
How Does Bayes Error Limit Probabilistic Robust Accuracy

Ruihan Zhang Jun Sun
School of Computing and Information Systems
Singapore Management University
{rhzhang, junsun}@smu.edu.sg

Abstract

Adversarial examples pose a security threat to many critical systems built on neural networks. Given that deterministic robustness often comes with significantly reduced accuracy, probabilistic robustness (*i.e.*, the probability of having the same label with a vicinity is $\geq 1 - \kappa$) has been proposed as a promising way of achieving robustness whilst maintaining accuracy. However, existing training methods for probabilistic robustness still experience non-trivial accuracy loss. It is unclear whether there is an upper bound on the accuracy when optimising towards probabilistic robustness, and whether there is a certain relationship between κ and this bound. This work studies these problems from a Bayes error perspective. We find that while Bayes uncertainty does affect probabilistic robustness, its impact is smaller than that on deterministic robustness. This reduced Bayes uncertainty allows a higher upper bound on probabilistic robust accuracy than that on deterministic robust accuracy. Further, we prove that with optimal probabilistic robustness, each probabilistically robust input is also deterministically robust in a smaller vicinity. We also show that voting within the vicinity always improves probabilistic robust accuracy and the upper bound of probabilistic robust accuracy monotonically increases as κ grows. Our empirical findings also align with our results.

1 Introduction

Neural networks (NNs) achieve remarkable success in various applications, including many security-critical systems [18, 37]. At the same time, several security vulnerabilities in NNs have been identified, including adversarial attacks that generate adversarial examples. Adversarial examples are inputs that are carefully crafted by adding human imperceptible perturbation to normal inputs to trigger wrong predictions [19]. Their presence is particularly concerning in security-critical NN applications.

To defend against adversarial examples, various methods for improving a model’s robustness have been proposed. Adversarial training works by training NNs with a mix of normal and adversarial examples, either pre-generated or generated during training. As it does not carry a formal guarantee on the achieved robustness [49], adversarially trained NNs are potentially vulnerable to new types of adversarial attacks [24, 43]. In contrast, certified training aims to provide a formal guarantee of robustness. A method in this category typically incorporates robustness verification techniques during training [47], *i.e.*, they aim to find a valuation of network parameters such that the model is provably robust to the training samples and some definition of vicinity. However, they are often impractical for multiple reasons including the irreducible errors due to Bayes error (the unavoidable inaccuracy in collecting or labelling training samples) which limits the level of accuracy that is achievable [7, 52].

Recent studies propose that probabilistic robustness, *i.e.*, the probability of having adversarial examples within a vicinity is no more than a certain tolerance level κ , could be sufficient for many practical applications [35, 22, 53]. Furthermore, it is shown to be achievable with a smaller accuracy

drop ($\sim 5\%$ when $\kappa = 0.1$) and less computational cost compared to certified training methods [35]. Probabilistic robustness thus balances between achieving high degrees of security and maintaining accuracy. However, it remains unclear if Bayes errors also limit the achievable performance when optimising towards probabilistic robustness *i.e.*, probabilistic robust accuracy. Further, if an upper limit does exist, how is such a limit correlated to the tolerance level κ ?

In this work, we aim to answer these questions. The Bayes error, in the context of statistics and machine learning, is a fundamental concept related to the inherent uncertainty in any classification system [15]. It represents the minimum error for any classifier on a given problem and is determined by the overlap in the probability distributions of different classes [10]. We remark that the relevance of Bayes error in simple classification tasks may occasionally be questioned given that many datasets, such as MNIST, provide a single, definite label for each input [20]. However, real-world data often lacks this clarity due to inevitable information loss, *e.g.*, during image capture or compression. For instance, the CIFAR-10H dataset showcases that over a third of CIFAR-10 inputs can be re-annotated with uncertain labels by human annotators (CIFAR-10H) [32]. Thus, this uncertainty leads to Bayes errors, which fundamentally constrain not only vanilla accuracy [15] but also deterministic robust accuracy [52] and probabilistic robust accuracy (as what will be shown in this work).

We study the limit on the probabilistic robust accuracy resulting from Bayes error. We first derive an optimal decision rule that maximises probabilistic robust accuracy. Similar to the Bayes classifier [9], the optimal decision rule for probabilistic robustness is also a Maximum A Posteriori (MAP [3]) probability decision, except that the posterior is regarding the vicinity, not a single input. Then, we show that the error from this optimal decision rule regarding probabilistic robustness is lower bounded by the Bayes error for deterministic robustness, but within a much smaller vicinity. After that, a relationship is established between the upper bound of probabilistic robust accuracy and the upper bound of vanilla accuracy or deterministic robust accuracy. We further show that the bound monotonically increases as κ grows. Empirically, we show that our bounds are consistent with what is observed on those probabilistically robust neural networks trained on various distributions.

2 Preliminary and Problem Definition

This section first reviews the background of robustness in machine learning. Then, we recall the Bayes error for the deterministic robustness of classification. Finally, we define our research problem.

2.1 Robustness in Neural Network Classification

We put the context in a K -class classification problem where a classifier $h : \mathbf{x} \mapsto y$ learns to fit a joint distribution D over input space \mathbb{R}^n and label space $\{0, 1, \dots, K - 1\}$. Let $h(\mathbf{x}) \in \{0, 1, \dots, K - 1\}$ denote prediction, and an error captures the difference between $h(\mathbf{x})$ and y . That is, vanilla accuracy is $\Upsilon_{\text{acc}}^+(D, h) = \mathbb{E}_{(\mathbf{x}, y) \sim D} [\mathbf{1}_{h(\mathbf{x})=y}]$ and thus its error is $\Upsilon_{\text{acc}}^-(D, h) = 1 - \Upsilon_{\text{acc}}^+(D, h)$. The capital Upsilon with a plus sign denotes accuracy itself, while a minus sign denotes its corresponding error.

Robustness $\Upsilon_{\text{rob}}^+(D, h, \mathbb{V})$ measures the change in prediction when a perturbation occurs on the input [41]. If the prediction changes when an input is perturbed, then this input is an adversarial example. Formally, an input \mathbf{x}' is an adversarial example of an input-label pair (\mathbf{x}, y) if $(h(\mathbf{x}) = y) \wedge (h(\mathbf{x}') \neq h(\mathbf{x})) \wedge (\mathbf{x}' \in \mathbb{V}(\mathbf{x}))$ [13, 19], where $\mathbb{V}(\mathbf{x})$ is the vicinity at \mathbf{x} . We define robustness to be the probability of *not* observing an adversarial example [23], as defined in Equation (1).

$$\Upsilon_{\text{rob}}^+(D, h, \mathbb{V}) = P_{(\mathbf{x}, y) \sim D} \left((h(\mathbf{x}) = y) \wedge \forall \mathbf{x}' \in \mathbb{V}(\mathbf{x}). h(\mathbf{x}') = h(\mathbf{x}) \right) \quad (1)$$

Remark 2.1 (Vicinity). \mathbf{x} -vicinity also has an equivalent distribution notation $\mathcal{V}(\mathbf{x})$. Let $(v : \mathbb{R}^n \rightarrow \mathbb{R})$ denote its probability density function (PDF) such that $v(\mathbf{x}' - \mathbf{x}) > 0$ means \mathbf{x}' is in \mathbf{x} -vicinity. v is an even and quasiconcave function *e.g.*, the Gaussian function. Appendix A.1 shows the equivalence.

Deterministic Robustness Deterministic robustness requires a zero probability of adversarial examples occurring in a vicinity. It is difficult as achieving $\forall \mathbf{x}' \in \mathbb{V}(\mathbf{x}), h(\mathbf{x}') = h(\mathbf{x})$ is challenging. Although adversarial training [13] empirically reduces adversarial examples [11], it lacks a formal guarantee. Meanwhile, certified training [28] guarantees deterministic robustness through integrating NN verification during training, often resulting in a significant accuracy drop (35%) [21].

Probabilistic Robustness While deterministic robustness is often infeasible without seriously compromising accuracy, probabilistic robustness claims to balance robustness and accuracy [35].

Probabilistic robustness is defined as in Equation (2), where a tolerance level κ limits the probability of having adversarial examples in a vicinity \mathcal{V} . Here, a small portion (such as 1% [53]-50% [8]) of adversarial examples within the vicinity is considered acceptable. Probabilistic robustness is often sufficient in practice [35]. Indeed, safety certification of many safety-critical domains such as aviation requires keeping safety violation probabilities below a non-zero threshold [14].

$$\Upsilon_{\text{prob}}^+(D, h, \mathcal{V}, \kappa) = P_{(\mathbf{x}, \mathbf{y}) \sim D} \left((h(\mathbf{x}) = \mathbf{y}) \wedge ((P_{\mathbf{x}' \sim \mathcal{V}(\mathbf{x})}(h(\mathbf{x}') \neq h(\mathbf{x})) \leq \kappa) \right) \quad (2)$$

2.2 An Upper Bound of Deterministic Robustness from Bayes Error

The Bayes Error In the presence of uncertainty in data distribution, a classifier (no matter how it is trained) inevitably makes some wrong predictions. Bayes error quantifies this inherent uncertainty and represents the irreducible error in accuracy [10, 12, 34], formally captured in Equation (3).

$$\min_{h \in \{\mathbb{R}^n \rightarrow \{0, 1, \dots, K-1\}\}} \Upsilon_{\text{acc}}^-(D, h) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim D} \left[1 - \max_k p(y = k | \mathbf{x}) \right] \quad (3)$$

A classifier achieves the Bayes error when its predictions correspond to the class with maximal posterior probability. Such a classifier is known as a Bayes classifier. The posterior of other classes thus contributes to the irreducible error. An example illustrating the Bayes error is shown in Figure 1a.

Bayes Error for Deterministic Robustness Prior work [52] shows that optimising towards deterministic robustness makes the Bayes error worse. Besides the posterior of other classes, forcing a prediction to be consistent with its neighbours constitutes another source of Bayes error [52]. As in Equation (4), the Bayes error for deterministic robustness can be derived from the Bayes error of a convolved distribution $D' = D * v$. In D' , $p(\mathbf{x}, \mathbf{y})$ is convolved from vicinity v and $p(\mathbf{x}, \mathbf{y})$ in D .

$$\min_{h \in \{\mathbb{R}^n \rightarrow \{0, 1, \dots, K-1\}\}} \Upsilon_{\text{rob}}^-(D, h, \mathbb{V}) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim D'} \left[1 - \max_k p(y = k | \mathbf{x}) \mathbf{1}_{\mathbf{x} \notin \mathbb{K}_{D^\dagger}} \right] \quad (4)$$

where $\Upsilon_{\text{rob}}^-(D, h, \mathbb{V}) = 1 - \Upsilon_{\text{rob}}^+(D, h, \mathbb{V})$. $D^\dagger = [D'] * v$ where $[D']$ is the ‘‘hardened’’ distribution of D' , *i.e.*, one-hot of Argmax posterior. $\mathbb{K}_{D^\dagger} = \{\mathbf{x} \mid (\mathbf{x}, \mathbf{y}) \sim D^\dagger, \max_k p(k | \mathbf{x} = \mathbf{x}) < 1\}$ represents a domain near the boundary where the marginal probability rather than joint probability contributes to the Bayes error for deterministic robustness. Therefore, the Bayes error for deterministic robustness of D is the Bayes error of D' plus the joint probability of non-max classes in \mathbb{K}_{D^\dagger} . As shown in [52], deterministic robustness $\Upsilon_{\text{rob}}^+(D, h, \mathbb{V})$ has an upper bound of 1 minus this irreducible error. Figure 1b illustrates the Bayes error for deterministic robustness.

2.3 Problem Definition

The fundamental problem in this study is finding an upper bound of probabilistic robust accuracy. Further, we aim to establish a relation between this upper bound and the tolerance level, *i.e.*, κ . Formally, we solve the minimisation problem in Equation (5) where $\Upsilon_{\text{prob}}^-(D, h, \mathcal{V}, \kappa) = 1 - \Upsilon_{\text{prob}}^+(D, h, \mathcal{V}, \kappa)$.

$$\min_{h \in \{\mathbb{R}^n \rightarrow \{0, 1, \dots, K-1\}\}} \Upsilon_{\text{prob}}^-(D, h, \mathcal{V}, \kappa) \quad (5)$$

3 Method

In the following, we study the upper bound of probabilistic robust accuracy. We first model the error when optimising towards probabilistic robustness and derive an optimal decision rule. Then, we study the Bayes error obtained from this rule. Further, we formally establish the relationship between the upper bounds of vanilla accuracy, probabilistic robust accuracy, and deterministic robust accuracy.

3.1 Error Modelling and Optimal Decision Rule for Probabilistic Robustness

To find Bayes error when optimising towards probabilistic robustness given distribution $(\mathbf{x}, \mathbf{y}) \sim D$, we first model $\Upsilon_{\text{prob}}^-(D, h, \mathcal{V}, \kappa)$. Intuitively, an error happens if the prediction is wrong or many of the samples in the vicinity are predicted wrongly. We denote the error from the former case as incorrectness and the latter as inconsistency. We analyze each type of error and their combined effect.

