi(b) & \text{otherwise.} \end{cases}$$

Since  $\phi$  is nondecreasing (cf. Lemma A.1),  $\hat{\phi}(z) \leq \phi(b)$ , and furthermore  $\hat{\phi}$  is Lipschitz with a constant that may be measured at  $b$ ; indeed, since  $\phi$  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  $z_1 \leq z_2$ ,

$$\begin{aligned} |\phi(z_2) - \phi(z_1)| &= \phi(z_2) - \phi(z_1) \\ &\leq \sup \{ \phi(z_2) - (\phi(z_2) + \langle g_2, z_1 - z_2 \rangle) : g_2 \in \partial\phi(z_2) \} \\ &= |z_2 - z_1| \sup \{ |g_2| : g_2 \in \partial\phi(z_2) \} \\ &< \infty; \end{aligned}$$

correspondingly, set a Lipschitz constant  $L_\phi := \sup \{ |g| : g \in \partial\phi(b) \}$ .

Note that for every  $f \in \text{span}(\mathcal{H}, b)$ ,  $\sup_{x \in \mathcal{X}} |f(x)| \leq b$ , and thus  $\mathcal{R}_\phi(f) = \mathcal{R}_{\hat{\phi}}(f)$ . Lastly, the desired constant  $c$ , which does not depend on  $\delta$ ,  $n$ , or  $m$ , will be  $c := \max\{2L_\phi b\sqrt{2}, \phi(b)\}$ .

Now let a sample of size  $m$  be given, and let  $R_m(\text{span}(\mathcal{H}, b))$  denote the Rademacher complexity of  $\text{span}(\mathcal{H}, b)$ . By properties of Rademacher complexity and a few appeals to McDiarmid's inequality, (Boucheron et al., 2005, Theorem 3.1, and the proof of Theorem 4.1), with probability at least  $1 - \delta$  over the draw of this sample,

$$\begin{aligned} \sup_{\|\lambda\|_1 \leq b} |\mathcal{R}_\phi(H\lambda) - \mathcal{R}_\phi^m(H\lambda)| &= \sup_{\|\lambda\|_1 \leq b} |\mathcal{R}_\phi(H\lambda) - \mathcal{R}_\phi^m(H\lambda)| \\ &\leq 2L_\phi R_m(\text{span}(\mathcal{H}, b)) + \sqrt{\frac{2 \ln(2/\delta)}{m}}. \end{aligned} \quad (\text{A.3})$$

Next, by  $R_m(\text{span}(\mathcal{H}, b)) = bR_m(\text{span}(\mathcal{H}, 1)) = bR_m(\mathcal{H})$  and an appeal to Massart's Finite Lemma (Boucheron et al., 2005, Theorem 3.3)

$$R_m(\text{span}(\mathcal{H}, b)) \leq \sqrt{\frac{2 \ln(n)}{m}}.$$

Plugging this into eq. (A.3) and recalling the choice  $c = \max\{2L_\phi b\sqrt{2}, \phi(b)\}$ , the result follows.  $\square$

**Lemma A.4.** *Let  $S \subset \mathbb{R}$  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  $z$ :

$$y = x \left( \frac{z - y}{z - x} \right) + z \left( \frac{y - x}{z - x} \right).$$

By convexity and  $f(y) > f(x)$ ,

$$\begin{aligned} f(y) \left( \frac{z - y}{z - x} \right) + f(z) \left( \frac{y - x}{z - x} \right) &> f(x) \left( \frac{z - y}{z - x} \right) + f(z) \left( \frac{y - x}{z - x} \right) \\ &\geq f \left( x \left( \frac{z - y}{z - x} \right) + z \left( \frac{y - x}{z - x} \right) \right) \\ &= f(y). \end{aligned}$$

Rearranging and using  $x < y$ , it follows that  $f(y) < f(z)$ .  $\square$

## B Convexity properties of $\mathcal{R}_\phi$

**Lemma B.1.** *Let finite measure  $\nu$  and  $\phi \in \Phi$  be given. Then the function*

$$L^\infty(\nu) \ni q \mapsto \int \phi(q) \in \mathbb{R}$$

*is well-defined, convex, and lower semi-continuous. Next,  $(L^\infty(\nu))^*$  can be written as the direct sum of two spaces, one being  $L^1(\nu)$ ; for any  $p \in (L^\infty(\nu))^*$ , let  $p_1 + p_2$  be the corresponding decomposition (with  $p_1 \in L^1(\nu)$ ). With this notation,  $\int qp_2 = 0$  for any  $q \in L^\infty(\nu)$ ; furthermore, the Fenchel conjugate to the above map is*

$$(L^\infty(\nu))^* \ni p \mapsto \int \phi^*(p_1),$$

*which is again well-defined, convex, and lower semi-continuous. Lastly, the subdifferential set to the first map may be obtained by simply passing the subdifferential operator through the integral,*

$$\partial \left( \int \phi \right) (q) = \{p \in (L^\infty(\nu))^* : p_1 \in \partial \phi(q) \text{ } \nu\text{-a.e.}\}.$$

*Proof.* The proof will proceed with heavy reliance upon results due to Rockafellar (1971). To start, note that  $\phi$ , being convex and continuous (cf. Lemma A.1), is a *normal convex integrand* (Rockafellar, 1971, Lemma 1).

Let  $Z : \mathcal{X} \rightarrow \mathbb{R}$  denote the zero map, i.e.  $Z(x) = 0$  everywhere. Note that  $\phi \circ Z \in L^1(\nu)$ , and similarly  $\phi^* \circ Z \in L^1(\nu)$  (since  $\phi(0) = 0$ ; cf. Lemma A.1); these facts provide the conjugacy formula

$$\left( \int \phi \right)^* (p) = \int \phi^*(p_1) + \sup \left\{ p_2(q) : q \in L^\infty(\nu), \int \phi(q) < \infty \right\}, \quad (\text{B.2})$$

where the decomposition  $p = p_1 + p_2$  is as in the lemma statement (Rockafellar, 1971, Theorem 1).

Next, notice that  $\text{dom}(\int \phi) = L^\infty(\nu)$ ; in particular, given any  $q \in L^\infty(\nu)$ ,

$$\int \phi(q) \leq \int \phi(\|q\|_\infty) = \phi(\|q\|_\infty)\nu(\mathcal{X}, \mathcal{Y}) < \infty.$$

As such, consider an arbitrary  $p_2$  and  $q \in L^\infty(\nu)$ . Since  $p$  is a continuous linear functional on  $L^\infty(\nu)$ , then so is  $p_2$  (otherwise the formula  $p = p_1 + p_2$  would not make sense). Next, as stated by Rockafellar (1971, introduction to Section 2), it is possible to choose sets  $S_k$  with  $\nu(S_k^c) < 1/k$ , and  $p_2(q) = 0$  over every  $S_k$  and  $q \in L^1(\nu)$ . Now define  $U_k = \cup_{i \leq k} S_i$ . By continuity of measures from below (Folland, 1999, Theorem 1.8c),  $\nu(U_k) \uparrow \nu(\mathcal{X} \times \mathcal{Y})$ . As such, by the dominated convergence theorem (Folland, 1999, Theorem 2.25), and setting  $U_0 = \emptyset$ ,

$$\begin{aligned} \int p_2 q &= \int_{\cup_{k=1}^\infty U_k} p_2 q \\ &= \sum_{k=1}^\infty \int_{U_k \setminus U_{k-1}} p_2 q \\ &= 0. \end{aligned}$$

That is to say, the supremum term in eq. (B.2) is simply zero; plugging this back into eq. (B.2), the desired conjugacy relation follows. Note that the same result, due to Rockafellar (1971, Theorem 1), provides the integrals are well-defined, and moreover that the pair of conjugate functions are both convex and lower semi-continuous (as a consequence of being mutually conjugate). Lastly, the above derivation has established that  $\int \phi$  is finite over  $L^\infty(\nu)$ , but it is possible that  $\int \phi^*$  is infinite, even over  $L^1(\nu)$  (i.e., and not just over  $(L^\infty(\nu))^*$ ).

For the subdifferential relation, a related result by Rockafellar (1971, Corollary 1A) provides that  $(L^\infty(\nu))^* \ni p \in \partial(\int \phi)(q)$  (for some  $q \in L^\infty(\nu)$ ) precisely when  $p_1 \in \partial\phi(q)$   $\nu$ -a.e., and the supremum in eq. (B.2) is attained for  $p_2$  at  $q$ . It was already established that the supremum is always zero, as is  $p_2(q)$ , and the result follows.  $\square$

**Corollary B.3.** *Let a finite measure  $\nu$  and  $\phi \in \Phi$  be given. The function*

$$\mathbb{R}^n \ni \lambda \quad \mapsto \quad \int \phi(-y(H\lambda)(x))d\nu(x, y) \in \mathbb{R}$$

*is convex and continuous.*

*Proof.* Note that

$$\lambda \mapsto -y(H\lambda)x$$

is a bounded linear operator (and thus continuous), and the latter object, taken as a function over  $\mathcal{X} \times \mathcal{Y}$ , is within  $L^\infty(\nu)$ . Combined with the lower semi-continuity and convexity of  $\int \phi$  as per Lemma B.1, it follows that the map in question is convex and lower semi-continuous. Since it is finite everywhere, it is in fact continuous (Rockafellar, 1970, Corollary 7.2.2).  $\square$

**Lemma B.4.** *Let a linear classification problem  $(\mathcal{H}, \nu)$  and any  $\phi \in \Phi$  be given. Then*

$$\inf \left\{ \int \phi(-y(H\lambda)x) d\nu(x, y) : \lambda \in \mathbb{R}^n \right\} = \max \left\{ \int -\phi^*(p) : \max\{p, 0\} \in \mathcal{D}(\mathcal{H}, \nu) \right\},$$

where the max is taken element-wise. Furthermore, if a primal optimum  $\bar{\lambda}$  exists, then there is a  $\bar{p} \in \mathcal{D}(\mathcal{H}, \nu)$  with  $\bar{p}(x, y) \in \partial\phi(-y(H\bar{\lambda})x)$   $\nu$ -a.e.

*Proof.* For convenience, define the linear operator

$$(A\lambda)(x, y) := -y(H\lambda)x.$$

Note that  $A$  is a bounded linear operator, and furthermore has transpose

$$A^\top p := \sum_{i=1}^n \mathbf{e}_i \int -y h_i(x) p(x, y) d\nu(x, y)$$

(this follows by checking  $\langle A\lambda, p \rangle = \langle \lambda, A^\top p \rangle$  for arbitrary  $\lambda \in \mathbb{R}^n$  and  $p \in (L^\infty(\nu))^*$ , which entails the formula above provides the unique transpose (Rudin, 1973, Theorem 4.10).)

Consider the following two Fenchel problems:

$$\begin{aligned} p &:= \inf \left\{ \int \phi(A\lambda) + \langle 0, \lambda \rangle : \lambda \in \mathbb{R}^n \right\}, \\ d &:= \sup \left\{ - \int \phi^*(p_1) - \iota_{\{0\}}(A^\top p) : p \in (L^\infty(\nu))^* \right\}, \end{aligned}$$

where  $\iota_{\{0\}}$  is the indicator for the set  $\{0\}$ ,

$$\iota_{\{0\}}(\lambda) = \begin{cases} 0 & \text{when } \lambda = 0, \\ \infty & \text{otherwise,} \end{cases}$$

and is the conjugate to  $\langle 0, \cdot \rangle$ ; additionally,  $p_1$  is as discussed in the statement of Lemma B.1. To show  $p = d$  and thus prove the desired result, an appropriate Fenchel duality rule will be applied (Zălinescu, 2002, Corollary 2.8.5 using condition (vii)).

