

# LIPSCHITZ REGULARIZED DEEP NEURAL NETWORKS CONVERGE AND GENERALIZE

ADAM M. OBERMAN AND JEFF CALDER

ABSTRACT. Generalization of deep neural networks (DNNs) is an open problem which, if solved, could impact the reliability and verification of deep neural network architectures. In this paper, we show that if the usual fidelity term used in training DNNs is augmented by a Lipschitz regularization term, then the networks converge and generalize. The convergence is in the limit as the number of data points,  $n \rightarrow \infty$ , while also allowing the network to grow as needed to fit the data. Two regimes are identified: in the case of clean labels, we prove convergence to the label function which corresponds to zero loss, in the case of corrupted labels which we prove convergence to a regularized label function which is the solution of a limiting variational problem. In both cases, a convergence rate is also provided.

## 1. INTRODUCTION

While deep neural networks (DNNs) give more accurate predictions than other machine learning methods (LeCun *et al.*, 2015), they lack some of the performance guarantees of these methods. One step towards performance guarantees for DNNs is a proof of convergence with a rate, which could lead to quantitative error estimates. In this paper, we present such a result, for Lipschitz regularized DNNs. A proof of generalization follows as a consequence of convergence.

We begin by establishing the notation for our problem. We consider the classification problem to fix ideas, although regularization can apply to other problems as well.

**Definition 1.1.** Assume the data is normalized so that the data space is  $X = [0, 1]^d$ . Write  $\mathcal{D}_n = x_1, \dots, x_n$  for the training data. Assume  $\mathcal{D}_n$  is a sequence of *i.i.d.* random variables on  $X$  sampled from the probability distribution  $\rho$ . We consider the classification problem with  $m$  labels which are imbedded into the probability simplex, the label space,  $Y \subset \mathbb{R}^m$ . Write  $u_0 : X \rightarrow Y$  for the map from data to label space, so that  $y_i = u_0(x_i)$ .

Our results point towards improved generalization results using Lipschitz regularization, which we define now.

**Definition 1.2.** Choose norms  $\|\cdot\|_Y$ , and  $\|\cdot\|_X$  on  $X$  and  $Y$ , respectively. The Lipschitz constant (in these norms) of a function  $u : X_0 \subset X \rightarrow Y$  is given by

$$\text{Lip}(u; X_0) = \sup_{x_1, x_2 \in X_0} \frac{\|u(x_1) - u(x_2)\|_Y}{\|x_1 - x_2\|_X}$$

---

*Date:* October 3, 2018.

research supported by: AFOSR FA9550-18-1-0167 (A.O.). NSF-DMS 1713691 (J.C.).

When  $X_0$  is all of  $X$ , we write  $\text{Lip}(u; X) = \text{Lip}(u)$ . The Lipschitz constant of the data is  $\text{Lip}(u_0; \mathcal{D}_n)$ .

Write  $u(x; w)$  for the last layer of the network.<sup>1</sup> We consider the variational problem with Lipschitz regularization

$$(1) \quad \min_{u: X \rightarrow Y} J^n[u] = \frac{1}{n} \sum_{i=1}^n \ell(u(x_i; w), y_i) + \lambda \max(\text{Lip}(u) - L_0, 0)$$

The first term in (1) is the usual averaged loss on the training data  $\mathcal{D}_n$ . The second term in (1) the Lipschitz regularization term: the excess Lipschitz constant of the map  $u$ , compared to the constant  $L_0$ . Choosing  $\lambda > 0$  in (1) introduces a Lipschitz regularization penalty to the standard learning problem, which corresponds to  $\lambda = 0$ . In theory we take  $L_0 = \text{Lip}(u_0)$ , the Lipschitz constant of the data on the whole data manifold (discussed below).

*Remark 1.3.* In practice,  $\text{Lip}(u_0)$  can be estimated by the Lipschitz constant of the data,  $\text{Lip}(u_0; \mathcal{D}_n)$ . In fact, for many common data sets, the Lipschitz constant is very small and the Lipschitz constants of networks is much larger (Finlay & Oberman, 2018), so in practice we can set  $L_0 = 0$ . However for clean data, if we wish to recover  $u_0$  exactly, we need  $L_0 > 0$ .

If  $u$  is differentiable<sup>2</sup>, then the boundedness of  $X$  implies

$$\text{Lip}(u; X) = \max_{x \in X} \|\nabla u(x; w)\|_{X, Y}$$

where  $\|M\|_{X, Y}$  is the matrix norm induced by the norms on  $X$  and  $Y$ . In practise, we take the  $\infty$ -norm for  $Y$  and the 2-norm for  $X$ , and use explicit formulas for the  $\|M\|_{\infty, 2}$  norm, see §4. The Lipschitz constant of  $u$ ,  $\text{Lip}(u; X)$ , can thus be estimated from below by the maximum of the norm of the gradient on a minibatch

$$(2) \quad \max_{i \in I} \|\nabla_x u(x_i; w)\|_{X, Y} \leq \text{Lip}(u; X)$$

Our analysis will apply to the problem (1) which is *convex* in  $u$ , and does not depend explicitly on the weights,  $w$ . Of course, once  $u$  is restricted to a fixed neural network architecture, the corresponding minimization problem becomes non-convex in the weights. Our analysis can avoid the dependence on the weights because we make the assumption that there are enough parameters so that  $u$  can perfectly fit the training data. The assumption is justified by Zhang *et al.* (2016). As we send  $n \rightarrow \infty$  for convergence, we require that the network also grow, in order to continue to satisfy this assumption. Our results apply to other non-parametric methods in this regime.

In Figure 1 we illustrate the solution of (1) (with  $L_0 = 0$ ), using synthetic one dimensional data. In this case, the labels  $\{-1, 0, 1\}$  are embedded naturally into  $Y = \mathbb{R}$ , and  $\lambda = 0.1$ . Notice that the solution matches the labels exactly on a subset of the data. In the second part of the figure, we show a solution with corrupted labels which introduce a large Lipschitz constant, in this case, the solution reduces the Lipschitz constant, thereby correcting the errors.

<sup>1</sup>We apologize for not using the standard notation  $f$  for the last layer!

<sup>2</sup>We follow the common practise and treat the the architecture as differentiable for training purposes. To be rigorous we can use the  $L^\infty$  norm and appeal to Rachevacher's Theorem.

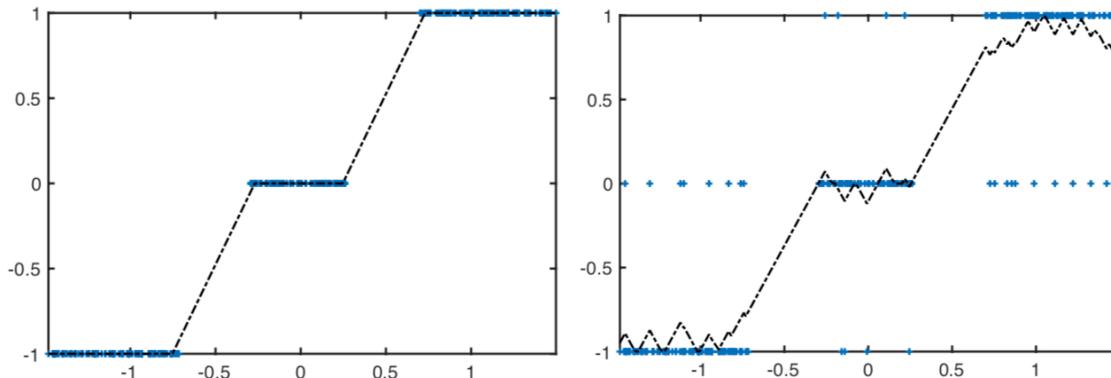


FIGURE 1. Synthetic labelled data and Lipschitz regularized solution  $u$ . Left: The solution value matches the labels perfectly on a large portion of the data set. Right: 10% of the data is corrupted by incorrect labels; the regularized solution corrects the errors.