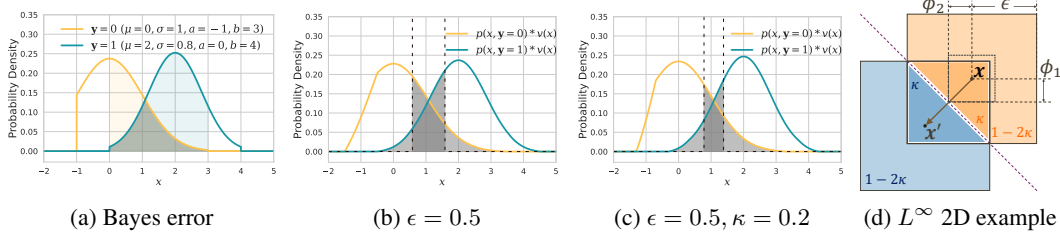


Figure 1: Two truncated normal distributions are used to visualise the Bayes error of (a) vanilla accuracy, (b) deterministic robust accuracy and (c) probabilistic robust accuracy. (d) Example of Corollary 3.6. The nearest adversarial example of \mathbf{x} is at the midpoint of \mathbf{x} and \mathbf{x}' . Both \mathbf{x} and \mathbf{x}' are probabilistically robust but $h(\mathbf{x}) \neq h(\mathbf{x}')$. The dashed box with side length $2\phi_i$ represents $\mathbb{V}_{\downarrow\kappa}(\mathbf{x})$.

Incorrectness The incorrectness for any example (\mathbf{x}, y) is simply $\mathbf{1}_{h(\mathbf{x}) \neq y}$. The incorrectness of a prediction at an input \mathbf{x} given all possible labels $y \in \{0, 1, \dots, K-1\}$ considers posterior at \mathbf{x} , as in Equation (6). Incorrectness is minimum when $h(\mathbf{x})$ equals the class with the highest posterior.

$$e_{\text{cor}}(\mathbf{x}, h; P(y | \mathbf{x})) = \sum_{y=0}^{K-1} P(y = y | \mathbf{x} = \mathbf{x}) \mathbf{1}_{h(\mathbf{x}) \neq y} = 1 - \sum_{y=0}^{K-1} P(y = y | \mathbf{x} = \mathbf{x}) \mathbf{1}_{h(\mathbf{x}) = y} \quad (6)$$

Inconsistency Inconsistency results from prediction at \mathbf{x} being not the same as some of its neighbours. Let $P_{\mathbf{t} \sim \mathcal{V}(\mathbf{x})}(h(\mathbf{t}) \neq h(\mathbf{x}))$ denote the probability of a neighbour of \mathbf{x} having a prediction different from \mathbf{x} . Since this probability is parameterised by $h(\mathbf{x}) \in \{0, 1, \dots, K-1\}$, it can be reformulated as $\sum_{k=0}^{K-1} \mathbf{1}_{h(\mathbf{x})=k} P_{\mathbf{t} \sim \mathcal{V}(\mathbf{x})}(h(\mathbf{t}) \neq k)$. Let $\mu_k(\mathbf{x}) = P_{\mathbf{t} \sim \mathcal{V}(\mathbf{x})}(h(\mathbf{t}) = k)$ and $\sum_{k=0}^{K-1} \mu_k(\mathbf{x}) = 1$. Intuitively, μ_k is the probability of a neighbour predicted as class- k . Thus, the probability of a neighbour of \mathbf{x} having a different prediction from \mathbf{x} can be written as Equation (7).

$$P_{\mathbf{t} \sim \mathcal{V}(\mathbf{x})}(h(\mathbf{t}) \neq h(\mathbf{x})) = \sum_{k=0}^{K-1} \mathbf{1}_{h(\mathbf{x})=k} (1 - \mu_k(\mathbf{x})) = 1 - \sum_{k=0}^{K-1} \mu_k(\mathbf{x}) \mathbf{1}_{h(\mathbf{x})=k} \quad (7)$$

Inconsistency exists when $P_{\mathbf{t} \sim \mathcal{V}(\mathbf{x})}(h(\mathbf{t}) \neq h(\mathbf{x})) > \kappa$. This thresholding can be represented by a unit step function (u) that takes an input $P_{\mathbf{t} \sim \mathcal{V}(\mathbf{x})}(h(\mathbf{t}) \neq h(\mathbf{x})) - \kappa$. Thus, inconsistency at \mathbf{x} is expressed as Equation (8). Also, Lemma 3.1 suggests that κ takes value from $[0, 1/2)$.

$$e_{\text{cns}}(\mathbf{x}, h; \mathcal{V}, \kappa) = u\left(P_{\mathbf{t} \sim \mathcal{V}(\mathbf{x})}(h(\mathbf{t}) \neq h(\mathbf{x})) - \kappa\right) = u\left(1 - \kappa - \sum_{k=0}^{K-1} \mu_k(\mathbf{x}) \mathbf{1}_{h(\mathbf{x})=k}\right) \quad (8)$$

Lemma 3.1. *For the prediction of input \mathbf{x} to be consistent, at most one class has a prediction probability $\geq 1 - \kappa$ in \mathbf{x} -vicinity. Thus, $\kappa < \frac{1}{2}$. (Proof is provided in Appendix B.1.)*

The overall error across the distribution Considering probabilistic robustness, the error at input \mathbf{x} is a combined error of e_{cor} and e_{cns} at \mathbf{x} . We need two intuitions to derive the combined error. First, if $e_{\text{cns}}(\mathbf{x}, h; \mathcal{V}, \kappa) = 1$, the combined error is always 1. Second, if $e_{\text{cns}}(\mathbf{x}, h; \mathcal{V}, \kappa) = 0$, the combined error equals $e_{\text{cor}}(\mathbf{x}, h; P(y|\mathbf{x}))$. Note that e_{cor} takes value from $[0, 1]$ and e_{cns} takes binary value from $\{0, 1\}$. The combined error e is expressed as Equation (9) whose derivation is in Appendix B.2.

$$\begin{aligned} e(\mathbf{x}, h; P(y|\mathbf{x}), \mathcal{V}, \kappa) &= (1 - e_{\text{cns}}(\mathbf{x}, h; \mathcal{V}, \kappa))e_{\text{cor}}(\mathbf{x}, h; P(y | \mathbf{x})) + e_{\text{cns}}(\mathbf{x}, h; \mathcal{V}, \kappa) \\ &= 1 - u\left(\kappa - 1 + \sum_{k=0}^{K-1} \mu_k(\mathbf{x}) \mathbf{1}_{h(\mathbf{x})=k}\right) \left(\sum_{y=0}^{K-1} P(y = y | \mathbf{x} = \mathbf{x}) \mathbf{1}_{h(\mathbf{x})=y}\right) \end{aligned} \quad (9)$$

Note that in general, the errors are functions of $\mathbf{x}, h, D, \mathcal{V}$, and κ . For simplicity, when h, D, \mathcal{V} , or κ can be inferred from the context, we simply omit them, e.g., the simplest case is written as $e_{\text{cor}}(\mathbf{x})$ to denote the incorrectness, $e_{\text{cns}}(\mathbf{x})$ to denote inconsistency, and $e(\mathbf{x})$ to denote the combined error.

To model the distribution-wise error of classifier h on $(\mathbf{x}, y) \sim D$, we compute the expectation of $e(\mathbf{x})$ across D . Formally, $\Upsilon_{\text{prob}}^-(D, h, \mathcal{V}, \kappa) = \int_{\mathbf{x} \in \mathbb{R}^n} e(\mathbf{x}) p(\mathbf{x} = \mathbf{x}) d\mathbf{x}$, where $p(\mathbf{x})$ is the marginal probability in D . Hereby, we get Υ_{prob}^- , the error when optimising towards probabilistic robustness.

To minimise Υ_{prob}^- of any measurable classification function h , we explore the optimal decision rules for probabilistic robustness. From Equation (9), we can establish a Maximum A Posteriori optimal decision rule, whose formal statement is given in Theorem 3.2.

Theorem 3.2. *If h^* is optimal for the probabilistic robustness on a given distribution, i.e., $h^* = \arg \min_h \int_{\mathbf{x} \in \mathbb{R}^n} e(\mathbf{x}) p(\mathbf{x} = \mathbf{x}) d\mathbf{x}$, we would always have $\forall \mathbf{x} \in \mathbb{R}^n, h^*(\mathbf{x}) = \arg \max_k \mu_k(\mathbf{x})$.*

Proof. Let h_1 and h_2 be two distinct classification functions such that $h_1(\mathbf{x}) = \arg \max_k \mu_k(\mathbf{x})$ and $h_2(\mathbf{x}) \neq h_1(\mathbf{x})$. If we can prove $e(\mathbf{x}, h_1) \leq e(\mathbf{x}, h_2)$, then we can know h_1 must be optimal for probabilistic robustness. First, we denote $k_1 = h_1(\mathbf{x})$ and $k_2 = h_2(\mathbf{x}) \neq k_1$. Then,

$$e(\mathbf{x}, h_1) - e(\mathbf{x}, h_2) = u(\kappa - 1 + \mu_{k_2}(\mathbf{x})) P(y = k_2 | \mathbf{x} = \mathbf{x}) - u(\kappa - 1 + \mu_{k_1}(\mathbf{x})) P(y = k_1 | \mathbf{x} = \mathbf{x}). \quad (10)$$

Since $\mu_{k_2}(\mathbf{x}) \leq \mu_{k_1}(\mathbf{x})$, we get $\mu_{k_2}(\mathbf{x}) \leq 1/2$. Recall $\kappa < 1/2$ from Lemma 3.1, we get $\kappa - 1 + \mu_{k_2}(\mathbf{x}) < 1/2 - 1 + 1/2 = 0$. Consequently, we have $u(\kappa - 1 + \mu_{k_2}(\mathbf{x})) = 0$. Therefore, $e(\mathbf{x}, h_1) - e(\mathbf{x}, h_2) = -u(\kappa - 1 + \mu_{k_1}(\mathbf{x})) P(y = k_1 | \mathbf{x} = \mathbf{x}) \leq 0$. This inequality applies to any input \mathbf{x} . Hence, a classification function like h_1 is optimal. An extended proof is in Appendix B.3. \square

Intuitively, the theorem states that when optimising towards probabilistic robustness, a Bayes (optimal) classifier would always classify a sample with the most popular label in the vicinity.

3.2 Bayes Error for Probabilistic Robustness from Bayes Error for Deterministic Robustness

The Bayes classifier, regarding probabilistic robustness, is closely related to the most popular label in its vicinity, leading us to study the properties of μ_k . Intuitively, μ_k is the probability of a neighbour predicted as class- k . Formally, μ_k has an equivalent convolutional form as

$$\mu_k(\mathbf{x}) = P_{\mathbf{t} \sim \mathcal{V}(\mathbf{x})}(h(\mathbf{t}) = k) = \int_{\mathbf{t} \in \mathbb{R}^n} \mathbf{1}_{h(\mathbf{t})=k} v(\mathbf{x} - \mathbf{t}) d\mathbf{t} = (\mathbf{1}_{h(\cdot)=k} * v)(\mathbf{x}), \quad (11)$$

where $*$ denotes convolution and $\mathbf{1}_{h(\cdot)=k}$ denotes an indicator function returning 1 if h of input equals k . v is the probability density function of vicinity distribution, e.g., uniform distribution.

Intuitively, convolution acts as a smoothing operation. Thus, $\mu_k(\mathbf{x})$ is expected to change *gradually* as \mathbf{x} moves in \mathbb{R}^n . Similarly, $\arg \max_k \mu_k(\mathbf{x})$ is unlikely to switch frequently. This implies that under the optimal probabilistic robustness condition, predictions do not change randomly or frequently (in \mathbb{R}^n) but exhibit a form of continuity. Lemma 3.3 and Theorem 3.4 formally states this intuition. Specifically, Lemma 3.3 states that $\mu_k(\mathbf{x})$ changes *gradually* as \mathbf{x} moves in \mathbb{R}^n . Moreover, Theorem 3.4 states that the Bayes classifier (for probabilistic robustness) achieves deterministic robustness with a much smaller vicinity at any input that achieves probabilistic robustness.

Lemma 3.3. *The change in μ_k resulting from shifting an input by a certain distance ϕ within the vicinity is bounded in any direction $\hat{\phi}$. Formally: where \mathbb{S}^{n-1} is the set of all unit vectors in \mathbb{R}^n ,*

$$\forall \mathbf{x} \in \mathbb{R}^n, \forall \phi \in \mathbb{R}, \left(\left(\forall \hat{\phi} \in \mathbb{S}^{n-1}, v\left(\frac{\phi}{2}\hat{\phi}\right) > 0 \right) \rightarrow \forall \hat{\phi} \in \mathbb{S}^{n-1}, \left(\left| \mu_k(\mathbf{x} + \phi\hat{\phi}) - \mu_k(\mathbf{x}) \right| \leq 1 - \min_{\hat{\phi}' \in \mathbb{S}^{n-1}} \int_{\mathbf{t} \in \mathbb{R}^n} \min(v(\mathbf{t} - \phi\hat{\phi}'), v(\mathbf{t})) d\mathbf{t} \right) \right). \quad (12)$$

Proof of Lemma 3.3 is given in Appendix B.4. Essentially, for all inputs shifting a distance ϕ , the μ_k value difference between the original input and the shifted input will be bounded by the 1 minus the minimum overlap between two vicinities that are ϕ apart. From Inequality (12), the correlation between the distance ϕ and maximum change of μ_k can be modelled as a function of ϕ expressed as

$$\Delta_{\max} \mu_k(\phi) = \max_{\hat{\phi} \in \mathbb{S}^{n-1}, \mathbf{x} \in \mathbb{R}^n} \left| \mu_k(\mathbf{x}) - \mu_k(\mathbf{x} + \phi\hat{\phi}) \right| \quad (13)$$

Note that $\Delta_{\max} \mu_k(\phi)$ is a monotonic function, i.e., a greater shift distance ϕ is required if the maximal change in μ_k needs to be increased. Theorem 3.4 leverages this monotonicity to formally show the connection between the Bayes error when optimising towards probabilistic robustness and that towards deterministic robustness.

