



Statistical Learning II

Lecture 9 - BSS & LASSO

Bruno Loureiro
@ CSD, DI-ENS & CNRS

brloureiro@gmail.com

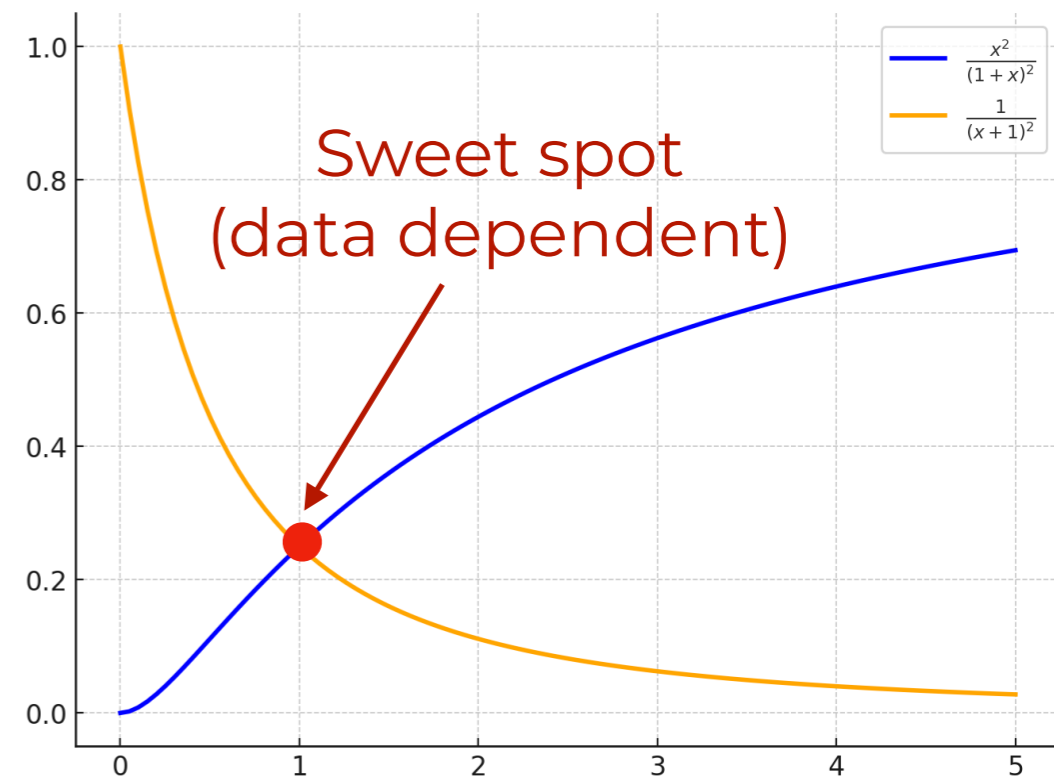
Risk of ridge

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

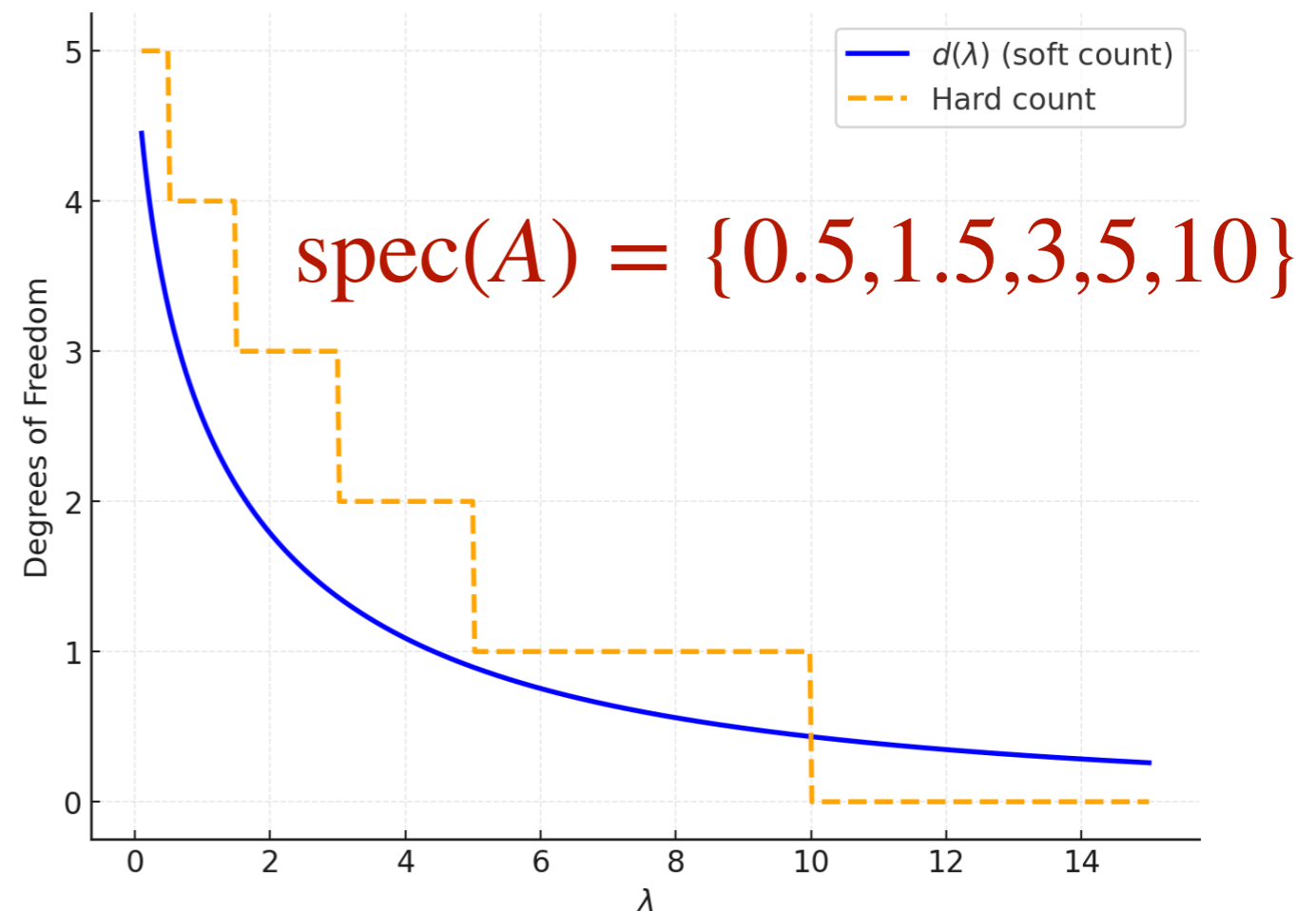
Let $A \in \mathbb{R}^{d \times d}$ be a positive definite matrix with decreasing eigenvalues $\text{spec}(A) = \{\lambda_k : k = 1, \dots, d\}$. Define the cumulative:

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

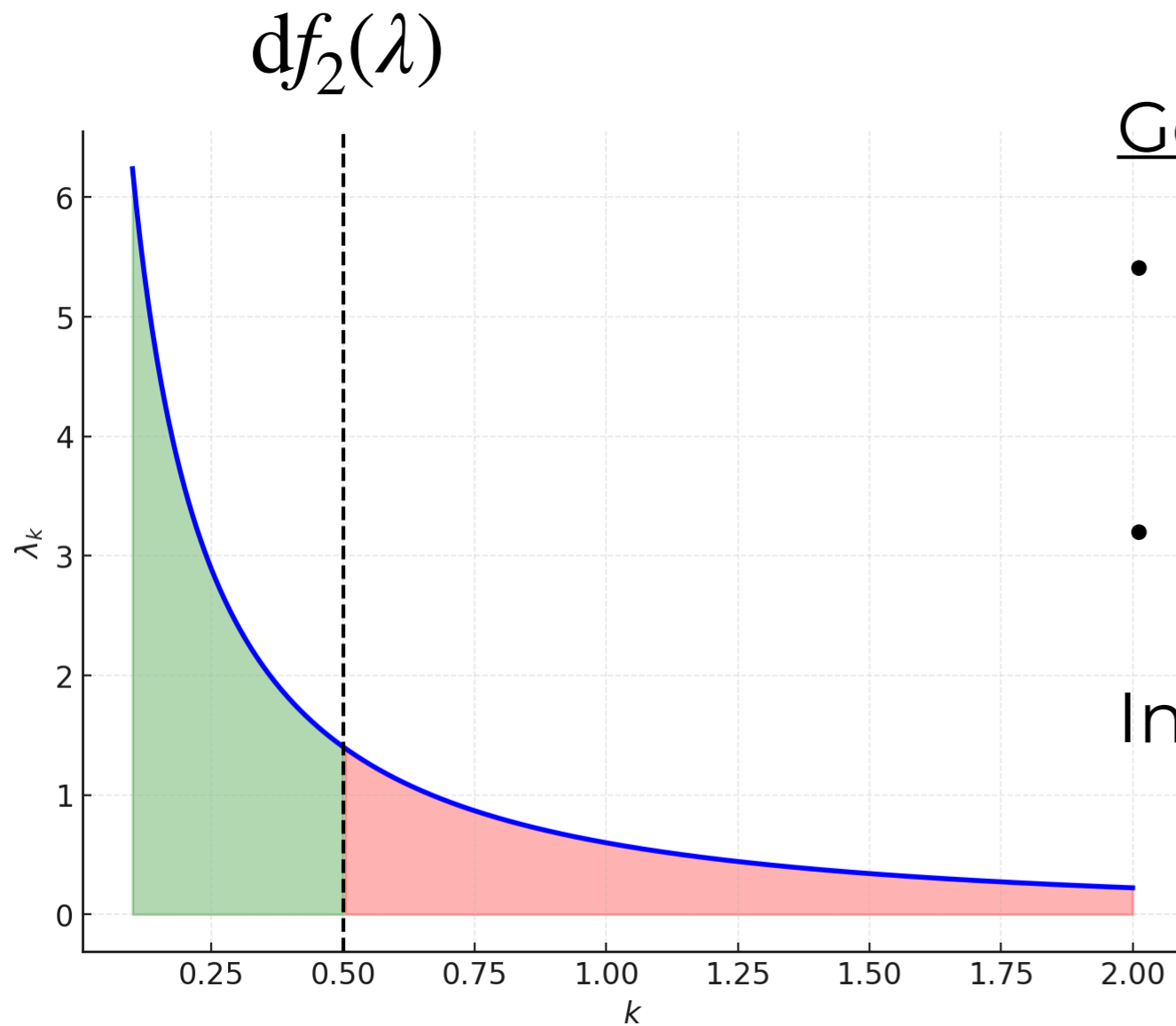
The variance of the ridge risk can be seen as a soft version:

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

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



Choosing regularisation



Goal: pick λ such that:

- directions in \mathbf{X} that better correlate with θ_\star are retained
- Shrink remaining directions

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

Low-frequency

High-frequency

Best subset selection & the LASSO

Pitfalls of ridge

The ridge estimation performs uniform shrinkage.

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

In other words: ℓ_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.

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

In other words: ℓ_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

$$\boldsymbol{\theta}_{\star, j} \neq 0 \quad i = 1, \dots, d$$

$$\boldsymbol{\theta}_{\star} = \begin{bmatrix} \color{teal}{\blacksquare} \\ \color{orange}{\blacksquare} \\ \color{blue}{\blacksquare} \\ \vdots \\ \color{red}{\blacksquare} \end{bmatrix}$$

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

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

$$\boldsymbol{\theta}_{\star} = \begin{bmatrix} \square \\ \color{orange}{\blacksquare} \\ \square \\ \vdots \\ \color{red}{\blacksquare} \end{bmatrix}$$