To start, note that  $\int \phi$  and  $\int \phi^*$  are conjugates, as provided by Lemma B.1. Next, also from Lemma B.1,  $\int \phi$  is finite everywhere over  $L^\infty(\nu)$ . As a result,

$$\text{Adom}(\langle 0, \cdot \rangle) - \text{dom}(\int \phi) = \text{Adom}(\langle 0, \cdot \rangle) - L^\infty(\nu) = L^\infty(\nu).$$

The significance of this fact is that it will act as the constraint qualification granting  $p = d$ .

Lastly,  $\mathbb{R}^n$  and  $L^\infty(\nu)$  are Banach and thus Fréchet spaces. As such, all conditions necessary for Fenchel duality are met (Zălinescu, 2002, Corollary 2.8.5 using condition (vii)), and it follows that  $p = d$  as desired, with attainment in the dual.

The next goal is to massage this duality expression into the one appearing in the lemma statement. To start, as provided by Lemma B.1,  $\int q p_2 = 0$  for any  $q \in L^\infty(\nu)$ , and in particular  $A^\top p_2 = 0$ ; consequently,  $p_2$  has no effect on either term in the dual objective, and the domain of the dual may be restricted to  $L^1(\nu)$ .

Next, Lemma A.1 grants  $\text{dom}(\phi^*) \subseteq \mathbb{R}_+$ , and so the domain of the dual problem may be safely restricted to  $p \geq 0$   $\nu$ -a.e. (since 0 is always dual feasible, and  $\nu([p < 0]) > 0$  entails an objective value of  $-\infty$ ). By the form of  $A^\top$ ,  $\iota_0(A^\top p)$  is finite iff

$$\int y h(x) p(x, y) d\nu(x, y) = 0$$

for all  $h$ ; it follows that  $\iota_0(A^\top p)$  is finite iff

$$\int (A\lambda)(x, y)p(x, y)d\nu(x, y) = 0$$

for all  $\lambda \in \mathbb{R}^n$ . Combining these facts, an equivalent form for the dual problem is

$$\max \left\{ - \int \phi^*(p) : \max\{p, 0\} \in \mathcal{D}(\mathcal{H}, \nu) \right\},$$

just as in the statement of the lemma.

Lastly, the Fenchel duality rule invoked above, as presented by Zălinescu (2002), also provides that a primal optimum  $\bar{\lambda}$  exists iff there is a  $p' \in (L^\infty(\nu))^*$  with  $-A^\top p' \in \partial(\langle 0, \cdot \rangle)(\bar{\lambda}) = 0$  and  $p' \in \partial(\int \phi)(A\bar{\lambda})$ . The first part simply states that  $\max\{p', 0\} \in \mathcal{D}(\mathcal{H}, \nu)$  as above. The second part, when combined with the subdifferential rule of Lemma B.1, gives  $p'_1 \in \partial\phi(A\bar{\lambda})$   $\nu$ -a.e. To obtain the desired statement, set  $\bar{p} := \max\{p'_1, 0\}$ , which satisfies all desired properties.  $\square$

## C Structure of $\mathcal{R}_\phi$ over $\mathcal{S}_\mathcal{D}(\mathcal{H}, \mu)$

The following theorem leads to a number of properties presented in Sections 4 and 5; it is easiest to prove them at once, as a ring of implications.

**Theorem C.1.** *Let a linear classification problem  $(\mathcal{H}, \mu)$  and a set  $D$  be given. The following statements are equivalent.*

1. *For every  $\lambda \in \mathbb{R}^n$ , either  $\mu(D \cap [y(H\lambda)x = 0]) = \mu(D)$  or  $\mu(D \cap [y(H\lambda)x < 0]) > 0$ .*
2. *Given any  $\rho$ , there exists a bound  $b$  and a null set  $N \subseteq \mathcal{X} \times \mathcal{Y}$  (i.e.,  $\mu(N) = 0$ ) so that for every  $\rho$ -suboptimal weighting  $\hat{\lambda}$  over  $D$ , meaning any weighting satisfying*

$$\mathcal{R}_{\phi; D}(H\hat{\lambda}) \leq \mathcal{R}_{\phi; D}(\text{span}(\mathcal{H})) + \rho,$$

*there exists  $\lambda'$  with  $\|\lambda'\|_1 \leq b$  and  $H\hat{\lambda} = H\lambda'$  over  $D \setminus N$ .*

3.  *$D \in \mathcal{S}_\mathcal{D}(\mathcal{H}, \mu)$ .*

The following structural lemma is crucial.

**Lemma C.2.** *Let  $(\mathcal{H}, \mu)$  and a set  $D$  be given. Define the set*

$$\mathcal{K} := \{\lambda \in \mathbb{R}^n : y(H\lambda)x = 0 \text{ for } \mu\text{-a.e. } (x, y) \in D\}.$$

*The following statements hold.*

1.  *$\mathcal{K}$  is a subspace.*
2. *There exists a set  $N$  with  $\mu(N) = 0$  so that, for any  $\lambda \in \mathbb{R}^n$ , the orthogonal projection  $\lambda \mapsto \lambda^\perp \in \mathcal{K}^\perp$  satisfies  $H\lambda = H\lambda^\perp$  everywhere over  $D \setminus N$ .*
3. *There exists a constant  $c > 0$  so that, for any  $\lambda \in \mathbb{R}^n$  with  $\mu(D \cap [H\lambda \neq 0]) > 0$ ,  $\|H\lambda\|_{L^\infty(\mu_D)} / \|\lambda^\perp\|_1 > c$ , where  $L^\infty(\mu_D)$  is the  $L^\infty$  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,  $\mathcal{K}$  is a subspace. Crucially, this means that  $\mathcal{K}^\perp$  is also a subspace, and the orthogonal projection  $\lambda \mapsto \lambda^\perp$  exists.

(Item 2) Given the subspace pair  $\mathcal{K}$  and  $\mathcal{K}^\perp$ , for any  $\lambda \in \mathbb{R}^n$ , there exists the decomposition  $\lambda \mapsto \lambda^\mathcal{K} + \lambda^\perp$ , where  $\lambda^\perp \in \mathcal{K}^\perp$ . By definition,  $H\lambda^\mathcal{K} = 0$   $\mu$ -a.e. over  $D$ , and thus  $H\lambda = H\lambda^\perp$   $\mu$ -a.e. over  $D$ .

Now let  $Q$  be any countable dense subset of  $\mathbb{R}^n$ . For each  $\lambda_i \in Q$ , define  $N_i := [H\lambda_i \neq H\lambda_i^\perp]$ , where the above provides  $\mu(N_i) = 0$ . Set  $N := \cup_i N_i$ , which is measurable since it is a countable union, and moreover  $\mu(N) = 0$  by  $\sigma$ -additivity. It will now be argued that the projections onto  $\mathcal{K}^\perp$  give equivalences over  $D \setminus N$ .

To this end, let any  $\lambda \in \mathbb{R}^n$ , any  $(x, y) \in D \setminus N$ , and any  $\tau > 0$  be given. Since  $Q$  is a countable dense subset of  $\mathbb{R}^n$ , there exists  $\lambda_i \in Q$  with  $\|\lambda_i - \lambda\|_1 \leq \tau/2$ . Now let  $P^\perp$  denote the orthogonal projection operator onto  $\mathcal{K}^\perp$ ; then

$$\begin{aligned} 0 &\leq |(H\lambda)(x) - (H\lambda^\perp)(x)| = |(H\lambda)(x) - (HP^\perp\lambda)(x)| \\ &= |(H(\lambda - \lambda_i + \lambda_i))(x) - (HP^\perp(\lambda - \lambda_i + \lambda_i))(x)| \\ &\leq |(H\lambda_i)(x) - (H\lambda_i^\perp)(x)| + |H(\lambda - \lambda_i)(x)| + |HP^\perp(\lambda - \lambda_i)(x)| \\ &\leq |0| + \|H\|_\infty \|\lambda - \lambda_i\|_1 + \|H\|_\infty \|P^\perp\|_\infty \|\lambda - \lambda_i\|_1 \\ &\leq 0 + \tau/2 + \tau/2 = \tau. \end{aligned}$$

Taking  $\tau \downarrow 0$ , it follows that  $H\lambda = H\lambda^\perp$  over  $D \setminus N$ .

(Item 3) For the final part, if every  $\lambda \in \mathbb{R}^n$  has  $\mu(D \cap [H\lambda \neq 0]) = 0$ , there is nothing to show, so suppose there exists  $\lambda \in \mathbb{R}^n$  with  $\mu_D([H\lambda \neq 0]) > 0$ . Consider the optimization problem

$$\inf \left\{ \frac{\|H\lambda\|_{L^\infty(\mu_D)}}{\|\lambda^\perp\|_1} : \lambda \in \mathbb{R}^n, \mu_D([H\lambda \neq 0]) > 0 \right\} = \inf \left\{ \|H\lambda\|_{L^\infty(\mu_D)} : \lambda \in \mathcal{K}^\perp, \|\lambda\|_1 = 1 \right\}.$$

The latter is a minimization of a continuous function over a nonempty compact set, and thus attains a minimizer  $\bar{\lambda}$ . But  $\bar{\lambda} \in \mathcal{K}^\perp$  and  $\|\bar{\lambda}\|_1 = 1$ , thus  $\|H\bar{\lambda}\|_{L^\infty(\mu_D)} > 0$ . The result follows with  $c := \|H\bar{\lambda}\|_{L^\infty(\mu_D)} > 0$ .  $\square$

*Proof of Theorem C.1.* (Item 1  $\implies$  Item 2.) Let  $\rho$  be given, and let  $N$  be the set, as provided by Lemma C.2, so that every  $\lambda \in \mathbb{R}^n$  has  $H\lambda = H\lambda^\perp$  everywhere on  $D \setminus N$ . Suppose contradictorily that the remainder of the desired statement is false; one way to say this is that there exists a sequence  $\{\lambda_i\}_{i=1}^\infty$  so that every equivalent representation over  $D \setminus N$  (i.e.,  $H\lambda_i = H\lambda_i'$  over this set) has  $\sup_i \|\lambda_i'\|_1 = \infty$ , but  $\mathcal{R}_{\phi;D}(H\lambda_i) \leq \mathcal{R}_{\phi;D}(\text{span}(\mathcal{H})) + \rho$ . (It can be taken without loss of generality that  $\lambda_i \neq 0$  for every  $i$ .)

To build the contradiction, choose representation  $\lambda_i^\perp$ , which satisfies  $H\lambda_i^\perp = H\lambda_i$  over  $D \setminus N$  via Lemma C.2. Note that  $\{\lambda_i^\perp / \|\lambda_i^\perp\|_1\}_{i=1}^\infty$  lies in a compact set (the unit  $l^1$  ball), and thus let  $\lambda_i^{(2)}$  be a subsequence with  $\lambda_i^{(2)} / \|\lambda_i^{(2)}\|_1 \rightarrow \bar{\lambda} \in \mathbb{R}^n$ . Since the assumed contradiction was that no representation is bounded,  $\lambda_i^{(2)}$  is unbounded; since there exists a  $c > 0$  with  $\|H\lambda_i^{(2)}\|_{L^\infty(\mu_D)} / \|\lambda_i^{(2)}\|_1 \geq c$  (cf. Lemma C.2), it follows by continuity of  $H$  and norms that  $\|H\bar{\lambda}\|_{L^\infty(\mu_D)} \geq c$ , and in particular  $\mu(D \cap [y(H\bar{\lambda})x \neq 0]) > 0$ .

