



# Statistical Learning II

Lecture 11 - LASSO & Feature maps

---

**Bruno Loureiro**  
@ CSD, DI-ENS & CNRS

[brloureiro@gmail.com](mailto:brloureiro@gmail.com)

# BSS: orthogonal covariates

Putting together, the solution of the BSS problem:

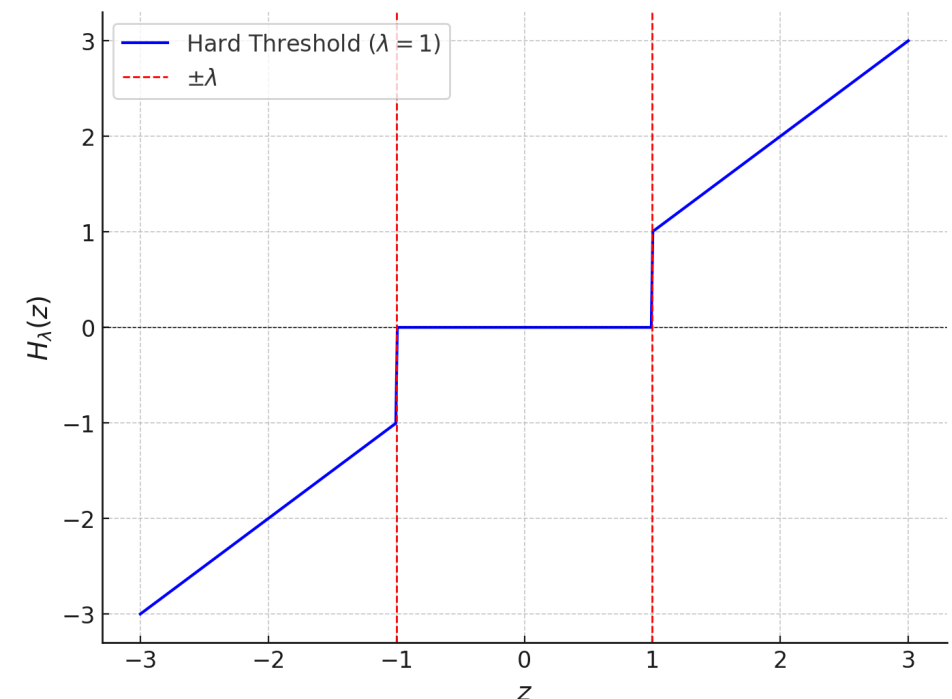
$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \frac{1}{2n} \sum_{i=1}^n (y_i - \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle)^2 + \lambda ||\boldsymbol{\theta}||_0$$

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:

$$y = X\theta_{\star} + \varepsilon$$

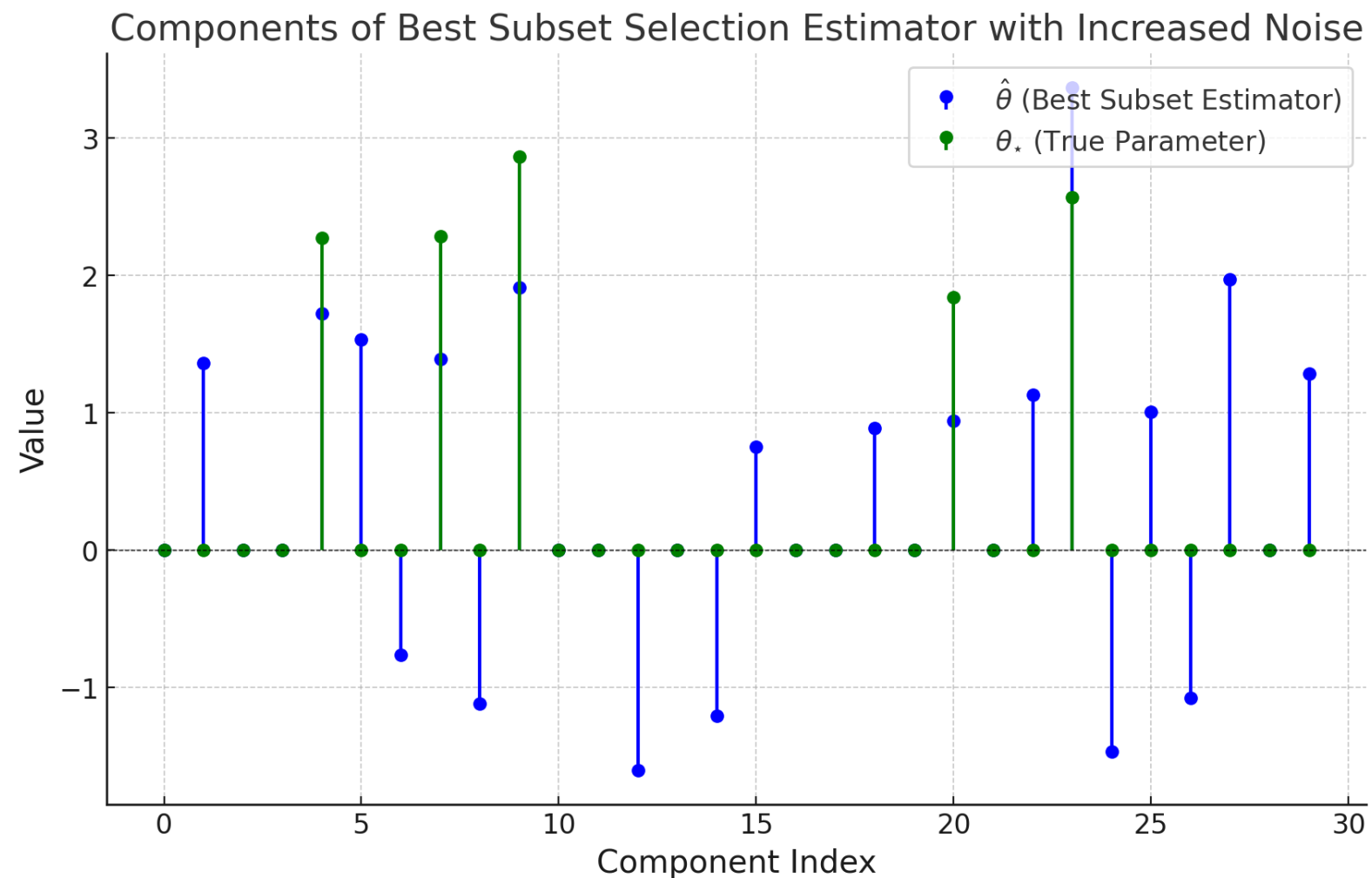
With  $\mathbb{E}[\varepsilon\varepsilon^{\top}] = \sigma\mathbf{I}_n$  and  $\theta_{\star}$  a  $k$ -sparse vector  
 $\mathbb{E}[\varepsilon] = 0$

