

## High Dimensional Linear Regression

### Recap: Setup and Underparametrized Regime

From Bach 2024 we have the following model

$$y_i = x_i^\top \theta_* + \epsilon_i \quad (1)$$

where  $x_i \sim \mathcal{N}(0, 1)$  and  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$ .

In vector form Equation 1 can be rewritten

$$\underbrace{\mathbf{y}}_{n \times 1} = \underbrace{X}_{n \times d} \underbrace{\boldsymbol{\theta}_*}_{d \times 1} + \underbrace{\boldsymbol{\epsilon}}_{n \times 1} \quad (2)$$

We are doing regression using random projections, that is

$$\begin{aligned} \underbrace{\hat{\mathbf{y}}}_{n \times 1} &= \underbrace{X}_{n \times d} \underbrace{S}_{d \times m} \underbrace{\hat{\boldsymbol{\eta}}}_{m \times 1}, \quad S_{ij} \sim \mathcal{N}(0, 1) \\ \hat{\boldsymbol{\theta}} &= S \hat{\boldsymbol{\eta}} \\ \hat{\boldsymbol{\eta}} &= \lim_{\lambda \rightarrow 0} \arg \min_{\boldsymbol{\eta}} \{ \|\mathbf{y} - X S \boldsymbol{\eta}\|^2 + \lambda \|\boldsymbol{\eta}\|^2 \} \end{aligned} \quad (3)$$

The risk  $R(\hat{\boldsymbol{\theta}}) = \|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_*\|^2$  can be decomposed into bias and variance components by taking the expectation over  $\boldsymbol{\epsilon}$  as follows:

$$\begin{aligned} \mathbb{E}_{\boldsymbol{\epsilon}}[R(\hat{\boldsymbol{\theta}})] &= \mathbb{E}_{\boldsymbol{\epsilon}}[\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_*\|^2] \\ &= \mathbb{E}_{\boldsymbol{\epsilon}}[\|\underbrace{S(S^\top X X^\top S + n\lambda I)^{-1} S^\top X^\top}_{M} \mathbf{y} - \boldsymbol{\theta}_*\|^2] \\ &= \mathbb{E}_{\boldsymbol{\epsilon}}[\|M \mathbf{y} - \boldsymbol{\theta}_*\|^2] \\ &= \mathbb{E}_{\boldsymbol{\epsilon}}[\|M(X \boldsymbol{\theta}_* + \boldsymbol{\epsilon}) - \boldsymbol{\theta}_*\|^2] \\ &= \mathbb{E}_{\boldsymbol{\epsilon}}[\|(MX - I)\boldsymbol{\theta}_* + M\boldsymbol{\epsilon}\|^2] \\ &= \|(MX - I)\boldsymbol{\theta}_*\|^2 + \mathbb{E}_{\boldsymbol{\epsilon}}[\|M\boldsymbol{\epsilon}\|^2] \\ &= R^{(\text{bias})}(\hat{\boldsymbol{\theta}}) + \mathbb{E}_{\boldsymbol{\epsilon}}[R^{(\text{var})}(\hat{\boldsymbol{\theta}})] \end{aligned} \quad (4)$$

By taking the limit  $d, n, m \rightarrow \infty$  such that  $\frac{d}{n} \rightarrow \gamma$  and  $\frac{m}{n} \rightarrow \delta$  we showed that in the underparametrized regime ( $\delta < 1, \gamma < 1$ ) we get

$$\mathbb{E}_{\epsilon}[R^{(\text{var})}(\hat{\theta})] \sim \frac{\sigma^2 \delta}{1 - \delta} \quad (5)$$

In addition one can show that

$$R^{(\text{bias})}(\hat{\theta}) \sim \frac{\gamma - \delta}{1 - \delta} \frac{\|\theta_*\|^2}{\gamma} \quad (6)$$

On this lecture we will see the derivation of the limit in the overparametrized case.

### Overparametrized Regime

Recall that  $\hat{\Sigma} = \frac{1}{n}XX^\top$  and  $\text{Tr}(\hat{\Sigma}(\hat{\Sigma} + \lambda I)^{-1}) \sim \text{Tr}(\Sigma(\Sigma + \kappa(\lambda)I)^{-1})$ , where  $\kappa(\lambda) = \frac{1}{\varphi(-\lambda)}$ . When taking  $\lambda \rightarrow 0$  we get

$$\kappa(\lambda) = \begin{cases} 0, \gamma \leq 1 & (\text{underparametrized}) \\ \gamma - 1, \gamma > 1 & (\text{overparametrized}) \end{cases} \quad (7)$$

This means that in the overparametrized regime, even when we have no regularization ( $\lambda \rightarrow 0$ ) we have asymptotic regularization ( $\kappa(\lambda) = \gamma - 1 > 0$ ).

