

---

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

---

Anonymous Author(s)

Affiliation

Address

email

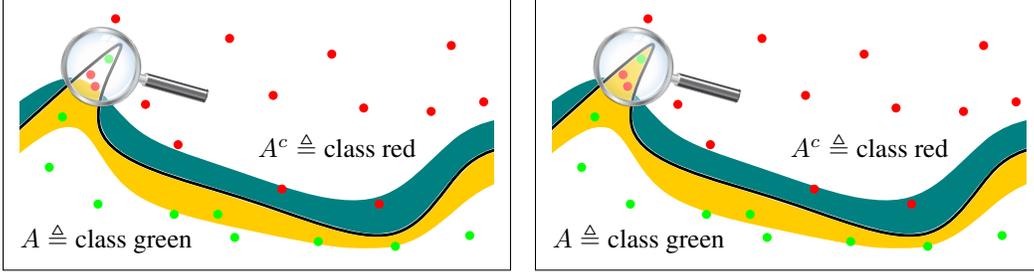
## Abstract

1 Although deep neural networks have achieved super-human performance on many  
2 classification tasks, they often exhibit a worrying lack of robustness towards ad-  
3 versarially generated examples. Thus, considerable effort has been invested into  
4 reformulating Empirical Risk Minimization (ERM) into an adversarially robust  
5 framework. Recently, attention has shifted towards approaches which interpolate  
6 between the robustness offered by adversarial training and the higher clean accu-  
7 racy and faster training times of ERM. In this paper, we take a fresh and geometric  
8 view on one such method—Probabilistically Robust Learning (PRL) [Robey et al.,  
9 2022]. We propose a geometric framework for understanding PRL, which allows  
10 us to identify a subtle flaw in its original formulation and to introduce a family of  
11 probabilistic nonlocal perimeter functionals to address this. We prove existence  
12 of solutions using novel relaxation methods and study properties as well as local  
13 limits of the introduced perimeters.

## 14 1 Introduction

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

26 Beginning with the Probabilistically Robust Learning (PRL) of Robey et al. [2022] we analyze such  
27 frameworks geometrically. This perspective reveals a subtle, paradoxical aspect of PRL: sometimes  
28 the adversary modeled by this framework corrects, instead of exploits, the learner! Fortunately, the  
29 geometric perspective we propose suggests a natural remedy which leads to an interpretation of the  
30 corrected PRL as regularized ERM where a certain nonlocal notion of length (or perimeter) of the  
31 decision boundary acts as a regularizer. We exemplify this correction in Figure 1. The interpretation  
32 of PRL as perimeter-regularized ERM leads us to further generalizations, and we provide a novel  
33 view of the Conditional Value at Risk (CVaR) relaxation of PRL proposed by Robey et al. [2022].



(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.

### 34 1.1 From empirical risk minimization to robustness

35 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  
 36  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ , and a hypothesis class  $\mathcal{H}$ , the standard risk minimization problem is

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

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

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

39 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  
 40 closed) ball of radius  $\varepsilon$  around  $x$ .

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

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

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

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

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

54 To better understand the model (3) we temporarily restrict our attention to binary classification (i.e.,  
 55  $\mathcal{Y} = \{0, 1\}$ ) using indicator functions of admissible sets (i.e.,  $\mathcal{H} := \{\mathbf{1}_A : A \in \mathcal{A}\}$ ). Note that we  
 56 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  
 57 one if  $y \neq \tilde{y}$  and zero otherwise. In this scenario (1) reduces to the geometric problem

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

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

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

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

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

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

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

## 65 1.2 Geometric modification of probabilistically robust learning

66 To motivate our geometric modification of the PRL model from Robey et al. [2022], it is insightful to  
67 investigate the regularization effect that PRL has compared to standard risk minimization. We let  
68  $\rho_i(\bullet) := \mu(\bullet \times \{i\})$  denote the non-normalized conditional distributions of the points with label  $i$ .  
69 Subtracting the standard risk in (4) from the one in (6) and disintegrating using  $\rho_0$  and  $\rho_1$  we obtain

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

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

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

## 82 1.3 Our contributions

83 Our main contributions are the following:

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

## 93 1.4 Related work

94 Adversarial training was developed by Goodfellow et al. [2014], Madry et al. [2017] as an approach  
95 to train networks that are less sensitive to adversarial attacks. Shafahi et al. [2019] reduced its  
96 computational complexity by reusing gradients from the backpropagation when training neural