**1.1. Related work and applications.** Generalization bounds have been obtained previously by using the Lipschitz constant of a network (Bartlett, 1997), as well as by using more general stability results (Bousquet & Elisseeff, 2002). More recently, (Bartlett *et al.*, 2017) proposed the Lipschitz constant of the network as a candidate measure for the Rademacher complexity, which a measure of generalization (Shalev-Shwartz & Ben-David, 2014, Chapter 26). However, our analysis is more direct and self-contained, and unlike other recent contributions such as (Hardt *et al.*, 2015), it does not depend on the training method.

The estimate of  $\text{Lip}(u; X)$  provided by (2) can be quite different from the the Tychonoff gradient regularization (Drucker & Le Cun, 1992),

$$\frac{1}{|I|} \sum_{i \in I} \|\nabla_x u(x_i)\|^2$$

since (2) corresponds to a maximum of the values of the norms, and the previous equation corresponds to average of the values. In fact, recent work on semi-supervised learning suggests that higher  $p$ -norms of the gradient are needed for generalization when the data manifold is not well approximated by the data (El Alaoui *et al.*, 2016; Calder, 2017; Kyng *et al.*, 2015; Slepcev & Thorpe, 2017). In Figure 2 we compare to the problems in Figure 1 using Tychonoff regularization. The Tychonoff regularization is less effective at correcting errors. The effect is more pronounced in higher dimensions.

An upper bound for the Lipschitz constant of the model is given by the norm of the product of the weight matrices (Szegedy *et al.*, 2013, Section 4.3). Let  $w = (w^1, \dots, w^J)$  be the weight matrices for each layer. Then

$$(3) \quad \text{Lip}(u; X) \leq \prod_{j=1}^J \|w^j\|.$$

Regularization of the network using methods based on (3) has been implemented recently in (Gouk *et al.*, 2018) and (Yoshida & Miyato, 2017). Because the upper bound in (3) does not take into account the coefficients in weight matrices which

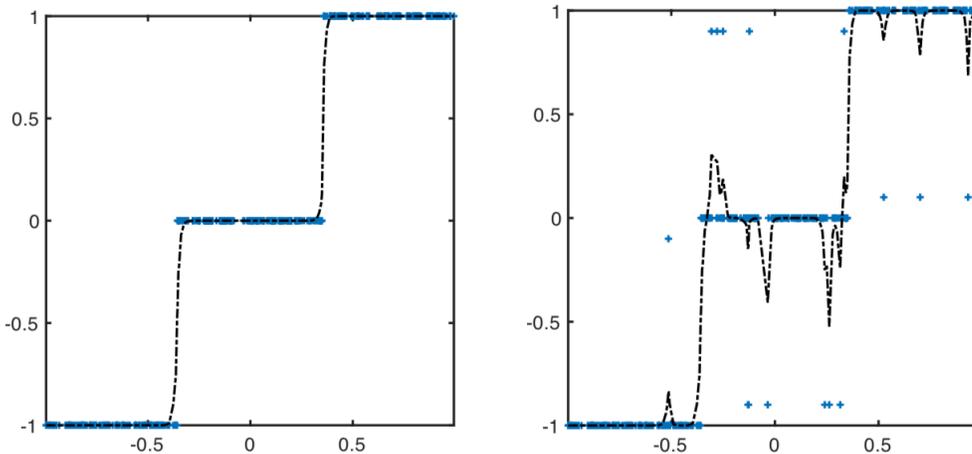


FIGURE 2. Synthetic labelled data and Tychonoff regularized solution  $u$ . Left: The solution value matches the labels perfectly on a large portion of the data set. Right: 10% of the data is corrupted by incorrect labels; the regularized solution is not as effective at correcting errors. The effect is more pronounced in higher dimensions.

are zero due to the activation functions, the gap in the inequality can be off by factors of many orders of magnitude for deep networks (Finlay & Oberman, 2018).

Implementing (2) can be accomplished using backpropagation in the  $x$  variable on each label, which can become costly for  $m$  large. Special architectures could also be used to implement Lipschitz regularization, for example, on a restricted architecture, Liao *et al.* (2018) renormalized the weight matrices of each layer to be norm 1. In practise, the computational cost can be reduced by regularizing the loss applied to the model  $\ell \circ u$ , instead

$$(4) \quad \max_{i \in I} \|\nabla_x \ell(u(x_i))\| \leq \text{Lip}(\ell \circ u; X)$$

which requires only one backpropagation. Lipschitz regularization may help with adversarial examples (Szegedy *et al.*, 2013) (Goodfellow *et al.*, 2014) which poses a problem for model reliability (Goodfellow *et al.*, 2018). Since the Lipschitz constant  $L_\ell$  of the loss,  $\ell$ , controls the norm of a perturbation

$$\|\ell(u(x_i + \epsilon v)) - \ell(u(x_i))\|_Y \leq \epsilon L_\ell \|v\|_X$$

maps with smaller Lipschitz constants may be more robust to adversarial examples. Finlay & Oberman (2018) implemented Lipschitz regularization of the loss, as in (4), and achieved better robustness against adversarial examples, compared to adversarial training (Goodfellow *et al.*, 2014) alone.

Lipschitz regularization may also improve stability of GANs. 1-Lipschitz networks with are also important for Wasserstein-GANs (Arjovsky *et al.*, 2017) (Arjovsky & Bottou, 2017). In (Wei *et al.*, 2018) the gradient penalty away from norm 1 is implemented, augmented by a penalty around perturbed points, with the

goal of improved stability. Spectral regularization for GANs was implemented in (Miyato *et al.*, 2018).

*Remark 1.4.* Consider the following problem inspired by (Zhang *et al.*, 2016). Given a labelled data set, which is Lipschitz continuous with constant  $L_0 = 1$ , say. Now consider making 100 copies of each data point, and adding a small  $\epsilon$  norm of noise to each image, but keeping the same labels. Call this data set  $\mathcal{D}_n$ . Now consider a second data set  $\tilde{\mathcal{D}}_n$  which is a copy of  $\mathcal{D}_n$ , but where with probability  $p$ , a given label is randomly changed. The Lipschitz constant of  $\tilde{\mathcal{D}}_n$  is  $O(1/\epsilon)$ , since two images a distance  $\epsilon$  apart will have different labels. When  $p$  is small, we expect that solving (1) with  $L_0 = 1$  will lead to a good approximation of the original map, for both data sets, but will result in an expected loss of approximately  $p$  for the second data set. See Figure 1.

## 2. LIPSCHITZ REGULARIZATION AND CONVERGENCE

**2.1. Limiting problem.** The variational problem (1) admits Lipschitz continuous minimizers, but in general the minimizers are not unique. When  $L_0 = \text{Lip}(u_0)$ , it is clear that  $u_0$  is a solution of (1): both the loss term and the regularization term are zero when applied to  $u_0$ . In addition, any  $L_0$ -Lipschitz extension of  $u_0|_{\mathcal{D}_n}$  is also a minimizer of (1), so solutions are not unique.

Let  $u_n$  be any solution of the Lipschitz regularized variational problem (1). We study the limit of  $u_n$  as  $n \rightarrow \infty$ . Since the empirical probability measures  $\rho_n$  converge to the data distribution  $\rho$ , we would expect the limit of  $u_n$  to be a solution of the continuum variational problem

$$(5) \quad \min_{u: X \rightarrow Y} J[u] \equiv L[u; \rho] + \lambda \max(\text{Lip}(u) - L_0, 0),$$

where in (5) we have introduced the following notation.