By assumption (i.e., by Item 1), since  $\mu(D \cap [y(H\bar{\lambda})x \neq 0]) > 0$ , then  $\mu(D \cap [y(H\bar{\lambda})x < 0]) > 0$ ; for convenience, define the set  $P := [y(H\bar{\lambda})x < 0]$ . Thus, for any  $\lambda \in \mathbb{R}^n$ , taking any  $g \in \partial\phi(0)$



(note  $g > 0$  via Lemma A.1),

$$\begin{aligned}
& \lim_{t \rightarrow \infty} \frac{\int_D \phi(-y(H(\lambda + t\bar{\lambda}))(x)) - \int_D \phi(-y(H\lambda)(x))}{t} \\
& \geq \lim_{t \rightarrow \infty} \frac{\int \phi(-y(H(\lambda + t\bar{\lambda}))(x)) \mathbf{1}((x, y) \in D \cap P) - \int_D \phi(-y(H\lambda)(x))}{t} \\
& \geq \lim_{t \rightarrow \infty} \frac{\int (\phi(0) + g(-y(H(\lambda + t\bar{\lambda}))(x))) \mathbf{1}((x, y) \in D \cap P) - \int_D \phi(-y(H\lambda)(x))}{t} \\
& = g \int -y(H\bar{\lambda})(x) \mathbf{1}((x, y) \in D \cap P) \\
& \quad + \lim_{t \rightarrow \infty} \frac{\int (\phi(0) + g(-y(H\lambda)(x))) \mathbf{1}((x, y) \in D \cap P) - \int_D \phi(-y(H\lambda)(x))}{t} \\
& > 0.
\end{aligned} \tag{C.3}$$

The above statement shows that  $\int_D \phi$  eventually grows in direction  $H\bar{\lambda}$ , and in particular must exit the desired  $\rho$ -sublevel set

$$C_\rho := \{\lambda \in \mathbb{R}^n : \mathcal{R}_{\phi;D}(H\lambda) \leq \mathcal{R}_{\phi;D}(\text{span}(\mathcal{H})) + \rho\}.$$

To develop the contradiction, it will be shown that the construction of  $\bar{\lambda}$  indicates it should be in this sublevel set  $C_\rho$ ; the proof will be similar to one due to Hiriart-Urruty and Lemaréchal (2001, Proposition A.2.2.3).

Since  $\int \phi$  and  $\int_D \phi$  are convex and lower semi-continuous (cf. Lemma B.1), sublevel sets, in particular  $C_\rho$ , are closed convex sets. By construction of  $\bar{\lambda}$ ,

$$H\lambda_j + tH\bar{\lambda} = \lim_{i \rightarrow \infty} \left( \left(1 - \frac{t}{\|\lambda_i^{(2)}\|_1}\right) H\lambda_j + \frac{t}{\|\lambda_i^{(2)}\|_1} H\lambda_i^{(2)} \right) \in C_\rho.$$

This holds for all  $t > 0$ , but since  $H\bar{\lambda} \neq 0$ , eq. (C.3) forces  $H\lambda_i + tH\bar{\lambda}$  to leave any sublevel set (for sufficiently large  $t$ ), and in particular  $C_\rho$ , a contradiction.

(Item 2  $\implies$  Item 3.) Choose  $\phi := \exp \in \Phi$ , and a minimizing sequence  $\lambda_i^{(1)}$  for  $\mathcal{R}_{\phi;D}$ , meaning  $\mathcal{R}_{\phi;D}(H\lambda_i^{(1)}) \rightarrow \mathcal{R}_{\phi;D}(\text{span}(\mathcal{H}))$ . Choose any suboptimality  $\rho$ , and produce  $\lambda_i^{(2)}$  by removing all  $\lambda_j^{(1)}$  with  $\mathcal{R}_{\phi;D}(H\lambda_j^{(1)}) > \mathcal{R}_{\phi;D}(\text{span}(\mathcal{H})) + \rho$  (this procedure must be possible, since otherwise  $\{\lambda_i^{(1)}\}_{i=1}^\infty$  is not a minimizing sequence). By the assumed statement, there exists  $b > 0$  and a null set  $N$  so that each  $\lambda_i^{(2)}$  may be replaced with  $\lambda_i^{(3)}$ , where  $\|\lambda_i^{(3)}\|_1 \leq b$ , and  $H\lambda_i^{(2)} = H\lambda_i^{(3)}$  over  $D \setminus N$ , which in particular means  $\lambda_i^{(3)}$  is also a minimizing sequence. But this is now a minimizing sequence lying within a compact set, so, perhaps by passing to a subsequence  $\lambda_i^{(4)}$ , it has a limit  $\bar{\lambda} \in \mathbb{R}^n$ . Since  $\lambda \mapsto \int \phi(-y(H\lambda)x)$  is continuous (cf. Corollary B.3), it follows that  $\bar{\lambda}$  attains the desired infimal value.

Applying the duality relation in Lemma B.4 to  $\mathcal{R}_{\phi;D}$  (i.e., using the measure  $\nu = \mu_D$ , meaning  $\nu(S) = \mu(D \cap S)$  for any measurable set  $S$ ), the existence of a primal minimum  $\bar{\lambda}$  grants the existence of a dual maximum  $\bar{p}$  satisfying  $\bar{p} \in \mathcal{D}(\mathcal{H}, \nu)$ , and moreover

$$\bar{p}(x, y) \in \partial \phi(-y(H\bar{\lambda})x) = \exp(-y(H\bar{\lambda})x)$$

$\nu$ -a.e. As such, the choice  $p'(x, y) := \exp(-y(H\bar{\lambda})x)$  satisfies  $p' := \bar{p}$   $\nu$ -a.e., and thus  $p' \in \mathcal{D}(\mathcal{H}, \nu)$ ; moreover  $p' > 0$  everywhere, since  $\exp > 0$  everywhere.

This reweighting  $p'$  was with respect to  $\nu$ , so to finish, define  $p^*(x, y) := p'(x, y) \mathbf{1}((x, y) \in D)$ .

By construction,  $[p^* > 0] = D$ . Finally, given any  $\lambda \in \mathbb{R}^n$ ,

$$\begin{aligned} \int y(H\lambda)(x)p^*(x,y)d\mu(x,y) &= \int y(H\lambda)(x)p'(x,y)\mathbf{1}((x,y) \in D)d\mu(x,y) \\ &= \int y(H\lambda)(x)p'(x,y)d\mu_D(x,y) \\ &= 0. \end{aligned}$$

It follows that  $p^* \in \mathcal{D}(\mathcal{H}, \mu)$ , and that  $D \in \mathcal{S}_{\mathcal{D}}(\mathcal{H}, \mu)$ .

(Item 3  $\implies$  Item 1.) Let  $p \in \mathcal{D}(\mathcal{H}, \mu)$  with  $D = [p > 0]$  be given, and take any  $\lambda \in \mathbb{R}^n$  satisfying  $\mu(D \cap [y(H\lambda)x > 0]) > 0$ . But notice then, since  $p$  decorrelates  $H\lambda$ ,

$$\begin{aligned} 0 &= \int p(x,y)y(H\lambda)(x)d\mu(x,y) \\ &= \int_{D, y(H\lambda)(x) > 0} p(x,y)y(H\lambda)(x)d\mu(x,y) + \int_{D, y(H\lambda)(x) < 0} p(x,y)y(H\lambda)(x)d\mu(x,y). \end{aligned}$$

From this it follows that

$$-\int_{D, y(H\lambda)(x) < 0} p(x,y)y(H\lambda)(x)d\mu(x,y) = \int_{D, y(H\lambda)(x) > 0} p(x,y)y(H\lambda)(x)d\mu(x,y) > 0,$$

where the inequality follows from  $\mu(D \cap [y(H\lambda)(x) > 0]) > 0$  (Folland, 1999, Proposition 2.23(b)). The result follows.  $\square$

## D Deferred material from Section 2

In order to invoke standard results for gradient descent, this proof will use material from Section 5 to establish the existence of minimizers. Although those results appear later in the text, they do not in turn depend on the material here.

*Proof of Proposition 2.6.* Suppose  $\mathcal{H}$ , a sample of size  $m$ , and suboptimality  $\rho > 0$  are given as specified. Before proceeding, note briefly that the results invoked below — those demonstrating  $\mathcal{O}(\text{poly}(1/\rho))$  iterations suffice — neglect to provide a mechanism to stop the algorithms, and thus provide a proper oracle. But this may be accomplished by measuring duality gap, for instance by specializing the duality relation in Lemma B.4 to the empirical measure.

First suppose  $\phi$  is Lipschitz continuous, attains its infimum, and subgradient descent is employed. Notice that  $\mathcal{R}_\phi^m \circ H$  is also Lipschitz continuous (since  $H$  is a bounded linear operator), so if it can be shown that the infimum is attained, the standard analysis of subgradient descent may be applied, which in particular grants a  $\mathcal{O}(1/\rho^2)$  convergence rate when a step size of  $\mathcal{O}(1/\sqrt{t})$  is employed, where  $t$  indexes the iterations (Nesterov, 2003, Theorem 3.2.2 and subsequent discussion on step sizes). To finish, it must be shown that the infimum is attained.

To this end, let  $\mu_m$  be the empirical measure of the training sample, and let  $\mathcal{C}$  be a corresponding hard core. By Theorem 5.1, since  $\mu_m$  is now a discrete measure, a single weighting  $\lambda_0 \in \mathbb{R}^n$  can be extracted out with  $y(H\lambda_0)(x) > 0$  over  $\mathcal{C}^c$  and  $y(H\lambda_0)(x) = 0$  over  $\mathcal{C}$ . Also by Theorem 5.1, every 1-suboptimal predictor to  $\mathcal{R}_\phi^m$  has a representation which lies in a compact set; thus, minimizing sequence lies in the compact set, and a minimizer  $\bar{\lambda}_0$  exists. To finish, since  $\lim_{z \rightarrow -\infty} \phi(z) = 0$  and  $\phi$  attains its infimum, necessarily there is a  $b$  with  $\phi(z) = 0$  for  $z \leq b$ . As such, it follows that

$$\lambda' := \bar{\lambda} + \lambda_0 \left( \frac{z + \|H\bar{\lambda}\|_\infty}{\min\{|y_i(H\lambda_0)(x_i)| : (x_i, y_i) \in \mathcal{C}^c\}} \right)$$

is an optimum to the full problem. First, it is zero over  $\mathcal{C}^c$ , since for any  $(x, y) \in \mathcal{C}^c$ ,

$$\begin{aligned} y(H\lambda')(x) &= y(H\bar{\lambda})(x) + y(H\lambda_0)(x) \left( \frac{z + \|H\bar{\lambda}\|_\infty}{\min\{|y_i(H\lambda_0)(x_i)| : (x_i, y_i) \in \mathcal{C}^c\}} \right) \\ &\geq -\|H\bar{\lambda}\|_\infty + (z + \|H\bar{\lambda}\|_\infty), \end{aligned}$$

and the choice of  $z$  (i.e.,  $\phi(-y(H\lambda')(x)) = 0$ ). Next,  $\lambda'$  is equivalent to  $\bar{\lambda}$  over  $\mathcal{C}$ . Finally, if there exists some  $\lambda^*$  which achieves a lower objective value than  $\lambda'$ , necessarily it would be better than  $\bar{\lambda}$  over  $\mathcal{C}$ , contradicting optimality of  $\bar{\lambda}$ . In particular, the infimum is attained, and the proof for this choice of  $\phi$  is complete.

Now suppose that  $\phi$  is in the convex cone generated by the logistic and exponential losses; if it can be shown that  $\phi$  is within  $\mathbb{G}$ , a class of losses known to possess  $\mathcal{O}(1/\rho)$  convergence rates for boosting (Telgarsky, 2012, Definition 19, Theorem 21, Theorem 23, Theorem 27), then the result follows.

To this end, first notice that  $\mathbb{G}$  is a cone: given any  $c > 0$  and  $g \in \mathbb{G}$  with certifying constants  $\eta, \beta$ , then  $cg \in \mathbb{G}$  with the exact same constants. Since the exponential and logistic losses are within  $\mathbb{G}$  (Telgarsky, 2012, Remark 46), then so are all rescalings.

