

Qian Qian, Xiaoyuan Qian

October 22, 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, β -smooth, and strictly decreasing to zero.

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

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

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

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