

Chapter 4

- Linear Regression
 - Ordinary Least Squares
 - Ridge Regression

Linear Regression:

Consider the problem

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n (y_i - x_i^T \theta)^2$$

We can think of this as approximating

$$\min_{\theta} \mathcal{L}(\theta) = \mathbb{E}_{(x,y)} (y - x^T \theta)^2$$

Let's write compactly

$$\min_{\theta \in \mathbb{R}^d} \mathcal{R}_n(\theta) := \frac{1}{n} \|y - X\theta\|_2^2$$

where $y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$, $X = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix}$

An optimal solution θ_n^{OLS} is called an ordinary least squares estimator.

We will assume for now that X has full column rank
So $d < n$

$$X = \begin{bmatrix} \text{---} \\ \text{---} \\ \text{---} \\ \vdots \\ \text{---} \end{bmatrix}$$

Lemma: Θ_n^{OLS} exists and is
unique and is given by

$$\Theta_n^{\text{OLS}} = (X^T X)^{-1} X^T y$$

proof: \hat{R} is coercive and continuous
and therefore has a minimizer.

Taking the gradient yields

$$0 = \nabla R_n(\theta) = \frac{2}{n} X^T (X\theta - y)$$

Solve for θ . \square

Lemma: The prediction

$$X\theta_n^{\text{OLS}} = X(X^T X)^{-1} X^T y$$
 is

the orthogonal projection of y onto $\text{Range}(X) \in \mathbb{R}^d$.

So θ_n^{OLS} is obtained by

① project y into $\text{Range}(X)$
to get \hat{y}

② Solve $X\theta = \hat{y}$ for θ .

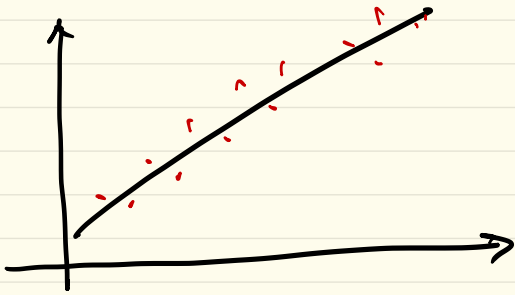
Next, let's estimate the generalization error:

$$R(\theta_n^{\text{OLS}}) - \underbrace{\min_{\theta} R(\theta)}_{R^*}$$

Two cases:

① Fixed design:

- $\exists \theta^*$ s.t. $y_i = x_i^T \theta^* + \varepsilon_i \quad \forall i$
- All ε_i are independent with $E[\varepsilon_i] = 0$ and $E[\varepsilon_i^2] \leq \sigma^2$.



② Random Design:

Both (x_i, y_i) are random.
More on this later.

Assume the fixed design setting. Set $\hat{\Sigma} = \frac{1}{n} X^T X$ and suppose $\hat{\Sigma}$ is invertible.

Define the inner product

$$\langle \theta, \theta' \rangle_{\hat{\Sigma}} = \theta^T \hat{\Sigma} \theta'$$

and induced norm

$$\|\theta\|_{\hat{\Sigma}}^2 = \langle \theta, \theta \rangle_{\hat{\Sigma}} = \|\hat{\Sigma}^{-1/2} \theta\|_2^2$$

Lemma:

$$R(\theta) - R^* = \|\theta - \theta^*\|_{\hat{\Sigma}}^2 \quad \forall \theta.$$

and $R^* = \sigma^2$?

pf:

$$\begin{aligned} R(\theta) &= \mathbb{E}_y \frac{1}{n} \|y - X\theta\|_2^2 \\ &= \frac{1}{n} \mathbb{E}_\varepsilon \|X(\theta^* - \theta) + \varepsilon\|_2^2 \\ &= \|\theta^* - \theta\|_{\hat{\Sigma}}^2 + \frac{1}{n} \mathbb{E}_\varepsilon \|\varepsilon\|_2^2 \\ &= \|\theta^* - \theta\|_{\hat{\Sigma}}^2 + \sigma^2 \end{aligned}$$

