



Statistical Learning II

Lecture 8 - Bias-Variance decomposition (Continued)

Bruno Loureiro
@ CSD, DI-ENS & CNRS

brloureiro@gmail.com

Risk of OLS

Therefore, we have the following final result for the excess risk of OLS

$$\mathbb{E}_{\varepsilon} \left[\mathcal{R}(\hat{\boldsymbol{\theta}}_{OLS}) \right] - \sigma^2 = \sigma^2 \frac{d}{n}$$

Remarks:

- Excess risk is proportional to the noise level $\mathbb{E}[\varepsilon^2] = \sigma^2$.
- Excess risk is proportional to the data dimension.
- To achieve excess risk $\Delta \mathcal{R} < \delta$, need:

$$n > \frac{\sigma^2 d}{\delta}$$

samples.

Bias-variance decomposition

Generally, if we have a data generative model for the training data $\mathcal{D} = \{(\mathbf{x}_i, y_i) \in \mathbb{R}^{d+1} : i = 1, \dots, n\}$:

$$y_i = f_{\star}(\mathbf{x}) + \varepsilon_i = \text{signal} + \text{noise}$$

With $\mathbb{E}[\varepsilon] = 0$ and $\mathbb{E}[\varepsilon^2] = \sigma^2$

Bias-variance decomposition

Generally, if we have a data generative model for the training data $\mathcal{D} = \{(\mathbf{x}_i, y_i) \in \mathbb{R}^{d+1} : i = 1, \dots, n\}$:

$$y_i = f_{\star}(\mathbf{x}) + \varepsilon_i = \text{signal} + \text{noise}$$

With $\mathbb{E}[\varepsilon] = 0$ and $\mathbb{E}[\varepsilon^2] = \sigma^2$, we can decompose the excess risk:

$$\mathbb{E}_{\varepsilon}[\mathcal{R}(\hat{\boldsymbol{\theta}})] - \sigma^2 = \mathbb{E} \left[(f_{\star}(\mathbf{x}) - f_{\hat{\boldsymbol{\theta}}}(\mathbf{x}))^2 \right]$$

Bias-variance decomposition

Generally, if we have a data generative model for the training data $\mathcal{D} = \{(\mathbf{x}_i, y_i) \in \mathbb{R}^{d+1} : i = 1, \dots, n\}$:

$$y_i = f_{\star}(\mathbf{x}) + \varepsilon_i = \text{signal} + \text{noise}$$

With $\mathbb{E}[\varepsilon] = 0$ and $\mathbb{E}[\varepsilon^2] = \sigma^2$, we can decompose the excess risk:

$$\begin{aligned}\mathbb{E}_{\varepsilon}[\mathcal{R}(\hat{\boldsymbol{\theta}})] - \sigma^2 &= \mathbb{E} \left[(f_{\star}(\mathbf{x}) - f_{\hat{\theta}}(\mathbf{x}))^2 \right] \\ &= \mathbb{E} \left[(f_{\star}(\mathbf{x}) - \mathbb{E}_{\varepsilon}[f_{\hat{\theta}}(\mathbf{x})] + \mathbb{E}_{\varepsilon}[f_{\hat{\theta}}(\mathbf{x})] - f_{\hat{\theta}}(\mathbf{x}))^2 \right]\end{aligned}$$

Bias-variance decomposition

Generally, if we have a data generative model for the training data $\mathcal{D} = \{(\mathbf{x}_i, y_i) \in \mathbb{R}^{d+1} : i = 1, \dots, n\}$:

$$y_i = f_{\star}(\mathbf{x}) + \varepsilon_i = \text{signal} + \text{noise}$$

With $\mathbb{E}[\varepsilon] = 0$ and $\mathbb{E}[\varepsilon^2] = \sigma^2$, we can decompose the excess risk:

$$\begin{aligned}\mathbb{E}_{\varepsilon}[\mathcal{R}(\hat{\boldsymbol{\theta}})] - \sigma^2 &= \mathbb{E} \left[(f_{\star}(\mathbf{x}) - f_{\hat{\theta}}(\mathbf{x}))^2 \right] \\ &= \mathbb{E} \left[(f_{\star}(\mathbf{x}) - \mathbb{E}_{\varepsilon}[f_{\hat{\theta}}(\mathbf{x})] + \mathbb{E}_{\varepsilon}[f_{\hat{\theta}}(\mathbf{x})] - f_{\hat{\theta}}(\mathbf{x}))^2 \right] \\ &= \mathbb{E} \left[(f_{\star}(\mathbf{x}) - \mathbb{E}_{\varepsilon}[f_{\hat{\theta}}(\mathbf{x})])^2 \right] + \mathbb{E} \left[(\mathbb{E}_{\varepsilon}[f_{\hat{\theta}}(\mathbf{x})] - f_{\hat{\theta}}(\mathbf{x}))^2 \right]\end{aligned}$$

Bias-variance decomposition

Generally, if we have a data generative model for the training data $\mathcal{D} = \{(\mathbf{x}_i, y_i) \in \mathbb{R}^{d+1} : i = 1, \dots, n\}$:

$$y_i = f_{\star}(\mathbf{x}) + \varepsilon_i = \text{signal} + \text{noise}$$

With $\mathbb{E}[\varepsilon] = 0$ and $\mathbb{E}[\varepsilon^2] = \sigma^2$, we can decompose the excess risk:

$$\begin{aligned}\mathbb{E}_{\varepsilon}[\mathcal{R}(\hat{\boldsymbol{\theta}})] - \sigma^2 &= \mathbb{E} \left[(f_{\star}(\mathbf{x}) - f_{\hat{\theta}}(\mathbf{x}))^2 \right] \\ &= \mathbb{E} \left[(f_{\star}(\mathbf{x}) - \mathbb{E}_{\varepsilon}[f_{\hat{\theta}}(\mathbf{x})] + \mathbb{E}_{\varepsilon}[f_{\hat{\theta}}(\mathbf{x})] - f_{\hat{\theta}}(\mathbf{x}))^2 \right] \\ &= \mathbb{E} \left[(f_{\star}(\mathbf{x}) - \mathbb{E}_{\varepsilon}[f_{\hat{\theta}}(\mathbf{x})])^2 \right] + \mathbb{E} \left[(\mathbb{E}_{\varepsilon}[f_{\hat{\theta}}(\mathbf{x})] - f_{\hat{\theta}}(\mathbf{x}))^2 \right] \\ &= \text{bias}^2 + \text{variance}\end{aligned}$$

Bias-variance decomposition

$$\mathbb{E}_{\boldsymbol{\varepsilon}}[\mathcal{R}(\hat{\boldsymbol{\theta}})] - \sigma^2 = \mathcal{B} + \mathcal{V}$$

$$\mathcal{B} = \mathbb{E} \left[(f_{\star}(\mathbf{x}) - \mathbb{E}_{\boldsymbol{\varepsilon}}[f_{\hat{\boldsymbol{\theta}}}(\mathbf{x})])^2 \right]$$

$$\mathcal{V} = \mathbb{E} \left[(\mathbb{E}_{\boldsymbol{\varepsilon}}[f_{\hat{\boldsymbol{\theta}}}(\mathbf{x})] - f_{\hat{\boldsymbol{\theta}}}(\mathbf{x}))^2 \right]$$

Bias-variance decomposition

$$\mathbb{E}_{\epsilon}[\mathcal{R}(\hat{\theta})] - \sigma^2 = \mathcal{B} + \mathcal{V}$$
$$\mathcal{B} = \mathbb{E} \left[(f_{\star}(\mathbf{x}) - \mathbb{E}_{\epsilon}[f_{\hat{\theta}}(\mathbf{x})])^2 \right]$$
$$\mathcal{V} = \mathbb{E} \left[(\mathbb{E}_{\epsilon}[f_{\hat{\theta}}(\mathbf{x})] - f_{\hat{\theta}}(\mathbf{x}))^2 \right]$$



Recall the the **approximation + estimation decomposition** from lecture 3:

$$\mathcal{R}(\theta) - \mathcal{R}_{\star} = \left(\mathcal{R}(\theta) - \inf_{\theta' \in \Theta} \mathcal{R}(\theta') \right) + \left(\inf_{\theta' \in \Theta} \mathcal{R}(\theta') - \mathcal{R}_{\star} \right)$$

