# Controlling Neural Level Sets

**Matan Atzmon, Niv Haim, Lior Yariv, Ofer Israelov, Haggai Maron, Yaron Lipman**
Weizmann Institute of Science
Rehovot, Israel

## Abstract

The level sets of neural networks represent fundamental properties such as decision boundaries of classifiers and are used to model non-linear manifold data such as curves and surfaces. Thus, methods for controlling the neural level sets could find many applications in machine learning.

In this paper we present a simple and scalable approach to directly control level sets of a deep neural network. Our method consists of two parts: (i) sampling of the neural level sets, and (ii) relating the samples' positions to the network parameters. The latter is achieved by a *sample network* that is constructed by adding a single fixed linear layer to the original network. In turn, the sample network can be used to incorporate the level set samples into a loss function of interest.

We have tested our method on three different learning tasks: improving generalization to unseen data, training networks robust to adversarial attacks, and curve and surface reconstruction from point clouds. For surface reconstruction, we produce high fidelity surfaces directly from raw 3D point clouds. When training small to medium networks to be robust to adversarial attacks we obtain robust accuracy comparable to state-of-the-art methods.

## 1 Introduction

The level sets of a Deep Neural Network (DNN) are known to capture important characteristics and properties of the network. A popular example is when the network $F(x;\theta) : \mathbb{R}^d \times \mathbb{R}^m \to \mathbb{R}^l$ represents a classifier, $\theta$ are its learnable parameters, $f_i(x;\theta)$ are its logits (the outputs of the final linear layer), and the level set

$$\mathcal{S}(\theta) = \left\{ x \in \mathbb{R}^d \,\middle|\, f_j - \max_{i \neq j} \{f_i\} = 0 \right\} \tag{1}$$

represents the decision boundary of the $j$-th class. Another recent example is modeling a manifold (e.g., a curve or a surface in $\mathbb{R}^3$) using a level set of a neural network (e.g., [26]). That is,

$$\mathcal{S}(\theta) = \left\{ x \in \mathbb{R}^d \,\middle|\, F = 0 \right\} \tag{2}$$

represents (generically) a manifold of dimension $d - l$ in $\mathbb{R}^d$.

The goal of this work is to provide practical means to directly control and manipulate *neural level sets* $\mathcal{S}(\theta)$, as exemplified in Equations 1, 2. The main challenge is how to incorporate $\mathcal{S}(\theta)$ in a differentiable loss function. Our key observation is that given a sample $p \in \mathcal{S}(\theta)$, its position can be associated to the network parameters: $p = p(\theta)$, in a *differentiable* and *scalable* way. In fact, $p(\theta)$ is itself a neural network that is obtained by an addition of a *single* linear layer to $F(x;\theta)$; we call these networks *sample networks*. Sample networks, together with an efficient mechanism for sampling the level set, $\{p_j(\theta)\} \subset \mathcal{S}(\theta)$, can be incorporated in general loss functions as a proxy for the level set $\mathcal{S}(\theta)$.

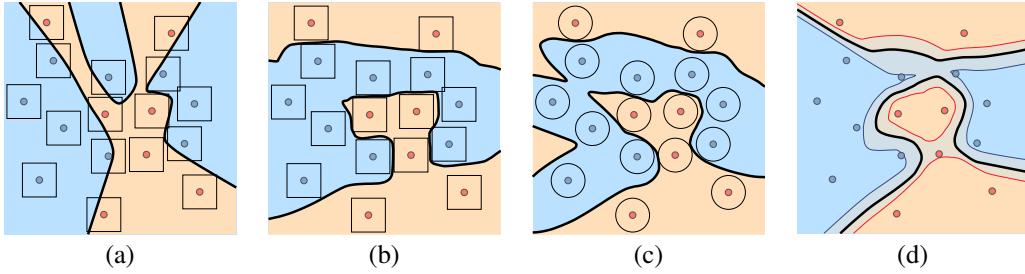|     |     |     |     |
| --- | --- | --- | --- |
| (a) | (b) | (c) | (d) |

Figure 1: Our method applied to binary classification in 2D. Blue and red dots represent positive and negative examples respectively. (a) standard cross-entropy loss baseline; (b) using our method to move the decision boundary at least $\varepsilon$ away from the training set in $L_\infty$ norm; (c) same for $L_2$ norm; (d) a geometrical adaptation of SVM soft-margin loss to neural networks using our method, the $+1, -1$ level sets are marked in light red and blue, respectively. Note that in (b),(c),(d) we achieve decision boundaries that seem to better explain the training examples compared to (a).

We apply our method of controlling the neural level sets to two seemingly different learning tasks: controlling decision boundaries of classifiers (Equation 1) and reconstructing curves and surfaces from point cloud data (Equation 2).

At first, we use our method to improve the generalization and adversarial robustness in classification tasks. In these tasks, the distance between the training examples and the decision boundary is an important quantity called the *margin*. Margins are traditionally desired to be as large as possible to improve generalization [7, 11] and adversarial robustness [11]. Usually, margins are only controlled indirectly by loss functions that measure network output values of training examples (e.g., cross entropy loss). Recently, researchers have been working on optimizations with more direct control of the margin using linearization techniques [14, 23, 11], regularization [28], output-layer margin [29], or using margin gradients [9]. We suggest controlling the margin by sampling the decision boundary, constructing the sample network, and measuring distances between samples and training examples directly. By applying this technique to train medium to small-size networks against adversarial perturbations we achieved comparable robust accuracy to state-of-the-art methods.

To improve generalization when learning from small amounts of data, we devise a novel geometrical formulation of the soft margin SVM loss to neural networks. This loss aims at directly increasing the *input space* margin, in contrast to standard deep network hinge losses that deal with output space margin [30, 29]. The authors of [11] also suggest to increase input space margin to improve generalization. Figure 1 shows 2D examples of training our adversarial robustness and geometric SVM losses for networks.

In a different application, we use our method for the reconstruction of manifolds such as curves and surfaces in $\mathbb{R}^3$ from point cloud data. The usage of neural networks for the modeling of surfaces has recently become popular [13, 31, 6, 2, 26, 24]. There are two main approaches: parametric and implicit. The parametric approach uses networks as parameterization functions to the manifolds [13, 31]. The implicit approach represents the manifold as a level set of a neural network [26, 6, 24]. So far, implicit representations were learned using regression, given a signed distance function or occupancy function computed directly from a ground truth surface. Unfortunately, for raw point clouds in $\mathbb{R}^3$, computing the signed distance function or an occupancy function is a notoriously difficult task [4]. In this paper we show that by using our sample network to control the neural level sets we can reconstruct curves and surfaces directly from point clouds in $\mathbb{R}^3$.

Lastly, to theoretically justify neural level sets for modeling manifolds or arbitrary decision boundaries, we prove a geometric version of the universality property of MLPs [8, 15]: any piecewise linear hyper-surface in $\mathbb{R}^d$ (i.e., a $d-1$ manifold built out of a finite number of linear pieces, not necessarily bounded) can be precisely represented as a neural level set of a suitable MLP.

## 2   Sample Network

Given a neural network $F(x; \theta) : \mathbb{R}^d \times \mathbb{R}^m \to \mathbb{R}^l$ its $0 \in \mathbb{R}^l$ level set is defined by

$$\mathcal{S}(\theta) := \left\{ x \mid F(x; \theta) = 0 \right\}. \tag{3}$$

We denote by $D_x F(p; \theta) \in \mathbb{R}^{l \times d}$ the matrix of partial derivatives of $F$ with respect to $x$. Assuming that $\theta$ is fixed, $F(p; \theta) = 0$ and that $D_x F(p; \theta)$ is of full rank $l$ ($l \ll d$), a corollary of the Implicit Function Theorem [17] implies that $\mathcal{S}(\theta)$ is a $d - l$ dimensional manifold in the vicinity of $p \in \mathcal{S}(\theta)$.

