# The perils of being unhinged:
# On the accuracy of classifiers minimizing
# a noise-robust convex loss

Philip M. Long
Google
plong@google.com

Rocco A. Servedio
Columbia University
rocco@cs.columbia.edu

**Abstract**

van Rooyen et al. [2015] introduced a notion of convex loss functions being robust to random classification noise, and established that the "unhinged" loss function is robust in this sense. In this note we study the accuracy of binary classifiers obtained by minimizing the unhinged loss, and observe that even for simple linearly separable data distributions, minimizing the unhinged loss may only yield a binary classifier with accuracy no better than random guessing.

## 1   Introduction

As van Rooyen et al. noted in the first sentence of the abstract of van Rooyen et al. [2015], "Convex potential minimisation is the *de facto* approach to binary classification." Given the ubiquity of this approach, it is natural to study its abilities and limitations in the presence of noise, and indeed this is the subject of many works [see Zhang, 2004, Bartlett et al., 2006, Long and Servedio, 2010, Manwani and Sastry, 2013, Natarajan et al., 2013, van Rooyen et al., 2015, Ghosh et al., 2017]).

The aim of this note is to clarify the connection between minimizing a convex potential function which is "robust to classification noise" in the sense of van Rooyen et al. [2015], and learning (i.e. performing accurate classification).

**Background.** Motivated by the observation that the popular AdaBoost algorithm (which works by minimizing the (convex) exponential potential function) can have empirically poor classification accuracy when run on noisy data [Dietterich, 2000, Freund and Schapire, 1996, Maclin and Opitz, 1997], Long and Servedio [2010] studied the performance of classification algorithms which work by minimizing a convex potential function in settings where linearly separable data is contaminated with random classification noise (RCN). The main result of Long and Servedio [2010] is a proof that for a certain simple learning problem corresponding to a "clean" data distribution $\mathcal{D}_1$ that is linearly separable with a margin, for *any* "convex potential function"

$\phi$, minimizing $\phi$ over all linear combinations of base features in the presence of random classification noise only yields a binary classifier with an error rate of $1/2$ under the clean distribution $\mathcal{D}_1$. (Here a "convex potential function" is a convex function $\phi : \mathbb{R} \to \mathbb{R}$ satisfying certain mild conditions which we detail in Definition 1.) This is in sharp contrast with the fact that, in the noise-free setting of a data distribution that is linearly separable with a margin, driving the potential to zero leads to a perfectly accurate binary classifier.

In an effort to address the discouraging negative result of Long and Servedio [2010], van Rooyen et al. [2015] considered a weakening of the [Long and Servedio, 2010] conditions for a convex potential function. In particular, they allow such functions $\phi$ to take negative values (which is disallowed by the definition of Long and Servedio). We refer to a function satisfying the condition of van Rooyen et al. [2015] as a "relaxed convex potential function."

The main result of van Rooyen et al. [2015] is that they propose a certain relaxed convex potential function, which we denote $\phi^\star$, and prove that it is "RCN-robust".[1] We give a formal definition of RCN-robustness in Section 2, but intuitively it means that a minimizer of this potential function (minimizing over all linear combinations of base features) under random classification noise performs no worse than a minimizer obtained with no random classification noise. van Rooyen et al. [2015] also define a notion of "strong RCN-robustness" and show that their $\phi^\star$ is the unique relaxed convex potential function which satisfies strong RCN-robustness.

**This note.** The purpose of the present note is to discuss the *accuracy* of the classifier obtained by minimizing the relaxed convex potential function $\phi^\star$ of van Rooyen et al. [2015]. Our main observation is that, for a simple learning problem corresponding to a certain "clean" data distribution $\mathcal{D}_2$ that is linearly separable with a margin, minimizing $\phi^\star$ over all bounded-norm linear combinations of base features *even when there is no random classification noise* only yields a binary classifier with an error rate of $1/2$. Since, as shown by van Rooyen et al. [2015], $\phi^\star$ is the unique strong RCN-robust relaxed convex potential function, this means that minimizing any strong RCN-robust relaxed convex potential function in this noise-free scenario may only yield a binary classifier with an error rate of $1/2$, which can be obtained through random guessing.

Our observation is consistent with the result of van Rooyen et al. [2015] that $\phi^\star$ is RCN-robust, since, informally, that condition only states that "you don't do any worse when there is RCN than when there is no RCN." Our example demonstrates even when there is no noise, the accuracy of the binary classifier obtained by minimizing $\phi^\star$ may be only $1/2$, and of course the accuracy is no worse than this when there actually is noise.

