

On the Implicit Bias of Adam

Matias D. Cattaneo*
Princeton University
cattaneo@princeton.edu

Jason M. Klusowski*
Princeton University
jason.klusowski@princeton.edu

Boris Shigida*
Princeton University
bs1624@princeton.edu

October 6, 2023

Abstract

In previous literature, backward error analysis was used to find ordinary differential equations (ODEs) approximating the gradient descent trajectory. It was found that finite step sizes implicitly regularize solutions because terms appearing in the ODEs penalize the two-norm of the loss gradients. We prove that the existence of similar implicit regularization in RMSProp and Adam depends on their hyperparameters and the training stage, but with a different “norm” involved: the corresponding ODE terms either penalize the (perturbed) one-norm of the loss gradients or, on the contrary, hinder its decrease (the latter case being typical). We also conduct numerical experiments and discuss how the proven facts can influence generalization.

1 Introduction

Gradient descent (GD) can be seen as a numerical method solving the ordinary differential equation (ODE) $\dot{\boldsymbol{\theta}} = -\nabla E(\boldsymbol{\theta})$, where $E(\cdot)$ is the loss function and $\nabla E(\boldsymbol{\theta})$ denotes its gradient. Starting at $\boldsymbol{\theta}^{(0)}$, it creates a sequence of guesses $\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}, \dots$, which lie close to the solution trajectory $\boldsymbol{\theta}(t)$ governed by aforementioned ODE. Since the step size h is finite, one could search for a modified differential equation $\dot{\tilde{\boldsymbol{\theta}}} = -\nabla \tilde{E}(\tilde{\boldsymbol{\theta}})$ such that $\boldsymbol{\theta}^{(n)} - \tilde{\boldsymbol{\theta}}(nh)$ is exactly zero, or at least closer to zero than $\boldsymbol{\theta}^{(n)} - \boldsymbol{\theta}(nh)$, that is, all the guesses of the descent lie exactly on the new solution curve or closer compared to the original curve. This approach to analysing properties of a numerical method is called backward error analysis in the numerical integration literature (see Chapter IX in [Ernst Hairer & Wanner \(2006\)](#)).

[Barrett & Dherin \(2021\)](#) first used this idea for full-batch gradient descent and found that the modified loss function $\tilde{E}(\tilde{\boldsymbol{\theta}}) = E(\tilde{\boldsymbol{\theta}}) + (h/4)\|\nabla E(\tilde{\boldsymbol{\theta}})\|^2$ makes the trajectory of the solution to $\dot{\tilde{\boldsymbol{\theta}}} = -\nabla \tilde{E}(\tilde{\boldsymbol{\theta}})$ approximate the sequence $\{\boldsymbol{\theta}^{(n)}\}_{n=0}^{\infty}$ one order of h better than the original differential equation, where $\|\cdot\|$ denotes the Euclidean norm. In related work, [Miyagawa \(2022\)](#) obtained the correction term for full-batch gradient descent up to any chosen order, also studying the global error (uniform in the iteration number) as opposed to the local (one-step) error.

The analysis was later extended to mini-batch gradient descent in [Smith et al. \(2021\)](#). Assume that the training set is split in batches of size B and there are m batches per epoch (so the training set size is mB), the cost function is rewritten $E(\boldsymbol{\theta}) = (1/m)\sum_{k=0}^{m-1} \hat{E}_k(\boldsymbol{\theta})$ with mini-batch costs denoted $\hat{E}_k(\boldsymbol{\theta}) = (1/B)\sum_{j=kB+1}^{kB+B} E_j(\boldsymbol{\theta})$. It was obtained in that work that after one epoch, the mean iterate of the algorithm, averaged over all possible shuffles of the batch indices, is close to the solution to $\dot{\boldsymbol{\theta}} = -\nabla \tilde{E}_{SGD}(\boldsymbol{\theta})$, where the modified loss is given by $\tilde{E}_{SGD}(\boldsymbol{\theta}) = E(\boldsymbol{\theta}) + h/(4m) \cdot \sum_{k=0}^{m-1} \|\nabla \hat{E}_k(\boldsymbol{\theta})\|^2$.

More recently, [Ghosh et al. \(2023\)](#) studied gradient descent with heavy-ball momentum iteration $\boldsymbol{\theta}^{(n+1)} = \boldsymbol{\theta}^{(n)} - h\nabla E(\boldsymbol{\theta}^{(n)}) + \beta(\boldsymbol{\theta}^{(n)} - \boldsymbol{\theta}^{(n-1)})$, where β is the momentum parameter. In the full-batch

*Equal contribution

setting, they proved that for n large enough it is close to the continuous trajectory of the first-order ODE

$$\dot{\boldsymbol{\theta}} = \frac{1}{1-\beta} \nabla E(\boldsymbol{\theta}) + h \underbrace{\frac{1+\beta}{4(1-\beta)^3} \nabla \|\nabla E(\boldsymbol{\theta})\|^2}_{\text{implicit regularization}}. \quad (1.1)$$

Their main theorem also provides the analysis for the general mini-batch case.

In another recent work, [Zhao et al. \(2022\)](#) introduce a regularization term $\lambda \cdot \|\nabla E(\boldsymbol{\theta})\|$ to the loss function as a way to ensure finding flatter minima, improving generalization. The only difference between their term and the first-order correction coming from backward error analysis (up to a coefficient) is that the norm is not squared and regularization is applied on a per-batch basis.

Using backward error analysis to approximate the discrete dynamics with a modified ODE for adaptive algorithms such as RMSProp ([Tieleman et al., 2012](#)) and Adam ([Kingma & Ba, 2015](#)) (which is an improvement over RMSProp and AdaGrad ([Duchi et al., 2011](#))) is currently missing in the literature. [Barrett & Dherin \(2021\)](#) note that “it would be interesting to use backward error analysis to calculate the modified loss and implicit regularization for other widely used optimizers such as momentum, Adam and RMSprop”. [Smith et al. \(2021\)](#) reiterate that they “anticipate that backward error analysis could also be used to clarify the role of finite learning rates in adaptive optimizers like Adam”. In the same context, [Ghosh et al. \(2023\)](#) agree that “RMSProp ... and Adam ..., albeit being powerful alternatives to SGD with faster convergence rates, are far from well-understood in the aspect of implicit regularization”. In a similar context, in Appendix G to [Miyagawa \(2022\)](#) it is mentioned that “its [Adam’s] counter term and discretization error are open questions”.

