# Statistical Learning II

## Lecture 10 - BSS

_____

**Bruno Loureiro**

@ CSD, DI-ENS & CNRS

brloureiro@gmail.com

*DL3 IASO, Université Paris Dauphine-PSL*
*19.11.2025*

# Risk of ridge

Considering the SVD of $X = \sum\limits_{k=1}^{\text{rank}(X)} \sigma_k \boldsymbol{u}_k \boldsymbol{v}_k^{\top}$, 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 $\lambda$



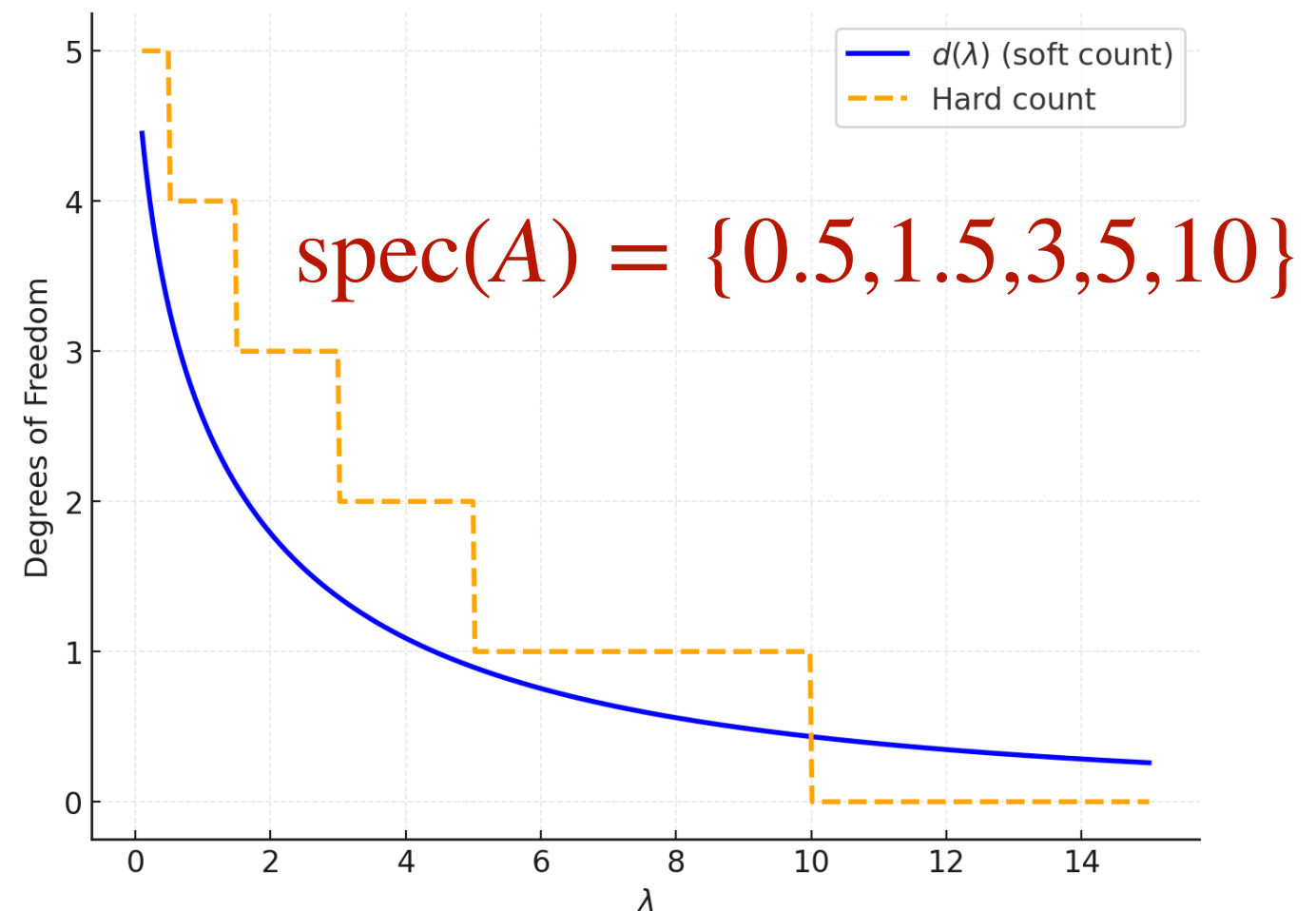Sweet spot
(data dependent)

