
It begins with a boundary: A geometric view on probabilistically robust learning

Leon Bungert

Technical University of Berlin
l.bungert@tu-berlin.de

Nicolás García Trillos

University of Wisconsin-Madison
garciatrillo@wisc.edu

Matt Jacobs

Purdue University
jacob225@purdue.edu

Daniel McKenzie

Colorado School of Mines
dmckenzie@mines.edu

Đorđe Nikolić

University of California Santa Barbara
nikolic@math.ucsb.edu

Qingsong Wang

University of Utah
qswang@math.utah.edu

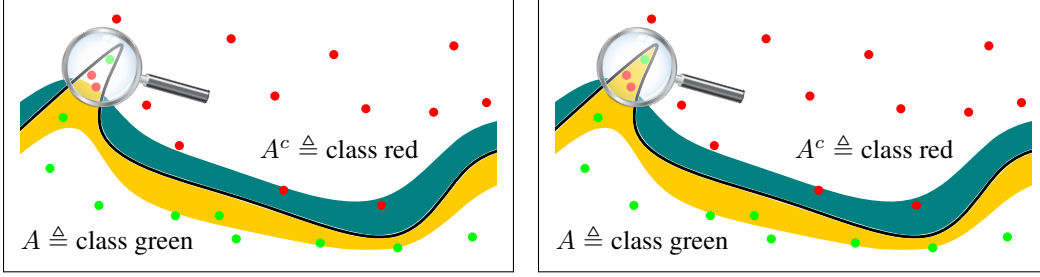
Abstract

Although deep neural networks have achieved super-human performance on many classification tasks, they often exhibit a worrying lack of robustness towards adversarially generated examples. Thus, considerable effort has been invested into reformulating Empirical Risk Minimization (ERM) into an adversarially robust framework. Recently, attention has shifted towards approaches which interpolate between the robustness offered by adversarial training and the higher clean accuracy and faster training times of ERM. In this paper, we take a fresh and geometric view on one such method—Probabilistically Robust Learning (PRL) [Robey et al., 2022]. We propose a geometric framework for understanding PRL, which allows us to identify a subtle flaw in its original formulation and to introduce a family of probabilistic nonlocal perimeter functionals to address this. We prove existence of solutions using novel relaxation methods and study properties as well as local limits of the introduced perimeters.

1 Introduction

The fragility of DNN-based classifiers in the face of adversarial examples [Goodfellow et al., 2014, Chen et al., 2017, Qin et al., 2019, Cai et al., 2021] and distributional shifts [Quinoñero Candela et al., 2008, Hendrycks et al., 2021] is by now nearly as familiar as their successes. In light of this, a multitude of works (see Section 1.4) propose replacing standard Empirical Risk Minimization (ERM) [Vapnik, 1999] with a more robust alternative (see, e.g., Madry et al. [2017]). Unfortunately there is no free lunch: robust classifiers frequently exhibit degraded performance on clean data and significantly longer training times [Tsipras et al., 2018]. Consequently, identifying frameworks which balance performance and robustness is of pressing interest to the Machine Learning (ML) community, and over the past several years many such frameworks have been proposed [Zhang et al., 2019, Wang et al., 2020, Robey et al., 2022]. Moreover, it is crucial that the mechanism by which such frameworks balance these competing aims be understood.

Beginning with the Probabilistically Robust Learning (PRL) of Robey et al. [2022] we analyze such frameworks geometrically. This perspective reveals a subtle, paradoxical aspect of PRL: sometimes the adversary modeled by this framework corrects, instead of exploits, the learner! Fortunately, the



(a) Robey et al. [2022]: The probabilistically non-robust region (**magnified**) reduces the loss.

(b) Our model: The probabilistically non-robust region is correctly identified and penalized.

Figure 1: Penalization effect of the original model [Robey et al., 2022] (**left**) and ours (**right**): The solid black is the decision boundary of a non-robust classifier induced by the set A . Both models penalize the numbers of green points in the yellow region and red points in the teal region. However, the original model *favors non-robust regions* of A for which most perturbations correct the class. Our model identifies this region as non-robust and penalizes it accordingly.

geometric perspective we propose suggests a natural remedy which leads to an interpretation of the corrected PRL as regularized ERM where a certain nonlocal notion of length (or perimeter) of the decision boundary acts as a regularizer. We exemplify this correction in Figure 1. The interpretation of PRL as perimeter-regularized ERM leads us to further generalizations, and we provide a novel view of the Conditional Value at Risk (CVaR) relaxation of PRL proposed by Robey et al. [2022].

1.1 From empirical risk minimization to robustness

Given an input space \mathcal{X} , an output space \mathcal{Y} , a probability measure $\mu \in \mathcal{P}(\mathcal{X} \times \mathcal{Y})$, a loss function $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$, and a hypothesis class \mathcal{H} , the standard risk minimization problem is

$$\inf_{h \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mu} [\ell(h(x), y)]. \quad (1)$$

For training classifiers which are robust against adversarial attacks Goodfellow et al. [2014], Madry et al. [2017] suggested adversarial training:

$$\inf_{h \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mu} \left[\sup_{x' \in B_\varepsilon(x)} \ell(h(x'), y) \right]. \quad (2)$$

Here \mathcal{X} is assumed to have the structure of a metric space and $B_\varepsilon(x)$ for $\varepsilon \geq 0$ denotes the (open or closed) ball of radius ε around x .

The recent work by Robey et al. [2022] offered an alternative to adversarial training in order to reduce the (in general) large trade-off between accuracy and robustness inherent in (2), see Tsipras et al. [2018], Robey et al. [2022] for discussion. Instead of requiring classifiers to be robust to *all* available attacks around a point x —as enforced through the supremum in (2)—one may consider a less stringent notion of robustness, only requiring classifiers to be robust to $100 \times (1 - p)\%$ of possible attacks when attacks are drawn from a certain distribution \mathfrak{p}_x centered at x . For this, the authors introduced the so-called p -ess sup operator for $p \in [0, 1)$ and suggested replacing (2) by

$$\inf_{h \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mu} \left[p\text{-ess sup}_{x' \sim \mathfrak{p}_x} \ell(h(x'), y) \right], \quad (3)$$

where $\{\mathfrak{p}_x\}_{x \in \mathcal{X}}$ is a family of probability distributions. The prototypical example to keep in mind for $\mathcal{X} = \mathbb{R}^d$ is the uniform distribution over the ε -ball around x , i.e., $\mathfrak{p}_x := \text{Unif}(B_\varepsilon(x))$, which is particularly relevant when dealing with adversarial attacks on image classifiers.

For a probability distribution \mathfrak{p} and a function f , the quantity $p\text{-ess sup}_{x' \sim \mathfrak{p}} f(x')$ is defined as the smallest value $t \in \mathbb{R}$ such that the probability of a randomly chosen point $x' \sim \mathfrak{p}$ satisfying $f(x') > t$ is smaller than p , which reduces to the usual essential supremum of f with respect to \mathfrak{p} if $p = 0$:

$$p\text{-ess sup}_{x' \sim \mathfrak{p}} f(x') := \inf \{t \in \mathbb{R} : \mathbb{P}_{x' \sim \mathfrak{p}} [f(x') > t] \leq p\}.$$

To better understand the model (3) we temporarily restrict our attention to binary classification (i.e., $\mathcal{Y} = \{0, 1\}$) using indicator functions of admissible sets (i.e., $\mathcal{H} := \{\mathbf{1}_A : A \in \mathcal{A}\}$). Note that we identify the two expressions $\mathbf{1}_A(x) = \mathbf{1}_{x \in A}$. We focus on the 0-1 loss $\ell(\tilde{y}, y) = \mathbf{1}_{\tilde{y} \neq y}$ which equals one if $y \neq \tilde{y}$ and zero otherwise. In this scenario (1) reduces to the geometric problem

$$\inf_{A \in \mathcal{A}} \left\{ R_{\text{std}}(A) := \mathbb{E}_{(x,y) \sim \mu} [y \mathbf{1}_{x \in A^c} + (1-y) \mathbf{1}_{x \in A}] \right\}, \quad (4)$$

and minimizers are called Bayes classifiers. Similarly, adversarial training (2) can be rewritten as

$$\inf_{A \in \mathcal{A}} \left\{ R_{\text{adv}}(A) := \mathbb{E}_{(x,y) \sim \mu} \left[y \mathbf{1}_{x \in (A^c)^{\oplus \varepsilon}} + (1-y) \mathbf{1}_{x \in A^{\oplus \varepsilon}} \right] \right\}, \quad (5)$$

where for a set $A \in \mathcal{A}$ its fattening by ε -balls is defined as $A^{\oplus \varepsilon} := \bigcup_{x \in A} B_\varepsilon(x)$. Hence (5) enforces that all points with distance at most ε to the decision boundary be adversarially robust.

On the other hand the PRL model (3) reduces to

$$\inf_{A \in \mathcal{A}} \left\{ R_{\text{prob}}(A) := \mathbb{E}_{(x,y) \sim \mu} \left[y \mathbf{1}_{\mathbb{P}_{x' \sim \mathbf{p}_x}[x' \in A^c] > p} + (1-y) \mathbf{1}_{\mathbb{P}_{x' \sim \mathbf{p}_x}[x' \in A] > p} \right] \right\}, \quad (6)$$

where $A^{\oplus \varepsilon}$ is replaced by a ‘‘probabilistic fattening’’, i.e., one considers the set of all x for which the probability that a neighboring point sampled from \mathbf{p}_x lies inside A is larger than p . To the best of our knowledge, existence of solutions for (6) or even (3) has not been proved so far.

1.2 Geometric modification of probabilistically robust learning

To motivate our geometric modification of the PRL model from Robey et al. [2022], it is insightful to investigate the regularization effect that PRL has compared to standard risk minimization. We let $\rho_i(\bullet) := \mu(\bullet \times \{i\})$ denote the non-normalized conditional distributions of the points with label i . Subtracting the standard risk in (4) from the one in (6) and disintegrating using ρ_0 and ρ_1 we obtain

$$\begin{aligned} & R_{\text{prob}}(A) - R_{\text{std}}(A) \\ &= \int_{\mathcal{X}} \mathbf{1}_{\mathbb{P}_{x' \sim \mathbf{p}_x}[x' \in A] > p} - \mathbf{1}_{x \in A} d\rho_0(x) + \int_{\mathcal{X}} \mathbf{1}_{\mathbb{P}_{x' \sim \mathbf{p}_x}[x' \in A^c] > p} - \mathbf{1}_{x \in A^c} d\rho_1(x). \end{aligned} \quad (7)$$

We highlight that this expression *does not constitute a non-negative functional of A* . Hence the loss function in (6) is not a regularized version of the standard risk (4) and in fact can be strictly smaller. This observation reveals a subtle flaw in the approach of Robey et al. [2022]: Points which lie in thin or spike-like regions of A penetrating the other class and that are more likely to have the label zero than the label one (meaning they lie in the set $\{\rho_0 > \rho_1\}$) yield negative contributions in (7) and are hence *favoured*. Such a scenario is visualized on the left side of Figure 1. From an adversarial perspective this means that points which are already misclassified are attacked nevertheless, which can lead to the bizarre situation that the adversary helps the learner by putting these points in the correct class with high probability, thereby reducing both adversarial robustness and clean accuracy.

We fix this by designing a probabilistically robust risk as non-negative regularization of the standard risk. For this we define probabilistic perimeter functionals which only penalize points which are classified correctly *and* admit a large portion of attacks around them, see the right side of Figure 1.

1.3 Our contributions

Our main contributions are the following:

- We address the geometric limitation of the model by Robey et al. [2022] by introducing a family of perimeter regularizations.
- We prove existence of soft and hard binary classifiers under weak conditions on the family of perimeters and hypothesis classes, using novel relaxation techniques.
- We investigate the relationship between the introduced family of perimeters and local perimeters in Euclidean space for small adversarial budgets.
- We extend our models to encompass general loss functions and hypothesis classes. Our numerical experiments demonstrate that our geometric correction can enhance the adversarial robustness of probabilistically robust classifiers without compromising clean accuracy.

