



Statistical Learning II

Lecture 8 - Ridge regression (continued)

Bruno Loureiro
@ CSD, DI-ENS & CNRS

brloureiro@gmail.com

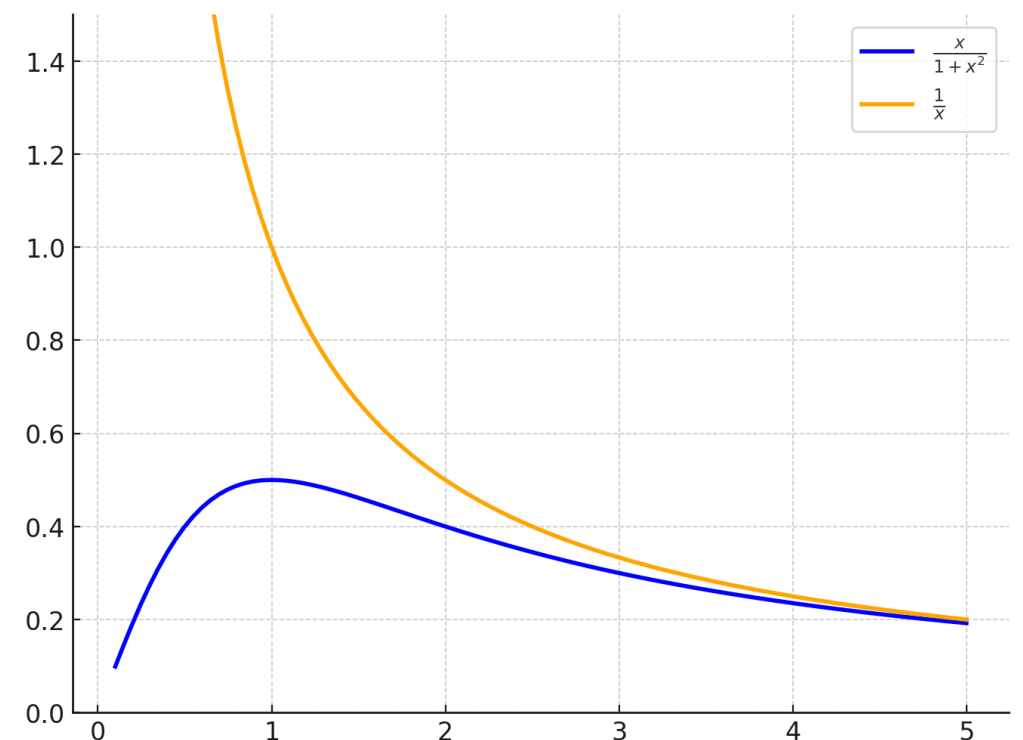
Ridge regression

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

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

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

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



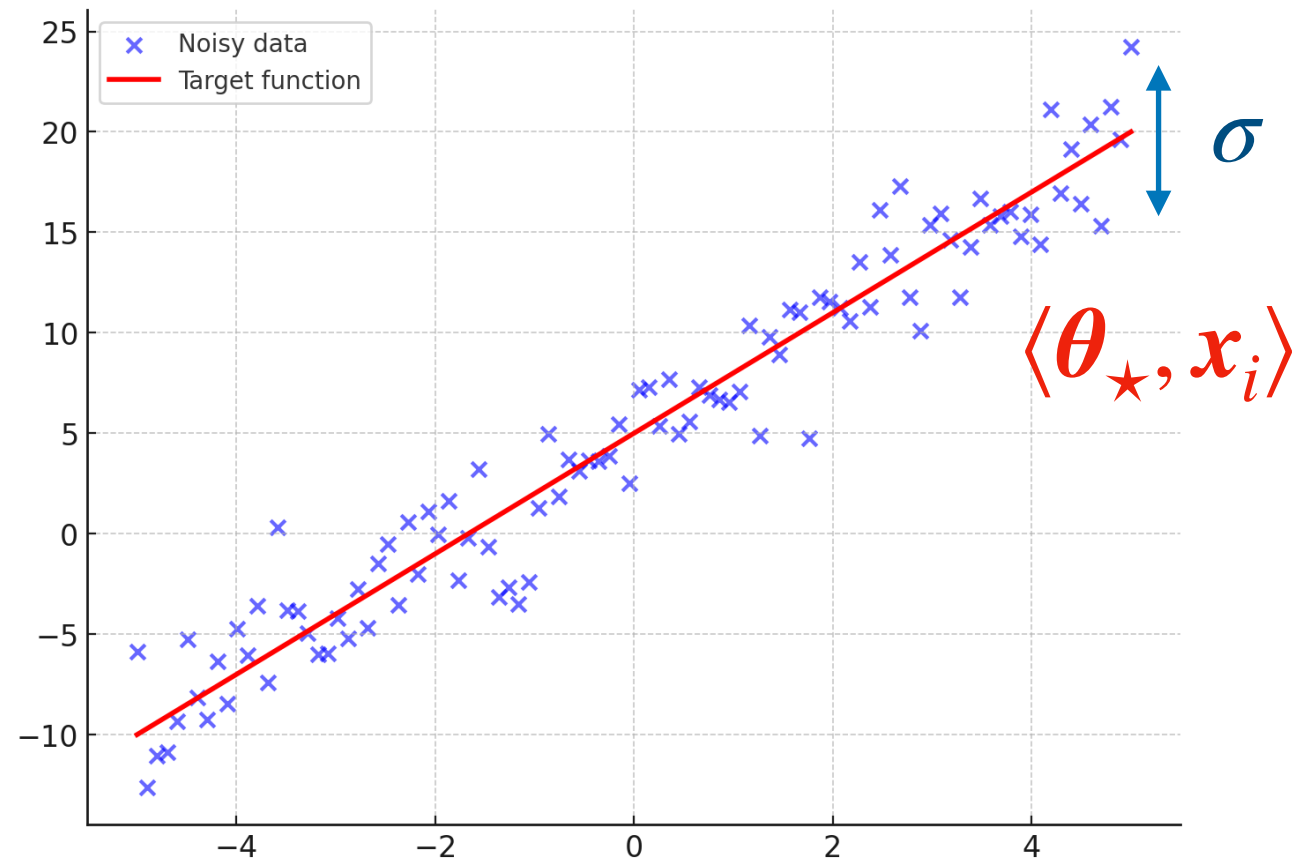
Statistical analysis of ridge regression