**Definition 2.1.** Given the loss function,  $\ell$ , a map  $u : X \rightarrow Y$ , and a probability measure,  $\mu$ , supported on  $X$ , define

$$L[u, \mu] = \mathbb{E}_{x \sim \mu}[\ell(u(x), u_0(x))] = \int_X \ell(u(x), u_0(x)) d\mu(x)$$

to be the expectation of the loss with respect to the measure. In particular, the *generalization loss* of the map  $u : X \rightarrow Y$  is given by  $L[u, \rho]$ . Write  $L[u, \mathcal{D}_n] := L[u, \rho_n]$  for the average loss on the data set  $\mathcal{D}_n$ , where  $\rho_n := \frac{1}{n} \sum \delta_{x_i}$  is the empirical measure corresponding to  $\mathcal{D}_n$ .

*Remark 2.2.* Generalization is defined in (Goodfellow *et al.*, 2016, Section 5.2) as the expected value of the loss function on a new input sampled from the data distribution. As defined, the full generalization error includes the training data, but it is of measure zero, so removing it does not change the value.

We would also expect the sequence of generalization losses  $L[u_n; \rho]$  to converge to zero in the case of perfect generalization.

We prove both of these results, including convergence rates, below. These results lead to an immediate proof that Lipschitz regularized DNN generalize.

**2.2. Loss function assumptions.** We introduce the following assumption on the loss function.

**Assumption 2.3** (Loss function). The function  $\ell : Y \times Y \rightarrow \mathbb{R}$  is a *loss* function if it satisfies (i)  $\ell \geq 0$ , (ii)  $\ell(y_1, y_2) = 0$  if and only if  $y_1 = y_2$ , and (iii)  $\ell$  is strictly convex in  $y_1$ .

*Example 2.4* ( $\mathbb{R}^m$  with  $L^2$  loss). Set  $Y = \mathbb{R}^m$ , and let each label be a basis vector. Set  $\ell(y_1, y_2) = \|y_1 - y_2\|_2^2$  to be the  $L^2$  loss.

*Example 2.5* (Classification). In classification problems, the output of the network is a probability vector on the labels. Thus  $Y = \Delta_p$ , the  $p$ -dimensional probability simplex, and each label is mapped to a basis vector. The cross-entropy loss is given by  $\ell^{KL}(y, z) = -\sum_{i=1}^p z_i \log(y_i/z_i)$ . For labels,  $\ell^{KL}(y, e_k) = -\log(y_k)$ .

*Example 2.6.* Define the *regularized cross entropy loss* with parameter  $\varepsilon > 0$  by

$$\ell_\varepsilon^{KL}(y, z) = -\sum_{i=1}^p z_i \log\left(\frac{y_i + \varepsilon}{z_i}\right).$$

For classification problems, where  $z = e_k$ , we have  $\ell_\varepsilon^{KL}(y, e_k) = -\log(y_k + \varepsilon)$ , which is Lipschitz and strongly convex for  $0 \leq y_i \leq 1$ .

In Theorems 2.11 and 2.14 which follow, the cross entropy loss  $\ell^{KL}$  does not satisfy the Lipschitz condition required. However they apply to  $\ell_\varepsilon^{KL}$  for any  $\varepsilon > 0$ .

**2.3. Convergence result for clean labels.** Here, we show that solutions of the random variational problem (1) converge to solutions of (5). We make the standard manifold assumption (Chapelle *et al.*, 2006), and assume the data distribution  $\rho$  is a probability density supported on a compact, smooth,  $m_0$ -dimensional manifold  $\mathcal{M}$  embedded in  $X = [0, 1]^d$ , where  $m_0 \ll d$ . We denote the probability density again by  $\rho : \mathcal{M} \rightarrow [0, \infty)$ . Hence, the data  $\mathcal{D}_n$  is a sequence  $x_1, \dots, x_n$  of *i.i.d.* random variables on  $\mathcal{M}$  with probability density  $\rho$ . Associated with the random sample we have the closet point projection map  $\sigma_n : X \rightarrow \{x_1, \dots, x_n\} \subset X$  that satisfies

$$\|x - \sigma_n(x)\|_X = \min_{1 \leq i \leq n} \{\|x - x_i\|_X\}$$

for all  $x \in X$ . We recall that  $W^{1,\infty}(X; Y)$  is the space of Lipschitz mappings from  $X$  to  $Y$ . Throughout this section,  $C, c > 0$  denote positive constants depending only on  $\mathcal{M}$ , and we assume  $C \geq 1$  and  $0 < c < 1$ .

We establish that that minimizers of (5) are unique on  $\mathcal{M}$ , which follows from the strict convexity of the loss restricted to the data manifold  $\mathcal{M}$ , in Theorem 3.1. See also Figure 3 which shows how the solutions need not be unique off the data manifold.

Our first convergence result is in the case where  $\text{Lip}[u_0] \leq L_0$ , and so the Lipschitz regularizer is not fully active. This corresponds to the case of clean labels.

**Theorem 2.7** (Convergence for clean labels). *Suppose that  $\text{Lip}[u_0] \leq L_0$  and  $\inf_{x \in \mathcal{M}} \rho(x) > 0$ . If  $u_n \in W^{1,\infty}(X; Y)$  is any sequence of minimizers of (1) then for any  $t > 0$*

$$\|u_0 - u_n\|_{L^\infty(\mathcal{M}; Y)} \leq CL_0 \left(\frac{t \log(n)}{n}\right)^{1/m}$$

*holds with probability at least  $1 - Ct^{-1}n^{-(ct-1)}$ .*

As an immediate corollary, we can prove that the generalization loss converges to zero, and so we obtain perfect generalization.

**Corollary 2.8.** *Assume that for some  $q \geq 1$  the loss  $\ell$  satisfies*

$$(6) \quad \ell(y, y_0) \leq C \|y - y_0\|_Y^q \text{ for all } y_0, y \in Y.$$

*Then under the assumptions of Theorem 2.7*

$$L[u_n, \rho] \leq CL_0^q \left( \frac{t \log(n)}{n} \right)^{q/m}$$

*holds with probability at least  $1 - Ct^{-1}n^{-(ct-1)}$ .*

*Proof.* By (6), we can bound the generalization loss as follows

$$L[u_n, \rho] = \int_{\mathcal{M}} \ell(u_n(x), u_0(x)) dVol(x) \leq C Vol(\mathcal{M}) \|u_n - u_0\|_{L^\infty(\mathcal{M}; Y)}^q.$$

The proof is completed by invoking Theorem 2.7.  $\square$

We now turn to the proof of Theorem 2.7, which requires a bound on the distance between the closest point projection  $\sigma_n$  and the identity. The result is standard in probability, and we include it for completeness in Lemma 2.9 proved in §3.1. We refer the interested reader to (Penrose *et al.*, 2003) for more details.

**Lemma 2.9.** *Suppose that  $\inf_{\mathcal{M}} \rho > 0$ . Then for any  $t > 0$*

$$\|Id - \sigma_n\|_{L^\infty(\mathcal{M}; X)} \leq C \left( \frac{t \log(n)}{n} \right)^{1/m}$$

*with probability at least  $1 - Ct^{-1}n^{-(ct-1)}$ .*

We now give the proof of Theorem 2.7.