This work fills the gap in the literature by conducting backward error analysis for (mini-batch, and full-batch as a special case) Adam and RMSProp. Our main contributions are listed below.

- In [Theorem 3.1](#), we provide a global second-order in h continuous ODE approximation to Adam in the general mini-batch setting. (A similar result for RMSProp is moved to the supplemental appendix). For the full-batch special case, it was shown in prior work [Ma et al. \(2022\)](#) that the continuous-time limit of both these algorithms is a (perturbed by ε) signGD flow $\dot{\boldsymbol{\theta}} = -\nabla E(\boldsymbol{\theta}) / (|\nabla E(\boldsymbol{\theta})| + \varepsilon)$ component-wise, where ε is the numerical stability parameter; we make this more precise by finding an additional “bias” term on the right (linearly depending on h).
- We analyze the full-batch case in more detail. We find that the bias term does something different from penalizing the two-norm of the loss gradient as in the case of gradient descent: it either penalizes the perturbed one-norm of the loss gradient, defined as $\|\mathbf{v}\|_{1,\varepsilon} = \sum_{i=1}^p \sqrt{v_i^2 + \varepsilon}$, or, on the contrary, hinders its decrease (depending on hyperparameters and the training stage). See the summary of our theoretical finding for the full-batch case in [Section 2](#). We also obtain the backward error analysis result for heavy-ball momentum gradient descent ([Ghosh et al., 2023](#)) as a special case ([Example 2.3](#)).
- We provide numerical evidence consistent with our results. In particular, we notice that often penalizing the perturbed one-norm appears to improve generalization, and hindering its decrease hurts it. The typical absence of implicit regularization appearing from backward error analysis in RMSProp and Adam (as opposed to GD) becomes one more previously unidentified possible explanation for poorer generalization of adaptive gradient algorithms compared to other methods.

Related work

Backward error analysis of first-order methods. We provide the history of finding ordinary differential equations approximating different algorithms above in the introduction. Recently, there have been other applications of backward error analysis related to machine learning. [Kunin et al. \(2020\)](#) show that the approximating continuous-time trajectories satisfy conservation laws that are broken in discrete time. [França et al. \(2021\)](#) use backward error analysis while studying how to discretize continuous-time dynamical systems preserving stability and convergence rates. [Rosca et al. \(2021\)](#) find continuous-time approximations of discrete two-player differential games.

Approximating gradient methods by differential equation trajectories. [Ma et al. \(2022\)](#) prove that the trajectories of Adam and RMSProp are close to signGD dynamics, and investigate different training regimes of these algorithms empirically. SGD is approximated by stochastic differential equations and novel adaptive parameter adjustment policies are devised in [Li et al. \(2017\)](#).

Implicit bias of first-order methods. Soudry et al. (2018) prove that GD trained to classify linearly separable data with logistic loss converges to the direction of the max-margin vector (the solution to the hard margin SVM). This result has been extended to different loss functions in Nacson et al. (2019b), to stochastic gradient descent in Nacson et al. (2019c) and more generic optimization methods in Gunasekar et al. (2018a), to the nonseparable case in Ji & Telgarsky (2018b), Ji & Telgarsky (2019). This line of research has been generalized to studying implicit biases of linear networks (Ji & Telgarsky, 2018a; Gunasekar et al., 2018b), homogeneous neural networks (Ji & Telgarsky, 2020; Nacson et al., 2019a; Lyu & Li, 2019). Woodworth et al. (2020) study the gradient flow of a diagonal linear network with squared loss and show that large initializations lead to minimum 2-norm solutions while small initializations lead to minimum 1-norm solutions. Even et al. (2023) extend this work to the case of non-zero step sizes and mini-batch training. Wang et al. (2021) prove that Adam and RMSProp maximize the margin of homogeneous neural networks.

Generalization of adaptive methods. Cohen et al. (2022) empirically investigate the edge-of-stability regime of adaptive gradient algorithms and the effect of sharpness (defined as the largest eigenvalue of the hessian) on generalization; Granzio (2020); Chen et al. (2021) observe that adaptive methods find sharper minima than SGD and Zhou et al. (2020); Xie et al. (2022) argue theoretically that it is the case. Jiang et al. (2022) introduce a statistic that measures the uniformity of the hessian diagonal and argue that adaptive gradient algorithms are biased towards making this statistic smaller. Keskar & Socher (2017) propose to improve generalization of adaptive methods by switching to SGD in the middle of training.

Notation

We denote the loss of the k th minibatch as a function of the network parameters $\theta \in \mathbb{R}^p$ by $E_k(\theta)$, and in the full-batch setting we omit the index and write $E(\theta)$. ∇E means the gradient of E , and ∇ with indices denotes partial derivatives, e.g. $\nabla_{ijs} E$ is a shortcut for $\frac{\partial^3 E}{\partial \theta_i \partial \theta_j \partial \theta_s}$. The norm without indices $\|\cdot\|$ is the two-norm of a vector, $\|\cdot\|_1$ is the one-norm and $\|\cdot\|_{1,\varepsilon}$ is the perturbed one-norm defined as $\|\mathbf{v}\|_{1,\varepsilon} = \sum_{i=1}^p \sqrt{v_i^2 + \varepsilon}$. (Of course, if $\varepsilon > 0$ the perturbed one-norm is not a norm, but $\varepsilon = 0$ makes it the one-norm.)

2 Implicit bias of full-batch Adam: an informal summary

To avoid ambiguity and to provide the names and notations for hyperparameters, we define the algorithm below.

Definition 2.1. The *Adam* algorithm is an optimization algorithm with numerical stability hyperparameter $\varepsilon > 0$, squared gradient momentum hyperparameter $\rho \in (0, 1)$, gradient momentum hyperparameter $\beta \in (0, 1)$, initialization $\theta^{(0)} \in \mathbb{R}^p$, $\nu^{(0)} = \mathbf{0} \in \mathbb{R}^p$, $\mathbf{m}^{(0)} = \mathbf{0} \in \mathbb{R}^p$ and the following update rule: for each $n \geq 0$, $j \in \{1, \dots, p\}$

$$\begin{aligned} \nu_j^{(n+1)} &= \rho \nu_j^{(n)} + (1 - \rho) (\nabla_j E_n(\theta^{(n)}))^2, & m_j^{(n+1)} &= \beta m_j^{(n)} + (1 - \beta) \nabla_j E_n(\theta^{(n)}), \\ \theta_j^{(n+1)} &= \theta_j^{(n)} - h \frac{m_j^{(n+1)} / (1 - \beta^{n+1})}{\sqrt{\nu_j^{(n+1)} / (1 - \rho^{n+1}) + \varepsilon}}. \end{aligned} \tag{2.1}$$