Fixed design assumption

As we did for the OLS, now let's assume:

$$y_i = \langle \theta_\star, x_i \rangle + \varepsilon_i$$

- With:
- Fixed $\theta_\star \in \mathbb{R}^d$ and $x_i \in \mathbb{R}^d$ “fixed design”
 - $\mathbb{E}[\varepsilon_i] = 0$ and $\mathbb{E}[\varepsilon_i^2] = \sigma^2 < \infty$



Decomposition of ridge

Given a batch of data sampled from this model:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta}_{\star} + \boldsymbol{\varepsilon} \in \mathbb{R}^n$$

The ridge estimator is given by:

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

Decomposition of ridge

Given a batch of data sampled from this model:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta}_\star + \boldsymbol{\varepsilon} \in \mathbb{R}^n$$

The ridge estimator is given by:

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

Decomposition of ridge

Given a batch of data sampled from this model:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta}_{\star} + \boldsymbol{\varepsilon} \in \mathbb{R}^n$$

The ridge estimator is given by:

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

Decomposition of ridge

Given a batch of data sampled from this model:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta}_{\star} + \boldsymbol{\varepsilon} \in \mathbb{R}^n$$

The ridge estimator is given by:

$$\begin{aligned}\hat{\boldsymbol{\theta}}_{\lambda}(\mathbf{X}, \mathbf{y}) &= \frac{1}{n} \left(\hat{\boldsymbol{\Sigma}}_n + \lambda \mathbf{I}_d \right)^{-1} \mathbf{X}^{\top} \mathbf{y} = \frac{1}{n} \left(\hat{\boldsymbol{\Sigma}}_n + \lambda \mathbf{I}_d \right)^{-1} \mathbf{X}^{\top} (\mathbf{X}\boldsymbol{\theta}_{\star} + \boldsymbol{\varepsilon}) \\ &= \left(\hat{\boldsymbol{\Sigma}}_n + \lambda \mathbf{I}_d \right)^{-1} \hat{\boldsymbol{\Sigma}}_n \boldsymbol{\theta}_{\star} + \frac{1}{n} \left(\hat{\boldsymbol{\Sigma}}_n + \lambda \mathbf{I}_d \right)^{-1} \mathbf{X}^{\top} \boldsymbol{\varepsilon}\end{aligned}$$

Decomposition of ridge

Given a batch of data sampled from this model:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta}_{\star} + \boldsymbol{\varepsilon} \in \mathbb{R}^n$$

The ridge estimator is given by:

