

Implicit Bias of Gradient Descent Toward Minimum-Norm Solution

Hossein Mobahi
Google Research
Mountain View, CA, USA
hmobahi@google.com

September 23, 2021

1 Introduction

You might have heard about implicit regularization of SGD; meaning preferences toward certain solutions that are not explicitly stated in the objective function. One of them is that SGD prefers solutions with small norm. In fact, such minimum norm bias is not a product of stochasticity of SGD, but the gradient descent itself. Where does this bias come from?

Here I analyze a simple scenario: linear regression with squared loss. We know gradient descent has min-norm bias in this case and we can simply prove that. I first derive the min-norm solution for this problem, and then show gradient descent converges to that.

2 Minimum-Norm Solution

Consider training data (\mathbf{x}_k, y_k) for $k = 1, \dots, n$, where each $\mathbf{x}_k \in \mathbb{R}^d$ and $y_k \in \mathbb{R}$. Suppose we are in over-parameterized regime, i.e. $d > n$. The goal is to find a \mathbf{w} that interpolates the data,

$$\mathbf{x}_k^T \mathbf{w} = y_k \quad \text{for } k = 1, \dots, n, \quad (1)$$

or equivalently in matrix form,

$$\mathbf{X}^T \mathbf{w} = \mathbf{y}, \quad (2)$$

where $\mathbf{X}_{d \times n} = [\mathbf{x}_1 \mid \dots \mid \mathbf{x}_n]$ and $\mathbf{y}_{n \times 1} = [y_1 \mid \dots \mid y_n]$. Since $d > n$, there are possibly infinitely many \mathbf{w} 's that satisfy the above. To keep the problem well-posed, we seek the \mathbf{w} that has the smallest norm,

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{s.t.} \quad \mathbf{X}^T \mathbf{w} = \mathbf{y}, \quad (3)$$

Using the method of Lagrange multiplier, it is straightforward¹ to show that the solution has the form,

$$\mathbf{w}^* = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{y}. \quad (11)$$

3 Gradient Descent

Suppose we wish to find the \mathbf{w} satisfying the $\mathbf{X}^T \mathbf{w} = \mathbf{y}$ using gradient descent with step size $\eta > 0$ by minimizing the squared loss J ,

$$J(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}^T \mathbf{w} - \mathbf{y}\|^2. \quad (12)$$

We assume the data matrix \mathbf{X} is full-rank, i.e. $\text{rank}(\mathbf{X}) = n$. Since $d > n$, there are infinitely many \mathbf{w} 's that can interpolate the data and achieve zero loss. Note that we are *not explicitly enforcing any regularization* here to prefer one solution versus the other. We now show that the solution of gradient descent converges to the min-norm solution \mathbf{w}^* in (11).

¹Denote the Lagrange multiplier vector by $\boldsymbol{\lambda} \in \mathbb{R}^n$. The Lagrangian can be expressed as,

$$L(\mathbf{w}, \boldsymbol{\lambda}) = \frac{1}{2} \|\mathbf{w}\|^2 + \langle \boldsymbol{\lambda}, \mathbf{X}^T \mathbf{w} - \mathbf{y} \rangle. \quad (4)$$

Since,

$$\nabla_{\mathbf{w}} L(\mathbf{w}, \boldsymbol{\lambda}) = \mathbf{w} + \mathbf{X} \boldsymbol{\lambda} \quad (5)$$

$$\nabla_{\boldsymbol{\lambda}} L(\mathbf{w}, \boldsymbol{\lambda}) = \mathbf{X}^T \mathbf{w} - \mathbf{y}. \quad (6)$$

By zero crossing the above we can obtain the optimality condition,

$$\mathbf{w} = -\mathbf{X} \boldsymbol{\lambda} \quad , \quad \mathbf{X}^T \mathbf{w} = \mathbf{y}. \quad (7)$$

Replacing \mathbf{w} implies,

$$-\mathbf{X}^T \mathbf{X} \boldsymbol{\lambda} = \mathbf{y}. \quad (8)$$

Since the $n \times n$ matrix $\mathbf{X}^T \mathbf{X}$ is full-rank (as \mathbf{X} has rank n and $d > n$), its inverse is *well-defined* and thus,

$$\boldsymbol{\lambda} = -(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{y}. \quad (9)$$

Replacing $\boldsymbol{\lambda}$ yields the optimal \mathbf{w} ,

$$\mathbf{w} = -\mathbf{X} \boldsymbol{\lambda} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{y}. \quad (10)$$

The gradient descent updates have the form,

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla J(\mathbf{w}_t) \quad (13)$$

$$= \mathbf{w}_t - \eta \mathbf{X} \left(\mathbf{X}^T \mathbf{w}_t - \mathbf{y} \right) \quad (14)$$

$$= \left(\mathbf{I} - \eta \mathbf{X} \mathbf{X}^T \right) \mathbf{w}_t + \eta \mathbf{X} \mathbf{y}. \quad (15)$$

From here it is easy to resolve the recurrence and obtain,

$$\mathbf{w}_t = (\mathbf{I} - \eta \mathbf{X} \mathbf{X}^T)^t \mathbf{w}_0 + \eta \left(\sum_{i=0}^{t-1} (\mathbf{I} - \eta \mathbf{X} \mathbf{X}^T)^i \right) \mathbf{X} \mathbf{y}. \quad (16)$$