# 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:

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

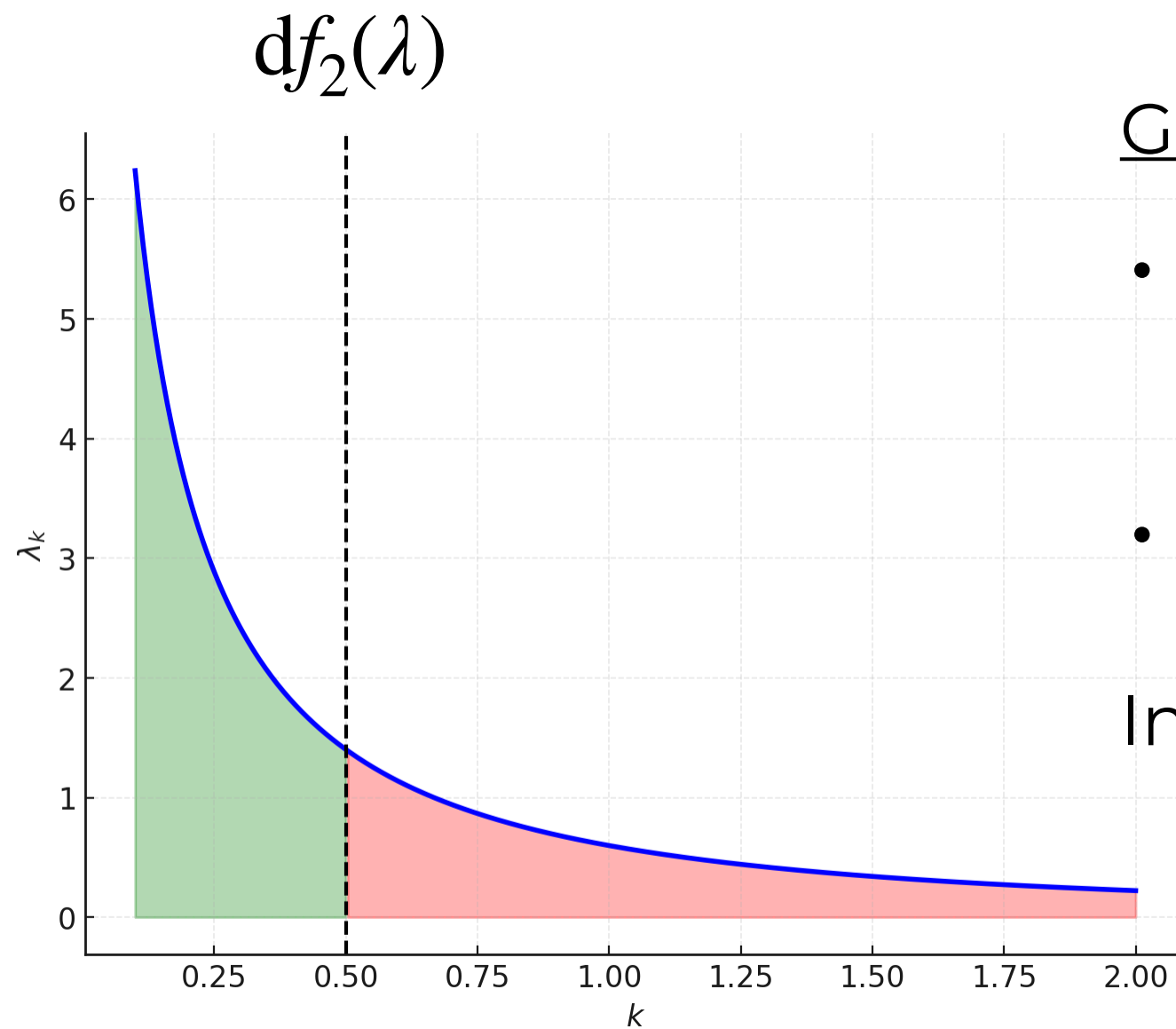"Count eigenvalues bigger than $\lambda$"

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

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

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



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

# 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...

# Best subset selection & the LASSO

# Pitfalls of ridge

The ridge estimation performs uniform shrinkage.

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

In other words: $\ell_2$ regularisation will control the overall norm $||\hat{\boldsymbol{\theta}}_{\lambda}||_2^2$ by reducing each entry equally

# Pitfalls of ridge

The ridge estimation performs uniform shrinkage.

