
Fast Adversarial Robustness Certification of Nearest Prototype Classifiers for Arbitrary Seminorms

Sascha Saralajew*

Dr. Ing. h.c. F. Porsche AG
Weissach, Germany
sascha.saralajew@porsche.de

Lars Holdijk*

University of Amsterdam
Amsterdam, Netherlands
larsholdijk@gmail.com

Thomas Villmann

UAS Mittweida
Mittweida, Germany,
villmann@hs-mittweida.de

Abstract

Methods for adversarial robustness certification aim to provide an upper bound on the test error of a classifier under adversarial manipulation of its input. Current certification methods are computationally expensive and limited to attacks that optimize the manipulation with respect to a norm. We overcome these limitations by investigating the robustness properties of Nearest Prototype Classifiers (NPCs) like learning vector quantization and large margin nearest neighbor. For this purpose, we study the hypothesis margin. We prove that if NPCs use a dissimilarity measure induced by a seminorm, the hypothesis margin is a tight lower bound on the size of adversarial attacks and can be calculated in constant time—this provides the first adversarial robustness certificate calculable in reasonable time. Finally, we show that each NPC trained by a triplet loss maximizes the hypothesis margin and is therefore optimized for adversarial robustness. In the presented evaluation, we demonstrate that NPCs optimized for adversarial robustness are competitive with state-of-the-art methods and set a new benchmark with respect to computational complexity for robustness certification.

1 Introduction

Adversarial robustness of a classifier describes its stability in classification under adversarial manipulations of the input. The adversarial setting has been studied extensively in numerous settings [1–3] but mainly found footing after the seminal paper by Szegedy et al. [4] that formalized the problem of *adversarial examples*. Since that, a wide line of research has sprung concerning both the construction of adversarial attacks [5–7] and the heuristic defense against them [8–10]. Unfortunately, while some progress has been made [11–13], most proposed defenses have been shown to be breakable by more advanced attacks [14–16] so that the adversarial robustness problem is far from being solved. With the heuristic defenses at the losing side of the metaphorical arms race, provable robustness guarantees for classifiers provide a welcome alternative [11, 12, 17]. Robustness guarantees aim to provide the so-called *robust test error* (or an upper bound) of a classifier under adversarial attacks and are therefore not dependent on the current state of adversarial attacks. However, current methods for determining the robust test error are computationally expensive [12, 17, 18] and are often limited to L^p -norm evaluations [11, 13, 18]. As an adverse effect, we cannot use them for the rejection of adversarial examples regarding an arbitrary seminorm without a huge computational overhead.

* Authors contributed equally.

To tackle these limitations, we extend the study of robustness guarantees beyond Neural Networks (NNs), on which most work has focused, by investigating *Nearest Prototype Classifiers* (NPCs) [19–22]. For that, given a data space \mathcal{X} , an NPC is defined by a set \mathcal{W} of *prototypes* selected from the data space \mathcal{X} and a dissimilarity measure $d : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_{\geq 0}$. Each prototype $\mathbf{w} \in \mathcal{W}$ has a predefined and fixed class label $c(\mathbf{w}) \in \mathcal{C} = \{1, 2, \dots, N_c\}$. The class assignment $c_{\mathcal{W}}^*(\mathbf{x}) \in \mathcal{C}$ of a data point $\mathbf{x} \in \mathcal{X}$ is determined by the class label of the closest prototype \mathbf{w}^* (1-nearest neighbor rule):

$$c_{\mathcal{W}}^*(\mathbf{x}) = c(\mathbf{w}^*) \text{ with } \mathbf{w}^* = \arg \min_{\mathbf{w} \in \mathcal{W}} \{d(\mathbf{x}, \mathbf{w})\}. \quad (1)$$

Due to the fundamental interpretability of the prototypes and dissimilarity measure used, NPCs are considered to be among the most interpretable machine learning models. This makes NPCs a preferred choice in the medical field, where the interpretability of models is a requirement for clinical trials [23, 24]. As a result, prototype-based principles have been adopted in a number of deep learning fields. Amongst the most notable of those is few-shot learning [25]. In addition to this, an empirical study has shown that NPCs are robust against adversarial attacks [26]. In summary—with the call for interpretable machine learning models increasing, the NPC principles adopted in deep learning, and the promising empirical robustness results—NPCs provide a worthwhile avenue for studying guaranteed adversarial robustness.

Contributions We analyze the adversarial robustness properties of NPCs in terms of the hypothesis margin. First, we show that if the dissimilarity measure $d(\mathbf{x}, \mathbf{w})$ is induced by a seminorm $\|\cdot\|$ (i. e., $d(\mathbf{x}, \mathbf{w}) = \|\mathbf{x} - \mathbf{w}\|$), the hypothesis margin regarding this seminorm can be computed in constant time during the inference—this even holds in the case of an uncountable set of prototypes. Second, we prove that this margin is a tight lower bound on the magnitude of an adversarial attack measured by the same seminorm—to the best of our knowledge, this presents the first robustness guarantee that holds for an arbitrary seminorm. Third, using this result, we show that every NPC that classifies by a seminorm and is trained by minimizing a triplet loss is inherently optimizing the adversarial robustness with respect to this seminorm. In an experimental section, we present how these results apply to different NPCs and that a violation of the assumptions (seminorm and triplet loss) does not necessarily lead to adversarially robust methods. The experimental results highlight that the derived robustness certificate is comparable with other methods in terms of guaranteed robustness but outperforms them all when computation speed is considered.

The following section discusses related work. After that, we give a brief introduction to NPCs and define the models we use in the evaluation. This section is followed by the main part where we define the hypothesis margin and investigate the relation to adversarial robustness. The subsequent experimental section shows how to apply these results and compares the derived robustness certificate to other methods. Finally, we finish with a discussion and an outlook of the presented results.

2 Related work

Besides general theoretical work about adversarial robustness [27–30], the defense investigations can be grouped into three areas: *robustification*, the research to improve the adversarial robustness of models [8, 10, 11, 31–33]; *verification* (complete or exact methods), the analysis of how to compute the robustness guarantees exact [17, 34, 35]; *certification* (incomplete methods), the study of fast calculable bounds for the robustness guarantees [36–38]. It is important to distinguish between verification methods and certification approaches. Naturally, verification approaches [17] have a combinatorial time complexity even though there are attempts to improve the computational complexity [39]. In contrast, certification methods [12, 40] try to return bounds for the robustness guarantees in polynomial time. Comparing verification and certification results only in terms of the returned robustness guarantees ignores the aspect of time complexity of the methods and is therefore not a fair comparison. We consider this thought carefully throughout the evaluation.

The presented work focuses on both the study of fast calculable certificates for NPCs and the use of the certificate to robustify the models. Currently, the most successful approach to robustify models is adversarial training [4, 41, 42]. However, being dependent on an adversarial attack used during training, adversarial training fails to provide a robustness certificate and increases the training time. Other empirical robustification approaches [8, 9, 43, 44] struggle with similar problems or fail to withstand stronger attacks [5, 7, 14–16, 41, 45–47]. It is therefore not surprising that several modern approaches strive—as the work presented in this article—towards combining robustification and

certification. Wong et al. [12, 40] studied an approach for NNs based on *convex outer adversarial polytopes*. Their method extends to arbitrary norms but is still computationally expensive and infeasibly complex to compute for deep NNs. Randomized Smoothing yields another approach mainly investigated for NNs despite being applicable for *arbitrary* classifiers [18, 48, 49]. This approach fails to scale to arbitrary norms though and introduces a heavy computational overhead through sampling. Besides, some work has also been done to study certification and verification for other classification approaches like decision trees [50, 51] and support vector machines [2, 36, 52] or to extend robustness certificates to an arbitrary L^p -norm knowing the results for a few L^p -norms [13].

Based on empirical observations, Saralajew et al. [26] discussed the relation between the margin maximization properties of Generalized Learning Vector Quantization (GLVQ) [53] and its adversarial robustness *without* providing a mathematical proof. The presented results rely on the hypothesis margin maximization properties of Learning Vector Quantization (LVQ) [54, 55], as originally studied by Crammer et al. [56]. There, the hypothesis margin was used to derive a generalization bound for LVQ with the Euclidean norm. Consequently, these results are not applicable to seminorms or an uncountable set of prototypes in arbitrary NPCs. Similarly, the work of Wang et al. [57] about the adversarial robustness of k-nearest neighbors cannot be scaled to an arbitrary NPC even though the results hold for an arbitrary norm. Wang et al. [58] presented a result that is similar to the hypothesis margin but limited to L^1 -, L^2 -, and L^∞ -norms for the attack and L^2 -norms for the classifier metric. However, their methods are orders of magnitude slower. Besides that, Yang et al. [59] proposed a generic defense by preprocessing the dataset before training a k-nearest neighbors classifier with adversarial pruning. In general, the method can be used to robustify NPCs, but for more than two classes it is time-consuming and it cannot be used to robustify seminorm-based NPCs. Brinkrolf et al. [60] studied Generalized Matrix LVQ (GMLVQ) [61] with reject options and derived an adversarial perturbation bound needed to fool the classifier. Compared to our work, the method requires the training of the classifier with reject options and does not provide a general framework for the evaluation of adversarial robustness and the robustification of NPCs.

3 Nearest prototype classifiers

As already mentioned in the introduction, an NPC consists of two main building blocks: a set \mathcal{W} of prototypes and a dissimilarity measure d . The prototypes $\mathbf{w} \in \mathcal{W}$ are elements of the data space \mathcal{X} and have a predefined class label $c(\mathbf{w}) \in \mathcal{C}$. Given a prototype \mathbf{w} and an input sample \mathbf{x} , we compute the “distance” between these two elements by the dissimilarity d . Based on this distance, we assign the class of the closest prototype \mathbf{w}^* to a given input, see Equation (1). The closest decision boundary to a given input \mathbf{x} is *implicitly* defined by the closest prototype \mathbf{w}^* and the closest prototype \mathbf{w}_* with a *different* class label than \mathbf{w}^* based on the given dissimilarity d .

In an NPC, the determination of the closest prototype of each class is considered as the *inference step*. During training, an NPC updates the prototypes and the maybe trainable dissimilarity measure. The training can be realized by heuristic methods and loss-based optimization and is an important difference between different realizations. For example, some NPCs optimize the selection of prototypes out of a given set of labeled data points (e. g., 1-nearest neighbor) or optimize the prototypes as free parameters (e. g., LVQ). Following the empirical observation of Saralajew et al. [26], we focus our analysis on the family of LVQ algorithms. These methods not only train the dissimilarity measure but the prototypes as well. In particular, compared to 1-nearest neighbor methods, LVQ does not only optimize the selection of prototypes out of a given training dataset but instead considers the prototypes as fully trainable parameters. We refer to the articles of Biehl et al. [22] and Nova and Estévez [21] for an overview and an in-depth introduction to LVQ and NPCs.

Generalized learning vector quantization GLVQ provides an LVQ version that can be trained by an empirical risk minimization and satisfies the convergence condition [53]. Because the method is trained by a gradient-based approach, the dissimilarity measure has to be differentiable almost everywhere. Except for this condition, the dissimilarity measure can be chosen freely (e. g., squared Euclidean distance). Given a training dataset \mathcal{T} of labeled inputs $(\mathbf{x}, c(\mathbf{x}))$, we fix the number of prototypes per class, initialize the prototypes, and optimize the prototypes by minimizing the following averaged loss function:

$$\frac{1}{\#\mathcal{T}} \sum_{(\mathbf{x}, c(\mathbf{x})) \in \mathcal{T}} \mu(\mathbf{x}, c(\mathbf{x})) = \frac{1}{\#\mathcal{T}} \sum_{(\mathbf{x}, c(\mathbf{x})) \in \mathcal{T}} \frac{d(\mathbf{x}, \mathbf{w}^+) - d(\mathbf{x}, \mathbf{w}^-)}{d(\mathbf{x}, \mathbf{w}^+) + d(\mathbf{x}, \mathbf{w}^-)}, \quad (2)$$

where \mathbf{w}^+ is the closest prototype to \mathbf{x} of the correct class $c(\mathbf{x})$ with respect to d and \mathbf{w}^- is the closest prototype to \mathbf{x} of an incorrect class. The expression $\mu(\mathbf{x}, c(\mathbf{x})) \in [-1, 1]$ is called the *relative distance difference* and returns negative values if and only if \mathbf{x} is correctly classified.

Generalized tangent learning vector quantization The Generalized Tangent LVQ (GTLVQ) algorithm is a version of GLVQ where a tangent distance is used [62]. The prototypes in GTLVQ are defined as elements of affine subspaces of the data space $\mathcal{X} = \mathbb{R}^n$. In other words, we can consider GTLVQ as an NPCs with infinitely many prototypes, approximating variations within the classes. The dissimilarity of a given data point \mathbf{x} to the k -th affine subspace is measured in terms of the smallest Euclidean distance d_E :

$$\min \{d_E(\mathbf{x}, \mathbf{t}_k + \mathbf{B}_k \boldsymbol{\theta}) \mid \boldsymbol{\theta} \in \mathbb{R}^m\} = d_E(\mathbf{x}, \mathbf{t}_k + \mathbf{B}_k \mathbf{B}_k^T (\mathbf{x} - \mathbf{t}_k)), \quad (3)$$

where $\mathbf{B}_k \in \mathbb{R}^{n \times m}$ is an m -dimensional orthonormal basis, $\mathbf{t}_k \in \mathbb{R}^n$ is a translation vector, and $\mathbf{t}_k + \mathbf{B}_k \mathbf{B}_k^T (\mathbf{x} - \mathbf{t}_k)$ determines the closest prototype at the k -th affine subspace. Given the number of affine subspaces and the dimension m , we optimize the bases \mathbf{B}_k and translations \mathbf{t}_k by minimizing Equation (2) with Equation (3) as the dissimilarity measure during training.

Robust soft learning vector quantization Robust Soft LVQ (RSLVQ) is a probabilistic version of LVQ where the prototypes are assumed as centers in a Gaussian mixture model [63].² The posterior probability $P(l \mid \mathbf{x})$ that an input \mathbf{x} belongs to a certain class $l \in \mathcal{C}$ is computed by

$$P(l \mid \mathbf{x}) = \frac{\sum_{\mathbf{w}:c(\mathbf{w})=l} \exp(-d_E^2(\mathbf{x}, \mathbf{w}))}{\sum_{\mathbf{w}} \exp(-d_E^2(\mathbf{x}, \mathbf{w}))}. \quad (4)$$

The prototypes are determined by minimizing the cross-entropy loss between the predicted probability vector $\mathbf{p}(\mathbf{x}) = (P(1 \mid \mathbf{x}), \dots, P(N_c \mid \mathbf{x}))^T$ and the true probability vector (one-hot encoding) of a labeled input $(\mathbf{x}, c(\mathbf{x}))$. Considering Equation (4), we observe that the calculation of the probabilities follows a softmax squashing. Together with the cross-entropy loss, this makes the RSLVQ algorithm highly similar to the design and training of classification layers in NNs.

