

# Implicit Regularization of AdaDelta

Anonymous authors

Paper under double-blind review

## Abstract

We consider the AdaDelta adaptive optimization algorithm on locally Lipschitz, positively homogeneous, and  $\alpha$ -minimally definable neural networks, with either the exponential or the logistic loss. We prove that, after achieving perfect training accuracy, the resulting adaptive gradient flows converge in direction to a Karush-Kuhn-Tucker point of the margin maximization problem, i.e. perform the same implicit regularization as the plain gradient flows. We also prove that the loss decreases to zero and the Euclidean norm of the parameters increases to infinity at the same rates as for the plain gradient flows. Moreover, we consider generalizations of AdaDelta where the exponential decay coefficients may vary with time and the numerical stability terms may be different across the parameters, and we obtain the same results provided the former do not approach 1 too quickly and the latter have isotropic quotients. Finally, we corroborate our theoretical results by numerical experiments on convolutional networks with MNIST and CIFAR-10 datasets.

## 1 Introduction

Understanding when, why, and how training overparameterized neural networks by gradient-based algorithms achieves good generalization (Zhang, Bengio, Hardt, Recht, and Vinyals, 2021; Belkin, Hsu, Ma, and Mandal, 2019) remains one of the central questions in machine learning, despite several years of vibrant research. Much progress on the question has been made by investigating *implicit regularization* (or implicit bias) (Neyshabur, Bhojanapalli, McAllester, and Srebro, 2017): the mysterious preference of the training algorithms for interpolators that perform well at test time.

Building on the seminal work of Soudry, Hoffer, Nacson, Gunasekar, and Srebro (2018), one of the most celebrated results in the field was obtained by Lyu and Li (2020); Ji and Telgarsky (2020): that after achieving perfect training accuracy, gradient flow implicitly regularizes locally Lipschitz, positively homogeneous, and  $\alpha$ -minimally definable networks so that their parameters converge in direction to a margin maximization KKT point. This precise bias towards margin maximization for this wide class of networks has been the basis of numerous theoretical works (cf. Vardi (2023)), as well as remarkable practical methods such as the reconstruction of training data by Haim, Vardi, Yehudai, Shamir, and Irani (2022); Buzaglo, Haim, Yehudai, Vardi, Oz, Nikankin, and Irani (2023).

Nevertheless, already Soudry et al. (2018) observed that algorithms such as Adam (Kingma and Ba, 2015) may not perform the same implicit regularization, and posed the research question:

*Can we characterize the bias of adaptive optimization algorithms for classification problems?*

The pertinence of this question was further attested by Wilson, Roelofs, Stern, Srebro, and Recht (2017), who demonstrated that, for several realistic deep learning models, the solutions found by adaptive methods often generalize significantly worse than those found by stochastic gradient descent, even when the former solutions have better training performance.

Extending the pioneering work of Gunasekar, Lee, Soudry, and Srebro (2018); Qian and Qian (2019), a major advance was made by Wang, Meng, Chen, and Liu (2021) who proved that, for the same wide class of networks as admitted by Lyu & Li (2020); Ji & Telgarsky (2020) and in the continuous setting corresponding to an infinitesimal learning rate, if the adapter of an algorithm without momentum can be shown to converge

to an isotropic vector without large fluctuations, then the implicit regularization is the same as for plain gradient flow. They also established that this holds for RMSProp (Hinton, Srivastava, and Swersky, 2012) and Adam without momentum, but fails for AdaGrad (Duchi, Hazan, and Singer, 2011).

**Our contributions.** In this work, we tackle the posed research question for AdaDelta (Zeiler, 2012), which has remained open. AdaDelta is one of the main adaptive optimization algorithms, implemented in PyTorch<sup>1</sup>, and known to perform well in many circumstances compared to other algorithms including RMSProp and Adam (cf. e.g. Ruder (2016)). Specifically:

- we overcome the technical challenge of AdaDelta having exponentially decaying averages in both the numerator and the denominator of the adapter (which makes it not readily amenable to the techniques of Wang et al. (2021)), and prove that the adapter has the required convergence properties, which enables us to conclude the same implicit regularization as for plain gradient flow;
- we show that the implicit regularization critically depends on the numerical stability terms in the numerator and the denominator of the adapter, and that it changes if they are permitted to have different components so that their quotient is not isotropic;
- we also investigate permitting the exponential decay coefficients to vary with time, and show that under a mild assumption on their integrals, the implicit regularization is not affected;
- we corroborate these theoretical results in three empirical settings, ranging from a simple visualization to a 14-layer convolutional network on CIFAR-10.

**Further related work.** Wang, Meng, Zhang, Sun, Chen, Ma, and Liu (2022) proved that momentum does not affect the implicit bias towards margin maximization for linear classification. Wang, Fu, Zhang, Zheng, and Chen (2023) proved a tight upper bound on Adam’s iteration complexity. Cattaneo, Klusowski, and Shigida (2023) studied the implicit bias of Adam and RMSProp by backward error analysis. Tarzanagh, Li, Zhang, and Oymak (2023) showed margin maximization when optimizing attention by gradient descent. Thilak, Littwin, Zhai, Saremi, Paiss, and Susskind (2024) identified and investigated a slingshot phenomenon in late-stage Adam and related it to grokking (Power, Burda, Edwards, Babuschkin, and Misra, 2022). Xie and Li (2024) showed that AdamW implicitly performs constrained optimization.

## 2 Preliminaries

**Basic notation.** We write:  $[n]$  for the set  $\{1, \dots, n\}$ ;  $\langle \mathbf{u}, \mathbf{v} \rangle$  for the inner product of vectors  $\mathbf{u}$  and  $\mathbf{v}$ ;  $\|\mathbf{v}\| = \sqrt{\langle \mathbf{v}, \mathbf{v} \rangle}$  for the Euclidean length of a vector  $\mathbf{v}$ ;  $\mathbf{v}_i$  for the  $i$ th component of a vector  $\mathbf{v}$ ;  $\mathbf{0}$ ,  $\mathbf{1}$ ,  $\infty$ , etc. for the vectors whose dimension is inferred from the context and whose all components are equal to the specified value, so that for all  $i$  we have  $\mathbf{0}_i = 0$ ,  $\mathbf{1}_i = 1$ ,  $\infty_i = \infty$ , etc.

Provided  $\mathbf{u}$  and  $\mathbf{v}$  are vectors of equal dimensions, we write  $\mathbf{u}\mathbf{v}$ ,  $\mathbf{u}/\mathbf{v}$ ,  $\mathbf{v}^2$ ,  $\sqrt{\mathbf{v}}$ , etc. for the component-wise product, quotient, square, square root, etc. operations, so that for all  $i$  we have  $(\mathbf{u}\mathbf{v})_i = \mathbf{u}_i\mathbf{v}_i$ ,  $(\mathbf{u}/\mathbf{v})_i = \mathbf{u}_i/\mathbf{v}_i$ ,  $(\mathbf{v}^2)_i = (\mathbf{v}_i)^2$ ,  $(\sqrt{\mathbf{v}})_i = \sqrt{\mathbf{v}_i}$ , etc. Similarly, we write  $\mathbf{u} < \mathbf{v}$ ,  $\mathbf{u} \leq \mathbf{v}$ , etc. for the component-wise less, less than or equal, etc. relations, so that we have  $\mathbf{u} < \mathbf{v} \Leftrightarrow \forall i: \mathbf{u}_i < \mathbf{v}_i$ ,  $\mathbf{u} \leq \mathbf{v} \Leftrightarrow \forall i: \mathbf{u}_i \leq \mathbf{v}_i$ , etc.