So $R^* = \sigma^2$. \square

Recall for any estimator $\hat{\theta}_n$,
we may write

$$\mathbb{E} \|\hat{\theta}_n - \theta^*\|_{\hat{\Sigma}}^2 = \underbrace{\mathbb{E} \|\hat{\theta}_n - \theta^*\|_{\hat{\Sigma}}^2}_{\text{Bias}(\hat{\theta}_n)} + \underbrace{\mathbb{E} \|\hat{\theta}_n - \theta^*\|_{\hat{\Sigma}}^2}_{\text{Var}(\hat{\theta}_n)}$$

Lemma: $\text{Bias}(\theta_n^{\text{OLS}}) = 0$

$$\mathbb{E}(\theta_n^{\text{OLS}} - \theta^*)(\theta_n^{\text{OLS}} - \theta^*)^T = \frac{\sigma^2}{n} \hat{\Sigma}^{-1}$$

pf: Recall $\theta_n^{OLS} = (X^T X)^{-1} X^T y$
 $= (X^T X)^{-1} X^T (\Phi \theta^* + \epsilon)$
 $= \theta^* + (X^T X)^{-1} X^T \epsilon$

So $\mathbb{E} \theta_n^{OLS} = \theta^*$ and

$$\mathbb{E} (\theta_n^{OLS} - \theta^*) (\theta_n^{OLS} - \theta^*)^T$$

$$= (X^T X)^{-1} X^T \underbrace{\mathbb{E} \epsilon \epsilon^T}_{\sigma^2 I} X (X^T X)^{-1}$$

$$= \sigma^2 (X^T X)^{-1} = \frac{\sigma^2}{n} \Sigma \quad \square$$

Cor: $\mathbb{E} \mathcal{R}(\theta_n^{OLS}) - \mathcal{R}(\theta^*) = \frac{\sigma^2 d}{n}$

pf: $\text{Var}(\theta_n^{OLS}) = \mathbb{E} \|\theta_n^{OLS} - \theta^*\|_{\Sigma}^2$
 $= \mathbb{E} \text{tr}(\Sigma (\theta_n^{OLS} - \theta^*) (\theta_n^{OLS} - \theta^*)^T)$
 $= \frac{\sigma^2}{n} \text{tr}(I) = \frac{d \sigma^2}{n} \quad \square$

Remark:

$$\begin{aligned} \mathbb{E} \mathcal{R}_n(\theta_n^{\text{OLS}}) &= \mathbb{E}_{y_1, \dots, y_n} \frac{1}{n} \|X\theta_n^{\text{OLS}} - y\|_2^2 \\ &= \mathbb{E}_{y_1, \dots, y_n} \frac{1}{n} \|(X(X^T X)^{-1} X^T - I)y\|_2^2 \\ &= \mathbb{E}_{\varepsilon} \frac{1}{n} \|(X(X^T X)^{-1} X^T - I)(\mathbb{E}\theta + \varepsilon)\|_2^2 \end{aligned}$$

$$= \mathbb{E}_{\varepsilon} \frac{1}{n} \|(X(X^T X)^{-1} X^T - I)\varepsilon\|_2^2$$

$$= \frac{\sigma^2}{n} \text{tr}((P - I)(P - I)^T)$$

$$= \frac{\sigma^2}{n} \text{tr}(P - I) = \frac{\sigma^2}{n} (n - d)$$

$$= \sigma^2 - \frac{\sigma^2 d}{n}$$

downward
bias.

Better estimators than OLS
trade bias for variance.

Defn: For $\lambda > 0$, define the
ridge least-squares estimator

$$\Theta_n^\lambda = \underset{\Theta}{\operatorname{argmin}} \frac{1}{n} \|y - X\Theta\|_2^2 + \lambda \|\Theta\|_2^2$$

Lemma: $\Theta_n^\lambda = \frac{1}{n} (\hat{\Sigma} + \lambda I)^{-1} X^T y$

pf: Set $\nabla = 0$ and solve.