4 Hypothesis margin maximization and adversarial robustness certification

In this section, we define the hypothesis margin for NPCs and derive some of its properties. For instance, we show that the hypothesis margin can be easily computed, lower bounds adversarial perturbations, and can be optimized during training. We refer the reader to Section A of the supplementary material for a visualization of the defined concepts and their properties in \mathbb{R}^2 .

4.1 Definition and calculation of the hypothesis margin

Margins are a common tool to measure the confidence of a classifier's decision. For example, the sample margin—the distance of a sample to the closest decision boundary—is used in the optimization of support vector machines. In the context of NPCs, the sample margin is defined as follows.

Definition 1 (sample margin). Given a set \mathcal{W} of prototypes and a dissimilarity d . The *sample margin* of \mathcal{W} with respect to a set \mathcal{S} of inputs is the maximum radius r such that the following condition holds: If we define a ball with radius r induced by d around each sample \mathbf{x} of \mathcal{S} , each point within a ball has the same assigned class label as the center \mathbf{x} of the ball. In symbols, we write $\text{margin}_s(\mathcal{S}, \mathcal{W})$.

As a decision boundary is implicitly defined by two prototypes of different classes, this margin definition is cumbersome for NPCs. In particular, given two prototypes \mathbf{w}_i and \mathbf{w}_j of different classes, a decision boundary is defined by those elements $\mathbf{x} \in \mathcal{X}$ for which $d(\mathbf{x}, \mathbf{w}_i) = d(\mathbf{x}, \mathbf{w}_j)$. Therefore, depending on the dissimilarity d , the calculation of the decision boundary can be difficult. The *hypothesis margin* offers an alternative and more suitable concept for NPCs.

²Compared to the other NPC methods, RSLVQ is a *probabilistic* version of a prototype-based classifier and *not* an NPC according to our definition. However, for simplicity, we consider RSLVQ as NPC in the following comparisons because it uses prototypes with fixed class labels and a dissimilarity. But note that it classifies and trains according to a probabilistic approach.

Definition 2 (hypothesis margin). Given a set \mathcal{W} of prototypes and a dissimilarity d . The *hypothesis margin* of \mathcal{W} with respect to a set \mathcal{S} of inputs is the maximum radius r such that the following condition holds: If we define a ball with radius r induced by d around each prototype, every change in the position of the prototypes within its ball does not change the class labels assigned to the inputs of \mathcal{S} . In symbols, we write $\text{margin}_h(\mathcal{S}, \mathcal{W})$.

We assume in these definitions—and in general—that the class label assignments to \mathcal{S} by \mathcal{W} are *unambiguous* and well-defined. This means that input samples do not lie on a decision boundary as otherwise the class assignments would be ill-defined. Additionally, a ball is always assumed as an *open* set. With the class assignments considered as unambiguous, this is not a limitation.

With the next theorem, we provide a formula to compute the hypothesis margin of a sample using only one additional floating-point operation after the inference of the NPC. This theorem is valid for seminorms and, thus, holds for all NPCs where the dissimilarity is induced by a seminorm.

Definition 3 (seminorm). Given a vector space \mathcal{V} over a field \mathcal{F} of the real or complex numbers, a *seminorm* $\|\cdot\|$ is a function $\|\cdot\| : \mathcal{V} \rightarrow \mathbb{R}$ that satisfies the following conditions for all $\mathbf{x}, \mathbf{y} \in \mathcal{V}$ and $\alpha \in \mathcal{F}$: $\|\mathbf{x}\| \geq 0$ (nonnegativity); $\|\alpha\mathbf{x}\| = |\alpha| \|\mathbf{x}\|$ (absolute homogeneity); $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ (triangle inequality).

Theorem 1. Let the data space \mathcal{X} be a vector space over a field of the real or complex numbers, $d(\mathbf{x}, \mathbf{w}) = \|\mathbf{x} - \mathbf{w}\|$ be a dissimilarity induced by a seminorm $\|\cdot\|$, and $\mathbf{x} \in \mathcal{X}$ be an input. Then, the hypothesis margin of the set \mathcal{W} of prototypes with respect to \mathbf{x} can be computed by

$$\text{margin}_h(\{\mathbf{x}\}, \mathcal{W}) = \frac{1}{2} (\|\mathbf{x} - \mathbf{w}_*\| - \|\mathbf{x} - \mathbf{w}^*\|), \quad (5)$$

where \mathbf{w}^* denotes the closest prototype to \mathbf{x} and \mathbf{w}_* denotes the closest prototype to \mathbf{x} with a different class label than the class label of \mathbf{w}^* .

A proof of the theorem can be found in the supplementary material Section B. Based on this theorem, we have a surprisingly simple rule to compute the hypothesis margin after the inference with only one additional floating-point operation and, therefore, in constant time.

4.2 Hypothesis margin lower bounds the adversarial perturbation

Given an input sample $\mathbf{x} \in \mathcal{X}$ with the class label $c(\mathbf{x})$, an *adversarial example* $\tilde{\mathbf{x}}$ of the sample \mathbf{x} is defined by an *adversarial perturbation* δ of \mathbf{x} such that $\tilde{\mathbf{x}} = \mathbf{x} + \delta$ is a point on the decision boundary or in the classification region of a different class than $c(\mathbf{x})$. Frequently, the computation of an adversarial perturbation is treated as an optimization problem by searching for a perturbation with minimum magnitude:

$$\min_{\delta} \|\delta\| \text{ such that } c_{\mathcal{W}}^*(\tilde{\mathbf{x}}) \neq c(\mathbf{x}) \text{ and } \tilde{\mathbf{x}} = \mathbf{x} + \delta \in \mathcal{X}. \quad (6)$$

We should note that the magnitude of the adversarial perturbation is measured in terms of a seminorm $\|\cdot\|$ and that the adversarial example has to be an element of the input space \mathcal{X} . Sometimes additional conditions are placed upon the adversarial example (e. g., it has to be a sample of a certain class, an element of a subset of \mathcal{X} , etc.). Frequently, the optimization task of Equation (6) is intractable and, hence, the optimal perturbation can only be approximated. A procedure that generates adversarial examples is called an *adversarial attack*. An adversarial attack that only creates adversarial examples with a maximum perturbation $\|\delta\|$ less than or equal a given bound $\epsilon > 0$ is called *ϵ -limited adversarial attack*.

The next theorem and corollary provide a statement about the relation between the sample and the hypothesis margin—see Section C of the supplementary material for the proofs.

Theorem 2. Let the data space \mathcal{X} be a vector space over a field of the real or complex numbers, $d(\mathbf{x}, \mathbf{w}) = \|\mathbf{x} - \mathbf{w}\|$ be a dissimilarity induced by a seminorm $\|\cdot\|$, and \mathcal{S} be a set of inputs. Then, the hypothesis margin of \mathcal{W} with respect to \mathcal{S} yields a lower bound on the sample margin of \mathcal{W} with respect to \mathcal{S} :

$$\text{margin}_h(\mathcal{S}, \mathcal{W}) \leq \text{margin}_s(\mathcal{S}, \mathcal{W}). \quad (7)$$

Corollary 1. Given a labeled data point $(\mathbf{x}, c(\mathbf{x}))$ that is correctly classified by an NPC according to Theorem 2 and a corresponding adversarial perturbation δ that changes the assigned class label.

Then, the following inequality is true and has tight bounds (existence of data points for which the equality is true):

$$\text{margin}_h(\{\mathbf{x}\}, \mathcal{W}) \leq \text{margin}_s(\{\mathbf{x}\}, \mathcal{W}) \leq \|\delta\|. \quad (8)$$

A direct result of this corollary is that we can use the hypothesis margin to reject adversarial examples with perfect recall during inference. Similar to Wong and Kolter [40, Corollary 2]: If a sample \mathbf{x} has a hypothesis margin greater than a certain threshold ϵ , then there exists no *original* sample that can be perturbed with an ϵ -limited adversarial attack such that its classification changes. Thus, \mathbf{x} is certified to be not an adversarial example. In Section 5, we investigate the related false rejection rate.

4.3 Hypothesis margin maximization leads to adversarially robust models

We can now define a signed version of the hypothesis margin that incorporates a given class label.

Definition 4 (signed hypothesis margin). Given a labeled input sample $(\mathbf{x}, c(\mathbf{x}))$ and a set \mathcal{W} of prototypes. The *signed hypothesis margin* of an input sample is

$$\text{margin}_h^c(\mathbf{x}, c(\mathbf{x}), \mathcal{W}) = \begin{cases} \text{margin}_h(\{\mathbf{x}\}, \mathcal{W}) & \text{if } \mathbf{x} \text{ is correctly classified,} \\ -\text{margin}_h(\{\mathbf{x}\}, \mathcal{W}) & \text{otherwise.} \end{cases} \quad (9)$$

Note that based on Theorem 1, the signed hypothesis margin is lower bounded by the *absolute distance difference* $\Delta(\mathbf{x})$:

$$\Delta(\mathbf{x}) = \|\mathbf{x} - \mathbf{w}^-\| - \|\mathbf{x} - \mathbf{w}^+\| \leq 2 \cdot \text{margin}_h^c(\mathbf{x}, c(\mathbf{x}), \mathcal{W}), \quad (10)$$

where the equality holds if \mathbf{w}^+ is equal to \mathbf{w}^* or \mathbf{w}_* (see Equation (2) for the definition of \mathbf{w}^+ and \mathbf{w}^-). Given a labeled test dataset \mathcal{T} and a set \mathcal{W} of prototypes—based on Equation (10)—we can calculate an upper bound on the robust test error under ϵ -limited adversarial attacks by

$$\text{error}_\epsilon(\mathcal{T}, \mathcal{W}) = \frac{\#\{(\mathbf{x}, c(\mathbf{x})) \in \mathcal{T} \mid \|\mathbf{x} - \mathbf{w}^-\| - \|\mathbf{x} - \mathbf{w}^+\| \leq 2\epsilon\}}{\#\mathcal{T}}. \quad (11)$$

Therefore, by maximizing $\Delta(\mathbf{x})$, we maximize the number of correctly classified samples with a large margin. Moreover, by Corollary 1, we can say that maximizing $\Delta(\mathbf{x})$ maximizes the robustness against adversarial examples or, in other words, minimizing $-\Delta(\mathbf{x})$ maximizes the adversarial robustness. Furthermore, the expression $-\Delta(\mathbf{x})$ is, in fact, a *triplet loss*. Consequently, to optimize an NPC for attacks less than ϵ -limited adversarial attacks one might optimize

$$\frac{1}{\#\mathcal{T}} \sum_{(\mathbf{x}, c(\mathbf{x})) \in \mathcal{T}} \text{ReLU}(\|\mathbf{x} - \mathbf{w}^+\| - \|\mathbf{x} - \mathbf{w}^-\| + 2\epsilon). \quad (12)$$

How does this result apply to the realizations proposed in Section 3? The result above states that as long as we optimize an NPC with a triplet loss, the NPC becomes adversarially robust. However, it is also known that optimizing an NPC with a triplet loss could lead to instabilities in training, which is the reason for the normalization in Equation (2). This normalization impacts the trade-off between large and small margins, but the loss still performs a margin maximization and optimizes for adversarial robustness. Additionally, this result remains true if squared seminorms are used to avoid the possible computation of the square root because the following holds:

$$\|\mathbf{x} - \mathbf{w}^+\|^2 - \|\mathbf{x} - \mathbf{w}^-\|^2 = (\|\mathbf{x} - \mathbf{w}^+\| + \|\mathbf{x} - \mathbf{w}^-\|) (\|\mathbf{x} - \mathbf{w}^+\| - \|\mathbf{x} - \mathbf{w}^-\|).$$

Consequently, we can expect that GLVQ and GTLVQ models become adversarially robust during training. In contrast, due to the softmax squashing and the cross-entropy loss, we cannot guarantee that RSLVQ models are robust against adversarial attacks.

5 Experiments

In order to verify the presented theoretical results, we performed an experimental analysis on the MNIST [64] and CIFAR-10 [65] datasets. By training several NPCs, we show that the certificate provided by Equation (11) is tight and applies to different L^p -norms. We compare the certificate

Table 1: Comparison of NPCs trained with the L^∞ -norm against state-of-the-art methods. Dashes “–” indicate that the quantity is not calculable or reported.

Dataset	Class	Model	CTE [%]	LRTE [%]	URTE [%]
MNIST $\epsilon = 0.3$	Certify	GLVQ (128 ppc)	3.66	16.39	20.58
		RSLVQ (128 ppc)	1.70	100.00	–
		RT [50, Table 3]	2.68	12.46	12.46
		CAP [12, Table 2 “Small”]	14.87	–	43.10
	Verify	RS [39, Table 3 “RS+”]	2.67	7.95	19.32
		IBP [33, Table 4]	1.66	6.12	8.05
CIFAR-10 $\epsilon = 8/255$	Certify	GLVQ (64 ppc)	59.35	79.54	79.62
		RSLVQ (128 ppc)	54.71	99.04	–
		RT [50, Table 3]	58.46	74.69	74.69
		CAP [12, Table 2 “Resnet”]	71.33	–	78.22
	Verify	RS [39, Table 3 “RS+”]	59.55	73.22	79.73
		IBP [33, Table 4]	50.51	65.23	67.96

with other methods for guaranteed robustness to highlight that NPCs have comparable guaranteed robustness against adversarial attacks. Particularly, the robust NPCs outperform all other methods in terms of the computational complexity for deriving the certificate. Based on this property, we show that NPCs are the first methods able to perform “real-time” adversarial rejection during inference.

For the L^∞ -norm, we compare GLVQ and RSLVQ with ReLU Stability training (RS) [39], Interval Bound Propagation (IBP) [33], Robust Trees (RT) [50], and Convex outer Adversarial Polytope (CAP) [12]. Out of these four, CAP and RT are certification methods. IBP and RS are robustification approaches analyzed by a verification approach.³ For the L^2 -norm, we compare GLVQ and GTLVQ with CAP, Stability Training with Noise (STN) [49], and randomized Smoothing (Smooth) [18]. Except for RT, all benchmark methods are based on NNs. Unless otherwise stated, all NPCs are trained by optimizing their respective loss function with stochastic gradient descent, see Section 3. We report the number of prototypes per class (ppc) and the subspace dimension m in the tables.

