# Statistical Learning II

## Lecture 9 - Ridge regression

_____

**Bruno Loureiro**

@ CSD, DI-ENS & CNRS

brloureiro@gmail.com

# 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:

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

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

💡 Key idea: penalise the norm.

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

# Ridge regression

Note the averaged norm of the OLS is given by:

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

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

💡 Key idea: penalise the norm.

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

Least squares empirical risk

Regularisation or "ridge" penalty

# Ridge regression

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

# Ridge regression

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

# Ridge regression

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

Remarks:

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

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

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

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

# Ridge regression

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

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

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

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

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

In other words, minimiser always exist and is unique.

# Ridge regression

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

$$\nabla_{\boldsymbol{\theta}} \hat{\mathscr{R}}_n^\lambda(\boldsymbol{\theta}) = -\frac{1}{n} \boldsymbol{X}^\top \left( \boldsymbol{y} - \boldsymbol{X}\boldsymbol{\theta} \right) + \lambda\boldsymbol{\theta} \overset{!}{=} \boldsymbol{0}$$

# Ridge regression

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

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

$$\Updownarrow$$

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

# Ridge regression

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

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

$$\Updownarrow$$

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

$$\Updownarrow$$

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

# Ridge regression

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

The unique solution is given by:

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

# Ridge regression

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

The unique solution is given by:

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

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

# Ridge regression

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

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

# Ridge regression

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

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

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

# Ridge regression

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

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

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

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:

$$y_i = \langle \boldsymbol{\theta}_\star, \boldsymbol{x}_i \rangle + \varepsilon_i$$

With:
- Fixed $\boldsymbol{\theta}_\star \in \mathbb{R}^d$ and $\boldsymbol{x}_i \in \mathbb{R}^d$   "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:

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

The ridge estimator is given by:

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

# Decomposition of ridge

Given a batch of data sampled from this model:

$$y = X\theta_\star + \varepsilon \in \mathbb{R}^n$$

The ridge estimator is given by:

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

# Decomposition of ridge

Given a batch of data sampled from this model:

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

The ridge estimator is given by:

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

# Decomposition of ridge

Given a batch of data sampled from this model:

$$y = X\theta_\star + \varepsilon \in \mathbb{R}^n$$

The ridge estimator is given by:

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

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

# Decomposition of ridge

Given a batch of data sampled from this model:

$$y = X\theta_\star + \varepsilon \in \mathbb{R}^n$$

The ridge estimator is given by:

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

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

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

# Decomposition of ridge

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

"signal"          "noise"

# Decomposition of ridge

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

"signal"                    "noise"

In particular:

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

# Decomposition of ridge

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

"signal"          "noise"

In particular:

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

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

# Decomposition of ridge

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

"signal"          "noise"

In particular:

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

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

⚠️
- 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 $\boldsymbol{\theta} \in \mathbb{R}^d$:

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

# 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 $\hat{\boldsymbol{\theta}}_\lambda(X, \boldsymbol{y})$:

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

$$+ \frac{1}{n^2} \boldsymbol{\varepsilon}^\top X (\hat{\boldsymbol{\Sigma}}_n + \lambda \boldsymbol{I}_d)^{-1} \hat{\boldsymbol{\Sigma}}_n (\hat{\boldsymbol{\Sigma}}_n + \lambda \boldsymbol{I}_d)^{-1} X^\top \boldsymbol{\varepsilon}$$

$$- \frac{\lambda}{n} \boldsymbol{\varepsilon}^\top X (\hat{\boldsymbol{\Sigma}}_n + \lambda \boldsymbol{I}_d)^{-2} \hat{\boldsymbol{\Sigma}}_n \boldsymbol{\theta}_\star$$

$$- \frac{\lambda}{n} \boldsymbol{\theta}_\star^\top (\hat{\boldsymbol{\Sigma}}_n + \lambda \boldsymbol{I}_d)^{-2} \hat{\boldsymbol{\Sigma}}_n X^\top \boldsymbol{\epsilon}$$

# Risk of ridge

Taking the expectation with respect to the noise:

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

# Risk of ridge

Taking the expectation with respect to the noise:

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

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

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

Where:

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

# Risk of ridge

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

$$\mathscr{B} = \frac{1}{n} \sum_{k=1}^{\text{rank}(X)} \frac{(n\lambda)^2 \sigma_k^2 \langle \boldsymbol{v}_k, \boldsymbol{\theta}_\star \rangle^2}{(\sigma_k^2 + n\lambda)^2} \qquad \mathscr{V} = \sigma^2 \sum_{k=1}^{\text{rank}(X)} \frac{\sigma_k^4}{(\sigma_k^2 + n\lambda)^2}$$
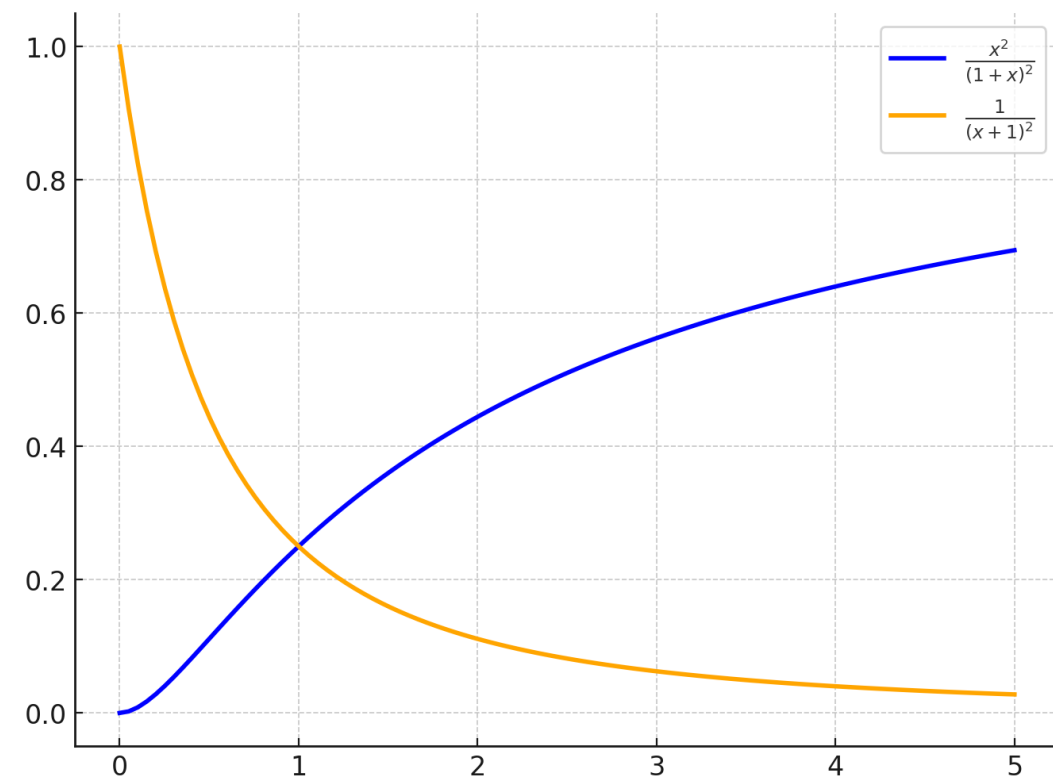
# Risk of ridge

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

$$\mathscr{B} = \frac{1}{n} \sum_{k=1}^{\text{rank}(X)} \frac{(n\lambda)^2 \sigma_k^2 \langle \boldsymbol{v}_k, \boldsymbol{\theta}_\star \rangle^2}{(\sigma_k^2 + n\lambda)^2} \qquad \mathscr{V} = \sigma^2 \sum_{k=1}^{\text{rank}(X)} \frac{\sigma_k^4}{(\sigma_k^2 + n\lambda)^2}$$