**Local Lipschitz continuity and the Clarke subdifferential.** Suppose a function  $f: \mathbb{R}^k \rightarrow \mathbb{R}$  is *locally Lipschitz*, i.e. every point  $\mathbf{v} \in \mathbb{R}^k$  has a neighborhood  $U$  such that  $f$  is Lipschitz continuous on  $U$ . By Rademacher’s theorem (cf. e.g. Borwein and Lewis (2010, Theorem 9.1.2)), then  $f$  is differentiable almost everywhere. The *Clarke subdifferential* of  $f$  at a point  $\mathbf{v}$  is the convex hull

$$\partial f(\mathbf{v}) := \text{conv} \left\{ \lim_{i \rightarrow \infty} \nabla f(\mathbf{v}^{(i)}) \mid \lim_{i \rightarrow \infty} \mathbf{v}^{(i)} = \mathbf{v} \text{ and } \nabla f(\mathbf{v}^{(i)}) \text{ exists for all } i \right\}.$$

It is nonempty and compact for all  $\mathbf{v}$ , and equals the singleton  $\{\nabla f(\mathbf{v})\}$  if  $f$  is continuously differentiable at  $\mathbf{v}$  (Clarke, 1975). It consists of *subgradients*, which we may refer to simply as gradients.

<sup>1</sup><https://pytorch.org/docs/stable/generated/torch.optim.Adadelta.html>

---

**Procedure 1** *Discrete generalized AdaDelta.* In step  $k + 1$ , using a current gradient  $\tilde{\partial}\mathcal{L}(\mathbf{w}_k) \in \partial\mathcal{L}(\mathbf{w}(t))$ , it computes the next exponentially decaying average  $\mathbf{g}_{k+1}$  of squared gradients, and computes the next adapted gradient  $\Delta_{k+1}$ . It then computes the next exponentially decaying average  $\mathbf{h}_{k+1}$  of squared adapted gradients, and computes the next parameter vector  $\mathbf{w}_{k+1}$  by subtracting  $\Delta_{k+1}$  scaled by the learning rate. The hyperparameters are: learning rate  $\eta > 0$ , exponential decay coefficients  $\rho_k \in [0, 1]^p$ , and numerical stability terms  $\delta > \mathbf{0}$  and  $\varepsilon > \mathbf{0}$ . The implementation of AdaDelta in PyTorch<sup>1</sup> corresponds to the special case where  $\rho_k = \varrho \mathbf{1}$  for all  $k$ , and  $\delta = \varepsilon = \epsilon \mathbf{1}$ . By specializing further to  $\eta = 1$ , we obtain the original AdaDelta (Zeiler, 2012).

---

$$\begin{aligned} \mathbf{g}_{k+1} &= \rho_k \mathbf{g}_k + (1 - \rho_k) \tilde{\partial}\mathcal{L}(\mathbf{w}_k)^2 \\ \mathbf{h}_{k+1} &= \rho_k \mathbf{h}_k + (1 - \rho_k) \Delta_{k+1}^2 \\ \mathbf{w}_{k+1} &= \mathbf{w}_k - \eta \Delta_{k+1} \end{aligned} \quad \Delta_{k+1} = \sqrt{\frac{\varepsilon + \mathbf{h}_k}{\delta + \mathbf{g}_{k+1}}} \tilde{\partial}\mathcal{L}(\mathbf{w}_k)$$


---

**O-minimal structures and definable functions.** An *o-minimal structure*  $\mathcal{S}$  is a family  $\{\mathcal{S}_k\}_{k=1}^\infty$  such that: each  $\mathcal{S}_k$  is a set of subsets of  $\mathbb{R}^k$ ;  $\mathcal{S}_1$  is the set of all finite unions of open intervals and points; each  $\mathcal{S}_k$  contains the zero sets of all polynomials on  $\mathbb{R}^k$ ; each  $\mathcal{S}_k$  is closed under finite union, finite intersection, and complement; each  $\mathcal{S}_{k+k'}$  contains the Cartesian products of all sets in  $\mathcal{S}_k$  and  $\mathcal{S}_{k'}$ ; each  $\mathcal{S}_k$  contains the projections of all sets in  $\mathcal{S}_{k+1}$  onto the first  $k$  components. A function  $f: \mathbb{R}^k \rightarrow \mathbb{R}^{k'}$  is *definable* in  $\mathcal{S}$  if and only if its graph is a set in  $\mathcal{S}_{k+k'}$ .

For every o-minimal structure, the collection of all definable functions is closed under algebraic operations, composition, inverse, maximum, minimum, etc. (cf. e.g. Ji & Telgarsky (2020, Appendix B)). Moreover, by Wilkie’s theorem (Wilkie, 1996), there exists an o-minimal structure in which the exponential function is definable.

**Predictor and loss functions.** We assume the following properties of the predictor function  $\Phi(\mathbf{w}, \mathbf{x})$  with parameters  $\mathbf{w} \in \mathbb{R}^p$ , inputs  $\mathbf{x} \in \mathbb{R}^d$ , and scalar outputs. The  $L$ -positive homogeneity means that scaling the parameters  $\mathbf{w}$  by any  $\alpha > 0$  scales the output by  $\alpha^L$ , i.e.  $\Phi(\alpha\mathbf{w}, \mathbf{x}) = \alpha^L \Phi(\mathbf{w}, \mathbf{x})$ .

**Assumption 1.** For some  $L > 0$  and some o-minimal structure  $\mathcal{S}$  in which the exponential function is definable, for each  $\mathbf{x}$  we have that  $\Phi(\mathbf{w}, \mathbf{x})$  as a function of  $\mathbf{w}$  is: (i) locally Lipschitz; (ii)  $L$ -positively homogeneous; (iii) definable in  $\mathcal{S}$ .

This assumption admits neural networks that are constructed from a wide variety of layer types, including fully connected, convolutional, ReLU, Leaky ReLU, and max-pooling, which may be composed arbitrarily. Points of nondifferentiability, such as at 0 in the case of the ReLU nonlinearity, are permitted because we assume only local Lipschitzness instead of continuous differentiability and we work with the Clarke subdifferential instead of the gradient. However, due to the  $L$ -positive homogeneity (ii), skip connections are excluded, and biases are excluded except at the first layer.

We consider minimizing the total loss  $\mathcal{L}(\mathbf{w}) := \sum_{i=1}^n \ell(y^{(i)} \Phi(\mathbf{w}, \mathbf{x}^{(i)}))$ , which is with respect to a finite dataset  $\{(\mathbf{x}^{(i)} \in \mathbb{R}^d, y^{(i)} \in \{\pm 1\})\}_{i=1}^n$  and where the individual loss function is:  $\ell(z) := e^{-z}$  for the exponential loss, and  $\ell(z) := \log(1 + e^{-z})$  for the logistic loss.

Since in both cases the individual loss functions  $\ell$  is locally Lipschitz and definable in the o-minimal structure  $\mathcal{S}$  from Assumption 1, the same is true of the total loss function  $\mathcal{L}$ .

**Generalized AdaDelta flow trajectories.** The starting point of our analysis is the adaptive gradient descent of Procedure 1, which is a generalization of AdaDelta (Zeiler, 2012) by allowing: the learning rate to be specified (this is already the case in the PyTorch implementation<sup>1</sup>), the exponential decay coefficients to vary with time (e.g. by following a specified schedule), and both those coefficients and the numerical stability terms to be specified differently across the vector components.

The focus of our theoretical study is the adaptive gradient flow that corresponds to the adaptive gradient descent with an infinitesimal learning rate. To arrive at its definition, we first restate the equations of

---

**Process 2** *Continuous generalized AdaDelta.* This adaptive gradient flow, where  $\tilde{\partial}\mathcal{L}(\mathbf{w}(t)) \in \partial\mathcal{L}(\mathbf{w}(t))$  is a gradient of the loss at the current parameters  $\mathbf{w}(t)$ , corresponds to the adaptive gradient descent of Procedure 1 with an infinitesimal learning rate. The hyperparameters are: exponential decay coefficients schedule  $\boldsymbol{\rho}: [0, \infty) \rightarrow [0, 1]^p$ , and numerical stability terms  $\boldsymbol{\delta} > \mathbf{0}$  and  $\boldsymbol{\varepsilon} > \mathbf{0}$ .