In Table 1 and Table 2, the results of the comparison are presented. We report the Clean Test Error (CTE) and an Upper bound on the Robust Test Error (URTE). For the L^∞ -norm, see Table 1, we also present a Lower bound on the Robust Test Error (LRTE), obtained using the Projected Gradient Descent (PGD) attack [41], and compare NPCs with certification and verification methods. Both LRTE and URTE are evaluated regarding ϵ -limited adversarial attacks—for example, for NPCs, the URTE regarding a certain ϵ is calculated by Equation (11). The ϵ ’s are selected in accordance with reported results in the literature. To compare space and time complexity of the certification methods, we present the number of trainable parameters (#param.) and the number of forward passes⁴ (forw. pass.) required to certify an input in the L^2 -norm setting, see Table 2—all methods presented there are *certification* methods. Further details about the experimental setting, together with additional results, can be found in the supplementary material Section D and Section E. Moreover, besides the evaluation on image datasets, we provide in the supplementary material Section F an evaluation and comparison of NPCs with robust tree-based methods on tabular data. The source code for training and evaluation is available at https://github.com/saraLajew/robust_NPCs.

Results of the comparison As expected, there is a large difference in robustness between the NPCs trained with a triplet loss (GLVQ and GTLVQ) and those trained with a different loss (RSLVQ), see Table 1. In addition to not being able to provide a guarantee, RSLVQ is not empirically robust—as shown by the trivial LRTE. GLVQ and GTLVQ, on the other hand, do provide a nontrivial robustness certificate comparable to, or even better than, the results of an NN trained with CAP as presented in Table 1 and Table 2. In combination with the small gap between LRTE and URTE in Table 1,

³For clarification, *IBP also provides a certificate* (see Table 3 in the referenced publication), but we refer to the robustness results achieved by a verification approach to compare with strong adversarial robustness results—note Footnote 4 in the referenced publication regarding possible overestimation of the reported errors.

⁴The CAP certificate is calculated using a dual network. Based on a discussion with the authors, we found that the number of forward passes in the full network can be approximated by the presented lower bounds.

Table 2: Comparison of NPCs trained with the L^2 -norm against state-of-the-art certification methods based on NNs. Values denoted with * were estimated from figures from the original publication. GTLVQ[†] was trained with the loss function from Equation (12) with an ϵ value of 1.58.

Dataset	Model	CTE [%]	URTE [%]	Forw. pass.	#param.
MNIST $\epsilon = 1.58$	GLVQ (256 ppc)	4.19	65.61	1	2.0 M
	GTLVQ [†] (10 ppc, $m = 12$)	2.92	55.32	1	1.0 M
	CAP [12, Table 4 “Large”]	11.88	55.47	≥ 749	2.0 M
	STN [49, Table 1]	1.10	31.00	100	0.7 M
CIFAR-10 $\epsilon = 36/255$	GLVQ (128 ppc)	51.41	61.90	1	3.9 M
	GTLVQ (1 ppc, $m = 100$)	40.53	55.96	1	3.1 M
	CAP [12, Table 2 “Resnet”]	38.80	48.04	≥ 3073	4.2 M
	STN [49, Table 1]	19.50	34.40	100	1.4 M
	Smooth [18, Figure 6 top, 0.12]	18*	27*	100100	1.7 M

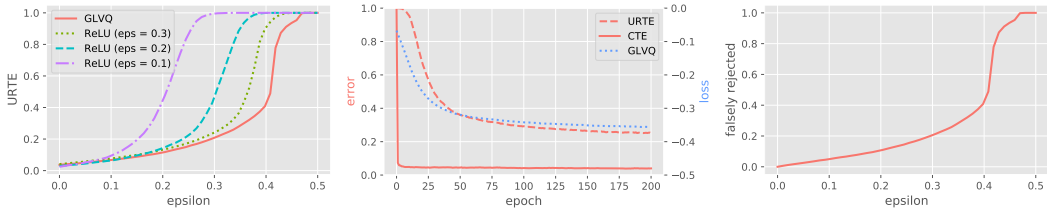


Figure 1: Different results on the MNIST test dataset of GLVQ models trained with the L^∞ -norm. Left: URTE after training with different losses. Middle: Evolution of CTE, GLVQ loss, and URTE ($\epsilon = 0.3$) during training. Right: Ratio of falsely rejected samples by the hypothesis margin.

this implies that the bound based on the hypothesis margin is not only formally but also empirically tight. However, the results from the same table show that the robust NPCs do not improve over the other models trained for high robustness but are reasonably close in terms of the URTE. Compared to current state-of-the-art robustness certification methods based on randomized smoothing (STN and Smooth) presented in Table 2, there is a considerable deficit in terms of URTE. However, it must be noted that these methods are not deterministic and require extensive sampling to compute the certificate (and prediction). To illustrate, Smooth [18, Section 4] takes *15 seconds to certify a sample from CIFAR-10*, while GLVQ *certifies the entire test dataset during this time*. This difference is also depicted in Table 2 by the number of forward passes needed to certify a single input.

Some notes on training adversarially robust NPCs We considered two triplet losses for training GLVQ and GTLVQ: the ReLU clipped absolute distance difference loss of Equation (12) and the GLVQ loss of Equation (2). While the ReLU loss optimizes for predetermined ϵ in ϵ -limited attacks, the GLVQ loss does not. Despite the former being more similar to other certification methods, we found that using the GLVQ loss often results in the highest guaranteed robustness for almost all ϵ 's, as shown in Figure 1 (left). Second, all NPCs are trained without early stopping based on plateauing CTE. Optimizing a triplet loss is directly tied to maximizing the guaranteed robustness. Hence, as long as the triplet loss is improving, the guarantee improves too, see Figure 1 (middle). This is an argument against early stopping based on plateauing CTE. Crucially, the URTE can be logged during training as it only requires one extra floating-point operation per sample.

Real-time adversarial rejection As stated in Section 4.2, the hypothesis margin can be used to reject adversarial examples with perfect recall. While this view on guaranteed adversarial robustness was voiced before, it suffered from the computational overhead of certifying a sample. With the hypothesis margin calculable with one floating-point operation after inference, it does not suffer from the same problem. To evaluate the rejection strategy, we apply it to the MNIST test dataset and the GLVQ model trained with the L^∞ -norm. The false positive rate for different values of ϵ is given in Figure 1 (right). Note that the adversarial rejection strategy is guaranteed to provide perfect recall, hence, investigating only the falsely rejected sample suffices. We find that with an ϵ of 0.3, the

adversarial rejection strategy falsely rejects only 20% of all samples. To emphasize, this is achieved without excessive overhead and perfect recall. Further investigation of the falsely rejected samples in supplementary Section E.3 shows that they are semantically close to the hypothetical original class. For example, most rejected samples from the class 9 are close to a prototype of the class 4.

6 Discussion

The experimental evaluation presents us with the following results: First of all, training NPCs with a triplet loss is an effective robustification strategy. The resulting models are empirically robust and their guaranteed robustness is tight, nontrivial, improves over other certification methods, and yields comparable performance even when compared with verification methods. More importantly, we showed a large difference in computational overhead to derive the guarantees between NN certifiers (verifiers) and the NPC certification.

There are also downsides to the work presented that we would like to clarify here. First, the method does not improve over state of the art and we expect issues in scaling to high dimensional datasets like ImageNet [66] without a proper feature engineering—this is a common challenge in NPCs. However, if resources and time are scarce, none of the current state-of-the-art methods are applicable. Hence, we do not consider not improving over these methods as a critical limitation. A second problem is that, in principle, it would be desired to guarantee robustness for several seminorms with one NPC. Unfortunately, the presented certificates are always related to the seminorm used to classify the data. But if we train a robust NPC with a *fixed* L^p -norm, Hölder’s inequality provides a simple bound on the size of adversarial perturbations δ regarding other L^q -norms: $\|\delta\|_q \geq \text{margin}_h(\{\mathbf{x}\}, \mathcal{W})$ if $q \leq p$ and $\|\delta\|_q \geq n^{\frac{1}{q} - \frac{1}{p}} \text{margin}_h(\{\mathbf{x}\}, \mathcal{W})$ otherwise. Preliminary results regarding this can be found in the supplementary material Section E.2.

We close this section with some concluding remarks considering the importance of the generality of the results as they hold for an arbitrary seminorm. A common subset of NPCs are those that use an adaptive dissimilarity measure—for example, GMLVQ and LMNN [19] with a 1-nearest neighbor rule. Both methods use the same dissimilarity measure: the quadratic form $d(\mathbf{x}, \mathbf{w}) = \|\mathbf{Q}(\mathbf{x} - \mathbf{w})\|_2$ with $\mathbf{Q} \in \mathbb{R}^{m \times n}$. The matrix \mathbf{Q} can be considered to encode the feature relevance with respect to the classification task, which is an important characteristic—for example, it is used to identify relevant biomarkers for malignancy of adrenal tumors [67, 68]. Based on the presented theory, both methods are hypothesis margin maximizers regarding the seminorm $\|\mathbf{x}\|_{\mathbf{Q}} = \|\mathbf{Q}\mathbf{x}\|_2$ since both optimize a triplet loss. Therefore, we can expect strong guaranteed robustness with respect to attacks that optimize $\|\delta\|_{\mathbf{Q}}$. However, being robust regarding $\|\delta\|_{\mathbf{Q}}$ optimized attacks does not necessarily imply robustness regarding L^p -norms. Hence, in commonly used empirical robustness evaluation settings, adaptive measure NPCs might seem to be non-robust—as empirically observed for GMLVQ [26].

7 Conclusion and outlook

In this work, we presented a theory to robustify and certify NPCs with respect to an *arbitrary* seminorm. To the best of our knowledge, this is the most universal and practically applicable approach with nontrivial robustness guarantees. Based on the hypothesis margin, we have proven an efficiently calculable and tight lower bound on the robust test error of an NPC. The numerical evaluation of this bound showed that the robustness guarantee of NPCs surpassed other NN-based certification methods and is close to verification methods. At the same time, it significantly improved the computational complexity. Together with their inherent interpretability [20, 22], NPCs are a great alternative for NNs in the adversarial setting and the superior choice when compute time is restricted.

To improve the presented results, we suggest the study of NPCs in the context of ensembles (similar to RTs [50]) and cascade models [12]. In the previous section, we discussed how robustness guarantees can be computed for various L^p -norms simultaneously. Since these guarantees are often too weak in practice, future work should examine whether the results of Croce and Hein [13] can be generalized to NPCs. Additionally, further studies should cover whether the idea of the hypothesis margin—*varying parameters instead of inputs* to derive a calculable margin—can be extended to NNs.

Broader impact

With the more widespread application of machine learning methods in our everyday life, the potential negative impact of adversarial attacks on society increases. As discussed in the introduction, neither current empirical robustness methods nor certification or verification methods are sufficient to eliminate this problem. For applied machine learning research in medium to large companies, the current state-of-the-art methods for certifying or verifying adversarial robustness require a too large investment in compute time to truly incorporate the guaranteed robustness of a model as a formal requirement for the productization of machine learning. The theoretical robustness bound presented in this work can however be parallelized with the accuracy evaluation of a model and can therefore be easily incorporated in the already existing evaluation pipelines. With the upper bound on the robust test error calculable in constant time, it is even possible to incorporate the certification of an NPC as a metric in the training procedure—outputting the certified adversarial robustness after each epoch. A potential application of the reduced impact on inference time is also discussed in Section 5.

Although deep neural networks frequently deliver excellent performances, the interpretability of those networks is difficult [69]. Recently, this has led to a wide line of research into the development of interpretable models, particularly for technical and medical applications [23, 24, 70, 71]. Having this in mind, the consideration of NPCs, as one of the most prominent interpretable paradigms [22], is of general interest. However, to ensure that interpretable models can be used without unwanted negative side effects, it is important to investigate their properties to the same extent as has been done for non-interpretable models. The investigation of the guaranteed adversarial robustness of NPCs is, therefore, a crucial step in this transition. In addition to this, the positive definiteness of norm-based distances can impose a significant restriction on the dissimilarity measure used in NPCs. By showing that this is not a requirement for constructing an adversarially robust NPC, this restriction is removed. Hence, more freedom is obtained in the selection of dissimilarity measures—for example, adaptive dissimilarity measures, as discussed in Section 6. This allows the application of NPCs as interpretable models in a wider variety of use cases.

To summarize, we foresee two potential areas where the theoretical work presented here could have a direct and lasting impact on society. First, with the upper bound on the robust test error calculable in constant-time, the certification method presented here is more suitable for direct incorporation in the development of machine learning methods. This has been extensively discussed and evaluated in previous sections. Second, as a side effect, with NPCs now proven to be robust against adversarial attacks, they are better suited and more widely applicable as an interpretable alternative to NNs in real-world applications.

Acknowledgments and disclosure of funding

We would like to thank Peter Schlicht for his valuable contribution to earlier versions of the manuscript and Eric Wong for his helpful discussion about CAP. Moreover, we would like to thank our attentive anonymous AC and reviewers whose comments have greatly improved this manuscript.

None of the authors received third party funding or have had any financial relationship with entities that could potentially be perceived to influence the submitted work during the 36 months prior to this submission.

References

- [1] Amir Globerson and Sam Roweis. Nightmare at test time: Robust learning by feature deletion. In William W. Cohen and Andrew W. Moore, editors, *Proceedings of the 23rd International Conference on Machine Learning – ICML 2006*, pages 353–360, Pittsburgh, PA, USA, 2006. ACM.
- [2] Huan Xu, Constantine Caramanis, and Shie Mannor. Robustness and regularization of support vector machines. *Journal of Machine Learning Research*, 10(51):1485–1510, 2009.
- [3] Battista Biggio, Iginio Corona, Blaine Nelson, Benjamin I. P. Rubinstein, Davide Maiorca, Giorgio Fumera, Giorgio Giacinto, and Fabio Roli. Security evaluation of support vector machines in adversarial environments. In *Support Vector Machines Applications*, pages 105–153. Springer, 2014.

