

Wasserstein Diffusion Tikhonov Regularization

Alex Tong Lin¹, Yonatan Dukler¹, Wuchen Li¹, and Guido Montúfar^{1,2,3}

¹Department of Mathematics and ²Department of Statistics, UCLA, CA 90095;

³Max Planck Institute for Mathematics in the Sciences, 04103 Leipzig, Germany

Abstract

We propose regularization strategies for learning discriminative models that are robust to in-class variations of the input data. We use the Wasserstein-2 geometry to capture semantically meaningful neighborhoods in the space of images, and define a corresponding input-dependent additive noise data augmentation model. Expanding and integrating the augmented loss yields an effective Tikhonov-type Wasserstein diffusion smoothness regularizer. This approach allows us to apply high levels of regularization and train functions that have low variability within classes but remain flexible across classes. We provide efficient methods for computing the regularizer at a negligible cost in comparison to training with adversarial data augmentation. Initial experiments demonstrate improvements in generalization performance under adversarial perturbations and also large in-class variations of the input data.

Key words: Wasserstein sample space; Gaussian distribution from Wasserstein Diffusion process; Tikhonov regularization; Adversarial robustness.

1 Introduction

The sensitivity of trained discriminative models to small perturbations of the input data has become a reason of concern and an important topic of research in recent years (Szegedy et al., 2014; Nguyen et al., 2015). In particular, it has been observed that neural networks which have been trained to have good test performance can be fooled when the inputs are slightly perturbed in a way that is imperceptible to humans. This indicates a poor generalization ability, and specifically, that the solutions found with naive training and validation techniques are not capturing appropriate smoothness priors over the input space. A number of recent works have proposed approaches to improve robustness to perturbations (Cissé et al., 2017; Liao et al., 2018; Samangouei et al., 2018; Wang et al., 2018a; Finlay et al., 2019), while a complementary line of work probes the limitations of trained networks (Salman et al., 2019) and develops strategies to generate adversarial attacks (Carlini et al., 2017; Moosavi-Dezfooli et al., 2016, 2017; Shafahi et al., 2019).

Intuitively, a smoother function at fixed training accuracy should be more robust to perturbations of the input, including adversarial attacks. Therefore, one strategy is to train the discriminative function with smoothness regularizers, such as noise added to the training examples (adversarial or random) or penalizing the norm of the gradient with respect to the inputs. We note, however, that the notion of a ‘small’ perturbation will strongly depend on how we decide to measure distances in the space of inputs. The gradient and its norm depend on the geometric structure that is laid on input space.

While it is convenient to use the L^2 metric (Euclidean), it is well understood that many data types of interest are not Euclidean. In particular, the L^2 metric does not measure distances between images in the way that we perceive them. Changes that humans consider small, might correspond to changes that the classifier considers to be large in this metric. Moreover, it is clear that a discriminative function on image data should be more stable in certain directions and more variable in other directions. This distinction is not well captured by isotropic smoothness regularizers.

To construct more effective smoothness regularizers, two general approaches come to mind: 1) Measure distances in a metric representation of the raw inputs, $d_\phi(x, y)^2 = \sum_j |\phi(x)_j - \phi(y)_j|^2$, where ϕ is some feature representation function that might be trained separately from or together with the discriminative task. Examples in this direction include preprocessing of the inputs by downsampling (Guo et al., 2018), autoencoders, and approaches that regularize intermediate representations within the neural network that is being trained for the discriminative task, such as injecting noise in the layers of a ResNet (Wang et al., 2018b). 2) Measure distances directly on the inputs (or following light preprocessing), but use a metric that is reflective of our perception of the data. Both approaches allow for data driven specifications and also direct incorporation of prior knowledge about the domain. We focus on the input space approach with the Wasserstein distance, in particular the Wasserstein-2 metric and geometry.

The Wasserstein distance is known to be an effective metric in the space of images, as demonstrated in image retrieval problems (Rubner et al., 2000; Solomon et al., 2015, 2014) and related applications (Peyré et al., 2019). In particular, the Wasserstein distance is robust to natural variations such as translations and independent noise added to the pixel values. Importantly, the Wasserstein distance exhibits a Riemannian metric structure. This allows us to define Wasserstein gradient penalties and effective Wasserstein Gaussian noise¹ in the space of images. The Wasserstein metric depends on the specific location at which it is being evaluated, and can define neighborhoods with a reasonable degree of semantic meaning. See, for example, the Wasserstein geodesics balls illustrated in Figure 1.

Our approach in this article is inspired by recent work (Dukler et al., 2019) which introduces a Wasserstein ground metric for Wasserstein GANs and demonstrates that this facilitates training of discriminators that are more stable to natural variations of image data. We suggest that regularization based on Wasserstein geometry can make a discriminative function noticeably smoother along the directions of natural variations of images, but without making it constant along the directions of semantic variation. Moreover, a generative perturbation model can be folded into the training objective (by computing the expectation value over perturbations of the Taylor expanded loss around each training example). This yields an effective penalty term that integrates (up to a given order in the expansion) a continuum of perturbations at once. The Wasserstein geometry on image space is described by a metric tensor whose inverse is a linear weighted Laplacian matrix. This fact allows us to compute the Wasserstein diffusion smoothness regularizer at a negligible cost.

In Section 2 we discuss adversarial attacks. In Section 3 we discuss training with input noise. We propose a Wasserstein Gaussian distribution in image space, which reflects the natural local variability of images. Then we compute the expectation of the perturbed objective by Taylor expansions in the Wasserstein space. This leads to a Tikhonov-type Wasserstein diffusion smoothness regularizer. In Section 5 we present preliminary experimental results. In Section 6 we discuss relations of the proposed methods to some of the existing literature, and in Section 7 we offer a discussion.

2 Adversarial training and ground truth geometry

An adversarial attack is a perturbed version $\pi(x)$ of an input example x , which alters the prediction of the classifier such that $f(\pi(x)) \neq f(x)$. According to this simplistic definition, every classifier can be successfully attacked, provided it has at least two possible output values. Taking a more refined perspective, consider $g(x)$ as the best possible classification (ground truth / Bayes classifier). Then a successful adversarial attack can be defined as a perturbation $\pi(x)$ of an input x such that $f(x) = g(x)$ but $f(\pi(x)) \neq g(\pi(x))$. This highlights that what we care about is not whether a classifier changes its prediction when the input is slightly perturbed, but rather in what scenario it changes its prediction.

In order to quantify the sensitivity to attacks, we need a measure of the size of the perturbation model and the effect that it has on the classifier. Consider a loss function of the form

$$E(f) = \mathbb{E}_{p(y|x)p(x)} [l(f(x), y)]. \quad (1)$$

We can measure the detriment of the loss when the data is perturbed in comparison with the unperturbed loss. For generality, we define a perturbation π at x as a random variable. For example, in the case of samples from a vector space, additive perturbations take the form $\pi = x + \xi$, where ξ can be,

¹Wasserstein Gaussians appear in the small time behavior of a process called Wasserstein diffusion, investigated in continuous (von Renesse and Sturm, 2009) and discrete (Li, 2018) states.