*Proof of Theorem 2.7.* Since  $J_n[u_n] = J_n[u_0] = 0$ , we must have  $\text{Lip}[u_n] \leq L_0$  and  $u_0(x_i) = u_n(x_i)$  for all  $1 \leq i \leq n$ . Then for any  $x \in X$  we have

$$\begin{aligned} \|u_0(x) - u_n(x)\|_Y &= \|u_0(x) - u_0(\sigma_n(x)) + u_0(\sigma_n(x)) - u_n(\sigma_n(x)) + u_n(\sigma_n(x)) - u_n(x)\|_Y \\ &\leq \|u_0(x) - u_0(\sigma_n(x))\|_Y + \|u_n(\sigma_n(x)) - u_n(x)\|_Y \\ &\leq 2L_0 \|x - \sigma_n(x)\|_X. \end{aligned}$$

Therefore, we deduce

$$\|u_0 - u_n\|_{L^\infty(\mathcal{M}; Y)} \leq 2L_0 \|Id - \sigma_n\|_{L^\infty(\mathcal{M}; X)}.$$

The proof is completed by invoking Lemma 2.9.  $\square$

**2.4. Convergence for noisy labels.** Our second result is in the setting where  $\text{Lip}(u_0) > L_0$ , so the regularizer is active and we do not expect  $u_n \rightarrow u_0$  as  $n \rightarrow \infty$ . This setting models the case where some labels have errors, and so the Lipschitz constant for the data is a large over-estimate of the Lipschitz constant of the clean labeling function. Our main result shows that minimizers of  $J^n$  converge to minimizers of  $J$ .

*Remark 2.10.* In the theorem which follows, the sequence  $u_n$  does not, in general, converge on the whole domain  $X$ . The important point is that the sequence converges on the data manifold  $\mathcal{M}$ , and solves the variational problem (5) off of the manifold, which ensures that the output of the DNN is stable with respect to the input. See Figure 3.

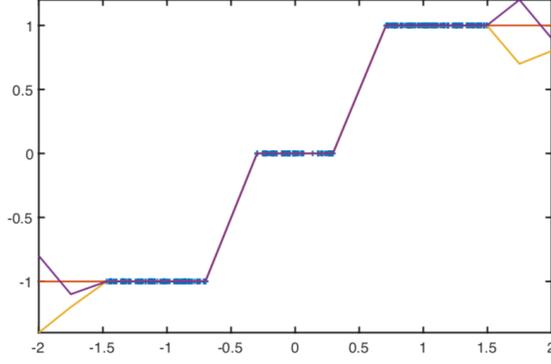


FIGURE 3. Non-uniqueness of the solution off the data manifold: in the middle areas off the data manifold where the Lipschitz constant is attained, the solution is unique. In the outer area off the data manifold, the solution is not unique.

**Theorem 2.11.** *Suppose that  $\inf_{\mathcal{M}} \rho > 0$ ,  $\ell : Y \times Y \rightarrow \mathbb{R}$  is Lipschitz, and let  $u^* \in W^{1,\infty}(X; Y)$  be any minimizer of (5). Then with probability one*

$$(7) \quad u_n \rightarrow u^* \text{ uniformly on } \mathcal{M} \text{ as } n \rightarrow \infty,$$

where  $u_n$  is any sequence of minimizers of (1). Furthermore, every uniformly convergent subsequence of  $u_n$  converges on  $X$  to a minimizer of (5).

The proof of Theorem 2.11 requires a preliminary Lemma. Let  $H_L(X; Y)$  denote the collection of  $L$ -Lipschitz functions  $w : X \rightarrow Y$ .

**Lemma 2.12.** *Suppose that  $\inf_{\mathcal{M}} \rho > 0$ , and  $\dim(\mathcal{M}) = m_0$ . Then for any  $t > 0$*

$$(8) \quad \sup_{w \in H_L(X; Y)} \left| \frac{1}{n} \sum_{i=1}^n w(x_i) - \int_{\mathcal{M}} w \rho dVol(x) \right| \leq CL \left( \frac{t \log(n)}{n} \right)^{\frac{1}{m_0+2}}$$

holds with probability at least  $1 - 2t^{-\frac{m_0}{m_0+2}} n^{-(ct-1)}$ .

The estimate (8) is called a discrepancy result (Talagrand, 2006; Györfi *et al.*, 2006), and is a uniform version of concentration inequalities. We include a simple proof in Section 3.2.

*Proof of Theorem 2.11.* By Lemma 2.12 the event that

$$(9) \quad \lim_{n \rightarrow \infty} \sup_{w \in H_L(X; Y)} |L[w, \rho_n] - L[w, \rho]| = 0$$

for all Lipschitz constants  $L > 0$  has probability one. For the rest of the proof we restrict ourselves to this event.

Let  $u_n \in W^{1,\infty}(X; Y)$  be a sequence of minimizers of (1), and let  $u^* \in W^{1,\infty}(X; Y)$  be any minimizer of (5). Then since

$$\lambda(\text{Lip}(u_n) - L_0) \leq J^n[u_n] \leq J^n[u_0] = \lambda(\text{Lip}(u_0) - L_0)$$

we have  $\text{Lip}(u_n) \leq \text{Lip}(u_0) =: L$  for all  $n$ . By the Arzelà-Ascoli Theorem (Rudin, 1976) there exists a subsequence  $u_{n_j}$  and a function  $u \in W^{1,\infty}(X; Y)$  such that

$u_{n_j} \rightarrow u$  uniformly as  $n_j \rightarrow \infty$ . Note we also have  $\text{Lip}(u) \leq \liminf_{j \rightarrow \infty} \text{Lip}(u_{n_j})$ . Since

$$\begin{aligned} |L[u_n, \rho_n] - L[u, \rho]| &\leq |L[u_n, \rho_n] - L[u, \rho_n]| + |L[u, \rho_n] - L[u, \rho]| \\ &\leq C \|u_n - u\|_{L^\infty(\mathcal{M}; Y)} + \sup_{w \in H_L(X; Y)} |L[w, \rho_n] - L[w, \rho]| \end{aligned}$$

it follows from (9) that  $L[u_{n_j}, \rho_{n_j}] \rightarrow L[u, \rho]$  as  $j \rightarrow \infty$ . It also follows from (9) that  $J^n[u^*] \rightarrow J[u^*]$  as  $n \rightarrow \infty$ . Therefore

$$\begin{aligned} J[u^*] &= \lim_{n \rightarrow \infty} J^n[u^*] \\ &\geq \liminf_{n \rightarrow \infty} J^n[u_n] \\ &= \liminf_{n \rightarrow \infty} L[u_n, \rho_n] + \lambda \max(\text{Lip}(u_n) - L_0, 0) \\ &= \lim_{n \rightarrow \infty} L[u_n, \rho_n] + \liminf_{n \rightarrow \infty} \lambda \max(\text{Lip}(u_n) - L_0, 0) \\ &\geq L[u, \rho] + \lambda \max(\text{Lip}(u) - L_0, 0) = J[u]. \end{aligned}$$

Therefore,  $u$  is a minimizer of  $J$ . By Theorem 3.1,  $u = u^*$  on  $\mathcal{M}$ , and so  $u_{n_j} \rightarrow u^*$  uniformly on  $\mathcal{M}$  as  $j \rightarrow \infty$ .

Now, suppose that (7) does not hold. Then there exists a subsequence  $u_{n_j}$  and  $\delta > 0$  such that

$$\max_{x \in \mathcal{M}} |u_{n_j}(x) - u^*(x)| > \delta$$

for all  $j \geq 1$ . However, we can apply the argument above to extract a further subsequence of  $u_{n_j}$  that converges uniformly on  $\mathcal{M}$  to  $u^*$ , which is a contradiction. This completes the proof.  $\square$

Finally, we prove a rate in the case where the loss  $\ell$  is strongly convex in the first variable.