## Remarks:

- For $\lambda \to 0^+$, we get the OLS excess risk

- $\mathscr{B}(\lambda)$ is an increasing function of $\lambda$

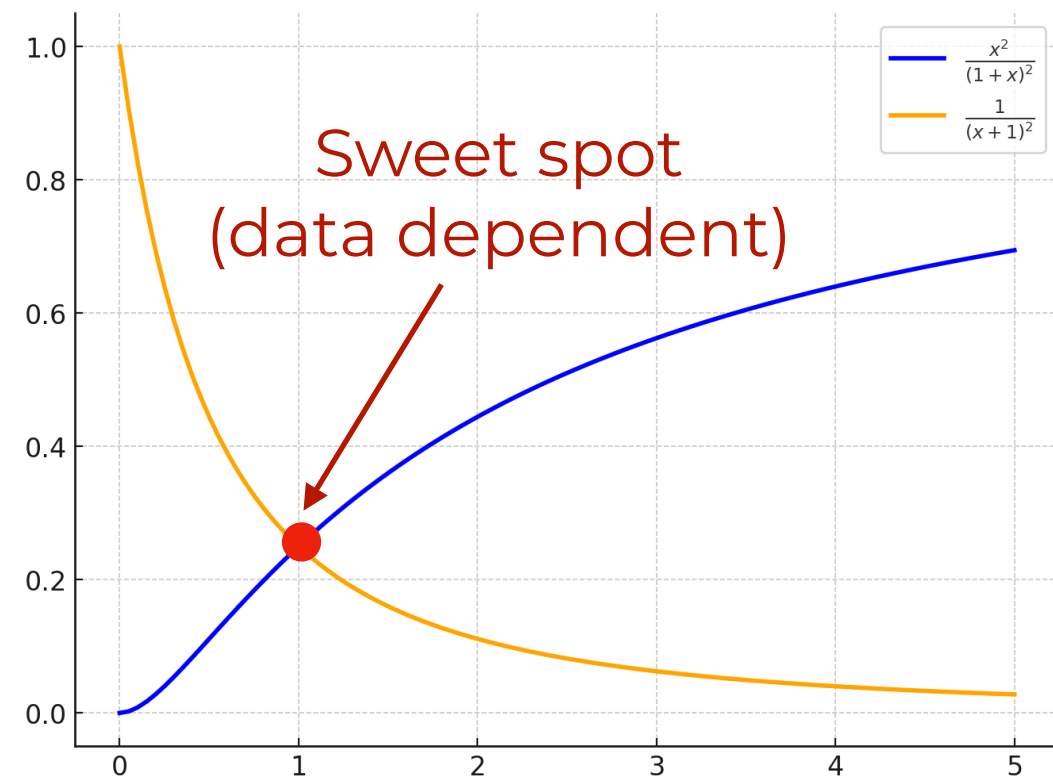- $\mathscr{V}(\lambda)$ is a decreasing function of $\lambda$

# Risk of ridge

Considering the SVD of $X = \sum\limits_{k=1}^{\text{rank}(X)} \sigma_k u_k v_k^{\top}$, we can also write:

$$\mathscr{B} = \frac{1}{n} \sum_{k=1}^{\text{rank}(X)} \frac{(n\lambda)^2 \sigma_k^2 \langle v_k, \theta_\star \rangle^2}{(\sigma_k^2 + n\lambda)^2} \quad \mathscr{V} = \sigma^2 \sum_{k=1}^{\text{rank}(X)} \frac{\sigma_k^4}{(\sigma_k^2 + n\lambda)^2}$$

## Remarks:

- For $\lambda \to 0^+$, we get the OLS excess risk

- $\mathscr{B}(\lambda)$ is an increasing function of $\lambda$

- $\mathscr{V}(\lambda)$ is a decreasing function of $\lambda$



Sweet spot (data dependent)