Our goal is to incorporate the neural level set $\mathcal{S}(\theta)$ in a differentiable loss. We accomplish that by performing the following procedure at each training iteration: (i) Sample $n$ points on the level set: $p_i \in \mathcal{S}(\theta)$, $i \in [n]$; (ii) Build the sample network $p_i(\theta)$, $i \in [n]$, by adding a fixed linear layer to the network $F(x; \theta)$; and (iii) Incorporate the sample network in a loss function as a proxy for $\mathcal{S}(\theta)$.

### 2.1   Level set sampling

To sample $\mathcal{S}(\theta)$ at some $\theta = \theta_0$, we start with a set of $n$ points $p_i$, $i \in [n]$ sampled from some probability measure in $\mathbb{R}^d$. Next, we perform generalized Newton iterations [3] to move each point $p$ towards $\mathcal{S}(\theta_0)$:

$$p^{\text{next}} = p - D_x F(p; \theta_0)^\dagger F(p; \theta_0), \tag{4}$$

where $D_x F(p, \theta_0)^\dagger$ is the Moore-Penrose pseudo-inverse of $D_x F(p, \theta_0)$. The generalized Newton step solves the under-determined ($l \ll d$) linear system $F(p; \theta_0) + D_x F(p; \theta_0)(p^{\text{next}} - p) = 0$. To motivate this particular solution choice we show that the generalized Newton step applied to a linear function is an orthogonal projection onto its zero level set (see proof in Appendix C):

**Lemma 1.** *Let $\ell(x) = Ax + b$, $A \in \mathbb{R}^{l \times d}$, $b \in \mathbb{R}^l$, $\ell < d$, and $A$ is of full rank $l$. Then Equation 4 applied to $F(x) = \ell(x)$ is an orthogonal projection on the zero level set of $\ell$, namely, on $\{x \mid \ell(x) = 0\}$.*

For $l > 1$ the computation of $D_x F(p; \theta_0)^\dagger$ requires inverting an $l \times l$ matrix; in this paper $l \in \{1, 2\}$. The directions $D_x F(p_i; \theta_0)$ can be computed efficiently using back-propagation where the points $p_i$, $i \in [n]$ are grouped into batches. We performed $10 - 20$ iterations of Equation 4 for each $p_i$, $i \in [n]$.

**Scalar networks.**   For scalar networks, i.e., $l = 1$, $D_x F \in \mathbb{R}^{1 \times d}$, a direct computation shows that

$$D_x F(p; \theta_0)^\dagger = \frac{D_x F(p; \theta_0)^T}{\|D_x F(p; \theta_0)\|^2}. \tag{5}$$

That is, the point $p$ moves towards the level set $\mathcal{S}(\theta_0)$ by going in the direction of the steepest descent (or ascent) of $F$.

It is worth mentioning that the projection-on-level-sets formula in the case of $l = 1$ has already been developed in [25] and was used to find adversarial perturbations; our result generalizes to the intersection of several level sets and shows that this procedure is an instantiation of the generalized Newton algorithm.

The generalized Newton method (similarly to Newton method) is not guaranteed to find a point on the zero level set. We denote by $c_i = F(p_i; \theta_0)$ the level set values of the final point $p_i$; in the following, we use also points that failed to be projected with their level set $c_i$. Furthermore, we found that the generalized Newton method usually does find zeros of neural networks but these can be far from the initial projection point. Other, less efficient but sometimes more robust ways to project on neural level sets could be using gradient descent on $\|F(x; \theta)\|$ or zero finding in the direction of a successful PGD attack [19, 22] as done in [9]. In our robust networks application (Section 3.2) we have used a similar approach with the false position method.

**Relation to Elsayed et al. [11].**   The authors of [11] replace the margin distance with distance to the linearized network, while our approach is to directly sample the actual network's level set and move it explicitly. Specifically, for $L_2$ norm ($p = 2$) Elsayed's method is similar to our method using a single generalized Newton iteration, Equation 4.

## 2.2 Differentiable sample position

Next, we would like to relate each sample $p$, belonging to some level set $F(p; \theta_0) = c$, to the network parameters $\theta$. Namely, $p = p(\theta)$. The functional dependence of a sample $p$ on $\theta$ is defined by $p(\theta_0) = p$ and $F(p(\theta); \theta) = c$, for $\theta$ in some neighborhood of $\theta_0$. The latter condition ensures that $p$ stays on the $c$ level set as the network parameters $\theta$ change. As only first derivatives are used in the optimization of neural networks, it is enough to replace this condition with its first-order version. We get the following two equations:

$$p(\theta_0) = p \quad ; \quad \frac{\partial}{\partial \theta}\Big|_{\theta=\theta_0} F(p(\theta); \theta) = 0. \tag{6}$$

Using the chain rule, the second condition in Equation 6 reads:

$$D_x F(p, \theta_0) D_\theta p(\theta_0) + D_\theta F(p, \theta_0) = 0. \tag{7}$$

This is a system of linear equations with $d \times m$ unknowns (the components of $D_\theta p(\theta_0)$) and $l \times m$ equations. When $d > l$, this linear system is under-determined. Similarly to what is used in the generalized Newton method, a minimal norm solution is given by the Moore-Penrose inverse:

$$D_\theta p(\theta_0) = -D_x F(p, \theta_0)^\dagger D_\theta F(p, \theta_0). \tag{8}$$

The columns of the matrix $D_\theta p(\theta_0) \in \mathbb{R}^{d \times m}$ describe the velocity of $p(\theta)$ w.r.t. each of the parameters in $\theta$. The pseudo-inverse selects the minimal norm solution that can be shown to represent, in this case, a movement in the orthogonal direction to the level set (see Lemma C2 for a proof). We reiterate that for scalar networks, where $l = 1$, $D_x F(p_i; \theta_0)^\dagger$ has a simple closed-form expression, as shown in Equation 5.

**The sample network.** A possible simple solution to Equation 6 would be to use the linear function $p(\theta) = p + D_\theta p(\theta_0)(\theta - \theta_0)$, with $D_\theta p(\theta_0)$ as defined in Equation 8. Unfortunately, this would require storing $D_\theta p(\theta_0)$, using at-least $\mathcal{O}(m)$ space (i.e., the number of network parameters), for every projection point $p$. (We assume the number of output channels $l$ of $F$ is constant, which is the case in this paper.) A much more efficient solution is

$$p(\theta) = p - D_x F(p; \theta_0)^\dagger [F(p; \theta) - c], \tag{9}$$

that requires storing $D_x F(p; \theta_0)^\dagger$, using only $\mathcal{O}(d)$ space, where $d$ is the input space dimension, for every projection point $p$. Furthermore, Equation 9 allows an efficient implementation with a single network

$$G(p, D_x F(p; \theta_0)^\dagger; \theta) := p(\theta).$$

We call $G$ the *sample network*. Note that a collection of samples $p_i$, $i \in [n]$ can be treated as a batch input to $G$.

## 3 Incorporation of samples in loss functions

Once we have the sample network $p_i(\theta)$, $i \in [n]$, we can incorporate it in a loss function to control the neural level set $\mathcal{S}(\theta)$ in a desired way. We give three examples in this paper.