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

In other words: $\ell_2$ regularisation will control the overall norm $||\hat{\boldsymbol{\theta}}_\lambda||_2^2$ by reducing each entry equally

- Good if $\boldsymbol{\theta}_\star$ is a dense vector
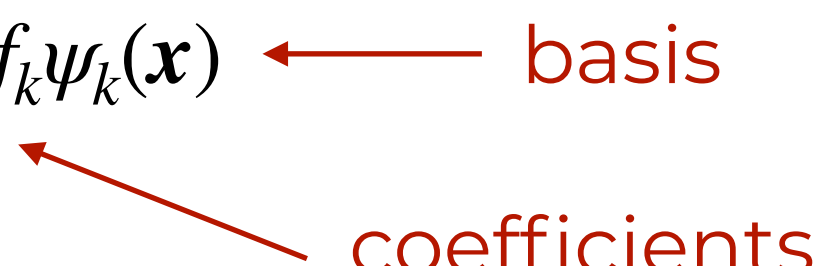
$$\boldsymbol{\theta}_{\star,j} \neq 0 \qquad i = 1, \cdots, d$$

- Bad if $\boldsymbol{\theta}_\star$ is a sparse vector

$$\boldsymbol{\theta}_{\star,j} = \begin{cases} 0 & j \in S \subset \{1,\ldots,d\} \\ \neq 0 & j \in \{1,\ldots,d\}\backslash S \end{cases}$$

$$\boldsymbol{\theta}_\star =$$

$$\boldsymbol{\theta}_\star =$$

# Sparsity is everywhere

Many signals of interest admit a sparse representation in a particular basis.

$$f(\boldsymbol{x}) = \sum_{k \geq 0} f_k \psi_k(\boldsymbol{x}) \quad \longleftarrow \text{basis}$$

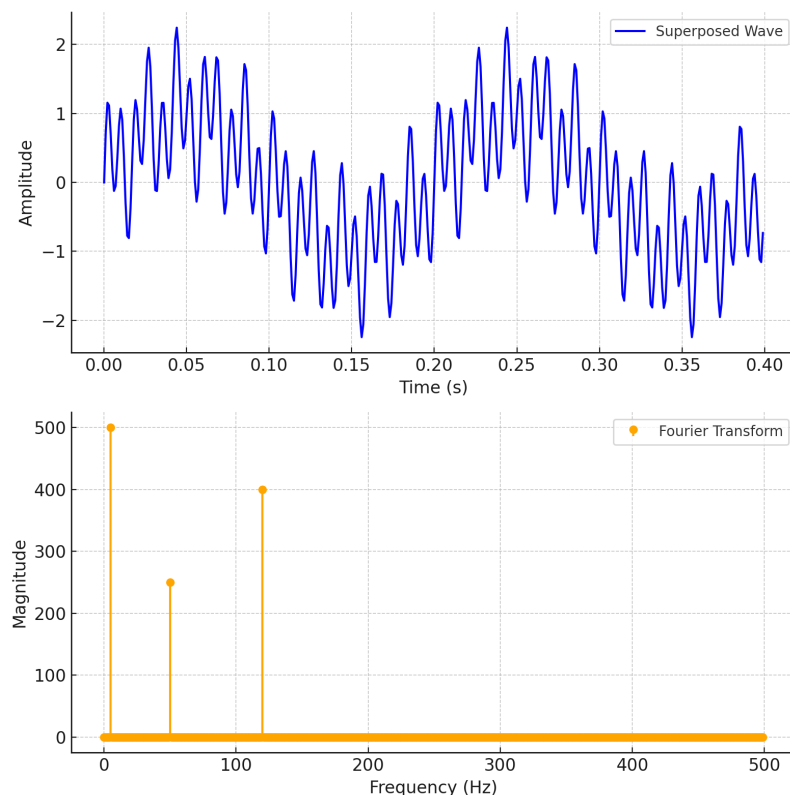coefficients

# Sparsity is everywhere

Many signals of interest admit a sparse representation in a particular basis.