1.4 Related work

Adversarial training was developed by Goodfellow et al. [2014], Madry et al. [2017] as an approach to train networks that are less sensitive to adversarial attacks. Shafahi et al. [2019] reduced its computational complexity by reusing gradients from the backpropagation when training neural networks. Wong et al. [2020] showed that training with noise perturbations followed by a single signed gradient ascent (FGSM) step can be on par with adversarial training while being much cheaper. This approach was picked up and improved upon by Andriushchenko and Flammarion [2020] based on gradient alignment. Different authors also investigated test-time robustification of pretrained classifiers using randomized smoothing [Cohen et al., 2019] or geometric / gradient-based approaches [Schwinn et al., 2021, 2022]. While some of the previous models use a combination of random perturbations and gradient-based adversarial attacks to robustify classifiers, Robey et al. [2022] proposed probabilistically robust learning, which is entirely based on random perturbations. PRL aims to interpolate between clean and adversarial accuracy and enjoys the favorable sample complexity of vanilla empirical risk minimization; see also Raman et al. [2023] for more insights on this issue. Connections between adversarial training and local perimeter regularization of decision boundaries were explored by García Trillos and Murray [2022] and then rigorously tied by Bungert and Stinson [2022]. Our work is in line with a series of papers [Pydi and Jog, 2021, Awasthi et al., 2021a,b, Frank and Niles-Weed, 2022, Frank, 2022, Bungert et al., 2023, García Trillos et al., 2023] that explore the existence of solutions to adversarial training problems in different settings. These existence proofs involve dealing with different kinds of measurability issues, depending on whether open or closed balls $B_\varepsilon(x)$ are used in the attack model. For open balls one can work with the Borel σ -algebra $\mathcal{A} = \mathfrak{B}(\mathcal{X})$ [Bungert et al., 2023], whereas closed balls require the use of the universal σ -algebra to make sure that $A^{\oplus\varepsilon}$ is measurable [Pydi and Jog, 2021, Awasthi et al., 2021a,b]. Recently, these results were improved by García Trillos et al. [2023] who also proved for the case of multi-class classification that even for the closed ball model Borel measurable classifiers (albeit not necessarily indicator functions of measurable sets) exist and that for all but countably many values of the adversarial budget $\varepsilon > 0$ the open and the closed ball models have the same minimal value.

2 Geometry and existence of probabilistically robust classifiers

2.1 The binary classification setting with 0-1 loss

In this section we shall introduce our baseline model, which is based on a suitable geometric regularization of the standard risk. Later we shall embed it into a family of models. For clarity we first discuss hard classifiers (characteristic functions of sets) and then soft classifiers (functions with values in $[0, 1]$). The generalization to general models and loss functions is postponed to Section 3.

We start by defining the *probabilistic perimeter* for $p \in [0, 1]$ of an admissible set $A \in \mathcal{A}$ as follows:

$$\begin{aligned} \text{ProbPer}(A) := & \rho_0(\{x \in A^c : \mathbb{P}_{x' \sim p_x}[x' \in A] > p\}) \\ & + \rho_1(\{x \in A : \mathbb{P}_{x' \sim p_x}[x' \in A^c] > p\}). \end{aligned} \quad (8)$$

$\text{ProbPer}(A)$ penalizes correctly classified points x for which more than $100 \times p$ % of their neighbors, sampled from p_x , constitute an attack. The perimeter can be rewritten in integral form:

$$\text{ProbPer}(A) = \int_{\mathcal{X}} \mathbf{1}_{x \in A \vee \mathbb{P}_{x' \sim p_x}[x' \in A] > p} - \mathbf{1}_{x \in A} d\rho_0(x) \quad (9)$$

$$\begin{aligned} & + \int_{\mathcal{X}} \mathbf{1}_{x \in A^c \vee \mathbb{P}_{x' \sim p_x}[x' \in A^c] > p} - \mathbf{1}_{x \in A^c} d\rho_1(x) \\ & = \int_{\mathcal{X}} \mathbf{1}_{x \in A^c} \mathbf{1}_{\mathbb{P}_{x' \sim p_x}[x' \in A] > p} d\rho_0(x) + \int_{\mathcal{X}} \mathbf{1}_{x \in A} \mathbf{1}_{\mathbb{P}_{x' \sim p_x}[x' \in A^c] > p} d\rho_1(x). \end{aligned} \quad (10)$$

The first reformulation (9) should be compared to (7), while the one in (10) will be useful later on. The use of the term perimeter to describe the functional ProbPer will become more apparent shortly in Section 2.4, and at this point it is worth highlighting that ProbPer is always a non-negative quantity. This motivates introducing the following regularized risk

$$\text{ProbR}(A) := R_{\text{std}}(A) + \text{ProbPer}(A), \quad A \in \mathcal{A}. \quad (11)$$

Our first theorem states that ProbR equals the expected maximum of the sample-wise standard risk and the probabilistically robust risk from Robey et al. [2022], cf. (4) and (6).

Theorem 1. For all $A \in \mathcal{A}$ it holds that

$$\text{ProbR}(A) = \mathbb{E}_{(x,y) \sim \mu} \left[\max \left\{ \mathbf{1}_{\mathbb{P}_{x' \sim p_x}[\mathbf{1}_A(x') \neq y] > p}, \mathbf{1}_{\mathbf{1}_A(x) \neq y} \right\} \right]. \quad (12)$$

The interpretation of the statement of this theorem in the light of Figure 1 is clear: Only if a point x is correctly classified—meaning $\mathbf{1}_{\mathbf{1}_A(x) \neq y} = 0$ —the probabilistically robust regularization kicks in through the first term in the maximum. Points which are incorrectly classified will always be penalized even if most attacks correct the label, i.e., if $\mathbf{1}_{\mathbb{P}_{x' \sim p_x}[\mathbf{1}_A(x') \neq y] > p} = 0$. Thus, minimizing ProbR instead of R_{prob} corrects the pathology identified in Section 1.2.

2.2 Extensions in the binary classification setting

Given the formula of ProbPer in (10), several natural extensions suggest themselves. E.g., one may replace the indicator function $\mathbf{1}_{t > p}$ with a different function $\Psi(t)$ to define other notions of *perimeter*

$$\begin{aligned} \text{ProbPer}_\Psi(A) := & \int_{\mathcal{X}} \mathbf{1}_{x \in A^c} \Psi(\mathbb{P}_{x' \sim p_x}[x' \in A]) \, d\rho_0(x) \\ & + \int_{\mathcal{X}} \mathbf{1}_{x \in A} \Psi(\mathbb{P}_{x' \sim p_x}[x' \in A^c]) \, d\rho_1(x) \end{aligned} \quad (13)$$

as well as their corresponding probabilistically robust losses

$$\text{ProbR}_\Psi(A) := R_{\text{std}}(A) + \text{ProbPer}_\Psi(A). \quad (14)$$

For $\Psi(t) := \mathbf{1}_{t > p}$ the perimeter ProbPer_Ψ reduces to ProbPer and so do the associated risks. Of particular interest is $\Psi_p(t) := \min\{t/p, 1\}$ —the smallest concave function that lies above $\Psi(t) = \mathbf{1}_{t > p}$ —which will allow us to develop deep connections between the theoretical and computational aspects of probabilistically robust learning. Our relaxation using the function Ψ is very similar to the one by Raman et al. [2023] who proved PAC learnability if Ψ is Lipschitz, see Appendix A.6 for more details. In order to rigorously study ProbR_Ψ we first make our setting precise.

Assumption 1. We let \mathcal{X} be a set and $\mathcal{A} \subset 2^{\mathcal{X}}$ be a σ -algebra. We assume that:

- $(\mathcal{X} \times \mathcal{Y}, \mathcal{A} \otimes 2^{\{0,1\}}, \mu)$ is a probability space;
- $(\mathcal{X}, \mathcal{A}, \rho)$ is a probability space, where we define $\rho(\bullet) := \mu(\bullet \times \{0, 1\})$;
- $\{p_x\}_{x \in \mathcal{X}}$ is a family such that $(\mathcal{X}, \mathcal{A}, p_x)$ is a probability space for ρ -almost every $x \in \mathcal{X}$.

The following theorem establishes existence of minimizers of the risk ProbR_Ψ for concave and non-decreasing functions Ψ . This existence result is astonishing since the standard method of calculus of variations is not directly applicable, with the reason being that problem (15) does not provide enough compactness for lower semicontinuity of the perimeter functional ProbPer_Ψ . Instead, the proof is based on convex relaxations to soft classifiers where we use a lower semicontinuous surrogate functional and a total variation defined through a coarea formula which—if Ψ is concave and non-decreasing—lower-bounds the surrogate.

Theorem 2. Suppose $\Psi : [0, 1] \rightarrow [0, 1]$ is concave and non-decreasing, and that Assumption 1 holds. Then, there exists a solution to the problem

$$\inf_{A \in \mathcal{A}} \text{ProbR}_\Psi(A). \quad (15)$$

Furthermore, ProbR_Ψ can also be interpreted as a sample-wise maximum, analogous to Theorem 1.

Theorem 3. For all $A \in \mathcal{A}$ and measurable $\Psi : [0, 1] \rightarrow [0, 1]$ it holds

$$\begin{aligned} \text{ProbR}_\Psi(A) &= R_{\text{std}}(A) + \text{ProbPer}_\Psi(A) \\ &= \mathbb{E}_{(x,y) \sim \mu} \left[\max \left\{ \Psi(\mathbb{P}_{x' \sim p_x}[\mathbf{1}_A(x') \neq y]), \mathbf{1}_{\mathbf{1}_A(x) \neq y} \right\} \right]. \end{aligned}$$

Note that for the non-concave function $\Psi(t) = \mathbf{1}_{t > p}$ an existence proof along the lines of Theorem 2 is not available since certain relaxation techniques therein rely on concavity of Ψ . However, in the next section we shall provide an existence theorem for soft classifiers which is valid for very general functions Ψ , including $\Psi(t) = \mathbf{1}_{t > p}$.

2.3 Extension to soft classifiers

Another natural extension features “soft classifiers” instead of indicator functions of admissible sets. Such classifiers are particularly relevant since they include the neural network based models with `Softmax` activation in the last layer which are used in practice. We start by defining a suitable regularization functional for soft classifiers. Given a \mathcal{A} -measurable function $u : \mathcal{X} \rightarrow [0, 1]$ we define

$$J_\Psi(u) := \int_{\mathcal{X}} (1 - u(x)) \Psi(\mathbb{E}_{x' \sim p_x} [u(x')]) \, d\rho_0(x) + \int_{\mathcal{X}} u(x) \Psi(\mathbb{E}_{x' \sim p_x} [1 - u(x')]) \, d\rho_1(x) \quad (16)$$

which satisfies $J_\Psi(\mathbf{1}_A) = \text{ProbPer}_\Psi(A)$ for every choice of Ψ . Hence, it is a natural generalization of the perimeter to soft classifiers and one could call J_Ψ a total variation. However, it is neither positively homogeneous nor convex so this name would be misleading. Instead, for the proof of Theorem 2 we shall construct a suitable total variation functional ProbTV_Ψ which upper-bounds J_Ψ .

The next theorem asserts existence of soft classifiers for the regularized risk minimization using J_Ψ for very general functions Ψ and hypothesis classes \mathcal{H} , requiring only that Ψ be lower semicontinuous. For example, every continuous function and also $\Psi(t) = \mathbf{1}_{t > p}$ for $p \in [0, 1]$ satisfies this. The existence theorem is valid for all hypotheses classes which are closed in a suitable sense.

Theorem 4. *Under Assumption 1, for every lower semicontinuous function $\Psi : [0, 1] \rightarrow [0, 1]$, and whenever \mathcal{H} is a weak-* closed hypothesis class of \mathcal{A} -measurable functions $u : \mathcal{X} \rightarrow [0, 1]$ in the sense of Definition 1 in the appendix, there exists a solution to the problem*

$$\inf_{u \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mu} [|u(x) - y|] + J_\Psi(u).$$

Example 1. Let us consider three interesting hypothesis classes of weak-* closed classifiers for which Theorem 4 applies. More detailed explanations are given in Appendix A.8.

1. The simplest such class \mathcal{H} is the class of all \mathcal{A} -measurable soft classifiers $u : \mathcal{X} \rightarrow [0, 1]$ which could be referred to as *agnostic* classifiers since they are not parametrized.
2. An example with more practical relevance is the class of (feedforward or residual) neural networks defined on the unit cube $\mathcal{X} := [-1, 1]^d$ with uniformly bounded parameters

$$\mathcal{H} := \left\{ \Phi_L \circ \dots \circ \Phi_1 : [-1, 1]^d \rightarrow [0, 1] : \Phi_l(\bullet) = A_l \bullet + \sigma_l(W_l \bullet + b_l), \right. \\ \left. \|(A_l, W_l, b_l)\| \leq C \, \forall l \in \{1, \dots, L\} \right\},$$

where we assume that the activations $\sigma_l : \mathbb{R} \rightarrow \mathbb{R}$ are continuous. Note that the boundedness of the weights cannot be relaxed. To see this, consider the (very simplistic) neural network $u_n(x) = \tanh(w_n x)$ for $x \in [-1, 1]$ and $w_n \in \mathbb{R}$. For $w_n \rightarrow \infty$ it is easy to see that u_n converges to $u(x) := \text{sign}(x)$ which does not lie in the same hypothesis class.

3. Finally, one can also consider the class of hard linear classifiers on \mathbb{R}^d . Letting $\theta(t) := \mathbf{1}_{t > 0}$ denote the Heaviside function, this class is given by

$$\mathcal{H} := \left\{ \theta(w \cdot x + b) : w \in \mathbb{R}^d, |w| = 1, b \in [-\infty, \infty] \right\},$$

where one interprets $u(x) := \theta(w \cdot x + b)$ as $u \equiv 1$ if $b = \infty$ and $u \equiv 0$ if $b = -\infty$. If the distributions ρ_0, ρ_1 , and p_x are sufficiently nice, then \mathcal{H} has the desired closedness property.