### 3.1 Geometric SVM

Support-vector machine (SVM) is a model which aims to train a linear binary classifier that would generalize well to new data by combining the hinge loss and a large margin term. It can be interpreted as encouraging large distances between training examples and the decision boundary. Specifically, the soft SVM loss takes the form [7]:

$$\mathrm{loss}(w, b) = \frac{1}{N} \sum_{j=1}^{N} \max\{0, 1 - y_j \ell(x_j; w, b)\} + \lambda \|w\|_2^2, \tag{10}$$

where $(x_j, y_j) \in \mathbb{R}^d \times \{-1, 1\}$, $j \in [N]$ is the binary classification training data, and $\ell(x; w, b) = w^T x + b$, $w \in \mathbb{R}^d$, $b \in \mathbb{R}$ is the linear classifier. We would like to generalize Equation 10 to a deep network $F : \mathbb{R}^d \times \mathbb{R}^m \to \mathbb{R}$ towards the goal of increasing the network's *input space margin*, which is defined as the minimal distance of training examples to the decision boundary $\mathcal{S}(\theta) = \{x \in \mathbb{R}^d \mid F(x; \theta) = 0\}$. Note that this is in strong contrast to standard deep network hinge

loss that works with the *output space margin* [30, 29], namely, measuring differences of output logits when evaluating the network on training examples. For that reason, this type of loss function does not penalize small input space margin, so long as it doesn't damage the output-level classification performance on the training data. Using the input margin over the output margin may also provide robustness to perturbations [11].

We now describe a new, geometric formulation of Equation 10, and use it to define the soft-SVM loss for general neural networks. In the linear case, the following quantity serves as the margin:

$$\|w\|_2^{-1} = d(\mathcal{S}(\theta), \mathcal{S}_1(\theta)) = d(\mathcal{S}(\theta), \mathcal{S}_{-1}(\theta))$$

where $\mathcal{S}_t(\theta) = \left\{ x \in \mathbb{R}^d \mid F(x;\theta) = t \right\}$, and $d(\mathcal{S}(\theta), \mathcal{S}_t(\theta))$ is the distance between the level sets, which are two parallel hyper-planes. In the general case, however, level sets are arbitrary hyper-surfaces which are not necessarily equidistant (i.e., the distance when traveling from $\mathcal{S}$ to $\mathcal{S}_t$ does not have to be constant across $\mathcal{S}$). Hence, for each data sample $x$, we define the following margin function:

$$\Delta(x;\theta) = \min \left\{ d\big(p(\theta), \mathcal{S}_1(\theta)\big), d\big(p(\theta), \mathcal{S}_{-1}(\theta)\big) \right\},$$

where $p(\theta)$ is the sample network of the projection of $x$ onto $\mathcal{S}(\theta_0)$. Additionally, note that in the linear case: $\left| w^T x + b \right| = d(x, \mathcal{S}(\theta))/\Delta(x;\theta)$. With these definitions in mind, Equation 10 can be given the geometric generalized form:

$$\text{loss}(w,b) = \frac{1}{N} \sum_{j=1}^{N} \max \left\{ 0, 1 - \text{sign}(y_j F(x_j;\theta)) \frac{d(x_j, p_j)}{\Delta(x_j;\theta)} \right\} + \frac{\lambda}{N} \sum_{j=1}^{N} \Delta(x_j;\theta)^\alpha, \qquad (11)$$

where $F(x;\theta)$ is a general classifier (such as a neural network, in our applications). Note that in the case where $F(x;\theta)$ is affine, $\alpha = -2$ and $d = L_2$, Equation 11 reduces back to the regular SVM loss, Equation 10. Figure 1d depicts the result of optimizing this loss in a 2D case, i.e., $x_j \in \mathbb{R}^2$; the light blue and red curves represent $\mathcal{S}_{-1}$ and $\mathcal{S}_1$.
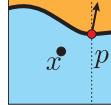
## 3.2 Robustness to adversarial perturbations

The goal of robust training is to prevent a change in a model's classification result when small perturbations are applied to the input. Following [22] the attack model is specified by some set $S \subset \mathbb{R}^d$ of allowed perturbations; in this paper we focus on the popular choice of $L_\infty$ perturbations, that is $S$ is taken to be the $\varepsilon$-radius $L_\infty$ ball, $\{x \mid \|x\|_\infty \leq \varepsilon\}$. Let $(x_j, y_j) \in \mathbb{R}^d \times \mathcal{L}$ denote training examples and labels, and let $\mathcal{S}^j(\theta) = \left\{ x \in \mathbb{R}^d \mid F^j(x;\theta) = 0 \right\}$, where $F^j(x;\theta) = f_j - \max_{i \neq j} f_i$, the decision boundary of label $j$. We define the loss

$$\text{loss}(\theta) = \frac{1}{N} \sum_{j=1}^{N} \lambda_j \max \left\{ 0, \varepsilon_j - \text{sign}(F^{y_j}(x_j;\theta)) d(x_j, \mathcal{S}^{y_j}(\theta)) \right\}, \qquad (12)$$

where $d(x, \mathcal{S}^j)$ is some notion of a distance between $x$ and $\mathcal{S}^j$, e.g., $\min_{y \in \mathcal{S}^j} \|x - y\|_p$ or $\int_{\mathcal{S}^j(\theta)} \|x - y\|_p \, d\mu(y)$, $d\mu$ is some probability measure on $\mathcal{S}^j(\theta)$. The parameter $\lambda_j$ controls the weighting between correct (i.e., $F^{y_j}(x_j;\theta) > 0$) and incorrect (i.e., $F^{y_j}(x_j;\theta) < 0$) classified samples. We fix $\lambda_j = 1$ for incorrectly classified samples and set $\lambda_j$ to be the same for correctly classified samples; The parameter $\varepsilon_j$ controls the desired target distances; Similarly to [11], the idea of this loss is: (i) if $x_j$ is incorrectly classified, pull the decision boundary $\mathcal{S}^{y_j}$ toward $x_j$; (ii) if $x_j$ is classified correctly, push the decision boundary $\mathcal{S}^{y_j}$ to be within a distance of at-least $\varepsilon_j$ from $x_j$.

In our implementation we have used $d(x, \mathcal{S}^j) = \rho(x, p(\theta))$, where $p = p(\theta_0) \in \mathcal{S}^j$ is a sample of this level set; $\rho(x, p) = \left| x^{i_*} - p^{i_*} \right|$, and $x^{i_*}, p^{i_*}$ denote the $i_*$-th coordinates of $x, p$ (resp.), $i_* = \arg\max_{i \in [d]} |D_x F(p;\theta_0)|$. This loss encourages $p(\theta)$ to move in the direction of the axis (i.e., $i_*$) that corresponds to the largest component in the gradient $D_x F(p;\theta_0)$. Intuitively, as $\rho(x,p) \leq \|x - p\|_\infty$, the loss pushes the level set $\mathcal{S}^j$ to leave the $\varepsilon_j$-radius $L_\infty$-ball in the direction of the axis that corresponds to the maximal speed of $p(\theta)$. The inset depicts an example where this distance measure is more effective than the standard $L_\infty$ distance: $D_x F(p;\theta)$ is shown as black arrow and the selected axis $i_*$ as dashed line.

In the case where the distance function $\min_{y \in \mathcal{S}(\theta)} \|x - y\|_p$ is used, the computation of its gradient using Equation 8 coincides with the gradient derivation of [9] up to a sign difference. Still, our derivation allows working with general level set points (i.e., not just the closest) on the decision boundary $\mathcal{S}^j$, and our sample network offers an efficient implementation of these samples in a loss function. Furthermore, we use the same loss for both correctly and incorrectly classified examples.