Bias-variance decomposition

$$\mathbb{E}_{\varepsilon}[\mathcal{R}(\hat{\theta})] - \sigma^2 = \mathcal{B} + \mathcal{V}$$
$$\mathcal{B} = \mathbb{E} \left[(f_{\star}(\mathbf{x}) - \mathbb{E}_{\varepsilon}[f_{\hat{\theta}}(\mathbf{x})])^2 \right]$$
$$\mathcal{V} = \mathbb{E} \left[(\mathbb{E}_{\varepsilon}[f_{\hat{\theta}}(\mathbf{x})] - f_{\hat{\theta}}(\mathbf{x}))^2 \right]$$



Recall the the **approximation + estimation decomposition** from lecture 3:

$$\mathcal{R}(\theta) - \mathcal{R}_{\star} = \left(\mathcal{R}(\theta) - \inf_{\theta' \in \Theta} \mathcal{R}(\theta') \right) + \left(\inf_{\theta' \in \Theta} \mathcal{R}(\theta') - \mathcal{R}_{\star} \right)$$

For the OLS setting from before ($\text{rank}(X) = d < n$):

$$\mathbb{E}[f_{\hat{\theta}}(\mathbf{x})] = \langle \boldsymbol{\theta}_{\star}, \mathbf{x} \rangle = f_{\star}(\mathbf{x}) \quad \Rightarrow \quad \mathcal{B} = 0 \quad \mathcal{V} = \sigma^2 \frac{d}{n}$$

Marvels and pitfalls of OLS

To summarise, the OLS estimator $\hat{\theta}_{\text{OLS}}(\mathbf{X}, \mathbf{y}) = \mathbf{X}^+ \mathbf{y}$:

- Can only fit **linear functions**.
- For $n > d$, **has low bias** $\mathcal{B} = 0$
- When, $n \gg d$, has **low variance** $\mathcal{V} = \sigma^2 \frac{d}{n}$

Marvels and pitfalls of OLS

To summarise, the OLS estimator $\hat{\theta}_{\text{OLS}}(X, y) = X^+y$:

- Can only fit **linear functions**.
- For $n > d$, **has low bias** $\mathcal{B} = 0$
- When, $n \gg d$, has **low variance** $\mathcal{V} = \sigma^2 \frac{d}{n}$

But what about $n \approx d$? Consider for instance $n = d$.

$$X \in \mathbb{R}^{d \times d} \text{ is invertible} \quad \Rightarrow \quad y = X\hat{\theta}_{\text{OLS}}$$

