

# The Implicit Bias of AdaGrad on Separable Data

Qian Qian, Xiaoyuan Qian

November 30, 2019

# Our contribution

- 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.

# 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  $y_n \in \{-1, 1\}$ .

Consider learning the logistic regression model over the empirical loss:

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

We focus on the following case, same as proposed in ?:

**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.

# Problem Setup

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

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

# Convergence of the Adaptive Learning Rates

## Theorem

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

# Main Results

**Theorem.** AdaGrad iterates has an asymptotic direction:

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

where

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

# Difference from the Asymptotic Direction of GD iterates

## Example.

Let  $\mathbf{x}_1 = (\cos \theta, \sin \theta)^T$  and  $\mathcal{L}(\mathbf{w}) = e^{-\mathbf{w}^T \mathbf{x}_1}$ . Suppose  $0 < \theta < \pi/2$ . Then  $\hat{\mathbf{w}} = \mathbf{x}_1$ . Selecting  $\mathbf{w}(0) = (a, b)^T$  and  $\epsilon = 0$ , then

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

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

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

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

# Difference from the Asymptotic Direction of GD iterates

**Example (continued).** Now

$$\tilde{\mathbf{w}} = \arg \min_{\mathbf{w}^T \mathbf{x}_1 \geq 1} (w_1^2 \cos \theta + w_2^2 \sin \theta) = \left( \frac{1}{\cos \theta + \sin \theta}, \frac{1}{\cos \theta + \sin \theta} \right),$$

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  $\theta = \pi/4$ .



# Sensitivity to Small Coordinate System Rotations

- If we consider the same setting as in the previous example, 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, 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.

## Cases that the Asymptotic Direction is Stable

**Proposition.** 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 and GD are equal.

# Cases that the Asymptotic Direction is Stable

**Proposition.** 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 and GD are equal.