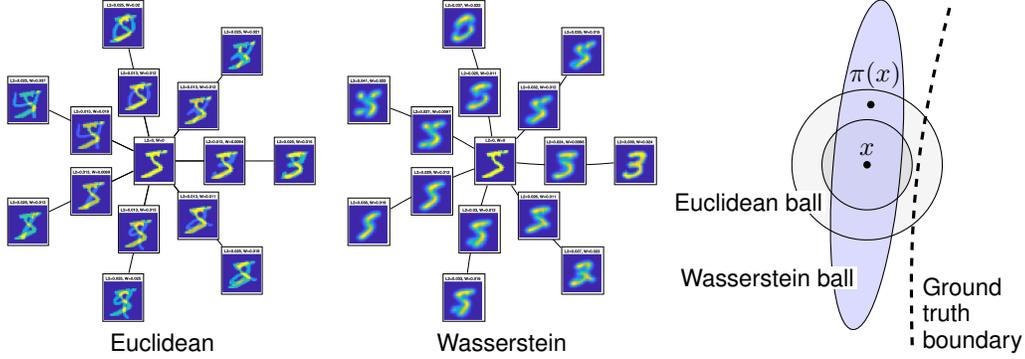


Figure 1: Shown are fixed-distance interpolates, measured by L^2 and Wasserstein metrics, between a source image and some other images in the MNIST data set. When training with smoothness regularizers towards achieving robustness against adversarial perturbations that are perceptually small, we need a good specification of semantically meaningful neighborhoods around which the discriminative function should not change much. As we see, the Wasserstein metric defines a more coherent neighborhood, with a natural interpretation in terms of continuous transportation of colors between pixels, which depends on the input example. The right part illustrates the intuition that an appropriate metric can capture larger in-class neighborhoods around an input image x , and hence allow us to apply stronger smoothness regularization against perturbations $\pi(x)$ without hurting test performance.

e.g., a zero mean Gaussian variable, or a deterministic value obtained by an attack strategy on the input x . The level of perturbation can be quantified, for example, in terms of the maximum or average size of the perturbations with respect to some norm, or in terms of the number of computations used to construct the attacks. The loss under perturbations is then

$$\mathbb{E}_{p(\pi|x)p(y|\pi(x))p(x)} [l(f(\pi(x)), y)]. \quad (2)$$

Unless we have access to the true labels given the perturbed inputs, i.e., the conditional probability $p(y|\pi(x))$ or the ground truth class $g(\pi(x))$, this measure is just theoretical. In practice the relation of labels to inputs is unknown, and for training we only have access to a training set $\{(x_i, y_i)\}_{i=1}^N$ and possibly the perturbation model π .

A natural approach to obtaining more robust classifiers is to train with perturbations. A special case is adversarial training, where the perturbations are constructed specifically to deceive the classifier. The problem with this approach is that typically the correct class labels for the perturbed input data are unknown, as mentioned above. Therefore, typically one considers simply

$$E_\pi(f) = \mathbb{E}_{p(\pi|x)p(y|x)p(x)} [l(f(\pi(x)), y)]. \quad (3)$$

For a given pair (x, y) with $y = g(x)$, in general $y \neq g(\pi(x))$. Therefore, the perturbation needs to be restricted in a way that ensures $g(\pi(x)) = g(x)$. The obvious and most common way to do this without using further information or prior knowledge about the ground truth, is to restrict the perturbations to be very small in some standard norm.

Adversarial examples are often constructed by minimizing the confidence of the discriminative function or increasing the training loss with respect to the input. Since this does not incorporate prior knowledge about the shape of the ground truth, usually the perturbations are restricted to lie within a very small L^p ball around the input example. Naive input noise regularization and gradient penalties suffer from a similar problem. When applied at the level that would be needed to prevent adversarial attacks, they tend to smear out the classifier in all directions around the input examples, leading to a significant detriment in test accuracy. A similar problem also arises in unsupervised adversarial training, where the base point is arbitrary (not necessarily a training example) and one requires that perturbations within a neighborhood are classified consistently. The situation is illustrated in the right part of Figure 1. Another difficulty is that data augmentation can be very expensive both in terms of the time it takes to compute each adversarial example and the number of examples that need to be added to the training data in order to obtain a sufficient level of robustness.

Instead of restricting the perturbations to be small in an L^p norm sense, we suggest to refine the metric on input space and the perturbation model. We propose to measure distances on input space using the Wasserstein metric and train with a corresponding Wasserstein Gaussian input noise. The Wasserstein metric assigns a small distance to natural local variations of an input image. This means that larger perturbations are more likely to remain within the class of the input example that is being perturbed. In turn, we can apply higher levels of noise, allow for larger size perturbations in adversarial training, or apply stronger gradient penalties during training. This is illustrated in Figure 1. In the next section we derive an effective regularizer for training with the Wasserstein metric on input space and which integrates the entire set of Wasserstein Gaussian noise perturbations (in a second order expansion) for each input example at once.

3 Perturbed loss and Wasserstein diffusion Tikhonov regularizer

It is well known that training with input noise can be related to training the original objective with an added penalty (Bishop, 1995). These derivations usually are based on Taylor expansion of the perturbed loss around a given example. By default, the inputs are considered to live in Euclidean space, with loss functions such as the mean square error or the log loss. Following the arguments from the previous sections, we model the input space of images as a Wasserstein space. We then derive the Wasserstein Taylor expansion of the perturbed loss and the corresponding regularization penalties. Once the input space is regarded as a Wasserstein metric space, our derivations follow Riemannian calculus therein.

We consider a perturbation model defined in terms of a ‘‘Wasserstein Gaussian’’, which at a given input image $x \in \mathcal{X} = \mathcal{P}(\Omega)$ has a density function of the form

$$p(\xi|x) = \exp(-d_W(x, x + \xi)^2/\eta^2)d(\xi),$$

with a scale parameter $\eta > 0$ and a given reference measure $d(\xi)$. Here the Wasserstein- q distance on image space can be defined in the linear programming formulation as

$$d_W(x, y) := \inf_{\Pi} (\mathbb{E}_{(a,b) \sim \Pi} [d_{\Omega}(a, b)^q])^{\frac{1}{q}}, \quad (4)$$

where Π is a joint distribution of pixel pairs (a, b) with marginals x, y . Here x, y are images viewed as histograms over the set of pixels Ω . The pixel ground metric $d_{\Omega}: \Omega \times \Omega \rightarrow \mathbb{R}_+$ assigns distances to pairs of pixels. The ground metric d_{Ω} can be defined in various ways that allow for efficient computations, and it can be trained from examples. We focus on the case $q = 2$. In this case, we can extract from the Wasserstein distance d_W a Riemannian metric for the space of images. Locally, the Wasserstein-2 distance can be expressed as

$$d_W(x, x + \xi)^2 = (\xi, G_W(x)\xi)_{L^2} + o(\|\xi\|^2), \quad (5)$$

where $\int \xi(a)da = 0$ and $G_W(x) = -(\nabla_a \cdot (x\nabla_a))^{-1}$ is the Wasserstein Riemannian metric tensor at x . Here ∇ and $\nabla \cdot$ are gradient and divergence operators in pixel space, with respect to $a \in \Omega$. In more details,

$$(\xi, G_W(x)\xi)_{L^2} = \int \|\nabla_a \Phi(a)\|^2 x(a)da,$$

where $\xi(a) = -\nabla \cdot (x(a)\nabla \Phi(a))$. Here the classical Wasserstein metric is defined on the space of images with equal total mass. This corresponds to the requirement $\int \xi(a)da = 0$. In general, we can apply the metric in unnormalized image spaces as well; see related studies on unnormalized optimal transport (Gangbo et al., 2019; Lee et al., 2019).