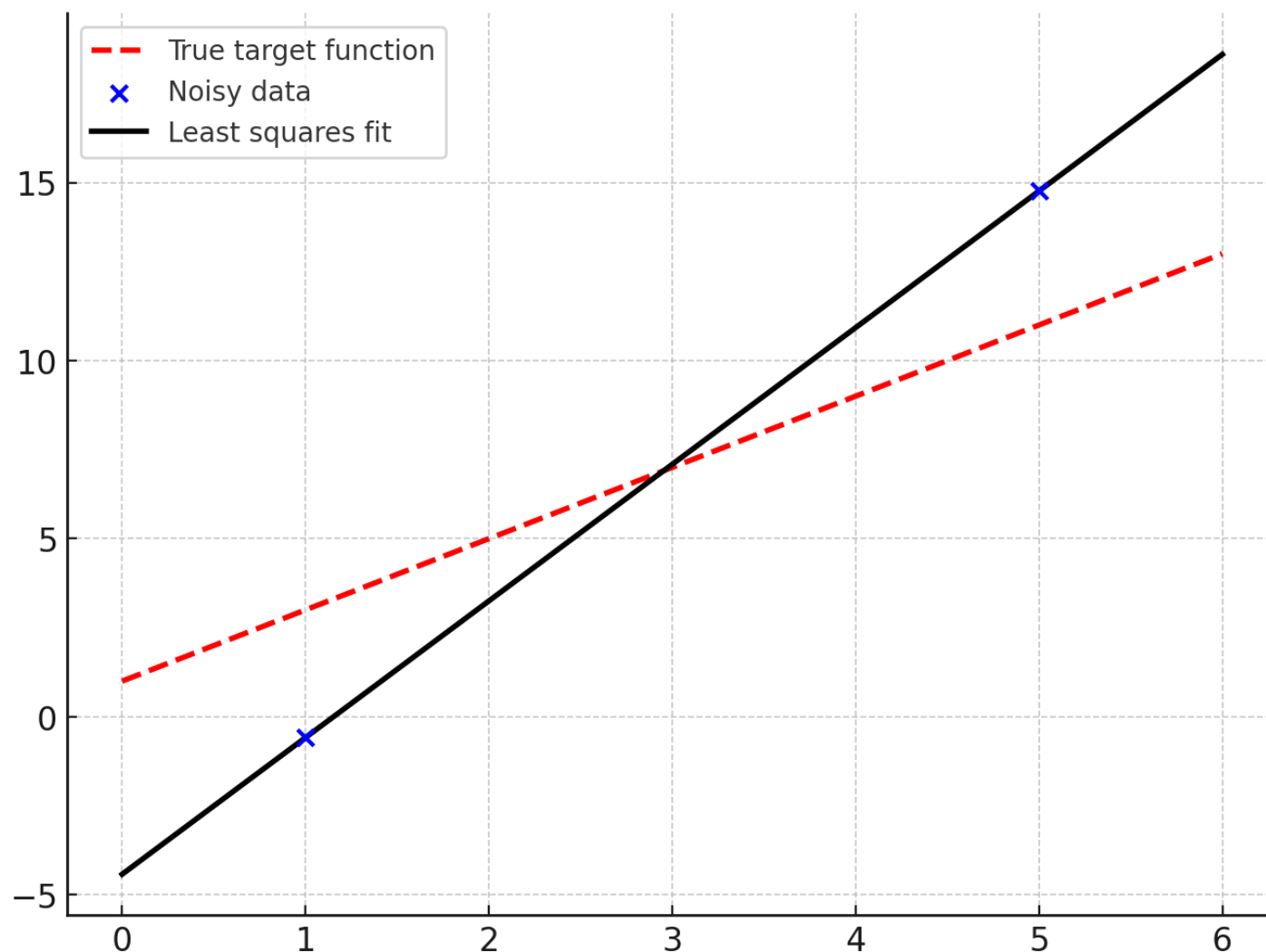
**interpolates the
training data.**

Marvels and pitfalls of OLS

But what about $n \approx d$? Consider for instance $n = d$.

$$X \in \mathbb{R}^{d \times d} \text{ is invertible} \Rightarrow y = X\hat{\theta}_{\text{OLS}}$$

interpolates the
training data.



$$\mathbb{E}_{\epsilon}[\mathcal{R}(\hat{\theta}_{\text{OLS}})] = 2\sigma^2$$

$$\hat{\mathcal{R}}_n(\hat{\theta}_{\text{OLS}}) = 0$$



The test error above is
valid for the fixed design.

Marvels and pitfalls of OLS

Recall that:

$$\hat{\theta}_{OLS}(X, y) = \theta_{\star} + \frac{1}{n} \hat{\Sigma}_n^{-1} X^{\top} \epsilon$$

Marvels and pitfalls of OLS

Recall that:

$$\begin{aligned}\hat{\boldsymbol{\theta}}_{OLS}(X, y) &= \boldsymbol{\theta}_{\star} + \frac{1}{n} \hat{\boldsymbol{\Sigma}}_n^{-1} X^{\top} \boldsymbol{\varepsilon} \\ &= \boldsymbol{\theta}_{\star} + \sum_{j=1}^d \frac{1}{\sigma_j} \langle \mathbf{u}_j, \boldsymbol{\varepsilon} \rangle \mathbf{v}_j\end{aligned}$$

Marvels and pitfalls of OLS

Recall that:

$$\begin{aligned}\hat{\boldsymbol{\theta}}_{OLS}(X, y) &= \boldsymbol{\theta}_{\star} + \frac{1}{n} \hat{\boldsymbol{\Sigma}}_n^{-1} X^{\top} \boldsymbol{\varepsilon} \\ &= \boldsymbol{\theta}_{\star} + \sum_{j=1}^d \frac{1}{\sigma_j} \langle \mathbf{u}_j, \boldsymbol{\varepsilon} \rangle \mathbf{v}_j\end{aligned}$$

- Hence:
- **signal** is stronger in directions with larger s.v.
 - **noise** dominates directions with smaller s.v.