---

$$\mathbf{g}'(t) = (\mathbf{1} - \boldsymbol{\rho}(t))(\tilde{\partial}\mathcal{L}(\mathbf{w}(t))^2 - \mathbf{g}(t)) \quad (1)$$

$$\mathbf{h}'(t) = (\mathbf{1} - \boldsymbol{\rho}(t))(\mathbf{w}'(t)^2 - \mathbf{h}(t)) \quad (2)$$

$$\mathbf{w}'(t) = -\sqrt{\frac{\boldsymbol{\varepsilon} + \mathbf{h}(t)}{\boldsymbol{\delta} + \mathbf{g}(t)}} \tilde{\partial}\mathcal{L}(\mathbf{w}(t)) \quad (3)$$


---

Procedure 1 as follows. We eliminate the auxiliary variable  $\boldsymbol{\Delta}_{k+1}$ , we suppose the hyperparameters  $\boldsymbol{\rho}_k$  obey that  $(\mathbf{1} - \boldsymbol{\rho}_k)/\eta$  is constant with respect to the hyperparameter  $\eta$ , and we replace steps  $k$  by times  $t = k\eta$ :

$$\begin{aligned} \frac{\mathbf{g}(t+\eta) - \mathbf{g}(t)}{\eta} &= (\mathbf{1} - \boldsymbol{\rho}(t)) \left( \tilde{\partial}\mathcal{L}(\mathbf{w}(t))^2 - \mathbf{g}(t) \right) \\ \frac{\mathbf{h}(t+\eta) - \mathbf{h}(t)}{\eta} &= (\mathbf{1} - \boldsymbol{\rho}(t)) \left( \left( \frac{\mathbf{w}(t+\eta) - \mathbf{w}(t)}{\eta} \right)^2 - \mathbf{h}(t) \right) \\ \frac{\mathbf{w}(t+\eta) - \mathbf{w}(t)}{\eta} &= -\sqrt{\frac{\boldsymbol{\varepsilon} + \mathbf{h}(t)}{\boldsymbol{\delta} + \mathbf{g}(t+\eta)}} \tilde{\partial}\mathcal{L}(\mathbf{w}(t)) . \end{aligned}$$

Now, regarding these equations as determining the endpoints  $\mathbf{g}(t+\eta), \mathbf{h}(t+\eta), \mathbf{w}(t+\eta)$  of the next line segments in some continuous polygonal curves  $\mathbf{g}, \mathbf{h}, \mathbf{w}: [0, \infty) \rightarrow \mathbb{R}^p$  born from the adaptive gradient descent of Procedure 1 with learning rate  $\eta$  (cf. e.g. [Elkabetz and Cohen \(2021\)](#)), letting  $\eta$  tend to 0 we obtain that the limits of their directions are given by the right-hand sides of eqs. (1) to (3), which we take to be the derivatives that determine the adaptive gradient flow we are seeking.

We therefore analyze *trajectories*  $\mathbf{g}, \mathbf{h}, \mathbf{w}: [0, \infty) \rightarrow \mathbb{R}^p$  of the two exponentially decaying averages and of the parameters, which are arcs (i.e. absolutely continuous on every compact subinterval) and which obey the adaptive gradient flow of Process 2 for almost all  $t \geq 0$ .

**Assumption 2.** (i)  $\int_0^\infty (\mathbf{1} - \boldsymbol{\rho}(t))dt = \infty$ . (ii)  $\mathbf{g}(0) \geq \mathbf{0}$  and  $\mathbf{h}(0) \geq \mathbf{0}$ . (iii) There exists a time  $t_0$  such that  $\mathcal{L}(\mathbf{w}(t_0)) < \ell(0)$ .

This is a mild regularity assumption. Part (i) ensures that the exponential decay coefficients are not scheduled to approach 1 so fast that they stop the learning, and for example it is satisfied by any constant coefficients smaller than 1. Part (ii) just requires that the exponentially decaying averages of squared gradients and squared adapted gradients are initialized as nonnegative (in original AdaDelta they are initialized to zero). Part (iii) assumes that the dataset is separable by the network, moreover that for every parameters trajectory  $\mathbf{w}$  that we consider arising from the generalized AdaDelta flow, there exists a time  $t_0$  at which a separation (i.e. the perfect training accuracy) is achieved; this commonly occurs when training overparameterized neural networks by gradient based algorithms (cf. e.g. [Zhang et al. \(2021\)](#)).

**Admittance of a chain rule.** Since the total loss function  $\mathcal{L}$  is locally Lipschitz and definable in  $\mathcal{S}$  (for both the exponential and the logistic individual loss functions), and the trajectory  $\mathbf{w}$  of the parameters is an arc, we are able to use the following fact applied to them.

**Proposition 1** ([Davis, Drusvyatskiy, Kakade, and Lee \(2020, Theorem 5.8\)](#)). *If  $f: \mathbb{R}^k \rightarrow \mathbb{R}$  is locally Lipschitz and definable in an o-minimal structure, then it admits a chain rule: for all arcs  $\mathbf{v}: [0, \infty) \rightarrow \mathbb{R}^k$ , almost all  $t \geq 0$ , and all  $\mathbf{u} \in \partial f(\mathbf{v}(t))$ , we have  $df(\mathbf{v}(t))/dt = \langle \mathbf{u}, d\mathbf{v}(t)/dt \rangle$ .*

**Directional convergence, KKT conditions, and margin maximization.** That a trajectory  $\mathbf{v}$  converges in direction to a vector  $\mathbf{u}$  means  $\lim_{t \rightarrow \infty} \mathbf{v}(t)/\|\mathbf{v}(t)\| = \mathbf{u}/\|\mathbf{u}\|$ .

Following [Dutta, Deb, Tulshyan, and Arora \(2013, Section 2.2\)](#), supposing  $f, g_1, \dots, g_n : \mathbb{R}^k \rightarrow \mathbb{R}$  are locally Lipschitz, we have that  $\mathbf{v} \in \mathbb{R}^k$  is a *Karush-Kuhn-Tucker point* of the problem

$$\text{minimize: } f(\mathbf{v}) \quad \text{subject to: } g_i(\mathbf{v}) \leq 0 \text{ for all } i \in [n]$$

if and only if there exist Lagrange multipliers  $\lambda_1, \dots, \lambda_n \geq 0$  such that:

$$\text{(feasibility)} \quad g_i(\mathbf{v}) \leq 0 \text{ for all } i \in [n];$$

$$\text{(equilibrium inclusion)} \quad \mathbf{0} \in \partial f(\mathbf{v}) + \sum_{i=1}^n \lambda_i \partial g_i(\mathbf{v});$$

$$\text{(complementary slackness)} \quad \lambda_i g_i(\mathbf{v}) = 0 \text{ for all } i \in [n].$$

For local optimality, these first-order conditions are necessary, however in general not sufficient.

The *margin* of a parameters vector  $\mathbf{w}$  is the smallest label-adjusted prediction for an input, i.e. formally  $\min_{i \in [n]} y^{(i)} \Phi(\mathbf{w}, \mathbf{x}^{(i)})$ . By the  $L$ -positive homogeneity of the predictor function (cf. Assumption 1.(ii)), the *normalized margin*  $\min_{i \in [n]} y^{(i)} \Phi(\mathbf{w}, \mathbf{x}^{(i)}) / \|\mathbf{w}\|^L$  depends only on the direction of  $\mathbf{w}$ , and it is straightforward to show that the directions that maximize it are also the optimal directions of the problem

$$\text{minimize: } \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{subject to: } y^{(i)} \Phi(\mathbf{w}, \mathbf{x}^{(i)}) \geq 1 \text{ for all } i \in [n]. \quad (4)$$

### 3 Main result

We prove that continuous generalized AdaDelta obeys the same tight rates for convergence of the loss and growth of the parameters as were established for plain gradient flow by [Lyu & Li \(2020, Corollary A.11\)](#), and also implicitly regularizes the parameters to converge in direction to a KKT point of a variant of the margin maximization problem in eq. (4), in which the objective function  $\frac{1}{2} \|\mathbf{w}\|^2$  is replaced by  $\frac{1}{2} \|\sqrt{\frac{\delta}{\epsilon}} \mathbf{w}\|^2$ . If the quotient  $\delta/\epsilon$  of the numerical stability hyperparameters is isotropic (i.e.  $\delta_j/\epsilon_j$  are equal for all  $j \in [p]$ ), then this modification does not alter the directions of the KKT points, so the implicit regularization is the same as was established for plain gradient flow by [Lyu & Li \(2020, Theorem A.8\)](#) and [Ji & Telgarsky \(2020, Theorem 3.1\)](#). Otherwise, the sets of points  $\mathbf{w}$  that have the same objective value  $\frac{1}{2} \|\sqrt{\frac{\delta}{\epsilon}} \mathbf{w}\|^2$  are ellipsoids which are not spheres, and the directions of the KKT points may be different from those for the margin maximization problem in eq. (4). Moreover, these conclusions are robust with respect to changing the exponential decay coefficients schedule  $\rho$  hyperparameter.