- [4] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In Yoshua Bengio and Yann LeCun, editors, *Proceedings of the 2nd International Conference on Learning Representations – ICLR 2014*, Banff, AB, Canada, 2014.
- [5] Thomas Brunner, Frederik Diehl, Michael Truong Le, and Alois Knoll. Guessing smart: Biased sampling for efficient black-box adversarial attacks. In *Proceedings of the IEEE International Conference on Computer Vision – ICCV 2019*, pages 4958–4966, Seoul, South Korea, 2019. IEEE.
- [6] Kevin Eykholt, Ivan Evtimov, Earlene Fernandes, Bo Li, Amir Rahmati, Chaowei Xiao, Atul Prakash, Tadayoshi Kohno, and Dawn Son. Robust physical-world attacks on deep learning visual classification. In *Proceedings of the 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition – CVPR 2018*, pages 1625–1634, Salt Lake City, UT, USA, 2018. IEEE.
- [7] Felix Assion, Peter Schlicht, Florens Gressner, Wiebke Günther, Fabian Hüger, Nico Schmidt, and Umair Rasheed. The attack generator: A systematic approach towards constructing adversarial attacks. In *Workshop proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition – CVPR 2019 Workshops*, Long Beach, CA, USA, 2019.
- [8] Nicolas Papernot, Patrick McDaniel, Xi Wu, Somesh Jha, and Ananthram Swami. Distillation as a defense to adversarial perturbations against deep neural networks. In *Proceedings of the 2016 IEEE Symposium on Security and Privacy – SP 2016*, pages 582–597, San Jose, CA, USA, 2016. IEEE.
- [9] Harini Kannan, Alexey Kurakin, and Ian Goodfellow. Adversarial logit pairing. *arXiv preprint arXiv:1803.06373 [cs.LG]*, 2018.
- [10] Lukas Schott, Jonas Rauber, Matthias Bethge, and Wieland Brendel. Towards the first adversarially robust neural network model on MNIST. In *Proceedings of the 7th International Conference on Learning Representations – ICLR 2019*, New Orleans, LA, USA, May 2019. OpenReview.net.
- [11] Francesco Croce, Maksym Andriushchenko, and Matthias Hein. Provable robustness of ReLU networks via maximization of linear regions. In Kamalika Chaudhuri and Masashi Sugiyama, editors, *Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics – AISTATS 2019*, volume 89 of the Proceedings of Machine Learning Research, pages 2057–2066, Naha, Okinawa, Japan, 2019. PMLR.
- [12] Eric Wong, Frank Schmidt, Jan Hendrik Metzen, and J. Zico Kolter. Scaling provable adversarial defenses. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32: Proceedings of the Neural Information Processing Systems Conference – NeurIPS 2019*, pages 8400–8409, Vancouver, BC, Canada, 2019. Curran Associates, Inc.
- [13] Francesco Croce and Matthias Hein. Provable robustness against all adversarial l_p -perturbations for $p \geq 1$. In *Proceedings of the 8th International Conference on Learning Representations – ICLR 2020*. OpenReview.net, 2020.
- [14] Nicholas Carlini and David A. Wagner. Towards evaluating the robustness of neural networks. In *Proceedings of the 2017 IEEE Symposium on Security and Privacy – SP 2017*, pages 39–57, San Jose, CA, USA, 2017. IEEE.
- [15] Chawin Sitawarin and David Wagner. On the robustness of deep k-nearest neighbors. In *Workshop proceedings of the 2019 IEEE Symposium on Security and Privacy Workshops – SP 2019 Workshops*, pages 1–7, San Francisco, CA, USA, 2019. IEEE.
- [16] Nicholas Carlini. Is Aml (Attacks Meet Interpretability) robust to adversarial examples? *arXiv preprint arXiv:1902.02322 [cs.LG]*, 2019.
- [17] Vincent Tjeng, Kai Y. Xiao, and Russ Tedrake. Evaluating robustness of neural networks with mixed integer programming. In *Proceedings of the 7th International Conference on Learning Representations – ICLR 2019*, New Orleans, LA, USA, 2019. OpenReview.net.

- [18] 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 – ICML 2019*, volume 97 of the Proceedings of Machine Learning Research, pages 1310–1320, Long Beach, CA, USA, 2019. PMLR.
- [19] Kilian Q. Weinberger and Lawrence K. Saul. Distance metric learning for large margin nearest neighbor classification. *Journal of Machine Learning Research*, 10:207–244, 2009.
- [20] Jacob Bien and Robert Tibshirani. Prototype selection for interpretable classification. *The Annals of Applied Statistics*, 5(4):2403–2424, 2011.
- [21] David Nova and Pablo A. Estévez. A review of learning vector quantization classifiers. *Neural Computing and Applications*, 25(3-4):511–524, 2014.
- [22] Michael Biehl, Barbara Hammer, and Thomas Villmann. Prototype-based models in machine learning. *Wiley Interdisciplinary Reviews Cognitive Science*, 7(2):92–111, 2016.
- [23] Cynthia Rudin. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence*, 1:206–215, 2019.
- [24] Alfredo Vellido. The importance of interpretability and visualization in machine learning for applications in medicine and health care. *Neural Computing and Applications*, 2019.
- [25] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30: Proceedings of the Neural Information Processing Systems Conference – NIPS 2017*, pages 4077–4087, Long Beach, CA, USA, 2017. Curran Associates, Inc.
- [26] Sascha Saralajew, Lars Holdijk, Maike Rees, and Thomas Villmann. Robustness of generalized learning vector quantization models against adversarial attacks. In Alfredo Vellido, Karina Gibert, Cecilio Angulo, and José David Martín-Guerrero, editors, *Advances in Self-Organizing Maps, Learning Vector Quantization, Clustering and Data Visualization: Proceedings of the 13th International Workshop, WSOM+ 2019*, volume 976 of the Advances in Intelligent Systems and Computing, pages 189–199, Barcelona, Spain, 2019. Springer, Cham.
- [27] Alhussein Fawzi, Hamza Fawzi, and Omar Fawzi. Adversarial vulnerability for any classifier. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31: Proceedings of the Neural Information Processing Systems Conference – NeurIPS 2018*, pages 1178–1187, Montréal, QC, Canada, 2018. Curran Associates, Inc.
- [28] Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry. Adversarial examples are not bugs, they are features. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32: Proceedings of the Neural Information Processing Systems Conference – NeurIPS 2019*, pages 125–136, Vancouver, BC, Canada, 2019. Curran Associates, Inc.
- [29] Carl-Johann Simon-Gabriel, Yann Ollivier, Leon Bottou, Bernhard Schölkopf, and David Lopez-Paz. First-order adversarial vulnerability of neural networks and input dimension. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning – ICML 2019*, volume 97 of the Proceedings of Machine Learning Research, pages 5809–5817, Long Beach, CA, USA, 2019. PMLR.
- [30] Zhuozhuo Tu, Jingwei Zhang, and Dacheng Tao. Theoretical analysis of adversarial learning: A minimax approach. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32: Proceedings of the Neural Information Processing Systems Conference – NeurIPS 2019*, pages 12280–12290, Vancouver, BC, Canada, 2019. Curran Associates, Inc.

- [31] Moustapha Cisse, Piotr Bojanowski, Edouard Grave, Yann Dauphin, and Nicolas Usunier. Parseval networks: Improving robustness to adversarial examples. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning – ICML 2017*, volume 70 of the Proceedings of Machine Learning Research, pages 854–863, Sydney, NSW, Australia, 2017. PMLR.
- [32] Gamaleldin F. Elsayed, Dilip Krishnan, Hossein Mobahi, Kevin Regan, and Samy Bengio. Large margin deep networks for classification. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31: Proceedings of the Neural Information Processing Systems Conference – NeurIPS 2018*, pages 842–852, Montréal, QC, Canada, 2018. Curran Associates, Inc.
- [33] Sven Gowal, Krishnamurthy Dvijotham, Robert Stanforth, Rudy Bunel, Chongli Qin, Jonathan Uesato, Relja Arandjelovic, Timothy Mann, and Pushmeet Kohli. On the effectiveness of interval bound propagation for training verifiably robust models. *arXiv preprint arXiv:1810.12715v4 [cs.LG]*, 2019.
- [34] Guy Katz, Clark Barrett, David L. Dill, Kyle Julian, and Mykel J. Kochenderfer. Reluplex: An efficient SMT solver for verifying deep neural networks. In Rupak Majumdar and Viktor Kunčák, editors, *29th International Conference on Computer Aided Verification – CAV 2017*, volume 10426 of the Lecture Notes in Computer Science, pages 97–117, Heidelberg, Germany, 2017. Springer International Publishing.
- [35] Matt Jordan, Justin Lewis, and Alexandros G. Dimakis. Provable certificates for adversarial examples: Fitting a ball in the union of polytopes. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32: Proceedings of the Neural Information Processing Systems Conference – NeurIPS 2019*, pages 14082–14092, Vancouver, BC, Canada, 2019. Curran Associates, Inc.
- [36] Matthias Hein and Maksym Andriushchenko. Formal guarantees on the robustness of a classifier against adversarial manipulation. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30: Proceedings of the Neural Information Processing Systems Conference – NIPS 2017*, pages 2266–2276, Long Beach, CA, USA, 2017. Curran Associates, Inc.
- [37] Aditi Raghunathan, Jacob Steinhardt, and Percy Liang. Certified defenses against adversarial examples. In *Proceedings of the 6th International Conference on Learning Representations – ICLR 2018*, Vancouver, BC, Canada, 2018. Vancouver, BC, Canada.
- [38] Huan Zhang, Tsui-Wei Weng, Pin-Yu Chen, Cho-Jui Hsieh, and Luca Daniel. Efficient neural network robustness certification with general activation functions. In Samy Bengio Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31: Proceedings of the Neural Information Processing Systems Conference – NeurIPS 2018*, pages 4939–4948, Montréal, QC, Canada, 2018. Curran Associates, Inc.
- [39] Kai Y. Xiao, Vincent Tjeng, Nur Muhammad (Mahi) Shafiullah, and Aleksander Madry. Training for faster adversarial robustness verification via inducing ReLU stability. In *Proceedings of the 7th International Conference on Learning Representations – ICLR 2019*, New Orleans, LA, USA, 2019. OpenReview.net.
- [40] Eric Wong and Zico Kolter. Provable defenses against adversarial examples via the convex outer adversarial polytope. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning – ICML 2018*, volume 80 of the Proceedings of Machine Learning Research, pages 5286–5295, Stockholm, Sweden, 2018. PMLR.
- [41] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *Proceedings of the 6th International Conference on Learning Representations – ICLR 2018*, Vancouver, BC, Canada, 2018. OpenReview.net.

- [42] Florian Tramer and Dan Boneh. Adversarial training and robustness for multiple perturbations. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32: Proceedings of the Neural Information Processing Systems Conference – NeurIPS 2019*, pages 5866–5876, Vancouver, BC, Canada, 2019. Curran Associates, Inc.
- [43] Nicolas Papernot and Patrick McDaniel. Deep k-nearest neighbors: Towards confident, interpretable and robust deep learning. *arXiv preprint arXiv:1803.04765 [cs.LG]*, 2018.
- [44] Guanhong Tao, Shiqing Ma, Yingqi Liu, and Xiangyu Zhang. Attacks meet interpretability: Attribute-steered detection of adversarial samples. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31: Proceedings of the Neural Information Processing Systems Conference – NeurIPS 2018*, pages 7717–7728, Montréal, QC, Canada, 2018. Curran Associates, Inc.
- [45] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. DeepFool: A simple and accurate method to fool deep neural networks. In *Proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition – CVPR 2016*, pages 2574–2582, Las Vegas, NV, USA, 2016. IEEE.
- [46] Wieland Brendel, Jonas Rauber, and Matthias Bethge. Decision-based adversarial attacks: Reliable attacks against black-box machine learning models. In *Proceedings of the 6th International Conference on Learning Representations – ICLR 2018*, Vancouver, BC, Canada, 2018. OpenReview.net.
- [47] Marius Mosbach, Maksym Andriushchenko, Thomas Trost, Matthias Hein, and Dietrich Klakow. Logit pairing methods can fool gradient-based attacks. *arXiv preprint arXiv:1810.12042*, 2018.
- [48] Guang-He Lee, Yang Yuan, Shiyu Chang, and Tommi Jaakkola. Tight certificates of adversarial robustness for randomly smoothed classifiers. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32: Proceedings of the Neural Information Processing Systems Conference – NeurIPS 2019*, pages 4910–4921, Vancouver, BC, Canada, 2019. Curran Associates, Inc.
- [49] Bai Li, Changyou Chen, Wenlin Wang, and Lawrence Carin. Certified adversarial robustness with additive noise. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32: Proceedings of the Neural Information Processing Systems Conference – NeurIPS 2019*, pages 9464–9474, Vancouver, BC, Canada, 2019. Curran Associates, Inc.
- [50] Maksym Andriushchenko and Matthias Hein. Provably robust boosted decision stumps and trees against adversarial attacks. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32: Proceedings of the Neural Information Processing Systems Conference – NeurIPS 2019*, pages 13017–13028, Vancouver, BC, Canada, 2019. Curran Associates, Inc.
- [51] Hongge Chen, Huan Zhang, Duane Boning, and Cho-Jui Hsieh. Robust decision trees against adversarial examples. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning – ICML 2019*, volume 97 of the Proceedings of Machine Learning Research, pages 1122–1131, Long Beach, CA, USA, 2019. PMLR.
- [52] Paolo Russu, Ambra Demontis, Battista Biggio, Giorgio Fumera, and Fabio Roli. Secure kernel machines against evasion attacks. In *Proceedings of the 2016 ACM Workshop on Artificial Intelligence and Security – ALSec 2016*, pages 59–69, Vienna, Austria, 2016. ACM.
- [53] Atsushi Sato and Keiji Yamada. Generalized learning vector quantization. In David S. Touretzky, Michael Mozer, and Michael E. Hasselmo, editors, *Advances in Neural Information Processing Systems 8: Proceedings of the Neural Information Processing Systems Conference – NIPS 1995*, pages 423–429, Denver, CO, USA, 1996. MIT Press.