2.4 Properties and asymptotics of ProbPer_Ψ

In this section we shall discuss the interpretation of the functional ProbPer_Ψ defined in (13) as a *perimeter*. We do this in two ways.

First, we focus on the case where Ψ is concave and non-decreasing and prove that ProbPer_Ψ is a *submodular functional*. If, in addition, Ψ is assumed to satisfy $\Psi(0) = 0$, then $\text{ProbPer}_\Psi(\mathcal{X}) = \text{ProbPer}_\Psi(\emptyset) = 0$. Following Chambolle et al. [2015], for Ψ satisfying these properties one can interpret ProbPer_Ψ as a generalized perimeter, i.e., a functional that can be used to measure the “size” of the boundary of a set. In Appendix A.3 we introduce ProbPer_Ψ ’s induced (generalized) total variation and use it in the proof of Theorem 2; note that, as discussed by Bungert et al. [2023], the adversarial problem (5) also induces a generalized perimeter with associated total variation.

Theorem 5. *If $\Psi(0) = 0$, then $\text{ProbPer}_\Psi(\mathcal{X}) = \text{ProbPer}_\Psi(\emptyset) = 0$. If Ψ is concave and non-decreasing, then the functional ProbPer_Ψ is submodular, meaning that*

$$\text{ProbPer}_\Psi(A \cup B) + \text{ProbPer}_\Psi(A \cap B) \leq \text{ProbPer}_\Psi(A) + \text{ProbPer}_\Psi(B) \quad \forall A, B \in \mathcal{A}.$$

Example 2. For $\Psi(t) = t$ our perimeter reduces to the perimeter on the *random walk space* $(\mathcal{X}, \mathfrak{p})$, introduced by Mazón et al. [2020]: $\text{ProbPer}_\Psi(A) = \int_{\mathcal{X} \setminus A} \int_A \text{d}\mathfrak{p}_x \text{d}\rho_0(x) + \int_A \int_{\mathcal{X} \setminus A} \text{d}\mathfrak{p}_x \text{d}\rho_1(x)$.

Second, we consider more general Ψ and show that ProbPer_Ψ is related to a standard *local* perimeter when the adversarial budget approaches zero; for the case of adversarial training such a connection was proved by Bungert and Stinson [2022] where the authors utilized the notion of Gamma-convergence of functionals. We take a first step in this direction by proving that for sufficiently smooth sets the probabilistic perimeter converges to a local one if the family of probability distributions \mathfrak{p}_x localizes suitably. For example, one could think of $\mathfrak{p}_x := \text{Unif}(B_\varepsilon(x))$, which converges to a point mass at x if $\varepsilon \rightarrow 0$. To make our setting precise, we pose the following general assumption:

Assumption 2. We assume that $\mathcal{X} = \mathbb{R}^d$, $\Psi(0) = 0$, Ψ is measurable and bounded, and ρ_1, ρ_0 have continuous densities with respect to the Lebesgue measure which we shall also denote as ρ_1, ρ_0 . Furthermore, we assume that there is $\varepsilon > 0$ and a measurable function $K : \mathcal{X} \times \mathbb{R}^d \rightarrow [0, \infty)$ such that for every $x \in \mathbb{R}^d$ we have the representation

$$\text{d}\mathfrak{p}_x(x') = \varepsilon^{-d} K\left(x, \frac{x' - x}{\varepsilon}\right) \text{d}x'.$$

We also assume that for every $x \in \mathcal{X}$ we have $K(x, \bullet) \in L^1(\mathbb{R}^d)$, $\int_{\mathbb{R}^d} K(x, z) \text{d}z = 1$, and $K(x, z) = 0$ if $|z| > 1$, and that for every $z \in \mathbb{R}^d$ the mapping $x \mapsto K(x, z)$ is C^1 .

Proposition 1. *Under Assumption 2, if A has a compact $C^{1,1}$ boundary and either Ψ is continuous or there exists a constant $c > 0$ such that $K(x, z) \geq c$ for all $x \in \mathcal{X}$ and $|z| \leq 1$, then*

$$\lim_{\varepsilon \rightarrow 0} \frac{1}{\varepsilon} \text{ProbPer}_\Psi(A) = \int_{\partial A} \sigma_{0, \Psi}[x, n(x)] \rho_0(x) + \sigma_{1, \Psi}[x, n(x)] \rho_1(x) \text{d}\mathcal{H}^{d-1}(x) \quad (17)$$

where we let $n(x)$ denote the normal to ∂A at a point $x \in \partial A$, and for any vector $v \in \mathbb{R}^d$ we define

$$\sigma_\Psi^0[x, v] := \int_0^1 \Psi\left(\int_{\{z \cdot v \leq -t\}} K(x, z) \text{d}z\right) \text{d}t, \quad \sigma_\Psi^1[x, v] := \int_0^1 \Psi\left(\int_{\{z \cdot v \geq t\}} K(x, z) \text{d}z\right) \text{d}t.$$

Remark 1. If K is radially symmetric and independent of $x \in \mathcal{X}$, then $\sigma_\Psi^0 = \sigma_\Psi^1 =: \sigma_\Psi$ is just a constant. E.g., for $K(x, z) := |B_1(0)|^{-1} \mathbf{1}_{|z| \leq 1}$ and $\Psi(t) = \mathbf{1}_{t > p}$ it is trivial that for $p = 0$ we have $\sigma_\Psi = 1$. However, for $p \geq \frac{1}{2}$ one easily sees $\sigma_\Psi = 0$, hence the limiting perimeter equals zero and there is no regularization effect. Using the function $\Psi(t) = \min\{t/p, 1\}$ corrects this degeneracy.

Notably, for radially symmetric K the limiting perimeter in (17) coincides, provided $\sigma_\Psi > 0$, with the one derived for adversarial training (problem (5)) by Bungert and Stinson [2022], although they considered more general (potentially discontinuous) densities ρ_i . In particular, our result indicates that for very small adversarial budgets the regularization effect of both probabilistically robust learning and adversarial training is dominated by the perimeter in (17). While Proposition 1 already completes half of the proof (namely the limsup inequality) of Gamma-convergence of $\frac{1}{\varepsilon} \text{ProbPer}_\Psi$ to the limiting perimeter, the remaining liminf inequality is beyond the scope of this paper. Proving that the convergence (17) does not only hold for sufficiently smooth sets as assumed in Proposition 1 but even in the sense of Gamma-convergence is an extremely important topic for future work since only Gamma-convergence allows to deduce from the convergence of the perimeters that also the solutions of probabilistically robust learning converge to certain regular Bayes classifiers as $\varepsilon \rightarrow 0$, see Bungert and Stinson [2022, Section 4.2].

3 General models

We now shift our attention to training general hypotheses $h \in \mathcal{H}$ using general loss functions $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$. Motivated by Theorems 1 and 3 we propose the following probabilistically robust optimization problem:

$$\inf_{h \in \mathcal{H}} \mathbb{E}_{(x, y) \sim \mu} \left[\max \left\{ p\text{-ess sup}_{x' \sim \mathfrak{p}_x} \ell(h(x'), y), \ell(h(x), y) \right\} \right]. \quad (18)$$

In the mathematical finance or economics literature the p -ess sup operator is better known as the value at risk (VaR) of a random variable at level p and it is notoriously hard to optimize. VaR is closely related to other risk measures like, for instance, the conditional value at risk (CVaR) which is convex and easier to optimize [Robey et al., 2022, Rockafellar et al., 2000]. For a function $f : \mathcal{X} \rightarrow \mathbb{R}$ and a probability distribution \mathbf{p} the CVaR at level p is defined as

$$\text{CVaR}_p(f; \mathbf{p}) := \inf_{\alpha \in \mathbb{R}} \alpha + \frac{\mathbb{E}_{x' \sim \mathbf{p}_x} [(f(x') - \alpha)_+]}{p}. \quad (19)$$

It is easy to see that p -ess sup $_{x' \sim \mathbf{p}} f(x') \leq \text{CVaR}_p(f; \mathbf{p})$. Using CVaR in place of the p -ess sup operator, a tractable version of (18) is

$$\inf_{h \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mu} \left[\max \left\{ \text{CVaR}_p(\ell(h(\bullet), y); \mathbf{p}_x), \ell(h(x), y) \right\} \right]. \quad (20)$$

We emphasize that, if the loss function $\ell(\bullet, \bullet)$ is convex in its first argument, then (20) is a convex function of the hypothesis h . Furthermore, CVaR is positively homogeneous and hence also (20) is positively homogeneous in the loss function. So, taking the maximum of the samplewise CVaR and standard risk is meaningful as both terms scale in the same way.

In the binary classification case we can prove the following interesting result that the CVaR relaxation corresponds precisely to using the risk ProbR_Ψ with a special piecewise linear and concave function Ψ for which our theory from Section 2.2 applies. In Appendix A.5 we prove a more general version of the following statement, replacing the $[\bullet]_+$ operation in (19) with a Leaky ReLU.

Theorem 6. *Let the function $\Psi_p : [0, 1] \rightarrow [0, 1]$ be defined as $\Psi_p(t) := \min \{t/p, 1\}$. Then it holds*

$$\text{CVaR}_p(\mathbf{1}_{A(\bullet) \neq y}; \mathbf{p}) = \Psi_p(\mathbb{P}_{x' \sim \mathbf{p}} [\mathbf{1}_A(x') \neq y])$$

and as a consequence for all $A \in \mathcal{A}$:

$$\mathbb{E}_{(x,y) \sim \mu} \left[\max \left\{ \text{CVaR}_p(\mathbf{1}_{A(\bullet) \neq y}; \mathbf{p}_x), \mathbf{1}_{A(x) \neq y} \right\} \right] = \text{ProbR}_{\Psi_p}(A).$$

An immediate consequence of Theorem 6 is that for binary classification (20) has a solution.

Corollary 1. *Under Assumption 1 and in the setting of Theorem 6 problem (20) has a solution.*

In Appendix A.5 we collect a few more observations concerning the CVaR, especially focussing on its behavior for $p > 1$. These geometric properties, the homogeneity with respect to the loss function, its potentially favorable sample complexity (see the discussion in Appendix A.6), and its versatility for algorithmic implementation make (20) a notable generalization of the adversarial training problem (2). Notice that when $p \rightarrow 0$ one formally recovers (2) from (20).

4 Numerical results

We build upon the code of Robey et al. [2022] and our implementation is available on GitHub.¹ The algorithmic realization of (20) is a straightforward adaptation of their algorithm, which alternately minimizes the inner optimization problem that defines CVaR and the outer optimization to find a suitable classifier, see Algorithm 1 in Appendix B. In our experiments, we conduct a comparative analysis between their algorithm (denoted as ‘‘Original’’ in Table 1) and Algorithm 1 in the appendix which is based on (20) (denoted as ‘‘Geometric’’). We report the clean, and adversarial accuracies (subject to PGD attacks), as well as accuracies on noise-augmented data and quantile accuracies for different values of p (see [Robey et al., 2022, (6.1)] for the definition) averaged over three runs; see Appendix B.2 for more training details. Our experiments are conducted on MNIST and CIFAR-10 and to ensure a fair comparison we adhere to the hyperparameter settings described by Robey et al. [2022], such that both the original and geometric algorithms utilize the same set of hyperparameters for each specified value of p . The corresponding results for several baseline algorithms including empirical risk minimization and adversarial training can be found in their paper. We perform model selection based on the best clean validation accuracy. The results in Table 1 show that for moderate values of p our geometric modification induces higher adversarial robustness than the original PRL without

¹https://github.com/DanielMckenzie/Begins_with_a_boundary

loss of clean accuracy (see, in particular, the results for MNIST with $p = 0.1$ and for CIFAR-10 with $p = 0.3$). In the noise augmented metrics as well as for extreme values of p close to 0 or equal to 0.5 both algorithms behave comparably. The latter can be expected from our theoretical results, in particular Proposition 1.

Note that the original or the geometric version of PRL should not be expected to match the adversarial robustness of classifiers trained with PGD attacks [Madry et al., 2017] or other worst-case optimization techniques. Instead, they shine with superior clean accuracies and easier training while maintaining probabilistic and a certain degree of adversarial robustness, as also observed by Robey et al. [2022].

We also remark that our sweep over different values of p confirms that increasing this parameter interpolates between low and high clean accuracies. However, it should be noted that it does not necessarily result in a direct interpolation between high and low adversarial or probabilistic accuracy, as claimed by Robey et al. [2022]. These observations hold true for both the original algorithm and our geometric modification, and despite utilizing their code and hyperparameters, we were unable to reproduce the exact results reported by Robey et al. [2022, Tables 1-4].