---

[1]van Rooyen et al. [2015] uses the term "SLN-robust", where the acronym stands for Symmetric Label Noise.

2

| Potential function | Reference | Satisfies Definition 1? |
|---|---|---|
| Exponential:<br>$\phi(z) = e^{-z}$ | [Freund and Schapire, 1997] | Yes |
| Mixed linear/exponential:<br>$\phi(z) = \begin{cases} 1 - z & \text{if } z \leq 0 \\ e^{-z} & \text{if } z > 0 \end{cases}$ | [Domingo and Watanabe, 2000] | Yes |
| Logistic:<br>$\phi(z) = \ln(1 + e^{-2z})$ | [Friedman et al., 1998] | Yes |
| Hinge:<br>$\phi(z) = \max\{0, 1 - z\}$ | [Gentile and Warmuth, 1998] | No |
| Unhinged:<br>$\phi(z) = 1 - z$ | [van Rooyen et al., 2015] | No |

Table 1: Some commonly used potential functions.

## 2 Preliminaries

### 2.1 Background: the negative result of Long and Servedio [2010] for convex potential functions

**Convex potential functions.** We recall the following definition which is central to the work of Long and Servedio [2010]:

**Definition 1** ([Long and Servedio, 2010], Definition 1). *A function $\phi : \mathbb{R} \to \mathbb{R}$ is a convex potential function if it satisfies the following:*

1. *$\phi \in C^1$ (i.e. $\phi$ is differentiable and $\phi'$ is continuous) and $\phi$ is convex and nonincreasing; and*

2. *$\phi'(0) < 0$ and $\lim_{x \to \infty} \phi(x) = 0$ (hence $\phi$ is everywhere non-negative).*

A number of potential functions used in the literature fit this definition, including the exponential potential function used by AdaBoost [Freund and Schapire, 1997], the mixed linear/exponential potential function used by MadaBoost [Domingo and Watanabe, 2000], and the logistic function used by LogitBoost [Friedman et al., 1998]; see Table 1.

**Linearly separable learning problems.** One of the simplest models for binary-labeled data over $\mathbb{R}^d$ is that of data which is *linearly separable with a margin.* A "clean" probability distribution $\mathcal{D}$ over $\mathbb{R}^d \times \{-1, 1\}$ is linearly separable with margin $\gamma > 0$ if there is a target weight vector $w = (w_1, \ldots, w_d) \in \mathbb{R}^d$ such that

$$\Pr_{(\boldsymbol{x}, \boldsymbol{y}) \sim \mathcal{D}} \left[ \frac{\boldsymbol{y}(w \cdot \boldsymbol{x})}{|w_1| + \cdots + |w_d|} < \gamma \right] = 0.$$

A very standard learning approach for such a setting is to choose a hypothesis weight vector $v = (v_1, \ldots, v_d) \in \mathbb{R}^d$ with the aim of minimizing the "global" potential function

$$P_{\phi,\mathcal{D}}(v) := \mathop{\mathbf{E}}_{(\boldsymbol{x},\boldsymbol{y}) \sim \mathcal{D}} [\phi(\boldsymbol{y}(v \cdot \boldsymbol{x}))]. \tag{1}$$

(Of course, given a finite sample of draws from $\mathcal{D}$, this is typically done by minimizing the corresponding expectation over the sample.) It is well known that if $\mathcal{D}$ is linearly separable with margin $\gamma > 0$, then for a range of different choices of the convex potential function $\phi$ (including the AdaBoost, MadaBoost and Logit-Boost potential functions described above), greedy iterative algorithms that perform coordinatewise gradient descent to minimize $P_{\phi,\mathcal{D}}$ will drive the misclassification error $\mathbf{Pr}_{(\boldsymbol{x},\boldsymbol{y}) \sim \mathcal{D}}[\boldsymbol{y} \neq \mathrm{sign}(v \cdot \boldsymbol{x})]$ to zero. Indeed, the AdaBoost [Freund and Schapire, 1997], MadaBoost [Domingo and Watanabe, 2000] and LogitBoost [Friedman et al., 1998] boosting algorithms correspond precisely to greedy coordinatewise gradient descent procedures of this sort; see the work of Mason et al. [1999] for details.