**Theorem 2.** *Under Assumptions 1 and 2, we have that  $\mathcal{L}(\mathbf{w}(t)) = \Theta\left(\frac{1}{t(\log t)^{2-2/L}}\right)$ ,  $\|\mathbf{w}(t)\| = \Theta((\log t)^{1/L})$ , and  $\mathbf{w}(t)$  converges in direction to a KKT point of the problem*

$$\text{minimize: } \frac{1}{2} \left\| \sqrt{\frac{\delta}{\epsilon}} \mathbf{w} \right\|^2 \quad \text{subject to: } y^{(i)} \Phi(\mathbf{w}, \mathbf{x}^{(i)}) \geq 1 \text{ for all } i \in [n].$$

A key role in the proof of Theorem 2 is played by the *skewed adapter*  $\beta(t) := \sqrt{\frac{1+h(t)/\epsilon}{1+g(t)/\delta}}$ , which is defined by multiplying the adapter  $\sqrt{\frac{\epsilon+h(t)}{\delta+g(t)}}$  of the continuous generalized AdaDelta (cf. Process 2) component-wise by the square rooted quotient  $\sqrt{\frac{\delta}{\epsilon}}$  of the numerical stability hyperparameters. Equation (3) that governs the parameters trajectory can then be restated as

$$\mathbf{w}'(t) = -\sqrt{\frac{\epsilon}{\delta}} \beta(t) \tilde{\partial} \mathcal{L}(\mathbf{w}(t)) \quad \text{for almost all } t \geq 0. \quad (5)$$

Our first lemma shows a useful expression for the derivative of  $\mathbf{1}$  minus the squared skewed adapter.

**Lemma 3.** *For almost all  $t \geq 0$  we have  $\frac{d(1-\beta(t)^2)}{dt} = -(1-\rho(t)) \frac{1-\beta(t)^2}{1+g(t)/\delta}$ .*

*Proof.* Observe that

$$\begin{aligned}
\frac{d(1 - \beta(t)^2)}{dt} &= -\frac{\mathbf{h}'(t)/\varepsilon + \mathbf{h}'(t)\mathbf{g}(t)/\varepsilon\delta - \mathbf{g}'(t)/\delta - \mathbf{g}'(t)\mathbf{h}(t)/\delta\varepsilon}{(1 + \mathbf{g}(t)/\delta)^2} && \text{by the def. of } \beta(t) \\
&= -(1 - \rho(t)) \frac{\mathbf{w}'(t)^2/\varepsilon - \mathbf{h}(t)/\varepsilon - \tilde{\partial}\mathcal{L}(\mathbf{w}(t))^2/\delta + \mathbf{g}(t)/\delta + \mathbf{w}'(t)^2\mathbf{g}(t)/\varepsilon\delta - \tilde{\partial}\mathcal{L}(\mathbf{w}(t))^2\mathbf{h}(t)/\delta\varepsilon}{(1 + \mathbf{g}(t)/\delta)^2} && \text{by eqs. (1) and (2)} \\
&= -(1 - \rho(t)) \frac{\mathbf{g}(t)/\delta - \mathbf{h}(t)/\varepsilon + (\mathbf{w}'(t)^2/\varepsilon)(1 + \mathbf{g}(t)/\delta) - (\tilde{\partial}\mathcal{L}(\mathbf{w}(t))^2/\delta)(1 + \mathbf{h}(t)/\varepsilon)}{(1 + \mathbf{g}(t)/\delta)^2} && \text{rearranging} \\
&= -(1 - \rho(t)) \frac{\mathbf{g}(t)/\delta - \mathbf{h}(t)/\varepsilon}{(1 + \mathbf{g}(t)/\delta)^2} && \text{by eq. (5)} \\
&= -(1 - \rho(t)) \frac{1 - \beta(t)^2}{1 + \mathbf{g}(t)/\delta} && \text{calculation}
\end{aligned}$$

for almost all  $t \geq 0$ .  $\square$

By Assumption 2.(ii), each component of the skewed adapter is initially positive. The second lemma shows that it never drops below that value or 1, whichever is smaller.

**Lemma 4.** *For all  $t \geq 0$  we have  $\beta(t) \geq \min\{\beta(0), 1\}$ .*

*Proof.* We first note that, if any component  $\mathbf{g}(t)_j$  or  $\mathbf{h}(t)_j$  of the exponentially decaying averages were negative at any time  $t \geq 0$ , then by eqs. (1) and (2) its derivative (provided it exists) would necessarily be positive. Recalling that  $\mathbf{g}(0) \geq \mathbf{0}$  and  $\mathbf{h}(0) \geq \mathbf{0}$  (cf. Assumption 2.(ii)), and that  $\mathbf{g}$  and  $\mathbf{h}$  are arcs, we infer:  $\mathbf{g}(t) \geq \mathbf{0}$  and  $\mathbf{h}(t) \geq \mathbf{0}$  for all  $t \geq 0$ .

Now, from the nonnegativity of  $\mathbf{g}$ , Lemma 3, and the fact that  $d\beta(t)^2/dt = 2\beta(t)\beta'(t)$ , for almost all  $t \geq 0$  and all  $j \in [p]$  we have that: if  $0 < \beta(t)_j < 1$  then  $\beta'(t)_j$  is nonnegative, if  $\beta(t)_j = 1$  then  $\beta'(t)_j$  is zero, and if  $\beta(t)_j > 1$  then  $\beta'(t)_j$  is nonpositive. These properties, together with the fact that the skewed adapter  $\beta$  is an arc, imply the lemma.  $\square$

Lemma 4 provides a lower bound for the skewed adapter, but it leaves open its asymptotic behavior. The next lemma fills that gap, establishing that all its components converge to 1. This equivalently means that the adapter  $\sqrt{\frac{\varepsilon + \mathbf{h}(t)}{\delta + \mathbf{g}(t)}}$  converges to  $\sqrt{\frac{\varepsilon}{\delta}}$ .

**Lemma 5.** *We have that  $\lim_{t \rightarrow \infty} \beta(t) = 1$ .*

*Proof.* We first use Lemma 4 to prove that the squared gradients have bounded integrals:

$$\begin{aligned}
\mathcal{L}(\mathbf{w}(0)) &\geq -\int_0^\infty (d\mathcal{L}(\mathbf{w}(t))/dt)dt && \text{since } \mathcal{L}(\mathbf{w}(t)) > 0 \text{ for all } t \geq 0 \\
&= -\int_0^\infty \langle \tilde{\partial}\mathcal{L}(\mathbf{w}(t)), \mathbf{w}'(t) \rangle dt && \text{by Proposition 1} \\
&= \int_0^\infty \langle \tilde{\partial}\mathcal{L}(\mathbf{w}(t)), \sqrt{\frac{\varepsilon}{\delta}} \beta(t) \tilde{\partial}\mathcal{L}(\mathbf{w}(t)) \rangle dt && \text{by eq. (5)} \\
&\geq \int_0^\infty \langle \sqrt{\frac{\varepsilon}{\delta}} \min\{\beta(0), 1\}, \tilde{\partial}\mathcal{L}(\mathbf{w}(t))^2 \rangle dt && \text{by Lemma 4} \\
&= \sum_{j \in [p]} \sqrt{\frac{\varepsilon_j}{\delta_j}} \min\{\beta(0)_j, 1\} \int_0^\infty \tilde{\partial}\mathcal{L}(\mathbf{w}(t))_j^2 dt && \text{rearranging ,}
\end{aligned}$$

so for each  $j \in [p]$  we have  $\int_0^\infty \tilde{\partial}\mathcal{L}(\mathbf{w}(t))_j^2 dt \leq \sqrt{\frac{\delta_j}{\varepsilon_j}} \frac{\mathcal{L}(\mathbf{w}(0))}{\min\{\beta(0)_j, 1\}} =: C_j$ .