$$\begin{aligned}\hat{\boldsymbol{\theta}}_{\lambda}(\mathbf{X}, \mathbf{y}) &= \frac{1}{n} \left(\hat{\boldsymbol{\Sigma}}_n + \lambda \mathbf{I}_d \right)^{-1} \mathbf{X}^{\top} \mathbf{y} = \frac{1}{n} \left(\hat{\boldsymbol{\Sigma}}_n + \lambda \mathbf{I}_d \right)^{-1} \mathbf{X}^{\top} (\mathbf{X}\boldsymbol{\theta}_{\star} + \boldsymbol{\varepsilon}) \\ &= \left(\hat{\boldsymbol{\Sigma}}_n + \lambda \mathbf{I}_d \right)^{-1} \hat{\boldsymbol{\Sigma}}_n \boldsymbol{\theta}_{\star} + \frac{1}{n} \left(\hat{\boldsymbol{\Sigma}}_n + \lambda \mathbf{I}_d \right)^{-1} \mathbf{X}^{\top} \boldsymbol{\varepsilon} \\ &= \boldsymbol{\theta}_{\star} - \lambda \left(\hat{\boldsymbol{\Sigma}}_n + \lambda \mathbf{I}_d \right)^{-1} \boldsymbol{\theta}_{\star} + \frac{1}{n} \left(\hat{\boldsymbol{\Sigma}}_n + \lambda \mathbf{I}_d \right)^{-1} \mathbf{X}^{\top} \boldsymbol{\varepsilon}\end{aligned}$$

Decomposition of ridge

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

“signal”

“noise”

Decomposition of ridge

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

“signal”

“noise”

In particular:

- Bias: $\mathbb{E}_{\epsilon} \left[\hat{\theta}_{\lambda}(X, y) \right] = \theta_{\star} - \lambda \left(\hat{\Sigma}_n + \lambda I_d \right)^{-1} \theta_{\star}$

Decomposition of ridge

$$\hat{\boldsymbol{\theta}}_{\lambda}(X, \mathbf{y}) = \underbrace{\boldsymbol{\theta}_{\star}}_{\text{"signal"}} - \lambda \left(\hat{\boldsymbol{\Sigma}}_n + \lambda \mathbf{I}_d \right)^{-1} \underbrace{\boldsymbol{\theta}_{\star}}_{\text{"signal"}} + \underbrace{\frac{1}{n} \left(\hat{\boldsymbol{\Sigma}}_n + \lambda \mathbf{I}_d \right)^{-1} X^{\top} \boldsymbol{\varepsilon}}_{\text{"noise"}}$$

“signal”

“noise”

In particular:

- Bias: $\mathbb{E}_{\boldsymbol{\varepsilon}} \left[\hat{\boldsymbol{\theta}}_{\lambda}(X, \mathbf{y}) \right] = \boldsymbol{\theta}_{\star} - \lambda \left(\hat{\boldsymbol{\Sigma}}_n + \lambda \mathbf{I}_d \right)^{-1} \boldsymbol{\theta}_{\star}$
- Variance: $\text{Var}_{\boldsymbol{\varepsilon}} \left[\hat{\boldsymbol{\theta}}_{\lambda}(X, \mathbf{y}) \right] = \frac{\sigma^2}{n} \left(\hat{\boldsymbol{\Sigma}}_n + \lambda \mathbf{I}_d \right)^{-2} \hat{\boldsymbol{\Sigma}}_n$

Decomposition of ridge

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

“signal”

“noise”

In particular:

- Bias: $\mathbb{E}_{\epsilon} \left[\hat{\theta}_{\lambda}(X, y) \right] = \theta_{\star} - \lambda \left(\hat{\Sigma}_n + \lambda I_d \right)^{-1} \theta_{\star}$
- Variance: $\text{Var}_{\epsilon} \left[\hat{\theta}_{\lambda}(X, y) \right] = \frac{\sigma^2}{n} \left(\hat{\Sigma}_n + \lambda I_d \right)^{-2} \hat{\Sigma}_n$



- Ridge is a **biased** estimator.
- Regularisation shrinks both signal and noise

Risk of ridge

Recall that in Lecture 5 we have shown that for any $\boldsymbol{\theta} \in \mathbb{R}^d$:

$$\mathcal{R}(\boldsymbol{\theta}) - \sigma^2 = (\boldsymbol{\theta} - \boldsymbol{\theta}_\star)^\top \hat{\boldsymbol{\Sigma}}_n (\boldsymbol{\theta} - \boldsymbol{\theta}_\star)$$

Risk of ridge

Recall that in Lecture 5 we have shown that for any $\boldsymbol{\theta} \in \mathbb{R}^d$:

$$\mathcal{R}(\boldsymbol{\theta}) - \sigma^2 = (\boldsymbol{\theta} - \boldsymbol{\theta}_\star)^\top \hat{\boldsymbol{\Sigma}}_n (\boldsymbol{\theta} - \boldsymbol{\theta}_\star)$$