# Interpretation of variance

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

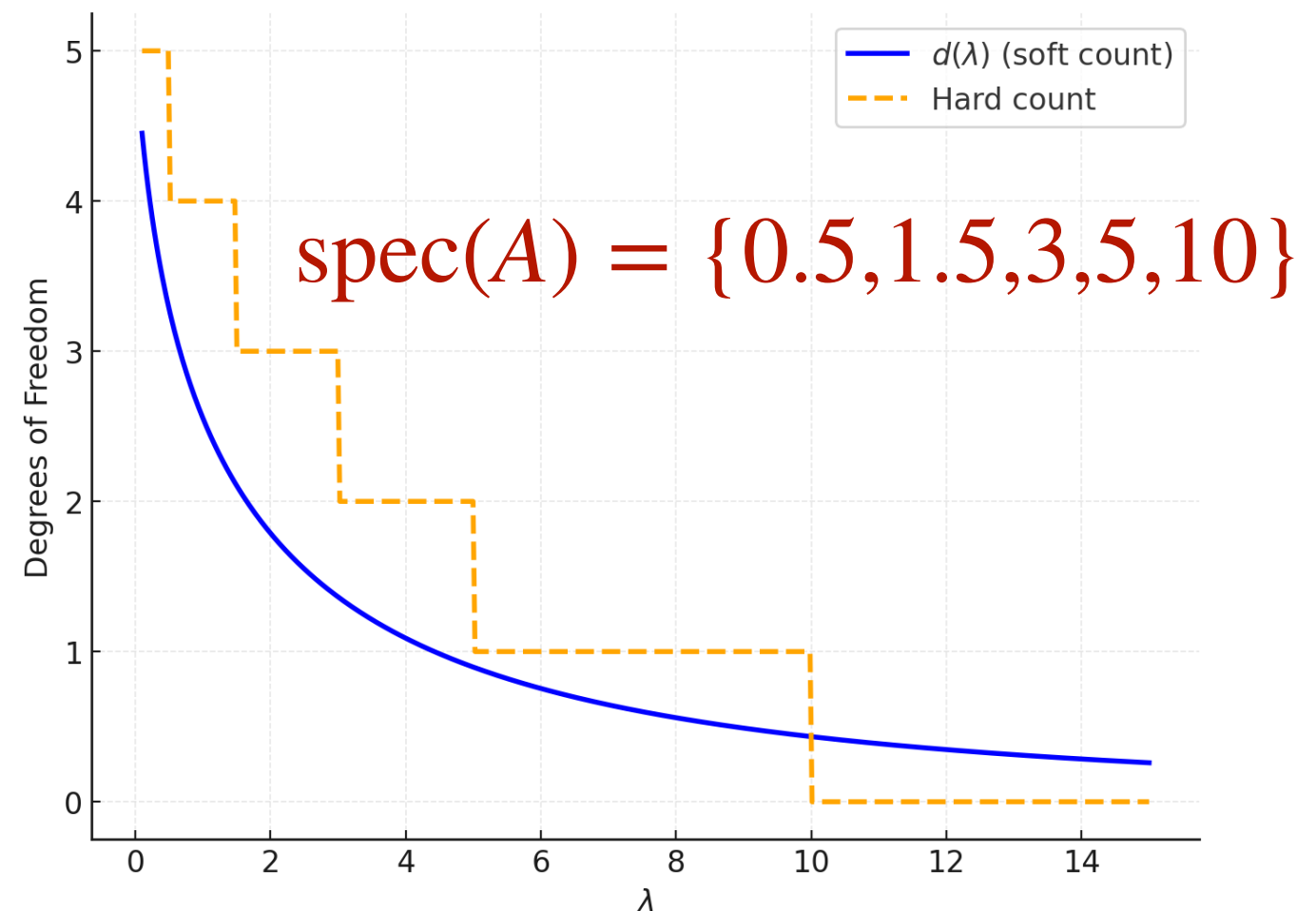$$\phi(\lambda) = \#\{k : \lambda_k > \lambda\}$$

"Count eigenvalues bigger than $\lambda$"

# Interpretation of variance

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

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

<span style="color:red">"Count eigenvalues bigger than $\lambda$"</span>

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

$$\mathrm{d}f_2(\lambda) = \sum_{k=1}^{d} \frac{\lambda_k^2}{(\lambda_k + \lambda)^2}$$



