

---

# The Implicit Bias of AdaGrad on Separable Data

---

**Qian Qian**

Department of Statistics  
Ohio State University  
Columbus, OH 43210, USA  
qian.216@osu.edu

**Xiaoyuan Qian**

School of Mathematical Sciences  
Dalian University of Technology  
Dalian, Liaoning 116024, China  
xyqian@dlut.edu.cn

## Abstract

We study the implicit bias of AdaGrad on separable linear classification problems. We show that AdaGrad converges to a direction that can be characterized as the solution of a quadratic optimization problem with the same feasible set as the hard SVM problem. We also give a discussion about how different choices of the hyperparameters of AdaGrad might impact this direction. This provides a deeper understanding of why adaptive methods do not seem to have the generalization ability as good as gradient descent does in practice.

## 1 Introduction

In recent years, implicit regularization from various optimization algorithms plays a crucial role in the generalization abilities in training deep neural networks (Salakhutdinov and Srebro [2015], Neyshabur et al. [2015], Keskar et al. [2016], Neyshabur et al. [2017], Zhang et al. [2017]). For example, in underdetermined problems where the number of parameters is larger than the number of training examples, many global optima fail to exhibit good generalization properties, however, a specific optimization algorithm (such as gradient descent) does converge to a particular optimum that generalize well, although no explicit regularization is enforced when training the model. In other words, the optimization technique itself "biases" towards a certain model in an implicit way (Soudry et al. [2018]). This motivates a line of works to investigate the implicit biases of various algorithms (Telgarsky [2013], Soudry et al. [2018], Gunasekar et al. [2017, 2018a,b]).

The choice of algorithms would affect the implicit regularization introduced in the learned models. In underdetermined least squares problems, where the minimizers are finite, we know that gradient descent yields the minimum  $L_2$  norm solution, whereas coordinate descent might give a different solution. Another example is logistic regression with separable data. While gradient descent converges in the direction of the hard margin support vector machine solution (Soudry et al. [2018]), coordinate descent converges to the maximum  $L_1$  margin solution (Telgarsky [2013], Gunasekar et al. [2018a]). Unlike the squared loss, the logistic loss does not admit a finite global minimizer on separable data: the iterates will diverge in order to drive the loss to zero. As a result, instead of characterizing the convergence of the iterates  $\mathbf{w}(t)$ , it is the asymptotic direction of these iterates i.e.,  $\lim_{t \rightarrow \infty} \mathbf{w}(t) / \|\mathbf{w}(t)\|$  that is important and therefore has been characterized (Soudry et al. [2018], Gunasekar et al. [2018b]).

Moreover, it has attracted much attention that different adaptive methods of gradient descent and hyperparameters of an adaptive method exhibit different biases, thus leading to different generalization performance in deep learning (Salakhutdinov and Srebro [2015], Keskar et al. [2016], Wilson et al. [2017], Hoffer et al. [2017]). Among those findings is that the vanilla SGD algorithm demonstrates better generalization than its adaptive variants (Wilson et al. [2017]), such as AdaGrad (Duchi et al. [2010]) and Adam (Kingma and Ba [2015]). Therefore it is important to precisely characterize how different adaptive methods induce difference biases. A natural question to ask is: can we explain this observation by characterizing the implicit bias of AdaGrad, which is a paradigm of adaptive

methods, in a binary classification setting with separable data using logistic regression? And how does the implicit bias depend on the choice of the hyperparameters of this specific algorithm, such as initialization, step sizes, etc?

## 1.1 Our Contribution

In this work we study AdaGrad applied to logistic regression with separable data. Our contribution is three-fold as listed as follows.

- We prove that the directions of AdaGrad iterates, with a constant step size sufficiently small, always converge.
- We formulate the asymptotic direction as the solution of a quadratic optimization problem. This achieves a theoretical characterization of the implicit bias of AdaGrad, which also provides insights about why and how the factors involved, such as certain intrinsic properties of the dataset, the initialization and the learning rate, affect the implicit bias.
- We introduce a novel approach to study the bias of AdaGrad. It is mainly based on a geometric estimation on the directions of the updates, which doesn't depend on any calculation on the convergence rate.

## 1.2 Paper Organization

This paper is organized as follows. In Section 2 we explain our problem setup. The main theory is developed in Section 3, including convergence of the adaptive learning rates of AdaGrad, existence of the asymptotic direction of AdaGrad iterates, and relations between the asymptotic directions of AdaGrad and gradient descent iterates. We conclude our paper in Section 4 with a review of our results and some questions left to future research.