**Learning problems with random classification noise.** Let $\mathcal{D}$ be a data distribution over $\mathbb{R}^d \times \{-1, 1\}$ as described above. The $\eta$-*RCN corrupted* version of $\mathcal{D}$ is the following distribution $\overline{\mathcal{D}}_\eta$ over $\mathbb{R}^d \times \{-1, 1\}$: a draw from $\overline{\mathcal{D}}_\eta$ is obtained by drawing $(\boldsymbol{x}, \boldsymbol{y}) \sim \mathcal{D}$ and flipping the label $\boldsymbol{y}$ with probability $\eta$.

**The negative result of Long and Servedio [2010].** The main result of Long and Servedio [2010] is that there is *no* convex potential function such that minimizing $\phi$ on the $\eta$-RCN corrupted distribution $\overline{\mathcal{D}}_\eta$ will succeed in achieving nontrivial classification accuracy:

**Theorem 2.** *Fix any noise rate $0 < \eta < 1/2$ and any convex potential function $\phi$. There is a distribution $\mathcal{D}$ over $\mathbb{R}^2 \times \{-1, 1\}$ (in fact the distribution $\mathcal{D}$ is supported on three points in the unit disc) and a margin parameter $\gamma > 0$ such that (a) $\mathcal{D}$ is linearly separable with margin $\gamma$, but (b) any weight vector $v$ which minimizes $\mathbf{E}_{(\boldsymbol{x},\boldsymbol{y}) \sim \overline{\mathcal{D}}_\eta} [\phi(\boldsymbol{y}(v \cdot \boldsymbol{x}))]$ has $\mathbf{Pr}_{(\boldsymbol{x},\boldsymbol{y}) \sim \mathcal{D}}[\boldsymbol{y} \neq \mathrm{sign}(v \cdot \boldsymbol{x})] = 1/2$.*

(See Appendix A for a proof that for any convex potential function $\phi$, the minimizer analyzed in Theorem 2 exists.)

## 2.2 Relaxed convex potential functions: a new hope?

Motivated by the goal of circumventing the negative result of Theorem 2, van Rooyen et al. [2015] consider a relaxed form of Definition 1:

**Definition 3.** *A function $\phi : \mathbb{R} \to \mathbb{R}$ is a* relaxed *convex potential function if it satisfies the following:*

1. *$\phi \in C^1$ (i.e. $\phi$ is differentiable and $\phi'$ is continuous) and $\phi$ is convex and nonincreasing; and*

2. *$\phi'(0) < 0$.*

The only difference between Definition 1 and Definition 3 is that the latter does not require $\lim_{x \to \infty} \phi(x) = 0$; a relaxed convex loss function may take (arbitrarily large magnitude) negative values. van Rooyen et al. [2015] exploit this flexibility by proposing the following simple potential function, which they call the "unhinged loss":

$$\phi^*(z) = 1 - z.$$

It is trivial to verify that $\phi^\star$ satisfies Definition 3 and hence is a valid relaxed convex potential function. (Note, also, that if the simpler $\phi(z) = -z$ is used instead, all gradients and minima are unaffected.) We note that the unhinged loss is a member of the class of *symmetric* loss functions, which satisfy $\phi(z) + \phi(-z) =$ constant; such loss functions have been studied by a number of authors, see e.g. Charoenphakdee et al. [2019a], Ghosh et al. [2015].

**RCN-robustness.** van Rooyen et al. [2015] analyze the relaxed convex potential function $\phi^\star$ through the lens of a new notion which we will call *RCN-robustness*. Their definition (Definition 1 of van Rooyen et al. [2015]) applies to a general pair $(\ell, \mathcal{F})$ where $\ell$ is a loss function and $\mathcal{F}$ is a class which may consist of any collection of functions mapping a domain $X$ to $\mathbb{R}$.

Informally, a pair $(\phi, \mathcal{F})$ is RCN-robust if minimizing $\phi$ over $\mathcal{F}$ on noise-free data gives the same binary classification performance as minimizing $\phi$ over $\mathcal{F}$ on RCN-contaminated data. More precisely, we have the following:

**Definition 4.** *Let $\mathcal{F}$ be a set of real-valued functions over $\mathbb{R}^d$ and let $\phi$ be a potential function. The pair $(\phi, \mathcal{F})$ is is said to be* RCN-robust *if the following holds: Let $\mathcal{D}$ be any distribution over $\mathbb{R}^d \times \{-1, 1\}$ and let $0 < \eta < 1/2$ be any noise rate. If $f$ is a minimizer of $\mathbf{E}_{(\boldsymbol{x},\boldsymbol{y}) \sim \mathcal{D}} [\phi(\boldsymbol{y}(f(\boldsymbol{x})))]$ over $f \in \mathcal{F}$ and $g$ is the minimizer of $\mathbf{E}_{(\boldsymbol{x},\boldsymbol{y}) \sim \overline{\mathcal{D}}_\eta} [\phi(\boldsymbol{y}(g(\boldsymbol{x})))]$ over $g \in \mathcal{F}$, then*