Theorem 3.4. *If h^* is optimal for the probabilistic robustness on a given distribution, i.e., $h^* = \arg \min_{h \in \{\mathbb{R}^n \rightarrow \{0, 1, \dots, K-1\}\}} \Upsilon_{\text{prob}}^-(D, h, \mathcal{V}, \kappa)$, then there is a lower bound on the distance between an input \mathbf{x} and any of its adversarial examples if probabilistic consistency is satisfied on \mathbf{x} . Formally,*

$$\forall \mathbf{x} \in \mathbb{R}^n. \left((\exists k \in \{0, 1, \dots, K-1\}. \mu_k(\mathbf{x}) > 1 - \kappa) \rightarrow \right. \\ \left. \forall \mathbf{x}' \in \mathbb{R}^n. ((h^*(\mathbf{x}') = k) \vee (|\mathbf{x} - \mathbf{x}'| \geq (\Delta_{\max} \mu_k)^{-1} (1/2 - \kappa))) \right). \quad (14)$$

Proof. Suppose an input \mathbf{x} whose prediction $h(\mathbf{x})$ is consistent, i.e., $\exists k. \mu_k(\mathbf{x}) \geq 1 - \kappa$. Let h^* denote this predicted class. Let ϕ_1 and ϕ_2 denote two scalar distances to shift \mathbf{x} and assume the following features of these two distances. First, $\Delta_{\max} \mu_{k^*}(\phi_1) = 1/2 - \kappa$ and $\Delta_{\mu_{k^*}}(\phi_2, \hat{\phi}_2) > 1/2 - \kappa$. The latter indicates that in some direction, moving the input by a distance of ϕ_2 results in a change in μ_{k^*} greater than $1/2 - \kappa$. Second, $\phi_1 > \phi_2$. From the first condition, we can derive that

$$\min_{\hat{\phi}_1 \in \mathbb{S}^{n-1}} \int_{\mathbf{t} \in \mathbb{R}^n} \min(v(\mathbf{t} - \phi_1 \hat{\phi}_1), v(\mathbf{t})) dt > \int_{\mathbf{t} \in \mathbb{R}^n} \min(v(\mathbf{t} - \phi_2 \hat{\phi}_2), v(\mathbf{t})) dt. \quad (15)$$

In other words, we can find some vector $\phi_2 \hat{\phi}_2$ such that this shift results in a vicinity overlap smaller than the minimum vicinity overlap caused by a ϕ_1 -magnitude shift.

Next, we further shift $\mathbf{x} + \phi_2 \hat{\phi}_2$ along $\hat{\phi}_2$ direction but with the magnitude $\phi_1 - \phi_2$. This new position, $\mathbf{x} + \hat{\phi}_2 \phi_1$, is farther away from \mathbf{x} than $\mathbf{x} + \phi_2 \hat{\phi}_2$ is because $\phi_1 > \phi_2$. Additionally, $\mathbf{x} + \hat{\phi}_2 \phi_1$ results in at most the same vicinity overlap (size) as $\mathbf{x} + \phi_2 \hat{\phi}_2$ does because v is quasiconcave. Formally,

$$\int_{\mathbf{t} \in \mathbb{R}^n} \min(v(\mathbf{t} - \phi_2 \hat{\phi}_2), v(\mathbf{t})) dt \geq \int_{\mathbf{t} \in \mathbb{R}^n} \min(v(\mathbf{t} - \phi_1 \hat{\phi}_2), v(\mathbf{t})) dt. \quad (16)$$

Observe that Inequality (16) contradicts (15). Therefore, the two assumptions cannot hold simultaneously. An adversarial example requires $\mathbf{x}' \in \mathbb{V}(\mathbf{x})$, $\mu_{k^*}(\mathbf{x}') \leq 1/2$. Thus, if $\Delta_{\max} \mu_{k^*}(\phi_1) = 1/2 - \kappa$, the distance between a consistent input and any of its adversarial examples is greater than (or equal to) ϕ_1 . The rationale of $1/2 - \kappa$ is proven in Appendix B.5. \square

Intuitively, at optimal probabilistic robust accuracy, if an input has a probabilistically consistent prediction, all its neighbours in some specific vicinity are predicted the same. Namely, the prediction at this input is deterministically robust within this (likely smaller) vicinity. Corollary 3.5 suggests that the bounds stated in Theorem 3.4 persist even as the shift approaches zero.

Corollary 3.5. *There exists a finite real value such that for all inputs, the directional derivative value with respect to any arbitrary nonzero vector $\hat{\phi}$ (unit vector) does not exceed this finite value and does not fall below the negative of this value. Formally: where \cdot denotes the dot product,*

$$\exists b \in \mathbb{R}, \forall \mathbf{x} \in \mathbb{R}^n, \forall \hat{\phi} \in \mathbb{S}^{n-1}, -b \leq \nabla \mu_k(\mathbf{x}) \cdot \hat{\phi} \leq b. \quad (17)$$

Proof of Corollary 3.5 is provided in Appendix B.6. A one-dimensional example in Appendix B.7 illustrates this idea. We then present an implication of Theorem 3.4 for L^∞ -norm in Corollary 3.6.

Corollary 3.6. *If h^* is optimal for the probabilistic robustness with respect to an L^∞ -vicinity on distribution D , then the vicinity size of the deterministically robust vicinity $\mathbb{V}^{\downarrow \kappa}$ around each probabilistically consistent input is $\epsilon(1 - (2\kappa)^{\frac{1}{n}})$.*

Proof of Corollary 3.6 is in Appendix B.8. We visualise its effect using a two-dimensional example in Figure 1d. Let \mathbf{x}, \mathbf{x}' be two probabilistically consistent inputs and $h(\mathbf{x}) \neq h(\mathbf{x}')$. To minimise $|\mathbf{x}' - \mathbf{x}|$, \mathbf{x}' must be in the diagonal direction (as in Corollary 3.6). The shift from \mathbf{x} to \mathbf{x}' is $-2(\phi_1 \hat{\mathbf{x}}_1 + \phi_2 \hat{\mathbf{x}}_2)$. Each triangle accounts for κ of the original vicinity volume. Box $\mathbb{V}^{\downarrow \kappa}$ has side length $2\phi_i$. Thus, solving $(2\epsilon - 2\phi_i)^2 = 2\kappa(2\epsilon)^2$, we get $\mathbb{V}^{\downarrow \kappa}$ has vicinity size $\phi_1 = \phi_2 = (1 - \sqrt{2\kappa})\epsilon$. Although Corollary 3.6 concerns L^∞ , we can analyse other types similarly, i.e., find the direction with the fastest vicinity overlap decrease and measure the distance to the nearest adversarial example.

In brief, an optimal probabilistic robust accuracy has the following implications. Theorem 3.4 suggests that each probabilistically robust input is also deterministically robust, but within a much smaller vicinity. Corollary 3.6 further quantifies the size of this vicinity. Particularly, the order $1/n$

lets vicinity shrink fast as n grows, making probabilistic robust accuracy higher. Overall, deterministic robust accuracy with the smaller vicinity bounds probabilistic robust accuracy. This is formally captured in Theorem 3.7. The effect of this reduced Bayes uncertainty is illustrated in Figure 1c.

Theorem 3.7. *Given distribution D , vicinity with size ϵ , and tolerance level κ , the probabilistic robust accuracy has an upper bound as shown in Equation (18), where $\mathbb{V}^{\downarrow\kappa}$ or $v^{\downarrow\kappa}$ is the “smaller” vicinity and \mathbb{K} follows the definition in Equation (4), denoting the domain near the boundary.*

$$\begin{aligned} \Upsilon_{\text{prob}}^+(D, h, \mathcal{V}, \kappa) &\leq \max_{h \in \{\mathbb{R}^n \rightarrow \{0, 1, \dots, K-1\}\}} \Upsilon_{\text{rob}}^+(D, h, \mathbb{V}^{\downarrow\kappa}) \\ &= \mathbb{E}_{(\mathbf{x}, y) \sim (D * v^{\downarrow\kappa})} \left[\max_k p(y = k | \mathbf{x}) \mathbf{1}_{\mathbf{x} \notin \mathbb{K}_{\lceil D \rceil * v^{\downarrow\kappa}}} \right] \end{aligned} \quad (18)$$

Proof. Theorem 3.4 infers $\forall \mathbf{x}, ((P_{\mathbf{x}' \sim \mathcal{V}(\mathbf{x})}(h(\mathbf{x}') \neq h(\mathbf{x})) \leq \kappa) \rightarrow \forall \mathbf{x}' \in \mathbb{V}^{\downarrow\kappa}(\mathbf{x}), h(\mathbf{x}') = h(\mathbf{x}))$ at optimal probabilistic robust accuracy. Consider this expression in the form of $\forall \mathbf{x}, \text{Event}_1(\mathbf{x}) \rightarrow \text{Event}_2(\mathbf{x})$. This implies that the occurrence rate of Event 2 (deterministic robust accuracy with $\mathbb{V}^{\downarrow\kappa}$) always upper bounds the rate of Event 1 (probabilistic robust accuracy with \mathbb{V}). Additionally, the deterministic robust accuracy with $\mathbb{V}^{\downarrow\kappa}$ has an upper bound, as what has been shown in [52]. \square

Besides the upper bound, we also find a relatively loose lower bound of probabilistic robust accuracy as shown in Theorem 3.8. This lower bound is the deterministic robust accuracy when the vicinity size is the same as the vicinity size assumed for probabilistic robustness.

Theorem 3.8. *The upper bound of probabilistic robust accuracy monotonically increases as κ grows. Further, for all tolerance levels κ , the upper bound of probabilistic robust accuracy lies between the upper bound of deterministic robust accuracy and the upper bound of vanilla accuracy. Formally,*

$$\begin{aligned} \forall \kappa_1, \kappa_2. \quad (\kappa_1 < \kappa_2) &\rightarrow \min_h \Upsilon_{\text{prob}}^+(D, h, \mathcal{V}, \kappa_1) \leq \min_h \Upsilon_{\text{prob}}^+(D, h, \mathcal{V}, \kappa_2) \\ \forall \kappa. \quad \min_h \Upsilon_{\text{rob}}^+(D, h, \mathbb{V}) &\leq \min_h \Upsilon_{\text{prob}}^+(D, h, \mathcal{V}, \kappa) \leq \min_h \Upsilon_{\text{acc}}^+(D, h) \end{aligned} \quad (19)$$

Proof. Let $\kappa_1 < \kappa_2$ and $k^* = \arg \max_k \mu_k(\mathbf{x})$, and then we get the sign of $e(\mathbf{x}, \kappa_1) - e(\mathbf{x}, \kappa_2)$ is the same as $u(\kappa_2 - 1 + \mu_{k^*}(\mathbf{x})) - u(\kappa_1 - 1 + \mu_{k^*}(\mathbf{x}))$. This is because the posterior probability $P(y | \mathbf{x} = \mathbf{x})$ is non-negative. Since $\kappa_1 < \kappa_2$, and unit step function monotonically increases, we get $e(\mathbf{x}, \kappa_1) - e(\mathbf{x}, \kappa_2) \geq 0$. A smaller κ leads to a lower or equal upper bound of probabilistic robust accuracy. For deterministic robust accuracy, $\kappa = 0$, which is the least value. Thus, deterministic robust accuracy has a smaller upper bound than probabilistic robust accuracy does.

On the other hand, from the intuition of error combination in Section 3.1, we get that the combined error is $e(\mathbf{x}) = 1 - (1 - e_{\text{cns}}(\mathbf{x}, h; \mathcal{V}, \kappa))(1 - e_{\text{cor}}(\mathbf{x}, h; P(y | \mathbf{x})))$. As $0 \leq e_{\text{cns}} \leq 1$, we get $e(\mathbf{x}) \geq 1 - (1 - e_{\text{cor}}(\mathbf{x}, h; \kappa))$. Note that the expectation of $e_{\text{cor}}(\mathbf{x}, h; \kappa)$ is the error in vanilla accuracy. Thus, $\Upsilon_{\text{prob}}^+(D, h, \mathcal{V}, \kappa) \leq \Upsilon_{\text{acc}}^+(D, h)$. So it is with their upper bounds. We also provide an extended proof for this theorem in Appendix B.9. \square

In summary, we show that probabilistic robust accuracy is both lower and upper bounded. We also show why the upper bound of probabilistic robust accuracy can be much greater than that of deterministic robust accuracy. Intuitively, our result suggests that probabilistic robustness indeed allows us to sacrifice much less accuracy, compared to that of deterministic robustness. Furthermore, adopting a larger (more relaxed) κ (up to 1/2) can effectively increase probabilistic robust accuracy.