In practice we will consider a discrete pixel space $\Omega = \{1, \dots, n\}$. In this case, the Wasserstein Riemannian metric $G_W(x)$ is an $n \times n$ symmetric positive definite matrix that we will discuss further in the next section. Practically, the Wasserstein noise model corresponds to Gaussian noise with a covariance matrix $\eta^2 G_W^{-1}(x)$ that depends on the input x and the choice of a ground metric d_{Ω} over pixels.

We are now ready to present our theorem that relates training with Wasserstein diffusion to training with an added penalty term.

Theorem 1 (Perturbed loss regularization). *Consider an input space (\mathcal{X}, g) with the Riemannian metric g represented by a matrix $G(x)$ depending on $x \in \mathcal{X}$, and consider the loss $E(f) =$*

$\mathbb{E}[l(f(x), y)]$ from equation 1 with some error function l that is twice differentiable in the first argument. Let ξ be a Gaussian noise variable with zero mean and covariance matrix $\eta^2 G^{-1}(x)$ depending on x . Then the perturbed loss from equation 3 takes the form

$$E_\xi(f) = E(f) + \frac{1}{2}\eta^2 E^R(f) + o(\eta^2),$$

where

$$E^R(f) = \mathbb{E}_{p(y|x)p(x)} \left[l''(f(x), y) \|\nabla_g f(x)\|_g^2 + l'(f(x), y) \Delta_g f(x) \right].$$

Here l' and l'' denote the first and second order ordinary partial derivatives of l with respect to the first argument, and ∇_g , $\|\cdot\|_g$, Δ_g are the gradient, norm, and Laplace-Beltrami operators on (\mathcal{X}, g) .

The proof is provided in Appendix A. In this paper, we focus on $G(x) = G_W(x)$. This tells us that training with Wasserstein Gaussian noise corresponds, to second order in the noise level, to training with the unperturbed loss plus a penalty to the squared Wasserstein gradient norm and Laplace operators of the discriminative function.

We are also interested in non-zero mean perturbations, such as adversarial perturbations used in adversarial training. In this case the proof of the theorem yields the expansion

$$E_\xi(f) = E(f) + E^R(f) + O(\|\xi\|^2),$$

where

$$E^R(f) = \mathbb{E}_{p(y|x)p(x)} \left[l'(f(x), y) \cdot (\mathbb{E}_{p(\xi|x)}[\xi])^\top \nabla f(x) \right].$$

If the perturbation is deterministic, $\mathbb{E}_{p(\xi|x)}[\xi]$ is simply $\xi(x)$. This suggests to regularize by penalizing the directional derivative of the discriminative function in the direction of the perturbation. If the perturbation is proportional to the Euclidean steepest descent direction of the discriminative function, $\xi \propto \nabla f(x)$, then the penalty is proportional to $\|\nabla f\|^2$. If the perturbation is proportional to the Riemannian steepest descent direction, $\xi \propto \nabla_g f(x) = G^{-1}(x) \nabla f(x)$, then the penalty is proportional to the Riemannian gradient norm squared, $\|\nabla_g f(x)\|_g^2$.

Example 2 (Square error). For the square error $l(f(x), y) = (f(x) - y)^2$ and a perturbation model as in Theorem 1, we obtain the regularizer

$$E^R(f) = \mathbb{E}_{p(y|x)p(x)} \left[\|\nabla_g f(x)\|_g^2 + (f(x) - y) \Delta_g f(x) \right].$$

For non-zero mean perturbations, we can consider an expansion to first order which gives

$$E^R(f) = \mathbb{E}_{p(y|x)p(x)} \left[2(f(x) - y) \cdot \mathbb{E}_{p(\xi|x)}[\xi]^\top \nabla f(x) \right].$$

Example 3 (Cross entropy error). For the cross entropy $l(f(x), y) = -y \ln(f(x)) - (1 - y) \ln(1 - f(x))$ and a perturbation model as in Theorem 1, we obtain the regularizer

$$E^R(f) = \mathbb{E}_{p(y|x)p(x)} \left[\left(\frac{y}{f^2(x)} + \frac{1-y}{(1-f(x))^2} \right) \|\nabla_g f(x)\|_g^2 + \left(-\frac{y}{f(x)} + \frac{1-y}{(1-f(x))} \right) \Delta_g f(x) \right].$$

In the case of k outputs (e.g., k -class classification), the loss function is simply the sum of the loss for each output times $1/k$. For non-zero mean perturbations we obtain

$$E^R(f) = \mathbb{E}_{p(y|x)p(x)} \left[\left(-\frac{y}{f(x)} + \frac{1-y}{(1-f(x))} \right) \cdot \mathbb{E}_{p(\xi|x)}[\xi]^\top \nabla f(x) \right].$$

Example 4 (Euclidean inputs). In the case of Euclidean inputs and uncorrelated zero mean Gaussian noise of variance η^2 , we recover some of the classic calculations by Bishop (1995). Consider as an example the square error function, for which the regularizer becomes

$$E^R(f) = \mathbb{E}_{p(y|x)p(x)} \left[\sum_i \left\{ \left(\frac{\partial f}{\partial x_i} \right)^2 + (f(x) - y) \frac{\partial^2 f}{\partial x_i^2} \right\} \right].$$

As pointed out by Bishop (1995), this is the Tikhonov regularizer that is usually added to the sum of squares error.

Theorem 1 shows that all noise perturbed versions of a given input example can be integrated (in a second order sense) into a single term. Formally, equivalence of the regularizer to training with noise is only valid for small values of η , since it is based on a second order Taylor expansion. The Wasserstein diffusion smoothness regularizer E^R also has the natural interpretation as decreasing the variability of the classifier in an anisotropic and input dependent way that is captured by the Wasserstein gradient norm and the Laplace-Beltrami operator. This interpretation remains valid for arbitrarily large values of η , even if in this case the regularized objective might no longer correspond to the integrated perturbed objective.

We note that the term involving the Laplace-Beltrami operator is premultiplied with the derivative of the error. For regular choices of l , if the classifier makes correct predictions on the training inputs x (which is often the case), the derivative $l'(f(x), y)$ will vanish. This suggests that for the purpose of regularization in settings where the training error vanishes, we can omit the Laplace-Beltrami term and consider only the gradient penalty.

Taking the perspective of smoothness suggests that we may also regularize by penalizing the gradient of the discriminator, instead of the gradient of the loss function. Finally, we point out that the Wasserstein metric can also be used to define the size constraints for adversarial training. Usually adversarial perturbations are constrained to have L^∞ norm (or some L^p norm) bounded by a small ϵ . Instead of using $\|\xi\|_{L^p} \leq \epsilon$, we can use $\|\xi\|_W \leq \epsilon$, or simply $\xi^\top G_W(x)\xi \leq \epsilon$.