$$\Pr_{(\boldsymbol{x},\boldsymbol{y}) \sim \mathcal{D}}[\boldsymbol{y} \neq \mathrm{sign}(f(\boldsymbol{x}))] = \Pr_{(\boldsymbol{x},\boldsymbol{y}) \sim \mathcal{D}}[\boldsymbol{y} \neq \mathrm{sign}(g(\boldsymbol{x}))]. \tag{2}$$

van Rooyen et al. [2015] specialize Definition 4 to the function class $\mathcal{F}_{\mathrm{lin}}$ of all linear functions $x \mapsto v \cdot x$ from $\mathbb{R}^d \to \mathbb{R}$ (see Section 3.2 of their paper). However, a problem with Definition 4 for this function class is that $\mathbf{E}_{(\boldsymbol{x},\boldsymbol{y}) \sim \mathcal{D}} [\phi(\boldsymbol{y}(v \cdot \boldsymbol{x}))]$ may not have a minimum. This is not merely a technicality. In fact, for standard loss functions such as the logistic loss or the exponential loss, for any linearly separable distribution $\mathcal{D}$, $\mathbf{E}_{(\boldsymbol{x},\boldsymbol{y}) \sim \mathcal{D}} [\phi(\boldsymbol{y}(v \cdot \boldsymbol{x}))]$ does not have a minimum, informally, because scaling up $v$ increases all of the margins, which decreases all of the losses.[2] van Rooyen et al. [2015] interpret Theorem 2 as saying that for $d \geq 2$, the pair $(\phi, \mathcal{F}_{\mathrm{lin}})$ cannot be RCN-robust for any convex potential function $\phi$ (see Proposition 1 of Section 3.2 of their paper), but the fact that the minimizer typically doesn't exist in the absence of noise interferes with this interpretation. The unhinged loss also cannot be minimized over $\mathcal{F}_{\mathrm{lin}}$, since

---

[2]Implicit bias research analyzes the effect of the algorithm that drives $\mathbf{E}_{(\boldsymbol{x},\boldsymbol{y}) \sim \mathcal{D}} [\phi(\boldsymbol{y}(v \cdot \boldsymbol{x}))]$ to zero on the classification behavior of the limiting classifier. Different algorithms lead to markedly different limiting classifiers [Telgarsky, 2013, Soudry et al., 2018, Ji and Telgarsky, 2019] .

by scaling up the weight vector of any linear separator, the unhinged loss can achieve an arbitrarily large negative value.

van Rooyen et al. [2015] also consider the class $\mathcal{F}_{\mathrm{lin},r}$ of all linear functions whose weight vector has length at most $r$. They prove the following:

**Theorem 5** ([van Rooyen et al., 2015], Section 5.1). *For all $d$ and all $r > 0$, $(\phi^\star, \mathcal{F}_{\mathrm{lin},r})$ is RCN-robust.*

van Rooyen et al. [2015] further establish a number of additional properties about the unhinged loss $\phi^\star$; most of these will not concern us, but one simple property, which we now explain, is relevant to our discussion in Section 3. As above let $v \in \mathbb{R}^d$ be the minimizer of $\mathbf{E}_{(\boldsymbol{x},\boldsymbol{y})\sim\mathcal{D}}\left[\phi(\boldsymbol{y}(v \cdot \boldsymbol{x}))\right]$ subject to $||v|| \leq r$ and let $v' \in \mathbb{R}^d$ be the minimizer of $\mathbf{E}_{(\boldsymbol{x},\boldsymbol{y})\sim\overline{\mathcal{D}}_\eta}\left[\phi(\boldsymbol{y}(v' \cdot \boldsymbol{x}))\right]$ subject to $||v'|| \leq r$. van Rooyen et al. [2015] make the straightforward but useful observation that $v$ is the vector corresponding to a "nearest centroid classifier" (see Servedio [2002], [Tibshirani et al., 2002], p. 181 of [Manning et al., 2008], and Section 5.1 of Shawe-Taylor and Cristianini [2004]), i.e. we have

$$v = \alpha \mathop{\mathbf{E}}_{(\boldsymbol{x},\boldsymbol{y})\sim\mathcal{D}}[\boldsymbol{y}\boldsymbol{x}] \tag{3}$$

for a suitable rescaling factor $\alpha$, and furthermore that $v = v'$ (this holds since the values of $\mathbf{E}_{(\boldsymbol{x},\boldsymbol{y})\sim\mathcal{D}}\left[\phi^\star(\boldsymbol{y}(w \cdot \boldsymbol{x}))\right]$ and $\mathbf{E}_{(\boldsymbol{x},\boldsymbol{y})\sim\overline{\mathcal{D}}_\eta}\left[\phi^\star(\boldsymbol{y}(w \cdot \boldsymbol{x}))\right]$ are linearly related with a slope of $1 - 2\eta > 0$).