Remark 2.2 (The ε hyperparameter is inside the square root). Note that the numerical stability hyperparameter $\varepsilon > 0$, which is introduced in these algorithms to avoid division by zero, is inside the square root in our definition. This way we avoid division by zero in the derivative too: the first derivative of $x \mapsto (\sqrt{x + \varepsilon})^{-1}$ is bounded for $x \geq 0$. This is useful for our analysis. In Theorems SA-2.4 and SA-4.4 in the appendix, the original versions of RMSProp and Adam are also tackled, though with an additional assumption which requires that no component of the gradient can come very close to zero in the region of interest. This is true only for the initial period of learning (whereas Theorem 3.1 tackles the whole period). Practitioners do not seem to make a distinction between the version with ε inside vs. outside the square root: tutorials with both versions abound on machine learning related websites. Moreover, the popular

Tensorflow variant of RMSProp has ε inside the square root¹ even though in the documentation² Kingma & Ba (2015) is cited, where ε is outside. While conducting experiments, we also noted that moving ε inside or outside the square root does not change the behavior of Adam or RMSProp qualitatively.

Summary of our main result (in the full-batch case)

We are ready to informally describe our theoretical result (in the full-batch special case). Assume $E(\boldsymbol{\theta})$ is the loss, whose partial derivatives up to the fourth order are bounded. Let $\{\boldsymbol{\theta}^{(n)}\}$ be iterations of Adam as defined in Definition 2.1. Our main result for this case is finding an ODE whose solution trajectory $\tilde{\boldsymbol{\theta}}(t)$ is h^2 -close to $\{\boldsymbol{\theta}^{(n)}\}$, meaning that for any positive time horizon $T > 0$ there exists a constant $C > 0$ such that for any step size $h \in (0, T)$ we have $\|\tilde{\boldsymbol{\theta}}(nh) - \boldsymbol{\theta}^{(n)}\| \leq Ch^2$ (for n between 0 and $\lfloor T/h \rfloor$). The ODE is written the following way (up to terms that rapidly go to zero as n grows): for the component number $j \in \{1, \dots, p\}$

$$\dot{\tilde{\theta}}_j(t) = -\frac{1}{\sqrt{|\nabla_j E(\tilde{\boldsymbol{\theta}}(t))|^2 + \varepsilon}} (\nabla_j E(\tilde{\boldsymbol{\theta}}(t)) + \text{bias}) \quad (2.2)$$

with initial conditions $\tilde{\boldsymbol{\theta}}_j(0) = \boldsymbol{\theta}_j^{(0)}$ for all j , where the bias term is

$$\text{bias} := \frac{h}{2} \left\{ \frac{1+\beta}{1-\beta} - \frac{1+\rho}{1-\rho} + \frac{1+\rho}{1-\rho} \cdot \frac{\varepsilon}{|\nabla_j E(\tilde{\boldsymbol{\theta}}(t))|^2 + \varepsilon} \right\} \nabla_j \|\nabla E(\tilde{\boldsymbol{\theta}}(t))\|_{1,\varepsilon}. \quad (2.3)$$

Depending on hyperparameter values and the training stage, the bias term can take two extreme forms, and during most of the training the reality is usually in between. The extreme cases are as follows.

- If $\sqrt{\varepsilon}$ is **small** compared to all components of $\nabla E(\tilde{\boldsymbol{\theta}}(t))$, i. e. $\min_j |\nabla_j E(\tilde{\boldsymbol{\theta}}(t))| \gg \sqrt{\varepsilon}$, which is the case during the initial learning stage, then

$$\text{bias} = \frac{h}{2} \left\{ \frac{1+\beta}{1-\beta} - \frac{1+\rho}{1-\rho} \right\} \nabla_j \|\nabla E(\tilde{\boldsymbol{\theta}}(t))\|_{1,\varepsilon}. \quad (2.4)$$

For small ε , the perturbed one-norm is indistinguishable from the usual one-norm, and for $\beta > \rho$ it is penalized (in much the same way as the squared two-norm is implicitly penalized in the case of GD), but for $\rho > \beta$ its decrease is actually hindered by this term (so the bias is opposite to penalization). The ODE in (2.2) can be approximately rewritten as

$$\dot{\tilde{\theta}}_j(t) = -\frac{\nabla_j \tilde{E}(\tilde{\boldsymbol{\theta}}(t))}{|\nabla_j E(\tilde{\boldsymbol{\theta}}(t))|}, \quad \tilde{E}(\boldsymbol{\theta}) = E(\boldsymbol{\theta}) + \frac{h}{2} \left\{ \frac{1+\beta}{1-\beta} - \frac{1+\rho}{1-\rho} \right\} \|\nabla E(\boldsymbol{\theta})\|_1. \quad (2.5)$$

- If $\sqrt{\varepsilon}$ is **large** compared to all gradient components, i. e. $\max_j |\nabla_j E(\tilde{\boldsymbol{\theta}}(t))| \ll \sqrt{\varepsilon}$, which may happen during the later learning stage, the fraction with ε is the numerator in (2.3) approaches one, the dependence on ρ cancels out, and

$$\|\nabla E(\tilde{\boldsymbol{\theta}}(t))\|_{1,\varepsilon} \approx \sum_{i=1}^p \sqrt{\varepsilon} \left(1 + \frac{|\nabla_i E(\tilde{\boldsymbol{\theta}}(t))|^2}{2\varepsilon} \right) = p\sqrt{\varepsilon} + \frac{1}{2\sqrt{\varepsilon}} \|\nabla E(\tilde{\boldsymbol{\theta}}(t))\|^2. \quad (2.6)$$

In other words, $\|\cdot\|_{1,\varepsilon}$ becomes $\|\cdot\|^2/(2\sqrt{\varepsilon})$ up to an additive constant (which is “eaten” by the gradient):

$$\text{bias} = \frac{h}{4\sqrt{\varepsilon}} \frac{1+\beta}{1-\beta} \nabla_j \|\nabla E(\tilde{\boldsymbol{\theta}}(t))\|^2.$$

The form of the ODE in this case is

$$\dot{\tilde{\theta}}_j(t) = -\nabla_j \tilde{E}(\tilde{\boldsymbol{\theta}}(t)), \quad \tilde{E}(\boldsymbol{\theta}) = \frac{1}{\sqrt{\varepsilon}} \left(E(\tilde{\boldsymbol{\theta}}(t)) + \frac{h}{4\sqrt{\varepsilon}} \frac{1+\beta}{1-\beta} \|\nabla E(\tilde{\boldsymbol{\theta}}(t))\|^2 \right). \quad (2.7)$$

These two extreme cases are summarized in Table 1. In Figure 1, we use the one-dimensional ($p = 1$) case to illustrate what kind of term is being implicitly penalized.