4 Training with the Wasserstein diffusion Tikhonov regularizer

In this section we describe the implementation of the regularization scheme introduced in Section 3. Theorem 1 suggests to replace the original objective function by a regularized objective $E + \eta^2 E^R$, where η corresponds to the strength of the regularization. For each training example x we add

$$l''(f(x), y) \cdot \|\nabla_W f(x)\|_W^2 + l'(f(x), y) \cdot \Delta_W f(x). \quad (6)$$

The first term of equation 6 involves the Wasserstein gradient of the discriminative function f with respect to the input. This type of calculation appeared recently in the context of Wasserstein Generative Adversarial Networks with Wasserstein ground metric (Dukler et al., 2019). We regard each input image as a histogram over pixels. Then we define a weighted graph $\mathcal{G} = (V, E, \omega)$, where the vertices $V = \{1, \dots, n\}$ correspond to pixels, edges E connect adjacent pixels, and ω is a symmetric matrix of weights ω_{ij} associated to the edges. The normalized volume form on node $i \in I$ is given by $d_i = \frac{\sum_{j \in N(i)} \omega_{ij}}{\sum_{i=1}^n \sum_{i' \in N(i)} \omega_{ii'}}$. For computational efficiency, we consider sparse graphs with a local connectivity structure that is invariant with respect to vertical and horizontal shifts in the pixel domain. In the experiments, we use local grids of radius $\text{rad} = 2, 4, 6, 8$ and constant weights.

The Laplacian matrix associated with the weighted graph \mathcal{G} is defined, depending on the input x , as

$$L(x)_{ij} = \begin{cases} \frac{1}{2} \sum_{k \in N(i)} \omega_{ik} \left(\frac{x_i}{d_i} + \frac{x_k}{d_k} \right), & \text{if } i = j \\ -\frac{1}{2} \omega_{ij} \left(\frac{x_i}{d_i} + \frac{x_j}{d_j} \right), & \text{if } j \in N(i) \\ 0, & \text{otherwise.} \end{cases}$$

The Wasserstein metric tensor is the matrix function given by the (pseudo) inverse of the weighted Laplacian operator,

$$G_W(x) = L(x)^{-1} \in \mathbb{R}^{n \times n}.$$

Written explicitly, the Wasserstein gradient norm squared is

$$\begin{aligned} \|\nabla_W f(x)\|_W^2 &= \nabla f(x)^\top G_W(x)^{-1} \nabla f(x) \\ &= \nabla f(x)^\top L(x) \nabla f(x) \\ &= \sum_{(i,j) \in E} \omega_{ij} \left(\frac{\partial}{\partial x_i} f(x) - \frac{\partial}{\partial x_j} f(x) \right)^2 \frac{x_i/d_i + x_j/d_j}{2}. \end{aligned} \quad (7)$$

An efficient implementation of equation 7 is described in Appendix C.

The second term of equation 6, the Wasserstein Laplace-Beltrami of the discriminative function, is computed as

$$\Delta_W f(x) = \text{tr}(L(x) \nabla^2 f(x)) + \nabla f(x)^\top L(x) \nabla \log \det(L(x))^{-\frac{1}{2}}.$$

Here $\det(L(x))^{-\frac{1}{2}}$ is the volume form for the Wasserstein manifold. Here $\det(L(x))$ is the product of non-singular eigenvalues. For well-posedness of the volume form, we consider a compact set in the interior of a finite dimensional probability simplex. In this case, $L(x)$ is a positive definite matrix, and $\det(L(x))^{-\frac{1}{2}}$ is well defined. In practice, we could also consider the reference density as the Lebesgue measure, which omits the volume term. Details on the implementation are provided in Appendix B. We leave a more systematic study and computation of the Wasserstein Laplace-Beltrami term for future work.

5 Experiments

In this section we present initial experimental results, leaving a more extensive experimental evaluation for future work. We evaluate the utility of Wasserstein smoothness regularization in terms of the robustness of the trained classifiers to small and large perturbations. We focus on regularization by the Wasserstein gradient norm penalty.

5.1 Stability to adversarial perturbations of the input data

In this experiment we test the effectiveness of the gradient penalty regularizer in terms of the test accuracy of the trained classifiers. We train a ResNet-20 on clean images from CIFAR-10 with gradient norm penalty computed under Euclidean and Wasserstein metrics. We run grid search for the regularization strength and the radius defining the ground metric on pixel space. The training error converges to zero in all cases. We consider two types of test data: the clean test data set (natural generalization) and the test data set with each test example perturbed by an adversarial attack (robust generalization). Following current literature, adversarial perturbations are computed by FGSM and I-FGSM (Goodfellow et al., 2015; Kurakin et al., 2017). More details on the implementation and hyperparameters are provided in Appendix D. Our results, reported in Table 1, compare well with the state of the art for CIFAR-10, where current works report robust test error for FGSM $\epsilon = 8/255$ of 37.52%, and for I-FGSM-20 $\alpha = 2/255, \epsilon = 8/255$ of 42.06% (Wang et al., 2018b).

Test data \ Regularizer	None	Euclidean grad.	Wasserstein grad.
Natural	16.29	15.61	15.35
FGSM $\epsilon = 8/255$	82.22	31.10	30.20
FGSM $\epsilon = 25/255$	89.72	66.83	44.32
I-FGSM-20 $\alpha = 2/255, \epsilon = 8/255$	90.15	40.06	32.12

Table 1: Robust test error percentage (lower is better) for a ResNet-20 network with softplus activation trained for 200 epochs on clean CIFAR-10 training images using gradient norm regularization with Euclidean and Wasserstein metric on image space. We run grid search over the regularization strength and the ground metric radius on pixel space.

5.2 Stability to large in-class variations of the input data

Most work on adversarial robustness focuses on small perturbations, with adversarial attacks restricted to have a small norm so that they remain imperceptible to humans. We are interested in generalization for all kinds of in-class variations of the data, including large perturbations that should not change the predicted class. In this experiment we train on the clean CIFAR-10 training set (no data augmentation), and compare between no regularization, Euclidean, and Wasserstein smoothness regularization. For testing, we randomly draw 1000 images from the natural test set and construct for each of them a sequence of translated versions with padding, as depicted in Figure 2. The semantic meaning of images should remain constant under relatively large translations, and therefore we expect a robust classifier to label all images in the sequence similarly. Quantitatively, this is measured by the number of label flips in the sequence. We report the average number of label flips over all sequences of test images in Table 2. As the table shows, Wasserstein gradient regularization improves the robustness of the classifier to translations.



Figure 2: Robust classifiers should be invariant to natural variations of the data. Shown are horizontal translations of an image from CIFAR-10.

Perturbation \ Regularizer	None	Euclidean grad.	Wasserstein grad.
Horizontal translation	10.009	7.898	6.488
Vertical translation	9.920	9.437	7.956

Table 2: Average number of prediction flips on sequences of translated test images from CIFAR-10. The classifiers were trained on the clean CIFAR-10 training set with no data augmentation, with either no regularization, Euclidean Tikhonov regularization, or Wasserstein Tikhonov regularization.

6 Related works

There are many works related to Wasserstein geometry, robustness, regularization. In this section we briefly mention some of the literature in relation to our discussion in this article.