# 3   A separable learning problem where minimizing the unhinged loss on clean data yields a poor classifier

In this section we observe that while the unhinged loss $\phi^\star$ is strongly robust, there are simple linearly separable data distributions for which minimizing $\phi^\star$ over all functions in $\mathcal{F}_{\mathrm{lin},r}$ *even in the absence of random classification noise* only yields a binary classifier with an error rate of $1/2$. So while (2) is satisfied, it holds because both error rates are equal to $1/2$.

We illustrate this with the linearly separable learning scenario which is depicted in Figure 1. The distribution $\mathcal{D}$ over $\mathbb{R}^2 \times \{-1, 1\}$ is as follows: given a parameter $0 < \gamma < 0.0901$,

- $\mathcal{D}$ puts weight $1/4$ on the labeled example $x^{(1)} := (1, 0), y^{(1)} := 1$;

- $\mathcal{D}$ puts weight $1/4$ on the labeled example $x^{(2)} := (\gamma, \sqrt{1 - \gamma^2}), y^{(2)} := 1$;

- $\mathcal{D}$ puts weight $1/2$ on the labeled example $x^{(3)} := (\gamma, -2\gamma), y^{(3)} := 1$.

It is clear that $\mathcal{D}$ is linearly separable with margin $\gamma$. By (3), the vector in $\mathbb{R}^2$ which minimizes $\mathbf{E}_{(\boldsymbol{x},\boldsymbol{y})\sim\mathcal{D}}\left[\phi(\boldsymbol{y}(v \cdot \boldsymbol{x}))\right]$ points in the direction of

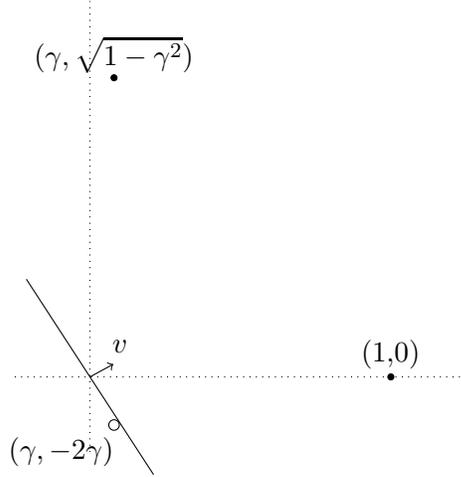$$v = (v_1, v_2) = \left(\frac{1}{4} + \frac{3\gamma}{4}, \frac{\sqrt{1 - \gamma^2}}{4} - \gamma\right).$$

Figure 1: The distribution $\mathcal{D}$ over $\mathbb{R}^2 \times \{-1, 1\}$. All examples have label $+1$; the two examples with weight $1/4$ are depicted with small filled circles and the example with weight $1/2$ is depicted with a larger unfilled circle.

For $0 < \gamma < 0.0901$ we have $v \cdot x^{(3)} < 0$ and hence $\text{sign}(v \cdot x^{(3)}) \neq y^{(3)}$, so the LHS of (2) is $1/2$. Since $v' = v$ the RHS of (2) is also $1/2$.

## 4    Implicit bias

This section includes a couple of observations about the implicit bias of algorithms that iteratively reduce the unhinged loss. Analogous results have been obtained for other loss functions [Telgarsky, 2013, Soudry et al., 2018, Ji and Telgarsky, 2019].

### 4.1    Gradient descent

Recall that the unhinged loss function is defined to be $\phi^*(z) = 1 - z$. If $\mathcal{D}$ is uniform over $(x_1, y_1), ..., (x_n, y_n) \in \mathbb{R}^d \times \{-1, 1\}$, then for any $v \in \mathbb{R}^d$ the gradient of the unhinged loss at $v$ is $-\sum_{i=1}^n y_i x_i$ (note that this does not depend on $v$). Thus, if the unhinged loss is minimized by gradient descent starting with an initial solution of 0, all iterates are multiples of $\sum_{i=1}^n y_i x_i$. If the initial solution is $v_0$, then, after $T$ updates with step size $\eta$, the weight vector is $v_0 + \eta T \sum_{i=1}^n y_i x_i$. As $T$ goes to infinity, the angle between this weight vector and $\sum_{i=1}^n y_i x_i$ goes to zero.