97 networks. Wong et al. [2020] showed that training with noise perturbations followed by a single  
 98 signed gradient ascent (FGSM) step can be on par with adversarial training while being much  
 99 cheaper. This approach was picked up and improved upon by Andriushchenko and Flammarion  
 100 [2020] based on gradient alignment. Different authors also investigated test-time robustification of  
 101 pretrained classifiers using randomized smoothing [Cohen et al., 2019] or geometric / gradient-based  
 102 approaches [Schwinn et al., 2021, 2022]. While some of the previous models use a combination  
 103 of random perturbations and gradient-based adversarial attacks to robustify classifiers, Robey et al.  
 104 [2022] proposed probabilistically robust learning, which is entirely based on random perturbations.  
 105 PRL aims to interpolate between clean and adversarial accuracy and enjoys the favorable sample  
 106 complexity of vanilla empirical risk minimization; see also Raman et al. [2023] for more insights on  
 107 this issue. Connections between adversarial training and local perimeter regularization of decision  
 108 boundaries were explored by García Trillos and Murray [2022] and then rigorously tied by Bungert  
 109 and Stinson [2022]. Our work is in line with a series of papers [Pydi and Jog, 2021, Awasthi et al.,  
 110 2021a,b, Frank and Niles-Weed, 2022, Frank, 2022, Bungert et al., 2023, García Trillos et al., 2023]  
 111 that explore the existence of solutions to adversarial training problems in different settings. These  
 112 existence proofs involve dealing with different kinds of measurability issues, depending on whether  
 113 open or closed balls  $B_\varepsilon(x)$  are used in the attack model. For open balls one can work with the  
 114 Borel  $\sigma$ -algebra  $\mathcal{A} = \mathfrak{B}(\mathcal{X})$  [Bungert et al., 2023], whereas closed balls require the use of the  
 115 universal  $\sigma$ -algebra to make sure that  $A^{\oplus\varepsilon}$  is measurable [Pydi and Jog, 2021, Awasthi et al., 2021a,b].  
 116 Recently, these results were improved by García Trillos et al. [2023] who also proved for the case of  
 117 multi-class classification that even for the closed ball model Borel measurable classifiers (albeit not  
 118 necessarily indicator functions of measurable sets) exist and that for all but countably many values of  
 119 the adversarial budget  $\varepsilon > 0$  the open and the closed ball models have the same minimal value.

## 120 2 Geometry and existence of probabilistically robust classifiers

### 121 2.1 The binary classification setting with 0-1 loss

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

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

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

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

$$\begin{aligned} \text{ProbPer}(A) = & \int_{\mathcal{X}} \mathbf{1}_{x \in A \vee \mathbb{P}_{x' \sim p_x}[x' \in A] > p} - \mathbf{1}_{x \in A} d\rho_0(x) \\ & + \int_{\mathcal{X}} \mathbf{1}_{x \in A^c \vee \mathbb{P}_{x' \sim p_x}[x' \in A^c] > p} - \mathbf{1}_{x \in A^c} d\rho_1(x) \\ = & \int_{\mathcal{X}} \mathbf{1}_{x \in A^c} \mathbf{1}_{\mathbb{P}_{x' \sim p_x}[x' \in A] > p} d\rho_0(x) + \int_{\mathcal{X}} \mathbf{1}_{x \in A} \mathbf{1}_{\mathbb{P}_{x' \sim p_x}[x' \in A^c] > p} d\rho_1(x). \end{aligned} \quad (9)$$

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

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

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

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

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

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

## 141 2.2 Extensions in the binary classification setting

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

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

144 as well as their corresponding probabilistically robust losses

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

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

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

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

155 The following theorem establishes existence of minimizers of the risk  $\text{ProbR}_\Psi$  for concave and  
 156 non-decreasing functions  $\Psi$ . This existence result is astonishing since the standard method of  
 157 calculus of variations is not directly applicable, with the reason being that problem (15) does not  
 158 provide enough compactness for lower semicontinuity of the perimeter functional  $\text{ProbPer}_\Psi$ . Instead,  
 159 the proof is based on convex relaxations to soft classifiers where we use a lower semicontinuous  
 160 surrogate functional and a total variation defined through a coarea formula which—if  $\Psi$  is concave  
 161 and non-decreasing—lower-bounds the surrogate.