## 2 Problem Setup

Let  $\{(\mathbf{x}_n, y_n) : n = 1, \dots, N\}$  be a training dataset with features  $\mathbf{x}_n \in \mathbb{R}^p$  and labels  $y_n \in \{-1, 1\}$ . To simplify the notation, we redefine  $y_n \mathbf{x}_n$  as  $\mathbf{x}_n$ ,  $n = 1, \dots, N$ , and consider learning the logistic regression model over the empirical loss:

$$\mathcal{L}(\mathbf{w}) = \sum_{n=1}^N l(\mathbf{w}^T \mathbf{x}_n), \quad \mathbf{w} \in \mathbb{R}^p,$$

where  $l : \mathbb{R}^p \rightarrow \mathbb{R}$ . We focus on the following case, same as proposed in Soudry et al. [2018]:

**Assumption 1.** There exists a vector  $\mathbf{w}_*$  such that  $\mathbf{w}_*^T \mathbf{x}_n > 0$  for all  $n$ .

**Assumption 2.**  $l(u)$  is continuously differentiable,  $\beta$ -smooth, and strictly decreasing to zero.

**Assumption 3.** There exist positive constants  $a, b, c$ , and  $d$  such that

$$|l'(u) + ce^{-au}| \leq e^{-(a+b)u}, \quad \text{for } u > d.$$

It is easy to see that the exponential loss  $l(u) = e^{-u}$  and the logistic loss  $l(u) = \log(1 + e^{-u})$  both meet these assumptions.

Given two hyperparameters  $\epsilon, \eta > 0$  and an initial point  $\mathbf{w}(0) \in \mathbb{R}^p$ , we consider the diagonal AdaGrad iterates

$$\mathbf{w}(t+1) = \mathbf{w}(t) - \eta \mathbf{h}(t) \odot \mathbf{g}(t), \quad t = 0, 1, 2, \dots, \quad (1)$$

where

$$\begin{aligned} \mathbf{g}(t) &= (g_1(t), \dots, g_p(t)), \\ g_i(t) &= \frac{\partial \mathcal{L}}{\partial w_i}(\mathbf{w}(t)), \\ \mathbf{h}(t) &= (h_1(t), \dots, h_p(t)), \\ h_i(t) &= \frac{1}{\sqrt{g_i(0)^2 + \dots + g_i(t)^2 + \epsilon}}, \quad i = 1, \dots, p, \end{aligned}$$

and  $\odot$  is the element-wise multiplication of two vectors, e.g.

$$\mathbf{a} \odot \mathbf{b} = (a_1 b_1, \dots, a_p b_p)^T$$

for  $\mathbf{a} = (a_1, \dots, a_p)^T$ ,  $\mathbf{b} = (b_1, \dots, b_p)^T$ .

To analyze the convergence of the algorithm, we put an additional restriction on the hyperparameter  $\eta$ .

**Assumption 4.** The hyperparameter  $\eta$  is not too large; specifically,

$$\eta < \frac{2 \min_{i \in \{1, \dots, p\}} \sqrt{g_i(0)^2 + \epsilon}}{\beta}. \quad (2)$$

We are interested in the asymptotic behavior of the AdaGrad iteration scheme in (1). The main problem is: does there exist some vector  $\mathbf{w}_A$  such that

$$\lim_{t \rightarrow \infty} \mathbf{w}(t) / \|\mathbf{w}(t)\| = \mathbf{w}_A?$$

We will provide an affirmative answer to this question in the following section.

### 3 The Asymptotic Direction of AdaGrad Iterates

#### 3.1 Convergence of the Adaptive Learning Rates

We first provide some elementary facts about AdaGrad iterates (1) with all assumptions (1-4) proposed in Section 2.

**Lemma 3.1.**  $\mathcal{L}(\mathbf{w}(t+1)) < \mathcal{L}(\mathbf{w}(t))$  ( $t = 0, 1, \dots$ ).

**Lemma 3.2.**  $\sum_{t=0}^{\infty} \|\mathbf{g}(t)\|^2 < \infty$ .

We notice that Gunasekar et al. [2018a] showed a similar result (Lemma 6, in Section 3.3 of their work) for exponential loss only, under slightly different assumptions. However, their approach depends on some specific properties of the exponential function, and thus cannot be extended to Lemma 3.2 in a trivial manner.

**Lemma 3.3.** The following statements hold:

- (i)  $\|\mathbf{g}(t)\| \rightarrow 0$  ( $t \rightarrow \infty$ ).
- (ii)  $\|\mathbf{w}(t)\| \rightarrow \infty$  ( $t \rightarrow \infty$ ).
- (iii)  $\mathcal{L}(\mathbf{w}(t)) \rightarrow 0$  ( $t \rightarrow \infty$ ).
- (iv)  $\forall n$ ,  $\lim_{t \rightarrow \infty} \mathbf{w}(t)^T \mathbf{x}_n = \infty$ .
- (v)  $\exists t_0$ ,  $\forall t > t_0$ ,  $\mathbf{w}(t)^T \mathbf{x}_n > 0$ .

**Theorem 3.1.** The sequence  $\{\mathbf{h}(t)\}_{t=0}^{\infty}$  converges as  $t \rightarrow \infty$  to a vector

$$\mathbf{h}_{\infty} = (h_{\infty,1}, \dots, h_{\infty,p})$$

satisfying  $h_{\infty,i} > 0$  ( $i = 1, \dots, p$ ).