Now let us derive the limit of the bias component of the risk in the overparametrized regime. Taking the expectation of the term defined in Equation 4 we get

$$\begin{aligned} \mathbb{E}_{\epsilon}[\|M\epsilon\|^2] &= \sigma^2 \text{Tr}(M^\top M) \\ &= \frac{\sigma^2}{n} \text{Tr}[S^\top S(S^\top \hat{\Sigma} S + \lambda I)^{-1} S \hat{\Sigma} S^\top (S^\top \hat{\Sigma} S + \lambda I)^{-1}] \end{aligned} \quad (8)$$

The trace expression surely looks unwieldy, but we can make use of the following result:

**Result (Proposition 2, Bach 2024).**

$$\begin{aligned} \text{Tr}[AZ^\top(Z\Sigma Z^\top - nzI)^{-1}ZBZ^\top(Z\Sigma Z^\top - nzI)^{-1}Z] &\sim \text{Tr}[A(\Sigma + \frac{1}{\varphi(z)}I)^{-1}B(\Sigma + \frac{1}{\varphi(z)}I)^{-1}] \\ &\quad + \frac{1}{\varphi(z)^2} \text{Tr}[A(\Sigma + \frac{1}{\varphi(z)}I)^{-2}] \text{Tr}[B(\Sigma + \frac{1}{\varphi(z)}I)^{-2}] \cdot \frac{1}{n - \text{df}_2(1/\varphi(z))} \end{aligned} \quad (9)$$

Using this result (taking  $Z \rightarrow S^\top, \Sigma \rightarrow \hat{\Sigma}, n \rightarrow m, z \rightarrow -\lambda, A \rightarrow I, B \rightarrow \hat{\Sigma}, \kappa(\lambda) \rightarrow \tilde{\kappa}(\lambda), \text{df}_2 \rightarrow \tilde{\text{df}}_2$  wrt  $\hat{\Sigma}$ ), we have the following asymptotic equivalence:

$$\mathbb{E}_\epsilon[R^{(\text{var})}(\hat{\theta})] \sim \frac{\sigma^2}{n} \text{Tr}[\hat{\Sigma}(\hat{\Sigma} + \tilde{\kappa}(\lambda)I)^{-2}] + \frac{\sigma^2 \tilde{\kappa}(\lambda)^2}{n(m - \tilde{\text{df}}_2(\tilde{\kappa}(\lambda)))} \text{Tr}[(\hat{\Sigma} + \tilde{\kappa}(\lambda)I)^{-2}] \text{Tr}[\hat{\Sigma}(\hat{\Sigma} + \tilde{\kappa}(\lambda)I)^{-2}] \quad (10)$$

To simplify this formula we will use the following results:

- $\lambda \rightarrow 0 \quad \tilde{\kappa}(\lambda) = \frac{1}{\varphi(-\lambda)} \rightarrow 0$
- $\tilde{\text{df}}_2(\tilde{\kappa}(\lambda)) = \text{Tr}[\hat{\Sigma}^2(\hat{\Sigma}^+)^2] = n$
- Using the push-through identity<sup>1</sup> we have  $\hat{\Sigma}(\hat{\Sigma} + \tilde{\kappa}I)^{-2} = nX^\top(XX^\top + n\tilde{\kappa}I)^{-2}X \rightarrow nX^\top(XX^\top)^{-2}X$
- Using the Woodbury matrix identity<sup>2</sup> we have

$$\begin{aligned} \tilde{\kappa}(\lambda)^2(\hat{\Sigma} + \tilde{\kappa}(\lambda)I)^{-2} &= (I - X^\top(XX^\top + n\tilde{\kappa}(\lambda)I)^{-1}X)^2 \\ &\rightarrow (I - \underbrace{X^\top(XX^\top)^{-1}X}_{\text{projection matrix } P_X})^2 \\ &= (I - P_X)^2 \\ &= I - P_X \end{aligned} \quad (11)$$