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

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

164 Furthermore,  $\text{ProbR}_\Psi$  can also be interpreted as a sample-wise maximum, analogous to Theorem 1.

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

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

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

170 **2.3 Extension to soft classifiers**

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

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

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

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

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

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

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

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

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

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

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

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

198 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  
 199 distributions  $\rho_0, \rho_1$ , and  $p_x$  are sufficiently nice, then  $\mathcal{H}$  has the desired closedness property.

200 **2.4 Properties and asymptotics of  $\text{ProbPer}_\Psi$**

201 In this section we shall discuss the interpretation of the functional  $\text{ProbPer}_\Psi$  defined in (13) as a  
 202 *perimeter*. We do this in two ways.

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

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

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

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

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

221 **Assumption 2.** We assume that  $\mathcal{X} = \mathbb{R}^d$ ,  $\Psi(0) = 0$ ,  $\Psi$  is measurable and bounded, and  $\rho_1, \rho_0$  have  
 222 continuous densities with respect to the Lebesgue measure which we shall also denote as  $\rho_1, \rho_0$ .  
 223 Furthermore, we assume that there is  $\varepsilon > 0$  and a measurable function  $K : \mathcal{X} \times \mathbb{R}^d \rightarrow [0, \infty)$  such  
 224 that for every  $x \in \mathbb{R}^d$  we have the representation

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

225 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  
 226  $K(x, z) = 0$  if  $|z| > 1$ , and that for every  $z \in \mathbb{R}^d$  the mapping  $x \mapsto K(x, z)$  is  $C^1$ .

227 **Proposition 1.** Under Assumption 2, if  $A$  has a compact  $C^{1,1}$  boundary and either  $\Psi$  is continuous  
 228 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

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

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

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

230 **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  
 231 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  
 232  $\sigma_\Psi = 1$ . However, for  $p \geq \frac{1}{2}$  one easily sees  $\sigma_\Psi = 0$ , hence the limiting perimeter equals zero and  
 233 there is no regularization effect. Using the function  $\Psi(t) = \min\{t/p, 1\}$  corrects this degeneracy.

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

### 246 3 General models

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

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

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

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

255 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  
 256 operator, a tractable version of (18) is

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

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

261 In the binary classification case we can prove the following interesting result that the CVaR relax-  
 262 ation corresponds precisely to using the risk  $\text{ProbR}_\Psi$  with a special piecewise linear and concave  
 263 function  $\Psi$  for which our theory from Section 2.2 applies. In Appendix A.5 we prove a more general  
 264 version of the following statement, replacing the  $[\bullet]_+$  operation in (19) with a Leaky ReLU.

265 **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*

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

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

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

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

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

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

## 274 4 Numerical results

275 We build upon the code of Robey et al. [2022]. The algorithmic realization of (20) is a straightforward  
 276 adaptation of their algorithm, which alternately minimizes the inner optimization problem that  
 277 defines CVaR and the outer optimization to find a suitable classifier, see Algorithm 1 in Appendix B.  
 278 In our experiments, we conduct a comparative analysis between their algorithm (denoted as ‘‘Original’’  
 279 in Table 1) and Algorithm 1 in the appendix which is based on (20) (denoted as ‘‘Geometric’’).  
 280 We report the clean, and adversarial accuracies (subject to PGD attacks), as well as accuracies on  
 281 noise-augmented data and quantile accuracies for different values of  $p$  (see [Robey et al., 2022,  
 282 (6.1)] for the definition) averaged over three runs; see Appendix B.2 for more training details. Our  
 283 experiments are conducted on MNIST and CIFAR-10 and to ensure a fair comparison we adhere to the  
 284 hyperparameter settings described by Robey et al. [2022], such that both the original and geometric  
 285 algorithms utilize the same set of hyperparameters for each specified value of  $p$ . The corresponding  
 286 results for several baseline algorithms including empirical risk minimization and adversarial training  
 287 can be found in their paper. We perform model selection based on the best clean validation accuracy.  
 288 The results in Table 1 show that for moderate values of  $p$  our geometric modification induces higher  
 289 adversarial robustness than the original PRL without loss of clean accuracy (see, in particular, the  
 290 results for MNIST with  $p = 0.1$  and for CIFAR-10 with  $p = 0.3$ ). In the noise augmented metrics as  
 291 well as for extreme values of  $p$  close to 0 or equal to 0.5 both algorithms behave comparably. The  
 292 latter can be expected from out theoretical results, in particular Proposition 1.