Now suppose  $j \in [p]$ . Equation (1) implies that

$$\mathbf{g}(t)_j \leq \mathbf{g}(0)_j + C_j =: C_j^\dagger \quad \text{for all } t \geq 0,$$

so recalling Lemmas 3 and 4 and setting  $C_j^\ddagger := \frac{1}{1+C_j^\dagger/\delta_j}$  we have

$$d \log |1 - \beta(t)_j^2|/dt \leq -C_j^\ddagger(1 - \rho(t)_j) \quad \text{for almost all } t \geq 0 \text{ such that } \beta(t)_j \neq 1.$$

Hence

$$|1 - \beta(t)_j^2| \leq |1 - \beta(0)_j^2| \exp \left( -C_j^\ddagger \int_0^t (1 - \rho(\tau)_j) d\tau \right) \quad \text{for all } t \geq 0,$$

which by positivity of  $C_j^\ddagger$  and unboundedness of the integrals of the complements of the exponential decay coefficients (cf. Assumption 2.(i)) establishes the lemma.  $\square$

Our final lemma shows that the components of the skewed adapter converge without large fluctuations.

**Lemma 6.** *The function  $d \log \beta(t)/dt$  is Lebesgue integrable.*

*Proof.* For each  $j \in [p]$ , from the proof of Lemma 4 we have that  $d \log \beta(t)_j/dt = \beta'(t)_j/\beta(t)_j$  is either nonnegative for almost all  $t \geq 0$ , or nonpositive for almost all  $t \geq 0$ . Thus by Lemma 5 we have

$$\int_0^\infty |d \log \beta(t)/dt| dt = \left| \int_0^\infty (d \log \beta(t)/dt) dt \right| = |\log \mathbf{1} - \log \beta(0)| = |\log \beta(0)| < \infty. \quad \square$$

Theorem 2 now follows from Lemmas 5 and 6 by applying Wang et al. (2021, Theorems 2, 3, and 10)<sup>2</sup> to the flow

$$\mathbf{v}'(t) = -\beta(t) \tilde{\partial} \hat{\mathcal{L}}(\mathbf{v}(t)) \quad \text{for almost all } t \geq 0,$$

which is a restatement of eq. (5) with  $\mathbf{v}(t) := \sqrt[4]{\frac{\delta}{\epsilon}} \mathbf{w}(t)$  and  $\hat{\mathcal{L}}(\mathbf{u}) := \mathcal{L}(\sqrt[4]{\frac{\epsilon}{\delta}} \mathbf{u})$ .

## 4 Experiments

To test our theoretical result for generalized AdaDelta, we evaluated comparatively five algorithms:

**SGD:** Stochastic gradient descent as implemented in PyTorch<sup>3</sup>.

**AdaDelta:** Its standard PyTorch implementation<sup>1</sup>, with the exponential decay coefficient  $\varrho = 0.9$ , and the numerical stability term  $\epsilon = 10^{-5}$ .

**AdaDeltaS:** This is AdaDelta amended to have the exponential decay coefficients follow the schedule  $\varrho_k = 1 - 0.1/(1 + \lfloor 100k/K \rfloor)$ , where  $K$  is the total number of steps. Thus  $1 - \varrho_k$  follows a harmonic sequence, increasing the coefficient from  $\varrho_0 = 0.9$  at the first step to  $\varrho_{K-1} = 0.999$  at the last step, which lessens aggressiveness of the decay in computing the averages of the squared gradients and the squared adapted gradients along the training.

**AdaDeltaN:** This is AdaDelta amended to have the numerical stability terms different in the numerator and the denominator of the adaptor, and different across the network parameters. At the start of the training, each component of  $\delta$  and of  $\epsilon$  is sampled independently from  $10^{-5+X}$ , where  $X$  is a centered Gaussian with standard deviation 1 for the smaller two networks we consider, and with standard deviation 0.25 for the largest network.

<sup>2</sup>Wang et al. (2021, Theorem 10) contains a typo: the rate  $\Theta\left(\frac{1}{(\log t)^{1/L}}\right)$  should be  $\Theta((\log t)^{1/L})$ .

<sup>3</sup><https://pytorch.org/docs/stable/generated/torch.optim.SGD.html>

**AdaDeltaNS:** This combines the generalizations in AdaDeltaS and AdaDeltaN, i.e. has exponential decay coefficients that follow the specified schedule as well as numerical stability terms whose components are initialized randomly as above.

We performed experiments in the following three gradually more complex settings. The total compute for the perceptron setting was around 10min on a mid-range CPU, whereas an experiment for a single algorithm on one of the two convolutional networks took around 2h on a mid-range GPU. In the plots for the experiments on the convolutional networks, the solid lines show the median values, and the shaded areas are between the 25th and 75th percentiles, over five runs of each experiment. Our experiments can be replicated by making: simple amendments to the PyTorch implementation of AdaDelta<sup>1</sup> to obtain AdaDeltaS, AdaDeltaN, and AdaDeltaNS; and straightforward extensions to the code of Wang et al. (2021) to implement the VGG network and load CIFAR-10.

**Two-layer Leaky ReLU perceptron on a two-dimensional dataset.** Following Wang et al. (2021), we first considered a perceptron with parameters  $v \in \mathbb{R}$  and  $\mathbf{w} \in \mathbb{R}^2$ , whose prediction for an input  $\mathbf{x} \in \mathbb{R}^2$  equals  $v \sigma(\langle \mathbf{w}, \mathbf{x} \rangle)$ , where  $\sigma$  is the Leaky ReLU nonlinearity with inactive gradient 0.5.

The dataset  $\{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(100)}, y^{(100)})\}$  consists of 50 points sampled from  $(\cos 0.5, \sin 0.5) + \mathbf{u}$  independently and labelled 1, and 50 points sampled from  $-(\cos 0.5, \sin 0.5) + \mathbf{u}$  independently and labelled  $-1$ , where  $\mathbf{u}$  is distributed uniformly on the square  $[-0.6, 0.6]^2$ . An instance is depicted in fig. 1.(a).

This toy setting is convenient for three-dimensional visualizations of the adapter’s reciprocal square roots. We trained the network for  $K = 5000$  full-batch epochs using the exponential loss, with learning rate 0.1 for SGD and learning rate 1 for the four variants of AdaDelta. We repeated the training 100 times, where in each round the five algorithms used the same randomly initialized parameters and numerical stability terms (if applicable). The results are shown in fig. 1.(b)–(g).

We observe that:

- all rounds achieved 100% training accuracy by around 256 epochs, in line with the separation Assumption 2.(iii);
- the final normalized margin, which here equals  $\min_{i \in [100]} \frac{y^{(i)} v_K \sigma(\langle \mathbf{w}_K, \mathbf{x}^{(i)} \rangle)}{v_K^2 + \|\mathbf{w}_K\|^2}$ , is within a higher and much narrower range for SGD and with isotropic numerical stability hyperparameters (i.e. for AdaDelta and AdaDeltaS), confirming the prediction of Theorem 2 that, without the isotropy, the implicit regularization may not be towards margin maximization;
- the final adapter’s reciprocal square root, which in the notations of Procedure 1 equals  $\sqrt[4]{\frac{\delta + \mathbf{g}_K}{\epsilon + \mathbf{h}_{K-1}}}$ , has a large variance of its direction when the numerical stability hyperparameters have random components (i.e. for AdaDeltaN and AdaDeltaNS), and it is the limit of this direction that according to the proof of Theorem 2 determines the nature of the implicit regularization.

**Four-layer convolutional network on MNIST.** Also following Wang et al. (2021), we then considered the network from Mądry, Makelov, Schmidt, Tsipras, and Vladu (2018) with biases removed, and in order consisting of: 32-channel  $5 \times 5$ -filter convolutional, ReLU, 2-kernel 2-stride max-pooling, 64-channel  $3 \times 3$ -filter convolutional, ReLU, 2-kernel 2-stride max-pooling, 1024-width fully connected, ReLU, and 10-width fully connected.