Table 1: Accuracies [%] of the geometric and original algorithm for different values of p .

| Data | p | Algorithm | Clean | Adv | Aug | Aug-0.1 | Aug-0.05 | Aug-0.01 |
|----------|------|-----------|--------------|--------------|-------|---------|----------|----------|
| MNIST | 0.01 | Geometric | 99.20 | 12.19 | 99.04 | 98.18 | 97.69 | 96.38 |
| | | Original | 99.19 | 10.76 | 98.90 | 97.94 | 97.38 | 95.67 |
| | 0.1 | Geometric | 99.28 | 14.20 | 99.22 | 98.70 | 98.45 | 97.86 |
| | | Original | 99.32 | 8.94 | 99.22 | 98.70 | 98.46 | 97.80 |
| | 0.3 | Geometric | 99.29 | 3.02 | 99.21 | 98.76 | 98.53 | 97.95 |
| | | Original | 99.27 | 3.02 | 99.22 | 98.77 | 98.55 | 98.01 |
| | 0.5 | Geometric | 99.27 | 1.80 | 99.21 | 98.72 | 98.44 | 97.93 |
| | | Original | 99.26 | 1.68 | 99.19 | 98.72 | 98.47 | 97.80 |
| CIFAR-10 | 0.01 | Geometric | 80.65 | 0.15 | 78.13 | 73.44 | 72.13 | 68.80 |
| | | Original | 81.73 | 0.24 | 79.16 | 74.61 | 73.19 | 69.96 |
| | 0.1 | Geometric | 88.15 | 0.14 | 85.96 | 82.55 | 81.46 | 78.81 |
| | | Original | 88.28 | 0.19 | 85.61 | 82.21 | 81.06 | 78.28 |
| | 0.3 | Geometric | 90.43 | 11.80 | 88.70 | 85.17 | 83.93 | 80.93 |
| | | Original | 89.97 | 7.20 | 88.62 | 85.07 | 83.75 | 80.87 |
| | 0.5 | Geometric | 91.51 | 1.93 | 88.94 | 85.53 | 84.18 | 81.21 |
| | | Original | 90.74 | 1.99 | 88.94 | 85.54 | 84.35 | 81.57 |

5 Discussion and Conclusion

In this paper we considered probabilistically robust learning (PRL), originally proposed by Robey et al. [2022]. We corrected a subtle but crucial theoretical flaw in the original formulation by introducing a regularization of the standard risk with nonlocal perimeters measuring the susceptibility of the decision boundary towards high-probability adversarial attacks. For binary classification we proved existence of optimal hard classifiers and of very general classes of soft classifiers including neural networks. We also provided an asymptotic expansion for smooth decision boundaries to show that for small adversarial budgets the probabilistic perimeters discussed in the paper induce the same regularization effect as adversarial training. For general (not necessarily binary) problems we showed that the natural loss function to choose is the sample-wise maximum of the standard loss and conditional value at risk (CVaR).

One limitation of PRL is that it does not completely solve the accuracy vs. robustness trade-off, which remains a challenging problem. Furthermore, while the formal limit of PRL as $p \rightarrow 0$ is the worst-case adversarial problem, the algorithms for solving PRL exhibit limitations for very small values of p (in the computation of CVaR_p). Still, the results for moderately large values of p are encouraging and future work should focus on understanding of this trade-off better.

The rich mathematical theory developed in this paper opens up new avenues for research, such as the explicit design of probabilistic regularizers for algorithms and exploring the variational convergence of the probabilistic perimeter and its implications for adversarial robustness.

Acknowledgments and Disclosure of Funding

This material is based upon work supported by the National Science Foundation under Grant Number DMS 1641020 and was started during the summer of 2022 as part of the AMS-MRC program *Data Science at the Crossroads of Analysis, Geometry, and Topology*. NGT was supported by the NSF grants DMS-2005797 and DMS-2236447.

References

- Maksym Andriushchenko and Nicolas Flammarion. Understanding and improving fast adversarial training. *Advances in Neural Information Processing Systems*, 33:16048–16059, 2020.
- Pranjal Awasthi, Natalie S Frank, and Mehryar Mohri. On the existence of the adversarial Bayes classifier. *Advances in Neural Information Processing Systems*, 34:2978–2990, 2021a.
- Pranjal Awasthi, Natalie S Frank, and Mehryar Mohri. On the existence of the adversarial Bayes classifier (extended version). *arXiv preprint arXiv:2112.01694*, 2021b.
- Leon Bungert and Kerrek Stinson. Gamma-convergence of a nonlocal perimeter arising in adversarial machine learning. *arXiv preprint arXiv:2211.15223*, 2022.
- Leon Bungert, Nicolás García Trillos, and Ryan Murray. The geometry of adversarial training in binary classification. *Information and Inference: A Journal of the IMA*, 12(2):921–968, 06 2023. ISSN 2049-8772. doi: 10.1093/imaia/iaac029.
- HanQin Cai, Yuchen Lou, Daniel McKenzie, and Wotao Yin. A zeroth-order block coordinate descent algorithm for huge-scale black-box optimization. In *International Conference on Machine Learning*, pages 1193–1203. PMLR, 2021.
- Antonin Chambolle, Massimiliano Morini, and Marcello Ponsiglione. Nonlocal curvature flows. *Archive for Rational Mechanics and Analysis*, 218:1263–1329, 2015.
- Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh. ZOO: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models. In *Proceedings of the 10th ACM workshop on artificial intelligence and security*, pages 15–26, 2017.
- Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. Certified adversarial robustness via randomized smoothing. In *international conference on machine learning*, pages 1310–1320. PMLR, 2019.
- Nelson Dunford and Jacob T Schwartz. *Linear Operators: General theory*. Linear Operators. Interscience Publishers, 1958. ISBN 9780470226056.
- Natalie S Frank. Existence and minimax theorems for adversarial surrogate risks in binary classification. *arXiv preprint arXiv:2206.09098*, 2022.
- Natalie S Frank and Jonathan Niles-Weed. The consistency of adversarial training for binary classification. *arXiv preprint arXiv:2206.09099*, 2022.
- Nicolás García Trillos and Ryan Murray. Adversarial classification: Necessary conditions and geometric flows. *Journal of Machine Learning Research*, 23(187):1–38, 2022.
- Nicolás García Trillos, Matt Jacobs, and Jakwang Kim. On the existence of solutions to adversarial training in multiclass classification. *arXiv preprint arXiv:2305.00075*, 2023.
- Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.

- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8340–8349, 2021.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- José M Mazón, Marcos Solera, and Julián Toledo. The total variation flow in metric random walk spaces. *Calculus of Variations and Partial Differential Equations*, 59:1–64, 2020.
- Muni Sreenivas Pydi and Varun Jog. The many faces of adversarial risk. *Advances in Neural Information Processing Systems*, 34:10000–10012, 2021.
- Yao Qin, Nicholas Carlini, Garrison Cottrell, Ian Goodfellow, and Colin Raffel. Imperceptible, robust, and targeted adversarial examples for automatic speech recognition. In *International conference on machine learning*, pages 5231–5240. PMLR, 2019.
- Joaquin Quinónero Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D Lawrence. *Dataset shift in machine learning*. Mit Press, 2008.
- Vinod Raman, Unique Subedi, and Ambuj Tewari. On proper learnability between average- and worst-case robustness. *arXiv preprint arXiv:2211.05656*, 2023.
- Alexander Robey, Luiz Chamon, George J Pappas, and Hamed Hassani. Probabilistically robust learning: Balancing average and worst-case performance. In *International Conference on Machine Learning*, pages 18667–18686. PMLR, 2022.
- R Tyrrell Rockafellar, Stanislav Uryasev, et al. Optimization of conditional value-at-risk. *Journal of risk*, 2:21–42, 2000.
- Leo Schwinn, An Nguyen, René Raab, Leon Bungert, Daniel Tenbrinck, Dario Zanca, Martin Burger, and Bjoern Eskofier. Identifying untrustworthy predictions in neural networks by geometric gradient analysis. In *Uncertainty in Artificial Intelligence*, pages 854–864. PMLR, 2021.
- Leo Schwinn, Leon Bungert, An Nguyen, René Raab, Falk Pulsmeier, Doina Precup, Björn Eskofier, and Dario Zanca. Improving robustness against real-world and worst-case distribution shifts through decision region quantification. In *International Conference on Machine Learning*, pages 19434–19449. PMLR, 2022.
- Ali Shafahi, Mahyar Najibi, Mohammad Amin Ghiasi, Zheng Xu, John Dickerson, Christoph Studer, Larry S Davis, Gavin Taylor, and Tom Goldstein. Adversarial training for free! *Advances in Neural Information Processing Systems*, 32, 2019.
- Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. Robustness may be at odds with accuracy. *arXiv preprint arXiv:1805.12152*, 2018.
- Vladimir Vapnik. *The nature of statistical learning theory*. Springer science & business media, 1999.
- Bao Wang, Binjie Yuan, Zuoqiang Shi, and Stanley J Osher. EnResNet: ResNets ensemble via the Feynman–Kac formalism for adversarial defense and beyond. *SIAM Journal on Mathematics of Data Science*, 2(3):559–582, 2020.
- Eric Wong, Leslie Rice, and J Zico Kolter. Fast is better than free: Revisiting adversarial training. *arXiv preprint arXiv:2001.03994*, 2020.
- Matthew D Zeiler. Adadelta: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*, 2012.
- Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan. Theoretically principled trade-off between robustness and accuracy. In *International conference on machine learning*, pages 7472–7482. PMLR, 2019.

A Proofs and theoretical aspects

A.1 Reformulation of the loss functions

In this section we will prove Theorems 1 and 3, which state that the proposed loss functionals equal the expected maximum of the sample-wise standard risk and the probabilistically robust risk from Robey et al. [2022]. Furthermore, we will prove Theorem 6, which states that the CVaR relaxation of the proposed regularization problem is equivalent to choosing a special piecewise linear function Ψ .

Proof of Theorem 1. We use the integral representations of the standard risk (4) and the proposed probabilistic perimeter (9) to express the risk functional $\text{ProbR}(A)$ as follows

$$\begin{aligned}
\text{ProbR}(A) &= \int_{\mathcal{X}} \mathbf{1}_{x \in A} d\rho_0(x) + \int_{\mathcal{X}} \mathbf{1}_{x \in A^c} d\rho_1(x) \\
&\quad + \int_{\mathcal{X}} \mathbf{1}_{x \in A \vee \mathbb{P}_{x' \sim p_x}[x' \in A] > p} - \mathbf{1}_{x \in A} d\rho_0(x) \\
&\quad + \int_{\mathcal{X}} \mathbf{1}_{x \in A^c \vee \mathbb{P}_{x' \sim p_x}[x' \in A^c] > p} - \mathbf{1}_{x \in A^c} d\rho_1(x). \\
&= \int_{\mathcal{X}} \mathbf{1}_{x \in A \vee \mathbb{P}_{x' \sim p_x}[x' \in A] > p} d\rho_0(x) + \int_{\mathcal{X}} \mathbf{1}_{x \in A^c \vee \mathbb{P}_{x' \sim p_x}[x' \in A^c] > p} d\rho_1(x) \\
&= \int_{\mathcal{X}} \max \left\{ \mathbf{1}_{\mathbb{P}_{x' \sim p_x}[x' \in A] > p}, \mathbf{1}_{x \in A} \right\} d\rho_0(x) + \int_{\mathcal{X}} \max \left\{ \mathbf{1}_{\mathbb{P}_{x' \sim p_x}[x' \in A^c] > p}, \mathbf{1}_{x \in A^c} \right\} d\rho_1(x),
\end{aligned}$$

where we used the fact that the indicator function of the union of two sets equals the maximum of the two indicator functions. Reverting the disintegration yields the claim:

$$\text{ProbR}(A) = \mathbb{E}_{(x,y) \sim \mu} \left[\max \left\{ \mathbf{1}_{\mathbb{P}_{x' \sim p_x}[\mathbf{1}_A(x') \neq y] > p}, \mathbf{1}_{\mathbf{1}_A(x) \neq y} \right\} \right].$$

□

Theorem 1 is a special case of the more general Theorem 3 which we prove in the following.