The, the solution is given by:

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

# BSS: orthogonal covariates

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



# 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  $\ell_1$ -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  $\ell_1$ -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:

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

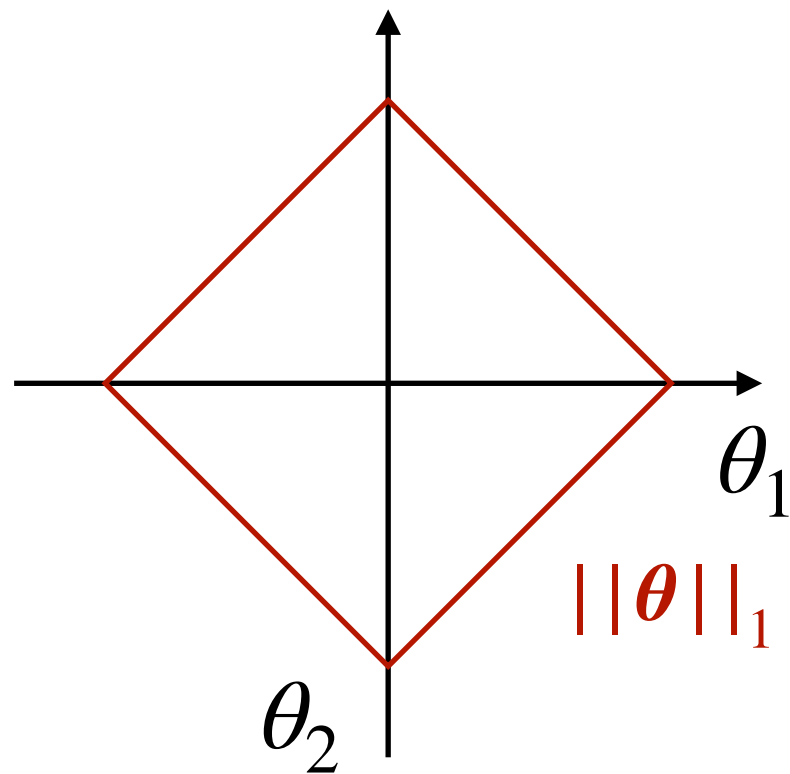
where  $||\cdot||_1 : \mathbb{R}^d \rightarrow \mathbb{R}_+$  is the  $\ell_1$ -norm:

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

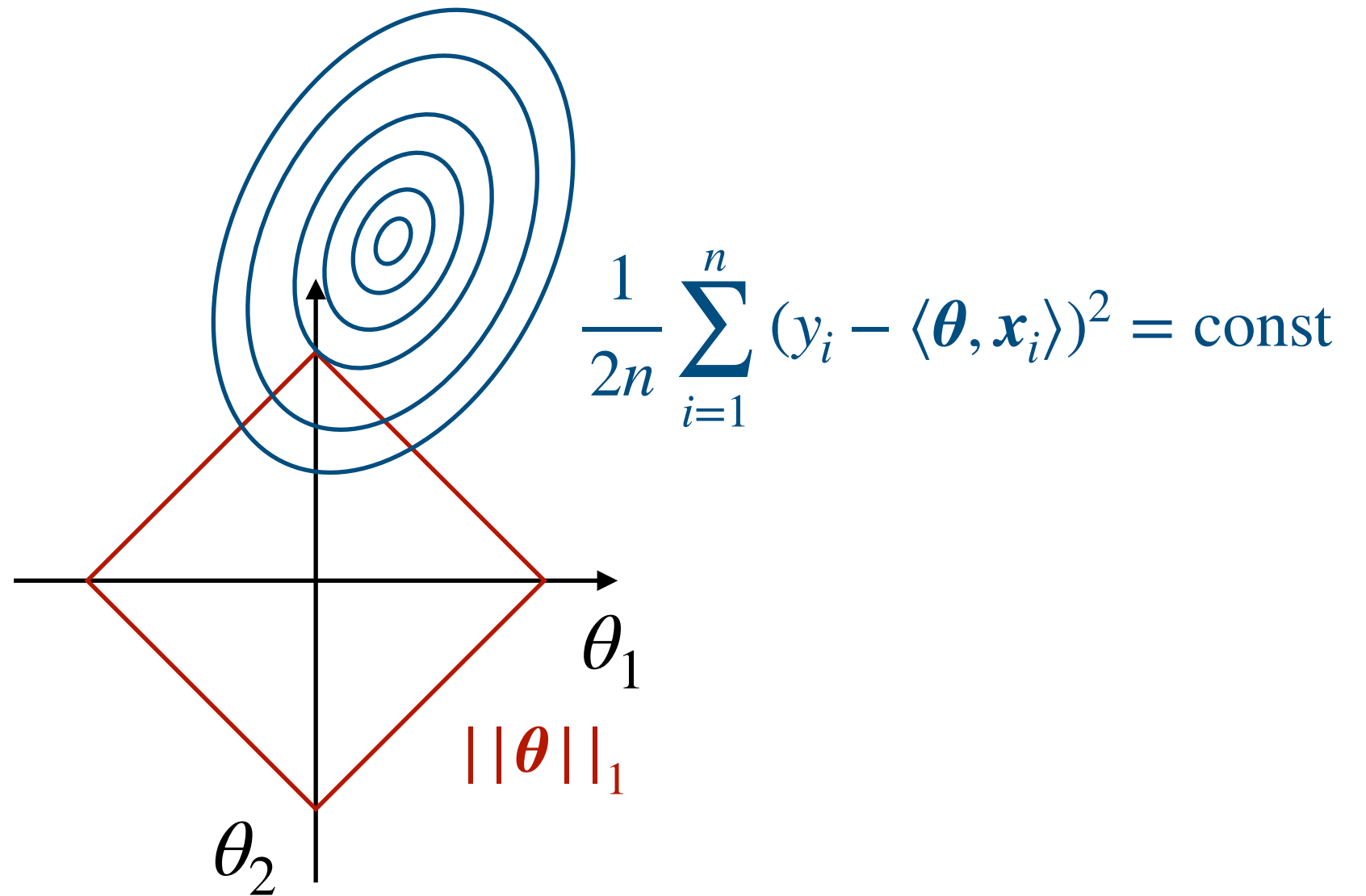
Moreover, this is a **convex** problem.

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

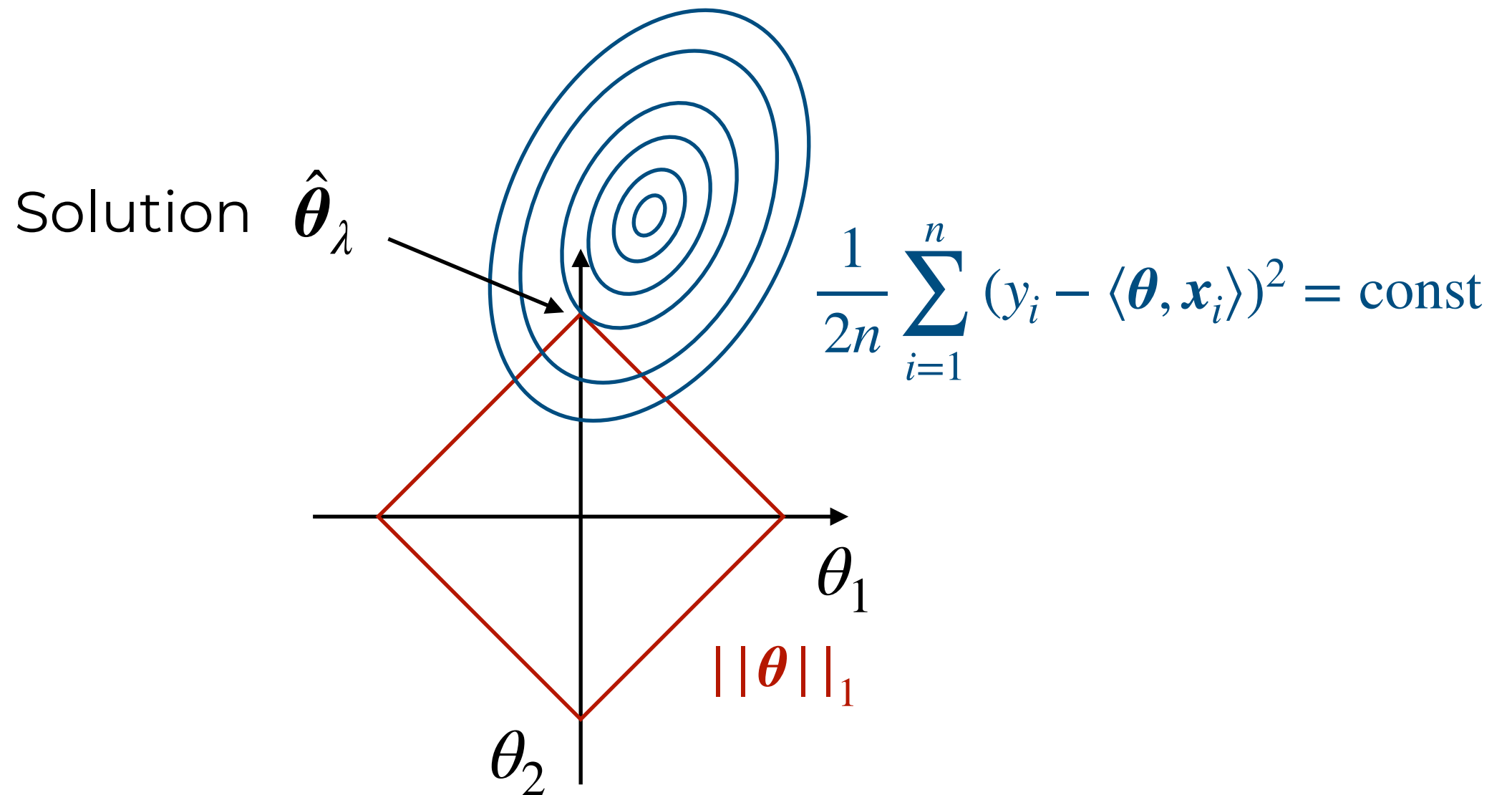
# LASSO: visualisation



# LASSO: visualisation



# LASSO: visualisation



Sharper corners favours sparser solutions!