To finish, let  $\phi_1$  and  $\phi_2$  respectively denote the logistic and exponential losses, and let any  $c_1, c_2 > 0$  be given; if it can be shown that  $c_1\phi_1 + c_2\phi_2 \in \mathbb{G}$ , then combined with the earlier cases, the proof is complete. First note that

$$\sum_{i=1}^m (c_1\phi_1(x_i) + c_2\phi_2(x_i)) \leq m(c_1\phi_1(0) + c_2\phi_2(0))$$

implies

$$\forall i. x_i \leq \ln \left( \frac{m(c_1\phi_1(0) + c_2\phi_2(0))}{c_2} \right);$$

henceforth define  $c := m(c_1\phi_1(0) + c_2\phi_2(0))/c_2$ , and as per the definition of  $\mathbb{G}$ , the constants  $\eta$  and  $\beta$  must be established under the assumption  $x \leq \ln(c)$ .

For any  $x \in (-\infty, \ln(c)]$ , since  $\ln$  is convex, there is a secant lower bound

$$\ln(1 + e^x) \geq \left( \frac{\ln(1 + c) - 0}{c - 0} \right) e^x;$$

as usual, there is also the upper bound  $\ln(1 + e^x) \leq e^x$ .

As such, for any  $x \in (-\infty, c]$ , since  $\phi'_1(x) = e^x/(1 + e^x)$ ,

$$\frac{c_1\phi_1(x) + c_2\phi_2(x)}{c_1\phi'_1(x) + c_2\phi'_2(x)} = \frac{c_1 \ln(1 + e^x) + c_2 e^x}{c_1 e^x/(1 + e^x) + c_2 e^x} \leq \frac{e^x(c_1 + c_2)}{e^x(c_1/(1 + c) + c_2)},$$

and so it suffices to set  $\beta := (c_1 + c_2)/(c_1/(1 + c) + c_2)$ . Furthermore, since  $\phi''_1(x) = e^x/(1 + e^x)^2$ ,

$$\frac{c_1\phi''_1(x) + c_2\phi''_2(x)}{c_1\phi_1(x) + c_2\phi_2(x)} = \frac{c_1 e^x/(1 + e^x)^2 + c_2 e^x}{c_1 \ln(1 + e^x) + c_2 e^x} \leq \frac{e^x(c_1 + c_2)}{e^x(c_1 \ln(1 + c)/c + c_2)},$$

thus  $\eta := (c_1 + c_2)/(c_1 \ln(1 + c)/c + c_2)$  suffices.  $\square$

## E Deferred material from Section 3

*Proof of Proposition 3.1.* As stated in the proposition, set  $\mathcal{X} = [-1, +1]^2$ , and  $\mathcal{H}$  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:

$$p_i = \left[ 1 - 0.5 \cdot 4^{2^{-i}} \right], \quad \mu(p_i) = 2^{-i-1}.$$

Here are the negative instances:

$$n_i = \left[ \frac{1}{1-0.3 \cdot 4^{2-i}} \right], \quad \mu(n_i) = 2^{-i-1}.$$

Notice that  $\mu$  has countable support, and  $\mu(\mathcal{X}) = 1$ . Furthermore, the vector  $\bar{\lambda} = (-1, +1)$  is a perfect separator: given any positive example  $p_i$ ,  $(H\bar{\lambda})(p_i) > 0$ , and given negative example  $n_i$ ,  $(H\bar{\lambda})(n_i) < 0$ . Note however that, as required by the proposition statement, the margins go to zero. However, given any  $\phi \in \Phi$ , since  $\lim_{z \rightarrow -\infty} \phi(z) = 0$ ,

$$0 \leq \inf_{\lambda} \mathcal{R}_{\phi}(H\lambda) \leq \lim_{c \uparrow \infty} \int \phi(-y_i(Hc\bar{\lambda})(z_i)) d\mu(z_i, y_i) = 0.$$

The key property of this construction is that the positive and negative examples are staggered; this will cause max margin solutions to avoid  $\bar{\lambda}$ . As such, let any finite sample of size  $m$  be given. If all drawn examples have the same class  $y$ , then  $\hat{\lambda} = (1 - y, 1 + y)$  (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,  $\mathcal{R}_{\phi}(cH\hat{\lambda}) > b$ .

As such, henceforth suppose there is at least one positive example, and at least one negative example. Suppose  $j$  and  $k$  respectively denote a sampled positive point  $p_j$  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 determine solely by  $p_j$  and  $n_k$ . To obtain one of them, define

$$\lambda := \left[ \frac{-(1+(n_k)_2)/(2+(p_j)_1+(n_k)_2)}{(1+(p_j)_1)/(2+(p_j)_1+(n_k)_2)} \right].$$

To verify that this is a max margin solution, note that for any sampled (positive or negative) point  $z_i$  with label  $y_i \in \{-1, +1\}$ ,

$$y_i(H\lambda)z_i \geq (H\lambda)(p_j) = -(H\lambda)(n_k) = -\langle \lambda, n_k \rangle = \frac{(p_j)_1(n_k)_2}{2 + (p_j)_1 + (n_k)_2} > 0.$$

By construction, however,  $(p_j)_1 \neq (n_k)_2$ , meaning  $\lambda$  is not a rescaling of  $\bar{\lambda}$ . As such,  $\lambda$  is wrong for either all large  $p_i$  or  $n_i$ , and taking  $\hat{\lambda} = q\lambda$  with  $q$  large, it follows that  $\mathcal{R}_{\phi}(H\hat{\lambda}) > b$ .  $\square$

## F Deferred material from Section 4

Throughout this section, the following notation for measures will be employed

**Definition F.1.** Given a measure  $\mu$  and a set  $P$ , let  $\mu_P$  be the restriction of  $\mu$  to  $P$ : for any measurable set  $S$ ,  $\mu_P(S) = \mu(P \cap S)$ . Note also that  $d\mu_P(x, y) = \mathbb{1}((x, y) \in P)d\mu(x, y)$ .  $\diamond$

### F.1 Proof of Theorem 4.2

In order to establish the existence of hard cores, this section first establishes a few properties of  $\mathcal{D}(\mathcal{H}, \mu)$  and  $\mathcal{S}_{\mathcal{D}}(\mathcal{H}, \mu)$ .

**Lemma F.2.** *Given any  $\{c_i\}_{i=1}^{\infty}$  with  $c_i \geq 0$  and  $\{p_i\}_{i=1}^{\infty}$  with  $p_i \in \mathcal{D}(\mathcal{H}, \mu)$  and  $\sum_i c_i \|p_i\|_1 < \infty$ , the limit object  $p_{\infty} := \sum_i c_i p_i$  exists, and satisfies  $p_{\infty} \in \mathcal{D}(\mathcal{H}, \mu)$ .*

*Proof.* Let  $\{c_i\}_{i=1}^{\infty}$  and  $\{p_i\}_{i=1}^{\infty}$  be given as specified. First, by the monotone convergence theorem, the function  $p_{\infty} = \sum_i c_i p_i$  exists (i.e., all limits converge pointwise), is measurable, and satisfies  $\int p_{\infty} = \sum_i \int c_i p_i < \infty$ , meaning  $p_{\infty} \in L^1(\mu)$  (Folland, 1999, Theorem 2.15). Now let any  $\lambda \in \mathbb{R}^n$

be given; note that  $\sum_i \int |c_i p_i(H\lambda)| \leq \|H\lambda\|_\infty \sum_i \|c_i p_i\|_1 < \infty$ . Thanks to this, by the dominated convergence theorem (Folland, 1999, Theorem 2.25),

$$\begin{aligned} \int p_\infty(x, y) y(H\lambda) x d\mu(x, y) &= \int \sum_{i=1}^{\infty} c_i p_i(x, y) y(H\lambda) x d\mu(x, y) \\ &= \sum_{i=1}^{\infty} \int c_i p_i(x, y) y(H\lambda) x d\mu(x, y) \\ &= \sum_{i=1}^{\infty} c_i \int p_i(x, y) y(H\lambda) x d\mu(x, y) \\ &= 0. \end{aligned}$$

□

**Lemma F.3.**  $\mathcal{S}_{\mathcal{D}}(\mathcal{H}, \mu)$  is closed under countable unions.

*Proof.* Let any collection  $\{C_i\}_{i=1}^{\infty}$  with  $C_i \in \mathcal{S}_{\mathcal{D}}(\mathcal{H}, \mu)$  and corresponding weighting  $p_i \in \mathcal{D}(\mathcal{H}, \mu)$  be given. Define

$$C := \bigcup_{i=1}^{\infty} C_i \quad \text{and} \quad p := \sum_{i=1}^{\infty} \frac{p_i}{2^i \max\{1, \|p_i\|_1\}}.$$

By Lemma F.2,  $p$  exists and satisfies  $p \in \mathcal{D}(\mathcal{H}, \mu)$ . Note further that  $C = [p > 0]$ , and thus  $C \in \mathcal{S}_{\mathcal{D}}(\mathcal{H}, \mu)$ . □

*Proof of Theorem 4.2.* Consider the optimization problem

$$d := \sup\{\mu(C) : C \in \mathcal{S}_{\mathcal{D}}(\mathcal{H}, \mu)\}.$$

Since  $\mathcal{S}_{\mathcal{D}}$  is nonempty (always contains  $\emptyset$  corresponding to  $p = 0 \in \mathcal{D}(\mathcal{H}, \mu)$ ) and  $\mu(\mathcal{X} \times \mathcal{Y}) < \infty$ , the supremum is finite. Let  $\{C_i\}_{i=1}^{\infty}$  be a maximizing sequence, and define  $D_j := \cup_{i \leq j} C_i$  and  $D := \cup_{j=1}^{\infty} D_j = \cup_{i=1}^{\infty} C_i$ . By Lemma F.3,  $D_j \in \mathcal{S}_{\mathcal{D}}(\mathcal{H}, \mu)$  for every  $j$ , and since  $\mu(D_j) \geq \mu(C_j)$ , it follows that  $\{D_j\}_{j=1}^{\infty}$  must also be a maximizing sequence to the above supremum. Finally, since Lemma F.3 also grants  $D \in \mathcal{S}_{\mathcal{D}}(\mathcal{H}, \mu)$ , then by continuity of measures from below (Folland, 1999, Theorem 1.8(c)),

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

Since  $D \in \mathcal{S}_{\mathcal{D}}(\mathcal{H}, \mu)$  attains the supremum, it is a dual hard core. □

## F.2 Primal hard cores

In light of the duality relationship for  $\mathcal{R}_\phi$  (cf. Lemma B.4), the definition for hard cores, provided in Section 4, is tied to the convex dual to  $\mathcal{R}_\phi$ . Analogously, it is possible to define a primal form of hard cores, which will lead to a proof of Theorem 4.3.

**Definition F.4.** Define  $\mathcal{S}_{\mathcal{P}}(\mathcal{H}, \mu)$  to contain all sets  $C$  for which there exists a sequence  $\{\lambda_i\}_{i=1}^{\infty}$  satisfying the following properties.

1. Every  $\lambda_i$  and  $(x, y) \in C$  satisfies  $y(H\lambda_i)x = 0$ .
2. For  $\mu$ -almost-every  $(x, y)$  in  $C^c$ ,  $y(H\lambda_i)x \uparrow \infty$ .