Therefore, inserting the solution $\hat{\boldsymbol{\theta}}_\lambda(X, \mathbf{y})$:

$$\begin{aligned} \mathcal{R}(\hat{\boldsymbol{\theta}}_\lambda) - \sigma^2 &= \lambda^2 \boldsymbol{\theta}_\star^\top (\hat{\boldsymbol{\Sigma}}_n + \lambda \mathbf{I}_d)^{-2} \boldsymbol{\theta}_\star \\ &\quad + \frac{1}{n^2} \boldsymbol{\varepsilon}^\top X (\hat{\boldsymbol{\Sigma}}_n + \lambda \mathbf{I}_d)^{-1} X^\top X (\hat{\boldsymbol{\Sigma}}_n + \lambda \mathbf{I}_d)^{-1} X^\top \boldsymbol{\varepsilon} \\ &\quad - \frac{\lambda}{n} \boldsymbol{\varepsilon}^\top X (\hat{\boldsymbol{\Sigma}}_n + \lambda \mathbf{I}_d)^{-2} \boldsymbol{\theta}_\star \end{aligned}$$

Risk of ridge

Taking the expectation with respect to the noise:

$$\mathbb{E}_{\boldsymbol{\varepsilon}}[\mathcal{R}(\hat{\boldsymbol{\theta}}_{\lambda})] - \sigma^2 = \lambda^2 \boldsymbol{\theta}_{\star}^{\top} (\hat{\boldsymbol{\Sigma}}_n + \lambda \mathbf{I}_d)^{-2} \hat{\boldsymbol{\Sigma}}_n \boldsymbol{\theta}_{\star} + \frac{\sigma^2}{n} \text{Tr} \hat{\boldsymbol{\Sigma}}_n^2 (\hat{\boldsymbol{\Sigma}}_n + \lambda \mathbf{I}_d)^{-2}$$

Risk of ridge

Taking the expectation with respect to the noise:

$$\mathbb{E}_{\epsilon}[\mathcal{R}(\hat{\boldsymbol{\theta}}_{\lambda})] - \sigma^2 = \lambda^2 \boldsymbol{\theta}_{\star}^{\top} (\hat{\boldsymbol{\Sigma}}_n + \lambda \mathbf{I}_d)^{-2} \hat{\boldsymbol{\Sigma}}_n \boldsymbol{\theta}_{\star} + \frac{\sigma^2}{n} \text{Tr } \hat{\boldsymbol{\Sigma}}_n^2 (\hat{\boldsymbol{\Sigma}}_n + \lambda \mathbf{I}_d)^{-2}$$

Alternatively, we can also write in terms of a bias-variance decomposition of the risk:

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

Where:

$$\mathcal{B} = \lambda^2 \boldsymbol{\theta}_{\star}^{\top} (\hat{\boldsymbol{\Sigma}}_n + \lambda \mathbf{I}_d)^{-2} \hat{\boldsymbol{\Sigma}}_n \boldsymbol{\theta}_{\star} \quad \mathcal{V} = \frac{\sigma^2}{n} \text{Tr } \hat{\boldsymbol{\Sigma}}_n^2 (\hat{\boldsymbol{\Sigma}}_n + \lambda \mathbf{I}_d)^{-2}$$

Risk of ridge

Considering the SVD of $X = \sum_{k=1}^{\text{rank}(X)} \lambda_k \mathbf{u}_k \mathbf{v}_k^\top$, we can also write:

$$\mathcal{B} = \sum_{k=1}^{\text{rank}(X)} \frac{(n\lambda)^2 \lambda_k \langle \mathbf{v}_k, \boldsymbol{\theta}_\star \rangle^2}{(\lambda_k + n\lambda)^2} \quad \mathcal{V} = \sum_{k=1}^{\text{rank}(X)} \frac{\sigma^2 \lambda_k^2}{(\lambda_k + n\lambda)^2}$$