OLS has larger variance for data with small “**effective dimension**”.

What to do?

Classical strategies to mitigate variance:

- Dimensionality reduction: PCA, random projections (sketching), etc.
- Variable subset selection: Stepwise selection, best Subset Selection, etc.
- Regularisation: ridge, LASSO, etc.

Ridge regression

Ridge regression

Note the averaged norm of the OLS is given by:

$$\mathbb{E}_{\boldsymbol{\varepsilon}} \left[||\hat{\boldsymbol{\theta}}_{OLS}||_2^2 \right] = ||\boldsymbol{\theta}_{\star}||_2^2 + \sigma^2 \sum_{j=1}^d \frac{1}{\sigma_j^2}$$

Therefore, small s.v.s lead to larger expected norm.

Ridge regression

Note the averaged norm of the OLS is given by:

$$\mathbb{E}_{\epsilon} \left[||\hat{\boldsymbol{\theta}}_{OLS}||_2^2 \right] = ||\boldsymbol{\theta}_{\star}||_2^2 + \sigma^2 \sum_{j=1}^d \frac{1}{\sigma_j^2}$$

Therefore, small s.v.s lead to larger expected norm.



Key idea: penalise the norm.

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \frac{1}{2n} ||\mathbf{y} - \mathbf{X}\boldsymbol{\theta}||_2^2 + \frac{\lambda}{2} ||\boldsymbol{\theta}||_2^2$$

Ridge regression

Note the averaged norm of the OLS is given by:

$$\mathbb{E}_{\varepsilon} \left[||\hat{\boldsymbol{\theta}}_{OLS}||_2^2 \right] = ||\boldsymbol{\theta}_{\star}||_2^2 + \sigma^2 \sum_{j=1}^d \frac{1}{\sigma_j^2}$$

Therefore, small s.v.s lead to larger expected norm.



Key idea: penalise the norm.

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \frac{1}{2n} ||\mathbf{y} - \mathbf{X}\boldsymbol{\theta}||_2^2 + \frac{\lambda}{2} ||\boldsymbol{\theta}||_2^2$$

Least squares
empirical risk

Regularisation or
“ridge” penalty

Ridge regression

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \frac{1}{2n} ||\mathbf{y} - \mathbf{X}\boldsymbol{\theta}||_2^2 + \frac{\lambda}{2} ||\boldsymbol{\theta}||_2^2$$

Ridge regression

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

Ridge regression

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

Remarks:

- The regularised empirical risk is a **strongly convex function** of $\boldsymbol{\theta} \in \mathbb{R}^d$

$$\nabla_{\boldsymbol{\theta}} \hat{\mathcal{R}}_n^\lambda(\boldsymbol{\theta}) = -\frac{1}{n} \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) + \lambda \boldsymbol{\theta}$$

$$\nabla_{\boldsymbol{\theta}}^2 \hat{\mathcal{R}}_n^\lambda(\boldsymbol{\theta}) = \frac{1}{n} \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d \succ 0$$

$$(= \hat{\boldsymbol{\Sigma}}_n + \lambda \mathbf{I}_n)$$

Ridge regression

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

Remarks:

- The regularised empirical risk is a **strongly convex function** of $\boldsymbol{\theta} \in \mathbb{R}^d$

$$\nabla_{\boldsymbol{\theta}} \hat{\mathcal{R}}_n^\lambda(\boldsymbol{\theta}) = -\frac{1}{n} \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) + \lambda \boldsymbol{\theta}$$

$$\nabla_{\boldsymbol{\theta}}^2 \hat{\mathcal{R}}_n^\lambda(\boldsymbol{\theta}) = \frac{1}{n} \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d \succ 0$$

$$(= \hat{\boldsymbol{\Sigma}}_n + \lambda \mathbf{I}_n)$$

In other words, **minimiser** always **exist** and is **unique**.