A *primal hard core*  $\mathcal{P}$  is a minimal set within  $\mathcal{S}_{\mathcal{P}}(\mathcal{H}, \mu)$ :

$$\mathcal{P} \in \mathcal{S}_{\mathcal{P}}(\mathcal{H}, \mu) \quad \text{and} \quad \forall C \in \mathcal{S}_{\mathcal{P}}(\mathcal{H}, \mu) \bullet \mu(\mathcal{P} \setminus C) = 0 \wedge \mu(C \setminus \mathcal{P}) \geq 0. \quad \diamond$$

**Lemma F.5.**  $\mathcal{S}_{\mathcal{P}}(\mathcal{H}, \mu)$  is closed under countable intersections.

*Proof.* To start, note that  $\mathcal{S}_{\mathcal{P}}(\mathcal{H}, \mu)$  is closed under finite intersections as follows. Let  $\{C_i\}_{i=1}^p$  be given with corresponding sequences  $\{\lambda_j^{(i)}\}_{j=1}^\infty$ . Define  $C := \cap C_i$  and  $\lambda_j := \sum_i \lambda_j^{(i)}$ . By construction, for every  $(x, y) \in C$  and pair  $(i, j)$ ,  $y(H\lambda_j^{(i)})x = 0$ , and thus  $y(H\lambda_j)x = 0$ . Next, for each  $C_i$ , define  $C'_i \subseteq C_i^c$  with  $\mu(C'_i) = \mu(C_i^c)$  so that, for every  $(x, y) \in C'_i$ ,  $y(H\lambda_j^{(i)})x \uparrow \infty$ . Correspondingly, define  $C' := \cup_i C'_i$ , where  $\mu(C') = \mu(C^c)$ . Now let any  $(x, y) \in C'$  and any  $B > 0$  be given. For each  $i$ , there are two cases: either this is an area where  $y(H\lambda_j^{(i)})x \uparrow \infty$ , or  $y(H\lambda_j^{(i)})x = 0$ . In the first case, let  $T_i$  denote an integer, as granted by  $y(H\lambda_j^{(i)})x \uparrow \infty$ , so that for all  $j \geq T_i$ ,  $y(H\lambda_j^{(i)})x > B$ . For those  $i$  where  $(x, y) \notin C'_i$  (but still  $(x, y) \in C'$ ), due to the ruled out nullsets,  $y(H\lambda_j^{(i)})x = 0$ , safely set  $T_i = 0$ . To finish, taking  $T := \max_i T_i$ , it follows that for every  $j > T$ ,  $y(H\lambda_j)x > B$ , whereby it follows that  $y(H\lambda_j)x \uparrow \infty$  over  $C'$ , and thus over  $C^c$   $\mu$ -a.e.

Now let a countable family  $\{D_i\}_{i=1}^\infty$  be given, and define  $D = \cap_i D_i$ . Consider the optimization problem

$$p := \inf \left\{ \int \exp(-y(H\lambda)x) d\mu_{D^c}(x, y) : \lambda \in \mathbb{R}^n, \forall (x, y) \in D \bullet y(H\lambda)x = 0 \right\}.$$

Define  $E_j := \cap_{i \leq j} D_i$ , whereby  $D := \cap_j E_j$ . Since  $\mu(\mathcal{X} \times \mathcal{Y}) < \infty$ , by continuity of measures from above (Folland, 1999, Theorem 1.8(d)), for any  $\tau > 0$  there exists  $E_k$  with  $\mu(D) > \mu(E_k) - \tau$ . Since it was shown above that  $\mathcal{S}_{\mathcal{P}}(\mathcal{H}, \mu)$  is closed under finite intersections,  $E_k = \cap_{i \leq k} D_i \in \mathcal{S}_{\mathcal{P}}(\mathcal{H}, \mu)$ ; consequently, let  $\{\lambda_i\}_{i=1}^\infty$  to be a sequence of predictors certifying that  $E_k \in \mathcal{S}_{\mathcal{P}}(\mathcal{H}, \mu)$ , as according to the definition. It follows that

$$p \leq \lim_{i \rightarrow \infty} \int \exp(-y(H\lambda_i)x) d\mu_{D^c}(x, y) = 0 + \int \exp(0) \mu_{E_k \setminus D} = \mu(E_k) - \mu(D) < \tau.$$

Since  $\tau$  was arbitrary, it follows that  $p = 0$ .

As such, for any  $n \in \mathbb{Z}_{++}$ , choose  $\bar{\lambda}_n \in \mathbb{R}^n$  with  $y(H\bar{\lambda}_n)x = 0$  over  $D$  satisfying

$$\int \exp(-y(H\bar{\lambda}_n)x) d\mu_{D^c}(x, y) < 1/n^2.$$

By Markov's inequality, it follows that

$$\mu_{D^c}([\exp(-y(H\bar{\lambda}_n)x) \geq 1/n]) \leq n \int \exp(-y(H\bar{\lambda}_n)x) d\mu_{D^c}(x, y) < 1/n.$$

As such, by definition,  $\exp(-y(H\bar{\lambda}_n)x)$  converges in measure to the function  $\mathbf{1}((x, y) \in D)$ . Consequently, there exists a subsequence  $\lambda_i^*$  with  $\exp(-y(H\lambda_i^*)x) \rightarrow \mathbf{1}(D)$   $\mu$ -a.e. (Folland, 1999, Theorem 2.30). This is only possible if  $y(H\lambda_i^*)x \uparrow \infty$  for  $\mu$ -a.e  $(x, y) \in D^c$ , and the result follows, with  $\{\lambda_i^*\}_{i=1}^\infty$  as the certifying sequence for  $D$ , since every  $y(H\lambda_i^*)x = 0$  for  $(x, y) \in D$  by construction.  $\square$

**Theorem F.6.** *Every linear classification problem  $(\mathcal{H}, \mu)$  has a primal hard core.*

*Proof.* Consider the optimization problem

$$p := \inf \{\mu(C) : C \in \mathcal{S}_{\mathcal{P}}(\mathcal{H}, \mu)\}.$$

Since  $\mathcal{S}_{\mathcal{P}}$  is nonempty (it always contains  $\mathcal{X} \times \mathcal{Y}$  with certifying sequence  $\lambda_i = 0$  for every  $i$ ) and  $\mu$  is a finite nonnegative measure, the infimum is finite. Let  $\{C_i\}_{i=1}^\infty$  be a minimizing sequence, and define  $D_j := \cap_{i \leq j} C_i$  and  $D := \cap_{j=1}^\infty D_j = \cap_{i=1}^\infty C_i$ . By Lemma F.5,  $D_j \in \mathcal{S}_{\mathcal{P}}(\mathcal{H}, \mu)$  for every  $j$ , and since  $\mu(D_j) \leq \mu(C_j)$ , it follows that  $\{D_j\}_{j=1}^\infty$  must also be a minimizing sequence to the above infimum. Finally, since  $\mu$  is finite and Lemma F.5 also grants  $D \in \mathcal{S}_{\mathcal{P}}(\mathcal{H}, \mu)$ , then by continuity of measures from above (Folland, 1999, Theorem 1.8(d)),

$$\mu(D) = \lim_{j \rightarrow \infty} \mu(D_j) = p.$$

Since  $D \in \mathcal{S}_{\mathcal{P}}(\mathcal{H}, \mu)$  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 F.7.** *Let a linear classification problem  $(\mathcal{H}, \mu)$  be given, along with a hard core  $\mathcal{C}$ , as well as a primal hard core  $\mathcal{P}$ . Then  $\mathcal{C}$  and  $\mathcal{P}$  agree on all but a null set.*

The proof needs the following lemma.

**Lemma F.8.** *Let a linear classification problem  $(\mathcal{H}, \mu)$ ,  $C_1 \in \mathcal{S}_{\mathcal{P}}(\mathcal{H}, \mu)$ , as well as a  $\lambda_2 \in \mathbb{R}^n$  be given, with  $y(H\lambda_2)x \geq 0$  for  $(x, y) \in C_1$  (but potentially  $y(H\lambda_2)x < 0$  elsewhere). Then  $C_1 \setminus [y(H\lambda_2)x > 0] \in \mathcal{S}_{\mathcal{P}}(\mathcal{H}, \mu)$ .*

*Proof.* Let  $C_1, \lambda_2$  be given as specified. Let  $\{\lambda_i^{(1)}\}_{i=1}^\infty$  be a certifying sequence for  $C_1$ . Define  $P := [y(H\lambda_2)x > 0]$  and  $C_3 := C_1 \setminus P = C_1 \setminus [y(H\lambda_2)x > 0]$ .

Now let  $i \in \mathbb{Z}_{++}$  be arbitrary; the following steps will construct  $\lambda_i^{(4)}$ , a certifying sequence for  $C_3$ , meaning  $C_3 \in \mathcal{S}_{\mathcal{P}}(\mathcal{H}, \mu)$ .

First, let  $c$  be sufficiently large so that  $\lambda_i^{(2)} := c\lambda_2$  satisfies

$$\int \exp(-y(H\lambda_i^{(2)})x) \mu_P(x, y) < 1/i^2.$$

By Markov's inequality, it follows that

$$\mu_P([\exp(-y(H\bar{\lambda}_2)x) \geq 1/i]) \leq i \int \exp(-y(H\bar{\lambda}_2)x) \mu_P(x, y) < 1/i. \quad (\text{F.9})$$

Consequently define  $P_i := [y(H\bar{\lambda}_2)x > \ln(i)]$ , where the above statements show  $\mu(P_i) > \mu(P) - 1/i$ .

Next, since  $\exp(-y(H\lambda_i^{(1)})x) \rightarrow \mathbb{1}(C_1)$   $\mu$ -a.e. and  $\mu(\mathcal{X} \times \mathcal{Y}) < \infty$ , by Egoroff's Theorem (Folland, 1999, Theorem 2.33), this convergence is uniform over a subset  $S_i$  with  $\mu(S_i) > \mu(\mathcal{X}, \mathcal{Y}) - 1/i$ . In particular, there exists an integer  $T_i$  so that, for any  $(x, y) \in S_i \cap C_1$ ,

$$y(H\lambda_{T_i}^{(1)})x > \|\lambda_i^{(2)}\|_1 + \ln(i).$$

As such, define  $\lambda_i^{(3)} := \lambda_{T_i}^{(1)} + \lambda_i^{(2)}$ . First, for any  $(x, y) \in C_3$  and any  $i$ ,

$$y(H\lambda_i^{(3)})x = 0 = y(H\lambda_i^{(1)})x = y(H\lambda_2)x.$$

On the other hand, for any  $(x, y) \in S_i \cap C_1$ ,

$$\begin{aligned} y(H\lambda_i^{(3)})x &= y(H\lambda_{T_i}^{(1)})x + y(H\lambda_i^{(2)})x \\ &> \|\lambda_i^{(2)}\|_1 + \ln(i) - \|\lambda_i^{(2)}\|_1 = \ln(i). \end{aligned}$$

Lastly, as shown above, for any  $(x, y) \in P_i$ ,

$$y(H\lambda_i^{(3)})x = 0 + y(H\lambda_i^{(2)})x \geq \ln(i).$$

Combining the above facts,

$$\mu([|\exp(-y(H\lambda_i^{(3)})x) - \mathbb{1}((x, y) \in C_3)| \geq 1/i]) < \mu(C_1^c \setminus S_i) + \mu(P \setminus P_i) \leq 2/i.$$

It follows that  $\exp(-y(H\lambda_i^{(3)})x) \rightarrow \mathbb{1}((x, y) \in C_3)$  in measure, and thus there is a subsequence  $\{\lambda_i^{(4)}\}_{i=1}^\infty$  which converges to  $\mathbb{1}((x, y) \in C_3)$   $\mu$ -a.e. (Folland, 1999, Theorem 2.30). It follows that  $\{\lambda_i^{(4)}\}_{i=1}^\infty$  is the desired sequence certifying that  $C_3 \in \mathcal{S}_{\mathcal{P}}(\mathcal{H}, \mu)$ .  $\square$

*Proof of Theorem F.7.* If  $\mu(\mathcal{P} \setminus \mathcal{C}) > 0$ , then by the maximality of  $\mathcal{C}$ ,  $\mathcal{P}$  is a set of positive measure away from any element of  $\mathcal{S}_{\mathcal{D}}(\mathcal{H}, \mu)$ , and in particular  $\mathcal{P} \notin \mathcal{S}_{\mathcal{D}}(\mathcal{H}, \mu)$ , and thus Theorem C.1 provides the existence of  $\lambda \in \mathbb{R}^n$  with  $\mu(\mathcal{P} \cap [y(H\lambda)x \geq 0]) = \mu(\mathcal{P})$  and  $\mu(\mathcal{P} \cap [y(H\lambda)x > 0]) > 0$ . But then, by Lemma F.8,  $\mathcal{P}$  can be reduced into a smaller element of  $\mathcal{S}_{\mathcal{P}}(\mathcal{H}, \mu)$ , contradicting its minimality.