- [54] Teuvo Kohonen. Improved versions of learning vector quantization. In *Proceedings of the 1990 International Joint Conference on Neural Networks – IJCNN 1990*, pages 545–550, San Diego, CA, USA, 1990. IEEE.
- [55] Teuvo Kohonen. *Self-Organizing Maps*, volume 30 of the Springer Series in Information Sciences, chapter Learning Vector Quantization, pages 175–189. Springer, Berlin, Heidelberg, 1995.
- [56] Koby Crammer, Ran Gilad-Bachrach, Amir Navot, and Naftali Tishby. Margin analysis of the LVQ algorithm. In Suzanna Becker, Sebastian Thrun, and Klaus Obermayer, editors, *Advances in Neural Information Processing Systems 15: Proceedings of the Neural Information Processing Systems Conference – NIPS 2002*, pages 479–486, Vancouver, BC, Canada, 2003. MIT Press.
- [57] Yizhen Wang, Somesh Jha, and Kamalika Chaudhuri. Analyzing the robustness of nearest neighbors to adversarial examples. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning – ICML 2018*, volume 80 of the Proceedings of Machine Learning Research, pages 5133–5142, Stockholm, Sweden, 2018. PMLR.
- [58] Lu Wang, Xuanqing Liu, Jinfeng Yi, Zhi-Hua Zhou, and Cho-Jui Hsieh. Evaluating the robustness of nearest neighbor classifiers: A primal-dual perspective. *arXiv preprint arXiv:1906.03972v1 [cs.LG]*, 2019.
- [59] Yao-Yuan Yang, Cyrus Rashtchian, Yizhen Wang, and Kamalika Chaudhuri. Robustness for non-parametric classification: A generic attack and defense. In Silvia Chiappa and Roberto Calandra, editors, *Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics – AISTATS 2020*, volume 108 of the Proceedings of Machine Learning Research, pages 941–951, Online, 2020. PMLR.
- [60] Johannes Brinkrolf and Barbara Hammer. Interpretable machine learning with reject option. *at - Automatisierungstechnik*, 66(4):283 – 290, 2018.
- [61] Petra Schneider, Michael Biehl, and Barbara Hammer. Adaptive relevance matrices in learning vector quantization. *Neural Computation*, 21(12):3532–3561, 2009.
- [62] Sascha Saralajew and Thomas Villmann. Adaptive tangent distances in generalized learning vector quantization for transformation and distortion invariant classification learning. In *Proceedings of the 2016 International Joint Conference on Neural Networks – IJCNN 2016*, pages 2672–2679, Vancouver, BC, Canada, 2016. IEEE.
- [63] Sambu Seo and Klaus Obermayer. Soft learning vector quantization. *Neural Computation*, 15(7):1589–1604, 2003.
- [64] Yann LeCun, Corinna Cortes, and Christopher J.C. Burges. The MNIST database of handwritten digits. 1998. <http://yann.lecun.com/exdb/mnist/>.
- [65] Alex Krizhevsky. Learning multiple layers of features from tiny images. techreport, Department of Computer Science, University of Toronto, Toronto, ON, Canada, 2009.
- [66] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *Proceedings of the 2009 IEEE Conference on Computer Vision and Pattern Recognition – CVPR 2009*, pages 248–255, Miami, Florida, USA, 2009. IEEE.
- [67] Irina Bancos, Angela E. Taylor, Vasileios Chortis, Alice J. Sitch, Carl Jenkinson, Caroline J. Davidge-Pitts, Katharina Lang, Stylianos Tsagarakis, Magdalena Macech, Anna Riestler, et al. Urine steroid metabolomics for the differential diagnosis of adrenal incidentalomas in the eurine-act study: a prospective test validation study. *The Lancet Diabetes & Endocrinology*, 8(9):773–781, 2020.
- [68] Wiebke Arlt, Michael Biehl, Angela E. Taylor, Stefanie Hahner, Rossella Libé, Beverly A. Hughes, Petra Schneider, David J. Smith, Han Stiekema, Nils Krone, Emilio Porfiri, Giuseppe

- Opocher, Jérôme Bertherat, Franco Mantero, Bruno Allolio, Massimo Terzolo, Peter Nightingale, Cedric H. L. Shackleton, Xavier Bertagna, Martin Fassnacht, and Paul M. Stewart. Urine steroid metabolomics as a biomarker tool for detecting malignancy in adrenal tumors. *The Journal of Clinical Endocrinology & Metabolism*, 96(12):3775–3784, 2011.
- [69] Wojciech Samek, Grégoire Montavon, Andrea Vedaldi, Lars Kai Hansen, and Klaus-Robert Müller, editors. *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*, volume 11700 of the Lecture Notes in Computer Science. Springer, 2019.
- [70] Diogo V. Carvalho, Eduardo M. Pereira, and Jaime S. Cardoso. Machine learning interpretability: A survey on methods and metrics. *Electronics*, 8(8):832, 2019.
- [71] Sascha Saralajew, Lars Holdijk, Maïke Rees, Ebubekir Asan, and Thomas Villmann. Classification-by-components: Probabilistic modeling of reasoning over a set of components. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32: Proceedings of the Neural Information Processing Systems Conference – NeurIPS 2019*, pages 2792–2803, Vancouver, BC, Canada, 2019. Curran Associates, Inc.
- [72] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun, editors, *Proceedings of the 3rd International Conference on Learning Representations – ICLR 2015*, San Diego, CA, USA, 2015.
- [73] Chih-Chung Chang and Chih-Jen Lin. LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2:27:1–27:27, 2011. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- [74] Dheeru Dua and Casey Graff. UCI machine learning repository. 2017.
- [75] William H Wolberg and Olvi L Mangasarian. Multisurface method of pattern separation for medical diagnosis applied to breast cytology. *Proceedings of the National Academy of Sciences of the United States of America*, 87(23):9193–9196, 1990.
- [76] Jack W. Smith, J. E. Everhart, W. C. Dickson, W. C. Knowler, and R. S. Johannes. Using the ADAP learning algorithm to forecast the onset of diabetes mellitus. In *Proceedings of the Annual Symposium on Computer Application in Medical Care*, pages 261–265, 1988.
- [77] Andrew V. Uzilov, Joshua M. Keegan, and David H. Mathews. Detection of non-coding RNAs on the basis of predicted secondary structure formation free energy change. *BMC Bioinformatics*, 7(173), 2006.

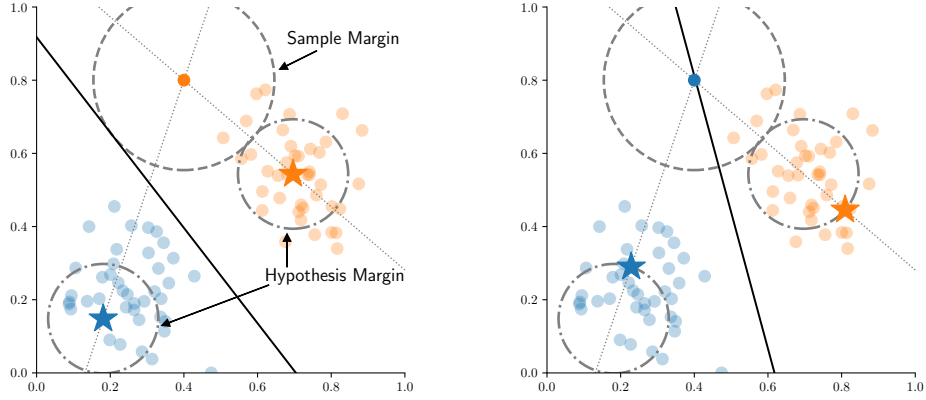


Figure 2: Visualization of the defined margins for a GLVQ model with the Euclidean distance in \mathbb{R}^2 . The softly colored dots are the training samples from the class “orange” or “blue”. The fully colored dot is a new sample that has to be classified by the model. The two stars represent the prototypes and thus the model weights trained by the GLVQ algorithm. The solid black line is the resulting decision boundary. Left: The final weight configuration of the trained GLVQ model including the margins (dashed circles) with respect to the given data sample. Right: The perturbed weight configuration that changes the assigned class label of the new data sample. Note that the prototypes have been shifted by a magnitude equal to the hypothesis margin.

A Visualization of the hypothesis and sample margins

In Figure 2, we visualize the effect of changing the weights of a GLVQ model by a magnitude equal to the hypothesis margin (with respect to a given sample). We note that if we span a ball with a radius equal to the hypothesis margin around the prototypes, the prototypes can be placed at an arbitrary position *inside* these balls without changing the assigned class label of the sample. Likewise, as the right visualization shows, we can systematically find a position for the prototypes such that they have been shifted by a magnitude *equal* to the hypothesis margin such that the prototypes assign a different class label to the given sample.

In the next sections, we use the observation that the prototypes have been shifted directly towards or away from the data point to formally prove the calculation of the hypothesis margin and the relation between the hypothesis and sample margin.

B Proof of Theorem 1: Calculation of the hypothesis margin

Theorem. *Let the data space \mathcal{X} be a vector space over a field of the real or complex numbers, $d(\mathbf{x}, \mathbf{w}) = \|\mathbf{x} - \mathbf{w}\|$ be a dissimilarity induced by a seminorm $\|\cdot\|$, and $\mathbf{x} \in \mathcal{X}$ be an input. Then, the hypothesis margin of the set \mathcal{W} of prototypes with respect to \mathbf{x} can be computed by*

$$\text{margin}_h(\{\mathbf{x}\}, \mathcal{W}) = \frac{1}{2} (\|\mathbf{x} - \mathbf{w}_*\| - \|\mathbf{x} - \mathbf{w}^*\|), \quad (13)$$

where \mathbf{w}^* denotes the closest prototype to \mathbf{x} and \mathbf{w}_* denotes the closest prototype to \mathbf{x} with a different class label than the class label of \mathbf{w}^* .

The proof is based on the ideas used by Crammer et al. [56].

Proof. The outline of the proof is as follows:

1. Given a set \mathcal{W} of prototypes, we define a set $\hat{\mathcal{W}}$ of shifted prototypes in which the prototypes are shifted to an arbitrary position within the corresponding balls of radius

$$r = \frac{1}{2} (\|\mathbf{x} - \mathbf{w}_*\| - \|\mathbf{x} - \mathbf{w}^*\|). \quad (14)$$

Using this set of shifted prototypes, we prove that

$$\text{margin}_h(\{\mathbf{x}\}, \mathcal{W}) \geq r.$$

2. For an arbitrary but sufficiently small $\varepsilon > 0$, we define a *second* set $\hat{\mathcal{W}}$ of shifted prototypes in which the prototypes are shifted by a vector of length $r + \frac{\varepsilon}{2}$. Based on this set, we prove that it assigns a different class label to the input \mathbf{x} . Consequently, we can conclude that

$$\text{margin}_h(\{\mathbf{x}\}, \mathcal{W}) \leq r.$$

Both steps together prove the theorem.

We use the following notations in the proof: Let \mathbf{x} , \mathbf{w}^* , and \mathbf{w}_* be defined as in the theorem, \mathbf{w} be an arbitrary prototype of \mathcal{W} , \mathbf{w}_\diamond be an arbitrary prototype of \mathcal{W} with a class label different than $c(\mathbf{w}^*)$, and r be defined as in Equation (14). With the dissimilarity induced by a seminorm, we use the notations $d(\mathbf{x}, \mathbf{w})$ and $\|\mathbf{x} - \mathbf{w}\|$ interchangeably.

Step 1: Let $\hat{\mathcal{W}}$ be a set of shifted prototypes constructed by repositioning each prototype in \mathcal{W} to an arbitrary position within its induced ball of radius r . Formally, this is realized by defining an arbitrary function $\mathbf{s} : \mathcal{W} \rightarrow \mathcal{X}$ such that $\|\mathbf{s}(\mathbf{w})\| < r$. The shifted prototype to \mathbf{w} is obtained by

$$\hat{\mathbf{w}}(\mathbf{w}) = \mathbf{w} + \mathbf{s}(\mathbf{w}).$$

All shifted prototypes combined, provide the set $\hat{\mathcal{W}}$ of shifted prototypes.

Given an arbitrary prototype \mathbf{w} , we conclude that the dissimilarity between the shifted and the non-shifted prototype is always less than r :

$$d(\mathbf{w}, \hat{\mathbf{w}}(\mathbf{w})) = \|\mathbf{w} - \hat{\mathbf{w}}(\mathbf{w})\| = \|\mathbf{s}(\mathbf{w})\| < r.$$

Using the triangle inequality, we can state that the dissimilarity between a shifted prototype and the input \mathbf{x} is less than $d(\mathbf{x}, \mathbf{w}) + r$:

$$\begin{aligned} d(\mathbf{x}, \hat{\mathbf{w}}(\mathbf{w})) &\leq d(\mathbf{x}, \mathbf{w}) + d(\mathbf{w}, \hat{\mathbf{w}}(\mathbf{w})), \\ &< d(\mathbf{x}, \mathbf{w}) + r. \end{aligned} \tag{15}$$

Similarly, using the triangle inequality again, we can state that

$$\begin{aligned} d(\mathbf{x}, \hat{\mathbf{w}}(\mathbf{w})) &\geq d(\mathbf{x}, \mathbf{w}) - d(\mathbf{w}, \hat{\mathbf{w}}(\mathbf{w})), \\ &> d(\mathbf{x}, \mathbf{w}) - r. \end{aligned} \tag{16}$$

Given a prototype $\mathbf{w}_\diamond \in \mathcal{W}$ with a different class label than $c(\mathbf{w}^*)$, it follows that

$$d(\mathbf{x}, \mathbf{w}_\diamond) \geq d(\mathbf{x}, \mathbf{w}_*).$$

Using Equation (16), we conclude

$$d(\mathbf{x}, \hat{\mathbf{w}}(\mathbf{w}_\diamond)) > d(\mathbf{x}, \mathbf{w}_\diamond) - r \geq d(\mathbf{x}, \mathbf{w}_*) - r. \tag{17}$$

By the definition of r , see Equation (14), it holds that

$$\begin{aligned} d(\mathbf{x}, \mathbf{w}^*) + r &= d(\mathbf{x}, \mathbf{w}^*) + \frac{1}{2}(d(\mathbf{x}, \mathbf{w}_*) - d(\mathbf{x}, \mathbf{w}^*)), \\ &= d(\mathbf{x}, \mathbf{w}_*) - r. \end{aligned} \tag{18}$$

Now, combining Equation (15) for the shifted prototype $\hat{\mathbf{w}}(\mathbf{w}^*)$ with Equation (18) and Equation (17), we obtain

$$d(\mathbf{x}, \hat{\mathbf{w}}(\mathbf{w}^*)) < d(\mathbf{x}, \mathbf{w}^*) + r = d(\mathbf{x}, \mathbf{w}_*) - r < d(\mathbf{x}, \hat{\mathbf{w}}(\mathbf{w}_\diamond)).$$

With \mathbf{w}_\diamond being an arbitrary prototype of \mathcal{W} with a class label different than $c(\mathbf{w}^*)$, this states that each shifted prototype $\hat{\mathbf{w}}(\mathbf{w}_\diamond)$ has a larger dissimilarity than $\hat{\mathbf{w}}(\mathbf{w}^*)$ to the input sample \mathbf{x} . In other words, the closest shifted prototype in $\hat{\mathcal{W}}$ to \mathbf{x} has the same class label as the closest prototype in \mathcal{W} —in symbols, $c_{\hat{\mathcal{W}}}^*(\mathbf{x}) = c_{\mathcal{W}}^*(\mathbf{x})$. Since the shift of the prototypes of \mathcal{W} to the new positions in $\hat{\mathcal{W}}$ was arbitrary, we conclude that *the hypothesis margin of \mathcal{W} is greater than or equal to r .*

Step 2: For an arbitrary $\varepsilon > 0$ that is less than or equal to $\|\mathbf{x} - \mathbf{w}_*\|$, we define a *new* set $\hat{\mathcal{W}}$ of shifted prototypes of \mathcal{W} by shifting each prototype by a magnitude of $r + \frac{\varepsilon}{2}$. Therefore, we relocate the prototypes inside balls of size $r + \varepsilon$. For each prototype $\mathbf{w} \in \mathcal{W}$, we define a unit vector with respect to $\|\cdot\|$ by

$$\mathbf{z}(\mathbf{w}) = \begin{cases} \mathbf{u} & \text{if } \|\mathbf{x} - \mathbf{w}\| = 0, \\ \frac{\mathbf{x} - \mathbf{w}}{\|\mathbf{x} - \mathbf{w}\|} & \text{otherwise,} \end{cases} \quad (19)$$

where \mathbf{u} is an arbitrary unit vector. Because the set \mathcal{W} of prototypes assigns a class label unambiguously to \mathbf{x} , the seminorm in Equation (19) vanishes only for prototypes of the class $c(\mathbf{w}^*)$.