4 Experiment

We conduct experiments¹ to validate the above established results empirically. Note that Theorem 3.2 infers voting is optimal. Section 3.2 establishes the probabilistic robust accuracy upper bound. Do these match empirical results? Further, ablation experiments on real-world distribution study how this upper bound changes as κ grows. The relationship between accuracy, probabilistic, and deterministic robust accuracy is also studied. In the following, we describe the setups and then answer these questions.

¹Available at <https://github.com/soumission-anonyme/irreducible.git>

Table 1: Probabilistic robustness of classifiers before and after voting. $\kappa = 0.1$ and $\epsilon = 0.15, 0.15, 0.1, 2/255$ for Moons, Chan, FashionMNIST, and CIFAR-10.

	Moons		Chan		FashionMNIST		CIFAR-10	
DA [39]	85.35	85.60 (+0.3%)	67.96	68.86 (+0.9%)	84.12	87.48 (+3.7%)	76.07	81.38 (+5.3%)
RS [8]	84.76	85.18 (+0.4%)	64.67	66.77 (+1.9%)	86.29	88.13 (+2.1%)	87.98	88.95 (+1.0%)
CVaR [35]	85.52	85.66 (+0.1%)	69.46	70.05 (+0.6%)	88.50	91.07 (+2.6%)	90.63	90.77 (+0.1%)

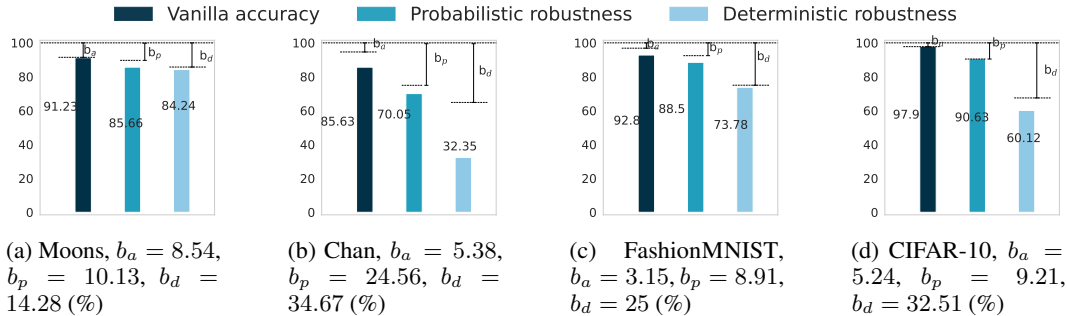


Figure 2: Comparing the SOTA classifier performance with upper bounds of vanilla accuracy (b_a), probabilistic robust accuracy (b_p), and deterministic robust accuracy (b_d).

Setup Our setup follows prior Bayes error studies [15, 52]. We include four datasets: Moons [31], Chan [6], FashionMNIST [46], and CIFAR-10 [17]. Given each dataset, we apply a direct method [15] to compute the Bayes error. L^∞ -vicinity is set with $\epsilon = 0.15, 0.15, 0.1, 2/255$ for defining robustness on respective distribution. For deterministic robustness, the Bayes error follows Equation (4) [52]. For probabilistic robustness, we set $\kappa = 0.1$ by default [35] and vary κ only for the ablation study. More details are in Appendix C.1. Statistic significance is included in Appendix C.2.

Does voting always increase probabilistic robust accuracy empirically? We first compute the probabilistic robust accuracy of some classifier h . We then compute that of a voting classifier h^\dagger , where $h^\dagger(x) = \arg \max_k P_{\mathbf{t} \sim \mathcal{V}(x)}(k = h(\mathbf{t}))$, with sample size 100. They are compared in Table 1. Training algorithms of h include data augmentation (DA [39]), randomised smoothing (RS [8]), and condition value-at-risk (CVaR [35]) which is state-of-the-art (SOTA) for probabilistic robust accuracy. Note that voting always improves probabilistic robust accuracy (at least by +0.1% or on average +1.58%). On DA and RS, the increase is significant (avg + 1.95%), partly because they are not designed specifically for probabilistic robustness. On CVaR, while modest (avg + 0.85%), we do observe an increase. This trend is maintained with a larger voting sample size (Appendix C.3).

Is our upper bound empirically valid on existing neural networks? To check if indeed all trained classifiers respect the theoretical upper bound of probabilistic robust accuracy on any distribution, we compare the SOTA CVaR training and the bound. The middle column of each plot in Figure 2 demonstrates this comparison. We observe that the SOTA probabilistic robust accuracy never exceeds our theoretical bound. Intriguingly, on certain distributions like CIFAR-10, SOTA training almost meets its upper bound with a small gap (0.2%), while on others, a gap remains (on average 4.06%). Theoretically, a negative gap may also occur when a classifier overfits the data samples [15]. Our upper bound is empirically useful in approximating the room for improvement.

How does probabilistic robust accuracy compare to vanilla accuracy and deterministic robust accuracy in terms of upper bounds? We observe in Figure 2 that invariably, the upper bound of probabilistic robust accuracy is lower than that of vanilla accuracy and higher than that of deterministic robust accuracy. In high-dimensional distributions, the upper bound of probabilistic robust accuracy is close to that of vanilla accuracy but over 27% higher than that of deterministic robust accuracy. This could be a result of the curse of dimensionality and much-reduced vicinity size according to Corollary 3.6. On Chan, the upper bound of probabilistic robustness is close to that of deterministic robust accuracy but over 20% lower than that of vanilla accuracy. This could be due to the high-frequency features in the distribution [52]. On Moons, these three bounds are close (at most 7% difference). The reason could be that this distribution is relatively smooth.

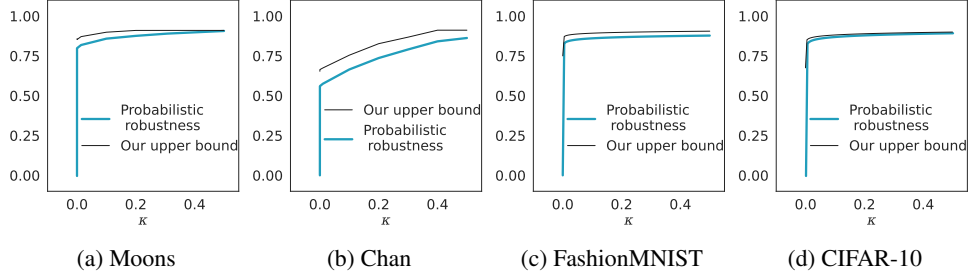


Figure 3: As κ increases, we plot the upper bounds of probabilistic robust accuracy as well as classifiers’ probabilistic robust accuracy change in the Moons and Chan dataset.

What is the effect of κ on the upper bound of probabilistic robust accuracy? Given different κ , the upper bound of probabilistic robust accuracy can be different. We vary κ in $[0, 0.5)$ increasing each time by 0.01. Figure 3 shows that this upper bound monotonically grows as κ grows, which matches Theorem 3.8. Besides, the growth is fast when κ is small (slope > 3 at $\kappa = 0.1$), and the growth rate decreases as κ grows (slope < 0.02 at $\kappa = 0.4$). Especially, for high-dimensional distributions, a small change in κ when κ is small, *e.g.* < 0.1 can significantly increase the upper bound. This could be explained by the $1/n$ order in Corollary 3.6. This is encouraging as it shows that by sacrificing deterministic robustness only slightly, we can already improve the accuracy significantly.

5 Related Work

This work is closely related to studies on Bayes errors and probabilistic robustness. Computing the Bayes error of a given distribution has been studied for over half a century [9], and one interesting topic is to derive or empirically estimate the upper and lower bounds of the Bayes error. Various f -divergences, such as the Bhattacharyya distance [10] or the Henze-Penrose divergence [4, 36], have been studied. Other approaches include directly estimating the Bayes error with f -divergence representation instead of using a bound [30], computing the Bayes error of generative models learned using normalizing flows [16, 42], or evaluate the Bayes error from data sample reassessment [15]. Recent studies apply Bayes error estimation to deterministic robustness beyond a vanilla accuracy perspective [52]. Our study extends this line of research and focuses on probabilistic robustness.

Improving robustness is a core topic in the recent decade [51, 44]. Adversarial training considers adversarial examples in the training phase [26, 13, 48, 45]. However, adversarially trained models do not come with a theoretical guarantee [50, 40, 2]. Certified training provides this guarantee by optimising bounds from formal verification during training but compromises performance on clean inputs [38, 28]. To mitigate these problems, probabilistic robustness methods such as PRoA [53] or CVaR [35] are proposed. Probabilistic robustness offers a desirable balance between robustness and accuracy, making it more applicable in real-world scenarios.

6 Conclusion

We investigate probabilistic robustness by formally deriving its upper bound. We find that the optimal prediction should be the Maximum A Posteriori of predictions in the vicinity. Then, we show that any probabilistically robust input is also deterministically robust within a smaller vicinity. Thus, the upper bound of probabilistic robust accuracy can be obtained from that of deterministic robust accuracy. We verify the empirical impact of this upper bound by comparing it with SOTA training and the upper bounds of vanilla accuracy or deterministic robust accuracy. Experiments match our theorems and show that our bounds could indicate room for improvement in practice.

Limitation A limitation of our upper bound is that Theorem 3.4 requires the posterior probability to calculate the probabilistic robust accuracy upper bound whereas it might be difficult to obtain the posterior for some cases. This generally occurs for Bayes uncertainty analyses [9] in finding bounds of accuracy [42, 29, 27] or deterministic robustness [52]. Yet, this can be compensated by various density estimation methods [33, 52] or reevaluating the probability from a dataset [15].

References

- [1] Anish Athalye, Logan Engstrom, Andrew Ilyas, and Kevin Kwok. Synthesizing robust adversarial examples. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 284–293. PMLR, 10–15 Jul 2018.
- [2] Mislav Balunovic and Martin T. Vechev. Adversarial training and provable defenses: Bridging the gap. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.
- [3] Robert Bassett and Julio Deride. Maximum a posteriori estimators as a limit of bayes estimators. *Mathematical Programming*, 174(1):129–144, Mar 2019.
- [4] Visar Berisha, Alan Wisler, Alfred O. Hero, and Andreas Spanias. Empirically estimable classification bounds based on a nonparametric divergence measure. *IEEE Transactions on Signal Processing*, 64(3):580–591, 2016.
- [5] Saswati Bhattacharya and Mousumi Gupta. A survey on: Facial emotion recognition invariant to pose, illumination and age. In *2019 Second International Conference on Advanced Computational and Communication Paradigms (ICACCP)*, pages 1–6. IEEE, 2019.
- [6] Qingqiang Chen, Fuyuan Cao, Ying Xing, and Jiye Liang. Evaluating classification model against bayes error rate. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(8):9639–9653, 2023.
- [7] Ping-yeh Chiang, Renkun Ni, Ahmed Abdelkader, Chen Zhu, Christoph Studer, and Tom Goldstein. Certified defenses for adversarial patches. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.
- [8] Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. Certified adversarial robustness via randomized smoothing. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 1310–1320. PMLR, PMLR, 09–15 Jun 2019.
- [9] K. Fukunaga and L. Hostetler. k-nearest-neighbor bayes-risk estimation. *IEEE Transactions on Information Theory*, 21(3):285–293, 1975.
- [10] Keinosuke Fukunaga. *Introduction to Statistical Pattern Recognition (2nd Ed.)*. Academic Press Professional, Inc., USA, 1990.
- [11] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario March, and Victor Lempitsky. Domain-adversarial training of neural networks. *Journal of Machine Learning Research*, 17(59):1–35, 2016.
- [12] F.D. Garber and A. Djouadi. Bounds on the bayes classification error based on pairwise risk functions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 10(2):281–288, 1988.
- [13] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- [14] Joris Guerin, Kevin Delmas, and Jérémie Guiochet. Certifying emergency landing for safe urban uav. In *2021 51st Annual IEEE/IFIP International Conference on Dependable Systems and Networks Workshops (DSN-W)*, pages 55–62, 2021.
- [15] Takashi Ishida, Ikko Yamane, Nontawat Charoenphakdee, Gang Niu, and Masashi Sugiyama. Is the performance of my deep network too good to be true? a direct approach to estimating the bayes error in binary classification. In *The Eleventh International Conference on Learning Representations*, 2023.
- [16] Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- [17] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.