Now suppose  $\mu(\mathcal{C} \setminus \mathcal{P}) > 0$ , and set  $\nu$  to be the restriction of  $\mu$  to  $\mathcal{C}$ : for any  $C$ ,  $\nu(C) := \mu(\mathcal{C} \cap C)$ . Consider the optimization problem

$$\inf \left\{ \int \exp(-y(H\lambda)(x)) d\nu(x, y) : \lambda \in \mathbb{R}^n \right\}.$$

Consider the sublevel set of 1-suboptimal points for this problem. By Theorem C.1, there exists  $B$  so that each  $\lambda$  in this sublevel set has  $\lambda'$  with  $H\lambda = H\lambda'$   $\mu$ -a.e. and  $\|\lambda'\|_1 \leq B$ . However, by the definition of  $\mathcal{P}$ , there exists a sequence  $\{\lambda_i\}_{i=1}^{\infty}$  which is zero over  $\mathcal{P}$  and approaches  $\infty$   $\mu$ -a.e. over  $\mathcal{P}^c$ , and in particular over the positive measure set  $\mathcal{C} \setminus \mathcal{P}$ . Thus, taking any  $\lambda$  in the 1-suboptimal set, notice that

$$\lim_{i \rightarrow \infty} \int \exp(-y(H(\lambda + \lambda_i))x) d\nu(x, y) = \int \exp(-y(H\lambda)(x)) \mathbb{1}((x, y) \notin \mathcal{P}) d\nu(x, y) =: p.$$

Since  $\lambda$  has a bounded representation,  $\exp(-y(H\lambda)x) \neq 0$ , and thus  $p < \mathcal{R}_{\phi}(H\lambda)$  (Folland, 1999, Theorem 2.23(b)). But since the objective function is continuous in  $\lambda$  (cf. Lemma B.1), there must exist a large  $j$  so that  $\mathcal{R}_{\phi}(H(\lambda + \lambda_j)) < \mathcal{R}_{\phi}(H\lambda)$ , and moreover  $y(H(\lambda + \lambda_j))(x) > B$  for a subset of  $\mathcal{C}$  with positive measure. But that means  $\lambda + \lambda_j$  is in the 1-sublevel set, but can not have a representation with norm at most  $B$  (since  $H$  is a bounded linear operator), contradicting Theorem C.1.  $\square$

### F.3 Proof of Theorem 4.3

This is now just a consequence of the equivalence to primal hard cores, and the structure over  $\mathcal{C}$  developed in Theorem C.1 (which was used to prove the equivalence to primal hard cores as well).

*Proof of Theorem 4.3.* The second property is direct from Theorem C.1. For the first property, since primal hard cores exist and are  $\mu$ -a.e. equivalent to hard cores (cf. Theorem F.7), and statement thus follows by taking the sequence provided by the definition of any primal hard core.  $\square$

## G Deferred material from Section 5

*Proof of Theorem 5.1.* (Item 1) Let  $\{\lambda_i\}_{i=1}^{\infty}$  be given as per Theorem 4.3. Automatically,  $y(H\lambda_i)x = 0$  for  $(x, y) \in \mathcal{C}$ . And since  $y'(H\lambda_i)x' \uparrow \infty$  for  $\mu$ -a.e.  $(x', y') \in \mathcal{C}^c$ , it follows from the definition of  $\Phi$  that  $\lim_{i \rightarrow \infty} \phi(-y'(H\lambda_i)x) = 0$ .

(Item 2) This is a consequence of Theorem C.1.  $\square$

*Proof of Theorem 5.2.* (Item 1) Let a sequence  $\{\lambda_i\}_{i=1}^{\infty}$  be given as provided by Theorem 4.3. In particular,  $\exp(-y(H\lambda_i)x) \rightarrow \mathbb{1}(\mathcal{C})$   $\mu$ -a.e. Now choose a finite sample size  $m$ ; by Egoroff's Theorem (Folland, 1999, Theorem 2.33), for any  $\tau > 0$ , there exists  $S_{\tau}$  with  $\mu(S_{\tau}) > \mu(\mathcal{X} \times \mathcal{Y}) - \tau/m$  over which this convergence is uniform. As such, choose  $\lambda_{\tau}$  so that  $\exp(-y(H\lambda_{\tau})x) < 1/2$  over  $S_{\tau} \cap \mathcal{C}^c$ , meaning in particular  $y(H\lambda_{\tau})x > 0$  for every  $(x, y) \in S_{\tau} \cap \mathcal{C}^c$ . The probability over a draw of  $m$  points that some within  $\mathcal{C}^c$  are misclassified by  $\lambda_{\tau}$  has upper bound bound

$$\mu^m(\exists (x_i, y_i) \in \mathcal{C}^c \cdot y(H\lambda_i)x \leq 0) \leq m\mu(\mathcal{C}^c \cap [y(H\lambda_i)x \leq 0]) < \tau.$$

Since  $\tau$  can be made arbitrarily small, the probability of failure is zero. Furthermore, since  $\lambda_{\tau}$  satisfies  $y(H\lambda_{\tau})(x) = 0$   $\mu$ -a.e. over  $\mathcal{C}$  (cf. Theorem 4.3), it also follows that, with probability 1,  $\lambda_{\tau}$  abstains on every example falling within  $\mathcal{C}$ .



(Item 2) Let  $\rho > 0$  and  $\phi \in \Phi$  be given. Choose  $b > 0$ , as provided by Theorem 5.1, so that every  $\lambda \in \mathbb{R}^n$  with  $\mathcal{R}_{\phi;\mathcal{C}}(H\lambda) \leq \mathcal{R}_{\phi;\mathcal{C}}(\text{span}(\mathcal{H})) + 4 + \rho$  has a representation  $\lambda'$  with  $\|\lambda'\|_1 \leq b$ , where  $H\lambda = H\lambda'$  everywhere along  $\mathcal{C} \setminus N$ , where  $\mu(N) = 0$ ; henceforth, rule out the event that any example falls within  $N$ . Additionally, choose  $c > 0$  as provided by Proposition A.2 so that, given  $m_{\mathcal{C}}$  i.i.d. points within  $\mathcal{C} \setminus N$ , every  $f \in \text{span}(\mathcal{H}, b)$  has

$$|\mathcal{R}_{\phi;\mathcal{C}}(f) - \mathcal{R}_{\phi;\mathcal{C}}^m(f)| \leq c \frac{\sqrt{\ln(n)} + \sqrt{\ln(2/\delta)}}{\sqrt{m_{\mathcal{C}}}}. \quad (\text{G.1})$$

Now consider any  $\lambda \in \mathbb{R}^n$  with no representation  $\|\lambda'\|_1 \leq b$  so that  $H\lambda = H\lambda'$  over  $\mathcal{C} \setminus N$ , which directly entails, by Theorem 5.1, that  $\mathcal{R}_{\phi;\mathcal{C}}(H\lambda) - \mathcal{R}_{\phi;\mathcal{C}}(\text{span}(\mathcal{H})) > \rho + 4$ . Additionally choose and any  $\bar{\lambda} \in \mathbb{R}^n$  with  $\mathcal{R}_{\phi;\mathcal{C}}(H\bar{\lambda}) - \mathcal{R}_{\phi;\mathcal{C}}(\text{span}(\mathcal{H})) < 1$ , whereby the choice of  $b > 0$  indicates that, without loss of generality,  $\|\bar{\lambda}\|_1 \leq b$ . Since  $\int \phi \circ H$  is continuous (cf. Corollary B.3), considering the line segment  $\{\alpha\lambda + (1 - \alpha)\bar{\lambda} : \alpha \in [0, 1]\}$ , there must exist  $\hat{\lambda}$  with

$$\rho + 3 \leq \mathcal{R}_{\phi;\mathcal{C}}(H\hat{\lambda}) - \mathcal{R}_{\phi;\mathcal{C}}(\text{span}(\mathcal{H})) \leq \rho + 4;$$

let  $\hat{\lambda}'$  be a representation with  $\|\hat{\lambda}'\|_1 \leq b$  and  $H\hat{\lambda} = H\hat{\lambda}'$  over  $\mathcal{C} \setminus N$  (and thus it holds for every example). Applying the deviation inequality in eq. (G.1) twice,

$$\begin{aligned} \mathcal{R}_{\phi;\mathcal{C}}^m(H\hat{\lambda}) - \mathcal{R}_{\phi;\mathcal{C}}^m(H\bar{\lambda}) &\geq \mathcal{R}_{\phi;\mathcal{C}}(H\hat{\lambda}') - \mathcal{R}_{\phi;\mathcal{C}}(H\bar{\lambda}) - 2c \frac{\sqrt{\ln(n)} + \sqrt{\ln(2/\delta)}}{\sqrt{m_{\mathcal{C}}}} \\ &= \mathcal{R}_{\phi;\mathcal{C}}(H\hat{\lambda}') - \mathcal{R}_{\phi;\mathcal{C}}(\text{span}(\mathcal{H})) - (\mathcal{R}_{\phi;\mathcal{C}}(H\bar{\lambda}) - \mathcal{R}_{\phi;\mathcal{C}}(\text{span}(\mathcal{H}))) \\ &\quad - 2c \frac{\sqrt{\ln(n)} + \sqrt{\ln(2/\delta)}}{\sqrt{m_{\mathcal{C}}}} \\ &> (\rho + 3) - (1) - 2c \frac{\sqrt{\ln(n)} + \sqrt{\ln(2/\delta)}}{\sqrt{m_{\mathcal{C}}}} \\ &\geq \rho, \end{aligned}$$

where the last step used the lower bound on  $m_{\mathcal{C}}$ . Returning to  $\lambda \in \mathbb{R}^n$  as specified above, convexity, in the form of Lemma A.4, grants that  $\mathcal{R}_{\phi;\mathcal{C}}^m(H\bar{\lambda}) < \mathcal{R}_{\phi;\mathcal{C}}^m(H\hat{\lambda})$  implies  $\mathcal{R}_{\phi;\mathcal{C}}^m(H\hat{\lambda}) \leq \mathcal{R}_{\phi;\mathcal{C}}^m(H\lambda)$ , and thus

$$\mathcal{R}_{\phi;\mathcal{C}}^m(H\lambda) - \mathcal{R}_{\phi;\mathcal{C}}^m(\text{span}(\mathcal{H})) \geq \mathcal{R}_{\phi;\mathcal{C}}^m(H\hat{\lambda}) - \mathcal{R}_{\phi;\mathcal{C}}^m(H\bar{\lambda}) > \rho.$$

Since  $\lambda$  was arbitrary, it follows that every  $\lambda$  with no representation  $\|\lambda'\|_1 > b$  that has agreement of  $H\lambda$  and  $H\lambda'$   $\mu$ -a.e. over  $\mathcal{C}$  does not lie in the empirical  $\rho$ -sublevel set. Since  $\mathcal{R}_{\phi;\mathcal{C}}^m$  is convex and continuous, the  $\rho$ -sublevel set is nonempty, and thus every  $\lambda'$  within it has a representation  $\|\lambda''\|_1 \leq b$ .  $\square$

## H Deferred material from Section 6

*Proof of Proposition 6.2.* This proof is essentially a repackaging of various results and comments due to Bartlett et al. (2006). Fix any  $\phi \in \Phi$ ;  $\phi$  is convex, increasing at 0, and differentiable at 0, which grants that the corresponding  $\psi$ -transform is classification calibrated (Bartlett et al., 2006, Theorem 6, although note losses in the present manuscript are increasing rather than decreasing). It follows that  $\psi(\mathcal{R}_{\mathcal{L}}(f) - \mathcal{R}_{\mathcal{L}}(\mathfrak{F})) \leq \mathcal{R}_{\phi}(f) - \mathcal{R}_{\phi}(\mathfrak{F})$ , (Bartlett et al., 2006, Theorem 3, part 3(c)).