Ridge regression

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

$$\nabla_{\boldsymbol{\theta}} \hat{\mathcal{R}}_n^\lambda(\boldsymbol{\theta}) = -\frac{1}{n} \boldsymbol{X}^\top (\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\theta}) + \lambda \boldsymbol{\theta} \stackrel{!}{=} \mathbf{0}$$

Ridge regression

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

$$\nabla_{\boldsymbol{\theta}} \hat{\mathcal{R}}_n^\lambda(\boldsymbol{\theta}) = -\frac{1}{n} \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) + \lambda \boldsymbol{\theta} \stackrel{!}{=} \mathbf{0}$$

$$\Leftrightarrow$$

$$\left(\frac{1}{n} \mathbf{X}^\top \mathbf{X} \boldsymbol{\theta} + \lambda \mathbf{I}_d \right) \boldsymbol{\theta} \stackrel{!}{=} \frac{1}{n} \mathbf{X}^\top \mathbf{y}$$

Ridge regression

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

$$\nabla_{\boldsymbol{\theta}} \hat{\mathcal{R}}_n^\lambda(\boldsymbol{\theta}) = -\frac{1}{n} \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) + \lambda \boldsymbol{\theta} \stackrel{!}{=} \mathbf{0}$$

$$\Leftrightarrow$$

$$\left(\frac{1}{n} \mathbf{X}^\top \mathbf{X} \boldsymbol{\theta} + \lambda \mathbf{I}_d \right) \boldsymbol{\theta} \stackrel{!}{=} \frac{1}{n} \mathbf{X}^\top \mathbf{y}$$

$$\Leftrightarrow$$

$$\hat{\boldsymbol{\theta}}_\lambda(\mathbf{X}, \mathbf{y}) = \frac{1}{n} \left(\frac{1}{n} \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d \right)^{-1} \mathbf{X}^\top \mathbf{y}$$

Ridge regression

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

The unique solution is given by:

$$\hat{\boldsymbol{\theta}}_\lambda(\mathbf{X}, \mathbf{y}) = \frac{1}{n} \left(\frac{1}{n} \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d \right)^{-1} \mathbf{X}^\top \mathbf{y}$$

Ridge regression

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

The unique solution is given by:

$$\hat{\boldsymbol{\theta}}_\lambda(\mathbf{X}, \mathbf{y}) = \frac{1}{n} \left(\frac{1}{n} \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d \right)^{-1} \mathbf{X}^\top \mathbf{y}$$



For $\lambda \rightarrow 0^+$, $\hat{\boldsymbol{\theta}}_\lambda \rightarrow \hat{\boldsymbol{\theta}}_{\text{OLS}}$

Ridge regression

$$\hat{\boldsymbol{\theta}}_{\lambda}(X, \mathbf{y}) = \frac{1}{n} \left(\frac{1}{n} X^{\top} X + \lambda \mathbf{I}_d \right)^{-1} X^{\top} \mathbf{y}$$

Remarks: • As before, consider s.v.d. of $X = \sum_{j=1}^{\text{rank}(X)} \sigma_j \mathbf{u}_j \mathbf{v}_j^{\top}$

Ridge regression

$$\hat{\boldsymbol{\theta}}_{\lambda}(X, \mathbf{y}) = \frac{1}{n} \left(\frac{1}{n} X^{\top} X + \lambda \mathbf{I}_d \right)^{-1} X^{\top} \mathbf{y}$$

Remarks: • As before, consider s.v.d. of $X = \sum_{j=1}^{\text{rank}(X)} \sigma_j \mathbf{u}_j \mathbf{v}_j^{\top}$

$$\hat{\boldsymbol{\theta}}_{\lambda}(X, \mathbf{y}) = \sum_{j=1}^{\text{rank}(X)} \frac{\sigma_j}{\sigma_j^2 + n\lambda} \langle \mathbf{u}_j, \mathbf{y} \rangle \mathbf{v}_j$$

Ridge regression

$$\hat{\theta}_{\lambda}(X, y) = \frac{1}{n} \left(\frac{1}{n} X^{\top} X + \lambda I_d \right)^{-1} X^{\top} y$$

Remarks: • As before, consider s.v.d. of $X = \sum_{j=1}^{\text{rank}(X)} \sigma_j \mathbf{u}_j \mathbf{v}_j^{\top}$

$$\hat{\theta}_{\lambda}(X, y) = \sum_{j=1}^{\text{rank}(X)} \frac{\sigma_j}{\sigma_j^2 + n\lambda} \langle \mathbf{u}_j, y \rangle \mathbf{v}_j$$

Ridge performs **shrinkage**:
small s.v.s are suppressed!