### 3.3 Manifold reconstruction

**Surface reconstruction.** Given a point cloud $\mathcal{X} = \{x_j\}_{j=1}^N \subset \mathbb{R}^d$ that samples, possibly with noise, some surface $\mathcal{M} \subset \mathbb{R}^3$, our goal is to find parameters $\theta$ of a network $F : \mathbb{R}^3 \times \mathbb{R}^m \to \mathbb{R}$, so that the neural level set $\mathcal{S}(\theta)$ approximates $\mathcal{M}$. Even more desirable is to have $F$ approximate the signed distance function to the unknown surface sampled by $\mathcal{X}$. To that end, we would like the neural level set $\mathcal{S}_t(\theta)$, $t \in \mathcal{T}$ to be of distance $|t|$ to $\mathcal{X}$, where $\mathcal{T} \subset \mathbb{R}$ is some collection of desired level set values. Let $\mathrm{d}(x, \mathcal{X}) = \min_{j \in [N]} \|x - x_j\|_2$ be the distance between $x$ and $\mathcal{X}$. We consider the reconstruction loss

$$\text{loss}(\theta) = \sum_{t \in \mathcal{T}} \left[ \int_{\mathcal{S}_t(\theta)} \left| \mathrm{d}(x, \mathcal{X}) - |t| \right|^p dv(x) \right]^{\frac{1}{p}} + \frac{\lambda}{N} \sum_{j=1}^N |F(x_j; \theta)|, \qquad (13)$$

where $dv(x)$ is the normalized volume element on $\mathcal{S}_t(\theta)$ and $\lambda > 0$ is a parameter. The first part of the loss encourages the $t$ level set of $F$ to be of distance $|t|$ to $\mathcal{X}$; note that for $t = 0$ this reconstruction error was used in level set surface reconstruction methods [36]. The second part of the loss penalizes samples $\mathcal{X}$ outside the zero level set $\mathcal{S}(\theta)$.

**Curve reconstruction.** In case of approximating a manifold $\mathcal{M} \subset \mathbb{R}^d$ with co-dimension greater than 1, e.g., a curve in $\mathbb{R}^3$, one cannot expect $F$ to approximate the signed distance function as no such function exists. Instead, we model the manifold via the level set of a vector-valued network $F : \mathbb{R}^d \times \mathbb{R}^m \to \mathbb{R}^l$ whose zero level set is an intersection of $l$ hyper-surfaces. As explained in Section 2, this generically defines a $d - l$ manifold. In that case we used the loss in Equation 13 with $\mathcal{T} = \{0\}$, namely, only encouraging the zero level set to be as close as possible to the samples $\mathcal{X}$.

## 4 Universality

To theoretically support the usage of neural level sets for modeling manifolds or controlling decision boundaries we provide a geometric universality result for multilayer perceptrons (MLP) with ReLU activations. That is, the level sets of MLPs can represent any watertight piecewise linear hyper-surface (i.e., manifolds of co-dimension 1 in $\mathbb{R}^d$ that are boundaries of $d$-dimensional polytopes). More specifically, we prove:

**Theorem 1.** *Any watertight, not necessarily bounded, piecewise linear hypersurface $\mathcal{M} \subset \mathbb{R}^d$ can be exactly represented as the neural level set $\mathcal{S}$ of a multilayer perceptron with ReLU activations, $F : \mathbb{R}^d \to \mathbb{R}$.*

The proof of this theorem is given in Appendix C. Note that this theorem is a geometrical version of Theorem 2.1 in [1], asserting that MLPs with ReLU activations can represent any piecewise linear continuous function.

## 5 Experiments

### 5.1 Classification generalization

In this experiment, we show that when training on small amounts of data, our geometric SVM loss (see Equation 11) generalizes better than the cross entropy loss and the hinge loss. Experiments were done on three datasets: MNIST [20], Fashion-MNIST [34] and CIFAR10 [18]. For all datasets we arbitrarily merged the labels into two classes, resulting



Figure 2: Generalization from small fractions of the data.

in a binary classification problem. We randomly sampled a fraction of the original training examples and evaluated on the original test set.

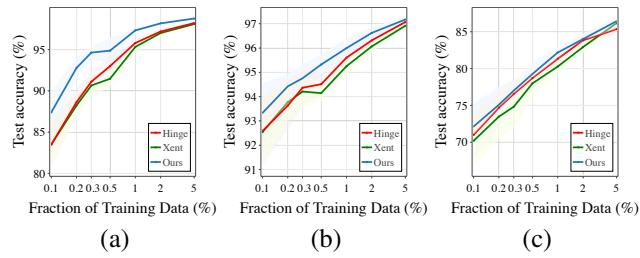| Method | Dataset | Arch. | Attack | $\varepsilon_{\text{train}}$ | Test Acc. | Rob. Acc. Xent | Rob. Acc. Margin |
|---|---|---|---|---|---|---|---|
| Standard | MNIST | ConvNet-4a | PGD$^{40}$($\varepsilon_{\text{attack}} = 0.3$) | - | 99.34% | 13.59% | 0.00% |
| Madry et al. [22] | MNIST | ConvNet-4a | PGD$^{40}$($\varepsilon_{\text{attack}} = 0.3$) | 0.3 | 99.35% | 96.04% | 96.11% |
| Madry et al. [22] | MNIST | ConvNet-4a | PGD$^{40}$($\varepsilon_{\text{attack}} = 0.3$) | 0.4 | 99.16% | 96.54% | 96.53% |
| TRADES [35] | MNIST | ConvNet-4a | PGD$^{40}$($\varepsilon_{\text{attack}} = 0.3$) | 0.3 | 98.97% | 96.75% | 96.74% |
| TRADES [35] | MNIST | ConvNet-4a | PGD$^{40}$($\varepsilon_{\text{attack}} = 0.3$) | 0.4 | 98.62% | 96.78% | 96.76% |
| Ours | MNIST | ConvNet-4a | PGD$^{40}$($\varepsilon_{\text{attack}} = 0.3$) | 0.4 | 99.35% | 99.23% | 97.35% |
| Standard | CIFAR10 | ConvNet-4b | PGD$^{20}$ ($\varepsilon_{\text{attack}} = 0.031$) | - | 83.67% | 0.00% | 0.00% |
| Madry et al. [22] | CIFAR10 | ConvNet-4b | PGD$^{20}$ ($\varepsilon_{\text{attack}} = 0.031$) | 0.031 | 71.86% | 39.84% | 38.18% |
| Madry et al. [22] | CIFAR10 | ConvNet-4b | PGD$^{20}$ ($\varepsilon_{\text{attack}} = 0.031$) | 0.045 | 63.66% | 41.53% | 39.13% |
| TRADES [35] | CIFAR10 | ConvNet-4b | PGD$^{20}$ ($\varepsilon_{\text{attack}} = 0.031$) | 0.031 | 71.24% | 41.89% | 38.4% |
| TRADES [35] | CIFAR10 | ConvNet-4b | PGD$^{20}$ ($\varepsilon_{\text{attack}} = 0.031$) | 0.045 | 68.24% | 42.04% | 38.18% |
| Ours | CIFAR10 | ConvNet-4b | PGD$^{20}$ ($\varepsilon_{\text{attack}} = 0.031$) | 0.045 | 71.96% | 38.45% | 38.54% |
| Standard | CIFAR10 | ResNet-18 | PGD$^{20}$ ($\varepsilon_{\text{attack}} = 0.031$) | - | 93.18% | 0.00% | 0.00% |
| Madry et al. [22] | CIFAR10 | ResNet-18 | PGD$^{20}$ ($\varepsilon_{\text{attack}} = 0.031$) | 0.031 | 81.0% | 47.29% | 46.58% |
| Madry et al. [22] | CIFAR10 | ResNet-18 | PGD$^{20}$ ($\varepsilon_{\text{attack}} = 0.031$) | 0.045 | 74.97% | 49.84% | 48.02% |
| TRADES [35] | CIFAR10 | ResNet-18 | PGD$^{20}$ ($\varepsilon_{\text{attack}} = 0.031$) | 0.031 | 83.04% | 53.31% | 51.36% |
| TRADES [35] | CIFAR10 | ResNet-18 | PGD$^{20}$ ($\varepsilon_{\text{attack}} = 0.031$) | 0.045 | 79.52% | 53.49% | 51.22% |
| Ours | CIFAR10 | ResNet-18 | PGD$^{20}$ ($\varepsilon_{\text{attack}} = 0.031$) | 0.045 | 81.3% | 79.74% | 43.17% |

Table 1: Results of different $L_\infty$-bounded attacks on models trained using our method (described in Section 3.2) compared to other methods.