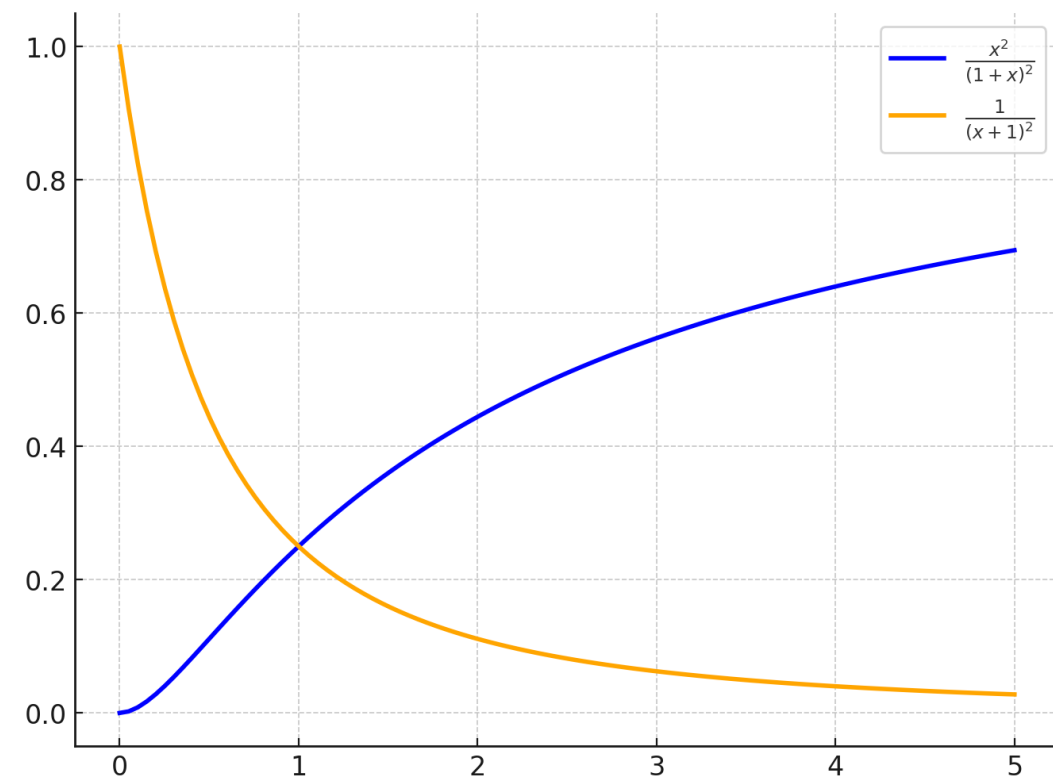
Risk of ridge

Considering the SVD of $X = \sum_{k=1}^{\text{rank}(X)} \lambda_k \mathbf{u}_k \mathbf{v}_k^\top$, we can also write:

$$\mathcal{B} = \sum_{k=1}^{\text{rank}(X)} \frac{(n\lambda)^2 \lambda_k \langle \mathbf{v}_k, \boldsymbol{\theta}_\star \rangle^2}{(\lambda_k + n\lambda)^2} \quad \mathcal{V} = \sum_{k=1}^{\text{rank}(X)} \frac{\sigma^2 \lambda_k^2}{(\lambda_k + n\lambda)^2}$$

Remarks:

- For $\lambda \rightarrow 0^+$, we get the OLS excess risk
- $\mathcal{B}(\lambda)$ is an increasing function of λ
- $\mathcal{V}(\lambda)$ is a decreasing function of λ



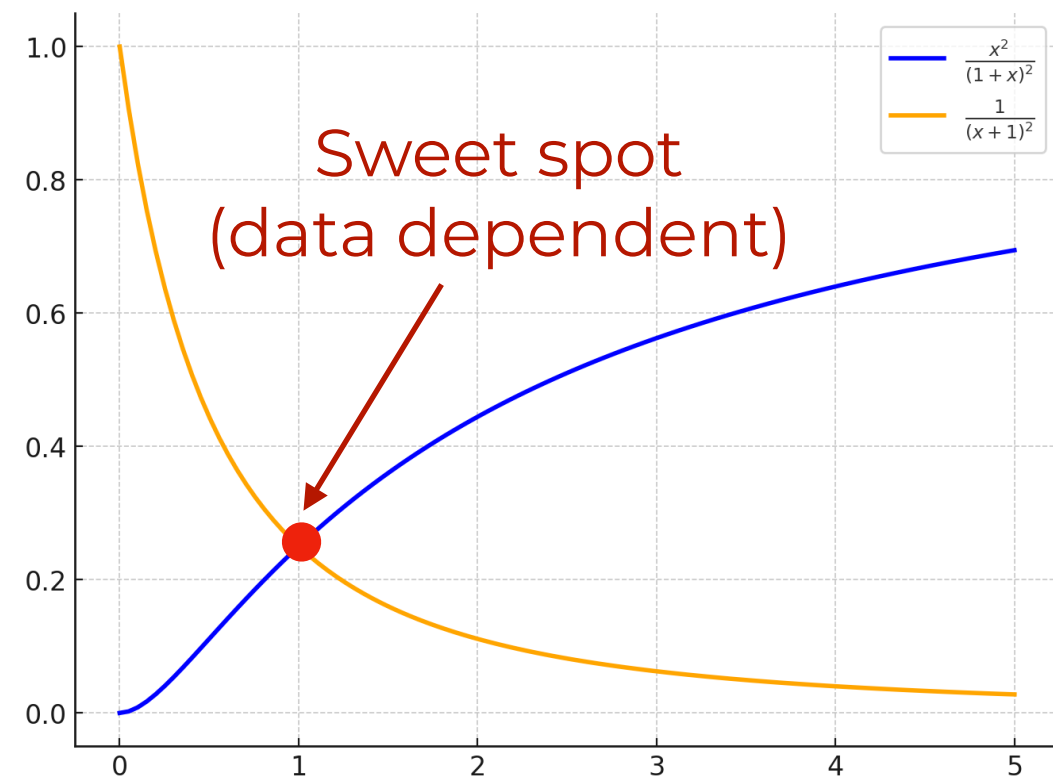
Risk of ridge

Considering the SVD of $X = \sum_{k=1}^{\text{rank}(X)} \lambda_k \mathbf{u}_k \mathbf{v}_k^\top$, we can also write:

$$\mathcal{B} = \sum_{k=1}^{\text{rank}(X)} \frac{(n\lambda)^2 \lambda_k \langle \mathbf{v}_k, \boldsymbol{\theta}_\star \rangle^2}{(\lambda_k + n\lambda)^2} \quad \mathcal{V} = \sum_{k=1}^{\text{rank}(X)} \frac{\sigma^2 \lambda_k^2}{(\lambda_k + n\lambda)^2}$$