- [18] Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. Adversarial examples in the physical world. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Workshop Track Proceedings*. OpenReview.net, 2017.
- [19] Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. Adversarial machine learning at scale. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017.
- [20] Yann LeCun, Corinna Cortes, and Chris Burges.
- [21] Linyi Li, Tao Xie, and Bo Li. Sok: Certified robustness for deep neural networks. In *44th IEEE Symposium on Security and Privacy, SP 2023, San Francisco, CA, USA, 22-26 May 2023*. IEEE, 2023.
- [22] Renjue Li, Pengfei Yang, Cheng-Chao Huang, Youcheng Sun, Bai Xue, and Lijun Zhang. Towards practical robustness analysis for dnns based on pac-model learning. In *Proceedings of the 44th International Conference on Software Engineering, ICSE '22*, page 2189–2201, New York, NY, USA, 2022. Association for Computing Machinery.
- [23] Wang Lin, Zhengfeng Yang, Xin Chen, Qingye Zhao, Xiangkun Li, Zhiming Liu, and Jifeng He. Robustness verification of classification deep neural networks via linear programming. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11418–11427, June 2019.
- [24] Hao Liu, Xiangyu Zhu, Zhen Lei, and Stan Z. Li. Adaptiveface: Adaptive margin and sampling for face recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [25] Xingjun Ma, Bo Li, Yisen Wang, Sarah M. Erfani, Sudanthi N. R. Wijewickrema, Grant Schoenebeck, Dawn Song, Michael E. Houle, and James Bailey. Characterizing adversarial subspaces using local intrinsic dimensionality. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018.
- [26] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018.
- [27] Kevin R. Moon, Alfred O. Hero, and B. Véronique Delouille. Meta learning of bounds on the bayes classifier error. In *2015 IEEE Signal Processing and Signal Processing Education Workshop (SP/SPE)*, pages 13–18, 2015.
- [28] Mark Niklas Müller, Franziska Eckert, Marc Fischer, and Martin T. Vechev. Certified training: Small boxes are all you need. *CoRR*, abs/2210.04871, 2022.
- [29] Frank Nielsen. Generalized bhattacharyya and chernoff upper bounds on bayes error using quasi-arithmetic means. *Pattern Recognition Letters*, 42:25–34, 2014.
- [30] Morteza Noshad, Li Xu, and Alfred Hero. Learning to benchmark: Determining best achievable misclassification error from training data. *arXiv preprint arXiv:1909.07192*, 2019.
- [31] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [32] Joshua C. Peterson, Ruairidh M. Battleday, Thomas L. Griffiths, and Olga Russakovsky. Human uncertainty makes classification more robust. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.
- [33] Cédric Renggli, Luka Rimanic, Nora Hollenstein, and Ce Zhang. Evaluating bayes error estimators on real-world datasets with feebee. In Joaquin Vanschoren and Sai-Kit Yeung, editors, *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual*, 2021.
- [34] Brian D. Ripley. *Pattern Recognition and Neural Networks*. Cambridge University Press, 1996.
- [35] Alexander Robey, Luiz Chamon, George J. Pappas, and Hamed Hassani. Probabilistically robust learning: Balancing average and worst-case performance. In Kamalika Chaudhuri, Stefanie

- Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 18667–18686. PMLR, 17–23 Jul 2022.
- [36] Salimeh Yasaei Sekeh, Brandon Oselio, and Alfred O. Hero. Learning to bound the multi-class bayes error. *IEEE Transactions on Signal Processing*, 68:3793–3807, 2020.
- [37] Mahmood Sharif, Sruti Bhagavatula, Lujo Bauer, and Michael K. Reiter. Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security, CCS ’16*, pages 1528–1540, New York, NY, USA, 2016. Association for Computing Machinery.
- [38] Zhouxing Shi, Yihan Wang, Huan Zhang, Jinfeng Yi, and Cho-Jui Hsieh. Fast certified robust training with short warmup. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 18335–18349. Curran Associates, Inc., 2021.
- [39] Connor Shorten and Taghi M. Khoshgoftaar. A survey on image data augmentation for deep learning. *Journal of Big Data*, 6(1):60, Jul 2019.
- [40] Gagandeep Singh, Timon Gehr, Markus Püschel, and Martin Vechev. An abstract domain for certifying neural networks. *Proc. ACM Program. Lang.*, 3(POPL), jan 2019.
- [41] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In Yoshua Bengio and Yann LeCun, editors, *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, 2014.
- [42] Ryan Theisen, Huan Wang, Lav R Varshney, Caiming Xiong, and Richard Socher. Evaluating state-of-the-art classification models against bayes optimality. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 9367–9377. Curran Associates, Inc., 2021.
- [43] Florian Tramer, Nicholas Carlini, Wieland Brendel, and Aleksander Madry. On adaptive attacks to adversarial example defenses. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1633–1645. Curran Associates, Inc., 2020.
- [44] Jingyi Wang, Jialuo Chen, Youcheng Sun, Xingjun Ma, Dongxia Wang, Jun Sun, and Peng Cheng. Robot: Robustness-oriented testing for deep learning systems. In *2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*, pages 300–311, 2021.
- [45] Yisen Wang, Difan Zou, Jinfeng Yi, James Bailey, Xingjun Ma, and Quanquan Gu. Improving adversarial robustness requires revisiting misclassified examples. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.
- [46] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *CoRR*, abs/1708.07747, 2017.
- [47] Kaidi Xu, Zhouxing Shi, Huan Zhang, Yihan Wang, Kai-Wei Chang, Minlie Huang, Bhavya Kaillkhura, Xue Lin, and Cho-Jui Hsieh. Automatic perturbation analysis for scalable certified robustness and beyond. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1129–1141. Curran Associates, Inc., 2020.
- [48] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan. Theoretically principled trade-off between robustness and accuracy. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 7472–7482. PMLR, 09–15 Jun 2019.
- [49] Huan Zhang, Hongge Chen, Zhao Song, Duane S. Boning, Inderjit S. Dhillon, and Cho-Jui Hsieh. The limitations of adversarial training and the blind-spot attack. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019.
- [50] Huan Zhang, Tsui-Wei Weng, Pin-Yu Chen, Cho-Jui Hsieh, and Luca Daniel. Efficient neural network robustness certification with general activation functions. In *Proceedings of the 32nd*

International Conference on Neural Information Processing Systems, volume 31 of *NIPS'18*, page 4944–4953, Red Hook, NY, USA, 2018. Curran Associates Inc.

- [51] Quan Zhang, Yongqiang Tian, Yifeng Ding, Shanshan Li, Chengnian Sun, Yu Jiang, and Jiaguang Sun. Coophance: Cooperative enhancement for robustness of deep learning systems. In *Proceedings of the 32nd ACM SIGSOFT International Symposium on Software Testing and Analysis, ISSTA 2023*, page 753–765, New York, NY, USA, 2023. Association for Computing Machinery.
- [52] Ruihan Zhang and Jun Sun. Certified robust accuracy of neural networks are bounded due to bayes errors. In *Computer Aided Verification*. Springer International Publishing, 2024.
- [53] Tianle Zhang, Wenjie Ruan, and Jonathan E. Fieldsend. Proa: A probabilistic robustness assessment against functional perturbations. In Massih-Reza Amini, Stéphane Canu, Asja Fischer, Tias Guns, Petra Kralj Novak, and Grigorios Tsoumakas, editors, *Machine Learning and Knowledge Discovery in Databases*, pages 154–170, Cham, 2023. Springer Nature Switzerland.

A Notations

A.1 Vicinity Notations

Definition A.1 (Vicinity). Inputs within the vicinity of an input \mathbf{x} are imperceptible from \mathbf{x} . We call this vicinity an \mathbf{x} -vicinity. To capture imperceptibility, an \mathbf{x} -vicinity can be denoted in (at least four) different but equivalent notations, *i.e.*, distance-threshold, set, distribution, and function notations. Occasionally, an input within \mathbf{x} -vicinity is called a neighbour of \mathbf{x} .

Distance-threshold Notation of Vicinity The distance-threshold notation is one of the earliest notations to depict vicinity [25, 1, 5]. Here, the neighbour of a sample $\mathbf{x} \in \mathbb{X}$ refers to an input that lies within a certain threshold distance from \mathbf{x} . Formally, \mathbf{x}' is within \mathbf{x} -vicinity if and only if

$$d(\mathbf{x}', \mathbf{x}) \leq \epsilon \quad (20)$$

where d measures the distance between two inputs, and this distance needs to be smaller than a threshold ϵ to be considered imperceptible.

Specifically, the distance function can be defined in a variety of ways *e.g.*, L^p norm ($p = 0, 1, 2$, or ∞) or domain-specific transformations that preserve labels, such as tilting or zooming.

$$d(\mathbf{x}', \mathbf{x}) = \|\mathbf{x}' - \mathbf{x}\|_p, \quad (\text{Additive in } L^p \text{ norm}), \text{ or} \quad (21)$$

$$d(\mathbf{x}', \mathbf{x}) = \begin{cases} |\epsilon|, & \text{if } f_{\text{transform}}(\mathbf{x}, \epsilon) = \mathbf{x}', \\ \epsilon + 1, & \text{otherwise} \end{cases}$$

where the transformation function $f_{\text{transform}} : \mathbb{X} \rightarrow \mathbb{X}$ can be, for example, an image rotation with a parameter determining the degree of rotation.

Set Notation of Vicinity Here, the vicinity of a sample $\mathbf{x} \in \mathbb{X}$ refers to a set containing all neighbours of \mathbf{x} . Given an input, all inputs whose distance to the given input is within a certain threshold form a set, defined as a vicinity of the given input. For any $\mathbf{x} \in \mathbb{X}$, its vicinity is expressed as

$$\mathbb{V}(\mathbf{x}) = \{\mathbf{x}' \mid d(\mathbf{x}, \mathbf{x}') \leq \epsilon\} \quad (22)$$

Since the corresponding distance function d can be a representation of different distance measures, the set notation $\mathbb{V}(\mathbf{x})$ can also be a representation of various sets, *i.e.*, $\mathbb{V}_1(\mathbf{x})$, $\mathbb{V}_2(\mathbf{x})$ could be \mathbf{x} -vicinities defined in two different ways.

Function Notation of Vicinity The set or distance representation may be inconvenient sometimes [52]. We may sometimes need the notion of $\mathbf{1}_{\mathbf{x}' \in \mathbb{V}(\mathbf{x})}$ to quantify if \mathbf{x}' is a neighbour of \mathbf{x} . For instance, if we would like to sum the marginal probability of all neighbours of \mathbf{x} , we can $\int_{\mathbb{X}} \mathbf{1}_{\mathbf{x}' \in \mathbb{V}(\mathbf{x})} p(\mathbf{x} = \mathbf{x}') d\mathbf{x}'$ instead of $\int_{\mathbb{V}(\mathbf{x})} p(\mathbf{x} = \mathbf{x}') d\mathbf{x}'$ to avoid a varying interval of integration.

In this case, a vicinity function, which is an equivalent form of the set, can be defined as

$$v_{\mathbf{x}}(\mathbf{x}') = \begin{cases} \left(\int_{\mathbb{V}(\mathbf{x})} d\mathbf{x}'' \right)^{-1}, & \text{if } \mathbf{x}' \in \mathbb{V}(\mathbf{x}) \\ 0, & \text{otherwise} \end{cases} \quad (23)$$

Essentially, Equation (23) can be viewed as a probability density function (PDF) uniformly defined over the vicinity around an input \mathbf{x} . Now we shift the x-coordinate by \mathbf{x} , we get

$$v_{\mathbf{0}}(\mathbf{x}' - \mathbf{x}) = \begin{cases} \left(\int_{\mathbb{V}(\mathbf{0})} d\mathbf{x}'' \right)^{-1}, & \text{if } \mathbf{x}' - \mathbf{x} \in \mathbb{V}(\mathbf{0}) \\ 0, & \text{otherwise} \end{cases} \quad (24)$$

Assuming that the vicinity function is translation invariant, we can drop the subscript $\mathbf{0}$, and use a positive constant ϵ_v to represent $\int_{\mathbb{V}(\mathbf{0})} d\mathbf{x}''$, *i.e.*, the volume of the vicinity. Thus, the vicinity function $v : \mathbb{X} \rightarrow \{0, \epsilon_v^{-1}\}$ can be expressed as

$$v(\mathbf{x}) = \begin{cases} \epsilon_v^{-1} & \text{if } \mathbf{x} \in \mathbb{V}(\mathbf{0}) \\ 0, & \text{otherwise} \end{cases} \quad (25)$$

An example of a one-dimensional input's vicinity is shown in Figure 4.

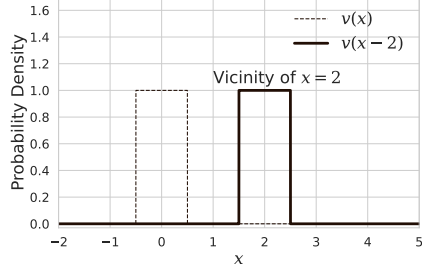


Figure 4: 1D visualizations of vicinity function. This vicinity function is a rectangular function that returns a constant value if an input is in the vicinity. Vicinity function $v(\mathbf{x})$ is shown in dashed line ($\epsilon = 0.5$). To get the vicinity at a specific input $\mathbf{x} = 2$, we simply shift $v(\mathbf{x})$ along the positive direction of the x-axis by 2.

Distribution Notation of Vicinity When we need to sample from the vicinity, we need its distribution notation. The distribution notation for \mathbf{x} -vicinity is $\mathcal{V}(\mathbf{x})$ whose PDF is denoted as $v : \mathbb{R}^n \rightarrow \mathbb{R}$. If $v(\mathbf{x}' - \mathbf{x}) > 0$, we say \mathbf{x}' is within \mathbf{x} -vicinity.