**Definition 2.13.** We say that  $\ell$  is strongly convex with parameter  $\theta > 0$  if

$$(10) \quad \ell(ty_1 + (1-t)y_2, y_0) + \frac{\theta}{2}t(1-t)\|y_1 - y_2\|_Y^2 \leq t\ell(y_1, y_0) + (1-t)\ell(y_2, y_0)$$

for all  $y_0, y_1, y_2 \in Y$  and  $0 \leq t \leq 1$ .

We note that when  $\ell$  is twice differentiable, this notion of strong convexity is equivalent to assuming  $\nabla_{y_1}^2 \ell \geq \theta I$ . The definition in equation (10) is useful for non-smooth functions, such as the Lipschitz semi-norm present in  $J[u]$ .

Our final result is the following  $L^2$  convergence rate in the strongly convex setting.

**Theorem 2.14.** *Suppose that  $\ell : Y \times Y \rightarrow \mathbb{R}$  is Lipschitz and strongly convex and let  $L = \text{Lip}(u_0)$ . Then for any  $t > 0$ , with probability at least  $1 - 2t^{-\frac{m}{m+2}}n^{-(ct-1)}$  all minimizing sequences  $u_n$  of (1) and all minimizers  $u^*$  of (5) satisfy*

$$\frac{\theta}{2} \int_{\mathcal{M}} \|u_n - u^*\|_Y^2 \rho d\text{Vol}(x) \leq CL \left( \frac{t \log(n)}{n} \right)^{\frac{1}{m+2}}.$$

Before proving Theorem 2.14, we require a preliminary lemma.

**Lemma 2.15.** *If  $u^* \in W^{1,\infty}(X; Y)$  is a minimizer of (5) and  $u \in W^{1,\infty}(X; Y)$  then*

$$\frac{\theta}{2} \int_{\mathcal{M}} \|u - u^*\|_Y^2 \rho d\text{Vol}(x) \leq J[u] - J[u^*].$$

*Proof.* We use Proposition 3.4 with  $u_1 = u^*$  and  $u_2 = u$  to obtain

$$J[tu^* + (1-t)u] + \frac{\theta}{2}t(1-t) \int_{\mathcal{M}} \|u^* - u\|_Y^2 \rho dVol(x) \leq tJ[u^*] + (1-t)J[u].$$

Since  $J[tu^* + (1-t)u] \geq J[u^*]$

$$J[u^*] + \frac{\theta}{2}t(1-t) \int_{\mathcal{M}} \|u^* - u\|_Y^2 \rho dVol(x) \leq tJ[u^*] + (1-t)J[u],$$

and so

$$\frac{\theta}{2}t \int_{\mathcal{M}} \|u^* - u\|_Y^2 \rho dVol(x) \leq J[u] - J[u^*].$$

Setting  $t = 1$  completes the proof.  $\square$

*Proof of Theorem 2.14.* Let  $L = \text{Lip}(u_0)$ . By Lemma 2.12

$$(11) \quad \sup_{w \in H_L(X;Y)} |L[w, \rho_n] - L[w, \rho]| \leq CL \left( \frac{t \log(n)}{n} \right)^{\frac{1}{m+2}}$$

holds with probability at least  $1 - 2t^{-\frac{m}{m+2}}n^{-(ct-1)}$  for any  $t > 0$ . Let us assume for the rest of the proof that (11) holds.

As in the proof of Theorem 2.11, we have  $\text{Lip}(u_n) \leq L$  and  $\text{Lip}(u^*) \leq L$ , and so

$$|J^n[u^*] - J[u^*]|, |J^n[u_n] - J[u_n]| \leq CL \left( \frac{t \log(n)}{n} \right)^{\frac{1}{m+2}}.$$

Therefore

$$J[u_n] - J[u^*] = J^n[u_n] - J[u^*] + J[u_n] - J^n[u_n] \leq CL \left( \frac{t \log(n)}{n} \right)^{\frac{1}{m+2}}.$$

By Lemma 2.15 we deduce

$$\frac{\theta}{2} \int_{\mathcal{M}} \|u_n - u^*\|_Y^2 \rho dVol(x) \leq CL \left( \frac{t \log(n)}{n} \right)^{\frac{1}{m+2}},$$

which completes the proof.  $\square$

### 3. PROOFS

**3.1. Proofs for clean labels.** In this section we provide the proof of results stated in §2.3.

**Theorem 3.1.** *Suppose the loss function satisfies Assumption 2.3. If  $u, v \in W^{1,\infty}(X;Y)$  are two minimizers of (5) and  $\inf_{\mathcal{M}} \rho > 0$  then  $u = v$  on  $\mathcal{M}$ .*

*Proof.* Let  $w = (u + v)/2$ . Then

$$\begin{aligned} J[w] &= \int_{\mathcal{M}} \ell \left( \frac{1}{2}u + \frac{1}{2}v, u_0 \right) \rho dVol(x) + \lambda \max \left( \text{Lip} \left( \frac{1}{2}u + \frac{1}{2}v \right), 0 \right) \\ &\leq \int_{\mathcal{M}} \left[ \frac{1}{2}\ell(u, u_0) + \frac{1}{2}\ell(v, u_0) \right] \rho dVol(x) + \lambda \max \left( \frac{1}{2} \text{Lip}(u) + \frac{1}{2} \text{Lip}(v), 0 \right) \\ &\leq \int_{\mathcal{M}} \left[ \frac{1}{2}\ell(u, u_0) + \frac{1}{2}\ell(v, u_0) \right] \rho dVol(x) + \lambda \left[ \frac{1}{2} \max(\text{Lip}(u), 0) + \frac{1}{2} \max(\text{Lip}(v), 0) \right] \\ &= \frac{1}{2}J[u] + \frac{1}{2}J[v] = \min_u J[u]. \end{aligned}$$

Therefore,  $w$  is a minimizer of  $J$  and so we have equality above, which yields

$$\int_{\mathcal{M}} \left[ \frac{1}{2} \ell(u, u_0) + \frac{1}{2} \ell(v, u_0) \right] \rho dVol(x) = \int_{\mathcal{M}} \ell\left(\frac{1}{2}u + \frac{1}{2}v, u_0\right) \rho dVol(x).$$

Since  $\ell$  is strictly convex in its first argument, it follows that  $u = v$  on  $\mathcal{M}$ .  $\square$

*Proof of Lemma 2.9 of §2.3.* There exists  $\varepsilon_{\mathcal{M}}$  such that for any  $0 < \varepsilon \leq \varepsilon_{\mathcal{M}}$ , we can cover  $\mathcal{M}$  with  $N$  geodesic balls  $B_1, B_2, \dots, B_N$  of radius  $\varepsilon$ , where  $N \leq C\varepsilon^{-m}$  and  $C$  depends only on  $\mathcal{M}$  (Györfi *et al.*, 2006). Let  $Z_i$  denote the number of random variables  $x_1, \dots, x_n$  falling in  $B_i$ . Then  $Z_i \sim B(n, p_i)$ , where  $p_i = \int_{B_i} \rho(x) dVol(x)$ . Since  $\rho \geq \theta > 0$  and  $Vol(B_i) \geq c\varepsilon^m$  we have  $p_i \geq c\varepsilon^m$ . Let  $A_n$  denote the event that at least one  $B_i$  is empty (i.e.,  $Z_i = 0$  for some  $i$ ). Then by the union bound we deduce

$$\begin{aligned} \mathbb{P}(A_n) &\leq \sum_{i=1}^N \mathbb{P}(Z_i = 0) \\ &\leq C\varepsilon^{-d}(1 - c\varepsilon^m)^n \\ &= C \exp(n \log(1 - c\varepsilon^m) - \log(\varepsilon^m)) \\ &\leq C \exp(-cn\varepsilon^m - \log(\varepsilon^m)). \end{aligned}$$