#### 3.2 Convergence of the Directions of AdaGrad Iterates

In the remainder of the paper we denote  $\mathbf{h}_{\infty} = \lim_{t \rightarrow \infty} \mathbf{h}(t)$  and  $\boldsymbol{\xi}_n = \mathbf{h}_{\infty}^{1/2} \odot \mathbf{x}_n$  ( $n = 1, \dots, N$ ). Since, by Theorem 3.1, the components of  $\mathbf{h}_{\infty}$  have a positive lower bound, we can define

$$\boldsymbol{\beta}(t) = \mathbf{h}_{\infty}^{-1} \odot \mathbf{h}(t) \quad (t = 0, 1, 2, \dots).$$

Here the squared root and the inverse of vectors are defined element-wise. We call the function

$$\mathcal{L}_{ind} : \mathbb{R}^p \rightarrow \mathbb{R}, \quad \mathcal{L}_{ind}(\mathbf{v}) = \sum_{n=1}^N l(\mathbf{v}^T \boldsymbol{\xi}_n)$$

the *induced loss* with respect to AdaGrad (1). Note that

$$\begin{aligned}\mathcal{L}(\mathbf{w}) &= \sum_{n=1}^N l(\mathbf{w}^T \mathbf{x}_n) = \sum_{n=1}^N l\left(\left(\mathbf{h}_\infty^{-1/2} \odot \mathbf{w}\right)^T \left(\mathbf{h}_\infty^{1/2} \odot \mathbf{x}_n\right)\right) \\ &= \sum_{n=1}^N l\left(\left(\mathbf{h}_\infty^{-1/2} \odot \mathbf{w}\right)^T \boldsymbol{\xi}_n\right) = \mathcal{L}_{ind}\left(\mathbf{h}_\infty^{-1/2} \odot \mathbf{w}\right).\end{aligned}$$

Thus if we set

$$\mathbf{v}(t) = \mathbf{h}_\infty^{-1/2} \odot \mathbf{w}(t) \quad (t = 0, 1, 2, \dots), \quad (3)$$

then  $\mathbf{v}(0) = \mathbf{h}_\infty^{-1/2} \odot \mathbf{w}(0)$ , and

$$\begin{aligned}\mathbf{h}_\infty^{1/2} \odot \mathbf{v}(t+1) &= \mathbf{w}(t+1) = \mathbf{w}(t) - \eta \mathbf{h}(t) \odot \nabla \mathcal{L}(\mathbf{w}(t)) \\ &= \mathbf{h}_\infty^{1/2} \odot \mathbf{v}(t) - \eta \mathbf{h}(t) \odot \mathbf{h}_\infty^{-1/2} \odot \nabla \mathcal{L}\left(\mathbf{h}_\infty^{1/2} \odot \mathbf{v}(t)\right) \\ &= \mathbf{h}_\infty^{1/2} \odot \mathbf{v}(t) - \eta \mathbf{h}(t) \odot \mathbf{h}_\infty^{-1/2} \odot \nabla \mathcal{L}_{ind}\left(\mathbf{h}_\infty^{-1/2} \odot \left(\mathbf{h}_\infty^{1/2} \odot \mathbf{v}(t)\right)\right) \\ &= \mathbf{h}_\infty^{1/2} \odot \left(\mathbf{v}(t) - \eta \boldsymbol{\beta}(t) \odot \nabla \mathcal{L}_{ind}(\mathbf{v}(t))\right),\end{aligned}$$

or

$$\mathbf{v}(t+1) = \mathbf{v}(t) - \eta \boldsymbol{\beta}(t) \odot \nabla \mathcal{L}_{ind}(\mathbf{v}(t)) \quad (t = 0, 1, \dots). \quad (4)$$

We refer to (4) as the *induced form* of AdaGrad (1).

The following result for the induced form is a simple corollary of Lemma 3.3.

**Lemma 3.4.** The following statements hold:

- (i)  $\|\nabla \mathcal{L}_{ind}(t)\| \rightarrow 0 \quad (t \rightarrow \infty)$ .
- (ii)  $\|\mathbf{v}(t)\| \rightarrow \infty \quad (t \rightarrow \infty)$ .
- (iii)  $\mathcal{L}_{ind}(\mathbf{v}(t)) \rightarrow 0 \quad (t \rightarrow \infty)$ .
- (iv)  $\forall n, \lim_{t \rightarrow \infty} \mathbf{v}(t)^T \boldsymbol{\xi}_n = \infty$ .
- (v)  $\exists t_0, \forall t > t_0, \mathbf{v}(t)^T \boldsymbol{\xi}_n > 0$ .