Remarks:

- For $\lambda \rightarrow 0^+$, we get the OLS excess risk
- $\mathcal{B}(\lambda)$ is an increasing function of λ
- $\mathcal{V}(\lambda)$ is a decreasing function of λ



Interpretation of variance

Let $\mathbf{A} \in \mathbb{R}^{d \times d}$ be a positive definite matrix with decreasing eigenvalues $\text{spec}(\mathbf{A}) = \{\lambda_k : k = 1, \dots, d\}$. Define the cumulative:

$$\phi(\lambda) = \#\{k : \lambda_k > \lambda\}$$

“Count eigenvalues bigger than λ ”

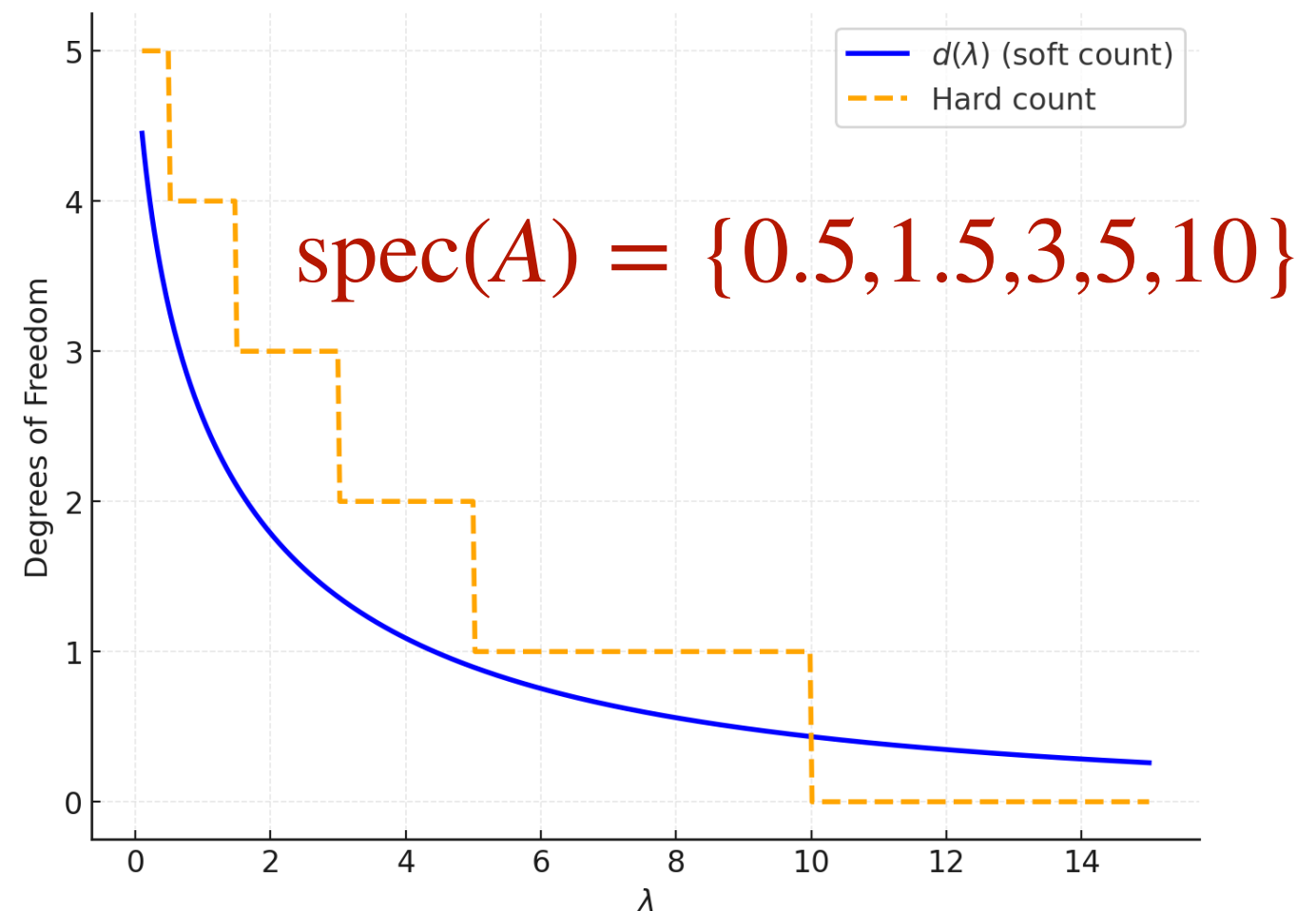
Interpretation of variance

Let $A \in \mathbb{R}^{d \times d}$ be a positive definite matrix with decreasing eigenvalues $\text{spec}(A) = \{\lambda_k : k = 1, \dots, d\}$. Define the cumulative:

$$\phi(\lambda) = \#\{k : \lambda_k > \lambda\} \quad \text{“Count eigenvalues bigger than } \lambda \text{”}$$