We trained the network on MNIST (LeCun, Bottou, Bengio, and Haffner, 1998) from the default PyTorch random initialization for 500 epochs using the cross-entropy loss: in a finer regime with batch size 100, learning rate 0.01 for SGD, and learning rate 0.1 for the four variants of AdaDelta; and in a coarser regime with batch size 1000, learning rate 0.1 for SGD, and learning rate 1 for the four variants of AdaDelta. The results are shown in fig. 2.

For both regimes, we observe that:

- all algorithms achieved perfect training accuracy by around 250 epochs, in line with the separation Assumption 2.(iii);



- the normalized margin, whose values here are relatively small partly due to the division by the fourth power of the norm of the network parameters, grows significantly higher for SGD and with isotropic numerical stability hyperparameters (i.e. for AdaDelta and AdaDeltaS), confirming the prediction of Theorem 2 that, without the isotropy, the implicit regularization may not be towards margin maximization;
- the test accuracy is consistently high for AdaDelta and AdaDeltaS, corroborating the link between the normalized margin and generalization (cf. e.g. Jiang, Neyshabur, Mobahi, Krishnan, and Bengio (2020));
- the training loss shrinks faster when the numerical stability hyperparameters have random components (i.e. for AdaDeltaN and AdaDeltaNS);
- the results are not substantially affected by whether the exponential decay coefficients follow the increasing schedule.

**VGG on CIFAR-10.** Following Lyu & Li (2020), we finally considered the 14-layer VGG-16 (Simonyan and Zisserman, 2015) with biases only at the first level, and in order consisting of: 64-channel convolutional then ReLU, repeated 2 times; max-pooling; 128-channel convolutional then ReLU, repeated 2 times; max-pooling; 256-channel convolutional then ReLU, repeated 3 times; max-pooling; 512-channel convolutional then ReLU, repeated 3 times; max-pooling; 256-channel convolutional then ReLU, repeated 3 times; max-pooling; 10-width fully connected. Each convolutional layer is  $3 \times 3$ -filter, and each max-pooling is 2-kernel 2-stride.

We trained the network on CIFAR-10 (Krizhevsky, 2009) from the default PyTorch random initialization for 1000 epochs using the cross-entropy loss: in a finer regime with batch size 100, and learning rate 0.1 for all five algorithms; and in a coarser regime with batch size 250, and learning rate 0.25 for all five algorithms. The results are shown in fig. 3.

Our observations are similar as for the previous setting, except that in the larger learning rate regime, the increasing schedule of the exponential decay coefficient (i.e. in AdaDeltaS and AdaDeltaNS) seems to help with both optimization and generalization.

## 5 Conclusion

Our main result, Theorem 2 which holds under Assumptions 1 and 2, indicates that AdaDelta may be used in practical methods that rely on the optimization algorithm implicitly regularizing the network parameters to converge in direction to a margin maximization KKT point.

Relaxing Assumption 1 on the predictor function is challenging even without considering adaptive algorithms, e.g. the definability excludes pathological examples such as based on the “Mexican hat” function (cf. Lyu & Li (2020, Appendix J)).

We focused on binary classification for simplicity of presentation; it is straightforward to extend Theorem 2 for logistic loss to an arbitrary number of classes (cf. Wang et al. (2021, Appendix E)).

An important future goal is to obtain a counterpart of Theorem 2 directly for adaptive gradient descent as in Procedure 1. This is also a challenge already without adaptivity: the directional convergence result (Ji & Telgarsky, 2020, Theorem 3.1) is only for gradient flow, and the characterization of directional limits for gradient descent (Lyu & Li, 2020, Theorem E.3) assumes  $C^2$ -smoothness of the predictor function which excludes nonlinearities such as ReLU and Leaky ReLU.

A potentially interesting direction for further empirical work is to evaluate more systematically the extensions of AdaDelta that we considered, in particular the benefits of anisotropic numerical stability terms and of scheduled exponential decay coefficients that we observed in some experiments.

### Broader Impact Statement

This is foundational research on a general algorithm for optimizing neural networks. Greater understanding of its implicit regularization properties may lead to machine learning models that have cheaper training, greater efficiency, and increased performance.

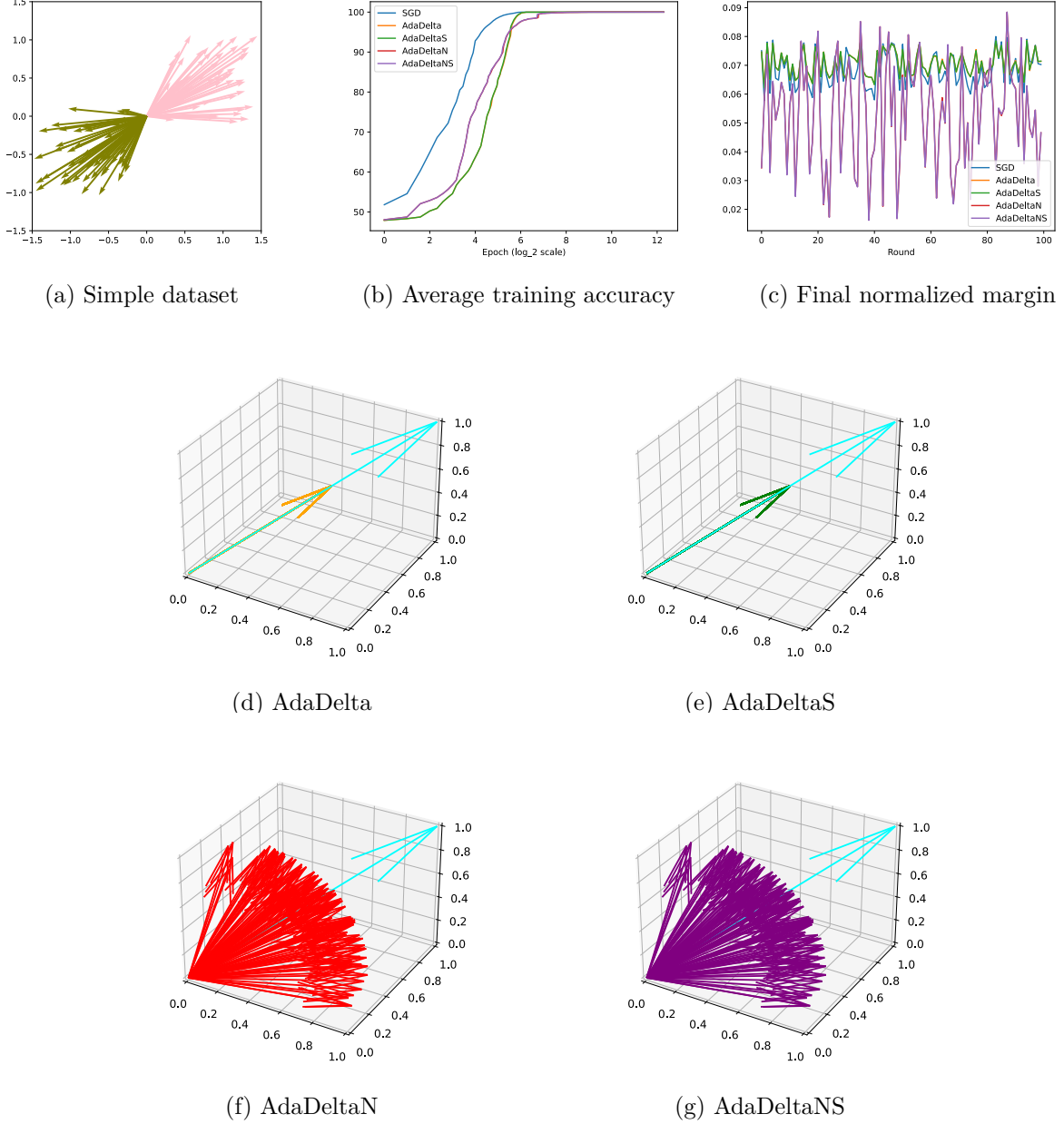


Figure 1: We trained a two-layer Leaky ReLU perceptron on a binary classification dataset depicted in plot (a) for 100 rounds. A round consisted of randomly initializing the network parameters and the numerical stability terms (if applicable), and then running separately each of the five algorithms for 5000 epochs. Plot (b) shows the training accuracies across the epochs averaged over the rounds, and plot (c) shows the final normalized margin across the rounds. In plots (b) and (c), the differences between AdaDelta and AdaDeltaS are small and thus hardly visible, and similarly for AdaDeltaN and AdaDeltaNS. Plots (d)–(g) visualize the direction of the final adapter’s reciprocal square root, where the isotropic direction is indicated by the longer cyan arrow.

