



Statistical Learning II

Lecture 7 - Ridge regression

Bruno Loureiro

@ CSD, DI-ENS & CNRS

brloureiro@gmail.com

Marvels and pitfalls of OLS

Recall that:

$$\hat{\boldsymbol{\theta}}_{OLS}(X, \mathbf{y}) = \boldsymbol{\theta}_{\star} + \frac{1}{n} \hat{\boldsymbol{\Sigma}}_{n}^{-1} X^{\mathsf{T}} \boldsymbol{\varepsilon}$$

$$= \boldsymbol{\theta}_{\star} + \sum_{j=1}^{d} \frac{1}{\sigma_{j}} \langle \boldsymbol{u}_{j}, \boldsymbol{\varepsilon} \rangle \boldsymbol{v}_{j}$$

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.

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.

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.

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

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

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

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

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

Remarks:

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

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

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

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

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

Remarks:

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

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

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

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

In other words, minimiser always exist and is unique.

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

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

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

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

$$\updownarrow$$

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

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

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

$$\updownarrow$$

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

$$\updownarrow$$

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

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

The unique solution is given by:

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

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

The unique solution is given by:

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



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

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

rank(X)

Remarks: • As before, consider s.v.d. of $X = \sum_{j=1}^{\infty} \sigma_j u_j v_j$

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

rank(X)

Remarks: • As before, consider s.v.d. of $X = \sum_{j=1}^{\infty} \sigma_j u_j v_j$

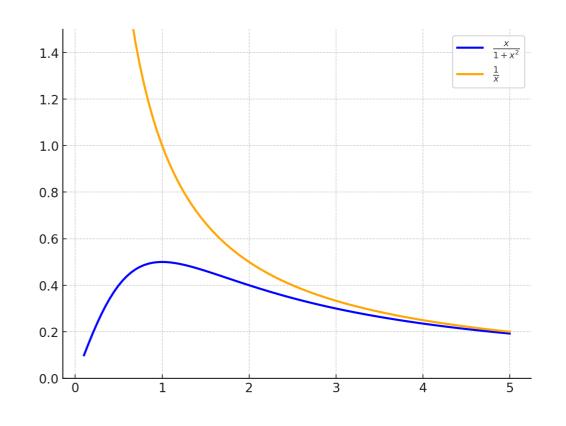
$$\hat{\boldsymbol{\theta}}_{\lambda}(\boldsymbol{X}, \boldsymbol{y}) = \sum_{j=1}^{\operatorname{rank}(\boldsymbol{X})} \frac{\sigma_{j}}{\sigma_{j}^{2} + n\lambda} \langle \boldsymbol{u}_{j}, \boldsymbol{y} \rangle \boldsymbol{v}_{j}$$

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

Remarks: • As before, consider s.v.d. of $X = \sum_{i=1}^{\infty} \sigma_i u_i v_i$

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

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



rank(X)