The variance of the ridge risk can be seen as a soft version:

$$\text{df}_2(\lambda) = \sum_{k=1}^d \frac{\lambda_k^2}{(\lambda_k + \lambda)^2}$$



Interpretation of variance

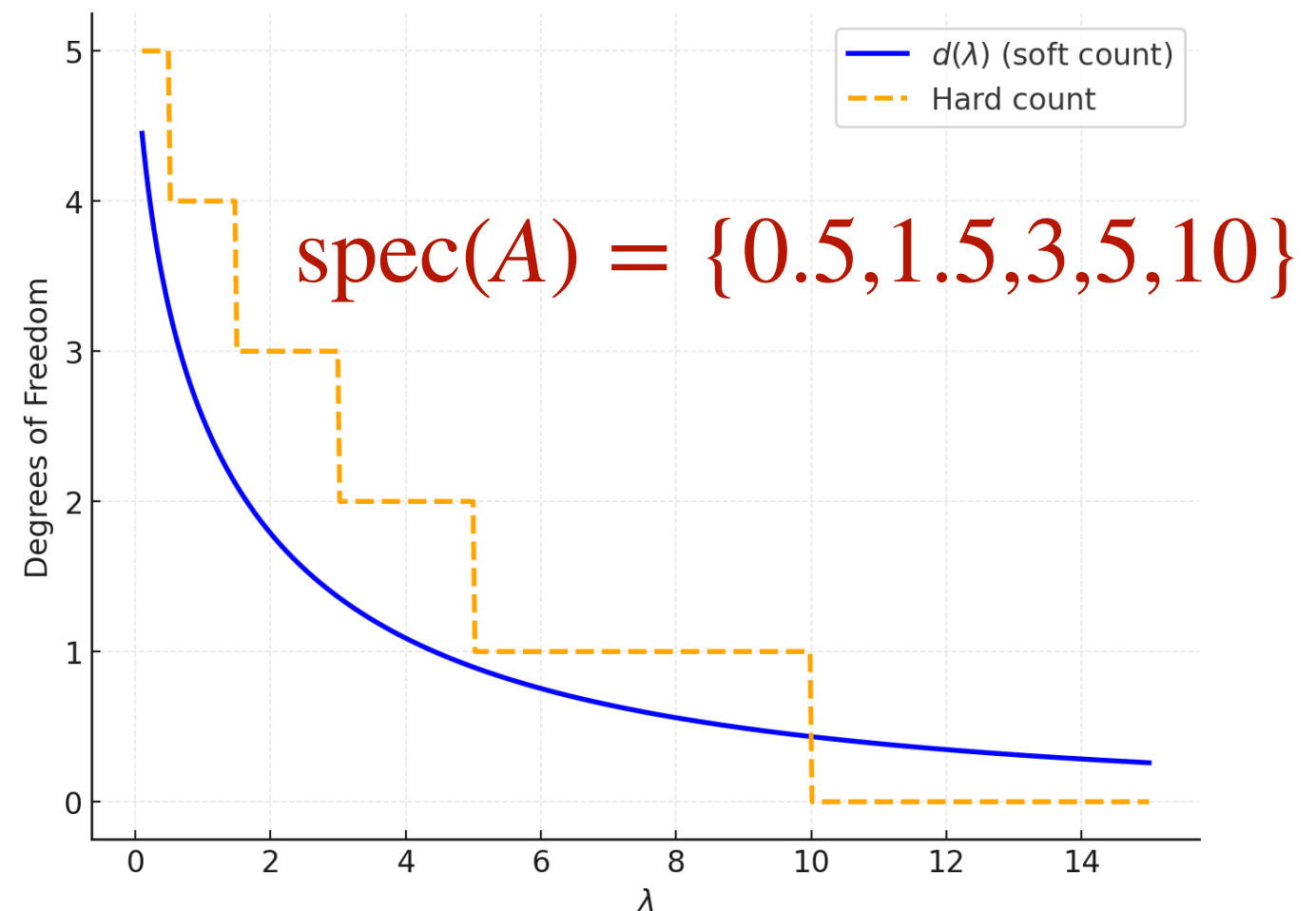
Let $A \in \mathbb{R}^{d \times d}$ be a positive definite matrix with decreasing eigenvalues $\text{spec}(A) = \{\lambda_k : k = 1, \dots, d\}$. Define the cumulative:

$$\phi(\lambda) = \#\{k : \lambda_k > \lambda\} \quad \text{“Count eigenvalues bigger than } \lambda \text{”}$$