$$f(\boldsymbol{x}) = \sum_{k \geq 0} f_k \psi_k(\boldsymbol{x})$$

← basis

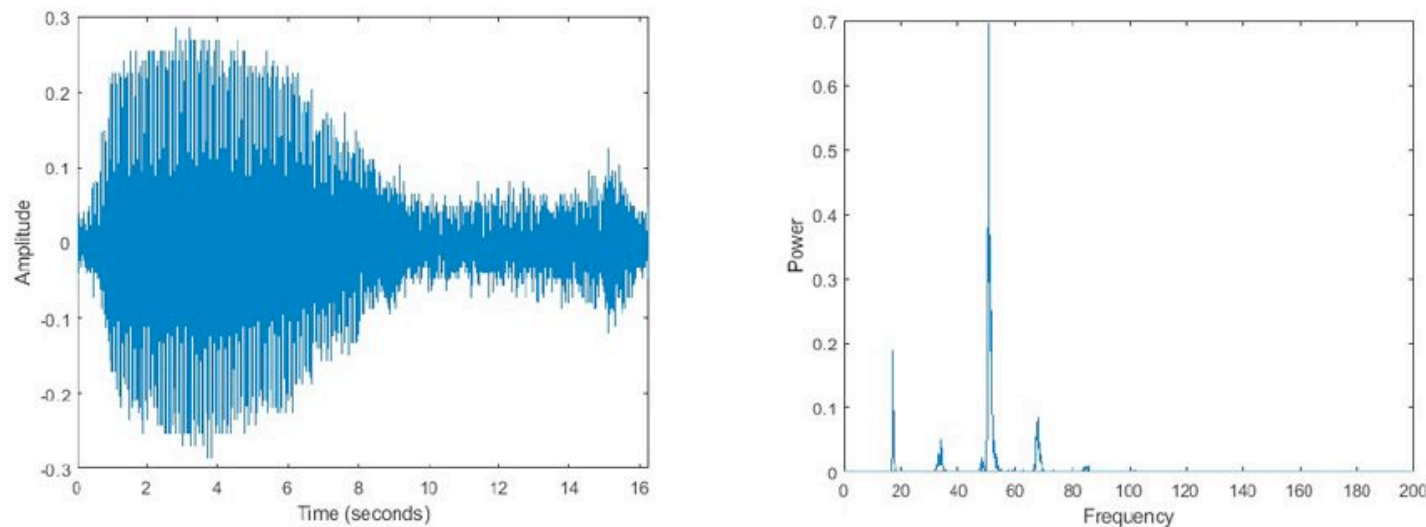← coefficients

Example: superposition of sine waves



$$f(t) = \sin(10\pi t) + 0.5\sin(100\pi t) + 0.8\sin(240\pi t)$$

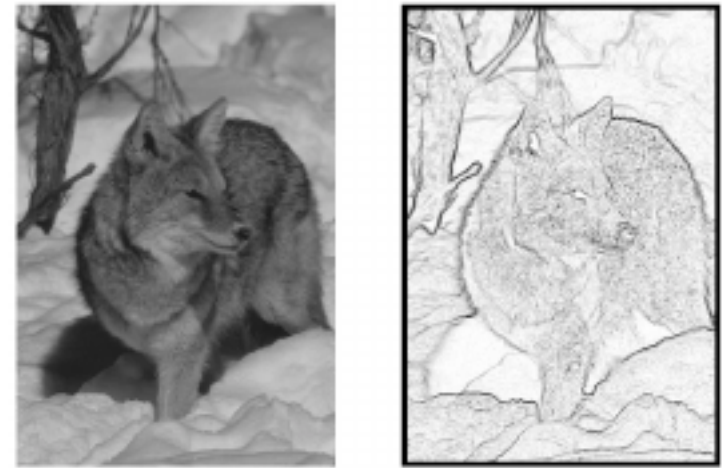$$\hat{f}(\omega) = \delta_5 + 0.5\ \delta_{50} + 0.8\ \delta_{120}$$
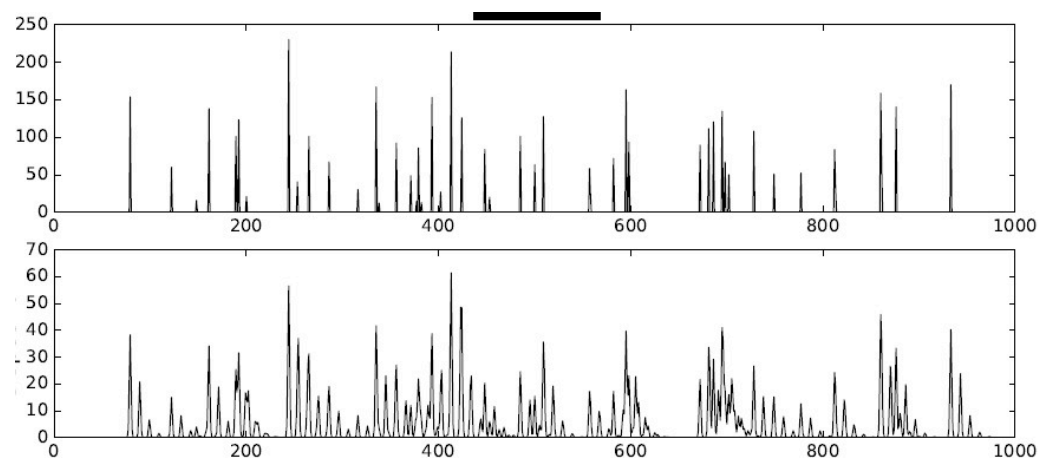
# Sparsity is everywhere