<span style="color:red">$\mathrm{spec}(A) = \{0.5, 1.5, 3, 5, 10\}$</span>
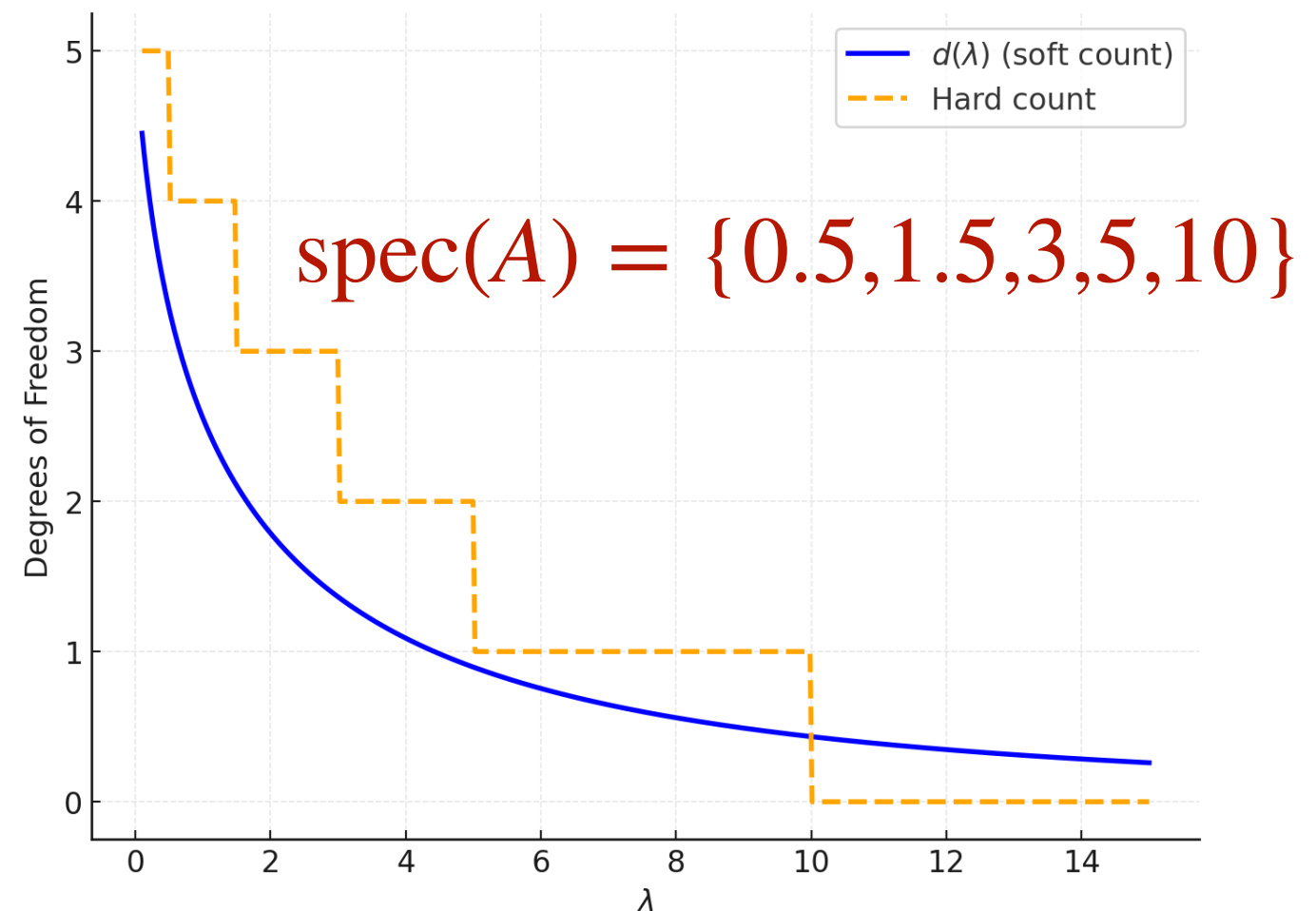
# Interpretation of variance