An imperceptible perturbation from any \mathbf{x} to \mathbf{x}' means that \mathbf{x}' is a ‘neighbour’ of \mathbf{x} , or equivalently, \mathbf{x}' is in the \mathbf{x} -vicinity. \mathbf{x} -vicinity is a (probabilistic) distribution $\mathcal{V}(\mathbf{x})$ centred at \mathbf{x} . A standard vicinity $\mathcal{V}(\mathbf{0})$ is centred at the origin and its PDF is denoted as $v : \mathbb{R}^n \rightarrow \mathbb{R}$. Thus, the PDF centred any specific \mathbf{x} would be $v(\mathbf{x}' - \mathbf{x})$.

v is typically an even and quasiconcave function. Formally,

$$\begin{aligned} v(\mathbf{x}) &= v(-\mathbf{x}) \\ \forall t \in [0, 1] \forall \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^n, \quad &v(t\mathbf{x}_1 + (1-t)\mathbf{x}_2) \geq \min(v(\mathbf{x}_1), v(\mathbf{x}_2)) \end{aligned} \quad (26)$$

A uniform L^p \mathbf{x} -vicinity assumes that all inputs outside an L^p -norm of \mathbf{x} are distinguishable from \mathbf{x} and all inputs within this norm are equally imperceptible from \mathbf{x} . This L^p vicinity function is captured in Equation (27) where parameter ϵ specifies a size.

$$v(\mathbf{x}' - \mathbf{x}) = \begin{cases} \frac{\Gamma(1+n/p)}{(2\epsilon\Gamma(1+1/p))^n}, & \text{if } \|\mathbf{x}' - \mathbf{x}\|_p \leq \epsilon \\ 0, & \text{otherwise} \end{cases} \quad (27)$$

The fraction in Equation (27) represents the inverse of the L^p -norm volume, where Γ denotes the gamma function. Vicinity functions assess the likelihood of \mathbf{x}' being a neighbour to \mathbf{x} . In the uniform L^p -norm context, all inputs within the norm are equally valid as neighbours and no inputs outside the norm are neighbours.

Remark A.2. Adversarial examples of an input \mathbf{x} always reside \mathbf{x} -vicinity.

B Complete Proofs and Derivations

This section provides detailed proofs for various lemmas, theorems, or corollaries. For each, we restate the original claim followed by a more comprehensive proof than what appears in the main text. Additionally, we include detailed derivations for certain equations not covered in the theorems.

B.1 Proof of Lemma 3.1

Lemma 3.1. *For the predictions within the vicinity (of an input \mathbf{x}) to be consistent, at most one class has a prediction probability exceeding $1 - \kappa$ in this vicinity. Thus, $\kappa < \frac{1}{2}$.*

Proof. Assume that $\kappa > \frac{1}{2}$ such that any $\sum_{k=0}^{K-1} \mu_k(\mathbf{x}) \mathbf{1}_{h(\mathbf{x})=k}$ that is greater than or equal to $\frac{1}{2}$ satisfies the consistency condition because $\frac{1}{2} > 1 - \kappa$. Hence, there may exist some $k_1 \neq k_2$ such that $\mu_{k_1} + \mu_{k_2} \leq 1$ and

$$\sum_{k=0}^{K-1} \mu_{k_1}(\mathbf{x}) \mathbf{1}_{h(\mathbf{x})=k_1} = \sum_{k=0}^{K-1} \mu_{k_2}(\mathbf{x}) \mathbf{1}_{h(\mathbf{x})=k_2} > 1 - \kappa. \quad (28)$$

The existence of such distinct indices k_1, k_2 implies that if prediction at \mathbf{x} is k_1 , it is consistent with its neighbours. Similarly, if prediction at \mathbf{x} is k_2 , it is consistent with its neighbours. However, since $k_1 \neq k_2$, it is not possible for the neighbours' predictions to simultaneously be consistent with both k_1 and k_2 . This scenario contradicts Inequality (28), and thus the initial assumption does not hold. \square

B.2 Derivation of the Combined Error Considering probabilistic robustness

Considering probabilistic robustness, the error at input \mathbf{x} is a combined error of e_{cor} and e_{cns} at \mathbf{x} . As discussed, we use two intuitions to derive the combined error. First, if $e_{\text{cns}}(\mathbf{x}, h; \kappa) = 1$, the combined error is always 1. Second, if $e_{\text{cns}}(\mathbf{x}, h; \mathcal{V}, \kappa) = 0$, the combined error equals $e_{\text{cor}}(\mathbf{x}, h; P(y | \mathbf{x}))$. Note that inconsistency is a binary value that takes either 0 or 1 while incorrectness takes a real value from 0 to 1 depending on the posterior at input. In the following, we derive the combined error $e(\mathbf{x}, h; P(y | \mathbf{x}), \kappa)$ as expressed in Equation (9).

$$\begin{aligned}
e(\mathbf{x}, h; P(y|\mathbf{x}), \kappa) &= (1 - e_{\text{cns}}(\mathbf{x}, h; \mathcal{V}, \kappa))e_{\text{cor}}(\mathbf{x}, h; P(y | \mathbf{x})) + e_{\text{cns}}(\mathbf{x}, h; \mathcal{V}, \kappa) \\
&= e_{\text{cor}}(\mathbf{x}, h; P(y | \mathbf{x})) - e_{\text{cor}}(\mathbf{x}, h; P(y | \mathbf{x}))e_{\text{cns}}(\mathbf{x}, h; \mathcal{V}, \kappa) + e_{\text{cns}}(\mathbf{x}, h; \mathcal{V}, \kappa) \\
&= e_{\text{cor}}(\mathbf{x}, h; P(y | \mathbf{x})) + e_{\text{cns}}(\mathbf{x}, h; \mathcal{V}, \kappa) - e_{\text{cor}}(\mathbf{x}, h; P(y | \mathbf{x}))e_{\text{cns}}(\mathbf{x}, h; \mathcal{V}, \kappa) \\
&= 1 - (1 - e_{\text{cns}}(\mathbf{x}, h; \mathcal{V}, \kappa))(1 - e_{\text{cor}}(\mathbf{x}, h; P(y | \mathbf{x}))) \\
&= 1 - \left(1 - u \left(1 - \kappa - \sum_{k=0}^{K-1} \mu_k(\mathbf{x}) \mathbf{1}_{h(\mathbf{x})=k}\right)\right) \left(\sum_{y=0}^{K-1} P(y = y | \mathbf{x} = \mathbf{x}) \mathbf{1}_{h(\mathbf{x})=y}\right) \\
&= 1 - u \left(\kappa - 1 + \sum_{k=0}^{K-1} \mu_k(\mathbf{x}) \mathbf{1}_{h(\mathbf{x})=k}\right) \left(\sum_{y=0}^{K-1} P(y = y | \mathbf{x} = \mathbf{x}) \mathbf{1}_{h(\mathbf{x})=y}\right)
\end{aligned} \tag{29}$$

B.3 Extended Proof of Theorem 3.2

Theorem 3.2. *If h^* is optimal for the probabilistic robustness on a given distribution, i.e., $h^* = \arg \min_h \int_{\mathbf{x} \in \mathbb{R}^n} e(\mathbf{x})p(\mathbf{x} = \mathbf{x})d\mathbf{x}$, we would always have $\forall \mathbf{x} \in \mathbb{R}^n, h^*(\mathbf{x}) = \arg \max_k \mu_k(\mathbf{x})$.*

Proof. Let h_1 and h_2 be two distinct classification functions such that $h_1(\mathbf{x}) = \arg \max_k \mu_k(\mathbf{x})$ and $h_2(\mathbf{x}) \neq h_1(\mathbf{x})$. We want to prove $e(\mathbf{x}, h_1) \leq e(\mathbf{x}, h_2)$, such that h_1 must be optimal for probabilistic robustness. First, we denote $k_1, k_2 \in \{0, 1, \dots, K-1\}$, $k_1 = h_1(\mathbf{x}), k_2 = h_2(\mathbf{x}) \neq k_1$. Then,

$$\begin{aligned}
e(\mathbf{x}, h_1) - e(\mathbf{x}, h_2) &= 1 - u \left(\kappa - 1 + \sum_{k=0}^{K-1} \mu_k(\mathbf{x}) \mathbf{1}_{h_1(\mathbf{x})=k}\right) \left(\sum_{y=0}^{K-1} P(y = y | \mathbf{x} = \mathbf{x}) \mathbf{1}_{h_1(\mathbf{x})=y}\right) \\
&\quad - 1 + u \left(\kappa - 1 + \sum_{k=0}^{K-1} \mu_k(\mathbf{x}) \mathbf{1}_{h_2(\mathbf{x})=k}\right) \left(\sum_{y=0}^{K-1} P(y = y | \mathbf{x} = \mathbf{x}) \mathbf{1}_{h_2(\mathbf{x})=y}\right) \\
&= u(\kappa - 1 + \mu_{k_2}(\mathbf{x}))P(y = k_2 | \mathbf{x} = \mathbf{x}) - u(\kappa - 1 + \mu_{k_1}(\mathbf{x}))P(y = k_1 | \mathbf{x} = \mathbf{x}).
\end{aligned} \tag{30}$$

Since $\mu_{k_2}(\mathbf{x}) \leq \mu_{k_1}(\mathbf{x})$, we get $\mu_{k_2}(\mathbf{x}) \leq 1/2$. Recall $\kappa < 1/2$ from Lemma 3.1, we get

$$\kappa - 1 + \mu_{k_2}(\mathbf{x}) < \frac{1}{2} - 1 + \frac{1}{2} = 0 \tag{31}$$

Consequently, when the input of a unit step function is negative, we have $u(\kappa - 1 + \mu_{k_2}(\mathbf{x})) = 0$. Therefore, we get the following expression where the value of a unit step function and the value of conditional probability are both non-negative.

$$e(\mathbf{x}, h_1) - e(\mathbf{x}, h_2) = -u(\kappa - 1 + \mu_{k_1}(\mathbf{x}))P(y = k_1 | \mathbf{x} = \mathbf{x}) \leq 0 \tag{32}$$

Hence, the error of h_1 is no greater than the error of h_2 . This inequality applies to any input \mathbf{x} . Hence, a classification function like h_1 is optimal. \square

B.4 Proof of Lemma 3.3

Lemma 3.3. *A change in μ_k results from shifting an input by a certain distance ϕ within the vicinity. This change is bounded in any direction $\hat{\phi}$. Formally,*

$$\forall \mathbf{x} \in \mathbb{R}^n, \forall \phi \in \mathbb{R}, \left(\left(\forall \hat{\phi} \in \mathbb{S}^{n-1}, v\left(\frac{\phi}{2}\hat{\phi}\right) > 0 \right) \rightarrow \right. \\ \left. \forall \hat{\phi} \in \mathbb{S}^{n-1}, \left(\left| \mu_k(\mathbf{x} + \phi\hat{\phi}) - \mu_k(\mathbf{x}) \right| \leq \left| 1 - \min_{\hat{\phi}' \in \mathbb{S}^{n-1}} \int_{\mathbf{t} \in \mathbb{R}^n} \min(v(\mathbf{t} - \phi\hat{\phi}'), v(\mathbf{t})) dt \right| \right) \right), \quad (33)$$

where \mathbb{S}^{n-1} denotes the set of all unit vectors in \mathbb{R}^n .

Proof. Each μ_k has a convolutional form as provided in Equation (11). Therefore, the change of μ_k resulting from a shift with magnitude ϕ in direction $\hat{\phi}$ can be expressed as

$$\mu_k(\mathbf{x} + \phi\hat{\phi}) - \mu_k(\mathbf{x}) = \int_{\mathbf{t} \in \mathbb{R}^n} \mathbf{1}_{h(\mathbf{t})=k} v(\mathbf{x} + \phi\hat{\phi} - \mathbf{t}) dt - \int_{\mathbf{t} \in \mathbb{R}^n} \mathbf{1}_{h(\mathbf{t})=k} v(\mathbf{x} - \mathbf{t}) dt \quad (34)$$

For simplicity, let $\phi = \phi\hat{\phi}$ for the moment. Then, according to the linearity of integration, we can put the two integrands under the same integral as

$$\mu_k(\mathbf{x} + \phi) - \mu_k(\mathbf{x}) = \int_{\mathbf{t} \in \mathbb{R}^n} (\mathbf{1}_{h(\mathbf{t})=k} v(\mathbf{x} + \phi - \mathbf{t}) - \mathbf{1}_{h(\mathbf{t})=k} v(\mathbf{x} - \mathbf{t})) dt \quad (35)$$

Observe that we can combine like terms $\mathbf{1}_{h(\mathbf{t})=k}$ shared by two parts of the integrands. Thus,

$$\mu_k(\mathbf{x} + \phi) - \mu_k(\mathbf{x}) = \int_{\mathbf{t} \in \mathbb{R}^n} \mathbf{1}_{h(\mathbf{t})=k} (v(\mathbf{x} + \phi - \mathbf{t}) - v(\mathbf{x} - \mathbf{t})) dt \quad (36)$$

Applying the symmetry of the vicinity function about axes, we can get

$$\mu_k(\mathbf{x} + \phi) - \mu_k(\mathbf{x}) = \int_{\mathbf{t} \in \mathbb{R}^n} \mathbf{1}_{h(\mathbf{t})=k} (v(\mathbf{t} - \mathbf{x} - \phi) - v(\mathbf{t} - \mathbf{x})) dt \quad (37)$$