	ε “small”	ε “large”
$\beta \geq \rho$	$\ \nabla E(\boldsymbol{\theta})\ _1$ -penalized	$\ \nabla E(\boldsymbol{\theta})\ _2^2$ -penalized
$\rho > \beta$	$-\ \nabla E(\boldsymbol{\theta})\ _1$ -penalized	$\ \nabla E(\boldsymbol{\theta})\ _2^2$ -penalized

Table 1: Implicit bias of Adam: special cases. “Small” and “large” are in relation to squared gradient components (Adam in the latter case is close to GD with momentum).

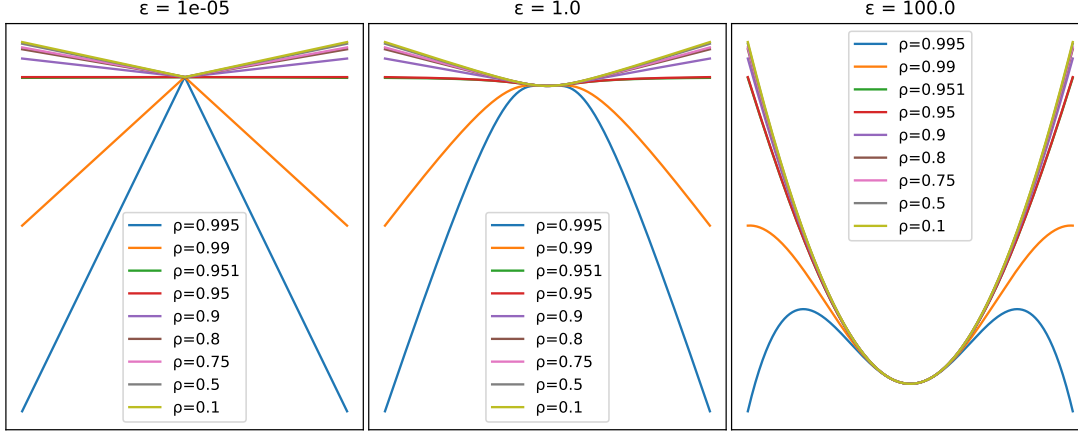


Figure 1: The graphs of $x \mapsto \int_0^x \left\{ \frac{1+\beta}{1-\beta} - \frac{1+\rho}{1-\rho} + \frac{1+\rho}{1-\rho} \cdot \frac{\varepsilon}{y^2+\varepsilon} \right\} d\sqrt{\varepsilon+y^2}$ with $\beta = 0.95$.

This overview also applies to RMSProp by setting $\beta = 0$. See Theorem SA-3.4 in the appendix for the formal result.

Example 2.3 (Backward Error Analysis for GD with Heavy-ball Momentum). Assume ε is very large compared to all squared gradient components during the whole training process, so that the form of the ODE is approximated by (2.7). Since Adam with a large ε and after a certain number of iterations approximates SGD with heavy-ball momentum with step size $h \frac{1-\beta}{\sqrt{\varepsilon}}$, a linear step size change (and corresponding time change) gives exactly the equations in Theorem 4.1 of Ghosh et al. (2023). Taking $\beta = 0$ (no momentum), we get the implicit regularization of GD from Barrett & Dherin (2021).

3 ODE approximating mini-batch Adam trajectories: full statement

We only make one assumption, which is standard in the literature: the loss for each mini-batch is 4 times continuously differentiable, and partial derivatives up to order 4 of each mini-batch loss E_k are bounded by constants, i. e. there exists a positive constant M such that for $\boldsymbol{\theta}$ in the region of interest

$$\sup_k \left\{ \sup_i |\nabla_i E_k(\boldsymbol{\theta})| \vee \sup_{i,j} |\nabla_{ij} E_k(\boldsymbol{\theta})| \vee \sup_{i,j,s} |\nabla_{ijs} E_k(\boldsymbol{\theta})| \vee \sup_{i,j,s,r} |\nabla_{ijsr} E_k(\boldsymbol{\theta})| \right\} \leq M. \quad (3.1)$$

We now state the main result for mini-batch Adam, whose proof is in the supplemental appendix (Theorem SA-5.4).

Theorem 3.1. For any sequence $\{a_k\}$, let $AV_\gamma^n[a.] := \frac{1}{1-\gamma^{n+1}} \sum_{k=0}^n \gamma^{n-k} (1-\gamma) a_k$ denote the exponential averaging operator. Assume (3.1) holds. Let $\{\boldsymbol{\theta}^{(n)}\}$ be iterations of Adam as defined in Definition 2.1,

¹<https://github.com/keras-team/keras/blob/f9336cc5114b4a9429a242deb264b707379646b7/keras/optimizers/rmsprop.py#L190>

²https://www.tensorflow.org/api_docs/python/tf/keras/optimizers/experimental/RMSprop

$\tilde{\boldsymbol{\theta}}(t)$ be the continuous solution to the piecewise ODE

$$\begin{aligned} \dot{\tilde{\boldsymbol{\theta}}}_j(t) &= -\frac{M_j^{(n)}(\tilde{\boldsymbol{\theta}}(t))}{R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t))} \\ &+ h \left(\frac{M_j^{(n)}(\tilde{\boldsymbol{\theta}}(t))(2P_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) + \bar{P}_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)))}{2R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t))^3} - \frac{2L_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) + \bar{L}_j^{(n)}(\tilde{\boldsymbol{\theta}}(t))}{2R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t))} \right). \end{aligned} \quad (3.2)$$

for $t \in [nh, (n+1)h]$ with the initial condition $\tilde{\boldsymbol{\theta}}(0) = \boldsymbol{\theta}^{(0)}$, where