# LASSO: orthogonal covariates

Again, we can get intuition by looking at the orthogonal covariate case:

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

# LASSO: orthogonal covariates

Again, we can get intuition by looking at the orthogonal covariate case:

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

Following exactly the same steps from before, in this case we need to solve the following coordinate wise problem:

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

# LASSO: orthogonal covariates

Again, we can get intuition by looking at the orthogonal covariate case:

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

Following exactly the same steps from before, in this case we need to solve the following coordinate wise problem:

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

As before, we note that:

$$L(\theta_j) = \begin{cases} \frac{1}{2n} (z_j - \theta_j)^2 + \lambda \theta_j & \text{for } \theta_j > 0 \quad \text{(a)} \\ \frac{z_j^2}{2n} & \text{for } \theta_j = 0 \quad \text{(b)} \\ \frac{1}{2n} (z_j - \theta_j)^2 - \lambda \theta_j & \text{for } \theta_j < 0 \quad \text{(c)} \end{cases}$$



# LASSO: orthogonal covariates

$$L(\theta_j) = \begin{cases} \frac{1}{2n}(z_j - \theta_j)^2 + \lambda\theta_j & \text{for } \theta_j > 0 \quad \text{(a)} \\ \frac{z_j^2}{2n} & \text{for } \theta_j = 0 \quad \text{(b)} \\ \frac{1}{2n}(z_j - \theta_j)^2 - \lambda\theta_j & \text{for } \theta_j < 0 \quad \text{(c)} \end{cases}$$

In case (a), solution is:  $\theta_j = z_j - n\lambda$  valid for  $z_j > n\lambda$

# LASSO: orthogonal covariates

$$L(\theta_j) = \begin{cases} \frac{1}{2n}(z_j - \theta_j)^2 + \lambda\theta_j & \text{for } \theta_j > 0 \quad \text{(a)} \\ \frac{z_j^2}{2n} & \text{for } \theta_j = 0 \quad \text{(b)} \\ \frac{1}{2n}(z_j - \theta_j)^2 - \lambda\theta_j & \text{for } \theta_j < 0 \quad \text{(c)} \end{cases}$$

In case (a), solution is:  $\theta_j = z_j - n\lambda$  valid for  $z_j > n\lambda$

In case (b), solution is:  $\theta_j = 0$

# LASSO: orthogonal covariates

$$L(\theta_j) = \begin{cases} \frac{1}{2n}(z_j - \theta_j)^2 + \lambda\theta_j & \text{for } \theta_j > 0 \quad \text{(a)} \\ \frac{z_j^2}{2n} & \text{for } \theta_j = 0 \quad \text{(b)} \\ \frac{1}{2n}(z_j - \theta_j)^2 - \lambda\theta_j & \text{for } \theta_j < 0 \quad \text{(c)} \end{cases}$$

In case (a), solution is:  $\theta_j = z_j - n\lambda$  valid for  $z_j > n\lambda$

In case (b), solution is:  $\theta_j = 0$

In case (c), solution is:  $\theta_j = z_j + n\lambda$  valid for  $z_j > -n\lambda$

# LASSO: orthogonal covariates

$$L(\theta_j) = \begin{cases} \frac{1}{2n}(z_j - \theta_j)^2 + \lambda\theta_j & \text{for } \theta_j > 0 \quad \text{(a)} \\ \frac{z_j^2}{2n} & \text{for } \theta_j = 0 \quad \text{(b)} \\ \frac{1}{2n}(z_j - \theta_j)^2 - \lambda\theta_j & \text{for } \theta_j < 0 \quad \text{(c)} \end{cases}$$

In case (a), solution is:  $\theta_j = z_j - n\lambda$  valid for  $z_j > n\lambda$

In case (b), solution is:  $\theta_j = 0$

In case (c), solution is:  $\theta_j = z_j + n\lambda$  valid for  $z_j > -n\lambda$

Putting together:  $\theta_j = \begin{cases} z_j - \text{sign}(z_j)n\lambda & \text{for } |z_j| > \lambda \\ 0 & \text{for } |z_j| \in [-\lambda, \lambda] \end{cases}$  Soft-thresholding function

# LASSO: orthogonal covariates

Putting together, the solution of the LASSO problem:

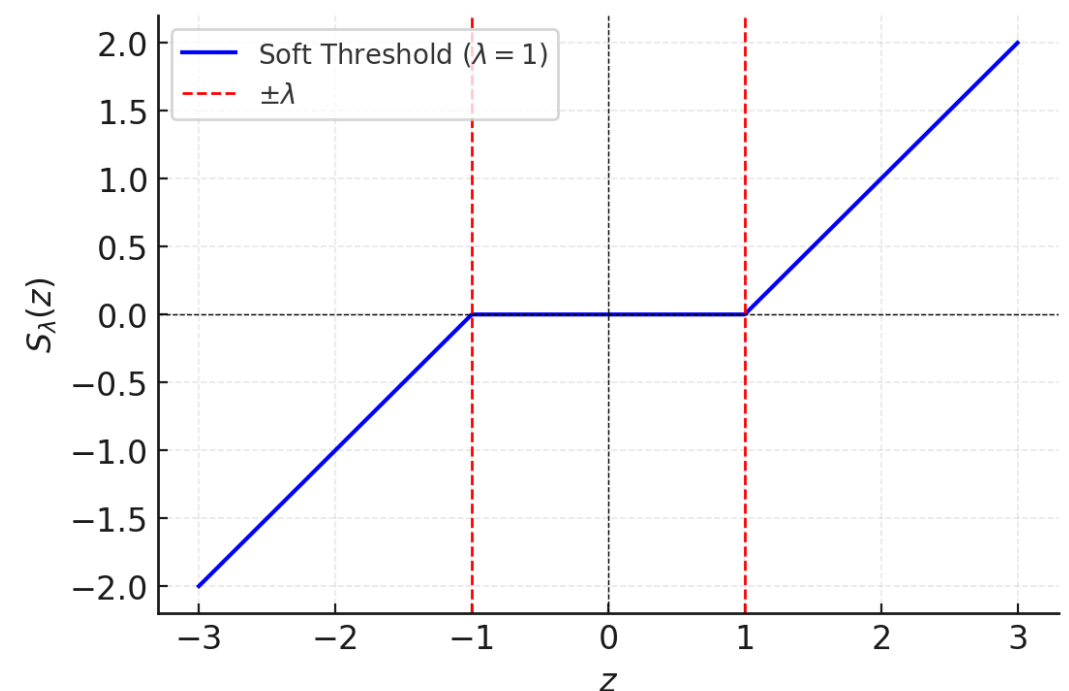
$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \frac{1}{2n} \sum_{i=1}^n (y_i - \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle)^2 + \lambda ||\boldsymbol{\theta}||_1$$

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

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

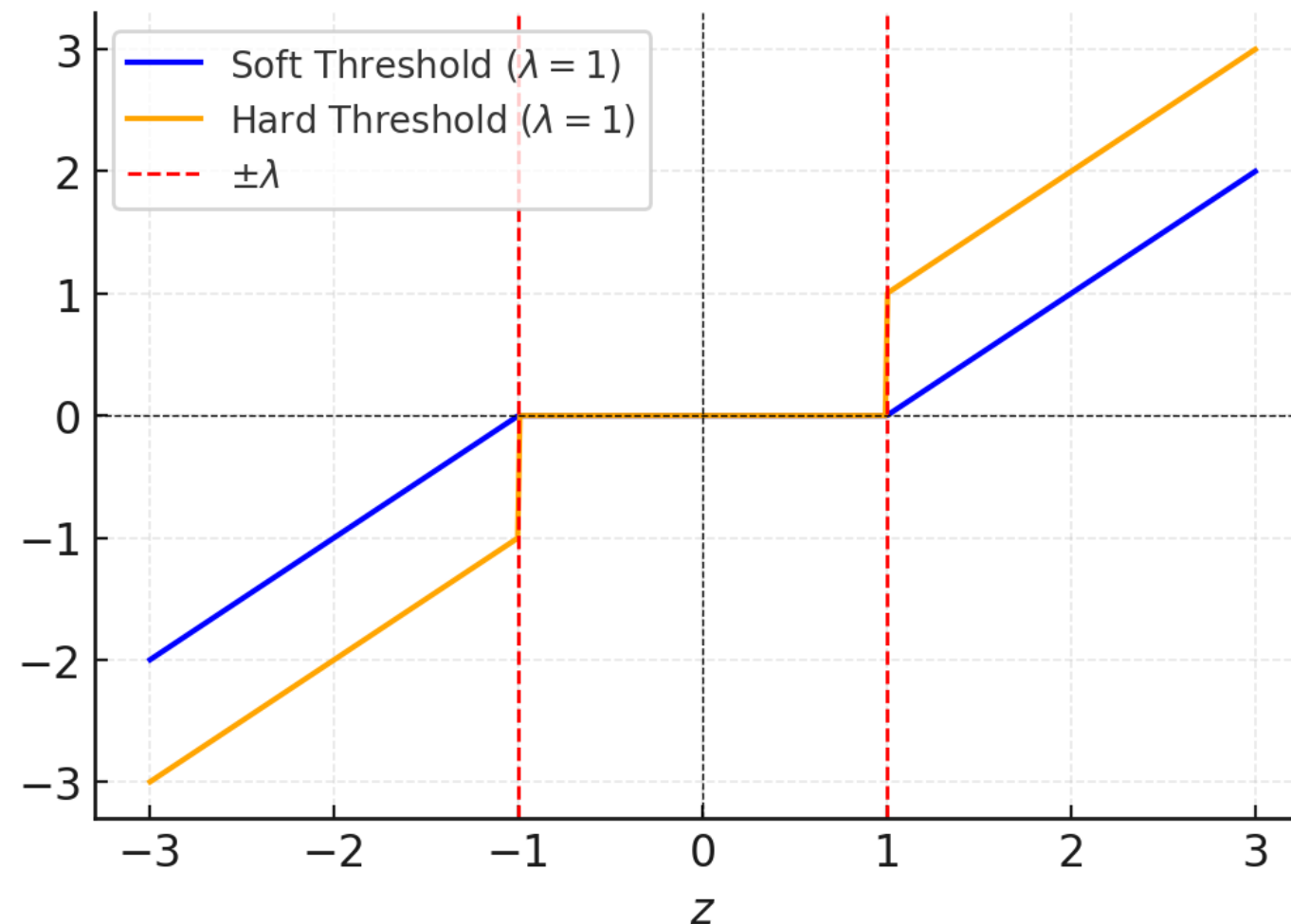
Where:

$$S_\lambda(z) = \begin{cases} z - \text{sign}(z)\lambda & \text{if } |z| > \lambda \\ 0 & \text{if } |z| < \lambda \end{cases}$$



