

Kernel and rich regimes in overparameterized models

Jinwon Park

May 20, 2021

Seoul National University

Table of Contents

- 1 Introduction
- 2 Simple 2-Homogeneous Model
- 3 D-Homogeneous Models

Implicit bias in Over-Parametrized models

- In many machine learning problems, the model is highly overparametrized
 - many possible parameters for which the training loss is zero
- Training algorithm (e.g Gradient descent) can provide "implicit regularization" towards certain solutions over others
- The scale of the initialization α controls the transition between the kernel and rich regimes

The Scale of Initialization

- Model: $f : \mathbb{R}^p \times \mathcal{X} \rightarrow \mathbb{R}$
 - Map parameters $w \in \mathbb{R}^p$ and examples $x \in \mathcal{X}$ to predictions $f(w, x) \in \mathbb{R}$
 - Much of focus will be on models that are linear in x (not in w)
 - D-homogeneous in w : $F(c \cdot w) = c^D F(w)$ for all $c > 0$
- Squared loss: $L(W) = \sum_{n=1}^N (f(w, x_n) - y_n)^2$
- Gradient descent/flow: $\dot{w}(t) = -\nabla L(w(t))$
- Scale of initialization: $\alpha \in (0, \infty)$, $w_{\alpha, w_0}(0) = \alpha w_0$

The Kernel Regime

- Gradient flow depends on first-order approximation w.r.t. w

$$f(w, x) = f(w(t), x) + \langle w - w(t), \nabla f(w(t), x) \rangle + O(\|w - w(t)\|^2)$$

- Gradient flow operates on model as if it were an affine model with feature map corresponding to tangent kernel
- Minimizing the loss of affine model reaches the solution nearest to the initialization where distance is measured w.r.t. the RKHS norm
- When does kernel regime happen?
 - "Width" $\rightarrow \infty$ leads to kernel regime
 - "Scale of initialization" $\rightarrow \infty$ leads to kernel regime

The Rich Regime

- Other studies have shown very different implicit biases
 - Matrix factorization with commutative measurements and $\alpha \rightarrow 0$ leads to implicit nuclear norm regularization
 - Deep linear convolutional networks \rightarrow implicit $L_{2/\text{depth}}$ regularization in frequency domain
 - Infinite-width, depth-2 ReLU networks with infinitesimal weight decay \rightarrow minimizes $\int |f''(w, x)| dx$ second order total variations
- These are not Hilbert norms, and cannot be captured by any kernel

The Transition Regime

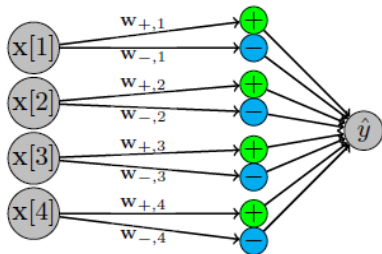
- Kernel regime: $\alpha \rightarrow \infty$
- Rich regime: $\alpha \rightarrow 0$
- Transition regime: finite α (i.e. the regime in which models are actually trained)

Table of Contents

- 1 Introduction
- 2 Simple 2-Homogeneous Model
- 3 D-Homogeneous Models

Simple 2-Homogeneous Model

- Diagonal linear neural network
 - linear model with unusual parametrization
 - $f(w, x) = \sum_{i=1}^d (w_{+,i}^2 - w_{-,i}^2)x_i = \langle \beta_w, x \rangle$



- Trained with gradient flow to minimize square loss
 - $\beta_\alpha^\infty := \lim_{t \rightarrow \infty} (w_+^2(t) - w_-^2(t))$ when $w_+(0) = w_-(0) = \alpha w_0$