For the induced loss  $\mathcal{L}_{ind}$ , consider GD iterates

$$\mathbf{u}(t+1) = \mathbf{u}(t) - \eta \nabla \mathcal{L}_{ind}(\mathbf{u}(t)) \quad (t = 0, 1, \dots). \quad (5)$$

According to Theorem 3 in Soudry et.al.(2018), we have

$$\lim_{t \rightarrow \infty} \frac{\mathbf{u}(t)}{\|\mathbf{u}(t)\|} = \frac{\hat{\mathbf{u}}}{\|\hat{\mathbf{u}}\|},$$

where

$$\hat{\mathbf{u}} = \arg \min_{\mathbf{u}^T \boldsymbol{\xi}_n \geq 1, \forall n} \|\mathbf{u}\|^2.$$

Noting that  $\boldsymbol{\beta}(t) \rightarrow \mathbf{1} \quad (t \rightarrow \infty)$  we can obtain GD iterates (5) by taking the limit of  $\boldsymbol{\beta}(t)$  in (4). Therefore it is reasonable to expect that these two iterative processes have similar asymptotic behaviors, especially a common limiting direction.

Different from the case of GD method discussed in Soudry et al. [2018], however, it is difficult to obtain an effective estimation about the convergence rate of  $\mathbf{w}(t)$ . Instead, we introduce an orthogonal decomposition approach to obtain the asymptotic direction of the original Adagrad process (1).

In the remainder of the paper, we denote by  $P$  the projection onto the 1-dimensional subspace spanned by  $\hat{\mathbf{u}}$ , and  $Q$  the projection onto the orthogonal complement. Without any loss of generality we may assume  $\|\hat{\mathbf{u}}\| = 1$ . Thus we have the orthogonal decomposition

$$\mathbf{v} = P\mathbf{v} + Q\mathbf{v} \quad (\mathbf{v} \in \mathbb{R}^p),$$

where  $P\mathbf{v} = \|\mathbf{v}\| \hat{\mathbf{u}} = (\mathbf{v}^T \hat{\mathbf{u}}) \hat{\mathbf{u}}$ . Moreover, we denote

$$\boldsymbol{\delta}(t) = -\eta \nabla \mathcal{L}_{ind}(\mathbf{v}(t)), \quad \mathbf{d}(t) = \boldsymbol{\beta}(t) \odot \boldsymbol{\delta}(t). \quad (6)$$

Using this notation we can rewrite the iteration scheme (4) as

$$\mathbf{v}(t+1) = \mathbf{v}(t) + \mathbf{d}(t) \quad (t = 0, 1, \dots).$$

By reformulating (6) as

$$\mathbf{d}(t) = \boldsymbol{\delta}(t) + (\beta(t) - \mathbf{1}) \odot \boldsymbol{\delta}(t),$$

where  $\beta(t) - \mathbf{1} \rightarrow \mathbf{0}$  as  $t \rightarrow \infty$ , we regard  $\boldsymbol{\delta}(t)$  as the decisive part of  $\mathbf{d}(t)$  and acquire properties of  $\mathbf{d}(t)$  through exploring analogues of  $\boldsymbol{\delta}(t)$ .

First, we can show a basic estimation:

$$\boldsymbol{\delta}(t)^T \hat{\mathbf{u}} = \|P\boldsymbol{\delta}(t)\| \geq \frac{\|\boldsymbol{\delta}(t)\|}{\max_n \|\boldsymbol{\xi}_n\|} \quad (t = 0, 1, 2, \dots).$$

The projection properties of  $\boldsymbol{\delta}(t)$  is easily passed on to  $\mathbf{d}(t)$ . In fact, for sufficiently large  $t$ ,

$$\mathbf{d}(t)^T \hat{\mathbf{u}} = \|P\mathbf{d}(t)\| \geq \frac{\|\mathbf{d}(t)\|}{4 \max_n \|\boldsymbol{\xi}_n\|}, \quad (7)$$

Inequality (7) provides a cumulative effect on the projection of  $\mathbf{v}(t)$  as  $t$  increases:

$$\|P\mathbf{v}(t)\| \geq \frac{\|\mathbf{v}(t)\|}{8 \max_n \|\boldsymbol{\xi}_n\|}, \quad \text{for sufficiently large } t.$$

The following lemma reveals a crucial characteristic of the iterative process (4): as  $t$  tends to infinity, the contribution of  $\boldsymbol{\delta}(t)$  to the increment of the deviation from the direction of  $\hat{\mathbf{u}}$ , compared to its contribution to the increment in the direction of  $\hat{\mathbf{u}}$ , becomes more and more insignificant.

**Lemma 3.5.** Given  $\varepsilon > 0$ . Let  $a, b, c$  be positive numbers as defined in Assumption 3 in Section 2. If  $\|Q\mathbf{v}(t)\| > 2N(c+1)(ace\varepsilon)^{-1}$ , then for sufficiently large  $t$ ,

$$Q\mathbf{v}(t)^T \boldsymbol{\delta}(t) < \varepsilon \|Q\mathbf{v}(t)\| \|\boldsymbol{\delta}(t)\|.$$

This property can be translated into a more convenient version for  $\mathbf{d}(t)$ .