The set $\hat{\mathcal{W}}$ of shifted prototypes is defined by the shifted prototypes $\hat{\mathbf{w}}(\mathbf{w})$ according to the following equation:

$$\hat{\mathbf{w}}(\mathbf{w}) = \begin{cases} \mathbf{w} + \left(r + \frac{\varepsilon}{2}\right) \mathbf{z}(\mathbf{w}) & \text{if } c(\mathbf{w}) \neq c(\mathbf{w}^*), \\ \mathbf{w} - \left(r + \frac{\varepsilon}{2}\right) \mathbf{z}(\mathbf{w}) & \text{otherwise.} \end{cases}$$

Keeping in mind the absolute homogeneity of a seminorm, we can state that for each prototype $\hat{\mathbf{w}}(\mathbf{w}) \in \hat{\mathcal{W}}$ that has the *same* class label as \mathbf{w}^* , the equality

$$\begin{aligned} \|\mathbf{x} - \hat{\mathbf{w}}(\mathbf{w})\| &= \left\| (\mathbf{x} - \mathbf{w}) \left(1 + \frac{r + \frac{\varepsilon}{2}}{\|\mathbf{x} - \mathbf{w}\|}\right) \right\|, \\ &= \left(1 + \frac{r + \frac{\varepsilon}{2}}{\|\mathbf{x} - \mathbf{w}\|}\right) \|\mathbf{x} - \mathbf{w}\|, \\ &= \|\mathbf{x} - \mathbf{w}\| + r + \frac{\varepsilon}{2} \end{aligned} \quad (20)$$

holds. Note that this equation is also valid if $\|\mathbf{x} - \mathbf{w}\| = 0$: From the triangle inequality, it follows that

$$\begin{aligned} \|\mathbf{x} - \hat{\mathbf{w}}(\mathbf{w})\| &= \left\| (\mathbf{x} - \mathbf{w}) + \left(r + \frac{\varepsilon}{2}\right) \mathbf{u} \right\|, \\ &\leq \|\mathbf{x} - \mathbf{w}\| + \left\| \left(r + \frac{\varepsilon}{2}\right) \mathbf{u} \right\|, \\ &\leq r + \frac{\varepsilon}{2}. \end{aligned} \quad (21)$$

Similarly, the triangle inequality implies that the seminorm of $\left(r + \frac{\varepsilon}{2}\right) \mathbf{u}$ is bounded by $\|\mathbf{x} - \hat{\mathbf{w}}(\mathbf{w})\|$:

$$\begin{aligned} \left\| \left(r + \frac{\varepsilon}{2}\right) \mathbf{u} \right\| &= \left\| \left(r + \frac{\varepsilon}{2}\right) \mathbf{u} + (\mathbf{x} - \mathbf{w}) - (\mathbf{x} - \mathbf{w}) \right\|, \\ &\leq \left\| (\mathbf{x} - \mathbf{w}) + \left(r + \frac{\varepsilon}{2}\right) \mathbf{u} \right\| + \|\mathbf{x} - \mathbf{w}\|, \\ &\leq \|\mathbf{x} - \hat{\mathbf{w}}(\mathbf{w})\|. \end{aligned}$$

Therefore, we obtain

$$\|\mathbf{x} - \hat{\mathbf{w}}(\mathbf{w})\| \geq r + \frac{\varepsilon}{2}. \quad (22)$$

The combination of Equation (21) and Equation (22) yields equality:

$$\|\mathbf{x} - \hat{\mathbf{w}}(\mathbf{w})\| = r + \frac{\varepsilon}{2}.$$

Consequently, Equation (20) is valid for *all* prototypes $\hat{\mathbf{w}}(\mathbf{w})$ with the *same* class label as \mathbf{w}^* .

Analogously, if $\hat{\mathbf{w}}(\mathbf{w}) \in \hat{\mathcal{W}}$ has a *different* class label than $c(\mathbf{w}^*)$, the seminorm of $\mathbf{x} - \hat{\mathbf{w}}(\mathbf{w})$ becomes

$$\begin{aligned} \|\mathbf{x} - \hat{\mathbf{w}}(\mathbf{w})\| &= \left\| (\mathbf{x} - \mathbf{w}) \left(1 - \frac{r + \frac{\varepsilon}{2}}{\|\mathbf{x} - \mathbf{w}\|}\right) \right\|, \\ &= \left| 1 - \frac{r + \frac{\varepsilon}{2}}{\|\mathbf{x} - \mathbf{w}\|} \right| \|\mathbf{x} - \mathbf{w}\|, \\ &= \left| \|\mathbf{x} - \mathbf{w}\| - r - \frac{\varepsilon}{2} \right|. \end{aligned} \quad (23)$$

With \mathbf{w}_* being the closest prototype with a class label other than $c(\mathbf{w}^*)$ and ε less than or equal to $\|\mathbf{x} - \mathbf{w}_*\|$, it follows that

$$\begin{aligned} \|\mathbf{x} - \mathbf{w}_\diamond\| - r - \frac{\varepsilon}{2} &\geq \|\mathbf{x} - \mathbf{w}_\diamond\| - \frac{1}{2} (\|\mathbf{x} - \mathbf{w}_*\| - \|\mathbf{x} - \mathbf{w}^*\|) - \frac{1}{2} \|\mathbf{x} - \mathbf{w}_*\|, \\ &\geq \|\mathbf{x} - \mathbf{w}_\diamond\| - \|\mathbf{x} - \mathbf{w}_*\| + \frac{1}{2} \|\mathbf{x} - \mathbf{w}^*\|, \end{aligned}$$

and, hence, we obtain

$$\|\mathbf{x} - \mathbf{w}_\diamond\| - r - \frac{\varepsilon}{2} \geq 0.$$

This implies that the argument of the absolute value function of Equation (23) is always positive so that

$$\|\mathbf{x} - \hat{\mathbf{w}}(\mathbf{w})\| = \|\mathbf{x} - \mathbf{w}\| - r - \frac{\varepsilon}{2} \quad (24)$$

for all prototypes with a different class label than $c(\mathbf{w}^*)$.

Using Equation (24) for the shifted prototype $\hat{\mathbf{w}}(\mathbf{w}_*)$, we get

$$\begin{aligned} \|\mathbf{x} - \hat{\mathbf{w}}(\mathbf{w}_*)\| &= \|\mathbf{x} - \mathbf{w}_*\| - \frac{1}{2} (\|\mathbf{x} - \mathbf{w}_*\| - \|\mathbf{x} - \mathbf{w}^*\|) - \frac{\varepsilon}{2}, \\ &= \frac{1}{2} (\|\mathbf{x} - \mathbf{w}_*\| + \|\mathbf{x} - \mathbf{w}^*\|) - \frac{\varepsilon}{2}. \end{aligned} \quad (25)$$

Similarly, using Equation (20) for the shifted prototype $\hat{\mathbf{w}}(\mathbf{w}^*)$, we obtain

$$\begin{aligned} \|\mathbf{x} - \hat{\mathbf{w}}(\mathbf{w}^*)\| &= \|\mathbf{x} - \mathbf{w}^*\| + \frac{1}{2} (\|\mathbf{x} - \mathbf{w}_*\| - \|\mathbf{x} - \mathbf{w}^*\|) + \frac{\varepsilon}{2}, \\ &= \frac{1}{2} (\|\mathbf{x} - \mathbf{w}_*\| + \|\mathbf{x} - \mathbf{w}^*\|) + \frac{\varepsilon}{2}. \end{aligned} \quad (26)$$

Comparing Equation (25) with Equation (26) and with ε being positive, we conclude that

$$\|\mathbf{x} - \hat{\mathbf{w}}(\mathbf{w}_*)\| < \|\mathbf{x} - \hat{\mathbf{w}}(\mathbf{w}^*)\|$$

holds. Since this implies that $\hat{\mathbf{w}}(\mathbf{w}_*)$ is the *closest* prototype in the set $\hat{\mathcal{W}}$, the set $\hat{\mathcal{W}}$ of shifted prototypes assigns the class label of $c(\mathbf{w}_*)$ to \mathbf{x} and thus labels \mathbf{x} other than \mathcal{W} —in symbols, $c_{\hat{\mathcal{W}}}^*(\mathbf{x}) \neq c_{\mathcal{W}}^*(\mathbf{x})$. Therefore, the hypothesis margin must be less than $r + \varepsilon$ because we constructed $\hat{\mathcal{W}}$ by means of prototype shifts with a magnitude of size $r + \frac{\varepsilon}{2}$. With $\varepsilon > 0$ arbitrarily chosen, it follows that *the hypothesis margin of \mathcal{W} is less than or equal to r .*

Now, combining the final result of Step 1 with the final result of Step 2, we obtain that the hypothesis margin is *exactly* r . \square

C Proof of Theorem 2: Relation to the sample margin

Theorem. *Let the data space \mathcal{X} be a vector space over a field of the real or complex numbers, $d(\mathbf{x}, \mathbf{w}) = \|\mathbf{x} - \mathbf{w}\|$ be a dissimilarity induced by a seminorm $\|\cdot\|$, and \mathcal{S} be a set of inputs. Then, the hypothesis margin of \mathcal{W} with respect to \mathcal{S} yields a lower bound on the sample margin of \mathcal{W} with respect to \mathcal{S} :*

$$\text{margin}_h(\mathcal{S}, \mathcal{W}) \leq \text{margin}_s(\mathcal{S}, \mathcal{W}).$$

The proof is based on the ideas used by Crammer et al. [56].

Proof. We prove the theorem by the following steps:

1. Given a set \mathcal{W} of prototypes and an arbitrary radius r that fulfills the requirement

$$\text{margin}_s(\mathcal{S}, \mathcal{W}) < r, \quad (27)$$

we define a set $\hat{\mathcal{W}}$ of shifted prototypes such that each prototype is shifted by a vector of length r .

2. We prove that this set $\hat{\mathcal{W}}$ of prototypes labels the inputs of \mathcal{S} differently than \mathcal{W} so that the hypothesis margin must be less than or equal to r . Because the radius r was arbitrarily chosen, it follows that the hypothesis margin must be less than or equal to the sample margin, which proves the theorem.

Step 1: Let r be an arbitrary radius such that Equation (27) holds. Because the sample margin $\text{margin}_s(\mathcal{S}, \mathcal{W})$ is less than r , there exists an element \mathbf{x} of \mathcal{S} such that there exists another element $\bar{\mathbf{x}}$ of \mathcal{X} with the distance

$$d(\mathbf{x}, \bar{\mathbf{x}}) = \|\mathbf{x} - \bar{\mathbf{x}}\| = r \quad (28)$$

that has an assigned class label different than the class label assigned to \mathbf{x} —in symbols, $c_{\mathcal{W}}^*(\mathbf{x}) \neq c_{\mathcal{W}}^*(\bar{\mathbf{x}})$. Because \mathbf{x} and $\bar{\mathbf{x}}$ are labeled differently, the closest prototypes must have different class labels too—we denote by \mathbf{w}^* the closest prototype to \mathbf{x} and by $\bar{\mathbf{w}}^*$ the closest prototype to $\bar{\mathbf{x}}$.