Next,  $\psi(0) = 0$  (Bartlett et al., 2006, Lemma 5, part 8),  $\psi(r) > 0$  when  $r > 0$  (Bartlett et al., 2006, Lemma 5, part 9(b)), and since  $\psi$  is convex by construction (Bartlett et al., 2006, Definition 2), it follows by Lemma A.4 that  $\psi$  is increasing. Since  $\psi$  is continuous as well, (Bartlett et al., 2006, Lemma 5, part 6), it follows that  $\psi$  has a well-defined inverse along the image  $\psi([0, 1])$ . Finally, the fact that  $\psi^{-1}(r) \downarrow 0$  as  $r \downarrow 0$  is due to Bartlett et al. (2006, Theorem 3, part 3(b)).  $\square$

*Proof of Theorem 6.4.* Throughout this proof,  $\delta' := \delta/8$  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  $\mu(\mathcal{C}) = 0$  or  $\mu(\mathcal{C}) = \mu(\mathcal{X} \times \mathcal{Y})$  (i.e., when terms depending on either appear in denominators); interpret these expressions as simply being  $\infty$ , whereby the bounds hold automatically.

(Item 1) Let  $\mathcal{S}_{\mathcal{C}}$  and  $\mathcal{S}_+$  respectively denote the set of samples landing in  $\mathcal{C}$  and  $\mathcal{C}^c$ , where the notation proposed in the theorem statement provides  $m_{\mathcal{C}} = |\mathcal{S}_{\mathcal{C}}|$  and  $m_+ = |\mathcal{S}_+|$ . By a Chernoff bound (Kearns and Vazirani, 1994, Theorem 9.2), basic deviations for these quantities are

$$\begin{aligned}\Pr^m[|\mathcal{S}_{\mathcal{C}}| < (\mu(\mathcal{C}) - \tau)m] &\leq \exp(-m\tau^2/2) = \delta', \\ \Pr^m[|\mathcal{S}_+| < (\mu(\mathcal{C}^c) - \tau)m] &\leq \exp(-m\tau^2/2) = \delta',\end{aligned}$$

where  $\tau = \sqrt{\frac{1}{2m} \ln(\frac{1}{\delta'})}$ , and  $\Pr^m$  denotes the product measure corresponding to  $\mu$ . Label these failure events  $F_1$  and  $F_2$ , and henceforth rule them out.

(Item 2) As provided by Theorem 5.2, there exists  $\bar{\lambda} \in \mathbb{R}^n$  with  $y_i(H\bar{\lambda})x_i > 0$  for all  $(x_i, y_i)$  falling in  $\mathcal{C}^c$ , and  $y_i(H\bar{\lambda})x_i = 0$  for those landing in  $\mathcal{C}$ . Consequently,

$$\begin{aligned}\mathcal{R}_{\phi}(\text{span}(\mathcal{H})) &= \inf_{\lambda} \inf_{c>0} \mathcal{R}_{\phi, \mathcal{C}}(H(\lambda + c\bar{\lambda})) + \mathcal{R}_{\phi, \mathcal{C}^c}(H(\lambda + c\bar{\lambda})) \\ &= \inf_{\lambda} \inf_{c>0} \mathcal{R}_{\phi, \mathcal{C}}(H\lambda) \\ &\leq \mathcal{R}_{\phi}(\text{span}(\mathcal{H})).\end{aligned}$$

Combining this with

$$\mathcal{R}_{\phi, \mathcal{C}^c}(H\lambda) + \mathcal{R}_{\phi, \mathcal{C}}(H\lambda) = \mathcal{R}_{\phi}(H\lambda) \leq \mathcal{R}_{\phi}(\text{span}(\mathcal{H})) + \epsilon,$$

it follows that

$$\mathcal{R}_{\phi, \mathcal{C}^c}(H\lambda) \leq \mathcal{R}_{\phi}(\text{span}(\mathcal{H})) - \mathcal{R}_{\phi, \mathcal{C}}(H\lambda) + \epsilon = \mathcal{R}_{\phi, \mathcal{C}}(\text{span}(\mathcal{H})) - \mathcal{R}_{\phi, \mathcal{C}}(H\lambda) + \epsilon \leq \epsilon.$$

Next, since  $\phi(0) > 0$  and  $\phi$  is nondecreasing (cf. Lemma A.1),

$$\mathcal{R}_{\mathcal{L}, \mathcal{C}^c}^m(H\lambda) \leq \frac{\mathcal{R}_{\phi, \mathcal{C}^c}^m(H\lambda)}{\phi(0)} = \frac{\epsilon}{\phi(0)}.$$

To obtain eq. (6.5) from here, first notice that  $\mathcal{S}_+$ , the portion of the sample falling within  $\mathcal{C}^c$ , can be interpreted as an i.i.d. sample from the probability measure  $\mu(\cdot \cap \mathcal{C})/\mu(\mathcal{C})$ . Next, the VC dimension of  $\text{span}(\mathcal{H})$  is the VC dimension of linear separators over the transformed space

$$\{(h_1(x), h_2(x), \dots, h_n(x)), y) : (x, y) \in \mathcal{X} \times \mathcal{Y}\};$$

namely, it is  $n$ . As such, eq. (6.5) follows by an application of a relative deviation version of the VC Theorem (Boucheron et al., 2005, discussion preceding Corollary 5.2).

To obtain eq. (6.6), note that  $\epsilon < \phi(0)/m$  means there are no mistakes over  $\mathcal{C}^c$ :

$$\begin{aligned}\phi(0) &> m\epsilon \geq m \left( \mathcal{R}_{\phi}^m(H\hat{\lambda}) - \bar{\mathcal{R}}_{\phi}^m(\text{span}(\mathcal{H})) \right) \\ &\geq m_+ \mathcal{R}_{\phi, \mathcal{C}^c}^m \\ &\geq \sum_{i=1}^{m_+} \phi(-y_i(H\hat{\lambda})x_i) \\ &\geq \max_{i \in [m_+]} \phi(-y_i(H\hat{\lambda})x_i);\end{aligned}$$

that is to say, for every  $(x_i, y_i) \in \mathcal{S}_+$ ,  $0 < y_i(H\hat{\lambda})x_i$ . Plugging  $\mathcal{R}_{\mathcal{L}}^m(H\lambda) = 0$  into the same relative deviation bound as before (Boucheron et al., 2005, discussion preceding Corollary 5.2), the second bound follows.

(Item 3) By Theorem 5.2, there exist constants  $b > 0$  and  $c \geq \phi(b)$ , depending on  $\mathcal{H}, \mu, \phi, \mathcal{C}$ , so that with probability at least  $1 - \delta'$ , if  $m_{\mathcal{C}} \geq c^2(\ln(n) + \ln(1/\delta'))$ , then every  $\rho$ -suboptimal predictor over  $\mathcal{C}$ , and in particular  $\lambda$ , has a representation  $\lambda'$  which is equivalent to  $\lambda$   $\mu$ -a.e. over  $\mathcal{C}$ , and satisfies  $\|\lambda'\|_1 \leq b$ . As such, since

$$\mathcal{R}_{\phi; \mathcal{C}}^m(H\lambda) = \mathcal{R}_{\phi; \mathcal{C}}^m(H\lambda') \quad \text{and} \quad \mathcal{R}_{\phi; \mathcal{C}}(H\lambda) = \mathcal{R}_{\phi; \mathcal{C}}(H\lambda'),$$

an application of Proposition A.2 grants

$$\begin{aligned} \mathcal{R}_{\phi; \mathcal{C}}(H\lambda) &= \mathcal{R}_{\phi; \mathcal{C}}(H\lambda') \\ &\leq \mathcal{R}_{\phi; \mathcal{C}}^m(H\lambda') + \frac{c \left( \sqrt{\ln(n)} + \sqrt{\ln(2/\delta')} \right)}{\sqrt{m_{\mathcal{C}}}} \\ &= \mathcal{R}_{\phi; \mathcal{C}}^m(H\lambda) + \frac{c \left( \sqrt{\ln(n)} + \sqrt{\ln(2/\delta')} \right)}{\sqrt{m_{\mathcal{C}}}} \\ &\leq \mathcal{R}_{\phi; \mathcal{C}}^m(\text{span}(\mathcal{H})) + \epsilon + \frac{c \left( \sqrt{\ln(n)} + \sqrt{\ln(2/\delta')} \right)}{\sqrt{m_{\mathcal{C}}}}. \end{aligned}$$

Next, noting that Theorem 5.1 provides that a minimizing sequence to  $\mathcal{R}_{\phi; \mathcal{C}}(\text{span}(\mathcal{H}))$  can be taken without loss of generality to lie within a compact set (e.g., points with  $l^1$  norm at most  $b$ ), it follows that a minimizer  $\bar{\lambda}$  exists; by an application of McDiarmid's inequality, with probability at least  $1 - \delta'$ ,

$$\mathcal{R}_{\phi; \mathcal{C}}^m(\text{span}(\mathcal{H})) \leq \mathcal{R}_{\phi; \mathcal{C}}^m(H\bar{\lambda}) \leq \mathcal{R}_{\phi; \mathcal{C}}(H\bar{\lambda}) + c \sqrt{\frac{2 \ln(1/\delta')}{m_{\mathcal{C}}}}.$$

(Note,  $\bar{\lambda}$  is independent of the sample, thus McDiarmid suffices, with constant  $c \geq \phi(b)$  since  $\bar{\lambda}$  is in this initial sublevel set.) Combining these two pieces, it follows that

$$\mathcal{R}_{\phi; \mathcal{C}}(H\lambda) - \mathcal{R}_{\phi; \mathcal{C}}(\text{span}(\mathcal{H})) \leq \epsilon + \frac{c \left( \sqrt{\ln(n)} + 4\sqrt{\ln(2/\delta')} \right)}{\sqrt{m_{\mathcal{C}}}},$$

which is precisely eq. (6.7).