Adversarial robustness. Previous works have investigated postprocessing with Jacobian regularization (Jakubovitz and Giryes, 2018) and cross Lipschitz regularization (Hein and Andriushchenko, 2017), whereby the input space was modeled as Euclidean space. The duality of attack norms and Lipschitz norms has been discussed as well (Finlay et al., 2019). Perturbation based regularization has been proposed (Yan et al., 2018), which penalizes the negative size of the deep fool attack in proportion to the size of the input. Gaussian data augmentation was proposed too (Zantedeschi et al., 2017), but evaluated by Monte Carlo samples and using Euclidean space. Recently, the tradeoff between natural and robust classification errors was studied, leading to a training objective with an added term of the form $\mathbb{E}_x[\max_{x' \in B(x, \epsilon)} \phi(f(x)f(x')/\lambda)]$ (Zhang et al., 2019). Similar to ours, this approach penalizes the variability of the classifier, but it is not incorporating priors about the geometry of the classes. While working on this article, we became aware of a work using modified Sinkhorn iterations to approximate projections of adversarial examples onto a Wasserstein ball (Wong et al., 2019). This is similar to the adversarial norm constraint that we suggested here. However, our approach is based on a Riemannian metric formulation, which allows us to obtain a very simple quadratic form approximation of the norm and also to integrate a generative noise model (Wasserstein diffusion) into an effective smoothness regularizer.

Wasserstein sample space. Optimal transport has been applied to the design of training objectives (Arjovsky et al., 2017). Recently this has been combined with a Wasserstein metric on image data space (Dukler et al., 2019). The Kantorovich duality of the Wasserstein training objective leads to a Lipschitz constraint on the discriminator networks, which itself is computed in Wasserstein space. This approach is known as the Wasserstein of Wasserstein GAN (WWGAN). Our regularization penalty derived from expanding a Wasserstein Gaussian noise variable in the input space of a classifier also includes a term that penalizes the Wasserstein gradient norm. However, our analysis has a different motivation and interpretation and also reveals higher order expansion terms.

Wasserstein diffusion. The noise in Wasserstein space has been studied in continuous (von Renesse and Sturm, 2009) and discrete (Li, 2018) state spaces. The present work seems to be the first to apply Wasserstein type noise in machine learning. In this trend, the definition and efficient computation of the Riemannian volume form on Wasserstein space remains an open problem for future work.

Wasserstein Information Geometry. The Wasserstein metric is gaining traction not only in the design of training objectives and in the definition of geometric structures on the sample space of generative models, but also in the development of natural gradients and optimizers (Li, 2018; Li and Montúfar, 2018; Li and Montúfar, 2018). Recently this has been applied to GANs (Lin et al., 2018). In this paper, we derive and apply second order calculus of Wasserstein geometry for improving generalization, especially improving the robustness to adversarial attacks.

Robustness and regularization. Wasserstein balls have appeared in the context of robust density estimation (Shafieezadeh-Abadeh et al., 2017), where they are also related to a form of Tikhonov

regularization. Wasserstein distributionally robust stochastic optimization has been related to regularization by certain empirical gradient norms (Gao et al., 2017). Close to our derivations, albeit not involving Wasserstein geometry, is the work by Bishop (1995), which shows that training with noise is equivalent to Tikhonov regularization.

7 Discussion

Training with input noise or data augmentation in general is known as an effective form of regularization to obtain classifiers that are more robust to natural variations of the data, or to reduce the sensitivity to perturbations. These methods usually have a high cost in terms of the number of examples needed and the cost of computing each of them (especially in the case of adversarial data augmentation obtained by iterated gradient methods). Another problem is that usually noise models and adversarial examples need to be restricted to tiny norm values to ensure that they remain within the class of the perturbed example. Smoothness regularizers based on L^p metrics are usually limited in the same way. In this paper we follow the idea that the space of inputs is not Euclidean and that smoothness priors should be implemented with respect to an appropriate metric, which in turn would allow us to apply higher levels of regularization without hurting test performance. We propose to use the Wasserstein-2 metric to capture semantically meaningful neighborhoods of images. As we show, the Wasserstein diffusion smoothness regularizer arises naturally by expanding and integrating the loss with respect to Wasserstein Gaussian noise on the inputs. We obtain an effective penalty that can be computed very efficiently, saving computation compared with adversarial data augmentation, and has a negligible overhead over L^2 gradient penalties. Preliminary experimental results indicate that our methods can significantly improve robust generalization performance on CIFAR-10. We obtain models that are robust not only to small perturbations (the usual setting in adversarial robustness literature), but also to large scale in-class perturbations, such as translations. We think that this is conceptually an important step towards learning models that generalize better in relation to all types of natural variations of the input data, not only small perturbations.

Acknowledgments

This research has received funding from AFOSR MURI FA9550-18-1-0502. YD has received funding from the National Science Foundation Graduate Research Fellowship under Grant No. DGE-1650604. Part of this research was performed at the Institute for Pure and Applied Mathematics (IPAM), which is supported by the National Science Foundation (Grant No. DMS-1440415). This project has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement n° 757983).

References

- M. Arjovsky, S. Chintala, and L. Bottou. Wasserstein GAN. *arXiv:1701.07875 [cs, stat]*, 2017.
- C. M. Bishop. Training with noise is equivalent to Tikhonov regularization. *Neural Computation*, 7(1):108–116, 1995.
- N. Carlini, G. Katz, C. Barrett, and D. L. Dill. Provably minimally-distorted adversarial examples. *CoRR*, abs/1709.10207, 2017. URL <http://arxiv.org/abs/1709.10207>.
- S.-N. Chow, W. Huang, Y. Li, and H. Zhou. Fokker–Planck Equations for a Free Energy Functional or Markov Process on a Graph. *Archive for Rational Mechanics and Analysis*, 203(3):969–1008, 2012.
- S.-N. Chow, W. Li, and H. Zhou. Entropy dissipation of Fokker-Planck equations on graphs. *Discrete & Continuous Dynamical Systems, series A*, 2018.
- M. Cissé, P. Bojanowski, E. Grave, Y. Dauphin, and N. Usunier. Parseval networks: Improving robustness to adversarial examples. In *ICML 34*, pages 854–863, 2017.
- Y. Dukler, W. Li, A. T. Lin, and G. Montúfar. Wasserstein of Wasserstein loss for learning generative models. In *ICML 36*, 2019.