Proof of Theorem 3. The proof is similar to that of Theorem 1 after noting that for $\Psi : [0, 1] \rightarrow [0, 1]$

$$\mathbf{1}_{x \in A} + \mathbf{1}_{x \in A^c} \Psi(\mathbb{P}_{x' \sim p_x}[x' \in A]) = \max \left\{ \mathbf{1}_{x \in A}, \Psi(\mathbb{P}_{x' \sim p_x}[x' \in A]) \right\}$$

which can easily be shown by checking cases. Then:

$$\begin{aligned}
\text{ProbR}(A) &= \int_{\mathcal{X}} \mathbf{1}_{x \in A} d\rho_0(x) + \int_{\mathcal{X}} \mathbf{1}_{x \in A^c} d\rho_1(x) + \int_{\mathcal{X}} \mathbf{1}_{x \in A^c} \Psi(\mathbb{P}_{x' \sim p_x}[x' \in A]) d\rho_0(x) \\
&\quad + \int_{\mathcal{X}} \mathbf{1}_{x \in A} \Psi(\mathbb{P}_{x' \sim p_x}[x' \in A^c]) d\rho_1(x) \\
&= \int_{\mathcal{X}} [\mathbf{1}_{x \in A} + \mathbf{1}_{x \in A^c} \Psi(\mathbb{P}_{x' \sim p_x}[x' \in A])] d\rho_0(x) \\
&\quad + \int_{\mathcal{X}} [\mathbf{1}_{x \in A^c} + \mathbf{1}_{x \in A} \Psi(\mathbb{P}_{x' \sim p_x}[x' \in A^c])] d\rho_1(x) \\
&= \int_{\mathcal{X}} \max \left\{ \mathbf{1}_{x \in A}, \Psi(\mathbb{P}_{x' \sim p_x}[x' \in A]) \right\} d\rho_0(x) \\
&\quad + \int_{\mathcal{X}} \max \left\{ \mathbf{1}_{x \in A^c}, \Psi(\mathbb{P}_{x' \sim p_x}[x' \in A^c]) \right\} d\rho_1(x)
\end{aligned}$$

and the claim follows via reverting the disintegration as in the proof of Theorem 1. □

Before proving Theorem 6 we will prove the following stronger theorem.

Theorem 7. For $\beta \geq 0$ define

$$\phi_\beta(t) := \begin{cases} t & t \geq 0, \\ \beta t & t < 0, \end{cases}$$

and

$$\text{CVaR}_{p,\beta}(f; \mathbf{p}) := \inf_{\alpha \in \mathbb{R}} \alpha + \frac{\mathbb{E}_{x' \sim \mathbf{p}} [\phi_\beta(f(x') - \alpha)]}{p}.$$

The for any set $A \in \mathcal{A}$ and for $\beta \leq p < 1$ it holds

$$\text{CVaR}_{p,\beta}(\mathbf{1}_A; \mathbf{p}) = \Psi_{p,\beta}(\mathbf{p}(A))$$

where the concave and non-decreasing function $\Psi_{p,\beta} : [0, 1] \rightarrow [0, \infty)$ is defined as $\Psi_{p,\beta}(t) = \min\{t/p, 1 - \beta/p(1-t)\}$.

Proof. Suppose that $f = \mathbf{1}_A$ for some set $A \in \mathcal{A}$. In that case we can write

$$\text{CVaR}_{p,\beta}(f; \mathbf{p}) = \inf_{\alpha \in \mathbb{R}} \alpha + \phi_\beta(1 - \alpha) \frac{\mathbf{p}(A)}{p} + \phi_\beta(-\alpha) \frac{1 - \mathbf{p}(A)}{p}.$$

Notice that the function

$$\zeta(\alpha) := \alpha + \phi_\beta(1 - \alpha) \frac{\mathbf{p}(A)}{p} + \phi_\beta(-\alpha) \frac{1 - \mathbf{p}(A)}{p}, \quad \alpha \in \mathbb{R},$$

is continuous and piecewise linear with kinks at $\alpha = 0$ and $\alpha = 1$. Moreover, since $\beta \leq p < 1$ it holds $\zeta(\alpha) \geq \zeta(1)$ for $\alpha > 1$ and $\zeta(\alpha) \geq \zeta(0)$ for $\alpha < 0$ such that the minimum of ζ is attained at either $\alpha = 0$ or $\alpha = 1$. Thus

$$\text{CVaR}_{p,\beta}(f; \mathbf{p}) = \min\{\zeta(0), \zeta(1)\} = \min\left\{\frac{\mathbf{p}(A)}{p}, 1 - \frac{\beta}{p}(1 - \mathbf{p}(A))\right\}.$$

□

Proof of Theorem 6. The first claim of Theorem 6 is a special case of Theorem 7 by choosing $\beta = 0$. The second claim follows by combining the first one with Theorem 3. □

A.2 Lower semicontinuity of the functional J_Ψ

An essential tool for the proof of Theorems 2 and 4 is lower semicontinuity of the functional J_Ψ , which we recall was defined in (16) as

$$J_\Psi(u) := \int_{\mathcal{X}} (1 - u(x)) \Psi(\mathbb{E}_{x' \sim \mathbf{p}_x}[u(x')]) \, d\rho_0(x) + \int_{\mathcal{X}} u(x) \Psi(\mathbb{E}_{x' \sim \mathbf{p}_x}[1 - u(x')]) \, d\rho_1(x)$$

for a measurable function $u : \mathcal{X} \rightarrow [0, 1]$. We have to construct a suitable topology for proving lower semicontinuity of this functional (recall Assumption 1) and define the probability measures

$$\rho := \rho_0 + \rho_1, \tag{21}$$

$$\nu(A) := \frac{1}{2} \int_{\mathcal{X}} \mathbf{p}_x(A) \, d\rho(x) + \frac{1}{2} \rho(A), \quad A \in \mathcal{A}. \tag{22}$$

The measure ρ equals the first marginal of μ and models the distribution of all data, irrespective of the label. The first summand of the measure ν is the convolution of ρ with the family of probability measures $\{\mathbf{p}_x\}_{x \in \mathcal{X}}$.

By construction we have the following two important absolute continuity properties which we shall use without further reference:

$$\begin{aligned} \nu(A) = 0 &\implies \left[\rho(A) = 0 \quad \text{and} \quad \mathbf{p}_x(A) = 0 \text{ for } \rho\text{-almost every } x \in \mathcal{X} \right], \\ \rho(A) = 0 &\implies \left[\rho_0(A) = 0 \quad \text{and} \quad \rho_1(A) = 0 \right]. \end{aligned}$$

A simple example for $\mathcal{X} = \mathbb{R}^d$ is $\rho := \frac{1}{N} \sum_{i=1}^N \delta_{x_i}$ and $\mathbf{p}_x := \text{Unif}(B_\varepsilon(x))$ in which case

$$\nu = \frac{1}{2N} \sum_{i=1}^N \left(\text{Unif}(B_\varepsilon(x_i)) + \delta_{x_i} \right)$$

is a sum of absolutely continuous measures on ball centered at x_i and the empirical measure of the points x_i .

The suitable topology in which we shall prove lower semicontinuity is the weak-* topology of $L^\infty(\mathcal{X}; \nu)$ which is the dual space of $L^1(\mathcal{X}; \nu)$ since ν is *a fortiori* a σ -finite measure [Dunford and Schwartz, 1958, IV.8.3, Theorem 5].

Definition 1. Under Assumption 1 we say that a sequence of functions $(u_n)_{n \in \mathbb{N}} \subset L^\infty(\mathcal{X}; \nu)$ converges to $u \in L^\infty(\mathcal{X}; \nu)$ in the weak-* sense (written $u_n \xrightarrow{*} u$) as $n \rightarrow \infty$ if

$$\lim_{n \rightarrow \infty} \int_{\mathcal{X}} u_n \varphi \, d\nu = \int_{\mathcal{X}} u \varphi \, d\nu \quad \forall \varphi \in L^1(\mathcal{X}; \nu). \quad (23)$$

The absolute continuity properties of ν allow us to deduce the following lemma

Lemma 1. Under Assumption 1 let $(u_n)_{n \in \mathbb{N}} \subset L^\infty(\mathcal{X}; \nu)$ satisfy $u_n \xrightarrow{*} u$ in the sense of Definition 1. Then it holds

$$\lim_{n \rightarrow \infty} \mathbb{E}_{x' \sim \mathbf{p}_x} [u_n(x')] = \mathbb{E}_{x' \sim \mathbf{p}_x} [u(x')] \quad \text{for } \rho\text{-almost every } x \in \mathcal{X}.$$

Proof. The Radon–Nikodým theorem and weak-* convergence imply that for ρ -almost every $x \in \mathcal{X}$ it holds

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbb{E}_{x' \sim \mathbf{p}_x} [u_n(x')] &= \lim_{n \rightarrow \infty} \int_{\mathcal{X}} u_n(x') \, d\mathbf{p}_x(x') = \lim_{n \rightarrow \infty} \int_{\mathcal{X}} u_n(x') \frac{d\mathbf{p}_x}{d\nu}(x') \, d\nu(x') \\ &= \int_{\mathcal{X}} u(x') \frac{d\mathbf{p}_x}{d\nu}(x') \, d\nu(x') = \int_{\mathcal{X}} u(x') \, d\mathbf{p}_x(x') = \mathbb{E}_{x' \sim \mathbf{p}_x} [u(x')] \end{aligned}$$

since $\frac{d\mathbf{p}_x}{d\nu} \in L^1(\mathcal{X}; \nu)$. □

Proposition 2 (Lower semicontinuity of J_Ψ). Under Assumption 1 let $(u_n)_{n \in \mathbb{N}} \subset L^\infty(\mathcal{X}; \nu)$ be a sequence of functions with values in $[0, 1]$ satisfying $u_n \xrightarrow{*} u$ in the sense of Definition 1, and let $\Psi : [0, 1] \rightarrow [0, 1]$ be lower semicontinuous. Then $0 \leq u \leq 1$ holds ν -almost everywhere and furthermore

$$J_\Psi(u) \leq \liminf_{n \rightarrow \infty} J_\Psi(u_n).$$

Proof. First we show that $0 \leq u \leq 1$. By the weak-* lower semicontinuity of the L^∞ -norm we get $u \leq 1$ from the fact that $0 \leq u_n \leq 1$. To show that $u \geq 0$ we assume that on a measurable set N with $\nu(N) > 0$ it holds $u < 0$. Then from the weak-* convergence and the fact that $u_n \geq 0$ we obtain

$$0 > \int_{\mathcal{X}} u \mathbf{1}_N \, d\nu = \lim_{n \rightarrow \infty} \int_{\mathcal{X}} u_n \mathbf{1}_N \, d\nu \geq 0$$

which is a contradiction. Therefore, $u \geq 0$ holds ν -almost everywhere.

Since both terms in the definition of J_Ψ are dealt with symmetrically, we assume without loss of generality and for an easier notation that $\rho_1 = 0$ and rewrite J_Ψ as

$$J_\Psi(u) = \int_{\mathcal{X}} (1 - u(x)) \Psi(\mathbb{E}_{x' \sim \mathbf{p}_x} [u(x')]) \, d\rho_0(x).$$

Since Ψ is lower semicontinuous there exists a sequence of continuous functions $\Psi_\delta : [0, 1] \rightarrow [0, 1]$ which converge to Ψ in the pointwise sense as $\delta \rightarrow 0$ and satisfy $\Psi_\delta \leq \Psi$. For instance, the functions

$$\Psi_\delta(t) := \inf_{s \in [0, 1]} \Psi(s) + \frac{1}{\delta} |s - t|, \quad t \in [0, 1],$$

do the job. Lemma 1 implies that $\mathbb{E}_{x' \sim p_x} [u_n(x')] \rightarrow \mathbb{E}_{x' \sim p_x} [u(x')]$ for ρ -almost every x as $n \rightarrow \infty$. Since Ψ_δ is continuous, we get $\Psi_\delta(\mathbb{E}_{x' \sim p_x} [u_n(x')]) \rightarrow \Psi_\delta(\mathbb{E}_{x' \sim p_x} [u(x')])$ for ρ -almost every x as $n \rightarrow \infty$. Since $0 \leq u_n \leq 1$ and hence $\Psi_\delta(\mathbb{E}_{x' \sim p_x} [u_n(x')])$ is uniformly bounded, the convergence even holds true in $L^1(\mathcal{X}; \rho_0)$ and therefore

$$\begin{aligned} & \lim_{n \rightarrow \infty} \int_{\mathcal{X}} (1 - u_n(x)) \Psi_\delta(\mathbb{E}_{x' \sim p_x} [u_n(x')]) \, d\rho_0(x) \\ &= \int_{\mathcal{X}} (1 - u(x)) \Psi_\delta(\mathbb{E}_{x' \sim p_x} [u(x')]) \, d\rho_0(x). \end{aligned} \tag{24}$$

Next we would like to use the Fatou lemma to take the limit as $\delta \rightarrow 0$ on both sides. For this we notice that the sequence of functions

$$f_\delta(x) := (1 - u_n(x)) \Psi_\delta(\mathbb{E}_{x' \sim p_x} [u_n(x')])$$

converges to $(1 - u_n(x)) \Psi(\mathbb{E}_{x' \sim p_x} [u_n(x')])$ pointwise as $\delta \rightarrow 0$ and satisfies the bounds $f_\delta \geq 0$. Thanks to the non-negativity we can apply the standard Fatou lemma. Using $\Psi_\delta \leq \Psi$ and (24) we get

$$\begin{aligned} J_\Psi(u) &= \int_{\mathcal{X}} (1 - u(x)) \Psi(\mathbb{E}_{x' \sim p_x} [u(x')]) \, d\rho_0(x) \\ &\leq \liminf_{\delta \rightarrow 0} \int_{\mathcal{X}} (1 - u(x)) \Psi_\delta(\mathbb{E}_{x' \sim p_x} [u(x')]) \, d\rho_0(x) \\ &= \liminf_{\delta \rightarrow 0} \lim_{n \rightarrow \infty} \int_{\mathcal{X}} (1 - u_n(x)) \Psi_\delta(\mathbb{E}_{x' \sim p_x} [u_n(x')]) \, d\rho_0(x) \\ &\leq \liminf_{n \rightarrow \infty} \int_{\mathcal{X}} (1 - u_n(x)) \Psi(\mathbb{E}_{x' \sim p_x} [u_n(x')]) \, d\rho_0(x) \\ &= \liminf_{n \rightarrow \infty} J_\Psi(u_n). \end{aligned}$$

□

A.3 The geometric problem for concave Ψ

We remind the reader of the definition of the following perimeter functional:

$$\text{ProbPer}_\Psi(A) := \int_{\mathcal{X}} \mathbf{1}_{x \in A^c} \Psi(\mathbb{P}_{x' \sim p_x} [x' \in A]) \, d\rho_0(x) + \int_{\mathcal{X}} \mathbf{1}_{x \in A} \Psi(\mathbb{P}_{x' \sim p_x} [x' \in A^c]) \, d\rho_1(x).$$

We first show that ProbPer_Ψ is a submodular function when Ψ is a concave non-decreasing function. For this, we first need a lemma.