### 4.2    Coordinate descent

As mentioned earlier, popular boosting algorithms can be viewed as coordinate descent on a convex potential function, which works by repeatedly finding the coordinate axis with the steepest descent direction and making an update in that direction. Informally, the unhinged loss rewards increasing the margin $yv \cdot x$ of a correctly classified example $(x, y)$ as much as increasing the negative margin of an incorrectly classified example, but

increasing the margin of a correctly classified example does not make progress towards overall classification accuracy. (In contrast, the tendency of the exponential loss to place more importance on misclassified examples is key to AdaBoost's ability to boost.) If we denote the components of $x_i$ by $x_{i,1}, ..., x_{i,d}$, since the gradient of the unhinged loss is the same for all $v$, when it is minimized by coordinate descent starting from the zero weight vector all of its iterates will only have nonzero components on members of $\text{argmax}_j \sum_i y_i x_{ij}$. (If there is not a tie for the best weak learner, this will be a single component.) From a boosting point of view, a boosting algorithm based on the unhinged loss allows a weak learner to keep returning the same (weak) hypothesis.

## 5 Discussion

While the unhinged loss is noise-tolerant in a sense, minimizing it can fail to find an accurate classifier on data that is linearly separable with a large margin. On the other hand, minimizing the unhinged loss has been found to yield reasonable accuracy on natural data [see Patrini et al., 2017, Charoenphakdee et al., 2019b]. This is not entirely unexpected, since, when learning linear models, minimizing the unhinged loss is closely related to performing Naive Bayes classification [Domingos and Pazzani, 1997, Ng and Jordan, 2001], using a spherical Gaussian to model the class-conditional distributions.

Given our results, one natural goal for future work is to study whether there are conditions on potential functions which achieve an attractive tradeoff between noise-robustness and usefulness for learning (in the sense that minimizing the potential function yields an accurate classifier). Tools developed for studying Fisher consistency [Fisher, 1922], consistent loss functions [Zhang, 2004], classification calibration [Bartlett et al., 2006] and $H$-consistency [Long and Servedio, 2013] may be useful for this. In particular, it would be interesting to investigate symmetric potential functions (see e.g. Charoenphakdee et al. [2019a], Ghosh et al. [2015]) and the multiclass setting (see e.g. [Ghosh et al., 2017, Zhang and Sabuncu, 2018]) in light of this question.

## A  A minimizer exists for noisy data

In this appendix we show that for all convex potential functions $\phi$, all finite-covariance distributions $\mathcal{D}$ over $\mathbb{R}^d$, and all $\eta \in (0, 1/2)$, the function $P_{\phi,\overline{\mathcal{D}}_\eta}$ has a minimum. (Recall from (1) that $P_{\phi,\overline{\mathcal{D}}_\eta}(v) := \mathbf{E}_{(\boldsymbol{x},\boldsymbol{y}) \sim \overline{\mathcal{D}}_\eta} \left[ \phi(\boldsymbol{y}(v \cdot \boldsymbol{x})) \right].$)

We recall some useful background.

**Definition 6.** *For any $a$, the set $\{x : f(x) \leq a\}$ is a* level set *for $f : \mathbb{R}^d \to \mathbb{R}$.*

**Definition 7.** *A nonzero vector $u \in \mathbb{R}^d$ is a* direction of recession *for a function $f$ if, for all nonempty level sets $L$ of $f$, there exists some $x_0$ such that $x_0 + \lambda u \in L$ for all $\lambda \geq 0$. (Informally, all non-empty level sets of $f$ extend infinitely in the $u$ direction.)*

**Lemma 8** ([Rockafellar, 2015], Theorem 27.1). *The set of minima of a continuous convex function $f$ is nonempty and bounded iff $f$ does not have any direction of recession.*

8

**Definition 9.** *Say that a probability distribution $\mathcal{D}$ over $\mathbb{R}^d \times \{-1, 1\}$ has* finite covariance *if, for all unit length $u \in \mathbb{R}^d$, $\mathbf{E}_{(\boldsymbol{x},\boldsymbol{y}) \sim \mathcal{D}}[(u \cdot \boldsymbol{x})^2]$ exists.*

Now we are ready to analyze $P_{\phi, \overline{\mathcal{D}}_\eta}$.