- C. Finlay, A. M. Oberman, and B. Abbasi. Improved robustness to adversarial examples using lipschitz regularization of the loss, 2019. URL <https://openreview.net/forum?id=HkxAisC9FQ>.
- W. Gangbo, W. Li, S. Osher, and M. Puthawala. Unnormalized Optimal Transport. *arXiv:1902.03367 [math]*, 2019.
- R. Gao, X. Chen, and A. J. Kleywegt. Wasserstein distributional robustness and regularization in statistical learning. *arXiv preprint arXiv:1712.06050*, 2017.
- I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. In *ICLR*, 2015.
- C. Guo, M. Rana, M. Cisse, and L. van der Maaten. Countering adversarial images using input transformations. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=SyJ7C1WCb>.
- M. Hein and M. Andriushchenko. Formal guarantees on the robustness of a classifier against adversarial manipulation. In *Advances in Neural Information Processing Systems 30*, pages 2266–2276. Curran Associates, Inc., 2017.
- D. Jakobovitz and R. Giryes. Improving DNN robustness to adversarial attacks using jacobian regularization. *CoRR*, abs/1803.08680, 2018. URL <http://arxiv.org/abs/1803.08680>.
- A. Kurakin, I. J. Goodfellow, and S. Bengio. Adversarial machine learning at scale. In *ICLR*. OpenReview.net, 2017.
- W. Lee, R. Lai, W. Li, and S. Osher. Fast algorithms for generalized unnormalized optimal transport. *CAM 19-40*, 2019.
- W. Li. Geometry of probability simplex via optimal transport. *arXiv:1803.06360 [math]*, 2018.
- W. Li and G. Montúfar. Natural gradient via optimal transport. *Information Geometry*, 1(2):181–214, 2018.
- W. Li and G. Montúfar. Ricci curvature for parametric statistics via optimal transport. *arXiv:1807.07095*, 2018.
- F. Liao, M. Liang, Y. Dong, T. Pang, X. Hu, and J. Zhu. Defense against adversarial attacks using high-level representation guided denoiser. In *CVPR*, pages 1778–1787, 2018.
- A. Lin, W. Li, S. Osher, and G. Montúfar. Wasserstein proximal of GANs. *CAM reports 18-53*, 2018.
- J. Maas. Gradient flows of the entropy for finite Markov chains. *Journal of Functional Analysis*, 261(8):2250–2292, 2011.
- A. Mielke. A gradient structure for reaction–diffusion systems and for energy-drift-diffusion systems. *Nonlinearity*, 24(4):1329, 2011.
- S. Moosavi-Dezfooli, A. Fawzi, and P. Frossard. Deepfool: A simple and accurate method to fool deep neural networks. In *CVPR*, pages 2574–2582. IEEE Computer Society, 2016.
- S. Moosavi-Dezfooli, A. Fawzi, O. Fawzi, and P. Frossard. Universal adversarial perturbations. In *CVPR*, pages 86–94. IEEE Computer Society, 2017.
- A. M. Nguyen, J. Yosinski, and J. Clune. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In *CVPR*, pages 427–436, 2015.
- G. Peyré, M. Cuturi, et al. Computational optimal transport. *Foundations and Trends® in Machine Learning*, 11(5-6):355–607, 2019.
- Y. Rubner, C. Tomasi, and L. J. Guibas. The earth mover’s distance as a metric for image retrieval. *Int. J. Comput. Vision*, 40(2):99–121, 2000.
- H. Salman, G. Yang, H. Zhang, C.-J. Hsieh, and P. Zhang. A convex relaxation barrier to tight robustness verification of neural networks. *ArXiv*, abs/1902.08722, 2019.

- P. Samangouei, M. Kabkab, and R. Chellappa. Defense-GAN: Protecting classifiers against adversarial attacks using generative models. In *ICLR*, 2018. URL <https://openreview.net/forum?id=BkJ3ibb0->.
- A. Shafahi, W. R. Huang, C. Studer, S. Feizi, and T. Goldstein. Are adversarial examples inevitable? In *ICLR*, 2019. URL <https://openreview.net/forum?id=r1lWUoA9FQ>.
- S. Shafieezadeh-Abadeh, D. Kuhn, and P. M. Esfahani. Regularization via mass transportation. *arXiv preprint arXiv:1710.10016*, 2017.
- J. Solomon, R. Rustamov, L. Guibas, and A. Butscher. Earth mover’s distances on discrete surfaces. *ACM Transactions on Graphics (TOG)*, 33(4):67, 2014.
- J. Solomon, F. De Goes, G. Peyré, M. Cuturi, A. Butscher, A. Nguyen, T. Du, and L. Guibas. Convolutional Wasserstein distances: Efficient optimal transportation on geometric domains. *ACM Transactions on Graphics (TOG)*, 34(4):66, 2015.
- C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus. Intriguing properties of neural networks. In *ICLR*, 2014. URL <http://arxiv.org/abs/1312.6199>.
- M.-K. von Renesse and K.-T. Sturm. Entropic measure and Wasserstein diffusion. *Ann. Probab.*, 37(3):1114–1191, 2009.
- B. Wang, A. Lin, Z. Shi, W. Zhu, P. Yin, A. L. Bertozzi, and S. J. Osher. Adversarial defense via data dependent activation function and total variation minimization. *CoRR*, abs/1809.08516, 2018a.
- B. Wang, B. Yuan, Z. Shi, and S. J. Osher. EnResNet: ResNet ensemble via the Feynman-Kac formalism. *CoRR*, abs/1811.10745, 2018b. URL <http://arxiv.org/abs/1811.10745>.
- E. Wong, F. Schmidt, and Z. Kolter. Wasserstein adversarial examples via projected Sinkhorn iterations. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 6808–6817, Long Beach, California, USA, 09–15 Jun 2019. PMLR. URL <http://proceedings.mlr.press/v97/wong19a.html>.
- Z. Yan, Y. Guo, and C. Zhang. Deep defense: Training dnns with improved adversarial robustness. In *Advances in Neural Information Processing Systems 31*, pages 419–428. Curran Associates, Inc., 2018.
- V. Zantedeschi, M.-I. Nicolae, and A. Rawat. Efficient defenses against adversarial attacks. In *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security, AISeC ’17*, pages 39–49, New York, NY, USA, 2017. ACM. ISBN 978-1-4503-5202-4.
- H. Zhang, Y. Yu, J. Jiao, E. P. Xing, L. E. Ghaoui, and M. I. Jordan. Theoretically principled trade-off between robustness and accuracy. *CoRR*, abs/1901.08573, 2019.

A Proof of Theorem 1

A Riemannian metric g defines an inner product between tangent vectors of the input space at each possible location. We choose standard coordinates for the input space $\mathcal{X} = \mathbb{R}^n$ and write $g(\xi, \zeta) = (\xi, \zeta)_g = \xi^\top G(x)\zeta$ for any pair $\xi, \zeta \in T_x\mathcal{X}$. We implicitly identify $T_x\mathcal{X}$ and \mathcal{X} so that adding a tangent vector $\xi \in T_x\mathcal{X}$ to an input vector $x \in \mathcal{X}$ makes sense. The Riemannian gradient with the metric g is given by $\nabla_g f(x) = G^{-1}(x)\nabla f(x)$, where ∇ is the ordinary gradient. This is also known as the natural gradient.