Lemma 2. *Let $\Psi : [0, \infty) \rightarrow \mathbb{R}$ be a concave and non-decreasing function, and let $0 \leq a \leq b \leq b' \leq a'$ be real numbers with $a + a' \leq b + b'$. Then*

$$\Psi(a) + \Psi(a') \leq \Psi(b) + \Psi(b').$$

Proof. Let a, a', b, b' be as stated. Since Ψ is concave and finite, it satisfies the fundamental theorem of calculus and thus it is possible to write

$$\Psi(s) = \Psi(a) + \int_a^s \Psi'(r) \, dr, \quad s \geq a$$

for a function Ψ' that is non-increasing and non-negative. It follows that

$$\Psi(b) - \Psi(a) = \int_a^b \Psi'(r) \, dr \geq (b - a)\Psi'(b) \geq (a' - b')\Psi'(b) \geq \int_{b'}^{a'} \Psi'(r) \, dr = \Psi(a') - \Psi(b'),$$

which is precisely what we wanted to show. □

We are ready to prove Theorem 5

Proof of Theorem 5. First, the fact that $\text{ProbPer}_\Psi(\mathcal{X}) = \text{ProbPer}_\Psi(\emptyset) = 0$ if $\Psi(0) = 0$ is easy to see from the definition of ProbPer_Ψ . Second, we trivially have

$$\mathbb{P}_{x' \sim p_x} [x' \in A \cup B] + \mathbb{P}_{x' \sim p_x} [x' \in A \cap B] \leq \mathbb{P}_{x' \sim p_x} [x' \in A] + \mathbb{P}_{x' \sim p_x} [x' \in B].$$

Define

$$\begin{aligned} a' &:= \mathbb{P}_{x' \sim p_x} [x' \in A \cup B], & b' &:= \mathbb{P}_{x' \sim p_x} [x' \in B] \\ b &:= \mathbb{P}_{x' \sim p_x} [x' \in A], & a &:= \mathbb{P}_{x' \sim p_x} [x' \in A \cap B]; \end{aligned}$$

without the loss of generality we can assume that b and b' defined above satisfy $b \leq b'$, for otherwise we can simply swap these labels. We can then use Lemma 2 to conclude that:

$$\begin{aligned} &\Psi(\mathbb{P}_{x' \sim p_x} [x' \in A \cup B]) + \Psi(\mathbb{P}_{x' \sim p_x} [x' \in A \cap B]) \\ &\leq \Psi(\mathbb{P}_{x' \sim p_x} [x' \in A]) + \Psi(\mathbb{P}_{x' \sim p_x} [x' \in B]). \end{aligned} \quad (25)$$

The submodularity follows directly once we have verified the following pointwise identity:

$$\begin{aligned} &\mathbf{1}_{x \in (A \cup B)^c} \Psi(\mathbb{P}_{x' \sim p_x} [x' \in A \cup B]) + \mathbf{1}_{x \in (A \cap B)^c} \Psi(\mathbb{P}_{x' \sim p_x} [x' \in A \cap B]) \\ &\leq \mathbf{1}_{x \in A^c} \Psi(\mathbb{P}_{x' \sim p_x} [x' \in A]) + \mathbf{1}_{x \in B^c} \Psi(\mathbb{P}_{x' \sim p_x} [x' \in B]). \end{aligned} \quad (26)$$

To do this we consider two complementary cases:

Case 1, $x \in (A \cup B)^c$: This is equivalent to $x \in A^c \cap B^c$. Furthermore, since $(A \cup B)^c \subset (A \cap B)^c$ we also have that $x \in (A \cap B)^c$. Hence, all indicator functions in (26) take the value one and (26) is the same as (25), which we have already verified.

Case 2, $x \in A \cup B$: In this case the first indicator function on the left hand side of (26) is zero.

Case 2.1, $x \in A \cap B$: In this subcase all indicator functions are equal to zero and the inequality is trivially satisfied.

Case 2.2, $x \in A \cup B \setminus (A \cap B)$: Without loss of generality we can assume that $x \in A \setminus B = A \cap B^c$. In this case only the second indicator function on the left hand side and the second one on the right hand side of (26) take the value one and the inequality reduces to the trivial inequality

$$\Psi(\mathbb{P}_{x' \sim p_x} [x' \in A \cap B]) \leq \Psi(\mathbb{P}_{x' \sim p_x} [x' \in B])$$

which is true since Ψ is non-decreasing and $A \cap B \subset B$. □

Motivated by Theorem 5 we define the associated total variation of a non-negative measurable function $u : \mathcal{X} \rightarrow [0, \infty)$ in terms of a coarea formula as

$$\text{ProbTV}_\Psi(u) := \int_0^\infty \text{ProbPer}_\Psi(\{u \geq t\}) dt. \quad (27)$$

By definition ProbTV_Ψ is positively homogeneous. It also satisfies $\text{ProbTV}_\Psi(\mathbf{1}_A) = \text{ProbPer}_\Psi(A) = J_\Psi(\mathbf{1}_A)$ but the functionals ProbTV_Ψ and J_Ψ do not coincide for general functions $u : \mathcal{X} \rightarrow [0, 1]$. Instead it holds $\text{ProbTV}_\Psi(u) \leq J_\Psi(u)$, as we prove in the following proposition.

Proposition 3. *If Ψ is concave and non-decreasing it holds*

$$\text{ProbTV}_\Psi(u) \leq J_\Psi(u).$$

Furthermore, for every sequence of measurable sets $(A_n)_{n \in \mathbb{N}} \subset \mathcal{X}$ such that $\mathbf{1}_{A_n}$ converges in the weak- sense to a function $u \in L^\infty(\mathcal{X}; \nu)$ as $n \rightarrow \infty$ it holds*

$$\text{ProbTV}_\Psi(u) \leq \liminf_{n \rightarrow \infty} \text{ProbPer}_\Psi(A_n).$$

Proof. As in the proof of Proposition 2 we assume without loss of generality that $\rho_1 = 0$. We compute

$$\begin{aligned}
\text{ProbTV}_\Psi(u) &= \int_0^1 \text{ProbPer}_\Psi(\{u \geq t\}) dt \\
&= \int_{\mathcal{X}} \int_0^1 \mathbf{1}_{u(x) < t} \Psi(\mathbb{P}_{x' \sim p_x}[u(x') \geq t]) dt d\rho_0(x) \\
&= \int_{\mathcal{X}} \int_0^1 \mathbf{1}_{u(x) < t} \Psi(\mathbb{E}_{x' \sim p_x}[\mathbf{1}_{\{u \geq t\}}(x')]) dt d\rho_0(x) \\
&= \int_{\mathcal{X}} \int_{u(x)}^1 \Psi(\mathbb{E}_{x' \sim p_x}[\mathbf{1}_{\{u \geq t\}}(x')]) dt d\rho_0(x). \tag{28}
\end{aligned}$$

Since Ψ is concave and non-decreasing, we get from (28) and Jensen's inequality that

$$\begin{aligned}
\text{ProbTV}_\Psi(u) &\leq \int_{\mathcal{X}} (1 - u(x)) \Psi\left(\frac{1}{1 - u(x)} \int_{u(x)}^1 \mathbb{E}_{x' \sim p_x}[\mathbf{1}_{\{u \geq t\}}(x')] dt\right) d\rho_0(x) \\
&= \int_{\mathcal{X}} (1 - u(x)) \Psi\left(\frac{1}{1 - u(x)} \mathbb{E}_{x' \sim p_x}\left[\int_{u(x)}^1 \mathbf{1}_{\{u \geq t\}}(x') dt\right]\right) d\rho_0(x) \\
&= \int_{\mathcal{X}} (1 - u(x)) \Psi\left(\frac{1}{1 - u(x)} \mathbb{E}_{x' \sim p_x}[u(x') - u(x)]\right) d\rho_0(x) \\
&\leq \int_{\mathcal{X}} (1 - u(x)) \Psi\left(\frac{1}{1 - u(x)} \mathbb{E}_{x' \sim p_x}[u(x')(1 - u(x))]\right) d\rho_0(x) \\
&= \int_{\mathcal{X}} (1 - u(x)) \Psi(\mathbb{E}_{x' \sim p_x}[u(x')]) d\rho_0(x) = J_\Psi(u).
\end{aligned}$$

The proof of the second statement of the proposition follows by combining the first one with Proposition 2, applied to the sequence $u_n := \mathbf{1}_{A_n}$, which satisfies $J_\Psi(u_n) = \text{ProbPer}_\Psi(A_n)$. \square

A remarkable consequence of this lower bound and the lower semicontinuity of J_Ψ from Proposition 2 is the following lower semicontinuity of ProbTV_Ψ for sequences of characteristic functions. For this sake, let $(A_n) \subset \mathcal{A}$ be a sequence of sets such that $\mathbf{1}_{A_n} \xrightarrow{*} u$ in $L^\infty(\mathcal{X}; \nu)$. Then it holds

$$\begin{aligned}
\text{ProbTV}_\Psi(u) &\leq J_\Psi(u) \leq \liminf_{n \rightarrow \infty} J_\Psi(\mathbf{1}_{A_n}) = \liminf_{n \rightarrow \infty} \text{ProbPer}_\Psi(A_n) \\
&= \liminf_{n \rightarrow \infty} \text{ProbTV}_\Psi(\mathbf{1}_{A_n}).
\end{aligned}$$

Remarkably, this observation suffices to prove Theorem 2 although there is no proof for lower semicontinuity of ProbTV_Ψ along *general sequences* of functions.

Proof of Theorem 2. Let $(A_n)_{n \in \mathbb{N}} \subset \mathcal{A}$ be a minimizing sequence such that

$$\lim_{n \rightarrow \infty} \text{R}_{\text{std}}(A_n) + \text{ProbPer}_\Psi(A_n) = \inf_{A \in \mathcal{A}} \text{R}_{\text{std}}(A) + \text{ProbPer}_\Psi(A). \tag{29}$$

The Banach–Alaoglu theorem implies that there exists $u \in L^\infty(\mathcal{X}; \nu)$ such that $\mathbf{1}_{A_n} \xrightarrow{*} u$ as $n \rightarrow \infty$. The standard risk is trivially lower semicontinuous because $\rho \ll \nu$. Combining this with Propositions 2 and 3 and using (29) we obtain

$$\begin{aligned}
\mathbb{E}_{(x,y) \sim \mu}[|u(x) - y|] + \text{ProbTV}_\Psi(u) &\leq \mathbb{E}_{(x,y) \sim \mu}[|u(x) - y|] + J_\Psi(u) \\
&\leq \liminf_{n \rightarrow \infty} \text{R}_{\text{std}}(A_n) + J_\Psi(\mathbf{1}_{A_n}) \\
&= \liminf_{n \rightarrow \infty} \text{R}_{\text{std}}(A_n) + \text{ProbPer}_\Psi(A_n) \\
&= \inf_{A \in \mathcal{A}} \text{R}_{\text{std}}(A) + \text{ProbPer}_\Psi(A).
\end{aligned}$$

For $t \in [0, 1]$ we define $A_t := \{u \geq t\}$. Trivially we have

$$\inf_{A \in \mathcal{A}} \text{R}_{\text{std}}(A) + \text{ProbPer}_\Psi(A) \leq \text{R}_{\text{std}}(A_t) + \text{ProbPer}_\Psi(A_t) \quad \forall t \in [0, 1].$$

If this inequality were strict for a set of $t \in [0, 1]$ with positive Lebesgue measure, integration and the coarea formula would give

$$\begin{aligned} \inf_{A \in \mathcal{A}} R_{\text{std}}(A) + \text{ProbPer}_{\Psi}(A) &< \int_0^1 R_{\text{std}}(A_t) + \text{ProbPer}_{\Psi}(A_t) dt \\ &= \mathbb{E}_{(x,y) \sim \mu} [|u(x) - y|] + \text{ProbTV}_{\Psi}(u) \\ &\leq \inf_{A \in \mathcal{A}} R_{\text{std}}(A) + \text{ProbPer}_{\Psi}(A), \end{aligned}$$

which is a contradiction. Hence, we have proved that

$$\inf_{A \in \mathcal{A}} R_{\text{std}}(A) + \text{ProbPer}_{\Psi}(A) = R_{\text{std}}(A_t) + \text{ProbPer}_{\Psi}(A_t)$$

for Lebesgue almost every $t \in [0, 1]$ and any such set A_t is a minimizer. \square

A.4 The soft classifier problems

The main issue with proving existence of solution for the original problem (i.e., $\Psi(t) = \mathbf{1}_{t > p}$ is not concave) is that

$$\text{ProbTV}_{\Psi}(u) \not\leq J_{\Psi}(u)$$

due to a failure of Jensen's inequality in this case. The validity of this inequality was central for the relaxation arguments in the proofs of Proposition 3 and Theorem 2.