Examples:

Sound



Images



Scientific signals
(mass spectrography)



And many more...

- Portfolio selection (finance)
- Networks (power grids)
- electroencephalogram
- Etc...

# Best subset selection

💡 <u>Idea</u>: encourage solutions which are sparse.

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \frac{1}{2n} \sum_{i=1}^{n} \left( y_i - \langle \boldsymbol{\theta}, \boldsymbol{x}_i \rangle \right)^2 + \lambda ||\boldsymbol{\theta}||_0$$

where $|| \cdot ||_0 : \mathbb{R}^d \to \{0, 1, \dots, d\}$ is the $\ell_0$-"norm": ⚠️ <span style="color:red">Strictly not a norm</span>

$$||\boldsymbol{\theta}||_0 = \sum_{j=1}^{d} \mathbb{I}(\theta_j \neq 0) = \quad \text{\# non-zero entries}$$

# Best subset selection

💡 <u>Idea</u>: encourage solutions which are sparse.

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \frac{1}{2n} \sum_{i=1}^{n} \left(y_i - \langle \boldsymbol{\theta}, \boldsymbol{x}_i \rangle\right)^2 + \lambda ||\boldsymbol{\theta}||_0$$

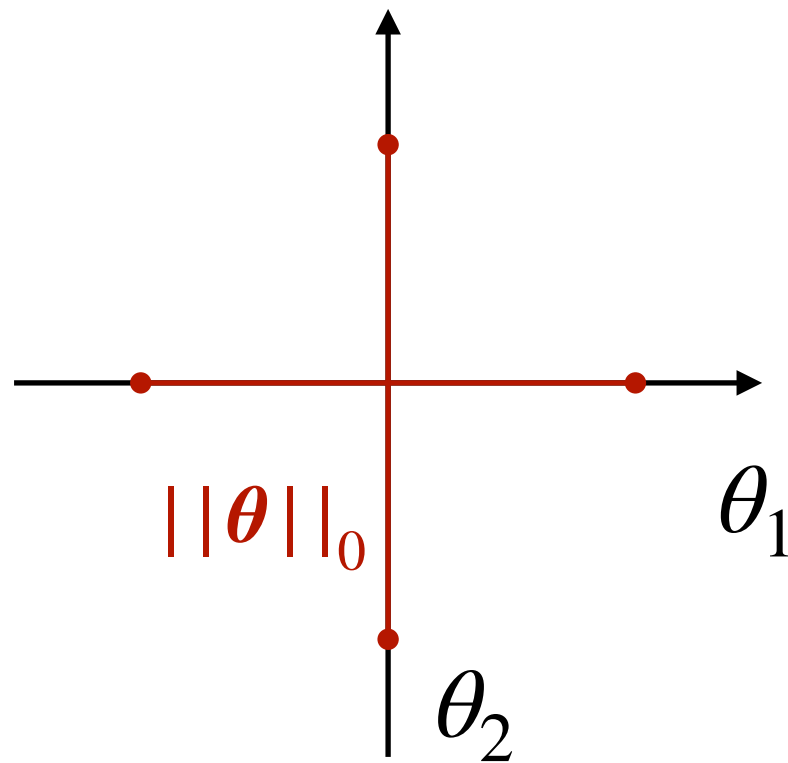where $|| \cdot ||_0 : \mathbb{R}^d \to \{0, 1, \ldots, d\}$ is the $\ell_0$-"norm": ⚠️ <span style="color:red">Strictly not a norm</span>

$$||\boldsymbol{\theta}||_0 = \sum_{j=1}^{d} \mathbb{I}(\theta_j \neq 0) = \quad \text{\# non-zero entries}$$

