

Implicit Bias of Gradient Descent on Linear Convolutional Networks

SeongSik Choi

May, 2021

Seoul National University

Table of Contents

- 1 Introduction
- 2 Multi-layer Linear Networks
- 3 Main Results
- 4 Discussion

Large scale neural networks are highly over-parameterized.

However, specific optimization algorithms take us some special global minima.

Example

Linear regression (under-determined model) - minimum l_2 solution

Linear logistic regression (linearly separable) with gradient descent
- hard margin support vector machine solution

Linear logistic regression (linearly separable) with coordinate descent
- maximum l_1 margin solution

Changing to a different parameterization of the same model class changes implicit bias.

(Ex 1. Fully connected, Convolutional)

(Ex 2. Optimizing w, β)

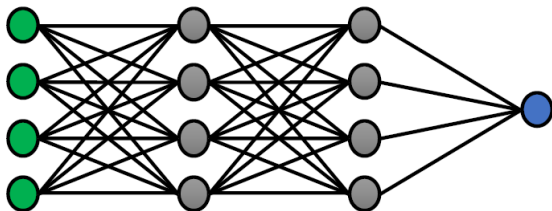
Changing the optimization algorithm changes implicit bias.

(Ex. Gradient descent, Coordinate descent)

Table of Contents

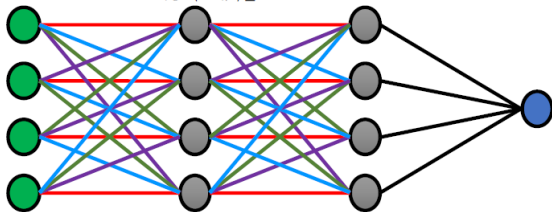
- ① Introduction
- ② Multi-layer Linear Networks**
- ③ Main Results
- ④ Discussion

Multi-layer Linear Networks



(a) Fully connected network of depth L

$$\bar{\beta}^\infty \propto \underset{\forall n, y_n(\mathbf{x}_n, \beta) \geq 1}{\operatorname{argmin}} \|\beta\|_2 \text{ (independent of } L)$$



(b) Convolutional network of depth L

$$\bar{\beta}^\infty \propto \text{first order stationary point of } \underset{\forall n, y_n(\mathbf{x}_n, \beta) \geq 1}{\operatorname{argmin}} \|\hat{\beta}\|_{2/L}$$

Multi-layer Linear Networks

Linear convolutional network

We consider 1-dim circular conv network. Each non-output layer has D units (same as the input dimensionality)

Circular convolutional operation parameterized by full width filters with weights $[w_l \in \mathbb{R}^D]_{l=1}^{L-1}$.

$$h_l[d] = \frac{1}{\sqrt{D}} \sum_{k=0}^{D-1} w_l[k] h_{l-1}[(d+k) \bmod D] := (h_{l-1} \star w_l)[d]$$

The output layer is fully connected.

Multi-layer Linear Networks

A mapping $\mathcal{P} : \mathcal{W} \rightarrow \mathbb{R}^D$ that maps the input parameters $w \in \mathcal{W}$ to a linear predictor in \mathbb{R}^D , such that the output of the network is given by $f_w(x) = \langle x, \mathcal{P}(w) \rangle$.

For fully connected networks, $\mathcal{P}_{\text{full}}(w) = w_1 w_2 \dots w_L$,

For convolutional networks, $\mathcal{P}_{\text{conv}}(w) = \left(\left(\left(w_L^\downarrow \star w_{L-1} \right) \star w_{L-2} \right) \dots \star w_1 \right)^\downarrow$, where w^\downarrow denotes the flipped vector corresponding to w given by $w^\downarrow[k] = w[D - k - 1]$ for $k = 0, 1, \dots, D - 1$.

Multi-layer Linear Networks

Objective for training network

$$\min_{\mathbf{w} \in \mathcal{W}} \mathcal{L}_{\mathcal{P}}(\mathbf{w}) := \sum_{n=1}^N \ell(\langle \mathbf{x}_n, \mathcal{P}(\mathbf{w}) \rangle, y_n) \quad (1)$$

is equivalent to the following optimization over $\beta = \mathcal{P}(\mathbf{w})$

$$\min_{\beta \in \mathbb{R}^D} \mathcal{L}(\beta) := \sum_{n=1}^N \ell(\langle \mathbf{x}_n, \beta \rangle, y_n) \quad (2)$$

But optimizing over w leads to different classifiers compared to optimizing over β directly.

Multi-layer Linear Networks

Consider problem (1) and (2) on a linearly separable dataset using the logistic loss.