Due to the variability in the results, we rerun the experiment 100 times for MNIST and 20 times for Fashion-MNIST and CIFAR10. We report the mean accuracy along with the standard deviation. Figure 2 shows the test accuracy of our loss compared to the cross-entropy loss and hinge loss over different training set sizes for MNIST (a), Fashion-MNIST (b) and CIFAR10 (c). Our loss function outperforms the standard methods.

For the implementation of Equation 11 we used $\alpha = -1$, d $= L_\infty$, and approximated $\mathrm{d}(x, \mathcal{S}_t) \approx \|x - p^t\|_\infty$, where $p^t$ denotes the projection of $p$ on the level set $\mathcal{S}_t$, $t \in \{-1, 0, 1\}$ (see Section 2.1). The approximation of the term $\Delta(x; \theta)$, where $x = x_j$ is a train example, is therefore $\min\left\{\|p^0 - p^{-1}\|_\infty, \|p^0 - p^1\|_\infty\right\}$. See Appendix B.2 for further implementation details.

## 5.2 Robustness to adversarial examples

In this experiment we used our method with the loss in Equation 12 to train robust models on MNIST [20] and CIFAR10 [18] datasets. For MNIST we used ConvNet-4a (312K params) used in [35] and for CIFAR10 we used two architectures: ConvNet-4b (2.5M params) from [33] and ResNet-18 (11.2M params) from [35]. We report results using the loss in Equation 12 with the choice of $\varepsilon_j$ fixed as $\varepsilon_{\text{train}}$ in Table 1, $\lambda_j$ to be 1, 10 for MNIST and CIFAR10 (resp.), and d $= \rho$ as explained in Section 3.2. We evaluated our trained networks on $L_\infty$ bounded attacks with $\varepsilon_{\text{attack}}$ radius using Projected Gradient Descent (PGD) [19, 22] and compared to networks with the same architectures trained using the methods of Madry et al. [22] and TRADES [35]. We found our models to be robust to PGD attacks based on the Xent loss; during the revision of this paper we discovered weakness of our trained models to PGD attack based on the margin loss, i.e., $\min\left\{F^j\right\}$, where $F^j = f_j - \max_{i \neq j} f_i$; we attribute this fact to the large gradients created at the level set. Consequently, we added margin loss attacks to our evaluation. The results are summarized in Table 1. Note that although superior in robustness to Xent attacks we are comparable to baseline methods for margin attacks. Furthermore, we believe the relatively low robust accuracy of our model when using the ResNet-18 architecture is due to the fact that we didn't specifically adapt our method to Batch-Norm layers.

In Appendix B.3 we provide tables summarizing robustness of our trained models (MNIST ConvNet-4a and CIFAR10 ConvNet-4b) to black-box attacks; we log black-box attacks of our and baseline methods [22, 35] in an all-versus-all fashion. In general, we found that all black-box attacks are less effective than the relevant white-box attacks, our method performs better when using standard model black-box attacks, and that all three methods compared are in general similar in their black-box robustness.
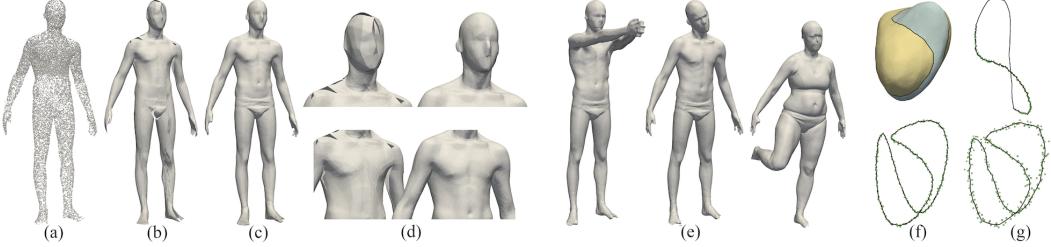
Figure 3: Point cloud reconstruction. Surface: (a) input point cloud; (b) AtlasNet [13] reconstruction; (c) our result; (d) blow-ups; (e) more examples of our reconstruction. Curve: (f) bottom image shows a curve reconstructed (black line) from a point cloud (green points) as an intersection of two scalar level sets (top image); (g) bottom shows curve reconstruction from a point cloud with large noise, where top image demonstrates the graceful completion of an open curve point cloud data.

## 5.3 Surface and curve reconstruction

In this experiment we used our method to reconstruct curves and surfaces in $\mathbb{R}^3$ using only incomplete point cloud data $\mathcal{X} \subset \mathbb{R}^3$, which is an important task in 3D shape acquisition and processing. Each point cloud is processed independently using the loss function described in Equation 13, which encourages the zero level set of the network to pass through the point cloud. For surfaces, it also moves other level sets to be of the correct distance to the point cloud.

For surface reconstruction, we trained on 10 human raw scans from the FAUST dataset [5], where each scan consists of $\sim$ 170K points in $\mathbb{R}^3$. The scans include partial connectivity information which we do not use. After convergence, we reconstruct the mesh using the marching cubes algorithm [21] sampled at a

|  | Chamfer L1 | Chamfer L2 |
|---|---|---|
| AtlasNet-1 sphere | $23.56 \pm 2.91$ | $17.69 \pm 2.45$ |
| AtlasNet-1 patch | $18.67 \pm 3.45$ | $13.38 \pm 2.66$ |
| AtlasNet-25 patches | $11.54 \pm 0.53$ | $7.89 \pm 0.42$ |
| Ours | $\mathbf{10.71} \pm 0.63$ | $\mathbf{7.32} \pm 0.46$ |

Table 2: Surface reconstruction results.

resolution of $[100]^3$. Table 2 compares our method with the recent method of [13] which also works directly with point clouds. Evaluation is done using the Chamfer distance [12] computed between 30K uniformly sampled points from our and [13] reconstructed surfaces and the ground truth registrations provided by the dataset, with both $L_1, L_2$ norms. Numbers in the table are multiplied by $10^3$. We can see that our method outperforms its competitor; Figure 3b-3e show examples of surfaces reconstructed from a point cloud (a batch of 10K points is shown in 3a) using our method (in 3c, 3d-right, 3e), and the method of [13] (in 3b, 3d-left). Importantly, we note that there are recent methods for implicit surface representation using deep neural networks [6, 26, 24]. These methods use signed distance information and/or the occupancy function of the ground truth surfaces and perform regression on these values. Our formulation, in contrast, allows working directly on the more common, raw input of point clouds.

For curve reconstruction, we took a noisy sample of parametric curves in $\mathbb{R}^3$ and used similar network to the surface case, except its output layer consists of two values. We trained the network with the loss Equation 13, where $\mathcal{T} = \{0\}$, using similar settings to the surface case. Figure 3f shows an example of the input point cloud (in green) and the reconstructed curve (in black) (see bottom image), as well as the two hyper-surfaces of the trained network, the intersection of which defines the final reconstructed curve (see top image); 3g shows two more examples: reconstruction of a curve from higher noise samples (see bottom image), and reconstruction of a curve from partial curve data (see top image); note how the network gracefully completes the curve.

## 6 Conclusions

We have introduced a simple and scalable method to incorporate level sets of neural networks into a general family of loss functions. Testing this method on a wide range of learning tasks we found the method particularly easy to use and applicable in different settings. Current limitations and interesting venues for future work include: applying our method with the batch normalization layer (requires generalization from points to batches); investigating control of intermediate layers' level sets; developing sampling conditions to ensure coverage of the neural level sets; and employing additional geometrical regularization to the neural level sets (e.g., penalize curvature).