Choose  $0 < \varepsilon \leq \varepsilon_{\mathcal{M}}$  in the form  $n\varepsilon^m = t \log(n)$  with  $t \leq n\varepsilon_{\mathcal{M}}^m / \log(n)$ . Then

$$\mathbb{P}(A_n) \leq Ct^{-1} \exp(-(ct - 1) \log(n)).$$

In the event that  $A_n$  does not occur, then each  $B_i$  has at least one point, and so  $|x - \sigma_n(x)| \leq C\varepsilon$  for all  $x \in \mathcal{M}$ . Therefore

$$\|\text{Id} - \sigma_n\|_{L^\infty(\mathcal{M}; X)} \leq C\varepsilon = C \left( \frac{t \log(n)}{n} \right)^{1/m}$$

with probability at least  $1 - Ct^{-1} \exp(-(ct - 1) \log(n))$ . Since  $\|\text{Id} - \sigma_n\|_{L^\infty(\mathcal{M}; X)} \leq C\sqrt{d}$ , the result holds for  $t \geq n\varepsilon_{\mathcal{M}}^m / \log(n)$ , albeit with a larger constant  $C$ .  $\square$

**3.2. Proofs for noisy labels.** Here, we give the proofs of lemmas and propositions required in Section 2.4. We first give the proof of Lemma 2.12. A key tool in the proof is Bernstein's inequality (Boucheron *et al.*, 2013), which we recall now for the reader's convenience. For  $X_1, \dots, X_n$  *i.i.d.* with variance  $\sigma^2 = \mathbb{E}[(X_i - \mathbb{E}[X_i])^2]$ , if  $|X_i| \leq M$  almost surely for all  $i$  then Bernstein's inequality states that for any  $\varepsilon > 0$

$$\mathbb{P} \left( \left| \frac{1}{n} \sum_{i=1}^n X_i - \mathbb{E}[X_i] \right| > \varepsilon \right) \leq 2 \exp \left( - \frac{n\varepsilon^2}{2\sigma^2 + 4M\varepsilon/3} \right).$$

*Proof of Lemma 2.12.* We note that it is sufficient to prove the result for  $w \in H_L(X; Y)$  with  $\int_{\mathcal{M}} w \rho dVol(x) = 0$ . In this case, we have  $w(x) = 0$  for some  $x \in \mathcal{M}$ , and so  $\|w\|_{L^\infty(X; Y)} \leq CL$ . We also write  $m$  in place of  $m_0$  in the proof for simplicity.

We first give the proof for  $\mathcal{M} = X = [0, 1]^m$ . We partition  $X$  into hypercubes  $B_1, \dots, B_N$  of side length  $h > 0$ , where  $N = h^{-m}$ . Let  $Z_j$  denote the number of

$x_1, \dots, x_n$  falling in  $B_j$ . Then  $Z_j$  is a Binomial random variable with parameters  $n$  and  $p_j = \int_{B_j} \rho dx \geq ch^m$ . By the Bernstein inequality we have for each  $j$  that

$$(12) \quad \mathbb{P} \left( \left| \frac{1}{n} Z_j - \int_{B_j} \rho dx \right| > \varepsilon \right) \leq 2 \exp(-cnh^{-m}\varepsilon^2)$$

provided  $0 < \varepsilon \leq h^m$ . Therefore, we deduce

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n w(x_i) &\leq \frac{1}{n} \sum_{j=1}^N Z_j \max_{B_j} w \\ &\stackrel{(12)}{\leq} \sum_{j=1}^N \left( \int_{B_j} \rho dx + \varepsilon \right) \max_{B_j} w \\ &\leq \sum_{j=1}^N \max_{B_j} w \int_{B_j} \rho dx + CLh^{-m}\varepsilon \\ &\leq \sum_{j=1}^N (\min_{B_j} w + CLh) \int_{B_j} \rho dx + CLh^{-m}\varepsilon \\ &\leq \sum_{j=1}^N \int_{B_j} w \rho dx + CLh^{-m}(h^{m+1} + \varepsilon) \\ &= \int_X w \rho dx + CL(h + h^{-m}\varepsilon) \end{aligned}$$

holds with probability at least  $1 - 2h^{-m} \exp(-cnh^{-m}\varepsilon^2)$  for any  $0 < \varepsilon \leq h^m$ . Choosing  $\varepsilon = h^{m+1}$  we have that

$$\left| \frac{1}{n} \sum_{i=1}^n w(x_i) - \int_X w \rho dx \right| \leq CLh$$

holds for all  $u \in H_L(X; Y)$  with probability at least  $1 - 2h^{-m} \exp(-cnh^{m+2})$ , provided  $h \leq 1$ . By selecting  $nh^{m+2} = t \log(n)$

$$\sup_{w \in H_L(X; Y)} \left| \frac{1}{n} \sum_{i=1}^n w(x_i) - \int_{\mathcal{M}} w \rho dVol(x) \right| \leq CL \left( \frac{t \log(n)}{n} \right)^{\frac{1}{m+2}}$$

holds with probability at least  $1 - 2t^{-\frac{m}{m+2}} n^{-(ct-1)}$  for  $t \leq n/\log(n)$ . Since we have  $\|w\|_{L^\infty(X; Y)} \leq CL$ , the estimate

$$\sup_{w \in H_L(X; Y)} \left| \frac{1}{n} \sum_{i=1}^n w(x_i) - \int_{\mathcal{M}} w \rho dVol(x) \right| \leq CL,$$

trivially holds, and hence we can allow  $t > n/\log(n)$  as well.

We sketch here how to prove the result on the manifold  $\mathcal{M}$ . We cover  $\mathcal{M}$  with  $k$  geodesic balls of radius  $\varepsilon > 0$ , denoted  $B_{\mathcal{M}}(x_1, \varepsilon), \dots, B_{\mathcal{M}}(x_k, \varepsilon)$ , and let  $\varphi_1, \dots, \varphi_k$  be a partition of unity subordinate to this open covering of  $\mathcal{M}$ . For  $\varepsilon > 0$  sufficiently small, the Riemannian exponential map  $\exp_x : B(0, \varepsilon) \subset T_x \mathcal{M} \rightarrow \mathcal{M}$  is a diffeomorphism between the ball  $B(0, r) \subset T_x \mathcal{M}$  and the geodesic ball  $B_{\mathcal{M}}(x, \varepsilon) \subset \mathcal{M}$ ,

where  $T_x\mathcal{M} \cong \mathbb{R}^m$ . Furthermore, the Jacobian of  $\exp_x$  at  $v \in B(0, r) \subset T_x\mathcal{M}$ , denoted by  $J_x(v)$ , satisfies (by the Rauch Comparison Theorem)

$$(1 + C|v|^2)^{-1} \leq J_x(v) \leq 1 + C|v|^2.$$

Therefore, we can run the argument above on the ball  $B(0, r) \subset \mathbb{R}^m$  in the tangent space, lift the result to the geodesic ball  $B_{\mathcal{M}}(x_i, \varepsilon)$  via the Riemannian exponential map  $\exp_x$ , and apply the bound

$$\left| \frac{1}{n} \sum_{i=1}^n w(x_i) - \int_{\mathcal{M}} w \rho dVol(x) \right| \leq \sum_{j=1}^k \left| \frac{1}{n} \sum_{i=1}^n \varphi_j(x_i) w(x_i) - \int_{\mathcal{M}} \varphi_j w \rho dVol(x) \right|$$

to complete the proof.  $\square$

*Remark 3.2.* The exponent  $1/(m_0 + 2)$  is not optimal, but affords a very simple proof. It is possible to prove a similar result with the optimal exponent  $1/m_0$  in dimension  $m_0 \geq 3$ , but the proof is significantly more involved. We refer the reader to (Talagrand, 2006) for details.