However, if we already consider the relaxed problem, optimizing over soft classifiers instead of characteristic functions and regularizing with J_{Ψ} , we obtain an existence proof straightforwardly for very general functions Ψ .

Proof of Theorem 4. The proof works precisely as the proof of Theorem 2, replacing the minimizing sequence $\mathbf{1}_{A_n}$ by functions $u_n \in L^{\infty}(\mathcal{X}; \nu)$ which satisfy $0 \leq u_n \leq 1$, and utilizing the weak-* closedness of \mathcal{H} as well as the lower semicontinuity of J_{Ψ} from Proposition 2. \square

A.5 CVaR relaxation: Existence for hard classifiers and some considerations

Proof of Corollary 1. Given the reformulation from Theorem 6 the existence result follows from Theorem 2. \square

We remark that if p was allowed to be larger than 1, then $\text{CVaR}_p(f; \mathfrak{p}) = -\infty$. Indeed, if $p > 1$, the function ζ defined in the proof of Theorem 7 would satisfy $\lim_{\alpha \rightarrow -\infty} \zeta(\alpha) = -\infty$. This insight allows us to show that the method of Robey et al. [2022] becomes meaningless for values $p > 1$, despite the authors using values bigger than one according to [Robey et al., 2022, Appendix C.4]. In contrast, our method just reduces to empirical risk minimization in this case.

Proposition 4. *For a non-negative function $f : \mathcal{X} \rightarrow [0, \infty]$ and for $p > 1$ it holds*

$$p\text{-ess sup}_{x' \sim \mathfrak{p}} f(x') = \inf \{t > 0 : \mathbb{P}_{x' \sim \mathfrak{p}} [f(x') > t] \leq p\} = -\infty$$

$$\text{CVaR}_p(f; \mathfrak{p}) = \inf_{\alpha \in \mathbb{R}} \alpha + \frac{\mathbb{E}_{x' \sim \mathfrak{p}} [(f(x') - \alpha)_+]}{p} = -\infty.$$

Proof. For $p > 1$ (in fact, even for $p = 1$) it holds for any $t \in \mathbb{R}$ that $\mathbb{P}_{x' \sim \mathfrak{p}} [f(x') > t] \leq 1 \leq p$. Hence, the infimum over such t equals $-\infty$.

Similarly, for $p > 1$ and $\alpha \leq 0$, we have, thanks to the non-negativity of f , that

$$\alpha + \frac{\mathbb{E}_{x' \sim \mathfrak{p}_x} [(f(x') - \alpha)_+]}{p} = \left(1 - \frac{1}{p}\right) \alpha + \frac{1}{p} \mathbb{E}_{x' \sim \mathfrak{p}} [f(x')].$$

If $p > 1$, the factor which multiplies α is positive and sending $\alpha \rightarrow -\infty$ shows the infimum equals $-\infty$. \square

Corollary 2. *Consequently, for $p > 1$ and non-negative loss functions $\ell(\bullet, \bullet)$ the objective in (3) (respectively, its CVaR relaxation [Robey et al., 2022, (P-CVaR)]) equals $-\infty$. In contrast, the proposed objective (18) and its CVaR relaxation (20) equal the standard risk $\mathbb{E}_{(x,y) \sim \mu} [\ell(h(x), y)]$ in this case.*

A.6 PAC learnability for Lipschitz continuous Ψ

Raman et al. [2023] considered PAC learnability of probabilistically robust learning models which depend on the *probabilistic margin* $|\mathbb{E}_{x' \sim p} [u(x')] - y|$ of u on $(x, y) \sim \mu$. More precisely (and in our notation) they consider the problem

$$\inf_{h \in \mathcal{H}} \mathbb{E}_{(x, y) \sim \mu} [\Psi (|\mathbb{E}_{x' \sim p_x} [u(x')] - y|)]. \quad (30)$$

Notably, for $\Psi(t) = \mathbf{1}_{t > p}$ this reduces precisely to the probabilistically robust learning model (3) by Robey et al. [2022] in case of using the loss function $\ell(y, \tilde{y}) = |y - \tilde{y}|$. However, they prove hardness results for such functions Ψ and also prove superior learnability properties if function Ψ is Lipschitz continuous. We expect that analogous statements carry over to the problems considered in Theorems 2 and 4 if one restricts Ψ to be Lipschitz. An interesting question for future investigation is whether concavity of Ψ (which is needed for the first theorem) would suffice to guarantee PAC learnability.

A.7 Pointwise consistency of the perimeter

Proof of Proposition 1. Under Assumption 2, a simple change of variables shows

$$\begin{aligned} \text{ProbPer}_\Psi(A) &= \int_{A^c} \Psi \left(\int_{\mathbb{R}^d} \mathbf{1}_A(x + \varepsilon z) K(x, z) dz \right) \rho_0(x) dx \\ &\quad + \int_A \Psi \left(\int_{\mathbb{R}^d} \mathbf{1}_{A^c}(x + \varepsilon z) K(x, z) dz \right) \rho_1(x) dx. \end{aligned}$$

Let $\tau : \mathbb{R}^d \rightarrow \mathbb{R}$ be the signed distance function to ∂A such that $\tau(x) \leq 0$ for $x \in A$. Using τ , we can rewrite the previous line as

$$\begin{aligned} \text{ProbPer}_\Psi(A) &= \int_{\{\tau(x) \geq 0\}} \Psi \left(\int_{\{\tau(x + \varepsilon z) \leq 0\}} K(x, z) dz \right) \rho_0(x) dx \\ &\quad + \int_{\{\tau(x) \leq 0\}} \Psi \left(\int_{\{\tau(x + \varepsilon z) \geq 0\}} K(x, z) dz \right) \rho_1(x) dx. \end{aligned}$$

Recalling that $K(x, z) = 0$ whenever $|z| > 1$, the previous line is equal to

$$\begin{aligned} \text{ProbPer}_\Psi(A) &= \int_{\{0 \leq \tau(x) \leq \varepsilon\}} \Psi \left(\int_{\{\tau(x + \varepsilon z) \leq 0\}} K(x, z) dz \right) \rho_0(x) dx \\ &\quad + \int_{\{-\varepsilon \leq \tau(x) \leq 0\}} \Psi \left(\int_{\{\tau(x + \varepsilon z) \geq 0\}} K(x, z) dz \right) \rho_1(x) dx. \end{aligned}$$

Since A has $C^{1,1}$ boundary, there exists some $\varepsilon_0 > 0$ such that τ is $C^{1,1}$ in an ε_0 neighborhood of ∂A . Furthermore, for any $\varepsilon < \varepsilon_0$ the mapping $T_\varepsilon : \partial A \times [-1, 1] \rightarrow \{x \in \mathbb{R}^d : \tau(x) \in [-\varepsilon, \varepsilon]\}$ given by $T_\varepsilon(y, t) = y + \varepsilon t n(y)$ is a bijection and $\lim_{\varepsilon \rightarrow 0} \frac{1}{\varepsilon} \det(DT_\varepsilon) = 1$ uniformly on $\partial A \times [-1, 1]$. Using this change of variables, we may write

$$\begin{aligned} &\int_{\{0 \leq \tau(x) \leq \varepsilon\}} \Psi \left(\int_{\{\tau(x + \varepsilon z) \leq 0\}} K(x, z) dz \right) \rho_0(x) dx \\ &= \int_{\partial A} \int_0^1 \det(DT_\varepsilon(y, t)) \Psi(a_\varepsilon(y, t)) \rho_0(y + \varepsilon t n(y)) dt d\mathcal{H}^{d-1}(y) \end{aligned}$$

and

$$\begin{aligned} &\int_{\{-\varepsilon \leq \tau(x) \leq 0\}} \Psi \left(\int_{\{\tau(x + \varepsilon z) \geq 0\}} K(x, z) dz \right) \rho_1(x) dx \\ &= \int_{\partial A} \int_0^1 \det(DT_\varepsilon(y, -t)) \Psi(b_\varepsilon(y, t)) \rho_1(y - \varepsilon t n(y)) dt d\mathcal{H}^{d-1}(y), \end{aligned}$$

where we abbreviate

$$a_\varepsilon(y, t) := \int_{\{\tau(y+\varepsilon(tn(y)+z)) \leq 0\}} K(y + \varepsilon tn(y), z) dz$$

$$b_\varepsilon(y, t) := \int_{\{\tau(y+\varepsilon(z-tn(y))) \geq 0\}} K(y - \varepsilon tn(y), z) dz$$

In the rest of what follows, we will focus on proving that

$$\lim_{\varepsilon \rightarrow 0} \int_{\partial A} \int_0^1 \frac{\det(DT_\varepsilon(y, t))}{\varepsilon} \Psi(a_\varepsilon(y, t)) \rho_0(y + \varepsilon tn(y)) dt d\mathcal{H}^{d-1}(y)$$

$$= \int_{\partial A} \sigma_{0, \Psi}[y, n(y)] \rho_0(y) d\mathcal{H}^{d-1}(y).$$

An identical argument will show that

$$\lim_{\varepsilon \rightarrow 0} \int_0^1 \frac{\det(DT_\varepsilon(y, -t))}{\varepsilon} \Psi(b_\varepsilon(y, t)) \rho_1(y - \varepsilon tn(y)) dt d\mathcal{H}^{d-1}(y)$$

$$= \int_{\partial A} \sigma_{1, \Psi}[y, n(y)] \rho_1(y) d\mathcal{H}^{d-1}(y).$$

Since we know that $\lim_{\varepsilon \rightarrow 0} \frac{\det(DT_\varepsilon(y, t))}{\varepsilon} = 1$ and $\lim_{\varepsilon \rightarrow 0} \rho_0(y + \varepsilon tn(y)) = \rho_0(y)$ pointwise almost everywhere, the main difficulty lies in passing to the limit in the term involving Ψ . For this we shall first prove convergence of $a_\varepsilon(y, t)$ to

$$a(y, t) := \int_{\{z \cdot n(y) \leq -t\}} K(y, z) dz$$

Since $\nabla \tau(y) = n(y)$ for any $y \in \partial A$, we have the expansion

$$\tau(y + \varepsilon(tn(y) + z)) = \varepsilon(t + z \cdot n(y)) + O(\varepsilon^2).$$

It now follows from our assumptions on K that for all $y \in \partial A$ and all $t \in [0, 1]$

$$\lim_{\varepsilon \rightarrow 0} a_\varepsilon(y, t) = a(y, t).$$

If Ψ is continuous, the result now follows from dominated convergence. If Ψ is not continuous, then we must work harder.