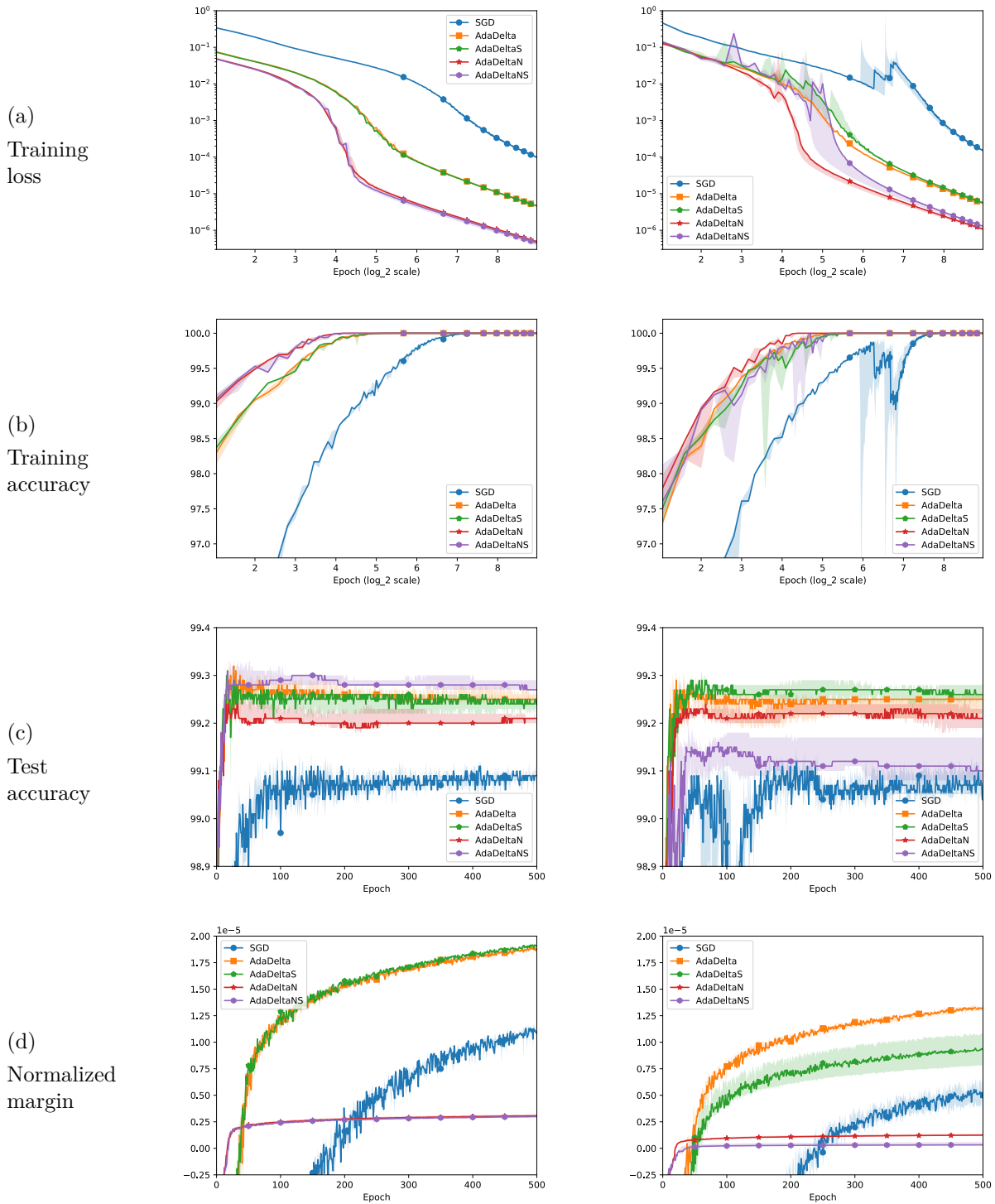
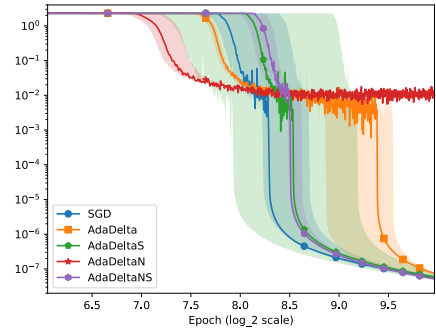
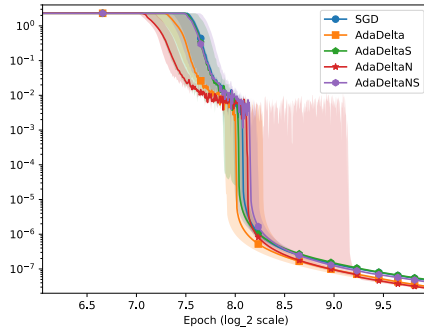
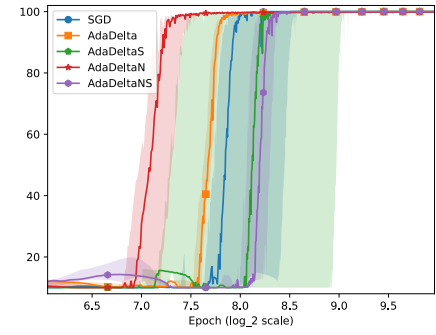
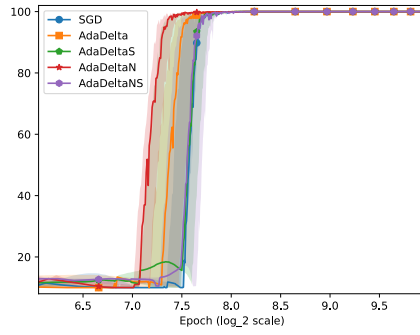


Figure 2: The plots show the results of training a 4-layer convolutional network on MNIST, with batch sizes and learning rates for the right-hand column that are 10 times larger than for the left-hand column.

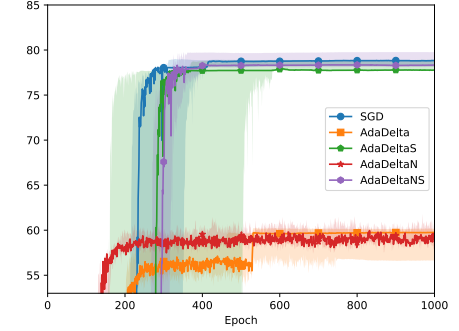
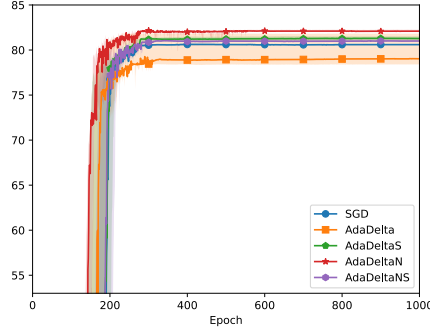
(a)  
Training  
loss



(b)  
Training  
accuracy



(c)  
Test  
accuracy



(d)  
Normalized  
margin

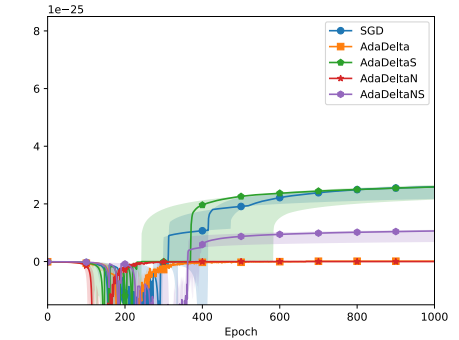
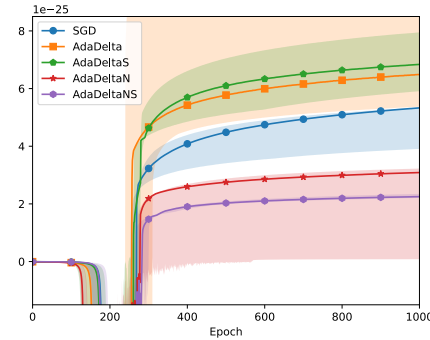


Figure 3: The plots show the results of training a 14-layer convolutional network on CIFAR-10, with batch sizes and learning rates for the right-hand column that are 2.5 times larger than for the left-hand column.