Proof of Theorem 1. We expand the error function l around a data point x with added noise ξ in the Riemannian space (\mathcal{X}, g) . We obtain

$$\begin{aligned}
 l(f(x + \xi), y) &= l(f(x), y) + l'(f(x), y)(\nabla_g f(x), \xi)_g \\
 &\quad + \frac{1}{2}l''(f(x), y)(\nabla_g f(x), \xi)_g^2 + \frac{1}{2}l'(f(x), y) \sum_{i,j} \xi_i \xi_j (\nabla_g^2 f(x))_{ij} + o(\|\xi\|_g^2).
 \end{aligned}$$

We discuss the individual terms in turn. The zero order term is just the unperturbed loss. On taking the expectation value with respect to ξ given x , the linear term vanishes when we assume that the perturbations have zero mean, $\mathbb{E}_{p(\xi|x)}[\xi] = 0$. If the perturbation does not have zero mean, we obtain

$$\mathbb{E}_{p(\xi|x)}[(\nabla_g f(x), \xi)_g] = \mathbb{E}_{p(\xi|x)}[(G^{-1}(x)\nabla f(x))^\top G(x)\xi] = \nabla f(x)^\top \mathbb{E}_{p(\xi|x)}[\xi].$$

For the first quadratic term, when $\mathbb{E}_{p(\xi|x)}[\xi\xi^\top] = \eta^2 G^{-1}(x)$, we obtain

$$\mathbb{E}_{p(\xi|x)}[(\nabla_g f(x)^\top G(x)\xi)^2] = \eta^2 \nabla_g f(x)^\top G(x)\nabla_g f(x) = \eta^2 \|\nabla_g f\|_g^2.$$

For the second quadratic term, again when $\mathbb{E}_{p(\xi|x)}[\xi\xi^\top] = \eta^2 G^{-1}(x)$, we obtain

$$\mathbb{E}_{p(\xi|x)}[\xi^\top \text{Hess } f(x)\xi] = \eta^2 \Delta_g f(x).$$

Here the Laplace-Beltrami operator is

$$\Delta_g f = \sum_{j,k} g^{jk} \frac{\partial^2 f}{\partial x^j \partial x^k} - g^{jk} \Gamma_{jk}^l \frac{\partial f}{\partial x^l},$$

where Γ_{jk}^l is the Christoffel symbol. □

B Riemannian calculus in Wasserstein space over discrete states

In this section, we review the definition of Wasserstein-2 Riemannian metric on discrete states founded in (Chow et al., 2012; Maas, 2011; Mielke, 2011) and developed in (Chow et al., 2018) and (Li and Montúfar, 2018). See Wasserstein Riemannian calculus in (Li, 2018). We also discuss practical implementations of the Laplace-Beltrami operator appearing in the Wasserstein Tikhonov regularizer.

Consider $I = \{1, \dots, n\}$. The probability simplex on I is the set

$$\mathcal{P}_+(I) = \left\{ (x_1, \dots, x_n) \in \mathbb{R}^n : \sum_{i=1}^n x_i = 1, \quad x_i > 0 \right\}.$$

Here $x = (x_1, \dots, x_n)$ is a probability vector with coordinates x_i corresponding to the probabilities assigned to each node $i \in I$.

We next define the Wasserstein-2 metric tensor on $\mathcal{P}_+(I)$. This is given in terms of a undirected graph with weighted edges, $\mathcal{G} = (I, E, \omega)$, where I is the vertex set, $E \subseteq \binom{I}{2}$ is the edge set, and $\omega = (\omega_{ij})_{i,j \in I} \in \mathbb{R}^{n \times n}$ is a matrix of edge weights satisfying

$$\omega_{ij} = \begin{cases} \omega_{ji} > 0, & \text{if } (i, j) \in E \\ 0, & \text{otherwise} \end{cases}.$$

The set of neighbors (adjacent vertices) of i is denoted by $N(i) = \{j \in V : (i, j) \in E\}$. The normalized volume form on node $i \in I$ is given by $d_i = \frac{\sum_{j \in N(i)} \omega_{ij}}{\sum_{i=1}^n \sum_{i' \in N(i)} \omega_{ii'}}$.

The graph structure $\mathcal{G} = (I, E, \omega)$ induces a graph Laplacian matrix function.

Definition 1 (Weighted Laplacian matrix). Given an undirected weighted graph $\mathcal{G} = (I, E, \omega)$, with $I = \{1, \dots, n\}$, the matrix function $L(\cdot) : \mathbb{R}^n \rightarrow \mathbb{R}^{n \times n}$ is defined as

$$L(p) = D^\top \Lambda(x) D, \quad x = (x_i)_{i=1}^n \in \mathbb{R}^n,$$

where

- $D \in \mathbb{R}^{|E| \times n}$ is the discrete gradient operator given by

$$D_{(i,j) \in E, k \in V} = \begin{cases} \sqrt{\omega_{ij}}, & \text{if } i = k, i > j \\ -\sqrt{\omega_{ij}}, & \text{if } j = k, i > j \\ 0, & \text{otherwise,} \end{cases}$$

- $-D^\top \in \mathbb{R}^{n \times |E|}$ is the oriented incidence matrix, and
- $\Lambda(x) \in \mathbb{R}^{|E| \times |E|}$ is a weight matrix depending on x ,

$$\Lambda(x)_{(i,j) \in E, (k,l) \in E} = \begin{cases} \frac{1}{2} \left(\frac{1}{d_i} x_i + \frac{1}{d_j} x_j \right), & \text{if } (i,j) = (k,l) \in E \\ 0, & \text{otherwise.} \end{cases}$$

We are now ready to present the Wasserstein-2 metric tensor. Consider the tangent space of $\mathcal{P}_+(I)$ at x ,

$$T_x \mathcal{P}_+(I) = \left\{ (\sigma_i)_{i=1}^n \in \mathbb{R}^n : \sum_{i=1}^n \sigma_i = 0 \right\}.$$

Definition 2 (Wasserstein-2 metric tensor). The inner product $g : T_x \mathcal{P}_+(I) \times T_x \mathcal{P}_+(I) \rightarrow \mathbb{R}$ takes any two tangent vectors $\sigma_1, \sigma_2 \in T_x \mathcal{P}_+(I)$ to

$$g_x(\sigma_1, \sigma_2) := \sigma_1^\top L(x)^{-1} \sigma_2, \quad \text{for any } \sigma_1, \sigma_2 \in T_x \mathcal{P}_+(I),$$

where $L(x)^{-1}$ is the pseudo inverse of $L(x)$.

Here $L(x)_{ij}^{-1}$ plays the role of g_{ij} in the last subsection A. Using this metric tensor, the Riemannian calculus in $(\mathcal{P}_+(I), g)$ has the following formulations.

(i) The Christoffel symbol $\Gamma_x^W : T_x \mathcal{P}_+(I) \times T_x \mathcal{P}_+(I) \rightarrow T_x \mathcal{P}_+(I)$ forms:

$$\Gamma_x^W(\sigma_1, \sigma_2) = -\frac{1}{2} [L(\sigma_1)L(\rho)^{-1}\sigma_2 + L(\sigma_2)L(\rho)^{-1}\sigma_1] + \frac{1}{2} L(\rho) \left(\nabla_G L(\rho)^{-1} \sigma_1 \circ \nabla_G L(\rho)^{-1} \sigma_2 \right),$$

where $\sigma_1, \sigma_2 \in T_x \mathcal{P}_+(I)$ and

$$\left(\nabla_G L(\rho)^{-1} \sigma_1 \circ \nabla_G L(\rho)^{-1} \sigma_2 \right) = \left(\frac{1}{2d_i} \sum_{j \in N(i)} (\nabla_{ij} L(\rho)^{-1} \sigma_1) (\nabla_{ij} L(\rho)^{-1} \sigma_2) \right)_{i=1}^n \in \mathbb{R}^n.$$

(ii) Given $F \in C^\infty(\mathcal{P}_+(I))$, the Riemannian gradient is

$$\text{grad}_W F(x) = \left(L(x)^{-1} \right)^{-1} \nabla F(x) = L(x) \nabla F(x),$$

where ∇ is the Euclidean L^2 derivative w.r.t. x .