$$\begin{aligned} R_j^{(n)}(\boldsymbol{\theta}) &:= \sqrt{\text{AV}_\rho^n [(\nabla_j E(\boldsymbol{\theta}))^2] + \varepsilon}, & M_j^{(n)}(\boldsymbol{\theta}) &:= \text{AV}_\beta^n [\nabla_j E(\boldsymbol{\theta})], \\ L_j^{(n)}(\boldsymbol{\theta}) &:= \text{AV}_\beta^n \left[\sum_{i=1}^p \nabla_{ij} E(\boldsymbol{\theta}) \sum_{l=1}^{n-1} \frac{M_i^{(l)}(\boldsymbol{\theta})}{R_i^{(l)}(\boldsymbol{\theta})} \right], & \bar{L}_j^{(n)}(\boldsymbol{\theta}) &:= \text{AV}_\beta^n \left[\sum_{i=1}^p \nabla_{ij} E(\boldsymbol{\theta}) \frac{M_i^{(n)}(\boldsymbol{\theta})}{R_i^{(n)}(\boldsymbol{\theta})} \right], \\ P_j^{(n)}(\boldsymbol{\theta}) &:= \text{AV}_\rho^n \left[\nabla_j E(\boldsymbol{\theta}) \sum_{i=1}^p \nabla_{ij} E(\boldsymbol{\theta}) \sum_{l=1}^{n-1} \frac{M_i^{(l)}(\boldsymbol{\theta})}{R_i^{(l)}(\boldsymbol{\theta})} \right], \\ \bar{P}_j^{(n)}(\boldsymbol{\theta}) &:= \text{AV}_\rho^n \left[\nabla_j E(\boldsymbol{\theta}) \sum_{i=1}^p \nabla_{ij} E(\boldsymbol{\theta}) \frac{M_i^{(n)}(\boldsymbol{\theta})}{R_i^{(n)}(\boldsymbol{\theta})} \right]. \end{aligned}$$

Then, for any fixed positive time horizon $T > 0$ there exists a constant C such that for any step size $h \in (0, T)$ we have $\|\boldsymbol{\theta}(nh) - \boldsymbol{\theta}^{(n)}\| \leq Ch^2$ for $n \in \{0, \dots, \lfloor T/h \rfloor\}$.

4 Discussion

First conclusion. Recall that from Ghosh et al. (2023) the ODE approximating the dynamics of full-batch heavy-ball momentum GD is close to (1.1). The correction term regularizes the training process by penalizing the two-norm of the gradient of the loss. We can conclude that *this* kind of regularization is typically absent in RMSProp (if ε is small) and Adam with $\rho > \beta$ (if ε is small). This may partially explain why these algorithms generalize worse than SGD, and it may be a previously unknown perspective on why they are biased towards higher-curvature regions and find “sharper” minima.

Second conclusion. However, the bias term in (2.3) does contain a kind of “norm” which is the perturbed one-norm $\|\mathbf{v}\|_{1,\varepsilon} = \sum_{i=1}^p \sqrt{v_i^2 + \varepsilon}$. If $\sqrt{\varepsilon}$ is small compared to gradient components, which is usually true except at the end of the training, we can conclude from (2.5) that it is only in the case $\beta > \rho$ that the perturbed norm *is* penalized, and decreasing ρ or increasing β moves the trajectory towards regions with lower “norm”.

Third conclusion. There is currently no theory indicating that penalizing the (perturbed) one-norm of the gradient improves generalization. However, reasoning by analogy (with the case of the two-norm), we can conjecture with lower confidence that at least in some stable regimes of training increasing β and decreasing ρ should improve the test error.

5 Illustration: simple bilinear model

We now analyze the effect of the first-order term for Adam in the same model as Barrett & Dherin (2021) and Ghosh et al. (2023) have studied. Namely, assume the parameter $\boldsymbol{\theta} = (\theta_1, \theta_2)$ is 2-dimensional, and the loss is given by $E(\boldsymbol{\theta}) := 1/2(y - \theta_1\theta_2x)^2$, where $x = 2$, $y = 3/2$ are fixed scalars. The loss is minimized on the hyperbola $\theta_1\theta_2 = y/x$. We graph the trajectories of Adam in this case: Figure 2 shows that increasing β forces the trajectory to the region with smaller $\|\nabla E(\boldsymbol{\theta})\|_1$, and increasing ρ does the opposite. Figure 3 shows that increasing the learning rate moves Adam towards the region with smaller $\|\nabla E(\boldsymbol{\theta})\|_1$ if $\beta > \rho$ (just like in the case of gradient descent, except the norm is different if ε is small compared to gradient components), and does the opposite if $\rho > \beta$. All these observations are exactly what Theorem 3.1 predicts.

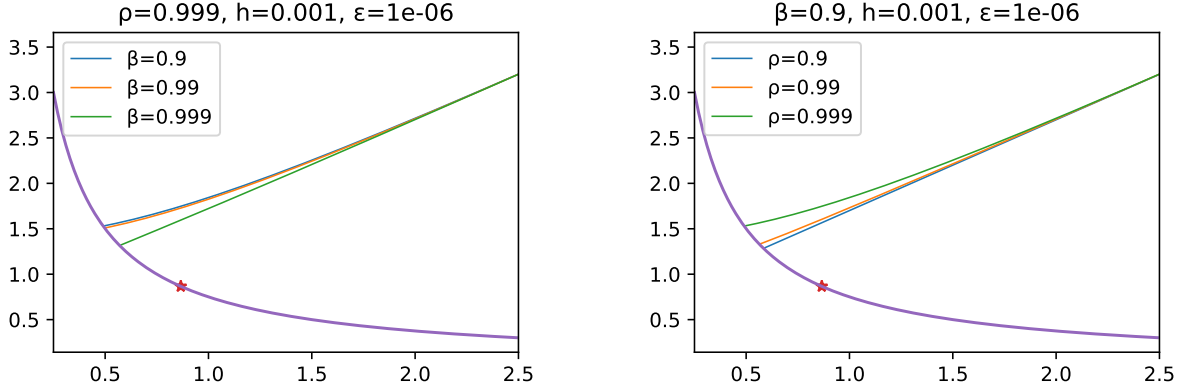


Figure 2: Increasing β moves the trajectory of Adam towards the regions with smaller one-norm of the gradient (if ε is sufficiently small); increasing ρ does the opposite. The violet line is the line of global minima, and the cross denotes the limiting point of minimal one-norm of the gradient. All Adam trajectories start at $(2.8, 3.5)$.