Sparsity is everywhere

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

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

← basis

← coefficients

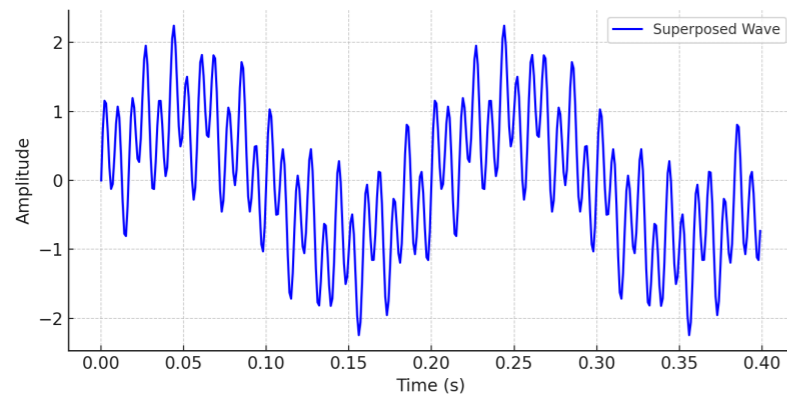
Sparsity is everywhere

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

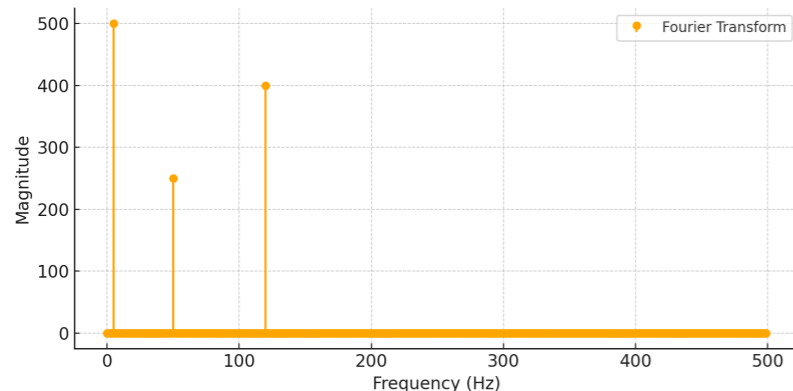
$$f(\mathbf{x}) = \sum_{k \geq 0} f_k \psi_k(\mathbf{x})$$

← basis
← coefficients

Example: superposition of sine waves



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

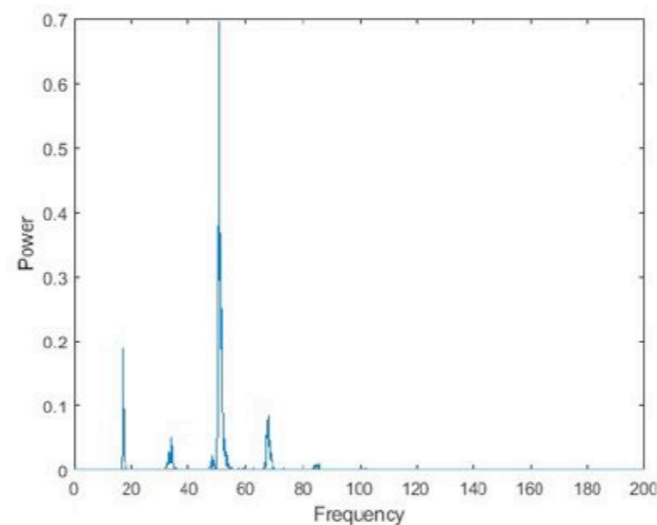
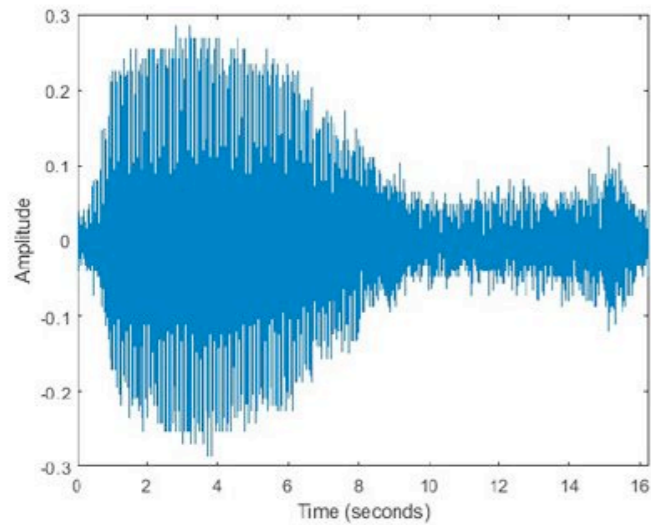


$$\hat{f}(\omega) = \delta_5 + 0.5 \delta_{50} + 0.8 \delta_{120}$$

Sparsity is everywhere

Examples:

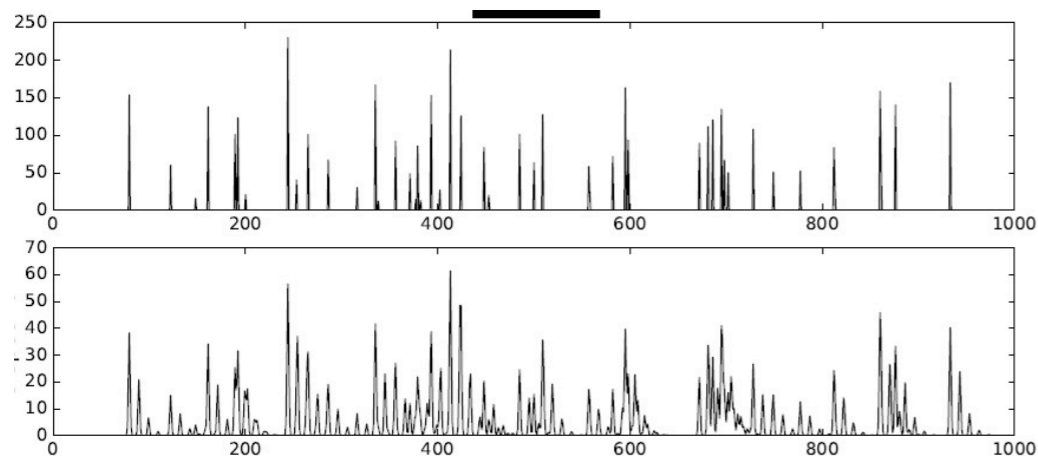
Sound



Images



Scientific signals (mass spectrography)



And many more...


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

Best subset selection



Idea: encourage solutions which are sparse.

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

where $\|\cdot\|_0 : \mathbb{R}^d \rightarrow \{0, 1, \dots, d\}$ is the ℓ_0 -“norm”:  Strictly not a norm

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

Best subset selection



Idea: encourage solutions which are sparse.