(iii) The Riemannian Hessian operator $\text{Hess}_W F(\rho) : T_\rho \mathcal{P}_+(M) \times T_\rho \mathcal{P}_+(M) \rightarrow \mathbb{R}$ is given by

$$\text{Hess}_W F(x)(\sigma_1, \sigma_2) = \sigma_1^\top \nabla^2 F(x) \sigma_2 + \sum_{i=1}^n \Gamma_x^W(\sigma_1, \sigma_2)_i \nabla_{x_i} F(x),$$

where $\sigma_1, \sigma_2 \in T_x \mathcal{P}_+(I)$.

(iv) The Riemannian volume is given by

$$\text{vol}_W(x) = \det(L(x))^{-\frac{1}{2}} = \prod_{i=1}^{n-1} \lambda_i(x)^{-\frac{1}{2}},$$

where $\lambda_i(x)$ are the positive eigenvalues of $L(x)$.

(v) The Laplacian–Beltrami operator is given by

$$\Delta_W F(x) = \text{tr} \left(L(x) \left(\frac{\partial^2 f}{\partial x_i \partial x_j} \right)_{1 \leq i, j \leq n} \right) + \nabla f(x)^\top L(x) \nabla \log \det(L(x))^{-\frac{1}{2}}.$$

If the reference measure is a Lebesgue measure in a simplex, the modified Wasserstein Laplacian operator satisfies

$$\begin{aligned} \tilde{\Delta}_W f(x) &= \text{tr} \left(L(x) \left(\frac{\partial^2 f}{\partial x_i \partial x_j} \right)_{1 \leq i, j \leq n} \right) \\ &= \sum_{(i,j) \in E} \omega_{ij} \left(\frac{\partial^2}{\partial x_i^2} f(x) - 2 \frac{\partial^2}{\partial x_i \partial x_j} f(x) + \frac{\partial^2}{\partial x_j^2} f(x) \right) \frac{x_i/d_i + x_j/d_j}{2}. \end{aligned} \quad (8)$$

C Efficient implementation of the Wasserstein gradient norm

To compute equation 7 in practice, we define a suitable similarity graph $\mathcal{G} = (V, E, \omega)$ for the space of images, displaying translation invariance and symmetries. First, there is the invariance with respect to pixel translations. Symmetries arise since the distance from a pixel to two spatially opposite pixels is equal. In addition, each term in the sum in equation 7 can be decomposed into the product of linear relations between values in nodes i and j . Each linear relation (e.g $\nabla_{X_i} f - \nabla_{X_j} f$) is computed via a convolution, in this case on $\nabla_X F$. Convolutions may replace a linear product owing it to the described symmetries and invariances. For each relative neighbor direction we define a convolutional filter with the number of filters equal to the number of neighbors. For the truncated similarity graph for the Wasserstein distance, the edge set E is sparse and the number of convolutional filters is reduced considerably from n^2 . We therefore can calculate all pairs $\nabla_{X_i} f - \nabla_{X_j} f$ with a given relative neighbor relation by passing X and $\nabla_X f$ through of kernels $K_{\mathcal{O}_1} \dots K_{\mathcal{O}_d}$. In this case a neighbor relation, is defined as the geometrical pattern between two pixels. A pixel located in position $(10, 10)$ satisfies a neighbor relation $(1, 2)$ with a pixel in location $(11, 12)$. The neighbor relation is indeed invariant to the global position of the pixel which allows for the use of convolutions. Given a ground metric, we enumerate all non-zero neighbor relations as $\mathcal{O}_1, \dots \mathcal{O}_d$, for truncated distances this is much smaller than the complete edge graph. For each neighbor relation \mathcal{O}_k we associate a zero-valued kernel that equals 1 and -1 in the corresponding \mathcal{O}_k pixels, we denote the gradient kernels as $K_{\mathcal{O}_k}$. Likewise, we apply the same \mathcal{O}_k pattern, now with $\frac{1}{2}, \frac{1}{2}$ in the corresponding neighbor pattern pixels to obtain the terms $\frac{X_i/d_i + X_j/d_j}{2}$. For each i, j we denote the input kernels as $M_{\mathcal{O}_k}$. Applying entry-wise multiplication (\odot) and a summation collapsing all pixel locations and channels then yields an efficient and general method of calculating the Wasserstein gradient $\|\text{grad } f\|_{W_{2,d}(\Omega)}$ for general local cost metrics on highly optimized convolution.

Algorithm 1 Wasserstein gradient norm $\|\text{grad } f(X)\|_W^2$.

Require: The pixel graph $\mathcal{G} = (V, E, \omega)$; local weights (w_{ij}) ; neighbor relation tuples arranged symmetrically $\mathcal{O}_1 \dots \mathcal{O}_d$

Require: Euclidean gradient $\nabla_X f$

- 1: *Wasserstein-grad* $\leftarrow 0$
- 2: **for** neighbor relations $k = 1, \dots, d$ **do**
- 3: Build kernel $K_{\mathcal{O}_k}$ to compute $\nabla_{X_i} f - \nabla_{X_{\mathcal{O}_k(i)}} f$
- 4: Build corresponding kernel $M_{\mathcal{O}_k}$ to compute $\frac{X_i}{2d_i} + \frac{X_{\mathcal{O}_k}}{2d_{\mathcal{O}_k}}$
- 5: $H \leftarrow K_{\mathcal{O}_k}(\nabla_X f)$
- 6: $V \leftarrow M_{\mathcal{O}_k}(X)$
- 7: $H \leftarrow H \odot H$ (entry-wise multiplication)
- 8: $W \leftarrow H \odot V$
- 9: *Wasserstein-grad* \leftarrow *Wasserstein-grad* + $\text{sum}(W)$
- 10: **end for**
- 11: **Return** $\|\text{grad } f(X)\|_W^2 =$ *Wasserstein-grad*

D Details on the experimental setup

For our experiments, we use the CIFAR-10 dataset, and perform white-box attacks on the ResNet20 network. For training, we fixed the batch size of 128, and used SGD with momentum and weight decay, where the momentum value is 0.9 and the weight-decay value is 10^{-4} . We start with a learning rate of 0.1, and at epoch 100 and 150 we divide the learning rate by 10 each time.

We examine the case of training the ResNet-20 network on the CIFAR-10 dataset, where the only data augmentation performed is normalization. This achieves a test accuracy of 83.71%. We then examine the effect of modifying the loss objective with either the Euclidean or Wasserstein gradient penalties of the original loss, namely we use the loss

$$\ell(f(x), y) + \eta^2 (\nabla_x \ell(f(x), y), G(x)^{-1} \nabla_x \ell(f(x), y)),$$

where ∇_x is the Euclidean gradient and $G(x) \in \mathbb{R}^{d \times d}$ represents the metric used in sample space. For the Wasserstein gradient norm, $G(x)^{-1} = L(x)$. For the Euclidean gradient norm, $G(x) = I$.