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

297 We also remark that our sweep over different values of  $p$  confirms that increasing this parameter  
 298 interpolates between low and high clean accuracies. However, it should be noted that it does not  
 299 necessarily result in a direct interpolation between high and low adversarial or probabilistic accuracy,  
 300 as claimed by Robey et al. [2022]. These observations hold true for both the original algorithm and  
 301 our geometric modification, and despite utilizing their code and hyperparameters, we were unable to  
 302 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	<b>99.20</b>	<b>12.19</b>	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	<b>14.20</b>	99.22	98.70	98.45	97.86
		Original	<b>99.32</b>	8.94	99.22	98.70	98.46	97.80
	0.3	Geometric	<b>99.29</b>	<b>3.02</b>	99.21	98.76	98.53	97.95
		Original	99.27	<b>3.02</b>	99.22	98.77	98.55	98.01
	0.5	Geometric	<b>99.27</b>	<b>1.80</b>	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	<b>81.73</b>	<b>0.24</b>	79.16	74.61	73.19	69.96
	0.1	Geometric	88.15	0.14	85.96	82.55	81.46	78.81
		Original	<b>88.28</b>	<b>0.19</b>	85.61	82.21	81.06	78.28
	0.3	Geometric	<b>90.43</b>	<b>11.80</b>	88.70	85.17	83.93	80.93
		Original	89.97	7.20	88.62	85.07	83.75	80.87
	0.5	Geometric	<b>91.51</b>	1.93	88.94	85.53	84.18	81.21
		Original	90.74	<b>1.99</b>	88.94	85.54	84.35	81.57

## 303 5 Discussion and Conclusion

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

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

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

## 322 References

- 323 Maksym Andriushchenko and Nicolas Flammarion. Understanding and improving fast adversarial  
324 training. *Advances in Neural Information Processing Systems*, 33:16048–16059, 2020.
- 325 Pranjali Awasthi, Natalie S Frank, and Mehryar Mohri. On the existence of the adversarial Bayes  
326 classifier. *Advances in Neural Information Processing Systems*, 34:2978–2990, 2021a.
- 327 Pranjali Awasthi, Natalie S Frank, and Mehryar Mohri. On the existence of the adversarial Bayes  
328 classifier (extended version). *arXiv preprint arXiv:2112.01694*, 2021b.
- 329 Leon Bungert and Kerrek Stinson. Gamma-convergence of a nonlocal perimeter arising in adversarial  
330 machine learning. *arXiv preprint arXiv:2211.15223*, 2022.
- 331 Leon Bungert, Nicolás García Trillos, and Ryan Murray. The geometry of adversarial training in  
332 binary classification. *Information and Inference: A Journal of the IMA*, 12(2):921–968, 06 2023.  
333 ISSN 2049-8772. doi: 10.1093/imaiai/iaac029.
- 334 HanQin Cai, Yuchen Lou, Daniel McKenzie, and Wotao Yin. A zeroth-order block coordinate  
335 descent algorithm for huge-scale black-box optimization. In *International Conference on Machine  
336 Learning*, pages 1193–1203. PMLR, 2021.
- 337 Antonin Chambolle, Massimiliano Morini, and Marcello Ponsiglione. Nonlocal curvature flows.  
338 *Archive for Rational Mechanics and Analysis*, 218:1263–1329, 2015.
- 339 Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh. ZOO: Zeroth order  
340 optimization based black-box attacks to deep neural networks without training substitute models.  
341 In *Proceedings of the 10th ACM workshop on artificial intelligence and security*, pages 15–26,  
342 2017.
- 343 Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. Certified adversarial robustness via randomized  
344 smoothing. In *international conference on machine learning*, pages 1310–1320. PMLR, 2019.
- 345 Nelson Dunford and Jacob T Schwartz. *Linear Operators: General theory*. Linear Operators.  
346 Interscience Publishers, 1958. ISBN 9780470226056.
- 347 Natalie S Frank. Existence and minimax theorems for adversarial surrogate risks in binary classifica-  
348 tion. *arXiv preprint arXiv:2206.09098*, 2022.
- 349 Natalie S Frank and Jonathan Niles-Weed. The consistency of adversarial training for binary  
350 classification. *arXiv preprint arXiv:2206.09099*, 2022.
- 351 Nicolás García Trillos and Ryan Murray. Adversarial classification: Necessary conditions and  
352 geometric flows. *Journal of Machine Learning Research*, 23(187):1–38, 2022.
- 353 Nicolás García Trillos, Matt Jacobs, and Jakwang Kim. On the existence of solutions to adversarial  
354 training in multiclass classification. *arXiv preprint arXiv:2305.00075*, 2023.
- 355 Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial  
356 examples. *arXiv preprint arXiv:1412.6572*, 2014.
- 357 Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul  
358 Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical  
359 analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF International  
360 Conference on Computer Vision*, pages 8340–8349, 2021.
- 361 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.  
362 Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*,  
363 2017.
- 364 José M Mazón, Marcos Solera, and Julián Toledo. The total variation flow in metric random walk  
365 spaces. *Calculus of Variations and Partial Differential Equations*, 59:1–64, 2020.
- 366 Muni Sreenivas Pydi and Varun Jog. The many faces of adversarial risk. *Advances in Neural  
367 Information Processing Systems*, 34:10000–10012, 2021.