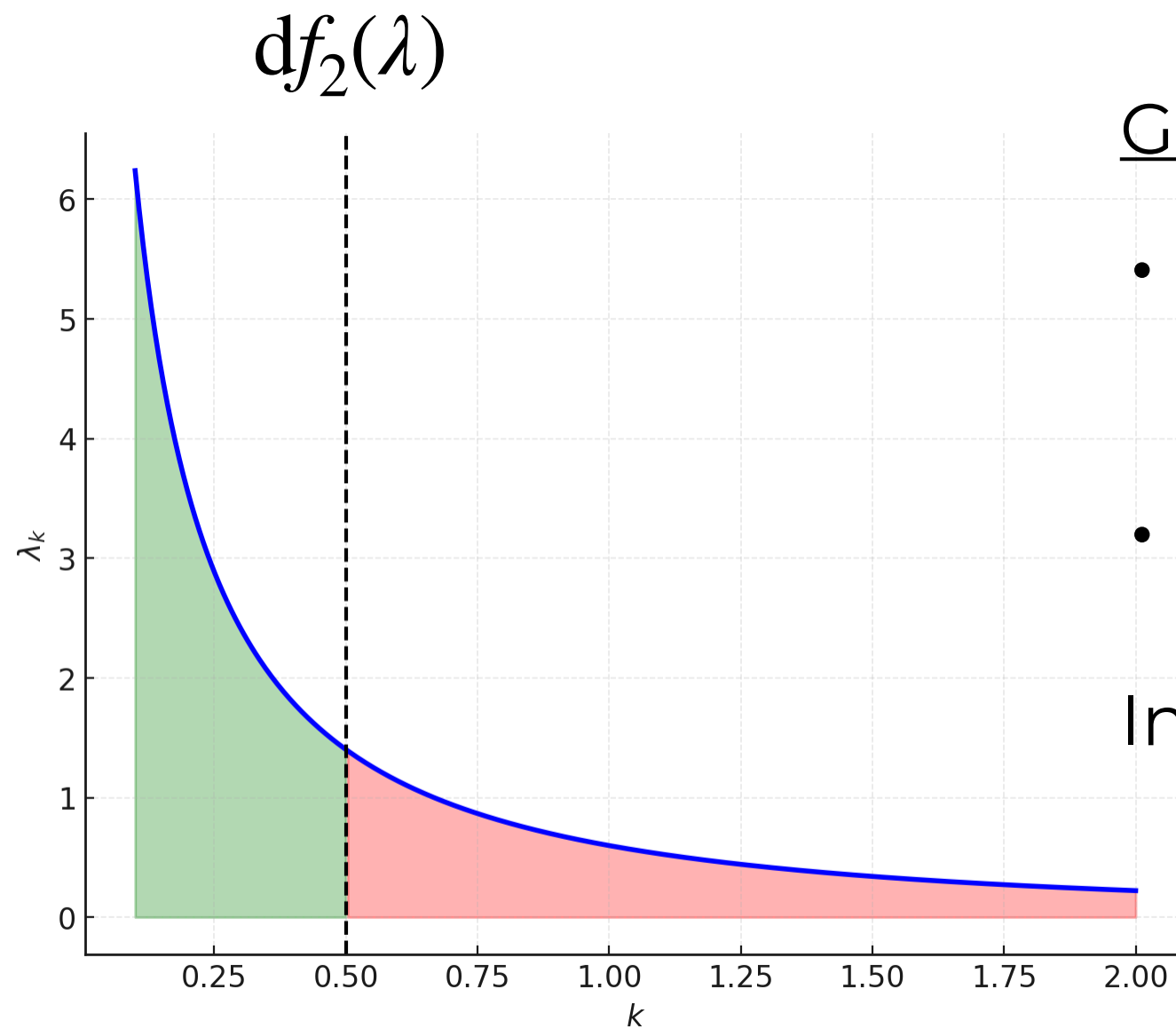
The variance of the ridge risk can be seen as a soft version:

$$\text{df}_2(\lambda) = \sum_{k=1}^d \frac{\lambda_k^2}{(\lambda_k + \lambda)^2}$$

- Fast decay: small λ
- Slow decay: large λ



Choosing regularisation



Goal: pick λ such that:

- directions in \mathbf{X} that better correlate with $\boldsymbol{\theta}_\star$ are retained
- Shrink remaining directions

In practice, **cross-validation**...

Low-frequency

High-frequency