*Remark 3.3.* The proof of Theorem 2.11 shows that (1)  $\Gamma$ -converges to (5) almost surely as  $n \rightarrow \infty$  in the  $L^\infty(X; Y)$  topology.  $\Gamma$ -convergence is a notion of convergence for functionals that ensures minimizers along a sequence of functionals converge to a minimizer of the  $\Gamma$ -limit. While we do not use the language of  $\Gamma$ -convergence here, the ideas are present in the proof of Theorem 2.11. We refer to (Braides, 2002) for details on  $\Gamma$ -convergence.

We finally give a proposition useful in the proof of Lemma 2.15.

**Proposition 3.4.** *If  $\ell$  is strongly convex with parameter  $\theta > 0$  then*

$$J[tu_1 + (1-t)u_2] + \frac{\theta}{2}t(1-t) \int_{\mathcal{M}} \|u_1 - u_2\|_Y^2 \rho dVol(x) \leq tJ[u_1] + (1-t)J[u_2]$$

for all  $u_1, u_2 \in W^{1,\infty}(X; Y)$  and  $0 \leq t \leq 1$ .

*Proof.* We compute

$$\begin{aligned} & J[tu_1 + (1-t)u_2] \\ &= \int_{\mathcal{M}} \ell(tu_1 + (1-t)u_2, u_0) \rho dVol(x) + \lambda \max(\text{Lip}(tu_1 + (1-t)u_2), 0) \\ &\leq tJ[u_1] + (1-t)J[u_2] - \frac{\theta}{2}t(1-t) \int_{\mathcal{M}} \|u_1 - u_2\|_Y^2 \rho dVol(x), \end{aligned}$$

which completes the proof.  $\square$

#### 4. INDUCED MATRIX NORMS

In some cases, we can take advantage of explicit formulas for matrix norms, which makes the estimates in (3) an explicit function of the weights. Define the induced matrix norm by

$$\|M\|_{p,q} = \sup_x \frac{\|Mx\|_q}{\|x\|_p}$$

Then the following matrix norms formulas hold (see (Horn *et al.*, 1990, Chapter 5.6.4))

$$\begin{aligned} \|M\|_{\infty,\infty} &= \max_i \sum_j |m_{ij}|, & \|M\|_{1,1} &= \max_j \sum_i |m_{ij}| \\ \|M\|_{1,\infty} &= \max_{i,j} |m_{ij}|, & \|M\|_{2,\infty} &= \max_i \sqrt{\sum_j m_{ij}^2} \end{aligned}$$

## 5. VARIATIONAL PROBLEMS IN IMAGE PROCESSING AND LIPSCHITZ EXTENSIONS

The variational problem (1) can be interpreted as a relaxation of the Lipschitz Extension problem.

$$(LE) \quad \begin{cases} \min_{u: X \rightarrow Y} \text{Lip}[u] \\ \text{subject to } u(x) = u_0(x) \text{ for } x \in \mathcal{D} \end{cases}$$

for  $\mathcal{D} \subset X$ . The problem (LE) has more than one solution. Two classical results giving explicit solutions in one dimension go back to Kirzbaum and to McShane (McShane, 1934). However solving (LE) is not practical for large scale problems. There has been extensive work on the Lipschitz Extension problem, see, (Johnson & Lindenstrauss, 1984), for example. More recently, optimal Lipschitz extensions have been studied, with connections to Partial Differential Equations, see (Aronsson *et al.*, 2004). We can interpret (1) as a relaxed version of (LE), where  $\lambda^{-1}$  is a parameter which replaces the unknown Lagrange multiplier for the constraint.

Variational problems are fundamental tools in mathematical approaches to image processing (Aubert & Kornprobst, 2006) and inverse problems more generally. Without regularization inverse problems can be ill-posed. The general form of the problem is

$$(13) \quad J[u] = L[u; u_0] + \lambda R[\nabla u]$$

which combines a loss or *fidelity* functional,  $L[u, u_0]$ , which depends on the values of  $u$  and the reference image  $u_0$ , and a *regularization* functional,  $R[\nabla u]$ , which depends on the gradient,  $\nabla u$ . The parameter  $\lambda$  determines the relative strength of the two terms which emphasize fidelity versus regularization.

*Example 5.1.* For example, a typical fidelity term is the standard least-squares  $L[u, u_0] = \|u - u_0\|_{L^2(D)}^2$ . The regularization  $\|\nabla u(x)\|_{L^2(D)}^2$  corresponds to the classical Tychonov regularization (Tikhonov & Arsenin, 1977),  $R[\nabla u] = \|\nabla u(x)\|_{L^1(D)}$  is the Total Variation regularization model of Rudin, Osher and Fatemi (Rudin *et al.*, 1992).

Lipschitz regularization is not nearly as common. It appears in image processing in (Pock *et al.*, 2010, §4.4) (Elion & Vese, 2007) and (Guillot & Le Guyader, 2009). Variational problems of the form (13) can be studied by the direct method in the calculus of variations (Dacorogna, 2007). The problem (13) can be discretized to obtain a finite dimensional convex optimization problem. The variational problem can also be studied by finding the first variation, which is a Partial Differential Equation (Evans, 1998), which can then be solved numerically. Both approaches are discussed in (Aubert & Kornprobst, 2006).

In Figure 4 we compare different regularization terms, in one dimension. The difference between the regularizers is more extreme in higher dimensions.

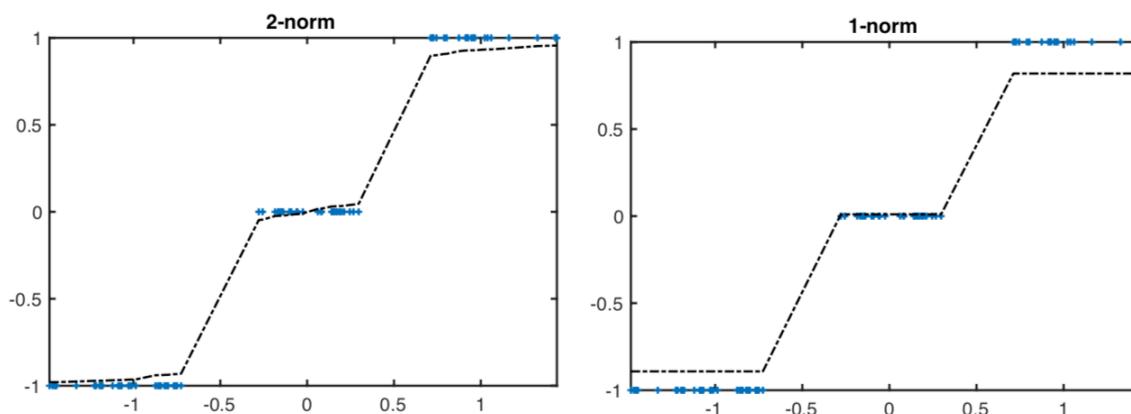


FIGURE 4. Comparison of different regularization methods. Lipschitz regularization preserves most of the labels (Figure 1). Tychonoff regularization smooths the solution (left). Total Variation regularization shifts the label values towards the mean (right).