- 368 Yao Qin, Nicholas Carlini, Garrison Cottrell, Ian Goodfellow, and Colin Raffel. Imperceptible, robust,  
369 and targeted adversarial examples for automatic speech recognition. In *International conference*  
370 *on machine learning*, pages 5231–5240. PMLR, 2019.
- 371 Joaquin Quinoñero Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D Lawrence. *Dataset*  
372 *shift in machine learning*. Mit Press, 2008.
- 373 Vinod Raman, Unique Subedi, and Ambuj Tewari. On proper learnability between average- and  
374 worst-case robustness. *arXiv preprint arXiv:2211.05656*, 2023.
- 375 Alexander Robey, Luiz Chamon, George J Pappas, and Hamed Hassani. Probabilistically robust  
376 learning: Balancing average and worst-case performance. In *International Conference on Machine*  
377 *Learning*, pages 18667–18686. PMLR, 2022.
- 378 R Tyrrell Rockafellar, Stanislav Uryasev, et al. Optimization of conditional value-at-risk. *Journal of*  
379 *risk*, 2:21–42, 2000.
- 380 Leo Schwinn, An Nguyen, René Raab, Leon Bungert, Daniel Tenbrinck, Dario Zanca, Martin Burger,  
381 and Bjoern Eskofier. Identifying untrustworthy predictions in neural networks by geometric  
382 gradient analysis. In *Uncertainty in Artificial Intelligence*, pages 854–864. PMLR, 2021.
- 383 Leo Schwinn, Leon Bungert, An Nguyen, René Raab, Falk Pulsmeier, Doina Precup, Björn Eskofier,  
384 and Dario Zanca. Improving robustness against real-world and worst-case distribution shifts  
385 through decision region quantification. In *International Conference on Machine Learning*, pages  
386 19434–19449. PMLR, 2022.
- 387 Ali Shafahi, Mahyar Najibi, Mohammad Amin Ghiasi, Zheng Xu, John Dickerson, Christoph Studer,  
388 Larry S Davis, Gavin Taylor, and Tom Goldstein. Adversarial training for free! *Advances in*  
389 *Neural Information Processing Systems*, 32, 2019.
- 390 Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry.  
391 Robustness may be at odds with accuracy. *arXiv preprint arXiv:1805.12152*, 2018.
- 392 Vladimir Vapnik. *The nature of statistical learning theory*. Springer science & business media, 1999.
- 393 Bao Wang, Binjie Yuan, Zuoqiang Shi, and Stanley J Osher. EnResNet: ResNets ensemble via the  
394 Feynman–Kac formalism for adversarial defense and beyond. *SIAM Journal on Mathematics of*  
395 *Data Science*, 2(3):559–582, 2020.
- 396 Eric Wong, Leslie Rice, and J Zico Kolter. Fast is better than free: Revisiting adversarial training.  
397 *arXiv preprint arXiv:2001.03994*, 2020.
- 398 Matthew D Zeiler. Adadelta: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*,  
399 2012.
- 400 Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan.  
401 Theoretically principled trade-off between robustness and accuracy. In *International conference on*  
402 *machine learning*, pages 7472–7482. PMLR, 2019.