

## HOMEWORK 3 - V0

CSC2532 WINTER 2024

*University of Toronto*

VERSION HISTORY: V0  $\rightarrow$  V1:

- **Deadline:** Apr 5, by 23:59.
- **Submission:** You need to submit your solutions through Crowdmark, including all your derivations, plots, and your code. You can produce the file however you like (e.g. L<sup>A</sup>T<sub>E</sub>X, Microsoft Word, etc), as long as it is readable. Points will be deducted if we have a hard time reading your solutions or understanding the structure of your code.

**1. Stieltjes Transform and Double descent - 30 pts.** In the lecture, as  $d/n \rightarrow \gamma$ , we proved that the risk of ridge regression can be written as

$$(1.1) \quad \text{Risk}(\lambda) = V(\lambda) + B(\lambda),$$

where the variance and the bias terms are given as

$$\begin{aligned} V(\lambda) &\rightarrow \sigma^2 \gamma \{s(-\lambda) - \lambda s'(-\lambda)\} \\ B(\lambda) &\rightarrow \lambda^2 s'(-\lambda) \end{aligned}$$

with  $s(z) = \int \frac{1}{x-z} d\mu(z)$  denoting the Stieltjes transform of the M-P law (explicit form given in lecture). Compute the risk of *ridgeless* regression as  $\lambda \rightarrow 0_+$  by deriving expressions for  $V(0_+)$  and  $B(0_+)$ . Plot the bias, variance and the risk as a function of  $\gamma$  (No need to submit code).

**2. Implicit bias and Double descent- 70 pts.** We have  $n$  data points  $\{(\mathbf{x}_i, y_i)\}$ , each of which is a pair of feature vector  $\mathbf{x}_i \in \mathbb{R}^d$  and corresponding label  $y_i$ , and our goal is to find some parameter vector  $\boldsymbol{\theta} \in \mathbb{R}^d$  that accurately predicts a linear relation between the features and the label. We do so by minimizing the squared difference between the predictions of our linear model and the labels, summed over  $n$  data points, i.e., the least squares objective

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \hat{R}(\boldsymbol{\theta}) := \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2$$

where  $\mathbf{y} = (y_i) \in \mathbb{R}^n$  is the response,  $\mathbf{X} = (\mathbf{x}_i) \in \mathbb{R}^{n \times d}$  is the feature matrix, and  $\boldsymbol{\theta}$  is the least squares parameter. We assume that the data matrix is not degenerate, i.e.,  $\text{rank}(\mathbf{X}) = \min\{n, d\}$ . This implies that when  $n > d$ , then  $\mathbf{X}^\top \mathbf{X}$  is invertible, and when  $n < d$ ,  $\mathbf{X} \mathbf{X}^\top$  is invertible.

Since we have  $n$  data points, and we aim to learn  $d$  parameters from data, we know that when  $n < d$ , the problem is underdetermined since we have more parameters than data points; we refer to this setting as the overparameterized regime; conversely, the underparameterized regime refers to the  $n > d$  setting.

We solve this problem with gradient flow

$$(2.1) \quad \frac{d}{dt} \boldsymbol{\theta}_t = -\nabla \hat{R}(\boldsymbol{\theta}_t), \quad \boldsymbol{\theta}_0 = \mathbf{0},$$

where  $\nabla \hat{R}(\boldsymbol{\theta}) = \mathbf{X}^\top (\mathbf{X}\boldsymbol{\theta} - \mathbf{y})$ .

1. **Underparametrized regime:** Assume  $n > d$  and set  $\lambda_{\min/\max} = \lambda_{\min/\max}(\mathbf{X}^\top \mathbf{X}) > 0$  so the problem is strongly convex. Prove that

$$\|\boldsymbol{\theta}_t - \hat{\boldsymbol{\theta}}\|^2 \leq e^{-\mu t} \|\hat{\boldsymbol{\theta}}\|^2 \quad \text{with} \quad \hat{\boldsymbol{\theta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

where  $\mu = 2\lambda_{\min}$ . Remark: A similar result also holds for the gradient descent.

2. **Overparametrized regime:** When  $n < d$ , we have that  $\lambda_{\min} = 0$ , thus the objective is no longer strongly convex (still convex). Since in this case, the equation  $\mathbf{X}\boldsymbol{\theta} = \mathbf{y}$  is underdetermined, there can be infinitely many solutions achieving zero loss:  $\hat{R}(\boldsymbol{\theta}) = 0$ . However, as it turns out, GF (starting from 0) has some implicit bias and does not return an arbitrary zero-loss solution. Prove that (2.1) at  $t = \infty$  returns the min-norm solution

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \|\boldsymbol{\theta}\|^2 \quad \text{such that} \quad \mathbf{X}\boldsymbol{\theta} = \mathbf{y}.$$

In other words, in the overparameterized setting, GF finds the zero-loss solution with the smallest Euclidean norm. This phenomenon is called *implicit bias*. Hint: GF solution is always spanned by the rows of  $\mathbf{X}$  for all  $t$ .