## Acknowledgments

# References

[1] R. Arora, A. Basu, P. Mianjy, and A. Mukherjee. Understanding deep neural networks with rectified linear units. *arXiv preprint arXiv:1611.01491*, 2016.

[2] H. Ben-Hamu, H. Maron, I. Kezurer, G. Avineri, and Y. Lipman. Multi-chart generative surface modeling. In *SIGGRAPH Asia 2018 Technical Papers*, page 215. ACM, 2018.

[3] A. Ben-Israel. A newton-raphson method for the solution of systems of equations. *Journal of Mathematical analysis and applications*, 15(2):243–252, 1966.

[4] M. Berger, A. Tagliasacchi, L. M. Seversky, P. Alliez, G. Guennebaud, J. A. Levine, A. Sharf, and C. T. Silva. A survey of surface reconstruction from point clouds. In *Computer Graphics Forum*, volume 36, pages 301–329. Wiley Online Library, 2017.

[5] F. Bogo, J. Romero, M. Loper, and M. J. Black. FAUST: Dataset and evaluation for 3D mesh registration. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, Piscataway, NJ, USA, June 2014. IEEE.

[6] Z. Chen and H. Zhang. Learning implicit fields for generative shape modeling. *arXiv preprint arXiv:1812.02822*, 2018.

[7] C. Cortes and V. Vapnik. Support-vector networks. *Machine learning*, 20(3):273–297, 1995.

[8] G. Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2(4):303–314, 1989.

[9] G. W. Ding, Y. Sharma, K. Y. C. Lui, and R. Huang. Max-margin adversarial (mma) training: Direct input space margin maximization through adversarial training. *arXiv preprint arXiv:1812.02637*, 2018.

[10] G. W. Ding, L. Wang, and X. Jin. AdverTorch v0.1: An adversarial robustness toolbox based on pytorch. *arXiv preprint arXiv:1902.07623*, 2019.

[11] G. Elsayed, D. Krishnan, H. Mobahi, K. Regan, and S. Bengio. Large margin deep networks for classification. In *Advances in Neural Information Processing Systems*, pages 842–852, 2018.

[12] H. Fan, H. Su, and L. J. Guibas. A point set generation network for 3d object reconstruction from a single image. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 605–613, 2017.

[13] T. Groueix, M. Fisher, V. G. Kim, B. C. Russell, and M. Aubry. A papier-mâché approach to learning 3d surface generation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 216–224, 2018.

[14] M. Hein and M. Andriushchenko. Formal guarantees on the robustness of a classifier against adversarial manipulation. In *Advances in Neural Information Processing Systems*, pages 2266–2276, 2017.

[15] K. Hornik, M. Stinchcombe, and H. White. Multilayer feedforward networks are universal approximators. *Neural networks*, 2(5):359–366, 1989.

[16] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[17] S. G. Krantz and H. R. Parks. *The implicit function theorem: history, theory, and applications*. Springer Science & Business Media, 2012.

[18] A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.

[19] A. Kurakin, I. Goodfellow, and S. Bengio. Adversarial machine learning at scale. *arXiv preprint arXiv:1611.01236*, 2016.

[20] Y. LeCun. The mnist database of handwritten digits. *http://yann. lecun. com/exdb/mnist/*, 1998.

[21] W. E. Lorensen and H. E. Cline. Marching cubes: A high resolution 3d surface construction algorithm. In *ACM siggraph computer graphics*, volume 21, pages 163–169. ACM, 1987.

[22] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.

[23] A. Matyasko and L.-P. Chau. Margin maximization for robust classification using deep learning. In *2017 International Joint Conference on Neural Networks (IJCNN)*, pages 300–307. IEEE, 2017.

[24] L. Mescheder, M. Oechsle, M. Niemeyer, S. Nowozin, and A. Geiger. Occupancy networks: Learning 3d reconstruction in function space. *arXiv preprint arXiv:1812.03828*, 2018.

[25] S.-M. Moosavi-Dezfooli, A. Fawzi, and P. Frossard. Deepfool: a simple and accurate method to fool deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2574–2582, 2016.

[26] J. J. Park, P. Florence, J. Straub, R. Newcombe, and S. Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. *arXiv preprint arXiv:1901.05103*, 2019.

[27] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer. Automatic differentiation in pytorch. 2017.

[28] J. Sokolić, R. Giryes, G. Sapiro, and M. R. Rodrigues. Robust large margin deep neural networks. *IEEE Transactions on Signal Processing*, 65(16):4265–4280, 2017.

[29] S. Sun, W. Chen, L. Wang, and T.-Y. Liu. Large margin deep neural networks: Theory and algorithms. *arXiv preprint arXiv:1506.05232*, 148, 2015.

[30] Y. Tang. Deep learning using support vector machines. *CoRR, abs/1306.0239*, 2, 2013.

[31] F. Williams, T. Schneider, C. Silva, D. Zorin, J. Bruna, and D. Panozzo. Deep geometric prior for surface reconstruction. *arXiv preprint arXiv:1811.10943*, 2018.

[32] E. Wong and J. Z. Kolter. Provable defenses against adversarial examples via the convex outer adversarial polytope. *arXiv preprint arXiv:1711.00851*, 2017.

[33] E. Wong, F. Schmidt, J. H. Metzen, and J. Z. Kolter. Scaling provable adversarial defenses. In *Advances in Neural Information Processing Systems*, pages 8400–8409, 2018.

[34] H. Xiao, K. Rasul, and R. Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.

[35] H. Zhang, Y. Yu, J. Jiao, E. P. Xing, L. E. Ghaoui, and M. I. Jordan. Theoretically principled trade-off between robustness and accuracy. *arXiv preprint arXiv:1901.08573*, 2019.

[36] H.-K. Zhao, S. Osher, and R. Fedkiw. Fast surface reconstruction using the level set method. In *Proceedings IEEE Workshop on Variational and Level Set Methods in Computer Vision*, pages 194–201. IEEE, 2001.

# A   Additional Experiments

| Initialization Method | Chamfer | Hausdorf |
|---|---|---|
| Uniform [-0.35,0.35] | 0.011 | 0.141 |
| Normal $\sigma = 0.01$ | 0.006 | 0.017 |
| Normal $\sigma = 0.05$ | 0.01 | 0.132 |



(a) Normal $\sigma = 0.05$      (b) Uniform
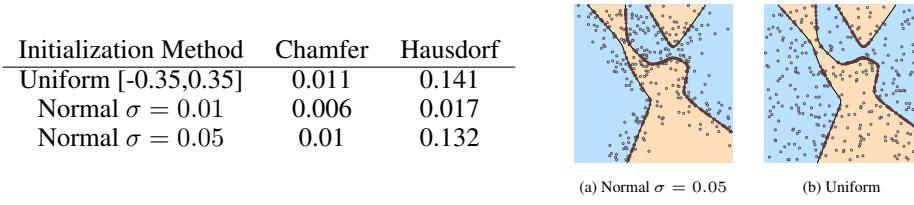
Figure A1: Distribution of samples on $d = 2$ neural level set.

## A.1   Distribution of points on the level set.

Achieving well distributed samples of neural level set is a challenge, especially for high dimensions. In the inset we quantify the quality of distribution in low dimension, $d = 2$, (where ground truth dense sampling of the level set is tractable). The table in Figure A1 logs the Chamfer and Hausdorff distances of the resulting sampling distribution and the level set of a neural network trained with Xent loss in 2 dimensions ((a) and (b)) where projected points (red) are initialized using a uniformly distributed points (gray, (b)) or normally perturbed level set samples (gray, (a)).
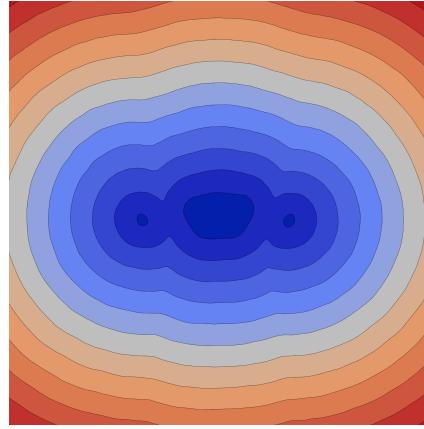


Figure A2: Level sets of a network $F$ from the experiment described in Section 5.3 shown along a cross-cut. Note how the iso-levels are equispaced, as encouraged by the loss in Equation 13.