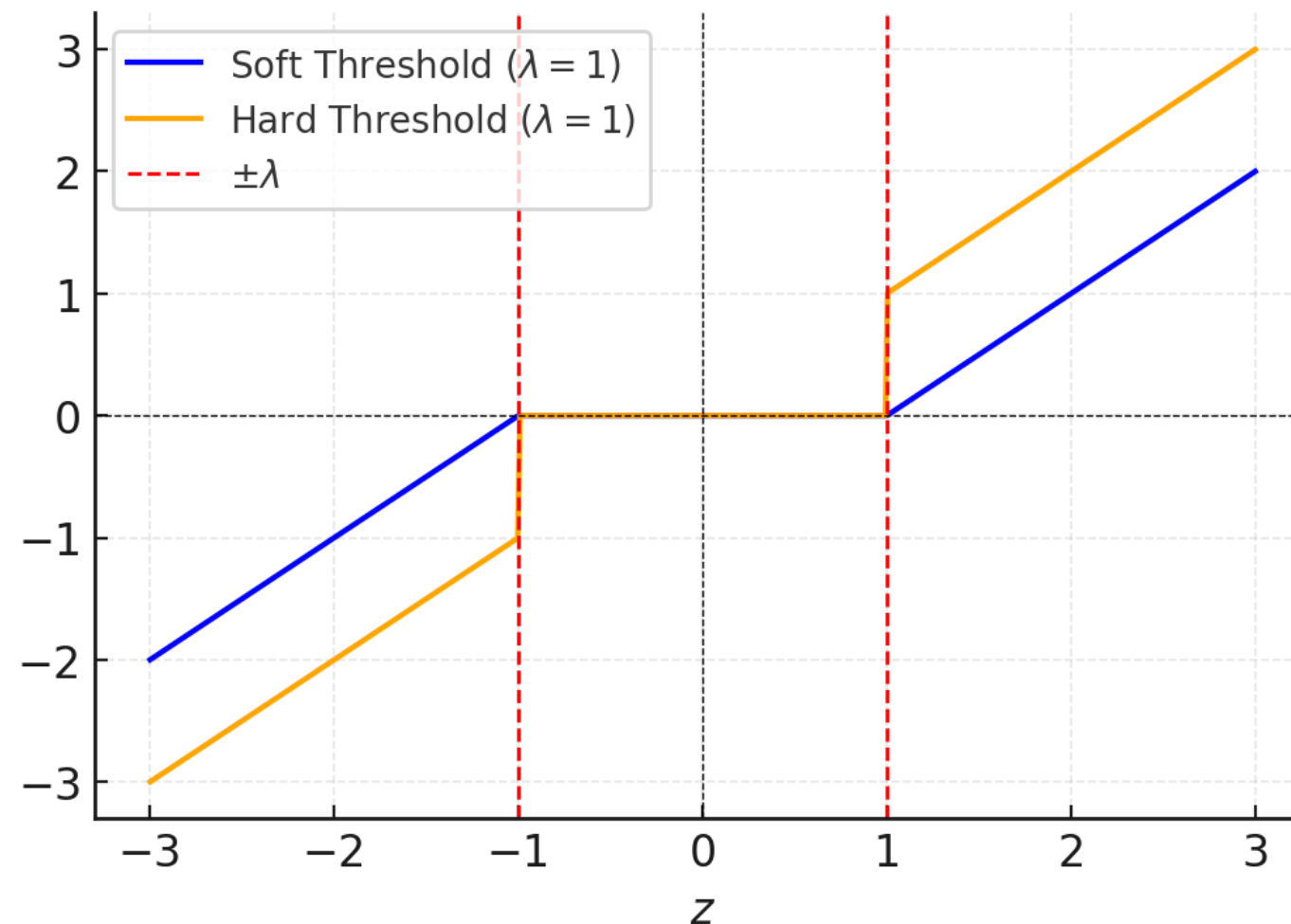
# BSS vs. LASSO

It is instructive to compare the BSS and LASSO solutions in the orthogonal covariate case