To produce eq. (6.8), the definition of the  $\psi$ -transform (cf. Proposition 6.2), combined with Equation (6.7), provides

$$\begin{aligned} \mathcal{R}_{\mathcal{L}; \mathcal{C}}(H\lambda) - \mathcal{R}_{\mathcal{L}; \mathcal{C}}(\mathfrak{F}) &\leq \psi^{-1}(\mathcal{R}_{\phi; \mathcal{C}}(H\lambda) - \mathcal{R}_{\phi; \mathcal{C}}(\mathfrak{F})) \\ &= \psi^{-1}(\mathcal{R}_{\phi; \mathcal{C}}(H\lambda) - \mathcal{R}_{\phi; \mathcal{C}}(\text{span}(\mathcal{H})) + \mathcal{R}_{\phi; \mathcal{C}}(\text{span}(\mathcal{H})) - \mathcal{R}_{\phi; \mathcal{C}}(\mathfrak{F})) \\ &\leq \psi^{-1} \left( \epsilon + \frac{c \left( \sqrt{\ln(n)} + 4\sqrt{\ln(2/\delta')} \right)}{\sqrt{m_{\mathcal{C}}}} + \mathcal{R}_{\phi; \mathcal{C}}(\text{span}(\mathcal{H})) - \mathcal{R}_{\phi; \mathcal{C}}(\mathfrak{F}) \right). \end{aligned}$$

(Item 4) Combining the lower bound on  $m$  with Item 1,

$$\begin{aligned} m_+ &\geq m\mu(\mathcal{C}^c)/2, \\ m_{\mathcal{C}} &\geq m\mu(\mathcal{C})/2 \geq c^2(\ln(n) + \ln(1/\delta')); \end{aligned}$$

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  $\mathcal{C}^c$  when  $\epsilon < \phi(0)/m$ ),

$$\begin{aligned}
\mathcal{R}_{\mathcal{L}}(H\lambda) - \mathcal{R}_{\mathcal{L}}(\mathfrak{F}) &= (\mathcal{R}_{\mathcal{L};\mathcal{C}}(H\lambda) - \mathcal{R}_{\mathcal{L};\mathcal{C}}(\mathfrak{F})) + (\mathcal{R}_{\mathcal{L};\mathcal{C}^c}(H\lambda) - \underbrace{\mathcal{R}_{\mathcal{L};\mathcal{C}^c}(\mathfrak{F})}_{=0}) \\
&\leq \psi^{-1} \left( \epsilon + \frac{c \left( \sqrt{\ln(n)} + 4\sqrt{\ln(2/\delta')} \right)}{\sqrt{m_{\mathcal{C}}}} + \mathcal{R}_{\phi;\mathcal{C}}(\text{span}(\mathcal{H})) - \mathcal{R}_{\phi;\mathcal{C}}(\mathfrak{F}) \right) \\
&\quad + \frac{4(n \ln(2m_+ + 1) + \ln(4/\delta'))}{m_+} \\
&\leq \psi^{-1} \left( \epsilon + \frac{c\sqrt{2} \left( \sqrt{\ln(n)} + 4\sqrt{\ln(2/\delta')} \right)}{\sqrt{m\mu(\mathcal{C})}} + \mathcal{R}_{\phi;\mathcal{C}}(\text{span}(\mathcal{H})) - \mathcal{R}_{\phi;\mathcal{C}}(\mathfrak{F}) \right) \\
&\quad + \frac{8(n \ln(m\mu(\mathcal{C}^c) + 1) + \ln(4/\delta'))}{m\mu(\mathcal{C}^c)}.
\end{aligned}$$

□

## I Deferred material from Section 7

*Proof of Theorem 7.1.* Let  $\mathcal{C}$  be a hard core for  $(\mathcal{H}, \mu)$ , set  $\rho := 1$ , and let  $b > 0$  and  $c > 0$  be the corresponding reals provided in the guarantee of Theorem 6.4. Note first that  $\mathcal{R}_{\phi}(\text{span}(\mathcal{H})) = \mathcal{R}_{\phi}(\mathfrak{F})$  implies  $\mathcal{R}_{\phi;\mathcal{C}}(\text{span}(\mathcal{H})) = \mathcal{R}_{\phi;\mathcal{C}}(\mathfrak{F})$ , since predictions are  $\mu$ -a.e. perfect off the hard core (cf. Theorem 5.1). Set  $\delta_i = 1/i^2$ , and choose  $m_i \uparrow \infty$  large enough and  $\epsilon_i \downarrow 0$  small enough so that the relevant finite sample bound from Theorem 6.4 holds, and goes to zero. (Note that all bounds go to zero as  $m_i \uparrow \infty$  and  $\epsilon_i \downarrow 0$ ; the word “relevant” refers to choosing a bound corresponding to the regime  $\mu(\mathcal{C}) = 0$ , or  $\mu(\mathcal{C}^c) = 0$ , or  $\min\{\mu(\mathcal{C}), \mu(\mathcal{C}^c)\} > 0$ .) Note, by the strong assumption, the term  $\mathcal{R}_{\phi;\mathcal{C}}(\text{span}(\mathcal{H})) - \mathcal{R}_{\phi;\mathcal{C}}(\mathfrak{F})$  may be dropped.

Now let  $F_i$  be the failure event of the corresponding finite sample guarantee; by choice of  $\delta_i$ ,  $\sum_i \Pr(F_i) = \sum_i i^{-2} = \pi^2/6 < \infty$ . Thus, by the Borel-Cantelli Lemma and de Morgan’s Laws,  $\Pr(\liminf_{i \rightarrow \infty} F_i^c) = 1$ , meaning  $\Pr(\exists j \cdot \forall i \geq j \cdot F_i^c) = 1$ . This means that the bounds hold for all large  $i$  (with probability 1), and the result follows by choice of  $m_i$  and  $\epsilon_i$ . □

*Proof of Proposition 7.3.* This proof will proceed in the following stages. First, it is shown that the infimal risk  $\mathcal{R}_{\phi}(\mathfrak{F})$  can be approximated arbitrarily well by bounded measurable functions. Next, Lusin’s theorem will allow this consideration to be restricted to a function which is continuous over a compact set. Finally, this function is approximated by a decision tree.

Let  $\mu, \phi, \{\mathcal{H}_i\}_{i=1}^{\infty}$ , and  $\epsilon > 0$  be given as specified. Since the infimum in  $\mathcal{R}_{\phi}(\mathfrak{F})$  is in general not attained, let  $g \in \mathfrak{F}$  be a measurable function satisfying

$$\mathcal{R}_{\phi}(g) \leq \epsilon/4 + \mathcal{R}_{\phi}(\mathfrak{F}).$$

Next let  $z > 0$  be a sufficiently large real so that  $\phi(-z) < \epsilon/4$ ; such a value must exist since  $\lim_{z \rightarrow -\infty} \phi(z) = 0$ . Correspondingly, define a truncation of  $g$  as

$$\hat{g}(x) := \min\{z, \max\{-z, g(x)\}\}.$$

There are three cases to consider. If  $|yg(x)| \leq z$ , then  $\phi(-y\hat{g}(x)) = \phi(-yg(x))$ . If  $-yg(x) > z$ , then by the nondecreasing property (cf. Lemma A.1),  $\phi(-yg(x)) \geq \phi(-y\hat{g}(x))$ . Lastly, if  $-yg(x) < -z$ ,

then  $\phi(-y\hat{g}(x)) \leq \phi(-yg(x)) + \epsilon/4$  by choice of  $z$ . Together, it follows that

$$\begin{aligned}\mathcal{R}_\phi(\hat{g}) &= \int \phi(-y\hat{g}(x))d\mu(x, y) \\ &\leq \int (\phi(-yg(x)) + \epsilon/4)d\mu(x, y) \\ &= \mathcal{R}_\phi(g) + \epsilon\mu(\mathcal{X} \times \mathcal{Y})/4 \\ &\leq \mathcal{R}_\phi(\mathfrak{F}) + \epsilon/2,\end{aligned}$$

which used the fact that  $\mu$  is a probability measure. Crucially,  $\hat{g}$  is now a bounded measurable function. Throughout the remained of this proof, let  $\|\cdot\|_u$  denote the uniform norm, meaning

$$\|f\|_u := \sup_x |f(x)|.$$

For example,  $\|\hat{g}\|_u < \infty$ .

In order to apply Lusin's Theorem and pass to continuous functions with compact support, a few properties must be verified. First, since  $\mu_{\mathcal{X}}$  is a Borel probability measure, it is finite on all compact Borel sets. Next,  $\mathbb{R}^d$  is a separable metric space, and thus second countable. Finally,  $\mathbb{R}^d$  is a locally compact Hausdorff space. It follows that  $\mu_{\mathcal{X}}$  is a Radon measure (Folland, 1999, Theorem 7.8).

Henceforth, set  $\tau := \epsilon/(8 \max\{1, \phi(\|\hat{g}\|_u)\})$ . By Lusin's Theorem, there exists a measurable function  $h$  which is continuous, has compact support, satisfies  $\mu_{\mathcal{X}}([\hat{g} \neq h]) < \tau$  and  $\|h\|_u \leq \|\hat{g}\|_u$  (Folland, 1999, Theorem 7.10, Lusin's Theorem). But continuity over a compact set implies uniform continuity. Furthermore, the convex function  $\phi$ , restricted to the domain  $[-z, z]$ , is necessarily Lipschitz. As such, it is possible to choose  $\delta > 0$  so that for any  $x, x'$  with  $\|x - x'\|_\infty < \delta$  and any  $y \in \{-1, +1\}$ , it follows that  $|\phi(-yh(x)) - \phi(-yh(x'))| < \tau$ . Notice that this in fact holds everywhere, since outside of its support  $h$  is just zero.

As such, let  $T$  be the smallest integer so that  $T > \sup\{\|x\|_\infty : h(x) \neq 0\}$  (which exists since  $h$  has compact support) and also  $1/T < \delta$ . For any  $t \geq T$ , construct a simple function approximation  $f$  to  $h$  as follows. Partition the cube  $[-t, t]^d$  into subcubes (formed as a product of half open intervals in order to correctly produce a partition) having side length  $1/t$  with vertices at the appropriate lattice points granting a correct partitioning. Let  $\{C_i\}_{i=1}^k$  index this family of subcubes, and let  $p_i$  be some point within each subcube. Define an approximant

$$f(x) := \sum_{i=1}^k h(p_i) \mathbb{1}(x \in C_i).$$

It follows that, for a point  $x \in C_i$  and any  $y \in \{-1, +1\}$ ,

$$|\phi(-yf(x)) - \phi(-yh(x))| = |\phi(-yh(p_i)) - \phi(-yh(x))| < \tau$$

by construction. Since  $C_i$  was arbitrary, this holds for every subcube; and it furthermore holds outside the support of  $f$ , where  $h$  and  $f$  are both guaranteed to be the constant 0.

Combining the various approximation components, it follows that

$$\begin{aligned}\mathcal{R}(f) &= \int \phi(-yf(x))d\mu(x, y) \\ &\leq \tau\mu(\mathcal{X} \times \mathcal{Y}) + \int \phi(-yh(x))d\mu(x, y) \\ &\leq \epsilon/8 + \int \phi(-y\hat{g}(x))\mathbb{1}(\hat{g}(x) = h(x))d\mu(x, y) + \int \phi(-yh(x))\mathbb{1}(\hat{g}(x) \neq h(x))d\mu(x, y) \\ &\leq \epsilon/8 + \mathcal{R}_\phi(\hat{g}) + \mu_{\mathcal{X}}([\hat{g} \neq h])\phi(\|\hat{g}\|_u) \\ &< \epsilon + \mathcal{R}_\phi(\mathfrak{F}).\end{aligned}$$

To finish, note by construction that  $f$ , which was formed from axis-aligned subcubes at lattice points within  $[-t, t)$ , satisfies  $f \in \text{span}(\mathcal{H}_t)$  (the indicator for each subcube can be modeled as an element of  $\mathcal{H}_t$ ).  $\square$

*Proof of Theorem 7.4.* Proceed as in the proof of Theorem 7.1, with one modification. First determine  $\epsilon_i$ . At each stage, choose  $j_i$  large enough so that  $\mathcal{H}_{j_i}$  satisfies  $\mathcal{R}_\phi(\text{span}(\mathcal{H}_{j_i})) < \mathcal{R}_\phi(\mathfrak{F}) + \epsilon_i$ ; the existence of such a  $j_i$  is straight from the definition of L-SRM families. Now choose  $m_i$  large enough to satisfy the necessary conditions in the proof of Theorem 7.1; meaning the relevant bound from Theorem 6.4 may be instantiated, and furthermore these bounds approach zero as  $i \rightarrow \infty$ . Now that  $m_i$  may be quite massive, as it must now smash the term  $n = |\mathcal{H}_{j_i}|$ . The proof is otherwise identical to before.  $\square$