## A.2   Level sets of reconstruction networks resemble signed distance function

Figure A2 shows iso-levels of one of the networks from the experiment described in Section 5.3. Note how the level sets resemble the level sets of a signed distance function.

# B Implementation details

All experiments are conducted on a Tesla V100 Nvidia GPU using PYTORCH framework [27].

| ConvNet-2a | ConvNet-2b | ConvNet-4a | ConvNet-4b | FC-1 | FC-2 |
|---|---|---|---|---|---|
| CONV 16 4x4+2 | CONV 32 5x5+1 | CONV 32 3x3+1 | CONV 32 3x3+1 | FC 512 | FC 509 |
| CONV 32 4x4+2 | MAXPOOL 2x2 | CONV 32 3x3+1 | CONV 32 4x4+2 | FC 512 | FC + SKIP 509 |
| FC 100 | CONV 64 5x5+1 | MAXPOOL 2x2 | CONV 64 3x3+1 | FC 512 | FC + SKIP 509 |
| FC 10 | MAXPOOL 2x2 | CONV 64 3x3+1 | CONV 64 4x4+2 | FC 1 | FC + SKIP 509 |
| | FC 512 | CONV 64 3x3+1 | FC 512 | | FC + SKIP 509 |
| | FC 10 | MAXPOOL 2x2 | FC 512 | | FC + SKIP 509 |
| | | FC 200 | FC 10 | | FC + SKIP 509 |
| | | FC 200 | | | FC 1 |
| | | FC 10 | | | |

Table B1: Our architectures. CONV $kw \times h + s$ corresponds to a convolution layer with $k$ channels, a kernel of size $w \times h$ and stride $s$. FC $n$ correspond to a fully connected layer with $n$ outputs. FC + SKIP indicates a skip connection to the input layer. Each CONV/FC layer is followed by a ReLU activation except for the last fully connected layer.

## B.1 Parameters of experiments shown in Figure 1

We train a 4-layer MLP $F(x; \theta) : \mathbb{R}^2 \times \mathbb{R}^m \to \mathbb{R}^2$, as in architecture FC-1, for 1000 epochs using the ADAM optimizer [16] with learning rate $0.001$. For the geometrical SVM loss we use $\lambda = 0.001$. Training set is composed of 16 points in $\mathbb{R}^2$, all of which lie inside $[0, 0.5]^2$. Batch size is 1. The sample network makes a maximum of 20 iterations for the projection procedure.

## B.2 Classification generalization

In Table B2 we summarize all hyper-parameters used in the generalization experiments (Section 5.1). For cross-entropy and hinge losses we checked learning rates of $0.001, 0.005, 0.01, 0.02$ and chose the ones that performed best. All models were trained using SGD (Nesterov) optimizer with momentum $0.9$ and weight decay $10^{-4}$.

| Dataset / Params | MNIST | Fashion-MNIST | CIFAR10 |
|---|---|---|---|
| Architecture | ConvNet-2b | ConvNet-2b | ConvNet-4b |
| Geometric SVM $\lambda$ | $\lambda$ grows linearly from 0.01 to 0.2 over 50 epochs | $\lambda$ grows linearly from 0.01 to 0.2 over 50 epochs | $\lambda$ grows linearly from 0 to 0.01 over 50 epochs |
| # Epochs | 200 | 200 | 100 |
| Batch size | 256 (32 for fraction $\leq 0.3$) | 32 | 32 |
| # Iterations for projection proc. | 20 | 20 | 20 |
| Initial learning rate | 0.02 | 0.02 | 0.01 |
| Learning rate decay | multipled by 0.5 at epochs 50, 100, 120, 140, 160, 180 | multipled by 0.5 at epochs 50, 100, 120 | multipled by 0.5 at epoch 50 |

Table B2: Generalization experiments hyperparameters

## B.3 Adversarial robustness

We describe the parameters used in the experiments shown in Secion 5.2.

**Training parameters**

We use the networks described in Table B1, labeled ConvNet-4a and ConvNet-4b (following [32]) for the MNIST and CIFAR10 experiments respectively. Additionally, for CIFAR10 we add an experiment with ResNet-18 architecture as in [35]. All networks are trained with batches of size 128. For the projection on the zero levelset procedure, we used the false-position method with a maximum of 40 iterations per batch. The standard models are trained using cross-entropy loss for 200 epochs on MNIST and CIFAR10 respectively (batch-size and learning rates are similar to the above mentioned models). All our models are trained using ADAM optimizer [16].

**Bounded Attack (Table 1)**   We use the *advertorch* library [10]. The attacks parameters are, for MNIST: $\varepsilon_{\text{attack}} = 0.3$, PGD-iterations 40 and 100 and step size 0.01. For CIFAR10: $\varepsilon_{\text{attack}} = 8/255$, PGD-iterations 20 and step size 0.003. All models are evaluated at epoch 200, except for Madry defense with ResNet-18 architecture evaluated at epoch 50.

|  (a) Robust Accuracy Xent | | | | |
|---|---|---|---|---|
| Source / Target | Standard | Madry | Trades | Ours |
| Madry | 98.96% | 96.04% | 97.76% | 99.21% |
| Trades | 98.57% | 97.46% | 96.78% | 98.87% |
| Ours | 99.04% | 97.78% | 97.95% | 99.23% |

|  (b) Robust Accuracy Margin | | | | |
|---|---|---|---|---|
| Source / Target | Standard | Madry | Trades | Ours |
| Madry | 98.95% | 96.11% | 97.81% | 98.78% |
| Trades | 98.56% | 97.5% | 96.74% | 98.46% |
| Ours | 99.04% | 97.87% | 97.99% | 97.35% |

Table B3: MNIST: Comparison of our method and baseline methods under black-box PGD[40] attack with $\varepsilon_{\text{attack}} = 0.3$. Rows (target) are the attacked models. All models are trained with ConvNet-4a architecture. Diagonal represents white-box attacks.

|  (a) Robust Accuracy Xent | | | | |
|---|---|---|---|---|
| Source / Target | Standard | Madry | Trades | Ours |
| Madry | 61.5% | 41.53% | 49.76% | 50.97% |
| Trades | 67.84% | 54.72% | 41.89% | 53.11% |
| Ours | 68.43% | 56.71% | 54.47% | 38.45% |

|  (b) Robust Accuracy Margin | | | | |
|---|---|---|---|---|
| Source / Target | Standard | Madry | Trades | Ours |
| Madry | 61.46% | 39.13% | 39.14% | 51.08% |
| Trades | 67.58% | 53.15% | 38.25% | 53.1% |
| Ours | 68.42% | 55.85% | 54.04% | 38.54% |

Table B4: CIFAR10: Comparison of our method and baseline methods under black-box PGD[20] attack with $\varepsilon_{\text{attack}} = 0.031$. Rows (target) are the attacked models. All models are trained with ConvNet-4b architecture. Diagonal represents white-box attacks.