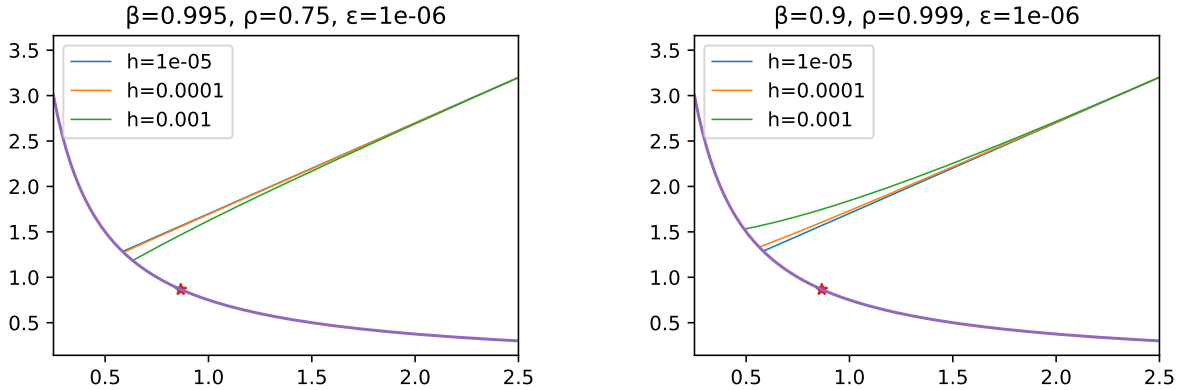


Figure 3: The setting is the same as in Figure 2. Increasing the learning rate moves the Adam trajectory towards the regions with smaller one-norm of the gradient if β is significantly larger than ρ and does the opposite if ρ is larger than β .

6 Numerical experiments

We offer some preliminary empirical evidence of how the bias term shows up in deep neural networks.

Ma et al. (2022) divides training regimes of Adam into three categories: the spike regime when ρ is much larger than β , in which the training loss curve contains very large spikes and the training is obviously unstable; the (stable) oscillation regime when ρ is sufficiently close to β , in which the loss curve contains fast and small oscillations; the divergence regime when β is much larger than ρ , in which Adam diverges. We of course exclude the last regime. Since it is very unlikely that an unstable Adam trajectory is close to the piecewise ODE emerging from backward error analysis, we exclude the spike regime as well, and confine ourselves to considering the oscillation regime (in which ρ and β should not be so far apart that spikes appear). This is the regime Ma et al. (2022) recommend to use in practice.

We train Resnet-50 on the CIFAR-10 dataset with full-batch Adam and investigate how the quantity $\|\nabla E(\theta)\|_{1,\varepsilon}$ and the test error are affected by increasing ρ or β . Figure 4 shows that in the stable oscillation regime increasing ρ seems to increase the perturbed one-norm (consistent with backward error analysis: the smaller ρ , the more this “norm” is penalized) and decrease the test accuracy. The opposite to the latter was noticed in Cohen et al. (2022), which we think is the case in the spike regime (where the trajectory of Adam is definitely far from the piecewise ODE trajectory at later stages of training). Figure 5 shows that increasing β seems to decrease the perturbed one-norm (consistent with backward error analysis:

the larger β , the more this norm is penalized) and increase the test accuracy. The picture confirms the finding in Ghosh et al. (2023) (for momentum gradient descent) that increasing the momentum parameter improves the test accuracy.

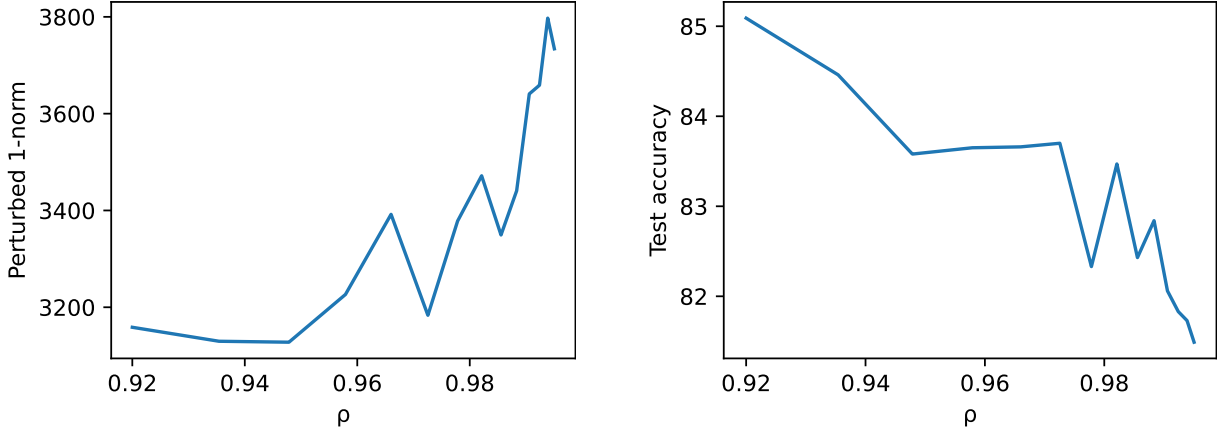


Figure 4: Resnet-50 on CIFAR-10 trained with full-batch Adam. The test accuracy seems to fall as ρ increases (in the stable “small oscillations” regime of training). The hyperparameters are as follows: $h = 7.5 \cdot 10^{-5}$, $\varepsilon = 10^{-8}$, $\beta = 0.99$. The test accuracies plotted here are maximal after more than 3600 epochs. The perturbed norms are calculated at the same epoch number 900. (It is fair to compare Adam with different parameters at one epoch since the effective learning rates are the same.)