Based on \mathbf{x} and $\bar{\mathbf{x}}$, we define the set $\hat{\mathcal{W}}$ of shifted prototypes of \mathcal{W} by shifting each prototype $\mathbf{w} \in \mathcal{W}$ to the position

$$\hat{\mathbf{w}}(\mathbf{w}) = \mathbf{w} + \mathbf{x} - \bar{\mathbf{x}}. \quad (29)$$

Due to Equation (28), the magnitude of the applied shift to a prototype is exactly r :

$$d(\mathbf{w}, \hat{\mathbf{w}}(\mathbf{w})) = \|\mathbf{w} - \mathbf{w} - \mathbf{x} + \bar{\mathbf{x}}\| = r. \quad (30)$$

Step 2: By Equation (29), it follows that $d(\bar{\mathbf{x}}, \mathbf{w})$ equals $d(\mathbf{x}, \hat{\mathbf{w}}(\mathbf{w}))$:

$$d(\bar{\mathbf{x}}, \mathbf{w}) = \|\underbrace{\bar{\mathbf{x}} - \mathbf{x} + \mathbf{x} - \mathbf{w}}_{=0}\| = \|\mathbf{x} - (\mathbf{w} + \mathbf{x} - \bar{\mathbf{x}})\| = d(\mathbf{x}, \hat{\mathbf{w}}(\mathbf{w})). \quad (31)$$

Because $\bar{\mathbf{w}}^*$ is the closest prototype to $\bar{\mathbf{x}}$, the dissimilarity of $\bar{\mathbf{x}}$ to an arbitrary prototype $\mathbf{w}_\diamond \in \mathcal{W}$ of a class *other* than $c(\bar{\mathbf{w}}^*)$ is larger:

$$d(\bar{\mathbf{x}}, \bar{\mathbf{w}}^*) < d(\bar{\mathbf{x}}, \mathbf{w}_\diamond). \quad (32)$$

Combining Equation (32) and Equation (31), we conclude that

$$d(\bar{\mathbf{x}}, \bar{\mathbf{w}}^*) = d(\mathbf{x}, \hat{\mathbf{w}}(\bar{\mathbf{w}}^*)) < d(\mathbf{x}, \hat{\mathbf{w}}(\mathbf{w}_\diamond)) = d(\bar{\mathbf{x}}, \mathbf{w}_\diamond),$$

which implies that $\hat{\mathbf{w}}(\bar{\mathbf{w}}^*)$ must be the closest prototype to \mathbf{x} regarding the set $\hat{\mathcal{W}}$ of prototypes. Therefore, the set $\hat{\mathcal{W}}$ of prototypes assigns a *different* class label than \mathcal{W} to \mathbf{x} —in symbols, $c_{\hat{\mathcal{W}}}^*(\mathbf{x}) \neq c_{\mathcal{W}}^*(\mathbf{x})$. Because all prototypes $\mathbf{w} \in \mathcal{W}$ have been shifted by the magnitude r , see Equation (30), the hypothesis margin must be less than or equal to r . Additionally, since this result holds for an arbitrary r according to Equation (27), the hypothesis margin must be less than or equal to the sample margin. \square

Corollary. *Given a labeled data point $(\mathbf{x}, c(\mathbf{x}))$ that is correctly classified by an NPC according to Theorem 2 and a corresponding adversarial perturbation δ that changes the assigned class label. Then, the following inequality is true and has tight bounds:*

$$\text{margin}_h(\{\mathbf{x}\}, \mathcal{W}) \leq \text{margin}_s(\{\mathbf{x}\}, \mathcal{W}) \leq \|\delta\|.$$

Proof. The inequality immediately follows from Theorem 2, Definition 1 of the sample margin, and the requirement that an adversarial perturbation changes the assigned class label. Moreover, by Definition 1 of the sample margin, the closest adversarial example is an element at a distance $\text{margin}_s(\{\mathbf{x}\}, \mathcal{W})$ to \mathbf{x} . Therefore, the minimum of Equation (6) is exactly $\text{margin}_s(\{\mathbf{x}\}, \mathcal{W})$ —which implies the tightness of the first lower bound. The tightness of the lower bound

$$\text{margin}_h(\{\mathbf{x}\}, \mathcal{W}) \leq \text{margin}_s(\{\mathbf{x}\}, \mathcal{W})$$

is proven by showing that there exists an element $\mathbf{x} \in \mathcal{X}$ such that the hypothesis margin is greater than or equal to the sample margin. Together with Equation (7) of the theorem, this implies equality and the tightness of the bound.

If \mathbf{x} is set to be equal to an arbitrary prototype $\mathbf{w} \in \mathcal{W}$, then the closest prototype \mathbf{w}^* to \mathbf{x} is \mathbf{w} . Additionally, we denote by \mathbf{w}_* the closest prototype with a different class label than $c(\mathbf{w}^*)$. The decision boundary between these two prototypes is defined by all elements $\mathbf{x}' \in \mathcal{X}$ that fulfill the following criterion:

$$\|\mathbf{x}' - \mathbf{w}^*\| = \|\mathbf{x}' - \mathbf{w}_*\|.$$

The vector $\bar{\mathbf{x}}$ defined as

$$\begin{aligned} \bar{\mathbf{x}} &= \mathbf{x} + \frac{1}{2}(\mathbf{w}_* - \mathbf{w}^*), \\ &= \mathbf{w}^* + \frac{1}{2}(\mathbf{w}_* - \mathbf{w}^*), \\ &= \frac{1}{2}(\mathbf{w}_* + \mathbf{w}^*) \end{aligned} \quad (33)$$

satisfies this decision-boundary criterion. Hence, the sample margin must be less than or equal to the dissimilarity from \mathbf{x} to the element $\bar{\mathbf{x}}$ from Equation (33):

$$\text{margin}_s(\{\mathbf{x}\}, \mathcal{W}) \leq d(\mathbf{x}, \bar{\mathbf{x}}) = \|\mathbf{x} - \bar{\mathbf{x}}\|. \quad (34)$$

On the other hand, the seminorm $\|\mathbf{x} - \bar{\mathbf{x}}\|$ is equal to the hypothesis margin, see Equation (5):

$$\begin{aligned} \|\mathbf{x} - \bar{\mathbf{x}}\| &= \left\| \mathbf{w}^* - \frac{1}{2}(\mathbf{w}_* + \mathbf{w}^*) \right\|, \\ &= \frac{1}{2} \|\mathbf{w}^* - \mathbf{w}_*\|, \\ &= \frac{1}{2} (\|\mathbf{w}^* - \mathbf{w}_*\| - \|\mathbf{w}^* - \mathbf{w}^*\|), \\ &= \frac{1}{2} (\|\mathbf{x} - \mathbf{w}_*\| - \|\mathbf{x} - \mathbf{w}^*\|), \\ &= \text{margin}_h(\{\mathbf{x}\}, \mathcal{W}). \end{aligned} \quad (35)$$

Combining Equation (34) and Equation (35), we conclude $\text{margin}_s(\{\mathbf{x}\}, \mathcal{W}) \leq \text{margin}_h(\{\mathbf{x}\}, \mathcal{W})$ and together with Equation (7) the equality:

$$\text{margin}_s(\{\mathbf{x}\}, \mathcal{W}) = \text{margin}_h(\{\mathbf{x}\}, \mathcal{W}).$$

This proves the tightness of the second lower bound and, therefore, the corollary. \square

D Experimental setup

This section presents all the required information for reproducing the results of Section 5. Additionally, we give a detailed explanation of how we selected the state-of-the-art methods.

D.1 Configuration, training, and evaluation of the selected NPCs

Model selection In the experimental evaluation in Section 5, we considered three different NPCs: GLVQ, GTLVQ, and RSLVQ. All methods belong to the family of LVQ algorithms. We focused on these methods because they optimize the prototypes as fully adjustable parameters and not merely select the prototypes from a given set of training points (e. g., k-nearest neighbors methods). By having access to the full data space and not only the data samples, the methods usually achieve better classification performances with fewer prototypes than prototype *selection* approaches. Below, we give a brief rationale for why we selected the specific LVQ approaches for the presented comparison.

GLVQ From the large LVQ family, GLVQ is by far the most commonly used variant. Therefore, the analysis of GLVQ provides results that are important for many applications. Moreover, the GLVQ loss is the standard loss function to train other LVQ variants.

GTLVQ Compared to NNs, GTLVQ is the closest in terms of CTE and was shown to be robust against adversarial attacks before [26]. Besides, GTLVQ generalizes LVQ to the case of an infinite number of prototypes.

RSLVQ As RSLVQ is not trained using a triplet loss and violates the seminorm assumption, the theorems discussed in Section 4 suggest that adversarial robustness cannot be guaranteed. Including RSLVQ in the comparison, allows for validation of this statement. Additionally, by being trained using the cross-entropy loss, RSLVQ is a variant of LVQ that is similar to NNs.

Initialization The prototypes of the GLVQ and RSLVQ networks were initialized by computing class-wise a k-means, where the number of means was equal to the number of prototypes in the respective class. For GTLVQ, we applied the following standard initialization strategy: The translation vectors are initialized by the same method as used for GLVQ. After that, we initialized each basis \mathbf{B}_k using the following procedure:

1. Determine all the training samples of the correct class for which \mathbf{t}_k is the closest prototype vector in terms of the Euclidean distance. Hence, we consider \mathbf{t}_k as a prototype vector of an ordinary LVQ approach and determine all the training samples that belong to the receptive field of \mathbf{t}_k .

Table 3: Configuration of the NPCs used in the evaluation. For GTLVQ, we report the number of prototypes per class and the subspace dimension. Moreover, the ReLU loss for the GTLVQ model refers to the loss of Equation (12).

Norm	Model	MNIST		CIFAR-10	
		Loss	Prototypes	Loss	Prototypes
L^∞	GLVQ	GLVQ	128 ppc	GLVQ	64 ppc
–	RSLVQ	cross-entropy	128 ppc	cross-entropy	128 ppc
L^2	GLVQ	GLVQ	256 ppc	GLVQ	128 ppc
	GTLVQ	ReLU ($\epsilon = 1.58$)	10 ppc, $m = 12$	GLVQ	1 ppc, $m = 100$

2. Compute the m eigenvectors that belong to the m largest eigenvalues of the estimated covariance matrix over these training samples.
3. Use these m eigenvectors as initialization for \mathbf{B}_k and orthonormalize the resulting matrix if necessary.

Training We trained each method by optimizing the reported loss function. In the case of GLVQ with the L^2 -norm and GTLVQ, we used the loss function with squared Euclidean distances. This avoids the computation of the square root. As discussed at the end of Section 4, this still optimizes for adversarially robust models.

Each method was trained for 1000 epochs without early stopping. The optimizer was Adam [72], with the default setting of the KERAS framework, a batch size of 128, and an initial learning rate of 0.001. During training, we monitored the validation loss and automatically adjusted the learning rate accordingly. If the validation loss did not decrease over 10 epochs, we reduced the learning rate by a factor of 0.9.

The datasets were normalized to the unit interval. During training, we applied basic data augmentations in the form of random shifts of up to ± 2 pixels and random rotations of up to ± 15 degrees.

Selected hyperparameters In Table 3, the hyperparameter settings for the NPCs are presented. To determine the number of prototypes for GLVQ and RSLVQ, we performed a grid search where we tested the following number of prototypes per class: 1, 2, 4, 8, 16, 32, 64, 128, and 256. The final models have been selected by the smallest URTE for GLVQ and the smallest CTE for RSLVQ. For GTLVQ trained on MNIST, we used the configuration reported by Saralajew et al. [26]. In contrast, the CIFAR-10 model of GTLVQ was selected by the following strategy: We defined the number of prototypes per class to be one and varied the subspace dimension from 1 to 256. For each configuration, we initialized the model and calculated the training accuracy. The final model was selected as the model where the usage of a higher subspace dimension showed no improvement of the accuracy.

Additionally, we considered both the ReLU loss function, see Equation (12), and the GLVQ loss function, see Equation 2, for GLVQ and GTLVQ, while we only considered the cross-entropy loss for RSLVQ. The free parameter of the ReLU loss was always set to be equal to the ϵ of the ϵ -limited attacks against which the model should be robust.

Configuration of the PGD attack In line with earlier work [33, 39], we used the PGD attack [41] to evaluate the empirical robustness of NPCs. We ran the PGD attack for 200 iterations and used random starts with Gaussian noise. For each sample, the worst-case adversary was selected from three starts. We found no evidence that performing more than three random starts had a significant effect on the reported LRTE.

Hardware and software frameworks used The implementation of the adversarial attack was supplied by the Python FOOLBOX⁵ library (version 2.4.0). All models were implemented using the

⁵<https://foolbox.readthedocs.io/en/v2.4.0/>

Python KERAS⁶ library (version 2.2.4) with TENSORFLOW⁷ back end (version 1.12.0). Evaluation and training were performed on an NVIDIA Tesla V100 32 GB GPU. However, the inference time of each model was obtained using an NVIDIA RTX 2080 Ti GPU. By using a commonly used GPU for this purpose, we hope to make a simple comparison possible.

D.2 Selection criteria for the state-of-the-art methods

Several important considerations were made when selecting the state-of-the-art certification and verification methods for the comparison. First of all, we wanted to make a comparison with the same number of verification and certification methods. Both approaches play an important role in the research of guaranteed adversarial robustness and therefore require equal consideration. However, as we explicitly focus on fast computation of the guaranteed robustness, a stronger focus was placed on certification methods when considering L^2 -norm limited attacks. The STN and Smooth methods for certification were selected based on their current state-of-the-art status among certification approaches. As already mentioned, these methods are not deterministic and require extensive sampling. Therefore, the derived certificates differ strongly from the proposed adversarial robustness certificates of NPCs. For this reason, CAP was included as a method more similar to the certification of NPCs. Both IBP and RS were chosen as verification methods because of their focus on fast verification and robustification. RT was chosen to add a second example of guaranteed adversarial robustness outside of NNs in addition to NPCs.

For the comparison of the results, we decided to present for each method the model with the best guaranteed robustness. We chose not to consider ensemble or cascade models (apart from RT) because these models accept a computational overhead for better performances. Considering our scope of fast adversarial robustness certification, we deemed this inappropriate. We presented the CTE, LRTE (if available), and URTE for each method, as reported in the respective papers. This was made possible by the widespread use of the PGD attack as a measure of empirical robustness.

E Additional experimental results for Section 5

In the following, we present additional experimental results for Section 5, including an evaluation of the presented NPCs for other ϵ -limited adversarial attacks, an analysis of the NPCs robustness regarding several norm-based adversarial attacks, and an extended discussion of the presented adversarial rejection strategy results.

E.1 Extended robustness evaluation for Section 5

In this section, we extend the robustness evaluation presented in the comparison of Section 5. For the L^∞ -norm, we present results for additional threshold values ϵ , see Table 4. The models were trained and selected like the models presented in the main results. We also include the verification and LRTE results for CAP provided by Gowal et al. [33] for the evaluation of IBP. These results are denoted as CAP-IBP. Furthermore, we included the results of evaluating a 1-nearest neighbor classifier (denoted as 1-NN) by Wang et al. [58], see L^∞ -norm with ϵ equal to 0.1.⁸ For the L^2 -norm, we extend the main comparison with RSLVQ, see Table 5.

In contrast to adversarial examples under the L^∞ -norm, no real standardized attack for evaluating the empirical robustness under the L^2 -norm is available. For the L^∞ -norm, the PGD attack plays this role. To provide some insights into the empirical robustness of NPCs trained to optimize L^2 -norm adversarial robustness, we present in Table 5 the results of an evaluation using the Carlini & Wagner (C&W) attack [14]. There, we present the LRTE obtained with the C&W attack for the GLVQ, RSLVQ, and GTLVQ models and compare it with the reported results of the STN certification method.

We used the C&W attack implementation of the Python FOOLBOX library (version 2.4.0). To determine the trade-off parameter—determining the trade-off between misclassification and perturbation distance in the C&W attack algorithm—a binary search with 10 steps was performed. The attack

⁶<https://www.keras.io/>

⁷<https://www.tensorflow.org/>

⁸Note that the LRTE is not obtained by the PGD attack.

Table 4: Comparison of NPCs trained with the L^∞ -norm against state-of-the-art methods. Dashes “-” indicate that the quantity is not calculable or reported.