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

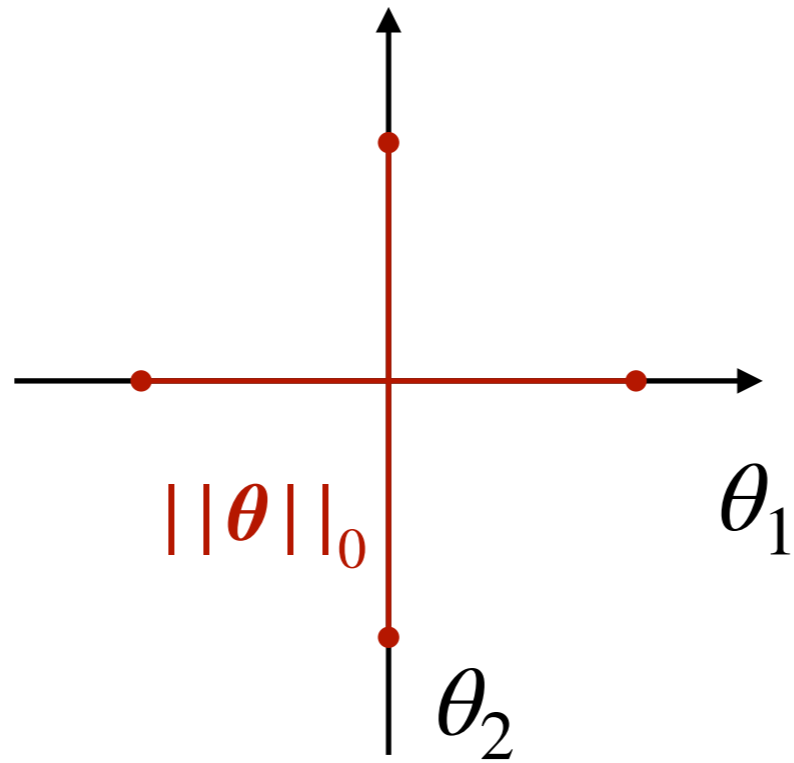
where $\|\cdot\|_0 : \mathbb{R}^d \rightarrow \{0, 1, \dots, d\}$ is the ℓ_0 -“norm”:  Strictly not a norm

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

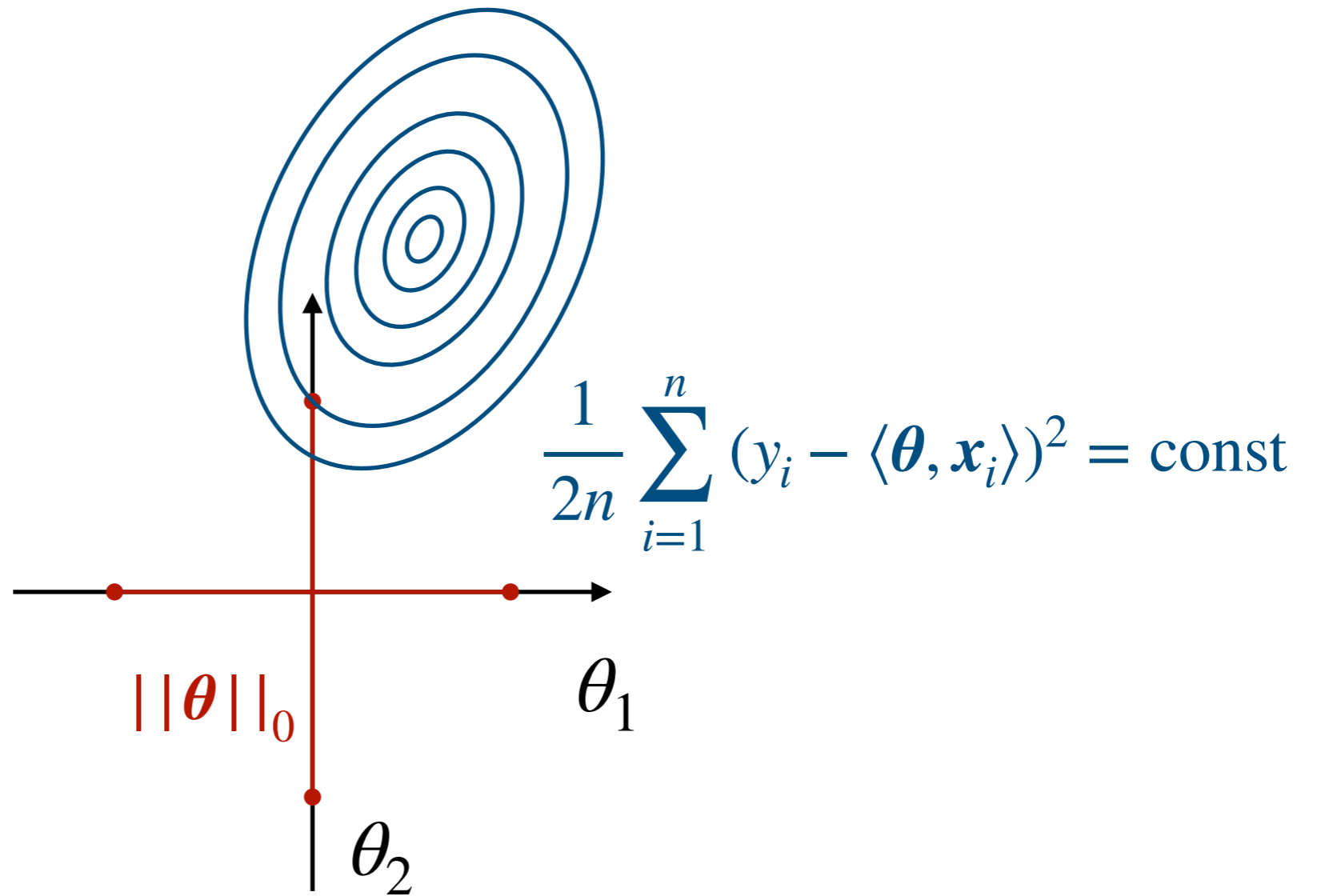
Hence, $\lambda \geq 0$ controls the desired sparsity level