## B.4   Surface Reconstruction

We describe the parameters used for the experiments in Section 5.3. For the Faust benchmark, the network architecture is set to FC-2 (similarly to [6, 26]) and is used both for our model and AtlasNet. The optimization is done using the ADAM optimizer, batch size set to 10 and the initial learning rate is set to 0.001 (decreased by half at epochs 500,1500,3500). Some additional implementation details are: first, we set the parameter $\lambda$ in our reconstruction loss to grow linearly from 1 to 5 over 1000 epochs. Next, to generate samples of $\mathcal{S}_t$ we add Gaussian noise ($\sigma = 0.1$) to the input batch, randomly sample half of the points and use it as initialization for the projection procedure to $\mathcal{S}_0$. The other half is used to sample various level sets $\mathcal{S}_t$ (see Equation 13). The number of iterations for the projection procedure is set to 10.

For the curve reconstruction experiment, the architecture used is FC-1 with the minor difference that the last layer output size is 2. The ground truth is generated by randomly sampling 6 points in space and generating a curve passing through the points, using cubic spline interpolation. We generate the input point cloud by sampling the ground truth curve and adding small Gaussian noise. The sample size is 300 and sample points are chosen using Farthest Point Sampling. We generate samples from

$\mathcal{S}_0$ using the same procedure described above with the minor difference that the entire batch is used as initialization for the projection procedure (other, non-zero level sets are not sampled).

## C   Proofs

**Lemma 1.** *Let $\ell(x) = Ax + b$, $A \in \mathbb{R}^{l \times d}$, $b \in \mathbb{R}^l$, $\ell < d$, and $A$ is of full rank $l$. Then Equation 4 applied to $F(x) = \ell(x)$ is an orthogonal projection on the zero level-set of $\ell$, namely, on $\{x \mid \ell(x) = 0\}$.*

*Proof.* Let $p \in \mathbb{R}^d$ be the starting point. A single generalized Newton iteration (Equation 4) is

$$p^{\text{next}} = p - A^\dagger(Ap + b). \tag{14}$$

First, $p^{\text{next}}$ is indeed on the level set because: $\ell(p^{\text{next}}) = A(p - A^\dagger(Ap + b)) + b = 0$, where we used the fact that $AA^\dagger A = A$, and $AA^\dagger = I$ (since $\text{rank}(A) = l$). Furthermore, from Equation 14 we read that $p^{\text{next}} - p \in \text{Im}A^\dagger$ and therefore $p^{\text{next}} - p \in \text{Im}A^T$. This implies that $p^{\text{next}} - p \perp \text{Ker}A$. But $\text{Ker}A$ is the tangent space of the level set $\{x \mid \ell(x) = 0\}$, so $p^{\text{next}}$ is the orthogonal projection of $p$ on the zero level set of $\ell$. $\qquad\square$

**Lemma 2.** *The columns of the solution in Equation 8, namely $D_\theta p(\theta_0)$, are in the orthogonal space to the level set $\mathcal{S}(\theta_0)$ at $p_0$.*

*Proof.* $D_\theta p \in \mathbb{R}^{d \times m}$ describes the speed of $p$ w.r.t. each of the parameters in $\theta$. If we assume $A := D_x F(p; \theta_0)$ is of full rank $l$, which is the generic case, then the Moore-Penrose inverse has the form $A^\dagger = A^T(AA^T)^{-1}$. This indicates that the columns of $D_\theta p(\theta_0) = -A^\dagger D_\theta F(p; \theta_0) \in \mathbb{R}^{d \times m}$ belong to $\text{Im}A^T$, which in turn implies that they are orthogonal to $\text{Ker}A$, which is the tangent space of the level set at the point $p_0$ $\qquad\square$

**Theorem 1.** *Any watertight, not necessarily bounded, piecewise linear hypersurface $\mathcal{M} \subset \mathbb{R}^d$ can be exactly represented as the neural level set $\mathcal{S}$ of a multilayer perceptron with ReLU activations, $F : \mathbb{R}^d \to \mathbb{R}$.*

*Proof.* Let $h_i(x) = a_i^T x + b_i = 0$, $i \in [k]$ denote the planes supporting the faces of $\mathcal{M}$ where $a_i$ are chosen to be the outward normals to $\mathcal{M}$. Since $\mathcal{M}$ is watertight, it is the boundary of a $d$-dimensional polytope $P$.

For each $\lambda \in \{-1, 0, 1\}^k$, let $P_\lambda = \cap_{i \in [k]} \{x \mid \lambda_i h_i(x) \geq 0\}$. Simply put, $P_\lambda$ is a polytope that is the intersection of closed half-spaces defined by the some of the hyperplanes $h_i$. Out of all the possible $P_\lambda$'s, we're only interested in those that are contained in $P$, so we define $\Lambda = \{\lambda \mid P_\lambda \subseteq P\}$. Now, we wish to show that every point in the interior of the large polytope necessarily also lies in the interior of some small polytope in our collection, i.e that $\cup_{\lambda \in \Lambda} \mathring{P}_\lambda = \mathring{P}$. So let $x \in \mathring{P}$. There are two cases:

Case 1: $h_i(x) \neq 0 \, \forall i \in [k]$. That is, $x$ does not lie exactly on a hyper-plane. We can take the following polytope $P_\lambda$ which contains $x$ in its interior: $\lambda_i = \text{sign}(h_i(x))$. We note that $\lambda \in \{-1, 1\}^k$, and we call such a polytope *minimal*. We argue that the interior of a minimal polytope is either completely inside $P$ or completely outside it. This is true because otherwise the minimal polytope will contain two points that are on two different sides of some hyper-plane, which is inconsistent with $\lambda \in \{-1, 1\}^k$. In our case, we know that $P_\lambda$ and $P$ both contain $x$ in their interior, so necessarily $P_\lambda \subseteq P$, which means that $\lambda \in \Lambda$.

Case 2: $\exists \{i_1, ..., i_l\} \subseteq [k]$ s.t. $h_i(x) = 0 \, \forall i \in \{i_1, ..., i_l\}$, and $h_i \neq 0 \, \forall i \in [k] \setminus \{i_1, ..., i_l\}$. In this case there is no minimal polytope that contains $x$ in its interior, so let us consider all of the minimal polytopes which contain $x$ on their boundary. Let $P_\mu$ be such a minimal polytope. As previously stated, the interior of $P_\mu$ is either completely inside $P$ or completely outside it, but since $x$ is both on the boundary of $P_\mu$ and in the interior of $P$ then necessarily $P_\mu$ is completely inside $P$, i.e, $P_\mu \subset P$. We are interested in the union of all such minimal polytopes. Note that for such a minimal polytope

$P_\mu$, necessarily $\mu_i = sign(h_i(x))\ \forall i \in [k] \setminus \{i_1, ..., i_l\}$. For $i \in \{i_1, ..., i_l\}$, $\mu_i$ may receive any value in $\{1, -1\}$. Thus, the union of all such minimal polytopes is $P_\lambda$ where:

$$\lambda_i = \begin{cases} 0 & ,\ i \in \{i_1, ..., i_l\} \\ \text{sign}(h_i(x)) & ,\ otherwise \end{cases}$$

which clearly contains $x$ in its interior and is itself contained in $P$ (because it is the union of minimal polytopes which are contained in $P$), i.e $\lambda \in \Lambda$.

We are now ready to define a function which will receive positive values on the interior of $P$, negative values outside of $P$, and will have $\mathcal{M}$ as its levelset:

$$f(x) = \max_{\lambda \in \Lambda} \min_{i \in [k]} \lambda_i h_i(x)$$

$f$ is a piecewise linear function and can, therefore, according to Theorem 2.1 in [1], be encoded as an MLP with ReLU activations. The idea is to build $\max$ operators using linear layers and ReLU via $\max\{a, b\} = \frac{\sigma(a-b)}{2} + \frac{\sigma(b-a)}{2} + \frac{a+b}{2}$, where $\sigma(x) = \max(0, s)$ is the ReLU activation. Using this binary $\max$, one can recursively build the $\max$ of a vector. $\min$ is treated similarly.

$\square$