**Proposition 10.** *For any convex potential function $\phi$, for any $\eta \in (0, 1/2)$, for any finite-covariance distribution $\mathcal{D}$ over $\mathbb{R}^d$, $P_{\phi, \overline{\mathcal{D}}_\eta}$ has a minimum.*

*Proof.* First, we may assume without loss of generality that, for all unit length $u \in \mathbb{R}^d$,

$$\mathbf{E}_{(\boldsymbol{x},\boldsymbol{y}) \sim \mathcal{D}}[|u \cdot \boldsymbol{x}|] > 0 \tag{4}$$

since otherwise $u \cdot \boldsymbol{x} = 0$ almost surely, and $P_{\phi, \overline{\mathcal{D}}_\eta}(v)$ is unaffected by projecting $v$ onto the subspace of $\mathbb{R}^d$ orthogonal to $u$.

Assume for contradiction that some $u$ is a direction of recession for $P_{\phi, \overline{\mathcal{D}}_\eta}$. For any $x_0 \in \mathbb{R}^d$ and $\lambda \geq 0$, we have

$$
\begin{aligned}
&P_{\phi, \overline{\mathcal{D}}_\eta}(x_0 + \lambda u) \\
&= \mathbf{E}_{(\boldsymbol{x},\boldsymbol{y}) \sim \overline{\mathcal{D}}_\eta}[\phi(\boldsymbol{y}((x_0 + \lambda u) \cdot \boldsymbol{x}))] && \text{(def. of } P_{\phi, \overline{\mathcal{D}}_\eta}) \\
&= \mathbf{E}_{(\boldsymbol{x},\boldsymbol{y}) \sim \mathcal{D}}[\eta\phi(-\boldsymbol{y}((x_0 + \lambda u) \cdot \boldsymbol{x})) + (1 - \eta)\phi(\boldsymbol{y}((x_0 + \lambda u) \cdot \boldsymbol{x}))] && \text{(def. of } \overline{\mathcal{D}}_\eta) \\
&\geq \mathbf{E}_{(\boldsymbol{x},\boldsymbol{y}) \sim \mathcal{D}}[\eta \max_{\tilde{y} \in \{-1,1\}} \phi(\tilde{y}((x_0 + \lambda u) \cdot \boldsymbol{x}))] && \text{(since } \eta < 1/2 \text{ and } \phi \geq 0) \\
&= \mathbf{E}_{(\boldsymbol{x},\boldsymbol{y}) \sim \mathcal{D}}[\eta\phi(-|(x_0 + \lambda u) \cdot \boldsymbol{x}|)] && \text{(since } \phi \text{ is nonincreasing)} \\
&\geq \eta \mathbf{E}_{(\boldsymbol{x},\boldsymbol{y}) \sim \mathcal{D}}[\phi(0) - \phi'(0) \cdot |(x_0 + \lambda u) \cdot \boldsymbol{x}|] && \text{(since } \phi \text{ is convex and } \phi'(0) < 0) \\
&\geq \eta\left(\phi(0) - \phi'(0)\lambda \mathbf{E}_{(\boldsymbol{x},\boldsymbol{y}) \sim \mathcal{D}}[|u \cdot \boldsymbol{x}|] + \phi'(0) \mathbf{E}_{(\boldsymbol{x},\boldsymbol{y})}[|x_0 \cdot \boldsymbol{x}|]\right). \\
&&& \text{(triangle inequality, } \phi'(0) < 0)
\end{aligned}
$$

Thus $\lim_{\lambda \to \infty} P_{\phi, \overline{\mathcal{D}}_\eta}(x_0 + \lambda u) = \infty$, which contradicts the assumption that $u$ is a direction of recession for $P_{\phi, \overline{\mathcal{D}}_\eta}$. $\square$

# References

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.

Nontawat Charoenphakdee, Jongyeong Lee, and Masashi Sugiyama. On symmetric losses for learning from corrupted labels. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 961–970. PMLR, 2019a.

Nontawat Charoenphakdee, Jongyeong Lee, and Masashi Sugiyama. On symmetric losses for learning from corrupted labels. In *International Conference on Machine Learning*, pages 961–970. PMLR, 2019b.

T.G. Dietterich. An experimental comparison of three methods for constructing ensembles of decision trees: bagging, boosting, and randomization. *Machine Learning*, 40 (2):139–158, 2000.

C. Domingo and O. Watanabe. MadaBoost: a modified version of AdaBoost. In *Proceedings of the Thirteenth Annual Conference on Computational Learning Theory (COLT)*, pages 180–189, 2000.