**Lemma 3.6.** For any  $\varepsilon > 0$ , there exist  $R > 0$  such that for sufficiently large  $t$  and  $\|Q\mathbf{v}(t)\| \geq R$ ,

$$\|Q\mathbf{v}(t+1)\| - \|Q\mathbf{v}(t)\| \leq \varepsilon \|\mathbf{d}(t)\|.$$

Therefore, over a long period, the cumulative increment of  $\mathbf{v}(t)$  in the direction of  $\hat{\mathbf{u}}$  will overwhelm the deviation from it, yielding the existence of an asymptotic direction for  $\mathbf{v}(t)$ .

**Lemma 3.7.**

$$\lim_{t \rightarrow \infty} \frac{\mathbf{v}(t)}{\|\mathbf{v}(t)\|} = \hat{\mathbf{u}}. \quad (8)$$

By the relation (3) between  $\mathbf{v}(t)$  and  $\mathbf{w}(t)$ , our main result directly follows from (8).

**Theorem 3.2.** AdaGrad iterates (1) has an asymptotic direction:

$$\lim_{t \rightarrow \infty} \frac{\mathbf{w}(t)}{\|\mathbf{w}(t)\|} = \frac{\tilde{\mathbf{w}}}{\|\tilde{\mathbf{w}}\|},$$

where

$$\tilde{\mathbf{w}} = \arg \min_{\mathbf{w}^T \mathbf{x}_n \geq 1, \forall n} \left\| \frac{1}{\sqrt{\mathbf{h}_\infty}} \odot \mathbf{w} \right\|^2. \quad (9)$$

### 3.3 Factors Affecting the Asymptotic Direction

Theorem 3.2 confirms that AdaGrad iterates (1) have an asymptotic direction  $\tilde{\mathbf{w}}/\|\tilde{\mathbf{w}}\|$ , where  $\tilde{\mathbf{w}}$  is the solution to the optimization problem (9). Since the objective function  $\left\| \frac{1}{\sqrt{\mathbf{h}_\infty}} \odot \mathbf{w} \right\|^2$  is determined by the limit vector  $\mathbf{h}_\infty$ , it is easy to see that the asymptotic direction may depend on the choices of the dataset  $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$ , the hyperparameters  $\varepsilon$ ,  $\eta$ , and the initial point  $\mathbf{w}(0)$ . In the following we will discuss this varied dependency in several respects.

### 3.3.1 Difference from the Asymptotic Direction of GD iterates

When the classic gradient descent method is applied to minimize the same loss, it is known (see Theorem 3, Soudry et al. [2018]) that GD iterates

$$\mathbf{w}_G(t+1) = \mathbf{w}_G(t) - \eta \nabla \mathcal{L}(\mathbf{w}_G(t)) \quad (t = 0, 1, 2, \dots), \quad (10)$$

have an asymptotic direction  $\hat{\mathbf{w}}/\|\hat{\mathbf{w}}\|$ , where  $\hat{\mathbf{w}}$  is the solution of the hard max-margin SVM problem

$$\arg \min_{\mathbf{w}^T \mathbf{x}_n \geq 1, \forall n} \|\mathbf{w}\|^2. \quad (11)$$

The two optimization problems (9) and (11) have the same feasible set

$$\{\mathbf{w} \in \mathbb{R}^p : \mathbf{w}^T \mathbf{x}_n \geq 1, \text{ for } n = 1, \dots, N\},$$

but they take on different objective functions. It is natural to expect that their solutions  $\tilde{\mathbf{w}}$  and  $\hat{\mathbf{w}}$  yield different directions, as shown in the following toy example.

**Example 3.1.** Let  $\mathbf{x}_1 = (\cos \theta, \sin \theta)^T$  and  $\mathcal{L}(\mathbf{w}) = e^{-\mathbf{w}^T \mathbf{x}_1}$ . Suppose  $0 < \theta < \pi/2$ . In this setting we simply have  $\hat{\mathbf{w}} = \mathbf{x}_1$ . Selecting  $\mathbf{w}(0) = (a, b)^T$  and  $\epsilon = 0$ , we have

$$-\mathbf{g}(0) = e^{-\mathbf{w}(0)^T \mathbf{x}_1} \mathbf{x}_1 = e^{-a \cos \theta - b \sin \theta} (\cos \theta, \sin \theta)^T,$$

$$\mathbf{h}(0) = (h_1(0), h_2(0))^T = e^{a \cos \theta + b \sin \theta} \left( \frac{1}{\cos \theta}, \frac{1}{\sin \theta} \right)^T.$$