Next, we shift the integral limit by $+\mathbf{x}$, the integrand becomes $\mathbf{1}_{h(\mathbf{t}+\mathbf{x})=k} (v(\mathbf{t} - \phi) - v(\mathbf{t}))$, and the interval of integration remains the same. To find the upper and lower bounds of $\mu_k(\mathbf{x} + \phi) - \mu_k(\mathbf{x})$, we want to find those for this integrand. Observe that $\mathbf{1}_{h(\mathbf{t}+\mathbf{x})=k}$ either takes value 0 or 1. Thus, letting $\mathbf{1}_{h(\mathbf{t}+\mathbf{x})=k} = 1$ if and only if $(v(\mathbf{t} - \phi) - v(\mathbf{t})) > 0$ will maximise the integrand, and letting $\mathbf{1}_{h(\mathbf{t}+\mathbf{x})=k} = 1$ if and only if $(v(\mathbf{t} - \phi) - v(\mathbf{t})) < 0$ will minimise the integrand. Formally,

$$\min(0, v(\mathbf{t} - \phi) - v(\mathbf{t})) \leq \mathbf{1}_{h(\mathbf{t}+\mathbf{x})=k} (v(\mathbf{t} - \phi) - v(\mathbf{t})) \leq \max(0, v(\mathbf{t} - \phi) - v(\mathbf{t})) \quad (38)$$

Substitute this inequality back into the integral gives rise to

$$|\mu_k(\mathbf{x} + \phi) - \mu_k(\mathbf{x})| \leq \int_{\mathbf{t} \in \mathbb{R}^n} \max(0, v(\mathbf{t} - \phi) - v(\mathbf{t})) dt \quad (39)$$

Now we add 1 minus 1 to the right-hand side. Note that integrating a probability density function v across the entire domain also equals 1.

$$\begin{aligned} \text{RHS} &= 1 + \int_{\mathbf{t} \in \mathbb{R}^n} \max(0, v(\mathbf{t} - \phi) - v(\mathbf{t})) dt - 1 \\ &= 1 + \int_{\mathbf{t} \in \mathbb{R}^n} (\max(0, v(\mathbf{t} - \phi) - v(\mathbf{t})) - v(\mathbf{t} - \phi)) dt \\ &= 1 + \int_{\mathbf{t} \in \mathbb{R}^n} \max(-v(\mathbf{t} - \phi), -v(\mathbf{t})) dt \\ &= 1 - \int_{\mathbf{t} \in \mathbb{R}^n} \min(v(\mathbf{t} - \phi), v(\mathbf{t})) dt \end{aligned} \quad (40)$$

Hence, the difference between $\mu_k(\mathbf{x} + \phi)$ and $\mu_k(\mathbf{x})$ is bounded by the complement of a vicinity from another vicinity shifted by ϕ .

Further, recall $\phi = \phi\hat{\phi}$, the integrand thus becomes $\min(v(\mathbf{t} - \phi\hat{\phi}), v(\mathbf{t}))$, and

$$\min(v(\mathbf{t} - \phi\hat{\phi}), v(\mathbf{t})) \leq \min\left(v(\mathbf{t}), \min_{\hat{\phi} \in \mathbb{S}^{n-1}} v(\mathbf{t} - \phi\hat{\phi})\right) \quad (41)$$

As a result, the problem of maximising the change in μ_k by a shifting magnitude ϕ is converted into the optimisation of finding the direction that results in the minimum overlap between the original vicinity function (v) and the vicinity function shifted by ϕ in that direction. The resulting upper bound can be expressed as Equation (13), which is re-displayed as follows.

$$\Delta_{\max} \mu_k(\phi) = \max_{\hat{\phi} \in \mathbb{S}^{n-1}} \mu_k(\mathbf{x} + \phi\hat{\phi}) - \mu_k(\mathbf{x}) = 1 - \min_{\hat{\phi} \in \mathbb{S}^{n-1}} \int_{\mathbf{t} \in \mathbb{R}^n} \min(v(\mathbf{t} - \phi\hat{\phi}), v(\mathbf{t})) dt \quad (42)$$

Similarly, the lower bound is the negative of the upper bound. Also, as long as two vicinities overlap, the change of μ_k is less than 1. \square

B.5 Extended Proof of Theorem 3.4

In the original proof of Theorem 3.4, a value $1/2 - \kappa$ is involved. Here, we explain how we get this value, and why it stands for the minimum required μ_k drop to allow an adversarial example.

Why $1/2 - \kappa$ marks the minimum μ_k change to have an adversarial example. For any consistent input \mathbf{x} , we have $\mu_{k^*}(\mathbf{x}) \geq 1 - \kappa$, where k^* is the major prediction in \mathbf{x} -vicinity. Consider \mathbf{x}' as a neighbour of \mathbf{x} . If $\mu_{k^*}(\mathbf{x}') > 1/2$, then we know that $h(\mathbf{x}') = k^*$ according to Theorem 3.2, i.e., \mathbf{x}' has the same prediction as \mathbf{x} does. In this way \mathbf{x}' can be possibly an adversarial example of \mathbf{x} only if $\mu_{k^*}(\mathbf{x}') \leq 1/2$. Thus, we get the minimum requirement of μ_k drop to allow an adversarial example to appear to be $1 - \kappa - 1/2 = 1/2 - \kappa$. \square

B.6 Proof of Corollary 3.5

Corollary 3.5. *There exists a finite real value such that for all inputs, the directional derivative value with respect to any arbitrary nonzero vector $\hat{\phi}$ (unit vector) does not exceed this finite value and does not fall below the negative of this value. Formally, \cdot denotes the dot product, and*

$$\exists b \in \mathbb{R}, \forall \mathbf{x} \in \mathbb{R}^n, \forall \hat{\phi} \in \mathbb{S}^{n-1}, -b \leq \nabla \mu_k(\mathbf{x}) \cdot \hat{\phi} \leq b. \quad (43)$$

Proof. The directional derivative of μ_k in the direction of $\hat{\phi} \in \mathbb{S}^{n-1}$ can be expressed as

$$\nabla \mu_k(\mathbf{x}) \cdot \hat{\phi} = \lim_{\delta \rightarrow 0} \frac{\mu_k(\mathbf{x} + \delta\hat{\phi}) - \mu_k(\mathbf{x})}{\delta} \quad (44)$$

According to Lemma 3.3, we can re-express the numerator such that the directional derivative is

$$\begin{aligned}
\nabla \mu_k(\mathbf{x}) \cdot \hat{\phi} &= \lim_{\delta \rightarrow 0} \frac{\int_{\mathbf{t} \in \mathbb{R}^n} \mathbf{1}_{h(\mathbf{t}+\mathbf{x})=k} \left(v(\mathbf{t} - \delta \hat{\phi}) - v(\mathbf{t}) \right) dt}{\delta} \\
&= \lim_{\delta \rightarrow 0} \int_{\mathbf{t} \in \mathbb{R}^n} \frac{\mathbf{1}_{h(\mathbf{t}+\mathbf{x})=k} \left(v(\mathbf{t} - \delta \hat{\phi}) - v(\mathbf{t}) \right)}{\delta} dt \\
&= \lim_{\delta \rightarrow 0} \int_{\mathbf{t} \in \mathbb{R}^n} \mathbf{1}_{h(\mathbf{t}+\mathbf{x})=k} \frac{v(\mathbf{t} - \delta \hat{\phi}) - v(\mathbf{t})}{\delta} dt \\
&= \int_{\mathbf{t} \in \mathbb{R}^n} \mathbf{1}_{h(\mathbf{t}+\mathbf{x})=k} \lim_{\delta \rightarrow 0} \frac{v(\mathbf{t} - \delta \hat{\phi}) - v(\mathbf{t})}{\delta} dt \\
&= - \int_{\mathbf{t} \in \mathbb{R}^n} \mathbf{1}_{h(\mathbf{t}+\mathbf{x})=k} \lim_{\delta \rightarrow 0} \frac{v(\mathbf{t}) - v(\mathbf{t} - \delta \hat{\phi})}{\delta} dt \tag{45} \\
&= - \int_{\mathbf{t} \in \mathbb{R}^n} \mathbf{1}_{h(\mathbf{t})=k} \lim_{\delta \rightarrow 0} \frac{v(\mathbf{t} - \mathbf{x}) - v(\mathbf{t} - \mathbf{x} - \delta \hat{\phi})}{\delta} dt \\
&= \int_{\mathbf{t} \in \mathbb{R}^n} \mathbf{1}_{h(-\mathbf{t}+\mathbf{x})=k} \lim_{\delta \rightarrow 0} \frac{v(-\mathbf{t}) - v(-\mathbf{t} - \delta \hat{\phi})}{\delta} dt \\
&= - \int_{\mathbf{t} \in \mathbb{R}^n} \mathbf{1}_{h(-\mathbf{t}+\mathbf{x})=k} \lim_{\delta \rightarrow 0} \frac{v(\mathbf{t} + \delta \hat{\phi}) - v(\mathbf{t})}{\delta} dt \\
&= - \int_{\mathbf{t} \in \mathbb{R}^n} \mathbf{1}_{h(\mathbf{t}+\mathbf{x})=k} \left(\nabla v(\mathbf{t}) \cdot \hat{\phi} \right) dt
\end{aligned}$$

The directional derivative of μ_k is maximised when the binary function takes 1 if and only if $\nabla v(\mathbf{t}) \cdot \hat{\phi} < 0$. Also, v is even in every dimension, and when $\nabla v(\mathbf{t}) \cdot \hat{\phi} < 0$, it is necessary that if $\nabla v(-\mathbf{t}) \cdot \hat{\phi} > 0$. Thus, the direction that maximises the directional derivative of μ_k can be expressed as

$$\phi^* = - \begin{bmatrix} \int_{\mathbf{t} \in \mathbb{R}^n} \left| \frac{\partial v(\mathbf{t})}{\partial t_1} \right| dt \\ \int_{\mathbf{t} \in \mathbb{R}^n} \left| \frac{\partial v(\mathbf{t})}{\partial t_2} \right| dt \\ \vdots \\ \int_{\mathbf{t} \in \mathbb{R}^n} \left| \frac{\partial v(\mathbf{t})}{\partial t_n} \right| dt \end{bmatrix} \tag{46}$$

Thus, the upper bound of the directional derivative of μ_k can be written as

$$\begin{aligned}
\nabla \mu_k(\mathbf{x}) \cdot \frac{\phi}{|\phi|} &\leq \frac{1}{2} \int_{\mathbf{t} \in \mathbb{R}^n} \left| \nabla v(\mathbf{t}) \cdot \frac{\phi^*}{|\phi^*|} \right| dt \\
&= \frac{1}{2 |\phi^*|} \int_{\mathbf{t} \in \mathbb{R}^n} \left| \sum_{i=1}^n \left(\frac{\partial v(\mathbf{t})}{\partial t_i} \int_{\boldsymbol{\tau} \in \mathbb{R}^n} \left| \frac{\partial v(\boldsymbol{\tau})}{\partial \tau_i} \right| d\boldsymbol{\tau} \right) \right| dt \tag{47}
\end{aligned}$$

Since function v is not a delta distribution PDF, this bound is always a finite number. \square

B.7 Example of Corollary 3.5

Example 1.

$$\begin{aligned}
&= \frac{1}{2 |\phi^*|} \int_{\mathbf{t} \in \mathbb{R}^n} \left| \sum_{i=1}^n \left(\frac{\partial v(\mathbf{t})}{\partial t_i} \int_{\boldsymbol{\tau} \in \mathbb{R}^n} \left| \frac{\partial v(\boldsymbol{\tau})}{\partial \tau_i} \right| d\boldsymbol{\tau} \right) \right| dt. \\
&= \frac{1}{2} \int_{\mathbf{t} \in \mathbb{R}^n} \left| \sum_{i=1}^n \left(\frac{\partial v(\mathbf{t})}{\partial t_i} \right) \right| dt. \tag{48} \\
&= \frac{1}{2} \int_{\mathbf{t} \in \mathbb{R}^n} \left| \left(\frac{\partial v(\mathbf{t})}{\partial t_0} \right) \right| dt. \\
&= \frac{1}{2} \int_{\mathbf{t} \in \mathbb{R}^n} \left| \left(\frac{\partial v(\mathbf{t})}{\partial t_0} \right) \right| dt.
\end{aligned}$$

Suppose we have input $\mathbf{x} \in \mathbb{R}$, $v : \mathbb{R} \rightarrow \mathbb{R}$, and v is a symmetric uniform distribution function. According to Corollary 3.5, the slope of μ_k in this example is within $\pm \frac{1}{2\epsilon}$. We can validate this value in Equation (49) using Leibniz's rule for differentiation under the integral sign.

$$\frac{d(\mathbf{1}_{h(\cdot)=k} * v)(\mathbf{x})}{d\mathbf{x}} = \frac{1}{2\epsilon} \frac{d}{d\mathbf{x}} \int_{-\epsilon}^{\epsilon} \mathbf{1}_{h(\mathbf{x}-t)=k} dt = \frac{1}{2\epsilon} (\mathbf{1}_{h(\mathbf{x}-\epsilon)=k} - \mathbf{1}_{h(\mathbf{x}+\epsilon)=k}) \quad (49)$$

The intuition of this example is that when a vicinity shifts from a region where all samples are labelled with one class to a region where all samples are labelled with another class, the slope reaches its maximum.