The Implicit Bias and the Scale of Initialization

- $\lim_{\alpha \rightarrow \infty} \beta_{\alpha}^{\infty} = \operatorname{argmin}_{\beta} \|\beta\|_2 \text{ s.t. } L(\beta) = 0$
- $\lim_{\alpha \rightarrow 0} \beta_{\alpha}^{\infty} = \operatorname{argmin}_{\beta} \|\beta\|_1 \text{ s.t. } L(\beta) = 0$
- Theorem: for any $\alpha \in (0, \infty)$ if the gradient flow solution $\beta_{\alpha, w_0}^{\infty}$ satisfies $X\beta_{\alpha, w_0}^{\infty} = y$, then

$$\beta_{\alpha, w_0}^{\infty} = \operatorname{argmin}_{\beta} Q_{\alpha, w_0}(\beta) \text{ s.t. } L(\beta) = 0$$

where $Q_{\alpha, w_0}(\beta) = \sum_{i=1}^d \alpha^2 w_{0,i}^2 q\left(\frac{\beta_i}{\alpha^2 w_{0,i}^2}\right)$

and $q(z) = 2 - \sqrt{4 + z^2} + z \operatorname{arcsinh}\left(\frac{z}{2}\right)$

Example: Sparse regression

- $y_n \sim N(\langle \beta^*, x_n \rangle, 0.01)$ for r^* -sparse β^* with non-zero entries
 - $N = \Omega(r^* \log d)$ samples suffice for $\beta_{\ell_1}^*$ to generalize well
 - $N = \Omega(d)$ samples needed for kernel regime solution $\beta_{\ell_2}^*$ to generalize

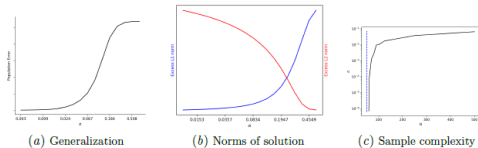


Figure 1: In (a) the population error of the gradient flow solution vs. α in the sparse regression problem described in Section 4. In (b), we plot $\|\beta_{\alpha,1}^{\infty}\|_1 - \|\beta_{\ell_1}^*\|_1$ in blue and $\|\beta_{\alpha,1}^{\infty}\|_2 - \|\beta_{\ell_2}^*\|_2$ in red vs. α . In (c), the largest α such that $\beta_{\alpha,1}^{\infty}$ achieves population error at most 0.025 is shown. The dashed line indicates the number of samples needed by $\beta_{\ell_1}^*$.

Table of Contents

- 1 Introduction
- 2 Simple 2-Homogeneous Model
- 3 D-Homogeneous Models**

D-Homogeneous Models

- $F_D(w) = \beta_{w,D} = W_+^D - W_-^D$ and $f_D(w, x) = \langle w_+^D - w_-^D, x \rangle$
- Theorem: For any $\alpha \in (0, \infty)$ and $D \geq 3$, if $X\beta_{\alpha,D}^\infty$, then

$$\beta_{\alpha,D}^\infty = \underset{\beta}{\operatorname{argmin}} Q_{\alpha,D}(\beta) \text{ s.t. } L(\beta) = 0$$

where $Q_\alpha^D(\beta) = \alpha^D \sum_{i=1}^d q_D \left(\frac{\beta_i}{\alpha^D} \right)$

and $q_D = \int h_D^{-1}$ for $h_D(z) = (1 - z)^{-\frac{D}{D-2}} - (1 + z)^{-\frac{D}{D-2}}$ /

D-Homogeneous Models

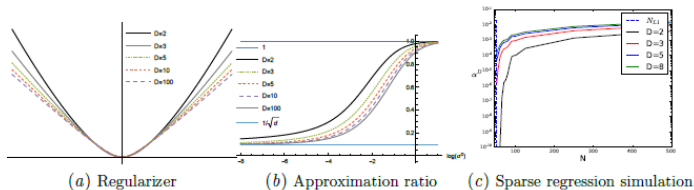


Figure 3: (a) $q_D(z)$ for several values of D . (b) The ratio $\frac{Q_\alpha^D(e_1)}{Q_\alpha^D(\mathbf{1}_d/\|\mathbf{1}_d\|_2)}$ as a function of α , where $e_1 = [1, 0, 0, \dots, 0]$ is the first standard basis vector and $\mathbf{1}_d = [1, 1, \dots, 1]$ is the all ones vector in \mathbb{R}^d . This captures the transition between approximating the ℓ_2 norm (where the ratio is 1) and the ℓ_1 norm (where the ratio is $1/\sqrt{d}$). (c) A sparse regression simulation as in Figure 1, using different order models. The y-axis is the largest α^D (the scale of β at initialization) that leads to recovery of the planted predictor to accuracy 0.025. The vertical dashed line indicates the number of samples needed in order for $\beta_{t_1}^*$ to approximate the plant.

Q&A