Since $d > n$ and $\text{rank}(\mathbf{X}) = n$ (full-rank assumption), $(\mathbf{X}^T \mathbf{X})^{-\frac{1}{2}}$ is well-defined. This enables us to have the identity²,

$$(\mathbf{I}_{d \times d} - \eta \mathbf{X} \mathbf{X}^T)^i \mathbf{X} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-\frac{1}{2}} (\mathbf{I}_{n \times n} - \eta \mathbf{X}^T \mathbf{X})^i (\mathbf{X}^T \mathbf{X})^{\frac{1}{2}}. \quad (19)$$

Applying this identity to the weight update yields,

$$\mathbf{w}_t = (\mathbf{I} - \eta \mathbf{X} \mathbf{X}^T)^t \mathbf{w}_0 + \eta \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-\frac{1}{2}} \left(\sum_{i=0}^{t-1} (\mathbf{I} - \eta \mathbf{X}^T \mathbf{X})^i \right) (\mathbf{X}^T \mathbf{X})^{\frac{1}{2}} \mathbf{y}. \quad (20)$$

The series $\sum_{i=0}^{\infty} (\mathbf{I} - \eta \mathbf{X}^T \mathbf{X})^i$ converges when $\|\mathbf{I} - \eta \mathbf{X}^T \mathbf{X}\| < 1$, which is equivalent to $|1 - \eta \lambda| < 1$ where λ is an eigenvalue of the matrix $\mathbf{X}^T \mathbf{X}$. This condition can be satisfied by choosing learning rate small enough: $\eta < \frac{2}{\lambda_{\max}}$ when $\lambda_{\min} \neq 0$ (the latter requirement is the reason we switched from the rank-deficient $d \times d$ matrix $\mathbf{X} \mathbf{X}^T$ to the full rank $n \times n$ matrix $\mathbf{X}^T \mathbf{X}$).

The series $\sum_{i=0}^{\infty} \mathbf{A}^i$ is called *Neumann series*, and when convergent ($\|\mathbf{A}\| < 1$), it converges to $(\mathbf{I} - \mathbf{A})^{-1}$. Thus, under the initialization assumption $\mathbf{w}_0 = \mathbf{0}$, it follows that, when $\eta < \frac{2}{\lambda_{\max}}$,

$$\mathbf{w}_{\infty} = (\mathbf{I} - \eta \mathbf{X} \mathbf{X}^T)^{\infty} \mathbf{w}_0 + \eta \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-\frac{1}{2}} (\eta \mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X})^{\frac{1}{2}} \mathbf{y} \quad (21)$$

$$= \eta \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-\frac{1}{2}} (\eta \mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X})^{\frac{1}{2}} \mathbf{y} \quad (22)$$

$$= \eta \mathbf{X} \frac{1}{\eta} (\mathbf{X}^T \mathbf{X})^{-\frac{1}{2} - 1 + \frac{1}{2}} \mathbf{y} \quad (23)$$

$$= \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{y}, \quad (24)$$

where in (21) we use the fact that $\lim_{i \rightarrow \infty} \|(\mathbf{I} - \eta \mathbf{X} \mathbf{X}^T)^i \mathbf{w}_0\| = 0$ because $\|\mathbf{I} - \eta \mathbf{X}^T \mathbf{X}\| \leq 1$ is bounded and $\mathbf{w}_0 = \mathbf{0}$. We just proved the following theorem.

Theorem 1 Consider linear regression with squared loss $\frac{1}{2} \|\mathbf{X}^T \mathbf{w} - \mathbf{y}\|^2$ in the over-parameterized regime $d > n$. Suppose the data matrix \mathbf{X} is full-rank, i.e. $\text{rank}(\mathbf{X}) = n$. Let λ_{\max} be the largest eigenvalue of the matrix $\mathbf{X}^T \mathbf{X}$. Then for any step satisfying $0 < \eta < \frac{2}{\lambda_{\max}}$ under initialization $\mathbf{w}_0 = \mathbf{0}$, \mathbf{w} converges to the solution of the problem $\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2$ s.t. $\mathbf{X}^T \mathbf{w} = \mathbf{y}$.

²This identity can be easily proved. Consider singular value decomposition of $\mathbf{X} = \mathbf{U} \mathbf{S} \mathbf{V}^T$. It is straightforward to verify that,

$$(\mathbf{I}_{d \times d} - \eta \mathbf{X} \mathbf{X}^T)^i \mathbf{X} = \mathbf{U} (\mathbf{I}_{d \times d} - \eta \mathbf{S} \mathbf{S}^T)^i \mathbf{S} \mathbf{V}^T \quad (17)$$

$$\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-\frac{1}{2}} (\mathbf{I}_{n \times n} - \eta \mathbf{X}^T \mathbf{X})^i (\mathbf{X}^T \mathbf{X})^{\frac{1}{2}} = \mathbf{U} \mathbf{S} (\mathbf{I}_{n \times n} - \eta \mathbf{S}^T \mathbf{S})^i \mathbf{V}^T. \quad (18)$$

It is easy to verify that $(\mathbf{I}_{d \times d} - \eta \mathbf{S} \mathbf{S}^T)^i \mathbf{S} = \mathbf{S} (\mathbf{I}_{n \times n} - \eta \mathbf{S}^T \mathbf{S})^i$ (recall that $\mathbf{S}_{d \times n}$ consists of a $n \times n$ diagonal block, and is zero elsewhere). This implies that the left hand side of (17) and (18) are identical.