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

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

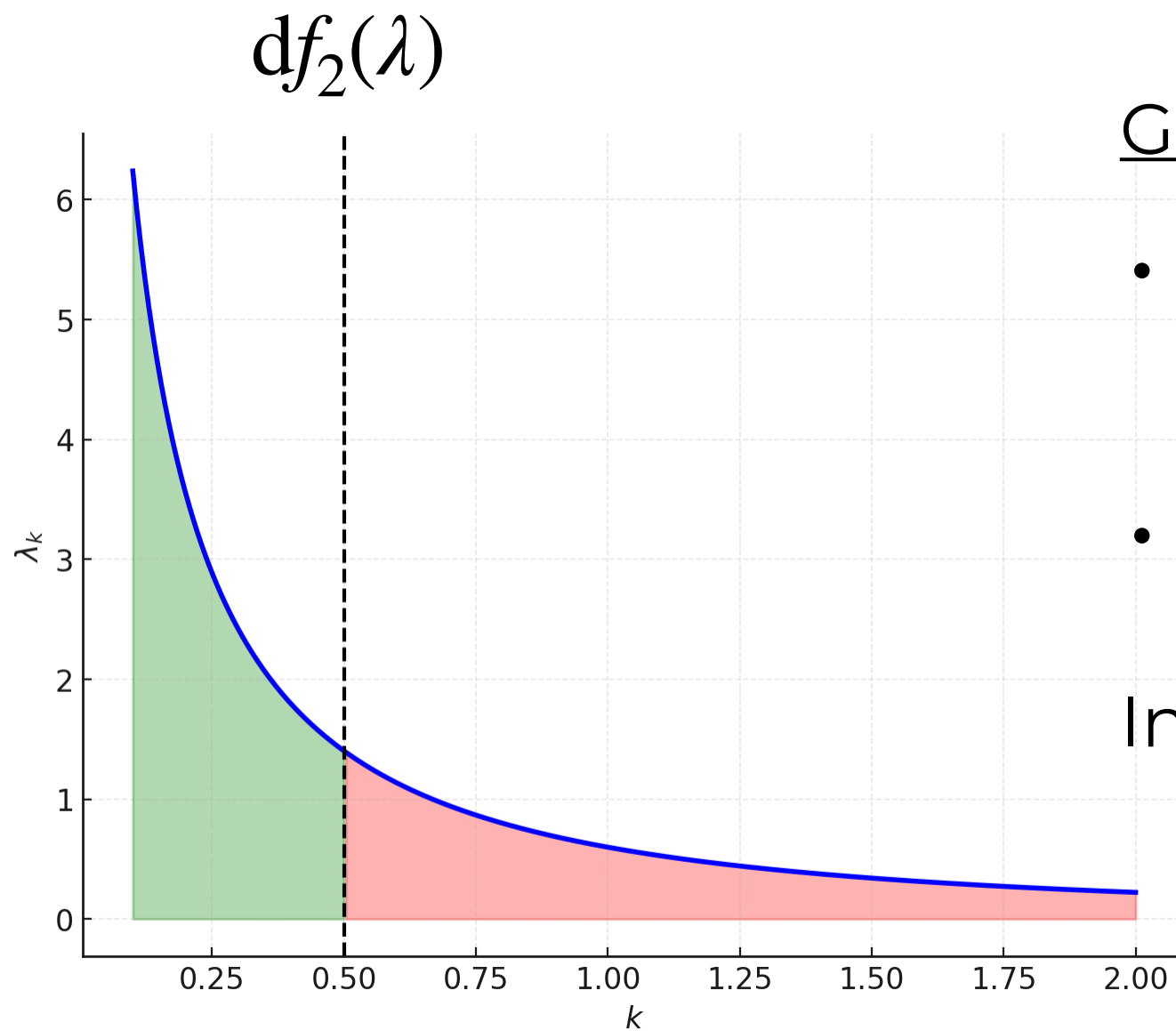"Count eigenvalues bigger than $\lambda$"

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

$$\mathrm{d}f_2(\lambda) = \sum_{k=1}^{d} \frac{\lambda_k^2}{(\lambda_k + \lambda)^2}$$



$$\mathrm{spec}(A) = \{0.5, 1.5, 3, 5, 10\}$$

- Fast decay: small $\lambda$
- Slow decay: large $\lambda$

# Choosing regularisation

$$\mathrm{d}f_2(\lambda)$$



Low-frequency    High-frequency

Goal: pick $\lambda$ such that:

- directions in $X$ that better correlate with $\theta_\star$ are retained

- Shrink remaining directions

In practice, cross-validation...