Pedro Domingos and Michael Pazzani. On the optimality of the simple bayesian classifier under zero-one loss. *Machine learning*, 29(2):103–130, 1997.

Ronald A Fisher. On the mathematical foundations of theoretical statistics. *Philosophical transactions of the Royal Society of London. Series A, containing papers of a mathematical or physical character*, 222(594-604):309–368, 1922.

Y. Freund and R. Schapire. Experiments with a new boosting algorithm. In *Proceedings of the Thirteenth International Conference on Machine Learning*, pages 148–156, 1996.

Y. Freund and R. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55(1):119–139, 1997.

J. Friedman, T. Hastie, and R. Tibshirani. Additive logistic regression: a statistical view of boosting. *Annals of Statistics*, 28(2):337–407, 1998.

Claudio Gentile and Manfred K. Warmuth. Linear hinge loss and average margin. In *Advances in Neural Information Processing Systems 11*, pages 225–231. The MIT Press, 1998.

Aritra Ghosh, Naresh Manwani, and P. S. Sastry. Making risk minimization tolerant to label noise. *Neurocomputing*, 160:93–107, 2015.

Aritra Ghosh, Himanshu Kumar, and PS Sastry. Robust loss functions under label noise for deep neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 31, 2017.

Ziwei Ji and Matus Telgarsky. The implicit bias of gradient descent on nonseparable data. In *Conference on Learning Theory*, pages 1772–1798. PMLR, 2019.

Phil Long and Rocco Servedio. Consistency versus realizable h-consistency for multiclass classification. In *International Conference on Machine Learning*, pages 801–809. PMLR, 2013.

Philip M Long and Rocco A Servedio. Random classification noise defeats all convex potential boosters. *Machine learning*, 78(3):287–304, 2010.

R. Maclin and D. Opitz. An empirical evaluation of bagging and boosting. In *Proceedings of the Fourteenth National Conference on Artificial Intelligence and Ninth Innovative Applications of Artificial Intelligence Conference (AAAI/IAAI)*, pages 546–551, 1997.

Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. *Introduction to Information Retrieval*. Cambridge University Press, 2008.

Naresh Manwani and P. S. Sastry. Noise tolerance under risk minimization. *IEEE Trans. Cybern.*, 43(3):1146–1151, 2013.

L. Mason, J. Baxter, P. Bartlett, and M. Frean. Boosting algorithms as gradient descent. In *Advances in Neural Information Processing Systems (NIPS)*, pages 512–518, 1999.

Nagarajan Natarajan, Inderjit S. Dhillon, Pradeep Ravikumar, and Ambuj Tewari. Learning with noisy labels. In *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013*, pages 1196–1204, 2013.

A. Y. Ng and M. I. Jordan. On discriminative vs. generative classifiers: A comparison of logistic regression and naive bayes. *NIPS*, 2001.

Giorgio Patrini, Alessandro Rozza, Aditya Krishna Menon, Richard Nock, and Lizhen Qu. Making deep neural networks robust to label noise: A loss correction approach. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1944–1952, 2017.

Ralph Tyrell Rockafellar. *Convex analysis*. Princeton university press, 2015.

Rocco A. Servedio. Perceptron, winnow, and PAC learning. *SIAM J. Comput.*, 31(5): 1358–1369, 2002.

Jonathan Shawe-Taylor and Nello Cristianini. *Kernel Methods for Pattern Analysis*. Cambridge University Press, 2004.

Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit bias of gradient descent on separable data. *The Journal of Machine Learning Research*, 19(1):2822–2878, 2018.

Matus Telgarsky. Margins, shrinkage, and boosting. In *ICML*, volume 28 of *JMLR Workshop and Conference Proceedings*, pages 307–315, 2013.

Robert Tibshirani, Trevor Hastie, Balasubramanian Narasimhan, and Gilbert Chu. Diagnosis of multiple cancer types by shrunken centroids of gene expression. *Proceedings of the National Academy of Sciences*, 99(10):6567–6572, 2002.

Brendan van Rooyen, Aditya Menon, and Robert C. Williamson. Learning with Symmetric Label Noise: The Importance of Being Unhinged. In *Advances in Neural Information Processing Systems 28*, 2015.

Tong Zhang. Statistical behavior and consistency of classification methods based on convex risk minimization. *The Annals of Statistics*, 32(1):56–85, 2004.

Zhilu Zhang and Mert Sabuncu. Generalized cross entropy loss for training deep neural networks with noisy labels. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.