Putting it all together we have

$$\begin{aligned} \mathbb{E}_\epsilon[R^{(\text{var})}(\hat{\theta})] &\sim \sigma^2 \text{Tr}[X^\top(XX^\top)^{-2}X] + \frac{\sigma^2}{m - n} \text{Tr}[I - P_X] \text{Tr}[(XX^\top)^{-1}] \\ &= \sigma^2 \text{Tr}[(XX^\top)^{-1}] + \frac{\sigma^2}{m - n} \text{Tr}[I - P_X] \text{Tr}[(XX^\top)^{-1}] \end{aligned} \quad (12)$$

Note that  $\text{Tr}[(XX^\top - n\lambda I)^{-1}] = \hat{\varphi}(z) \xrightarrow{n \rightarrow \infty} \varphi(z) \xrightarrow{z \rightarrow 0} \frac{1}{\kappa(0)} = \frac{1}{\gamma - 1}$ . Also,  $\text{Tr}[I - P_X] = d - n$ , because  $\text{rank}(P_X) = n$ . Thus

$$\begin{aligned} \mathbb{E}_\epsilon[R^{(\text{var})}(\hat{\theta})] &\sim \sigma^2 \text{Tr}[(XX^\top)^{-1}] + \frac{\sigma^2}{m - n} \text{Tr}[I - P_X] \text{Tr}[(XX^\top)^{-1}] \\ &\sim \sigma^2 \frac{1}{\gamma - 1} + \frac{\sigma^2}{m - n} (d - n) \frac{1}{\gamma - 1} \\ &= \sigma^2 \frac{1}{\gamma - 1} \left(1 + \frac{d - n}{m - n}\right) \\ &= \sigma^2 \left(\frac{1}{\gamma - 1} + \frac{1}{\delta - 1}\right) \end{aligned} \quad (13)$$

<sup>1</sup>Push-through identity:  $X(X^\top X + \kappa I)^{-1} = (XX^\top + \kappa I)^{-1}X$

<sup>2</sup>Woodbury matrix identity:  $\kappa(X^\top X + \kappa I)^{-1} = I - X^\top(XX^\top + \kappa I)^{-1}X$

Therefore we can see that  $\mathbb{E}_\epsilon[R^{(\text{var})}(\hat{\boldsymbol{\theta}})]$  decreases as  $\delta$  or  $\gamma$  increases.

The result for the bias component of the risk is the following (we will not prove it):

$$R^{(\text{bias})} \sim \left(1 - \frac{1}{\gamma}\right) \frac{\delta}{1-\delta} \|\boldsymbol{\theta}_*\|^2 \quad (14)$$

Thus  $R^{(\text{bias})}$  increases with  $\gamma$  and decreases with  $\delta$ .

## Double Descent

Combining the results from the previous lecture (underparametrized regime) and today's results (overparametrized regime), we illustrate the double descent phenomenon in Figure 1.

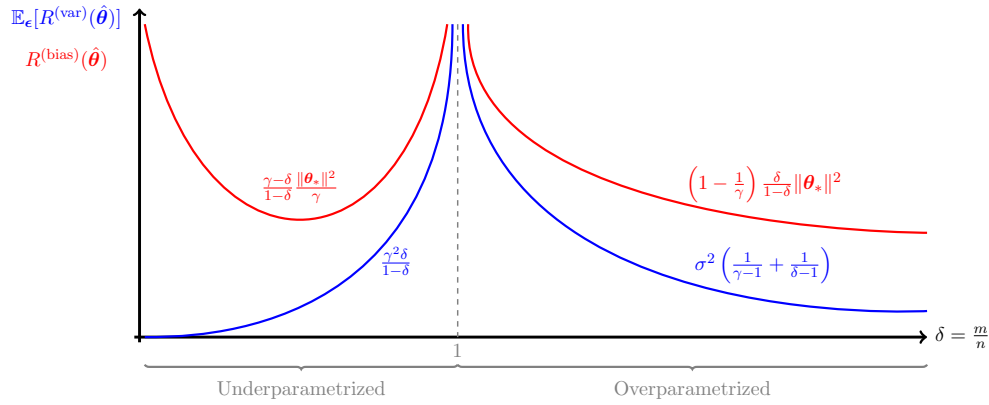


Figure 1: Illustration of the double descent phenomenon.

## References

Bach, Francis (2024). “High-Dimensional Analysis of Double Descent for Linear Regression with Random Projections”. In: *SIAM Journal on Mathematics of Data Science* 6.1, pp. 26–50.