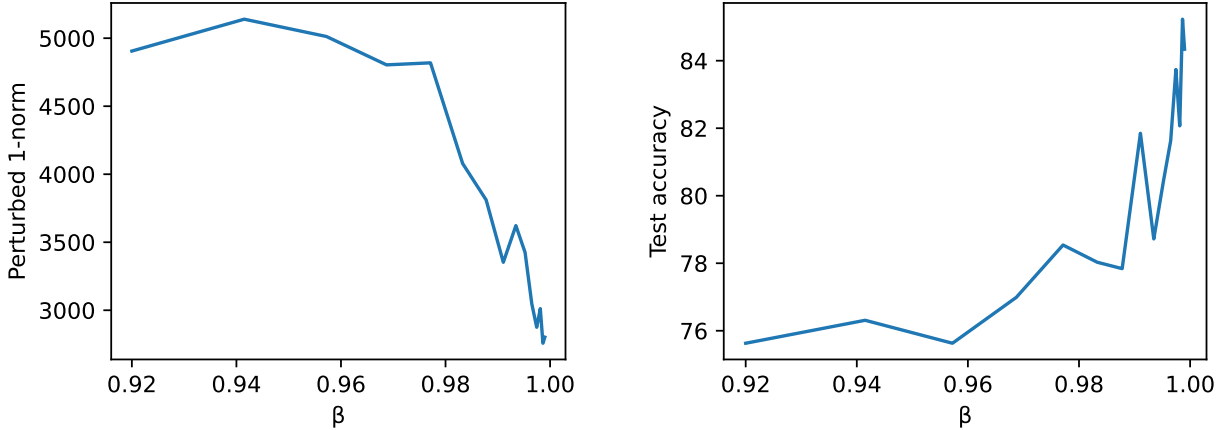


Figure 5: Resnet-50 on CIFAR-10 trained with full-batch Adam. The perturbed one-norm seems to fall as β increases (in the stable oscillation regime of training), and the test accuracy seems to rise. The hyperparameters are as follows: $h = 10^{-4}$, $\rho = 0.999$, $\varepsilon = 10^{-8}$. Both metrics are calculated when the loss first drops below the threshold 0.1.

We obtain a more detailed picture of the perturbed norm’s behavior by training Resnet-101 on CIFAR-10 and CIFAR-100 with full-batch Adam. Figure 6 shows the graphs of $\|\nabla E\|_{1,\varepsilon}$ as functions of the epoch number. The “norm” decreases, then rises again, and then decreases further until it flatlines. Throughout most of the training, the larger β the smaller the “norm”. The “hills” of the “norm” curves are higher with smaller β and larger ρ . This is completely consistent with backward analysis because the larger ρ compared to β , the more $\|\nabla E\|_{1,\varepsilon}$ is prevented from falling by the bias term.

7 Future directions

As far as we know, the assumption similar to (3.1) is explicitly or implicitly present in all previous work on backward error analysis of gradient-based machine learning algorithms. Apart from the technicality

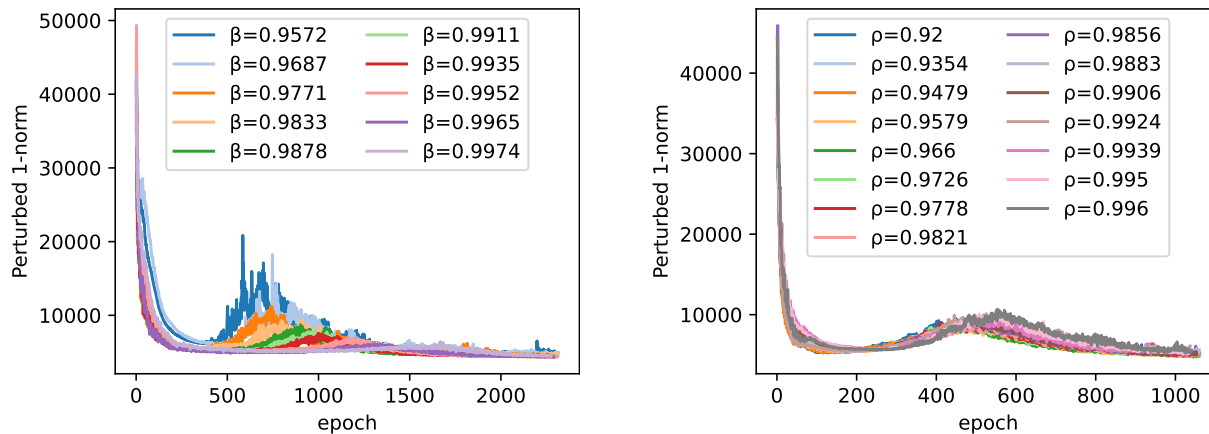


Figure 6: Curves plotting $\|\nabla E\|_{1,\varepsilon}$ after each epoch for a full-batch Adam. On the left: Resnet-101 on CIFAR-10, $h = 10^{-4}$, $\rho = 0.999$, $\varepsilon = 10^{-8}$. On the right: Resnet-101 on CIFAR-100, $h = 10^{-4}$, $\beta = 0.97$, $\varepsilon = 10^{-8}$.

that ReLU activations cause the loss to not be differentiable everywhere (though it is very common to ignore this), there is evidence that large-batch algorithms often operate at the edge of stability (Cohen et al., 2021, 2022), in which the largest eigenvalue of the hessian can be quite large, making it unclear whether the higher-order partial derivatives can safely be assumed bounded near optimality. However, as Smith et al. (2021) point out, in the mini-batch setting backward error analysis can be more accurate. We leave a qualitative analysis of the behavior (average or otherwise) of first-order terms in Theorem 3.1 in the mini-batch case as a future direction.

Also, the constant C in Theorem 3.1 depends on ε and goes to infinity as ε goes to zero. Theoretically, our proof does not exclude the case where for very small ε the trajectory of the piecewise ODE is only close to the Adam trajectory for small, suboptimal learning rates, at least at later stages of learning. (For the initial learning period, this is not a problem.) It appears to also be true of Proposition 1 in Ma et al. (2022) (zeroth-order approximation by sign-GD). This is especially noticeable in the large-spike regime of training (see Section 6 and Ma et al. (2022)) which, despite being obviously pretty unstable, can still minimize the training loss well and lead to acceptable test errors. It would be interesting to investigate this regime in connection with Theorem 3.1 in detail.

We believe these considerations can fruitfully guide future work in this area.

Acknowledgments

We specially thank Boris Hanin and Sam Smith for their insightful comments and suggestions. Cattaneo gratefully acknowledges financial support from the National Science Foundation through DMS-2210561 and SES-2241575. Klusowski gratefully acknowledges financial support from the National Science Foundation through CAREER DMS-2239448, DMS-2054808, and HDR TRIPODS CCF-1934924.

References

- David Barrett and Benoit Dherin. Implicit gradient regularization. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=3q5IqUrkcF>.
- Xiangning Chen, Cho-Jui Hsieh, and Boqing Gong. When vision transformers outperform resnets without pre-training or strong data augmentations. *arXiv preprint arXiv:2106.01548*, 2021.
- Jeremy Cohen, Simran Kaur, Yuanzhi Li, J Zico Kolter, and Ameet Talwalkar. Gradient descent on neural networks typically occurs at the edge of stability. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=jh-rTtvkGeM>.
- Jeremy M Cohen, Behrooz Ghorbani, Shankar Krishnan, Naman Agarwal, Sourabh Medapati, Michal Badura, Daniel Suo, David Cardoze, Zachary Nado, George E Dahl, et al. Adaptive gradient methods at the edge of stability. *arXiv preprint arXiv:2207.14484*, 2022.

- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- Christian Lubich Ernst Hairer and Gerhard Wanner. *Geometric numerical integration*. Springer-Verlag, Berlin, 2 edition, 2006. ISBN 3-540-30663-3.
- Mathieu Even, Scott Pesme, Suriya Gunasekar, and Nicolas Flammarion. (s) gd over diagonal linear networks: Implicit regularisation, large stepsizes and edge of stability. *arXiv preprint arXiv:2302.08982*, 2023.
- Guilherme França, Michael I Jordan, and René Vidal. On dissipative symplectic integration with applications to gradient-based optimization. *Journal of Statistical Mechanics: Theory and Experiment*, 2021 (4):043402, 2021.
- Avrajit Ghosh, He Lyu, Xitong Zhang, and Rongrong Wang. Implicit regularization in heavy-ball momentum accelerated stochastic gradient descent. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=ZzdBhtEH9yB>.
- Diego Granzio. Flatness is a false friend. *arXiv preprint arXiv:2006.09091*, 2020.
- Suriya Gunasekar, Jason Lee, Daniel Soudry, and Nathan Srebro. Characterizing implicit bias in terms of optimization geometry. In *International Conference on Machine Learning*, pp. 1832–1841. PMLR, 2018a.
- Suriya Gunasekar, Jason D Lee, Daniel Soudry, and Nati Srebro. Implicit bias of gradient descent on linear convolutional networks. *Advances in neural information processing systems*, 31, 2018b.
- Ziwei Ji and Matus Telgarsky. Gradient descent aligns the layers of deep linear networks. *arXiv preprint arXiv:1810.02032*, 2018a.
- Ziwei Ji and Matus Telgarsky. Risk and parameter convergence of logistic regression. *arXiv preprint arXiv:1803.07300*, 2018b.
- Ziwei Ji and Matus Telgarsky. The implicit bias of gradient descent on nonseparable data. In *Conference on Learning Theory*, pp. 1772–1798. PMLR, 2019.
- Ziwei Ji and Matus Telgarsky. Directional convergence and alignment in deep learning. *Advances in Neural Information Processing Systems*, 33:17176–17186, 2020.
- Kaiqi Jiang, Dhruv Malik, and Yuanzhi Li. How does adaptive optimization impact local neural network geometry? *arXiv preprint arXiv:2211.02254*, 2022.
- Nitish Shirish Keskar and Richard Socher. Improving generalization performance by switching from adam to sgd. *arXiv preprint arXiv:1712.07628*, 2017.
- Diederick P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.
- Daniel Kunin, Javier Sagastuy-Brena, Surya Ganguli, Daniel LK Yamins, and Hidenori Tanaka. Neural mechanics: Symmetry and broken conservation laws in deep learning dynamics. *arXiv preprint arXiv:2012.04728*, 2020.
- Qianxiao Li, Cheng Tai, and Weinan E. Stochastic modified equations and adaptive stochastic gradient algorithms. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 2101–2110. PMLR, 8 2017. URL <https://proceedings.mlr.press/v70/li17f.html>.
- Kaifeng Lyu and Jian Li. Gradient descent maximizes the margin of homogeneous neural networks. *arXiv preprint arXiv:1906.05890*, 2019.
- Chao Ma, Lei Wu, and E Weinan. A qualitative study of the dynamic behavior for adaptive gradient algorithms. In *Mathematical and Scientific Machine Learning*, pp. 671–692. PMLR, 2022.

- Taiki Miyagawa. Toward equation of motion for deep neural networks: Continuous-time gradient descent and discretization error analysis. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=qq84D17BPu>.
- Mor Shpigel Nacson, Suriya Gunasekar, Jason Lee, Nathan Srebro, and Daniel Soudry. Lexicographic and depth-sensitive margins in homogeneous and non-homogeneous deep models. In *International Conference on Machine Learning*, pp. 4683–4692. PMLR, 2019a.
- Mor Shpigel Nacson, Jason Lee, Suriya Gunasekar, Pedro Henrique Pamplona Savarese, Nathan Srebro, and Daniel Soudry. Convergence of gradient descent on separable data. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 3420–3428. PMLR, 2019b.
- Mor Shpigel Nacson, Nathan Srebro, and Daniel Soudry. Stochastic gradient descent on separable data: Exact convergence with a fixed learning rate. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 3051–3059. PMLR, 2019c.
- Mihaela C Rosca, Yan Wu, Benoit Dherin, and David Barrett. Discretization drift in two-player games. In *International Conference on Machine Learning*, pp. 9064–9074. PMLR, 2021.
- Samuel L Smith, Benoit Dherin, David Barrett, and Soham De. On the origin of implicit regularization in stochastic gradient descent. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=rq_Qr0c1Hyo.
- Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit bias of gradient descent on separable data. *The Journal of Machine Learning Research*, 19(1):2822–2878, 2018.
- Tijmen Tieleman, Geoffrey Hinton, et al. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural networks for machine learning*, 4(2):26–31, 2012.
- Bohan Wang, Qi Meng, Wei Chen, and Tie-Yan Liu. The implicit bias for adaptive optimization algorithms on homogeneous neural networks. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 10849–10858. PMLR, 7 2021. URL <https://proceedings.mlr.press/v139/wang21q.html>.
- Blake Woodworth, Suriya Gunasekar, Jason D Lee, Edward Moroshko, Pedro Savarese, Itay Golan, Daniel Soudry, and Nathan Srebro. Kernel and rich regimes in overparametrized models. In *Conference on Learning Theory*, pp. 3635–3673. PMLR, 2020.
- Zeke Xie, Xinrui Wang, Huishuai Zhang, Issei Sato, and Masashi Sugiyama. Adaptive inertia: Disentangling the effects of adaptive learning rate and momentum. In *International conference on machine learning*, pp. 24430–24459. PMLR, 2022.
- Yang Zhao, Hao Zhang, and Xiuyuan Hu. Penalizing gradient norm for efficiently improving generalization in deep learning. In *International Conference on Machine Learning*, pp. 26982–26992. PMLR, 2022.
- Pan Zhou, Jiashi Feng, Chao Ma, Caiming Xiong, Steven Chu Hong Hoi, et al. Towards theoretically understanding why sgd generalizes better than adam in deep learning. *Advances in Neural Information Processing Systems*, 33:21285–21296, 2020.