In general we can show there is a sequence of positive numbers  $p(t)$  such that

$$-\mathbf{g}(t) = p(t) (\cos \theta, \sin \theta)^T,$$

and

$$\mathbf{h}_\infty = \lim_{t \rightarrow \infty} \frac{1}{\sqrt{p(0)^2 + p(1)^2 + \dots + p(t)^2}} \left( \frac{1}{\cos \theta}, \frac{1}{\sin \theta} \right)^T = \frac{1}{\rho} \left( \frac{1}{\cos \theta}, \frac{1}{\sin \theta} \right)^T.$$

Now

$$\begin{aligned} \tilde{\mathbf{w}} &= \arg \min_{\mathbf{w}^T \mathbf{x}_1 \geq 1} \left\| \mathbf{h}_\infty^{-1/2} \odot \mathbf{w} \right\|^2 = \arg \min_{\mathbf{w}^T \mathbf{x}_1 \geq 1} \rho (w_1^2 \cos \theta + w_2^2 \sin \theta) \\ &= \arg \min_{\mathbf{w}^T \mathbf{x}_1 \geq 1} (w_1^2 \cos \theta + w_2^2 \sin \theta) = \left( \frac{1}{\cos \theta + \sin \theta}, \frac{1}{\cos \theta + \sin \theta} \right), \end{aligned}$$

and we have  $\tilde{\mathbf{w}}/\|\tilde{\mathbf{w}}\| = (\sqrt{2}/2, \sqrt{2}/2)^T$ . Note that this direction is invariant when  $\theta$  ranges between 0 and  $\pi/2$ , i.e., irrelevant to  $\mathbf{x}_1$ . These two directions coincide only when  $\theta = \pi/4$ .

### 3.3.2 Sensitivity to Small Coordinate System Rotations

If we consider the same setting as in Example 3.1, but taking  $\theta \in (\pi/2, \pi)$ . Then the asymptotic direction  $\tilde{\mathbf{w}}/\|\tilde{\mathbf{w}}\|$  will become  $(-\sqrt{2}/2, \sqrt{2}/2)^T$ . This implies, however, if  $\mathbf{x}_1$  is close to the direction of  $y$ -axis, then a small rotation of the coordinate system may result in a large change of the asymptotic direction reaching a right angle, i.e., in this case the asymptotic direction is highly unstable even for a small perturbation of its  $x$ -coordinate.

### 3.3.3 Effects of the Initialization and Hyperparameter $\eta$

It is reasonable to believe that the asymptotic direction of AdaGrad depends on the initial conditions, including initialization and step size (see Section 3.3, Gunasekar et al. [2018a]). Theorem 3.2 yields a geometric interpretation for this dependency as shown in Figure 1, where the red arrows indicate  $\mathbf{x}_1 = (\cos(3\pi/8), \sin(3\pi/8))$  and  $\mathbf{x}_2 = (\cos(9\pi/20), \sin(9\pi/20))$ , and the cyan arrow indicates the max-margin separator, which points at  $\mathbf{m}$ , the corner of the feasible set  $\{\mathbf{w} \mid \mathbf{w}^T \mathbf{x}_n \geq 1, \forall n = 1, 2\}$  (the yellow shadowed area).

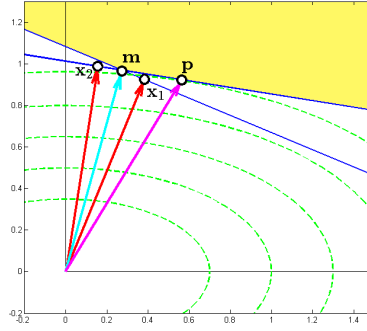


Figure 1: A case that the asymptotic directions of AdaGrad and GD are different.

Since the isolines of the function  $\|\mathbf{h}_\infty^{-1/2} \odot \mathbf{w}\|^2$  are ellipses (the green dashed curves) centered at the origin, the unique minimizer of the function in the feasible set must be the tangency point  $\mathbf{p}$  (pointed at by the magenta arrow) between the tangent ellipse and the boundary of the feasible set. If  $\mathbf{h}_\infty$  varies, then the eccentricity of the tangent ellipses may change. It makes the tangency point move along the boundary, indicating the change of the asymptotic direction.

Numerical simulations also reveal the differences among the asymptotic directions of AdaGrad iterates with various learning rates, as shown in Figure 2. On the left-hand diagram,  $\mathbf{x}_1 = (\cos(\pi/8), \sin(\pi/8))$  and  $\mathbf{x}_2 = (\cos(\pi/20), \sin(\pi/20))$  are two support vectors.  $\mathbf{d}_m$  denotes the direction of the max-margin separator.  $\mathbf{d}_{01}$  and  $\mathbf{d}_{05}$  denote the directions of AdaGrad iterates computed after  $10^8$  steps, with  $\eta = 0.1$  and  $0.5$ , respectively. The small angle between the two may indicate that the asymptotic direction depends on  $\eta$ . However, all the asymptotic directions apparently diverge from the max-margin separator. On the right-hand diagram, the red and blue curves plot  $\|\mathbf{w}(t)/\|\mathbf{w}(t)\| - \mathbf{d}_m\|$  vs. the number of the iterates with  $\eta = 0.1$  and  $0.5$ , respectively. It illustrates that the two sequences of the directions of AdaGrad iterates slowly converge to their own asymptotic directions, slightly different from each other.