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

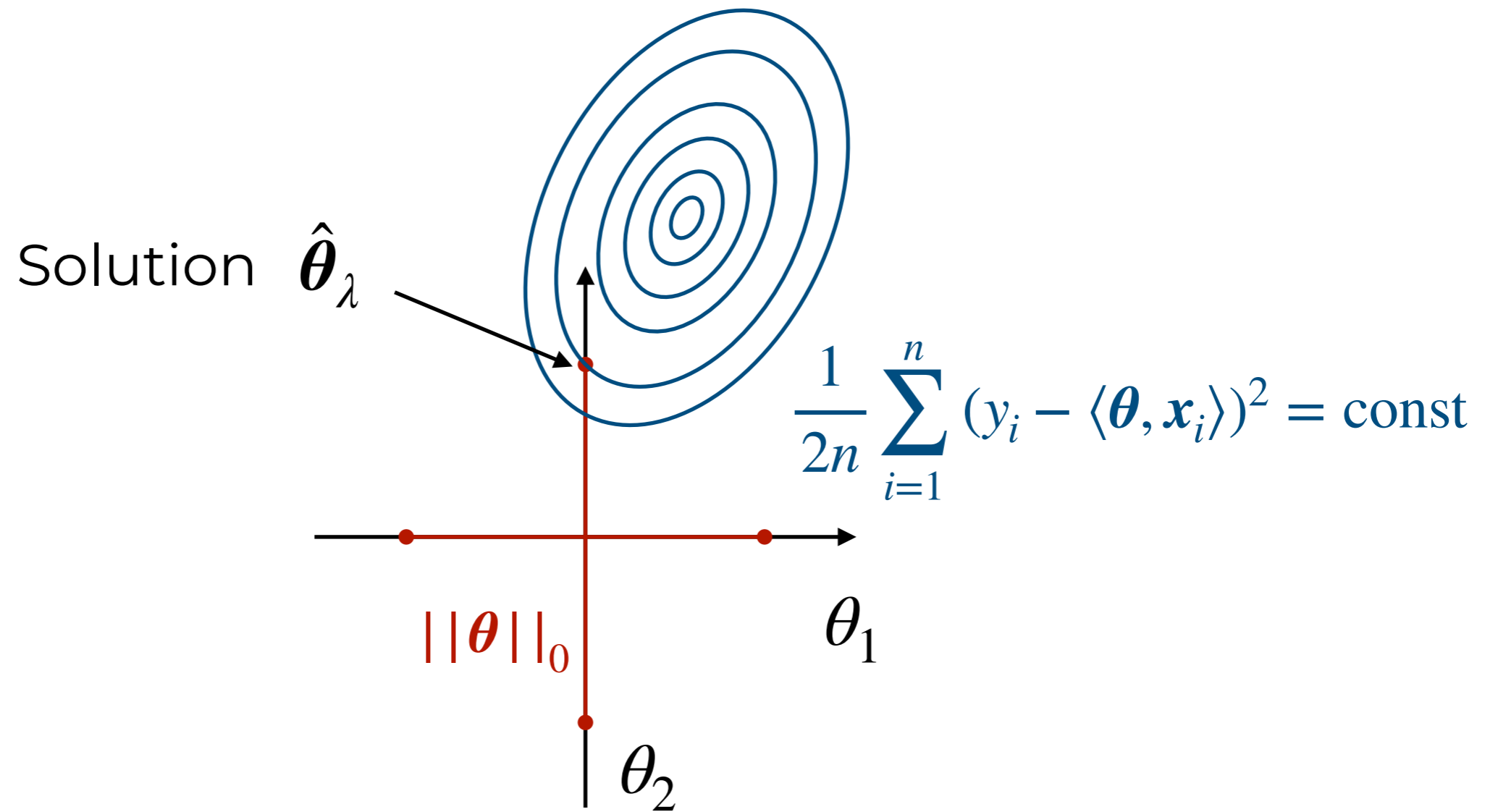
BSS: visualisation



BSS: visualisation



BSS: visualisation



BSS: orthogonal covariates

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

$$\mathbf{X}^\top \mathbf{X} = \mathbf{I}_d \quad (n \geq d)$$

BSS: orthogonal covariates

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

$$\mathbf{X}^\top \mathbf{X} = \mathbf{I}_d \quad (n \geq d)$$

Then, we can rewrite:

$$\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 = \|\mathbf{y}\|_2^2 + \boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\theta} - 2\boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{y}$$

BSS: orthogonal covariates

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

$$\mathbf{X}^\top \mathbf{X} = \mathbf{I}_d \quad (n \geq d)$$

Then, we can rewrite:

$$\begin{aligned} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 &= \|\mathbf{y}\|_2^2 + \boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\theta} - 2\boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{y} \\ &= \|\mathbf{y}\|_2^2 + \|\boldsymbol{\theta}\|_2^2 - 2\boldsymbol{\theta}^\top \mathbf{z} \quad (\mathbf{z} = \mathbf{X}^\top \mathbf{y} \in \mathbb{R}^d) \end{aligned}$$

BSS: orthogonal covariates

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

$$\mathbf{X}^\top \mathbf{X} = \mathbf{I}_d \quad (n \geq d)$$

Then, we can rewrite:

$$\begin{aligned} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 &= \|\mathbf{y}\|_2^2 + \boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\theta} - 2\boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{y} \\ &= \|\mathbf{y}\|_2^2 + \|\boldsymbol{\theta}\|_2^2 - 2\boldsymbol{\theta}^\top \mathbf{z} \quad (\mathbf{z} = \mathbf{X}^\top \mathbf{y} \in \mathbb{R}^d) \\ &= \|\mathbf{y}\|_2^2 + \|\mathbf{z}\|_2^2 - \|\mathbf{z} - \boldsymbol{\theta}\|_2^2 \end{aligned}$$

BSS: orthogonal covariates

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

$$\mathbf{X}^\top \mathbf{X} = \mathbf{I}_d \quad (n \geq d)$$

Therefore, under the above:

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

Is equivalent to:

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

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

BSS: orthogonal covariates

Coordinate-wise, we need to solve

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

BSS: orthogonal covariates

Coordinate-wise, we need to solve

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

Note that:

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

BSS: orthogonal covariates

Coordinate-wise, we need to solve

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

Note that:

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

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

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

Note that:

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

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

Which one to pick? The one with minimal loss.

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

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.

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

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:

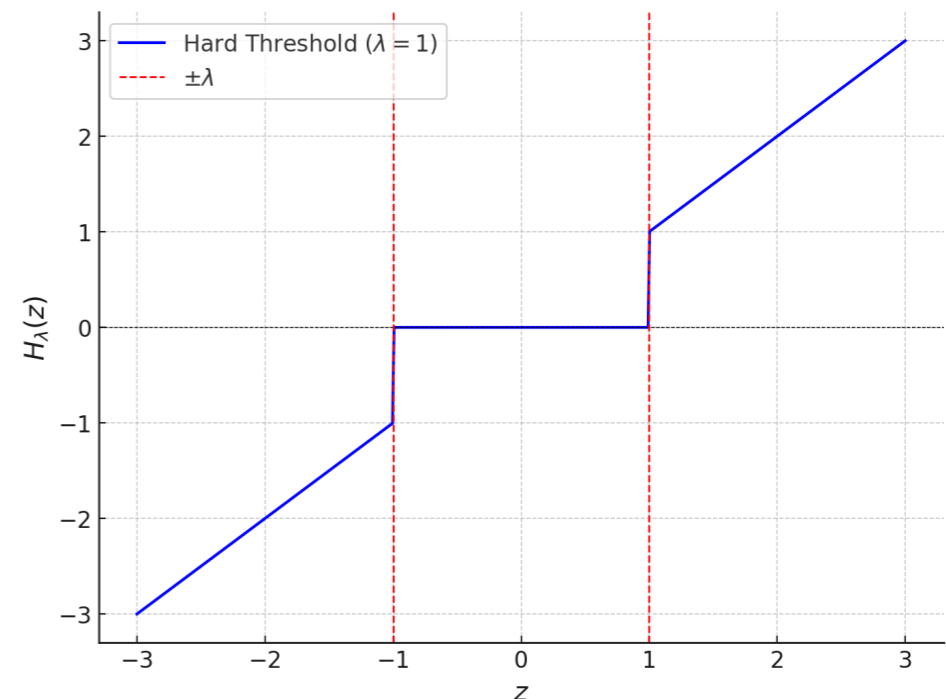
$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \frac{1}{2n} \sum_{i=1}^n (y_i - \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle)^2 + \lambda \|\boldsymbol{\theta}\|_0$$