## References

- Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. [Reconciling modern machine-learning practice and the classical bias–variance trade-off](#). *Proc. Natl. Acad. Sci.*, 116(32):15849–15854, 2019. [1](#)
- Jonathan Borwein and Adrian Lewis. *Convex Analysis and Nonlinear Optimization: Theory and Examples*. Springer, 2nd edition, 2010. [2](#)
- Gon Buzaglo, Niv Haim, Gilad Yehudai, Gal Vardi, Yakir Oz, Yaniv Nikankin, and Michal Irani. [Deconstructing Data Reconstruction: Multiclass, Weight Decay and General Losses](#). In *NeurIPS*, 2023. [1](#)
- Matias D. Cattaneo, Jason M. Klusowski, and Boris Shigida. [On the Implicit Bias of Adam](#). *CoRR*, abs/2309.00079, 2023. [2](#)
- Frank H. Clarke. [Generalized gradients and applications](#). *Trans. Amer. Math. Soc.*, 205:247–262, 1975. [2](#)
- Damek Davis, Dmitriy Drusvyatskiy, Sham M. Kakade, and Jason D. Lee. [Stochastic Subgradient Method Converges on Tame Functions](#). *Found. Comput. Math.*, 20(1):119–154, 2020. [4](#)
- John C. Duchi, Elad Hazan, and Yoram Singer. [Adaptive Subgradient Methods for Online Learning and Stochastic Optimization](#). *J. Mach. Learn. Res.*, 12:2121–2159, 2011. [2](#)
- Joydeep Dutta, Kalyanmoy Deb, Rupesh Tulshyan, and Ramnik Arora. [Approximate KKT points and a proximity measure for termination](#). *J. Glob. Optim.*, 56(4):1463–1499, 2013. [5](#)
- Omer Elkabetz and Nadav Cohen. [Continuous vs. Discrete Optimization of Deep Neural Networks](#). In *NeurIPS*, pp. 4947–4960, 2021. [4](#)
- Suriya Gunasekar, Jason D. Lee, Daniel Soudry, and Nathan Srebro. [Characterizing Implicit Bias in Terms of Optimization Geometry](#). In *ICML*, pp. 1827–1836, 2018. [1](#)
- Niv Haim, Gal Vardi, Gilad Yehudai, Ohad Shamir, and Michal Irani. [Reconstructing Training Data From Trained Neural Networks](#). In *NeurIPS*, 2022. [1](#)
- Geoffrey Hinton, Nitish Srivastava, and Kevin Swersky. [Overview of mini-batch gradient descent](#), 2012. *Neural Networks for Machine Learning: Lecture 6a*. [2](#)
- Ziwei Ji and Matus Telgarsky. [Directional convergence and alignment in deep learning](#). In *NeurIPS*, 2020. [1](#), [3](#), [5](#), [9](#)
- Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. [Fantastic Generalization Measures and Where to Find Them](#). In *ICLR*, 2020. [9](#)
- Diederik P. Kingma and Jimmy Ba. [Adam: A Method for Stochastic Optimization](#). In *ICLR*, 2015. [1](#)
- Alex Krizhevsky. [Learning Multiple Layers of Features from Tiny Images](#). Master’s thesis, University of Toronto, 2009. [9](#)
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. [Gradient-based learning applied to document recognition](#). *Proc. IEEE*, 86(11):2278–2324, 1998. [8](#)
- Kaifeng Lyu and Jian Li. [Gradient Descent Maximizes the Margin of Homogeneous Neural Networks](#). In *ICLR*, 2020. [1](#), [5](#), [9](#)
- Aleksander Mądry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. [Towards Deep Learning Models Resistant to Adversarial Attacks](#). In *ICLR*, 2018. [8](#)
- Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nati Srebro. [Exploring Generalization in Deep Learning](#). In *NeurIPS*, pp. 5947–5956, 2017. [1](#)

- Alethea Power, Yuri Burda, Harrison Edwards, Igor Babuschkin, and Vedant Misra. [Grokking: Generalization Beyond Overfitting on Small Algorithmic Datasets](#). *CoRR*, abs/2201.02177, 2022. 2
- Qian Qian and Xiaoyuan Qian. [The Implicit Bias of AdaGrad on Separable Data](#). In *NeurIPS*, pp. 7759–7767, 2019. 1
- Sebastian Ruder. [An overview of gradient descent optimization algorithms](#). *CoRR*, abs/1609.04747, 2016. 2
- Karen Simonyan and Andrew Zisserman. [Very Deep Convolutional Networks for Large-Scale Image Recognition](#). In *ICLR*, 2015. 9
- Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. [The Implicit Bias of Gradient Descent on Separable Data](#). *J. Mach. Learn. Res.*, 19(70):1–57, 2018. 1
- Davoud Ataee Tarzanagh, Yingcong Li, Xuechen Zhang, and Samet Oymak. [Max-Margin Token Selection in Attention Mechanism](#). In *NeurIPS*, 2023. 2
- Vimal Thilak, Etai Littwin, Shuangfei Zhai, Omid Saremi, Roni Paiss, and Joshua M. Susskind. [The Slingshot Effect: A Late-Stage Optimization Anomaly in Adaptive Gradient Methods](#). *Trans. Mach. Learn. Res.*, 2024. 2
- Gal Vardi. [On the Implicit Bias in Deep-Learning Algorithms](#). *Commun. ACM*, 66(6):86–93, 2023. 1
- Bohan Wang, Qi Meng, Wei Chen, and Tie-Yan Liu. [The Implicit Bias for Adaptive Optimization Algorithms on Homogeneous Neural Networks](#). In *ICML*, pp. 10849–10858, 2021. 1, 2, 7, 8, 9
- Bohan Wang, Qi Meng, Huishuai Zhang, Ruoyu Sun, Wei Chen, Zhi-Ming Ma, and Tie-Yan Liu. [Does Momentum Change the Implicit Regularization on Separable Data?](#) In *NeurIPS*, 2022. 2
- Bohan Wang, Jingwen Fu, Huishuai Zhang, Nanning Zheng, and Wei Chen. [Closing the gap between the upper bound and lower bound of Adam’s iteration complexity](#). In *NeurIPS*, 2023. 2
- Alex Wilkie. [Model completeness results for expansions of the ordered field of real numbers by restricted Pfaffian functions and the exponential function](#). *J. Am. Math. Soc.*, 9(4):1051–1094, 1996. 3
- Ashia C. Wilson, Rebecca Roelofs, Mitchell Stern, Nati Srebro, and Benjamin Recht. [The Marginal Value of Adaptive Gradient Methods in Machine Learning](#). In *NeurIPS*, pp. 4148–4158, 2017. 1
- Shuo Xie and Zhiyuan Li. [Implicit Bias of AdamW:  \$\ell\_\infty\$  Norm Constrained Optimization](#). In *ICML*, pp. 54488–54510, 2024. 2
- Matthew D. Zeiler. [ADADELTA: An Adaptive Learning Rate Method](#). *CoRR*, abs/1212.5701, 2012. 2, 3
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. [Understanding deep learning \(still\) requires rethinking generalization](#). *Commun. ACM*, 64(3):107–115, 2021. 1, 4