Dataset	ϵ	Class	Model	CTE [%]	LRTE [%]	URTE [%]	Notes
MNIST	0.1	Certify	GLVQ	3.66	6.42	6.67	GLVQ loss, 128 ppc
			RSLVQ	1.70	93.97	-	cross-entropy, 128 ppc
			1-NN	3.41	27.06	27.06	[58, Table 3]
		Verify	CAP	1.08	-	3.67	[12, Table 2 “Large”]
			CAP-IBP	1.08	2.89	3.01	[33, Table 4]
			RS	1.32	4.87	5.66	[39, Table 3 “RS+”]
	0.2	Certify	GLVQ	3.66	10.63	11.53	GLVQ loss, 128 ppc
			RSLVQ	1.70	100.00	-	cross-entropy, 128 ppc
		Verify	CAP-IBP	3.22	6.93	7.27	[33, Table 4]
			RS	1.90	6.86	10.21	[39, Table 3 “RS+”]
			IBP	1.66	3.90	4.48	[33, Table 4]
		0.3	Certify	GLVQ	3.66	16.39	20.58
	RSLVQ			1.70	100.00	-	cross-entropy, 128 ppc
	RT			2.68	12.46	12.46	[50, Table 3]
	Verify		CAP	14.87	-	43.10	[12, Table 2 “Small”]
CAP-IBP			13.52	26.16	26.92	[33, Table 4]	
RS			2.67	7.95	19.32	[39, Table 3 “RS+”]	
CIFAR-10	$2/255$	Certify	GLVQ	59.35	65.50	65.52	GLVQ loss, 64 ppc
			RSLVQ	54.71	89.54	-	cross-entropy, 128 ppc
			CAP	31.72	-	46.11	[12, Table 2 “Resnet”]
		Verify	CAP-IBP	36.01	45.11	49.96	[33, Table 4]
			RS	38.88	51.08	54.07	[39, Table 3 “RS+”]
			IBP	29.84	45.09	49.98	[33, Table 4]
	$8/255$	Certify	GLVQ	59.35	79.54	79.62	GLVQ loss, 64 ppc
			RSLVQ	54.71	99.04	-	cross-entropy, 128 ppc
			RT	58.46	74.69	74.69	[50, Table 3]
		Verify	CAP	71.33	-	78.22	[12, Table 2 “Resnet”]
			CAP-IBP	71.03	78.14	79.21	[33, Table 4]
			RS	59.55	73.22	79.73	[39, Table 3 “RS+”]
IBP	50.51	65.23	67.96	[33, Table 4]			

used the Adam optimizer with a learning rate of 0.05 during the 1000 optimization steps, except for RSLVQ, which ran for 3000 steps.

Results L^∞ -norm With respect to the additional ϵ values, the comparison sways in favor of the NPCs for larger ϵ values as can be seen in Table 4 for the L^∞ -norm: For $\epsilon = 0.1$ and $\epsilon = 0.2$ on MNIST, the URTE of GLVQ is higher than the URTE of the certification method CAP and is significantly higher than the robustness guarantees of the verification methods. This is no longer the case for $\epsilon = 0.3$, GLVQ now improves over CAP and is closer to verification methods with regard to URTE. Since large ϵ -limited attacks are a more realistic attack model, we deem this to be an advantageous behavior. For the CIFAR-10 dataset, the same trend is observable—regarding URTE, GLVQ is closer to CAP and the verification methods for $\epsilon = 8/255$ than for $\epsilon = 2/255$.

The same holds for CAP-IBP concerning both the URTE and LRTE. Most notably, with $\epsilon = 0.3$ on MNIST, GLVQ outperforms CAP-IBP by a large margin in terms of LRTE. Despite being verified using an exact method, this model was still trained to optimize certification. Hence, we can conclude that in terms of LRTE, NPCs provide exceptional empirical adversarial robustness compared to other certification methods.

Table 5: Comparison of NPCs trained with the L^2 -norm against state-of-the-art methods. Dashes “–” indicate that the quantity is not calculable or reported. Values denoted with * were estimated from figures from the original publication.

Dataset	ϵ	Model	CTE [%]	LRTE [%]	URTE [%]	Notes
MNIST	1.58	GLVQ	4.19	32.42	65.61	GLVQ loss, 256 ppc
		GTLVQ	2.92	25.86	55.32	ReLU loss, 10 ppc, $m = 12$
		RSLVQ	1.70	84.26	–	cross-entropy, 128 ppc
		CAP	11.88	–	55.47	[12, Table 4 “Large”]
		STN	1.10	10*	31.00	[49, Table 1]
CIFAR-10	$\frac{3}{255}$	GLVQ	51.41	58.40	61.90	GLVQ loss, 128 ppc
		GTLVQ	40.53	49.02	55.96	GLVQ loss, 1 ppc, $m = 100$
		RSLVQ	54.71	75.31	–	cross-entropy, 128 ppc
		CAP	38.80	–	48.04	[12, Table 2 “Resnet”]
		STN	19.50	30*	34.40	[49, Table 1]
		Smooth	18*	–	27*	[18, Figure 5c]

Table 6: Certified robustness for the L^∞ -norm of NPCs trained with the L^2 -norm.

Dataset	ϵ	Model	CTE [%]	LRTE [%]	URTE [%]	Notes
MNIST	0.1	GLVQ (L^2)	4.19	18.45	97.15	GLVQ loss, 256 ppc
		GTLVQ	2.92	13.98	100.00	ReLU loss, 10 ppc, $m = 12$
	0.3	GLVQ (L^2)	4.19	89.47	100.00	GLVQ loss, 256 ppc
		GTLVQ	2.92	99.2	100.00	ReLU loss, 10 ppc, $m = 12$
CIFAR-10	$\frac{2}{255}$	GLVQ (L^2)	51.41	64.45	74.46	GLVQ loss, 128 ppc
		GTLVQ	40.53	59.69	78.10	GLVQ loss, 1 ppc, $m = 100$
	$\frac{3}{255}$	GLVQ (L^2)	51.41	84.21	95.96	GLVQ loss, 128 ppc
		GTLVQ	40.53	90.2	99.61	GLVQ loss, 1 ppc, $m = 100$

In comparison to the 1-NN classifier results, which is also an NPC, GLVQ has a much better adversarial robustness even though the CTE of GLVQ is slightly worse. However, it must be noted that the 1-NN method was not purposefully trained to be adversarial robust.













































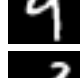





Results L^2 -norm Regarding RSLVQ and the L^2 -norm, we find the same results as presented in the paper for the L^∞ -norm: RSLVQ is not robust against adversarial attacks, see Table 5. However, the LRTE score is not trivial. On the other hand, the NPCs that optimize a triplet loss—GLVQ and GTLVQ—are robust against adversarial examples generated by the C&W attack. It must, however, be noted that the LRTE of NPCs is significantly higher than for state-of-the-art methods.

E.2 Multiple seminorm robustness: Evaluation of the L^∞ -norm robustness for L^2 -norm robustified models

As discussed in Section 6, we can use Hölder’s inequality to extend the robustness guarantees between different L^p -norms. In this section, we present some preliminary results demonstrating this. Given the NPCs trained using the L^2 -norm, as presented in Section 5 and Section E.1, Table 6 presents their adversarial robustness with regard to the L^∞ -norm, both empirical and guaranteed.

Results The NPCs trained to classify the MNIST dataset have trivial guaranteed robustness and poor empirical robustness. However, if we consider the CIFAR-10 dataset, we find different results. For small values of ϵ , both the empirical and guaranteed robustness against adversarial examples are nontrivial. Interestingly, GLVQ and GTLVQ have a better LRTE with respect to the L^∞ -norm when trained using the L^2 -norm than GLVQ has when trained using the L^∞ -norm directly—as can be seen by comparing the LRTE of GLVQ on CIFAR-10 in Table 4 with the LRTE of GLVQ and GTLVQ on CIFAR-10 in Table 6 (the latter has lower LRTE than the former).

Table 7: The first 10 falsely rejected samples from the MNIST test dataset.

id	$c(\mathbf{x})$	$c_{\mathcal{W}}^*(\mathbf{x}) = c(\mathbf{w}^*)$	$c(\mathbf{w}_*)$	\mathbf{x}	\mathbf{w}^*	\mathbf{x}^*	\mathbf{w}_*	\mathbf{x}_*
63	3	2	3					
92	9	9	4					
115	4	6	9					
124	7	4	7					
125	9	9	4					
149	2	2	9					
151	9	9	2					
193	9	9	4					
195	3	3	9					
219	5	5	3					

E.3 Extended adversarial rejection evaluation for Section 5

In Table 7, we show the first 10 falsely rejected samples from the MNIST test dataset by the adversarial rejection strategy presented in Section 5. For each sample, we visualized the following images: the sample \mathbf{x} , the closest prototype \mathbf{w}^* , a closest representation \mathbf{x}^* of the closest prototype \mathbf{w}^* in the training dataset, the closest prototype \mathbf{w}_* from another class, and a closest representation \mathbf{x}_* of the closest prototype \mathbf{w}_* from another class in the training dataset. Additionally, we present the following information: the sample identifier (id) in the MNIST test dataset, the true label $c(\mathbf{x})$ of the sample \mathbf{x} , the predicted class label $c_{\mathcal{W}}^*(\mathbf{x})$ by GLVQ, and the class label $c(\mathbf{w}_*)$ of the prototype \mathbf{w}_* . Note that the predicted class label $c_{\mathcal{W}}^*(\mathbf{x})$ is equal to $c(\mathbf{w}^*)$ and that the samples of class $c(\mathbf{w}_*)$ are the closest samples that can alter the prediction $c_{\mathcal{W}}^*(\mathbf{x})$. Consequently, if we consider \mathbf{x} as a potentially adversarially manipulated sample, the class $c(\mathbf{w}_*)$ is the class that contains the original sample (the non-adversarially manipulated sample) with the highest probability.

We used a rejection threshold of 0.1. Therefore, we rejected all samples with a hypothesis margin of less than 0.1. The rejection threshold reflects the amount of perturbation that can be expected to be used by a malicious attacker without being detectable by other means. According to the original definition of adversarial examples (being perturbed by an imperceptible amount of noise to humans), the threshold can also be considered as the amount of noise that would be *imperceptible* to the human eye. The threshold of 0.1 equates to roughly 6% of the examples from the MNIST test dataset being falsely rejected.

Results We can make several distinct considerations within the set of falsely rejected samples. First of all, several samples are rejected despite being classified correctly by the NPC. Sample 92, for example, is classified correctly as a nine and is still rejected. Second, we find some samples that are

Table 8: Comparison of NPCs trained with the L^∞ -norm against robust boosted decision trees and robust boosted decision stumps on tabular data. For Stumps-A, we took the versions trained with the exact robust loss. The URTEs for Stumps-C and RT-C are robust test errors and therefore the best possible URTEs.

Dataset	ϵ	Model	CTE [%]	URTE [%]	Notes
breast-cancer	0.3	Stumps-A	5.1	10.9	[50, Table 1]
		Stumps-C	8.8	16.8	
		RT	0.7	6.6	[50, Table 2]
		RT-C	0.7	13.1	
		GLVQ	0.0	7.3	
diabetes	0.05	Stumps-A	27.3	31.8	[50, Table 1]
		Stumps-C	23.4	30.5	
		RT	27.3	35.7	[50, Table 2]
		RT-C	22.1	40.3	
		GLVQ	25.3	31.8	
cod-rna	0.025	Stumps-A	11.2	22.6	[50, Table 1]
		Stumps-C	11.6	23.2	
		RT	6.9	21.4	[50, Table 2]
		RT-C	10.2	24.2	
		GLVQ	7.8	21.4	

misclassified and rejected. For obvious reason, these false rejections are not as severe as the rejection of correctly classified samples. Within this category, a further division can be made between samples for which the class $c(\mathbf{w}_*)$ is the ground truth class (e. g., sample 63) and samples for which this is not the case (e. g., sample 115).

F Robustness evaluation on tabular data

Additionally to the results on image datasets, we present a robustness evaluation on tabular data and compare it with state-of-the-art robust boosted tree and stump methods. The datasets and the threat models (ϵ parameters) were chosen in accordance with the evaluation of robust boosted tree-based methods [50]. In particular, we used the source code provided by Andriushchenko and Hein [50] to load, preprocess, and split the following datasets (all available in the LIBSVM library [73]):

breast-cancer a binary classification dataset (also denoted as Breast Cancer Wisconsin (Original) dataset [74]) consisting of 546 training and 137 test samples with 10 features after removing samples with missing values [75];

diabetes a binary classification dataset (also denoted as Pima Indians Diabetes dataset) consisting of 614 training and 154 test samples with 8 features [76];

cod-rna a binary classification task consisting of 59535 training and 271617 test samples with 8 features [77].

Compared to the image datasets that are used in the main part of the paper, these datasets are low-dimensional and do not provide (except for cod-rna) a large number of training samples.

For each dataset, we trained a GLVQ classifier with the L^∞ -norm with the setting specified in Table 8. These settings are the result of a hyperparameter optimization by grid search on the specified parameters. With the respective setting, each model was trained several times. From all training runs, the model with the best URTE was selected. In general, the models were trained by the following setting:

- initialization by applying class-wise a k-means algorithm;
- optimizing the models by ADAM with a learning rate scheduler (similar to the one used for the image datasets) and without augmentation;
- stopping the training after 1000 training epochs.

In addition to RT models (see Section 5), we compare the achieved performances with:

- robust boosted decision stumps (denoted by Stumps-A) of Andriushchenko and Hein [50],
- robust boosted decision trees (denoted by RT-C) of Chen et al. [51], and
- robust boosted decision stumps (denoted by Stumps-C) of Chen et al. [51].

For all models, we used the results reported by Andriushchenko and Hein [50]. All models, except for GLVQ, are ensembles.

In Table 8, we report the results of the experiments. Summarily, there is no superior method. For example, on the breast-cancer dataset, the RT achieves the best URTE, but GLVQ has a better accuracy and is just slightly worse in terms of URTE. On the other hand, the Stumps-C outperforms all other models on the diabetes dataset with respect to URTE, but RT and GLVQ are close behind it. Moreover, on the cod-rna dataset, the RT and the GLVQ classifier achieve the same and best URTE and outperform all other models. With no clear superior method for either guaranteed robustness or clean accuracy, GLVQ can therefore be considered as state of the art for guaranteed adversarial robustness on tabular datasets, together with the stump- and tree-based methods.