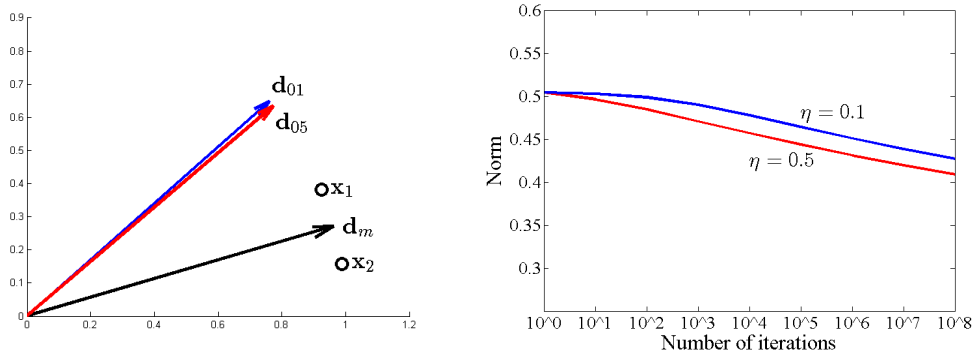


Figure 2: Numerical simulations with  $\eta = 0.1$  and  $0.5$ .

### 3.3.4 Cases that the Asymptotic Direction is Stable

Above we have observed that the asymptotic direction of AdaGrad iterates can be very different from the asymptotic direction of GD iterates, which is robust with respect to different choices of initialization and learning rate  $\eta$ . It is natural to ask what are the conditions under which the two asymptotic directions coincide. The following proposition provides a sufficient one.

**Proposition 3.1.** Let  $\mathbf{a} = (a_1, \dots, a_p)^T$  be a vector satisfying  $\mathbf{a}^T \mathbf{x}_n \geq 1$  ( $n = 1, \dots, N$ ) and  $a_1 \cdots a_p \neq 0$ . Suppose that  $\mathbf{w} = (w_1, \dots, w_p)^T$  satisfies  $\mathbf{w}^T \mathbf{x}_n \geq 1$  ( $n = 1, \dots, N$ ) and

$$a_i (w_i - a_i) \geq 0 \quad (i = 1, \dots, p).$$

Then for any  $\mathbf{b} = (b_1, \dots, b_p)^T$  such that  $b_1 \cdots b_p \neq 0$ ,

$$\arg \min_{\mathbf{w}^T \mathbf{x}_n \geq 1, \forall n} \|\mathbf{b} \odot \mathbf{w}\|^2 = \arg \min_{\mathbf{w}^T \mathbf{x}_n \geq 1, \forall n} \|\mathbf{w}\|^2 = \mathbf{a},$$

and therefore the asymptotic directions of AdaGrad (1) and GD (10) are equal.

Such a condition seems at first sight quite harsh to be satisfied. However, there is a significant proportion of the chances that a dataset  $\{\mathbf{x}_n : n = 1, \dots, N\}$  meets the requirement, as shown in the following result.

**Proposition 3.2.** Suppose  $N \geq p$  and  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{p \times N}$  is sampled from any distribution whose density function is nonzero almost everywhere. Then with a positive probability the asymptotic directions of AdaGrad (1) and GD (10) are equal.

**Example 3.2.** Let  $r_1, r_2 > 0$ ,

$$\mathbf{x}_1 = r_1 (\cos \theta_1, \sin \theta_1)^T, \quad \frac{\pi}{2} \leq \theta_1 < \pi,$$

$$\mathbf{x}_2 = r_2 (\cos \theta_2, \sin \theta_2)^T, \quad \theta_1 - \pi < \theta_2 \leq 0,$$

and  $\mathcal{L}(\mathbf{w}) = l(\mathbf{w}^T \mathbf{x}_1) + l(\mathbf{w}^T \mathbf{x}_2)$ . The system of equations

$$\mathbf{w}^T \mathbf{x}_i = 1 \quad (i = 1, 2)$$

has a unique solution  $(\alpha, \beta)^T$ , where

$$\alpha = \frac{r_2^{-1} \sin \theta_1 - r_1^{-1} \sin \theta_2}{\sin(\theta_1 - \theta_2)} > 0, \quad \beta = \frac{r_1^{-1} \cos \theta_2 - r_2^{-1} \cos \theta_1}{\sin(\theta_1 - \theta_2)} > 0.$$

It is easy to check that if  $\mathbf{w} = (w_1, w_2)^T$  satisfies  $\mathbf{w}^T \mathbf{x}_i \geq 1$  ( $i = 1, 2$ ), then  $w_1 \geq \alpha$ ,  $w_2 \geq \beta$ . Thus any quadratic form  $b_1 w_1^2 + b_2 w_2^2$  ( $b_1, b_2 > 0$ ) takes its minimum at  $(\alpha, \beta)^T$  over the feasible set  $\{\mathbf{w} : \mathbf{w}^T \mathbf{x}_i \geq 1 \text{ (} i = 1, 2)\}$ . Hence the asymptotic direction of AdaGrad (1) applying to this problem is always equal to  $(\alpha, \beta)^T / \|(\alpha, \beta)\|$ , which is also the asymptotic direction of GD (10).