3. Conclude that GF finds the following solutions to the least squares objective

$$(2.2) \quad \boldsymbol{\theta}^\infty = \begin{cases} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}, & n > d \\ \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{y}, & n < d. \end{cases}$$

4. (Digression) Prove that the ridge regression solution  $\boldsymbol{\theta}(\lambda) = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d)^{-1} \mathbf{X}^\top \mathbf{y}$  in the overparametrized regime converges to the same minimum norm solution in the limit  $\lambda \rightarrow 0_+$ . This is what we analyzed in the lecture as well as Problem 1 above.
5. The above calculations do not rely on a particular statistical model. In what follows, we will assume that the data generating process satisfies

$$y_i = \langle \mathbf{x}_i, \boldsymbol{\theta}_* \rangle + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2).$$

where  $\epsilon_i$  is independent of  $\mathbf{x}_i$ . If we assume that the features are Gaussian  $\mathbf{x}_i \sim \mathcal{N}(0, \mathbf{I}_d)$ , show that the population risk  $\mathcal{R}(\boldsymbol{\theta}) = \mathbb{E}[(y - \langle \mathbf{x}, \boldsymbol{\theta} \rangle)^2]$  of any (possibly random)  $\hat{\boldsymbol{\theta}}$  is

$$\mathcal{R}(\hat{\boldsymbol{\theta}}) = \mathbb{E}[\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_*\|^2] + \sigma^2.$$

Thus the excess risk is

$$\mathcal{ER}(\hat{\boldsymbol{\theta}}) = \mathbb{E}[\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_*\|^2].$$

6. Using the explicit form of the GF solution (2.2), prove that

$$\mathcal{ER}(\boldsymbol{\theta}^\infty) = \mathbb{E}[\|\boldsymbol{\theta}^\infty - \boldsymbol{\theta}_*\|^2] = \begin{cases} \sigma^2 \mathbb{E}[\text{Tr}((\mathbf{X}^\top \mathbf{X})^{-1})], & n > d + 1 \\ \frac{d-n}{d} \|\boldsymbol{\theta}_*\|^2 + \sigma^2 \mathbb{E}[\text{Tr}((\mathbf{X} \mathbf{X}^\top)^{-1})], & n < d - 1 \end{cases}$$

Hint: In the case  $d > n$ , you will need to compute  $\mathbb{E}[\mathbf{P}_R]$  where  $\mathbf{P}_R = \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{X}$  is the projection matrix to the row space of the Gaussian matrix  $\mathbf{X}$ . Note that Gaussian

matrices are rotationally invariant, i.e.  $\mathbf{X} \stackrel{d}{=} \mathbf{X}\mathbf{Q}$  for any unitary matrix  $\mathbf{Q} \in \mathbb{R}^{d \times d}$ . Due to this property, if we write the EVD of  $\mathbf{X}^\top \mathbf{X} = \mathbf{V}\mathbf{D}\mathbf{V}^\top$ , the (diagonal) matrix  $\mathbf{D}$  containing the eigenvalues is independent of the matrix  $\mathbf{V}$ . This in hand, show that the projection matrix is given as  $\mathbf{P}_R = \mathbf{V}\mathbf{S}\mathbf{V}^\top$  where  $\mathbf{S}$  is a  $d \times d$  diagonal matrix with entries either 0 or 1, with trace  $n$ . Argue that  $\mathbf{S}$  and  $\mathbf{V}$  are independent, and by symmetry,  $\mathbb{E}[\mathbf{S}] = \frac{n}{d}\mathbf{I}_d$ .

7. Using the properties of the inverse Wishart matrices<sup>1</sup>, show

$$\mathcal{ER}(\boldsymbol{\theta}^\infty) = \begin{cases} \sigma^2 \frac{d}{n-d-1}, & n > d + 1 \\ \frac{d-n}{d} \|\boldsymbol{\theta}_*\|^2 + \sigma^2 \frac{n}{d-n-1}, & n < d - 1. \end{cases}$$

8. Compare this non-asymptotic result to the asymptotic result obtained via M-P law. You may assume  $\boldsymbol{\theta}_*$  is a multivariate Gaussian. Do you observe the same asymptotic behavior as  $n, d \rightarrow \infty$  and  $d/n \rightarrow \gamma$ ?

- 3. Course evaluations - 0 pts.** Can you please fill out the course evaluations?

---

<sup>1</sup>[https://en.wikipedia.org/wiki/Inverse-Wishart\\_distribution](https://en.wikipedia.org/wiki/Inverse-Wishart_distribution)