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

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

Where:

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



BSS: orthogonal covariates

To understand better this solution, consider a linear model for the data:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta}_\star + \boldsymbol{\varepsilon}$$

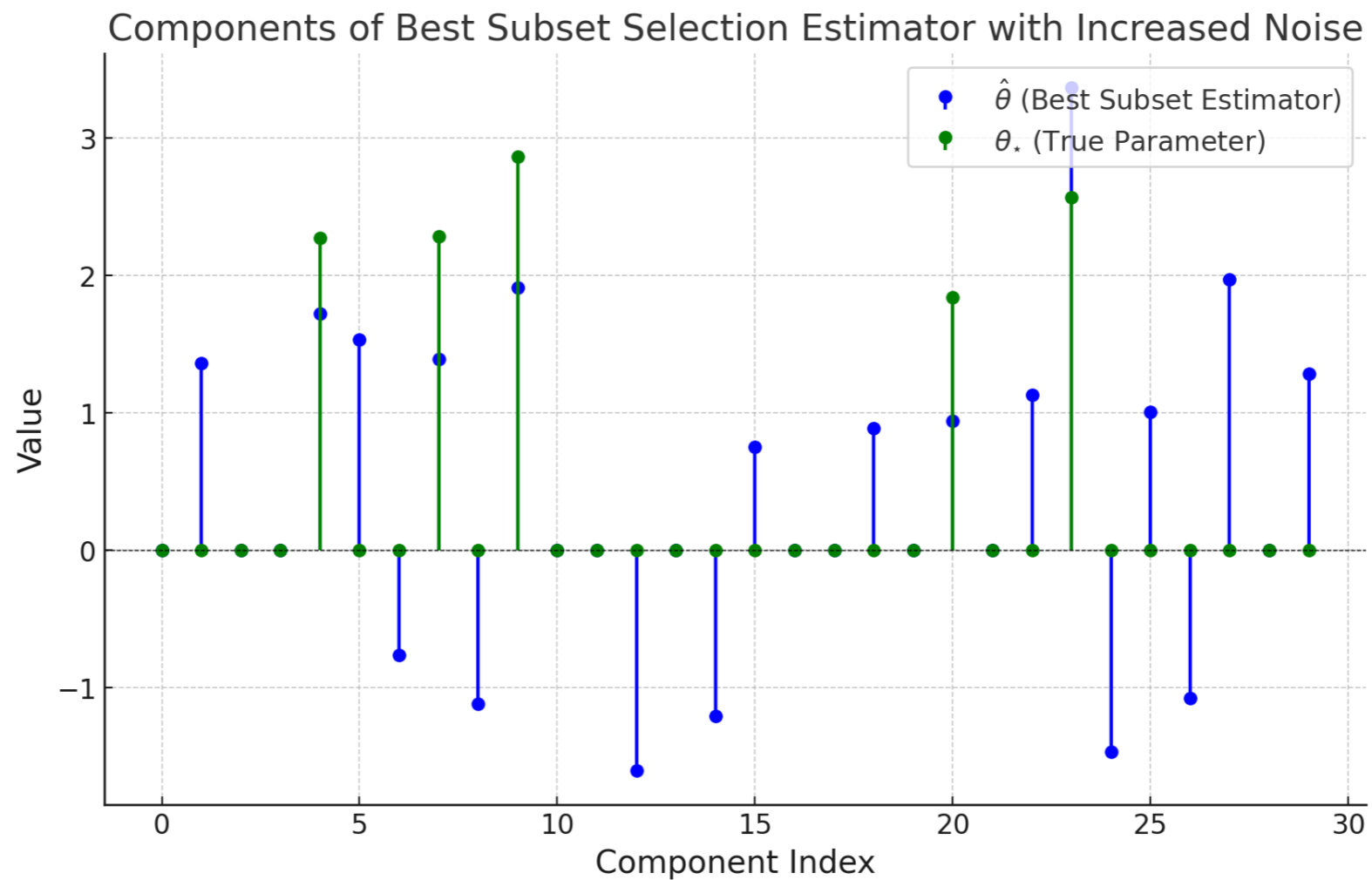
With $\mathbb{E}[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top] = \sigma\mathbf{I}_n$ and $\boldsymbol{\theta}_\star$ a k -sparse vector
 $\mathbb{E}[\boldsymbol{\varepsilon}] = 0$

The, the solution is given by:

$$\hat{\boldsymbol{\theta}}_\lambda = H_{\sqrt{2n\lambda}}(\boldsymbol{\theta}_\star + \mathbf{X}^\top\boldsymbol{\varepsilon})$$

BSS: orthogonal covariates

Example: $n = 40$ $\lambda = 0.5$ θ_{\star} 5-sparse
 $d = 30$ $\sigma^2 = 1$ $\|\theta_{\star}\|_2^2 = 5.35$



BSS: beyond orthogonal

When the covariates are not orthogonal, an explicit solution is not available. Nevertheless, we can partially characterise it.

BSS: beyond orthogonal

When the covariates are not orthogonal, an explicit solution is not available. Nevertheless, we can partially characterise it.

Let $S = \{j \in [d] : \hat{\theta}_{\lambda,j} \neq 0\}$ denote the support of the BSS solution

Denoting:

- $\hat{\boldsymbol{\theta}}_S \in \mathbb{R}^{|S|}$ the non-zero entries of $\hat{\boldsymbol{\theta}}_\lambda \in \mathbb{R}^d$
- $\mathbf{X}_S \in \mathbb{R}^{n \times |S|}$ the corresponding covariates

BSS: beyond orthogonal

When the covariates are not orthogonal, an explicit solution is not available. Nevertheless, we can partially characterise it.