A geometric perspective of this example is given in Figure 2, where the red arrows indicate  $\mathbf{x}_1 = (\cos(5\pi/8), \sin(5\pi/8))$  and  $\mathbf{x}_2 = (\cos(-\pi/8), \sin(-\pi/8))$ , and the magenta arrow indicates  $(\alpha, \beta)^T$ . It is easy to see that the isoline (the thick ellipse drawn in green) along which the function  $\|\mathbf{h}_\infty^{-1/2} \odot \mathbf{w}\|^2$  equals its minimum must intersect with the feasible set (the grey shaded area) at the corner  $(\alpha, \beta)^T$ , no matter what  $\mathbf{h}_\infty$  is.

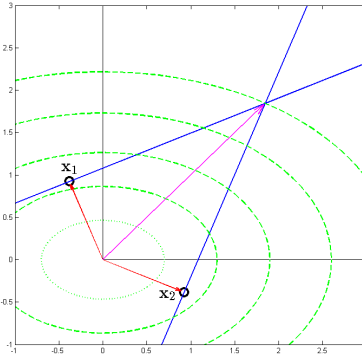


Figure 3: A case that the asymptotic directions of AdaGrad and GD are equal.



## 4 Conclusion

We proved that the basic diagonal AdaGrad, when minimizing a smooth monotone loss function with an exponential tail, has an asymptotic direction, which can be characterized as the solution of a quadratic optimization problem. In this respect AdaGrad is similar to GD, even though their asymptotic directions are usually different. The difference between them also lies in the stability of their asymptotic directions. The asymptotic direction of GD is uniquely determined by the predictors  $x_n$ 's and independent of initialization and the learning rate, as well as the rotation of coordinate system, while the asymptotic direction of AdaGrad is likely to be affected by those factors.

In spite of all these findings, we still do not know whether the asymptotic direction of AdaGrad will change for various initialization or different learning rates. Furthermore, we hope our approach can be applied to the research on the implicit biases of other adaptive methods such as AdaDelta, RMSProp, and Adam.

## References

- B. Neyshabur and R. R. Salakhutdinov and N. Srebro. Path-sgd: Path-normalized optimization in deep neural networks. *In Advances in Neural Information Processing Systems*, page 2422–2430, 2015.
- B. Neyshabur, R. Tomioka, and N. Srebro. In search of the real inductive bias: On the role of implicit regularization in deep learning. *In International Conference on Learning Representations*, 2015.
- Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. *ICLR*, 2016.
- B. Neyshabur, R. Tomioka, R. Salakhutdinov, and N. Srebro. Geometry of optimization and implicit regularization in deep learning, 2017. URL <https://arxiv.org/pdf/1705.03071.pdf>.
- C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals. Understanding deep learning requires rethinking generalization. *In International Conference on Learning Representations*, 2017.
- D. Soudry, E. Hoffer, M. S. Nacson, S. Gunasekar, and N. Srebro. The implicit bias of gradient descent on separable data, 2018.
- M. Telgarsky. Margins, shrinkage and boosting. *Proceedings of the 30th International Conference on Machine Learning, PMLR*, 28(2):307–315, 2013.
- Suriya Gunasekar, Blake E. Woodworth, Srinadh Bhojanapalli, Behnam Neyshabur, and Nati Srebro. Implicit regularization in matrix factorization, 2017.
- Suriya Gunasekar, Jason Lee, Daniel Soudry, and Nathan Srebro. Characterizing implicit bias in terms of optimization geometry. *In Proceedings of the 35th International Conference on Machine Learning*, 2018a.
- S. Gunasekar, J. Lee, D. Soudry, and N. Srebro. Implicit bias of gradient descent on linear convolutional networks. *In Proceedings of the 35th International Conference on Machine Learning*, 2018b.
- Ashia C. Wilson, Rebecca Roelofs, Mitchell Stern, Nathan Srebro, and Benjamin Recht. The marginal value of adaptive gradient methods in machine learning. *arXiv*, pages 1–14, 2017.
- E. Hoffer, I. Hubara, and D. Soudry. Train longer, generalize better: closing the generalization gap in large batch training of neural networks. *In Advances in Neural Information Processing Systems*, page 1–13, 2017.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12:2121 – 2159, 2010.
- Diederik P. Kingma and Jimmy Lei Ba. Adam: a method for stochastic optimization. *International Conference on Learning Representations*, pages 1–13, 2015.