Lemma: $\operatorname{Bias}(\Theta_n^\lambda) = \lambda^2 \Theta^{*T} (\hat{\Sigma} + \lambda I)^{-2} \hat{\Sigma} \Theta^{*}$

$$\operatorname{Var}(\Theta_n^\lambda) = \frac{\sigma^2}{n} \operatorname{tr}(\hat{\Sigma}^2 (\hat{\Sigma} + \lambda I)^{-2})$$

$$\begin{aligned}
 \text{pf: } E\theta_n^\lambda &= \frac{1}{n} (\hat{\Sigma} + \lambda I)^{-1} X^T X \theta^* \\
 &= (\hat{\Sigma} + \lambda I)^{-1} \hat{\Sigma} \theta^* \\
 &= \theta^* + ((\hat{\Sigma} + \lambda I)^{-1} \hat{\Sigma} - I) \theta^* \\
 &= \theta^* + ((\hat{\Sigma} + \lambda I)^{-1} (\hat{\Sigma} + \lambda I - \lambda I) - I) \theta^* \\
 &= \theta^* - \lambda (\hat{\Sigma} + \lambda I)^{-1} \theta^*
 \end{aligned}$$

$$\begin{aligned}
 \text{So Bias}(\theta_n^\lambda) &= \lambda^2 \left\| (\hat{\Sigma} + \lambda I)^{-1} \theta^* \right\|_{\hat{\Sigma}}^2 \\
 &= \lambda^2 \theta^{*\top} (\hat{\Sigma} + \lambda I)^{-2} \hat{\Sigma} \theta^*
 \end{aligned}$$

$$\begin{aligned}
 \text{Var} \theta_n^\lambda &= E \left\| \theta_n^\lambda - E \theta_n^\lambda \right\|_{\hat{\Sigma}}^2 \\
 &= E \left\| \frac{1}{n} (\hat{\Sigma} + \lambda I)^{-1} X^T y - \frac{1}{n} (\hat{\Sigma} + \lambda I)^{-1} X^T X \theta^* \right\|_{\hat{\Sigma}}^2 \\
 &= E \left\| \frac{1}{n} (\hat{\Sigma} + \lambda I)^{-1} X^T \varepsilon \right\|_{\hat{\Sigma}}^2 \\
 &= \frac{\sigma^2}{n} \text{tr} \left(\hat{\Sigma} (\hat{\Sigma} + \lambda I)^{-1} \hat{\Sigma} (\hat{\Sigma} + \lambda I)^{-1} \right)
 \end{aligned}$$

The quantity $\text{evals of } \hat{\Sigma}$

$$\text{tr}(\hat{\Sigma}^2(\hat{\Sigma} + \lambda I)^{-2}) = \sum_{i=1}^d \frac{\lambda_i^2}{(\lambda_i + \lambda)^2} \leq d$$

is called the degrees of freedom

Summary:

$$\mathbb{E} \mathcal{R}(\hat{\theta}_\lambda) - \mathcal{R}^* = \lambda^2 \theta^{*\top} (\hat{\Sigma} + \lambda I)^{-2} \hat{\Sigma} \theta^* + \frac{\sigma^2}{n} \text{tr}(\hat{\Sigma}^2(\hat{\Sigma} + \lambda I)^{-2})$$

Remarks:

- No explicit dependence on d
- Optimizing over λ can lead to a lower value than θ_n^{OLS} .