# BSS vs. LASSO

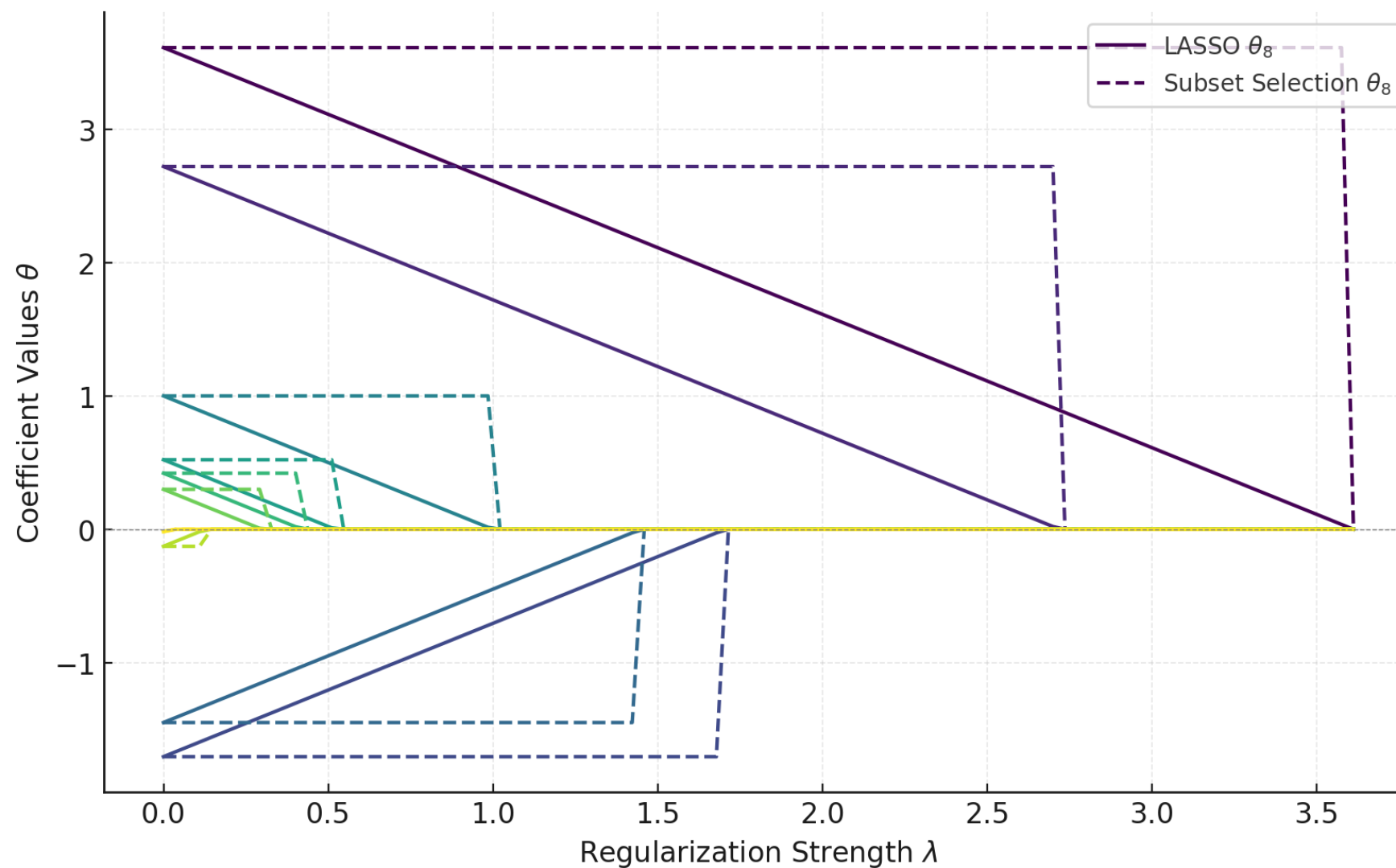
It is instructive to compare the BSS and LASSO solutions in the orthogonal covariate case



- Key similarity: both solutions induce sparsity
- Key differences: LASSO is convex and induce shrinkage (e.g.  $z - \lambda$  for  $z > \lambda$ )

# BSS vs. LASSO

$n = 20$        $d = 10$        $y_i = \langle \boldsymbol{\theta}_\star, \mathbf{x}_i \rangle + \varepsilon_i$        $\varepsilon_i \sim \mathcal{N}(0,1)$        $\mathbf{X}^\top \mathbf{X} = \mathbf{I}_{10}$ ,       $\boldsymbol{\theta}_\star$  is 5-sparse



- BSS is discontinuous
- LASSO is piece-wise continuous



For general design, non-zero path not simply a line



# LASSO: beyond orthogonal

Again, when the covariates are not orthogonal, an explicit solution for the LASSO is not available. Nevertheless, we can partially characterise it.

# LASSO: beyond orthogonal

Again, when the covariates are not orthogonal, an explicit solution for the LASSO 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 LASSO solution



For this to be well-defined, assume  $\hat{\theta}_{\lambda}$  unique

# LASSO: beyond orthogonal

Again, when the covariates are not orthogonal, an explicit solution for the LASSO 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 LASSO solution



For this to be well-defined, assume  $\hat{\theta}_\lambda$  unique

- Denote:
- $\hat{\theta}_S \in \mathbb{R}^{|S|}$  the non-zero entries of  $\hat{\theta}_\lambda \in \mathbb{R}^d$
  - $X_S \in \mathbb{R}^{n \times |S|}$  the corresponding covariates
  - $s_S = \text{sign}(\hat{\theta}_\lambda) \in \{-1, +1\}^{|S|}$  the signs.

# LASSO: beyond orthogonal

Again, when the covariates are not orthogonal, an explicit solution for the LASSO 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 LASSO solution



For this to be well-defined, assume  $\hat{\theta}_\lambda$  unique

- Denote:
- $\hat{\theta}_S \in \mathbb{R}^{|S|}$  the non-zero entries of  $\hat{\theta}_\lambda \in \mathbb{R}^d$
  - $X_S \in \mathbb{R}^{n \times |S|}$  the corresponding covariates
  - $s_S = \text{sign}(\hat{\theta}_\lambda) \in \{-1, +1\}^{|S|}$  the signs.