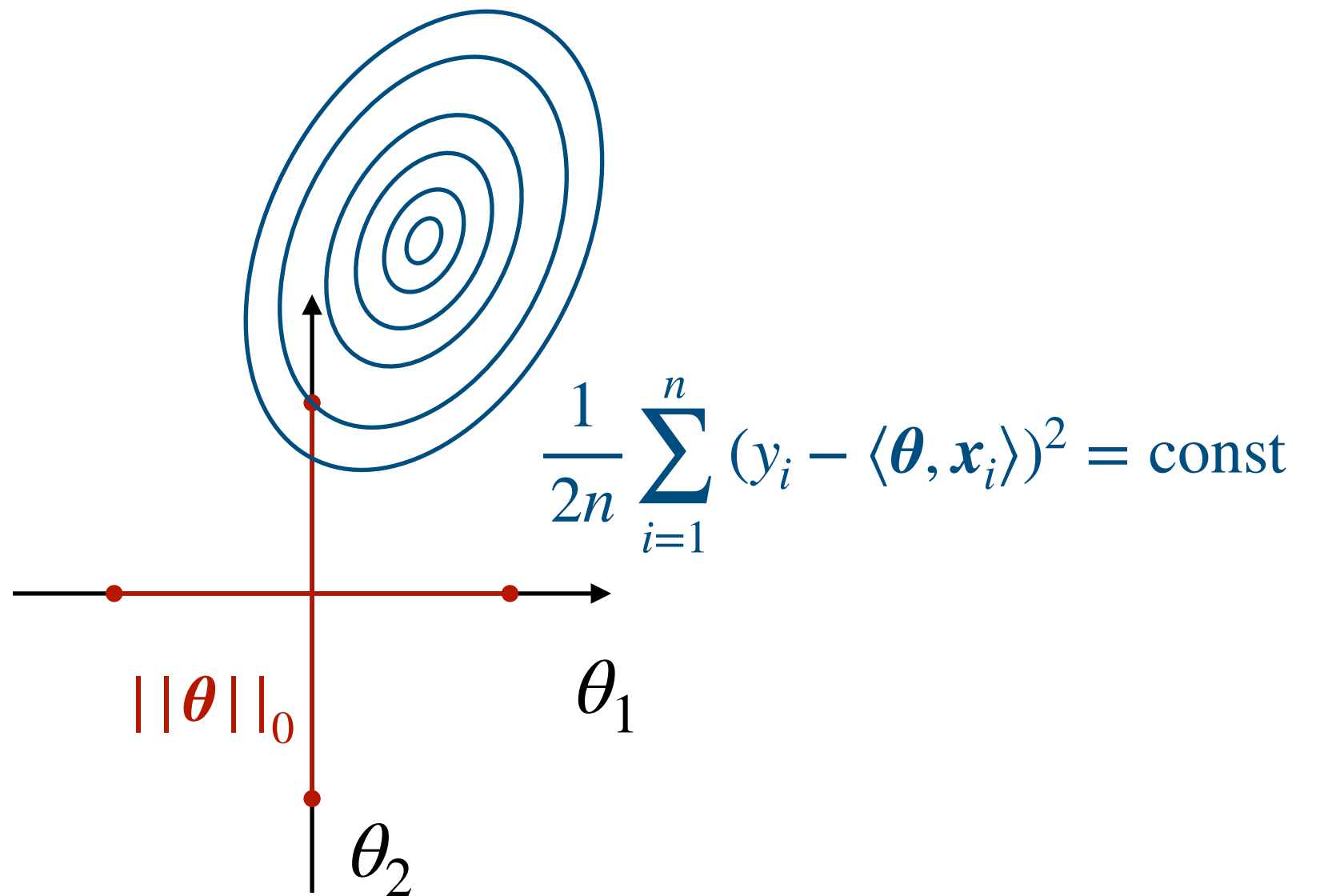Hence, $\lambda \geq 0$ controls the desired sparsity level

- Large $\lambda \gg 1$: encourage more sparsity
- Small $\lambda \ll 1$: encourage less sparsity

# BSS: visualisation

# BSS: visualisation



$$\frac{1}{2n}\sum_{i=1}^{n}(y_i - \langle \boldsymbol{\theta}, \boldsymbol{x}_i \rangle)^2 = \text{const}$$