Let $S = \{j \in [d] : \hat{\theta}_{\lambda,j} \neq 0\}$ denote the support of the BSS solution

Denoting:

- $\hat{\theta}_S \in \mathbb{R}^{|S|}$ the non-zero entries of $\hat{\theta}_\lambda \in \mathbb{R}^d$
- $X_S \in \mathbb{R}^{n \times |S|}$ the corresponding covariates

We can write:

$$\hat{\theta}_S = X_S^+ \mathbf{y}$$

In other words, BSS = OLS in the support!

The hard part is to find S as a function of $X, \mathbf{y}, \lambda \dots$

Pitfalls of BSS

More generally, BSS is that it is a **non-convex** problem

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

In particular, for general covariates it is **hard to optimise**.
(it is actually a **NP-hard problem** in the worst case)

Pitfalls of BSS

More generally, BSS is that it is a **non-convex** problem

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

In particular, for general covariates it is **hard to optimise**.
(it is actually a **NP-hard problem** in the worst case)



Question: $\|\cdot\|_0$ is what makes this non-convex. Can we find another regularisation with similar properties but convex?

Pitfalls of BSS

More generally, BSS is that it is a **non-convex** problem

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

In particular, for general covariates it is **hard to optimise**.
(it is actually a **NP-hard problem** in the worst case)



Question: $\|\cdot\|_0$ is what makes this non-convex. Can we find another regularisation with similar properties but convex?



That's the key idea of the LASSO.

LASSO

The Least Absolute Shrinkage and Selection Operator (LASSO) is defined as the solution of the following problem:

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \frac{1}{2n} \sum_{i=1}^n (y_i - \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle)^2 + \lambda \|\boldsymbol{\theta}\|_1$$

where $\|\cdot\|_1 : \mathbb{R}^d \rightarrow \mathbb{R}_+$ is the ℓ_1 -norm:

$$\|\boldsymbol{\theta}\|_1 = \sum_{j=1}^d |\theta_j|$$

LASSO

The Least Absolute Shrinkage and Selection Operator (LASSO) is defined as the solution of the following problem:

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \frac{1}{2n} \sum_{i=1}^n (y_i - \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle)^2 + \lambda \|\boldsymbol{\theta}\|_1$$

where $\|\cdot\|_1 : \mathbb{R}^d \rightarrow \mathbb{R}_+$ is the ℓ_1 -norm:

$$\|\boldsymbol{\theta}\|_1 = \sum_{j=1}^d |\theta_j|$$

Moreover, this is a **convex** problem.

LASSO

The Least Absolute Shrinkage and Selection Operator (LASSO) is defined as the solution of the following problem:

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \frac{1}{2n} \sum_{i=1}^n (y_i - \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle)^2 + \lambda \|\boldsymbol{\theta}\|_1$$

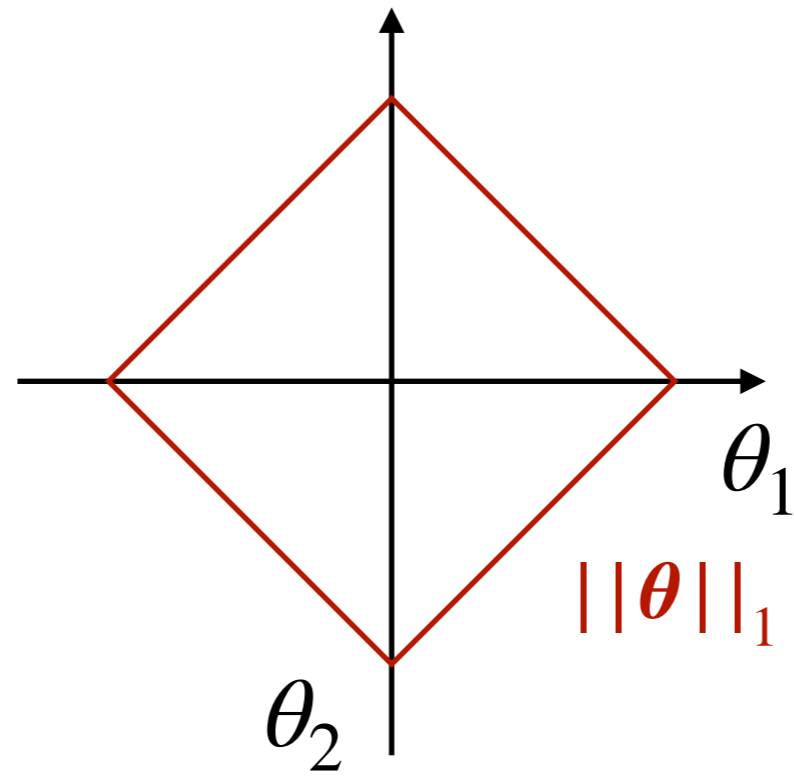
where $\|\cdot\|_1 : \mathbb{R}^d \rightarrow \mathbb{R}_+$ is the ℓ_1 -norm:

$$\|\boldsymbol{\theta}\|_1 = \sum_{j=1}^d |\theta_j|$$

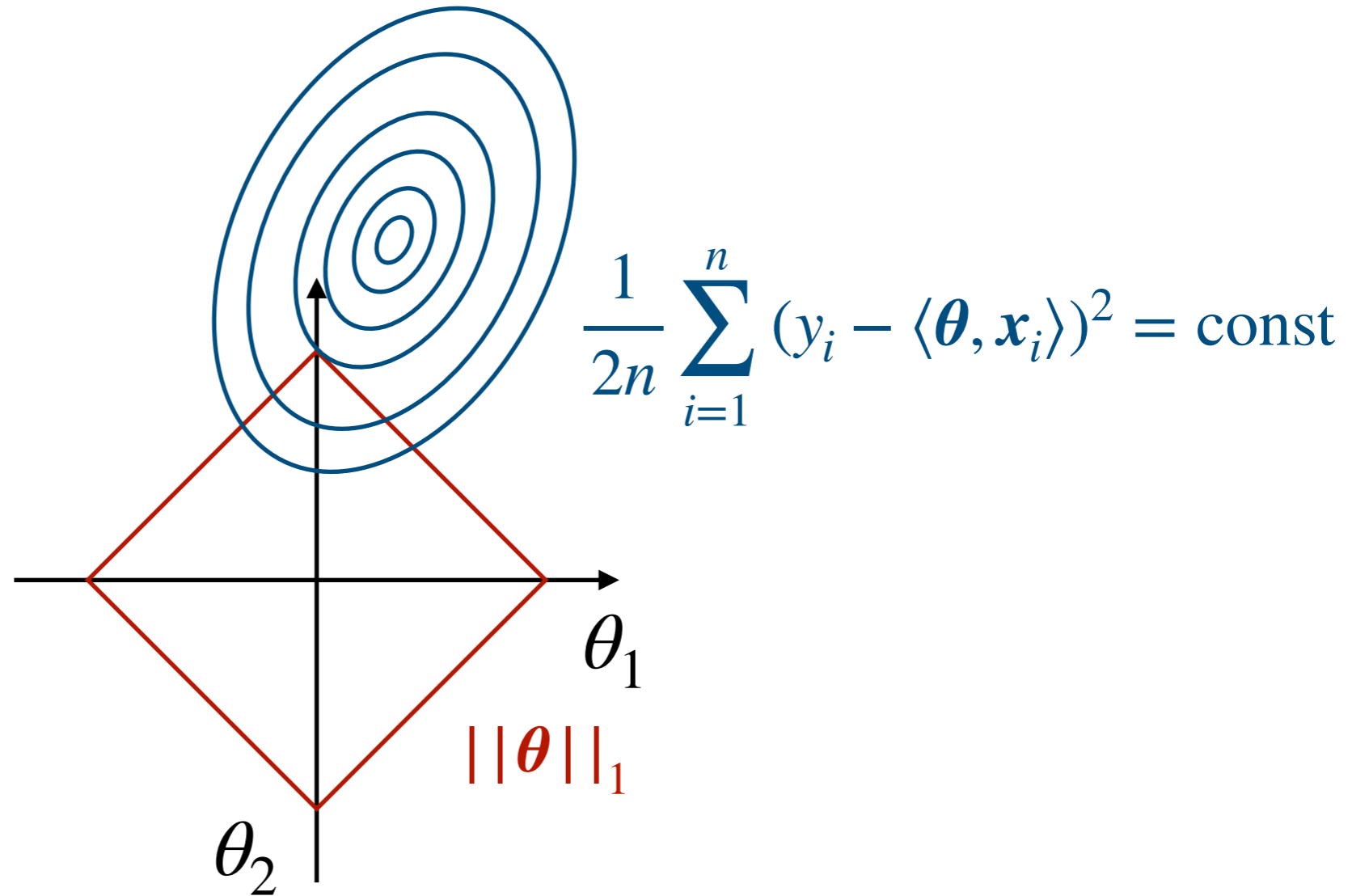
Moreover, this is a **convex** problem.

Note that both $\|\cdot\|_1$ and $\|\cdot\|_2$ are small for sparse vectors... why this is different?

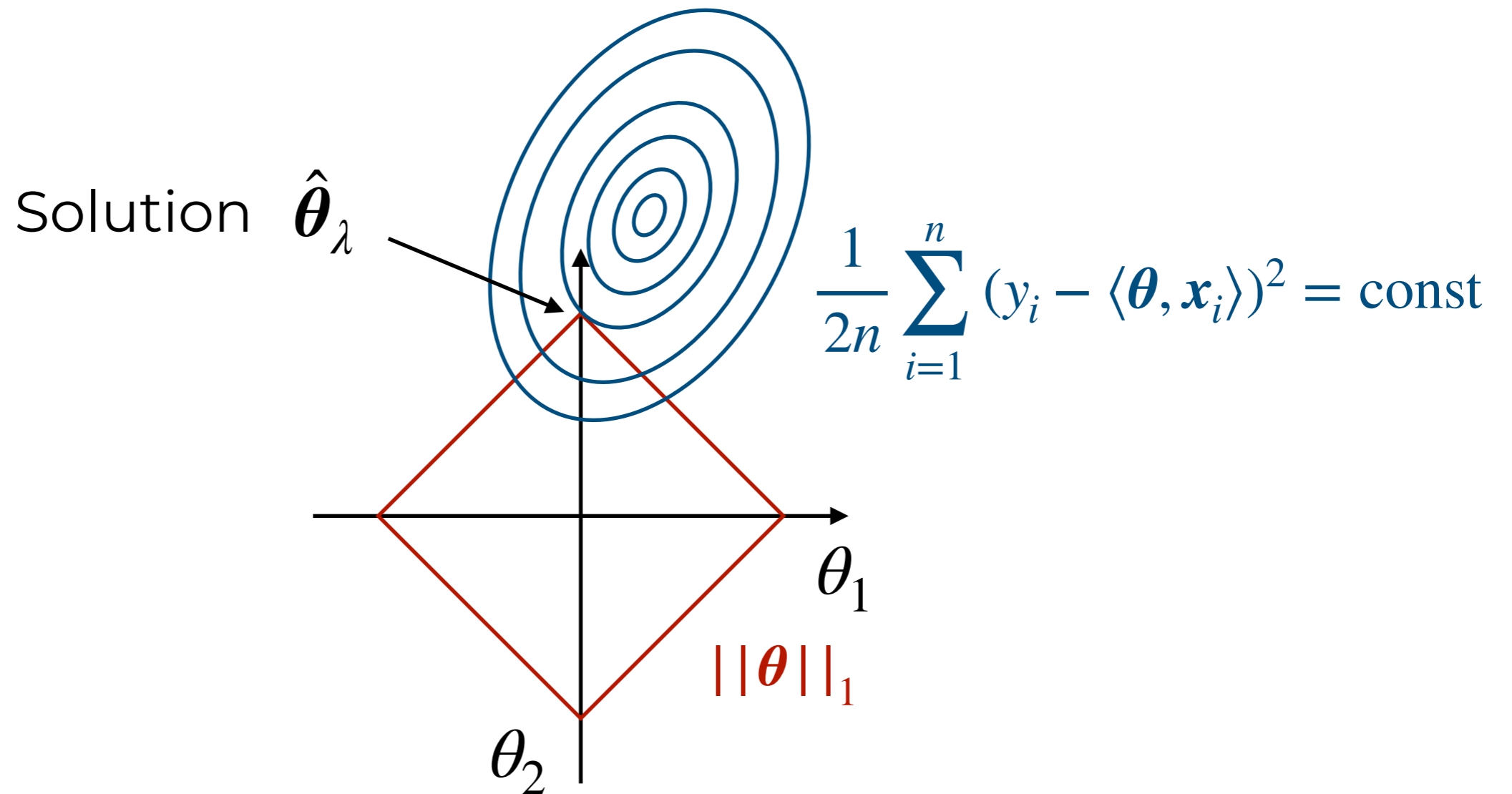
LASSO: visualisation



LASSO: visualisation



LASSO: visualisation



Sharper corners favours sparser solutions!