Global infimum of $L(\beta)$ is 0, but not attainable by any finite β .

Want to find the direction $\bar{\beta}^\infty = \lim_{t \rightarrow \infty} \frac{\beta^{(t)}}{\|\beta^{(t)}\|}$.

If this limit exist we say that $\beta^{(t)}$ converges in direction to the limit direction $\bar{\beta}^\infty$.

Multi-layer Linear Networks

Soudry et al. studied this implicit bias of gradient descent on (2) over the direct parameterization of β .

In this paper, we study the behavior of gradient descent on (1) for linear fully connected or convolutional networks.

Want to find the direction $\bar{\beta}^\infty = \lim_{t \rightarrow \infty} \frac{\mathcal{P}(\mathbf{w}^{(t)})}{\|\mathcal{P}(\mathbf{w}^{(t)})\|}$

Table of Contents

- ① Introduction
- ② Multi-layer Linear Networks
- ③ Main Results**
- ④ Discussion

Assumptions.

In the following theorems, we characterize the limiting predictor $\bar{\beta}^\infty = \lim_{t \rightarrow \infty} \frac{\beta^{(t)}}{\|\beta^{(t)}\|}$ under the following assumptions:

1. $w^{(t)}$ minimize the objective, i.e., $\mathcal{L}_{\mathcal{P}}(w^{(t)}) \rightarrow 0$.
 2. $w^{(t)}$, and consequently $\beta^{(t)} = \mathcal{P}(w^{(t)})$, converge in direction to yield a separator $\bar{\beta}^\infty = \lim_{t \rightarrow \infty} \frac{\beta^{(t)}}{\|\beta^{(t)}\|}$ with positive margin, i.e., $\min_n y_n \langle x_n, \bar{\beta}^\infty \rangle > 0$.
 3. $\nabla_{\beta} \mathcal{L}(\beta^{(t)})$ converge in direction.
- (+) The phase of the Fourier coefficients $\hat{\beta}^{(t)}$ of the linear predictors $\beta^{(t)}$ converge coordinate-wise.

Theorem 1 (Linear fully connected networks).

For any depth L , almost all linearly separable datasets $\{\mathbf{x}_n, y_n\}_{n=1}^N$, almost all initializations $\mathbf{w}^{(0)}$, and any bounded sequence of step sizes $\{\eta_t\}_t$, consider the sequence gradient descent iterates $\mathbf{w}^{(t)}$ for minimizing $\mathcal{L}_{\mathcal{P}_{\text{full}}}(\mathbf{w})$ in (1) with exponential loss over L -layer fully connected linear networks.

$$\overline{\beta}^\infty = \lim_{t \rightarrow \infty} \frac{\mathcal{P}_{\text{full}}(\mathbf{w}^{(t)})}{\|\mathcal{P}_{\text{full}}(\mathbf{w}^{(t)})\|} = \frac{\beta_{\ell_2}^*}{\|\beta_{\ell_2}^*\|},$$

where $\beta_{\ell_2}^* := \underset{\mathbf{w}}{\operatorname{argmin}} \|\beta\|_2^2$ s.t. $\forall n, y_n \langle \mathbf{x}_n, \beta \rangle \geq 1$

Theorem 2a (Linear Convolutional Networks of any Depth).

For any depth L , the limit direction $\bar{\beta}^\infty = \lim_{t \rightarrow \infty} \frac{\mathcal{P}_{\text{conv}}(w^{(t)})}{\|\mathcal{P}_{\text{conv}}(w^{(t)})\|}$ is a scaling of a first order stationary point of the following optimization problem,

$$\min_{\beta} \|\hat{\beta}\|_{2/L} \text{ s.t. } \forall n, y_n \langle \beta, x_n \rangle \geq 1$$

where $\hat{\beta} \in \mathbb{C}^D$ denote the Fourier coefficients of β , and the ℓ_p penalty given by $\|z\|_p = \left(\sum_{i=1}^D |z[i]|^p\right)^{1/p}$ is a norm for $p = 1$ and a quasi-norm for $p < 1$.

Table of Contents

- ① Introduction
- ② Multi-layer Linear Networks
- ③ Main Results
- ④ Discussion

Discussion

- Merely changing to a convolutional parameterization introduces radically different bias.
- For convenience, we studied one dimensional circular convolutions.
- These results can be directly extended to higher dimensional input signals and convolutions.
- When using convolutions as part of a larger network, with multiple parallel filters, max pooling, and non-linear activations, the situation is of course more complex, and we do not expect to get the exact same bias.
- Another important direction for future study is understanding the implicit bias for networks with multiple outputs.

Thank you for listening.

btd63@snu.ac.kr