$||\boldsymbol{\theta}||_0$

$\theta_1$

$\theta_2$

# BSS: visualisation

Solution $\hat{\boldsymbol{\theta}}_{\lambda}$

$$\frac{1}{2n}\sum_{i=1}^{n}(y_i - \langle \boldsymbol{\theta}, \boldsymbol{x}_i \rangle)^2 = \text{const}$$

$||\boldsymbol{\theta}||_0$

$\theta_1$

$\theta_2$

# BSS: orthogonal covariates

To get some intuition about this problem, let's consider a simplified setting: assume the covariates are orthogonal

$$X^\top X = I_d \qquad (n \geq d)$$

# BSS: orthogonal covariates

To get some intuition about this problem, let's consider a simplified setting: assume the covariates are orthogonal

$$X^\top X = I_d \qquad (n \geq d)$$

Then, we can rewrite:

$$||y - X\theta||_2^2 = ||y||_2^2 + \theta^\top X^\top X \theta - 2\theta^\top X^\top y$$

# BSS: orthogonal covariates

To get some intuition about this problem, let's consider a simplified setting: assume the covariates are orthogonal

$$X^\top X = I_d \qquad (n \geq d)$$

Then, we can rewrite:

$$||y - X\theta||_2^2 = ||y||_2^2 + \theta^\top X^\top X\theta - 2\theta^\top X^\top y$$

$$= ||y||_2^2 + ||\theta||_2^2 - 2\theta^\top z \quad (z = X^\top y \in \mathbb{R}^d)$$

# BSS: orthogonal covariates

To get some intuition about this problem, let's consider a simplified setting: assume the covariates are orthogonal

$$X^\top X = I_d \qquad (n \geq d)$$

Then, we can rewrite:

$$||y - X\theta||_2^2 = ||y||_2^2 + \theta^\top X^\top X\theta - 2\theta^\top X^\top y$$

$$= ||y||_2^2 + ||\theta||_2^2 - 2\theta^\top z \quad (z = X^\top y \in \mathbb{R}^d)$$

$$= ||y||_2^2 + ||z||_2^2 - ||z - \theta||_2^2$$

# BSS: orthogonal covariates

To get some intuition about this problem, let's consider a simplified setting: assume the covariates are orthogonal

$$X^\top X = I_d \qquad (n \geq d)$$

Therefore, under the above:

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \frac{1}{2n} \sum_{i=1}^{n} \left( y_i - \langle \boldsymbol{\theta}, x_i \rangle \right)^2 + \lambda ||\boldsymbol{\theta}||_0$$

Is equivalent to:

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \frac{1}{2n} ||z - \boldsymbol{\theta}||_2^2 + \lambda ||\boldsymbol{\theta}||_0$$

Which is a simpler problem since it factorises coordinate-wise.

# BSS: orthogonal covariates

Coordinate-wise, we need to solve

$$\min_{\theta_j \in \mathbb{R}} L(\theta_j) := \left\{ \frac{1}{2n}(z_j - \theta_j)^2 + \lambda \mathbb{I}(\theta_j \neq 0) \right\}$$

# BSS: orthogonal covariates

Coordinate-wise, we need to solve

$$\min_{\theta_j \in \mathbb{R}} L(\theta_j) := \left\{ \frac{1}{2n}(z_j - \theta_j)^2 + \lambda \mathbb{I}(\theta_j \neq 0) \right\}$$

Note that:

$$L(\theta_j) = \frac{1}{2n}(z_j - \theta_j)^2 + \lambda \mathbb{I}(\theta_j \neq 0) = \begin{cases} \frac{1}{2n} z_j^2 & \text{if } \theta_j = 0 \text{ (a)} \\ \frac{1}{2n}(z_j - \theta_j)^2 + \lambda & \text{if } \theta_j \neq 0 \text{ (b)} \end{cases}$$

# BSS: orthogonal covariates

Coordinate-wise, we need to solve

$$\min_{\theta_j \in \mathbb{R}} L(\theta_j) := \left\{ \frac{1}{2n}(z_j - \theta_j)^2 + \lambda \mathbb{I}(\theta_j \neq 0) \right\}$$

Note that:

$$L(\theta_j) = \frac{1}{2n}(z_j - \theta_j)^2 + \lambda \mathbb{I}(\theta_j \neq 0) = \begin{cases} \frac{1}{2n}z_j^2 & \text{if } \theta_j = 0 \text{ (a)} \\ \frac{1}{2n}(z_j - \theta_j)^2 + \lambda & \text{if } \theta_j \neq 0 \text{ (b)} \end{cases}$$