B.8 Proof of Corollary 3.6

Corollary 3.6. *If h^* is optimal for the probabilistic robustness with respect to an L^∞ -vicinity on a given distribution, then the vicinity size for the deterministically robust region around each consistent is $\epsilon \in \left(1 - (2\kappa)^{\frac{1}{n}}\right)$.*

Proof. An L^∞ norm looks like a n -dimensional cube. A two-dimensional illustration is given in Figure 1d. Generally, the vicinity function can be expressed as

$$v(\mathbf{x}) = \begin{cases} (2\epsilon)^{-n} & \text{for } \|\mathbf{x}\|_\infty \leq \epsilon \\ 0 & \text{otherwise.} \end{cases} \quad (50)$$

According to Theorem 3.4, if we would like to find the closest adversarial example to a consistent input at the upper right centre of the yellow vicinity, we first need to find a shift magnitude that causes as large as a μ_{k^*} drop by $1/2 - \kappa$. Suppose this drop goes further and the closest (probabilistically) consistent input is found. In this way, the minimum μ_k drop from consistent point \mathbf{x} and its (probabilistically) consistent adversarial example is $1 - \kappa - \kappa = 1 - 2\kappa$. In an L^∞ vicinity scenario, the magnitude of shift can be solved based on Equation (51).

$$1 - 2\kappa = 1 - \min_{\hat{\phi}} \prod_{i=1}^n \max(0, 2\epsilon - \phi \hat{\phi}_i) \quad (51)$$

where $\hat{\phi}_i$ is each element of the unit directional vector $\hat{\phi}$. Since the shift is within the vicinity, we can write $\max(0, 2\epsilon - \phi \hat{\phi}_i)$ as $2\epsilon - \phi \hat{\phi}_i$. For two identical n -dimensional cubes, the fastest way to reduce the overlap is to move one of them in the diagonal direction, such that $\hat{\phi}_i = \frac{1}{\sqrt{n}}$. Then the overlap volume becomes $(2\epsilon - \phi/\sqrt{n})^n$. Then, Equation (51) can be simplified as $2\kappa(2\epsilon)^n = (2\epsilon - |\phi|/\sqrt{n})^n$. Solving this equation, we get $\epsilon(1 - (2\kappa)^{\frac{1}{n}})$. This $\epsilon(1 - (2\kappa)^{\frac{1}{n}})$ also serves as the vicinity size for deterministic robustness at this input. As n grows, this vicinity size decreases. As κ grows, this vicinity size may grow. \square

We may visualise the effect of Corollary 3.6 using a two-dimensional example in Figure 1d. Let \mathbf{x}, \mathbf{x}' be two probabilistically consistent inputs and $h(\mathbf{x}) \neq h(\mathbf{x}')$. Suppose \mathbf{x}' is the nearest adversarial example (of \mathbf{x}) that achieves its own probabilistic consistency. Thus, \mathbf{x}' must be in the direction (that travels away from \mathbf{x}) that fastest decreases the vicinity overlap (suggested by Theorem 3.4), *i.e.*, the diagonal (suggested by Corollary 3.6). Consequently, the nearest adversarial example of \mathbf{x} would occur on the midpoint between \mathbf{x} and \mathbf{x}' . The shift from \mathbf{x} to \mathbf{x}' is $-2(\phi_1 \hat{\mathbf{x}}_1 + \phi_2 \hat{\mathbf{x}}_2)$. Each triangle accounts for a κ portion of the original vicinity volume, and the vicinity overlap is 2κ . The dashed box $\nabla^{\downarrow\kappa}$ has side length $2\phi_i$. Thus, solving $(2\epsilon - 2\phi_i)^2 = 2\kappa(2\epsilon)^2$, we get $\nabla^{\downarrow\kappa}$ has vicinity size $\phi_1 = \phi_2 = (1 - \sqrt{2\kappa})\epsilon$. Although Corollary 3.6 specifically captures the L^∞ scenario, we remark that other vicinity types can be analysed similarly. First, the direction that fastest decreases the vicinity overlap needs to be found. Then, the distance between the input and its nearest adversarial example can be measured.

B.9 Extended Proof of Theorem 3.8

Theorem 3.8. *The upper bound of probabilistic robust accuracy monotonically increases as κ grows. Further, for all tolerance levels κ , the upper bound of probabilistic robust accuracy lies between the*

upper bound of deterministic robust accuracy and vanilla accuracy. Formally,

$$\begin{aligned} \forall \kappa_1, \kappa_2. \quad (\kappa_1 < \kappa_2) &\rightarrow \min_h \Upsilon_{\text{prob}}^+(D, h, \mathcal{V}, \kappa_1) \leq \min_h \Upsilon_{\text{prob}}^+(D, h, \mathcal{V}, \kappa_2) \\ \forall \kappa. \quad \min_h \Upsilon_{\text{rob}}^+(D, h, \mathbb{V}) &\leq \min_h \Upsilon_{\text{prob}}^+(D, h, \mathcal{V}, \kappa) \leq \min_h \Upsilon_{\text{acc}}^+(D, h) \end{aligned} \quad (52)$$

Proof. Let $\kappa_1 < \kappa_2$, such that for κ_1, κ_2 , we have their corresponding error as the following equation.

$$\begin{aligned} e(\mathbf{x}, \kappa_1) - e(\mathbf{x}, \kappa_2) &= 1 - u \left(\kappa_1 - 1 + \sum_{k=0}^{K-1} \mu_k(\mathbf{x}) \mathbf{1}_{h(\mathbf{x})=k} \right) \left(\sum_{y=0}^{K-1} P(y = y | \mathbf{x} = \mathbf{x}) \mathbf{1}_{h(\mathbf{x})=y} \right) \\ &\quad - 1 + u \left(\kappa_2 - 1 + \sum_{k=0}^{K-1} \mu_k(\mathbf{x}) \mathbf{1}_{h(\mathbf{x})=k} \right) \left(\sum_{y=0}^{K-1} P(y = y | \mathbf{x} = \mathbf{x}) \mathbf{1}_{h(\mathbf{x})=y} \right) \end{aligned} \quad (53)$$

Note the cancelled 1 and the like terms in the rest of the terms, we can further write the above equation as Equation (54), where we let $T = \sum_{k=0}^{K-1} \mu_k(\mathbf{x}) \mathbf{1}_{h(\mathbf{x})=k}$. Note that T is just a temporary substituting variable, and we do not mean to use it to denote any particular quantity.

$$e(\mathbf{x}, \kappa_1) - e(\mathbf{x}, \kappa_2) = \left(\sum_{y=0}^{K-1} P(y = y | \mathbf{x} = \mathbf{x}) \mathbf{1}_{h(\mathbf{x})=y} \right) (u(\kappa_2 - 1 + T) - u(\kappa_1 - 1 + T)) \quad (54)$$

Now that $e(\mathbf{x}, \kappa_1) - e(\mathbf{x}, \kappa_2)$ is a product of two expressions. Since the posterior probability $P(y | \mathbf{x} = \mathbf{x})$ is non-negative, the sum of posteriors, *i.e.*, the former expression, is non-negative. Thus, the sign of $e(\mathbf{x}, \kappa_1) - e(\mathbf{x}, \kappa_2)$ is the same as $u(\kappa_2 - 1 + T) - u(\kappa_1 - 1 + T)$.

Since $\kappa_1 < \kappa_2$, we get $\kappa_1 - 1 + T < \kappa_2 - 1 + T$. Further, the unit step function (u) is monotonically increasing, we get $e(\mathbf{x}, \kappa_1) - e(\mathbf{x}, \kappa_2) \geq 0$. A more stringent κ (*i.e.*, smaller) leads to a lower or equal upper bound of probabilistic robust accuracy. For deterministic robust accuracy, $\kappa = 0$, which is the least value. Thus, deterministic robust accuracy has a lower upper bound than probabilistic robust accuracy does.

On the other hand, from the intuition of error combination in Section 3.1, we get that the combined error is $e(\mathbf{x}) = 1 - (1 - e_{\text{cns}}(\mathbf{x}, h; \mathcal{V}, \kappa))(1 - e_{\text{cor}}(\mathbf{x}, h; P(y | \mathbf{x})))$. As $0 \leq e_{\text{cns}} \leq 1$, we get $e(\mathbf{x}) \geq 1 - (1 - e_{\text{cor}}(\mathbf{x}, h; \kappa))$. Note that the expectation of $e_{\text{cor}}(\mathbf{x}, h; \kappa)$ is the error in accuracy. Thus, $\Upsilon_{\text{prob}}^+(D, h, \mathcal{V}, \kappa) \leq \Upsilon_{\text{acc}}^+(D, h)$. So it is with their upper bounds. \square

C Additional Experiments and Plots

In this section, we present the results of some additional experiments in which we investigate the effect of the upper bounds and decision rules.

C.1 Setup Details

The experiments are conducted with four data sets: two synthetic ones (*i.e.*, Moons and Chan [6], whose distributions are illustrated in Figure 5) and two standard benchmarks (*i.e.*, FashionMNIST [46] and CIFAR-10 [17]). Moons is used for binary classification with two-dimensional features, where each class’s distribution is described analytically with specific likelihood equations, and uses a three-layer Multi-Layer Perceptron (MLP) neural network for classification. The Chan data set, also for binary classification with two-dimensional features, differs in that it does not follow a standard PDF pattern, requiring kernel density estimation (KDE) for non-parametric PDF estimation, and also uses the three-layer MLP. FashionMNIST, a collection of fashion item images, involves a 10-class classification task with 784-dimensional inputs (28×28 pixel grayscale images). Each class has an equal prior probability, and their conditional distributions are estimated non-parametrically using KDE. CIFAR-10 uses images with a resolution of 32×32 pixels. Similar to FashionMNIST, it has a balanced class distribution and is estimated using KDE. We use a seven-layer convolutional neural network (CNN-7) [38] as the classifier of both FashionMNIST and CIFAR-10. We adopt a direct approach [15, 52] to compute the original Bayes error and deterministic Bayes error of both real-world data sets [15].

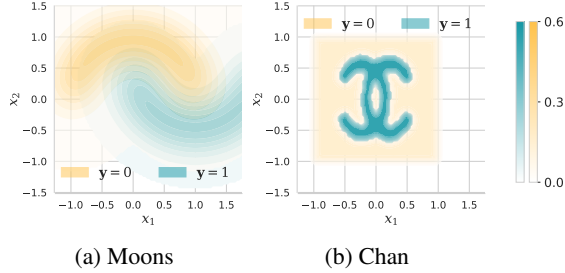


Figure 5: This figure illustrates the conditional distribution for (a) Moons [31] and (b) Chan [6].

Table 2: Probabilistic robustness of classifiers before voting. The value in parentheses represents the 95% confidence level range when we repeat the same program 100 times.

	Moons	Chan	FashionMNIST	CIFAR-10
DA [39]	85.35 ($\pm 0.020\%$)	67.96 ($\pm 0.023\%$)	84.12 ($\pm 0.000\%$)	76.07 ($\pm 0.021\%$)
RS [8]	84.76 ($\pm 0.020\%$)	64.67 ($\pm 0.002\%$)	86.29 ($\pm 0.012\%$)	87.98 ($\pm 0.038\%$)
CVaR [35]	85.52 ($\pm 0.017\%$)	69.46 ($\pm 0.034\%$)	88.50 ($\pm 0.028\%$)	90.63 ($\pm 0.007\%$)

C.2 Statistical Significance of Experiments

We include confidence levels of performances of trained classifiers Table 2 and their corresponding voting classifiers Table 3. We provide their statistical significance to support the claims addressed by our research questions.

Besides, the train/validation splits follow that of the data set [46, 17] if there is already a split guideline. For Moons [31] and Chan [6], we follow the setup in [52].

C.3 How does the sample size affect the voting effectiveness?

The voting process of classifier h^\dagger can be viewed as a process of taking the expected value of prediction within the vicinity. The law of large numbers is key to understanding the relationship between sample size and expectation. It states that as the sample size increases, the sample mean converges to the true expectation.

To validate this effect empirically, as well as to verify a suitable sample size at which the performance of the voting classifiers can be properly represented, we gradually increase the sample size from 10 to 10,000. As demonstrated in Figure 6, at a small sample size, the voting may result in an increase or a decrease in the performance. However, as the sample exceeds 100, its positive impact on the probabilistic robust accuracy becomes more noticeable. Thus, our empirical intuition matches Theorem 3.2.

Table 3: Probabilistic robustness of classifiers after voting. The value in parentheses represents the 95% confidence level range when we repeat the same program 100 times.

	Moons	Chan	FashionMNIST	CIFAR-10
DA [39]	85.60 ($\pm 0.034\%$)	68.86 ($\pm 0.006\%$)	87.48 ($\pm 0.022\%$)	81.38 ($\pm 0.026\%$)
RS [8]	85.18 ($\pm 0.005\%$)	66.77 ($\pm 0.008\%$)	88.13 ($\pm 0.037\%$)	88.95 ($\pm 0.008\%$)
CVaR [35]	85.66 ($\pm 0.011\%$)	70.05 ($\pm 0.005\%$)	91.07 ($\pm 0.028\%$)	90.77 ($\pm 0.011\%$)

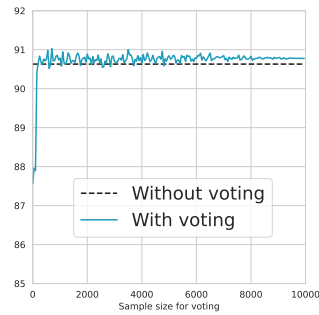


Figure 6: Probabilistic robust accuracy of voting classifier (CVaR [35]) tested on CIFAR-10 [17]) as the sample size grows to 10,000.