#### REFERENCES

- Arjovsky, Martin, & Bottou, Léon. 2017. Towards principled methods for training generative adversarial networks. *arXiv preprint arXiv:1701.04862*.
- Arjovsky, Martin, Chintala, Soumith, & Bottou, Léon. 2017. Wasserstein gan. *arXiv preprint arXiv:1701.07875*.
- Aronsson, Gunnar, Crandall, Michael, & Juutinen, Petri. 2004. A tour of the theory of absolutely minimizing functions. *Bulletin of the American mathematical society*, **41**(4), 439–505.
- Aubert, Gilles, & Kornprobst, Pierre. 2006. *Mathematical problems in image processing: partial differential equations and the calculus of variations*. Vol. 147. Springer Science & Business Media.
- Bartlett, Peter L. 1997. For valid generalization the size of the weights is more important than the size of the network. *Pages 134–140 of: Advances in neural information processing systems*.
- Bartlett, Peter L, Foster, Dylan J, & Telgarsky, Matus J. 2017. Spectrally-normalized margin bounds for neural networks. *Pages 6240–6249 of: Advances in Neural Information Processing Systems*.
- Boucheron, Stéphane, Lugosi, Gábor, & Massart, Pascal. 2013. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press.
- Bousquet, Olivier, & Elisseeff, André. 2002. Stability and generalization. *Journal of machine learning research*, **2**(Mar), 499–526.
- Braides, Andrea. 2002. *Gamma-convergence for Beginners*. Vol. 22. Clarendon Press.
- Calder, Jeff. 2017. Consistency of Lipschitz learning with infinite unlabeled data and finite labeled data. *arXiv preprint arXiv:1710.10364*.
- Chapelle, Olivier, Scholkopf, Bernhard, & Zien, Alexander. 2006. *Semi-supervised learning*. MIT.

- Dacorogna, Bernard. 2007. *Direct methods in the calculus of variations*. Vol. 78. Springer Science & Business Media.
- Drucker, Harris, & Le Cun, Yann. 1992. Improving generalization performance using double backpropagation. *IEEE Transactions on Neural Networks*, **3**(6), 991–997.
- El Alaoui, Ahmed, Cheng, Xiang, Ramdas, Aaditya, Wainwright, Martin J, & Jordan, Michael I. 2016. Asymptotic behavior of  $\ell_p$ -based laplacian regularization in semi-supervised learning. *Pages 879–906 of: Conference on Learning Theory*.
- Elion, Christopher, & Vese, Luminita A. 2007. An image decomposition model using the total variation and the infinity Laplacian. *Page 64980W of: Computational Imaging V*, vol. 6498. International Society for Optics and Photonics.
- Evans, Lawrence C. 1998. *Partial differential equations*. Graduate Studies in Mathematics, vol. 19. American Mathematical Society.
- Finlay, Chris, & Oberman, Adam M. 2018. Improved robustness to adversarial examples using Lipschitz regularization of the loss. *arXiv preprint arXiv:1810.00953*.
- Goodfellow, Ian, Bengio, Yoshua, & Courville, Aaron. 2016. *Deep Learning*. MIT Press. <http://www.deeplearningbook.org>.
- Goodfellow, Ian, McDaniel, Patrick, & Papernot, Nicolas. 2018. Making machine learning robust against adversarial inputs. *Communications of the ACM*, **61**(7), 56–66.
- Goodfellow, Ian J, Shlens, Jonathon, & Szegedy, Christian. 2014. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*.
- Gouk, Henry, Frank, Eibe, Pfahringer, Bernhard, & Cree, Michael. 2018. Regularisation of Neural Networks by Enforcing Lipschitz Continuity. *arXiv preprint arXiv:1804.04368*.
- Guillot, Laurence, & Le Guyader, Carole. 2009. Extrapolation of Vector Fields Using the Infinity Laplacian and with Applications to Image Segmentation. *Pages 87–99 of: Tai, Xue-Cheng, Mørken, Knut, Lysaker, Marius, & Lie, Knut-Andreas (eds), Scale Space and Variational Methods in Computer Vision*. Berlin, Heidelberg: Springer Berlin Heidelberg.
- Györfi, László, Kohler, Michael, Krzyzak, Adam, & Walk, Harro. 2006. *A distribution-free theory of nonparametric regression*. Springer Science & Business Media.
- Hardt, Moritz, Recht, Benjamin, & Singer, Yoram. 2015. Train faster, generalize better: Stability of stochastic gradient descent. *arXiv preprint arXiv:1509.01240*.
- Horn, Roger A, Horn, Roger A, & Johnson, Charles R. 1990. *Matrix analysis*. Cambridge university press.
- Johnson, William B, & Lindenstrauss, Joram. 1984. Extensions of Lipschitz mappings into a Hilbert space. *Contemporary mathematics*, **26**(189-206), 1.
- Kyng, Rasmus, Rao, Anup, Sachdeva, Sushant, & Spielman, Daniel A. 2015. Algorithms for Lipschitz learning on graphs. *Pages 1190–1223 of: Conference on Learning Theory*.
- LeCun, Yann, Bengio, Yoshua, & Hinton, Geoffrey. 2015. Deep learning. *nature*, **521**(7553), 436.
- Liao, Qianli, Miranda, Brando, Banburski, Andrzej, Hidary, Jack, & Poggio, Tomaso. 2018. A Surprising Linear Relationship Predicts Test Performance in Deep Networks. *arXiv preprint arXiv:1807.09659*.

- McShane, Edward James. 1934. Extension of range of functions. *Bulletin of the American Mathematical Society*, **40**(12), 837–842.
- Miyato, Takeru, Kataoka, Toshiki, Koyama, Masanori, & Yoshida, Yuichi. 2018. Spectral normalization for generative adversarial networks. *arXiv preprint arXiv:1802.05957*.
- Penrose, Mathew, *et al.* . 2003. *Random geometric graphs*. Oxford university press.
- Pock, Thomas, Cremers, Daniel, Bischof, Horst, & Chambolle, Antonin. 2010. Global solutions of variational models with convex regularization. *SIAM Journal on Imaging Sciences*, **3**(4), 1122–1145.
- Rudin, Leonid I, Osher, Stanley, & Fatemi, Emad. 1992. Nonlinear total variation based noise removal algorithms. *Physica D: nonlinear phenomena*, **60**(1-4), 259–268.
- Rudin, Walter. 1976. *Principles of mathematical analysis*. McGraw-hill New York.
- Shalev-Shwartz, Shai, & Ben-David, Shai. 2014. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press.
- Slepcev, Dejan, & Thorpe, Matthew. 2017. Analysis of p-Laplacian Regularization in Semi-Supervised Learning. *arXiv preprint arXiv:1707.06213*.
- Szegedy, Christian, Zaremba, Wojciech, Sutskever, Ilya, Bruna, Joan, Erhan, Dumitru, Goodfellow, Ian, & Fergus, Rob. 2013. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*.
- Talagrand, Michel. 2006. *The generic chaining: upper and lower bounds of stochastic processes*. Springer Science & Business Media.
- Tikhonov, AN, & Arsenin, V Ya. 1977. *Solutions of Ill-Posed Problems*. Winston and Sons, New York.
- Wei, Xiang, Gong, Boqing, Liu, Zixia, Lu, Wei, & Wang, Liqiang. 2018. Improving the Improved Training of Wasserstein GANs: A Consistency Term and Its Dual Effect. *arXiv preprint arXiv:1803.01541*.
- Yoshida, Yuichi, & Miyato, Takeru. 2017. Spectral Norm Regularization for Improving the Generalizability of Deep Learning. *arXiv preprint arXiv:1705.10941*.
- Zhang, Chiyuan, Bengio, Samy, Hardt, Moritz, Recht, Benjamin, & Vinyals, Oriol. 2016. Understanding deep learning requires rethinking generalization. *arXiv:1611.03530*.

DEPARTMENT OF MATHEMATICS AND STATISTICS, MCGILL UNIVERSITY  
*E-mail address:* `adam.oberman@mcgill.ca`

SCHOOL OF MATHEMATICS, UNIVERSITY OF MINNESOTA  
*E-mail address:* `jcalder@umn.edu`