Note the solution of the problem is not unique:

- In case (a), solution is $\hat{\theta}_{\lambda,j}^{(1)} = 0$
- In case (b), solution is $\hat{\theta}_{\lambda,j}^{(2)} = z_j$

# BSS: orthogonal covariates

Coordinate-wise, we need to solve

$$\min_{\theta_j \in \mathbb{R}} L(\theta_j) := \left\{ \frac{1}{2n}(z_j - \theta_j)^2 + \lambda \mathbb{I}(\theta_j \neq 0) \right\}$$

Note that:

$$L(\theta_j) = \frac{1}{2n}(z_j - \theta_j)^2 + \lambda \mathbb{I}(\theta_j \neq 0) = \begin{cases} \frac{1}{2n} z_j^2 & \text{if } \theta_j = 0 \text{ (a)} \\ \frac{1}{2n}(z_j - \theta_j)^2 + \lambda & \text{if } \theta_j \neq 0 \text{ (b)} \end{cases}$$

Note the solution of the problem is not unique:

- In case (a), solution is $\hat{\theta}_{\lambda,j}^{(1)} = 0$

- In case (b), solution is $\hat{\theta}_{\lambda,j}^{(2)} = z_j$

Which one to pick? The one with minimal loss.

# BSS: orthogonal covariates

Note the solution of the problem is not unique:

- In case (a), solution is $\hat{\theta}_{\lambda,j}^{(1)} = 0$

- In case (b), solution is $\hat{\theta}_{\lambda,j}^{(2)} = z_j$

Which one to pick? The one with minimal loss.

$$L\left(\hat{\theta}_{\lambda,j}^{(2)}\right) - L\left(\hat{\theta}_{\lambda,j}^{(1)}\right) = -\frac{z_j^2}{2n} + \lambda \underset{?}{\geq} 0$$

# BSS: orthogonal covariates

Note the solution of the problem is not unique:

- In case (a), solution is $\hat{\theta}^{(1)}_{\lambda,j} = 0$

- In case (b), solution is $\hat{\theta}^{(2)}_{\lambda,j} = z_j$

Which one to pick? The one with minimal loss.

$$L\left(\hat{\theta}^{(2)}_{\lambda,j}\right) - L\left(\hat{\theta}^{(1)}_{\lambda,j}\right) = -\frac{z_j^2}{2n} + \lambda \underset{?}{\geq} 0 \quad \Leftrightarrow \quad 2n\lambda \geq z_j^2$$

# BSS: orthogonal covariates

Note the solution of the problem is not unique:

- In case (a), solution is $\hat{\theta}_{\lambda,j}^{(1)} = 0$

- In case (b), solution is $\hat{\theta}_{\lambda,j}^{(2)} = z_j$

Which one to pick? The one with minimal loss.

$$L\left(\hat{\theta}_{\lambda,j}^{(2)}\right) - L\left(\hat{\theta}_{\lambda,j}^{(1)}\right) = -\frac{z_j^2}{2n} + \lambda \underset{?}{\geq} 0 \quad \Leftrightarrow \quad 2n\lambda \geq z_j^2$$

Hence, the solution is given by:

$$\hat{\theta}_{\lambda,j} = \begin{cases} 0 & \text{if } z_j^2 < 2n\lambda \\ z_j & \text{if } z_j^2 \geq 2n\lambda \end{cases}$$

"Hard threshold" function

# BSS: orthogonal covariates

Putting together, the solution of the BSS problem:

$$\min_{\boldsymbol{\theta}\in\mathbb{R}^d} \frac{1}{2n}\sum_{i=1}^{n}\left(y_i - \langle\boldsymbol{\theta},\boldsymbol{x}_i\rangle\right)^2 + \lambda\|\boldsymbol{\theta}\|_0$$

Under the assumption of $\boldsymbol{X}^\top\boldsymbol{X} = \boldsymbol{I}_d$ is given by:

$$\hat{\boldsymbol{\theta}}_\lambda = H_{\sqrt{2n\lambda}}(\boldsymbol{X}^\top\boldsymbol{y})$$

Where:

$$H_\lambda(z) = \begin{cases} 0 & \text{if } |z| < \lambda \\ z & \text{otherwise} \end{cases}$$