Let us first assume that K is C^1 in both variables. Changing variables $z' = z + \varepsilon tn(y)$ we can write

$$a_\varepsilon(y, t) = \int_{\{\tau(y+\varepsilon z') \leq 0\}} K(y + \varepsilon tn(y), z' - tn(y)) dz',$$

and hence

$$\partial_t a_\varepsilon(y, t) = \int_{\{\tau(y+\varepsilon z') \leq 0\}} \varepsilon \nabla_y K(y + \varepsilon tn(y), z' - tn(y)) \cdot n(y) dz'$$

$$- \int_{\{\tau(y+\varepsilon z') \leq 0\}} \nabla_{z'} K(y + \varepsilon tn(y), z' - tn(y)) \cdot n(y) dz'.$$

Since the second term is the complete derivative with respect to z' , we can integrate by parts to obtain

$$\partial_t a_\varepsilon(y, t) = \int_{\{\tau(y+\varepsilon z') \leq 0\}} \varepsilon \nabla_y K(y + \varepsilon tn(y), z' - tn(y)) \cdot n(y) dz'$$

$$- \int_{\{\tau(y+\varepsilon z') = 0\}} K(y + \varepsilon tn(y), z' - tn(y)) n(y) \cdot \nabla \tau(y + \varepsilon z') d\mathcal{H}^{d-1}(z')$$

Expanding $\nabla \tau(y + \varepsilon z) = n(y) + O(\varepsilon)$, we see that

$$\partial_t a_\varepsilon(y, t) = - \int_{\{\tau(y+\varepsilon(z+tn(y))) = 0\}} K(y + \varepsilon tn(y), z) d\mathcal{H}^{d-1}(z) + O(\varepsilon),$$

where we note that the constant in the big O bound does not depend on the differentiability of K with respect to the z variable. If we define $f : [0, 1] \rightarrow \mathbb{R}$ by setting

$$f(t) = \mathcal{H}^{d-1}(\{|z| < 1\} \cap \{z : \tau(y + \varepsilon(z + tn(y))) = 0\})$$

then f is uniformly bounded away from zero on any compact subset of $[0, 1)$ and our assumptions on K give us

$$- \int_{\{\tau(y + \varepsilon(z + tn(y))) = 0\}} K(y + \varepsilon tn(y), z) d\mathcal{H}^{d-1}(z) \leq -cf(t).$$

Thus, there exists some $B > 0$ such that for all $y \in \partial A$,

$$\partial_t a_\varepsilon(y, t) \leq B\varepsilon - cf(t).$$

Hence, if we let $\mathcal{L}_{[a,b]}$ denote the Lebesgue measure on the interval $[a, b]$, then for any $\delta > 0$ sufficiently large, we have the following pushforward bound $a_\varepsilon(y, \cdot) \# \mathcal{L}_{[0,1-\delta]} \leq \frac{1}{cf_\delta - B\varepsilon} \mathcal{L}_{[0,1]}$ where $f_\delta = \inf_{t \in [0,1-\delta]} f(t)$.

When K is not differentiable with respect to z , we can first approximate K with a sequence of smooth kernels to obtain the same pushforward bound as above. Since the constant B does not depend on the differentiability of K with respect to z , we can pass to the limit to conclude that the pushforward bound

$$a_\varepsilon(y, \cdot) \# \mathcal{L}_{[0,1-\delta]} \leq \frac{1}{cf_\delta - B\varepsilon} \mathcal{L}_{[0,1]}$$

holds whenever $cf_\delta > B\varepsilon$.

Now let Ψ_n be a sequence of smooth functions converging to Ψ in $L^1([0, 1])$ whose L^∞ norms do not exceed $\|\Psi\|_{L^\infty([0,1])}$. Given any bounded function ϕ , and any $\delta > 0$ such that $cf_\delta > B\varepsilon$ we have, using a change of variables and the pushforward bound, that

$$\begin{aligned} & \left| \int_{\partial A} \int_0^1 \phi(t, y) (\Psi(a_\varepsilon(t, y)) - \Psi_n(a_\varepsilon(t, y))) dt d\mathcal{H}^{d-1}(y) \right| \\ & \leq \|\phi\|_{L^\infty(\partial A \times [0,1])} \left(\frac{1}{cf_\delta - B\varepsilon} \|\Psi - \Psi_n\|_{L^1([0,1])} + 2\delta \|\Psi\|_{L^\infty([0,1])} \right) \mathcal{H}^{d-1}(\partial A). \end{aligned}$$

Hence,

$$\lim_{\delta \rightarrow 0} \lim_{n \rightarrow \infty} \limsup_{\varepsilon \rightarrow 0} \left| \int_{\partial A} \int_0^1 \phi(t, y) (\Psi(a_\varepsilon(t, y)) - \Psi_n(a_\varepsilon(t, y))) dt d\mathcal{H}^{d-1}(y) \right| = 0.$$

This estimate shows us that we can replace Ψ by a continuous approximation. As a result, we can simply argue as we did above in the case where Ψ was continuous. \square

A.8 Examples for weak-* closed hypothesis classes

In this section we continue the discussion in Example 1 and argue why the hypothesis classes considered there are indeed closed in the weak-* topology of $L^\infty(\mathcal{X}; \nu)$. In fact, all these classes are even weak-* compact.

1. The class of all \mathcal{A} -measurable functions $u : \mathcal{X} \rightarrow [0, 1]$ is a bounded subset of $L^\infty(\mathcal{X}; \nu)$ and therefore, by the Banach–Alaoglu theorem, it is weak-* compact.
2. To argue why neural networks with bounded parameters and continuous activation functions are weak-* compact, let $(u_n)_{n \in \mathbb{N}} \subset \mathcal{H}$ be a sequence. Thanks to finite-dimensional compactness a subsequence of the associated parameters converge to some limiting parameters. The continuity of the activations implies that the associated neural networks converge (uniformly in the space of continuous functions on the unit cube $[-1, 1]^d$) to a limiting neural network $u \in \mathcal{H}$. In particular, the convergence is true in the weak-* sense, which shows that \mathcal{H} is weak-* compact.

3. Finally, we consider the class of hard linear classifiers of the form $u(x) = \theta(w \cdot x + b)$ where we assume that the distributions ρ_0, ρ_1 , and \mathfrak{p}_x are such that ν defined in (22) has a density with respect to the Lebesgue measure. A sufficient condition for this to hold is that ρ_0, ρ_1 , and \mathfrak{p}_x have densities with respect to the Lebesgue measure.

If $(u_n)_{n \in \mathbb{N}} \subset \mathcal{H}$ is a sequence of linear classifiers, thanks to finite-dimensional compactness a subsequence (which we do not relabel) of the associated parameters (w_n, b_n) will converge to $w \in \mathbb{R}^d$ with $|w| = 1$ and $b \in [-\infty, \infty]$. For simplicity we only consider the case where $b \neq \pm\infty$. In this case one can define the half-spaces

$$\begin{aligned} A_n &:= \{x \in \mathbb{R}^d : w_n \cdot x + b > 0\}, \\ A &:= \{x \in \mathbb{R}^d : w \cdot x + b > 0\}, \end{aligned}$$

where u_n and u are supported. Then for any $\phi \in L^1(\mathbb{R}^d; \nu)$ it holds

$$\left| \int_{\mathbb{R}^d} (u_n - u) \phi \, d\nu \right| = \left| \int_{A_n} \phi \, d\nu - \int_A \phi \, d\nu \right| \leq \int_{A_n \Delta A} |\phi| \, d\nu$$

where we used the symmetric difference $A_n \Delta A := (A_n \setminus A) \cup (A \setminus A_n)$. Note that this set is either a double cone (if $w_n \neq w$) or a strip of width $|b_n - b|$ (if $w_n = w$).

Since $\phi \in L^1(\mathbb{R}^d; \nu)$ and ν is a probability measure, for every $\varepsilon > 0$ there exists a compact set $K \subset \mathbb{R}^d$ such that $\int_{\mathbb{R}^d \setminus K} |\phi| \, d\nu < \varepsilon$. Using this, we can compute

$$\left| \int_{\mathbb{R}^d} (u_n - u) \phi \, d\nu \right| \leq \int_{A_n \Delta A} |\phi| \, d\nu \leq \int_{(A_n \Delta A) \cap K} |\phi| \, d\nu + \varepsilon.$$

Using that ν has a density with respect to the Lebesgue measure \mathcal{L}^d and using also that $\mathcal{L}^d(A_n \Delta A \cap K) \rightarrow 0$ as $n \rightarrow \infty$, we obtain

$$\lim_{n \rightarrow \infty} \left| \int_{\mathbb{R}^d} (u_n - u) \phi \, d\nu \right| \leq \varepsilon$$

and since $\varepsilon > 0$ was arbitrary we get

$$\lim_{n \rightarrow \infty} \int_{\mathbb{R}^d} (u_n - u) \phi \, d\nu = 0,$$

which implies the weak-* convergence of u_n to $u \in \mathcal{H}$ and hence the weak-* compactness of \mathcal{H} .

Note that for general measures ν the above argument fails. For instance, the sequence of linear classifiers $u_n(x) = \mathbf{1}_{x_1 > -1/n}$ has the natural limit $u(x) = \mathbf{1}_{x_1 > 0}$. However, if $\nu = \delta_0$ then $\int_{\mathbb{R}^d} u_n \, d\nu = 1$ for all $n \in \mathbb{N}$ but $\int_{\mathbb{R}^d} u \, d\nu = 0$, meaning that u is not the weak-* limit of u_n .

B Computational aspects

B.1 Pseudocode for geometric probabilistically robust learning

In Algorithm 1 we provide a pseudocode for parametrized classifiers $f \equiv f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$ based on stochastic gradient descent with batch size B . Furthermore, it involves a sample size of M samples from a distribution \mathfrak{p}_x around an input $x \in \mathcal{X}$, a learning rate η_α for the inner optimization in CVaR, and a learning rate η for the parameter updates. The pseudocode is a straightforward generalization of Robey et al. [2022, Algorithm 1] and we implemented it in their code framework.² The code which can be used to reproduce our results is part of the supplementary material of this paper.

B.2 Training details

Hyperparameter values specific to Algorithm 1 used in training are presented in Table 2. Following Robey et al. [2022], we use AdaDelta [Zeiler, 2012] for MNIST experiments and SGD with

²<https://github.com/arobey1/advbench>

Algorithm 1 Proposed algorithm for solving (20) for $p \in (0, 1)$.

```

1: for minibatch  $(x_j, y_j)_{j=1}^B$  do
2:   for  $T$  steps do                                     ▷ Approximate solution of inner problem
3:     Draw  $x'_k \sim \mathfrak{p}_{x_j}, k = 1, \dots, M$ 
4:      $g_{\alpha_j} \leftarrow 1 - \frac{1}{pM} \sum_{k=1}^M \mathbf{1}_{\ell(f_{\theta}(x'_k), y_j) \geq \alpha_j}$ 
5:      $\alpha_j \leftarrow \alpha_j - \eta_{\alpha} g_{\alpha_j}$ 
6:   end for
7:    $S_j \leftarrow \alpha_j + \frac{1}{pM} \sum_{k=1}^M (\ell(f_{\theta}(x'_k), y_j) - \alpha_j)_+$    ▷ Approximate value of CVaR $_p$ 
8:   if  $S_j > \ell(f_{\theta}(x_j), y_j)$  then                   ▷ If CVaR $_p$  kicks in
9:      $g_j \leftarrow \frac{1}{pM} \sum_{k=1}^M \nabla_{\theta} (\ell(f_{\theta}(x'_k), y_j) - \alpha_j)_+$ 
10:  else                                                 ▷ If it doesn't
11:     $g_j \leftarrow \nabla_{\theta} \ell(f_{\theta}(x_j), y_j)$ 
12:  end if
13:   $g \leftarrow \frac{1}{B} \sum_{j=1}^B g_j$                        ▷ Compute full  $\theta$ -gradient
14:   $v \leftarrow \text{optimizer}(g)$                          ▷ AdaDelta or SGD(+M)
15:   $\theta \leftarrow \theta - \eta v$                          ▷ Update parameters
16: end for

```

momentum for CIFAR-10 experiments. The MNIST experiments use a CNN architecture with two convolutional layers (32 and 64 filters, size 3x3), two dropout layers (dropout probabilities 0.25, 0.5), and two fully connected layers (dimensions 9216 to 128 and 128 to 10). A ResNet-18 [He et al., 2016] is used in the CIFAR-10 experiments. The hyperparameter values used for these algorithms are contained in `hparams_registry.py` in the accompanying code.

Table 2: Hyperparameters used for training. The probability distribution \mathfrak{p}_x is always taken to be the uniform distribution over the ball $B_{\varepsilon}(x)$. Note that p is called beta in the code. For consistency we always used the same hyperparameter values for the ‘‘Geometric’’ and ‘‘Original’’ versions.

| Data | p | ε | η_{α} | M | T |
|----------|------|---------------|-----------------|----|---|
| MNIST | 0.01 | 0.3 | 0.1 | 20 | 5 |
| | 0.1 | 0.3 | 1.0 | 20 | 5 |
| | 0.3 | 0.3 | 1.0 | 20 | 5 |
| | 0.5 | 0.3 | 1.0 | 20 | 5 |
| CIFAR-10 | 0.01 | 8/255 | 0.1 | 20 | 5 |
| | 0.1 | 8/255 | 1.0 | 20 | 5 |
| | 0.3 | 8/255 | 1.0 | 20 | 5 |
| | 0.5 | 8/255 | 1.0 | 20 | 5 |

B.3 Computational resources used

We performed the majority of the prototyping and some experimentation on a LambdaLabs Vector workstation equipped with 3 NVIDIA A6000 GPUs. We estimate that we used approximately 500 GPU-hours on this machine. We supplemented this with 550 GPU-hours of cloud compute—using the Lambda GPU cloud—predominately on instances equipped with a single A10 GPU.