Note: $\lambda^2 \theta^{*\top} (\hat{\Sigma} + \lambda I)^{-2} \hat{\Sigma} \theta^* \leq \lambda \underbrace{\|(\hat{\Sigma} + \lambda I)^{-2} \lambda \hat{\Sigma}\|_p}_{\leq \lambda \|\theta^*\|_2} \|\theta^*\|_2$

$$\frac{\sigma^2}{n} \text{tr}(\hat{\Sigma}^2 (\hat{\Sigma} + \lambda I)^{-2}) \leq \frac{\sigma^2}{\lambda n} \text{tr}(\hat{\Sigma} \lambda \hat{\Sigma} (\hat{\Sigma} + \lambda I)^{-2})$$

$$\leq \frac{\sigma^2}{2\lambda n} \text{tr}(\hat{\Sigma})$$

Optimizing over λ gives

$$\bar{\lambda} = \frac{\sigma}{\|\theta^*\|_2} \cdot \sqrt{\frac{\text{tr}(\hat{\Sigma})}{n}}$$

Then $\mathbb{E} \mathcal{R}(\theta_n^\top) - \mathcal{R}^* \leq \sigma \|\theta^*\|_2 \sqrt{\frac{\text{tr}(\hat{\Sigma})}{n}}$

Note $\text{tr}(\hat{\Sigma}) = \frac{1}{n} \sum_{i=1}^n \|x_i\|^2$

typically independent of d .

• So worse rate in n , but better in σ and no dependence on d !

Random Design:

Suppose $\exists \theta^* \in \mathbb{R}^d$ s.t.

$$y = x^T \theta^* + \varepsilon$$

- ε is independent from x ,
and $\mathbb{E}\varepsilon = 0$, $\mathbb{E}\varepsilon^2 = \sigma^2$

Define $\Sigma = \mathbb{E}xx^T$

Lemma: $\mathcal{R}(\theta) - \mathcal{R}^* = \|\theta - \theta^*\|_{\Sigma}^2$

and $\mathcal{R}^* = \sigma^2$.

Pf:

$$\begin{aligned}\mathcal{R}(\theta) &= \mathbb{E}(y - x^T \theta)^2 \\ &= \mathbb{E}(x^T(\theta^* - \theta) + \varepsilon)^2 \\ &= (\theta^* - \theta)^T \mathbb{E}xx^T(\theta^* - \theta) + \sigma^2.\end{aligned}$$

□

Lemma:

$$E R(\theta_n^{OLS}) - R^* = \frac{\sigma^2}{n} E \text{tr}(\Sigma \hat{\Sigma}^{-1})$$

pf: Recall $\theta_n^{OLS} = \theta^* + \frac{1}{n} \hat{\Sigma}^{-1} X^T \varepsilon$

So

$$\begin{aligned} E R(\theta_n^{OLS}) - R^* &= E \|\theta_n^{OLS} - \theta^*\|_{\Sigma}^2 \\ &= \frac{1}{n^2} E \varepsilon^T X \hat{\Sigma}^{-1} \Sigma \hat{\Sigma}^{-1} X^T \varepsilon \\ &= \frac{1}{n^2} \text{tr} (X \hat{\Sigma}^{-1} \Sigma \hat{\Sigma}^{-1} X^T E(\varepsilon \varepsilon^T)) \\ &= \frac{\sigma^2}{n^2} \text{tr} (\hat{\Sigma}^{-1} \Sigma \hat{\Sigma}^{-1} X^T X) \\ &= \frac{\sigma^2}{n} \text{tr} (\hat{\Sigma}^{-1} \Sigma) \quad \square \end{aligned}$$

So we must compute

$$\mathbb{E} \text{tr}(\Sigma \hat{\Sigma}^{-1}) = \mathbb{E} \text{tr}(\Sigma^{1/2} \Sigma^{\lambda-1} \Sigma^{1/2})$$

$$\Sigma^{1/2} \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^T \right)^{-1} \Sigma^{1/2}$$
$$= \left(\frac{1}{n} \sum_{i=1}^n (\Sigma^{-1/2} x_i) (\Sigma^{-1/2} x_i)^T \right)^{-1}$$

$$= n \cdot \left(\sum_{i=1}^n z_i z_i^T \right)^{-1}$$

where $z_i = \Sigma^{-1/2} x_i$

If $x_i \sim N(0, \Sigma)$, then

$z_i \sim N(0, I)$ and

$(Z^T Z)^{-1}$ has a special distribution.

Can compute $\mathbb{E} (Z^T Z)^{-1} = \frac{1}{n-d-1} I$