Then, the by the optimality condition  $\hat{\theta}_S \in \mathbb{R}^{|S|}$  satisfies:

$$X_S^\top (\mathbf{y} - X_S \hat{\theta}_S) = n\lambda s_S$$

# LASSO: beyond orthogonal

Then, the by the optimality condition  $\hat{\boldsymbol{\theta}}_S \in \mathbb{R}^{|S|}$  satisfies:

$$\mathbf{X}_S^\top (\mathbf{y} - \mathbf{X}_S \hat{\boldsymbol{\theta}}_S) = n\lambda \mathbf{s}_S$$

Therefore, the LASSO solution satisfies:

$$\hat{\boldsymbol{\theta}}_S = (\mathbf{X}_S^\top \mathbf{X}_S)^{-1} (\mathbf{X}_S^\top \mathbf{y} - n\lambda \mathbf{s}_S) \quad \hat{\boldsymbol{\theta}}_{-S} = \mathbf{0}_{d-|S|}$$



It can be shown uniqueness imply  $\mathbf{X}_S^\top \mathbf{X}_S \succ 0$

# LASSO: beyond orthogonal

Then, the by the optimality condition  $\hat{\boldsymbol{\theta}}_S \in \mathbb{R}^{|S|}$  satisfies:

$$\mathbf{X}_S^\top (\mathbf{y} - \mathbf{X}_S \hat{\boldsymbol{\theta}}_S) = n\lambda \mathbf{s}_S$$

Therefore, the LASSO solution satisfies:

$$\hat{\boldsymbol{\theta}}_S = \underbrace{(\mathbf{X}_S^\top \mathbf{X}_S)^{-1}}_{\text{OLS}} \underbrace{(\mathbf{X}_S^\top \mathbf{y} - n\lambda \mathbf{s}_S)}_{\text{Shrinkage}} \quad \hat{\boldsymbol{\theta}}_{-S} = \mathbf{0}_{d-|S|}$$



It can be shown uniqueness imply  $\mathbf{X}_S^\top \mathbf{X}_S \succ 0$

# LASSO: beyond orthogonal

Then, the by the optimality condition  $\hat{\boldsymbol{\theta}}_S \in \mathbb{R}^{|S|}$  satisfies:

$$\mathbf{X}_S^\top (\mathbf{y} - \mathbf{X}_S \hat{\boldsymbol{\theta}}_S) = n\lambda \mathbf{s}_S$$

Therefore, the LASSO solution satisfies:

$$\hat{\boldsymbol{\theta}}_S = \underbrace{(\mathbf{X}_S^\top \mathbf{X}_S)^{-1}}_{\text{OLS}} \underbrace{(\mathbf{X}_S^\top \mathbf{y} - n\lambda \mathbf{s}_S)}_{\text{Shrinkage}} \quad \hat{\boldsymbol{\theta}}_{-S} = \mathbf{0}_{d-|S|}$$



It can be shown uniqueness imply  $\mathbf{X}_S^\top \mathbf{X}_S \succ 0$

In particular, note that:

$$||\hat{\boldsymbol{\theta}}_\lambda||_1 = \mathbf{s}_S^\top \hat{\boldsymbol{\theta}}_S$$

# LASSO: beyond orthogonal

Then, the by the optimality condition  $\hat{\boldsymbol{\theta}}_S \in \mathbb{R}^{|S|}$  satisfies:

$$\mathbf{X}_S^\top (\mathbf{y} - \mathbf{X}_S \hat{\boldsymbol{\theta}}_S) = n\lambda \mathbf{s}_S$$

Therefore, the LASSO solution satisfies:

$$\hat{\boldsymbol{\theta}}_S = \underbrace{(\mathbf{X}_S^\top \mathbf{X}_S)^{-1}}_{\text{OLS}} \underbrace{(\mathbf{X}_S^\top \mathbf{y} - n\lambda \mathbf{s}_S)}_{\text{Shrinkage}} \quad \hat{\boldsymbol{\theta}}_{-S} = \mathbf{0}_{d-|S|}$$



It can be shown uniqueness imply  $\mathbf{X}_S^\top \mathbf{X}_S \succ 0$

In particular, note that:

$$||\hat{\boldsymbol{\theta}}_\lambda||_1 = \mathbf{s}_S^\top \hat{\boldsymbol{\theta}}_S = \mathbf{s}_S^\top (\mathbf{X}_S^\top \mathbf{X}_S)^{-1} \mathbf{X}_S^\top \mathbf{y} - n\lambda \mathbf{s}_S^\top (\mathbf{X}_S^\top \mathbf{X}_S)^{-1} \mathbf{s}_S$$



# LASSO: beyond orthogonal

Then, the by the optimality condition  $\hat{\boldsymbol{\theta}}_S \in \mathbb{R}^{|S|}$  satisfies:

$$\mathbf{X}_S^\top (\mathbf{y} - \mathbf{X}_S \hat{\boldsymbol{\theta}}_S) = n\lambda \mathbf{s}_S$$

Therefore, the LASSO solution satisfies:

$$\hat{\boldsymbol{\theta}}_S = \underbrace{(\mathbf{X}_S^\top \mathbf{X}_S)^{-1}}_{\text{OLS}} \underbrace{(\mathbf{X}_S^\top \mathbf{y} - n\lambda \mathbf{s}_S)}_{\text{Shrinkage}} \quad \hat{\boldsymbol{\theta}}_{-S} = \mathbf{0}_{d-|S|}$$



It can be shown uniqueness imply  $\mathbf{X}_S^\top \mathbf{X}_S \succ 0$

In particular, note that:

$$\begin{aligned} ||\hat{\boldsymbol{\theta}}_\lambda||_1 &= \mathbf{s}_S^\top \hat{\boldsymbol{\theta}}_S = \mathbf{s}_S^\top (\mathbf{X}_S^\top \mathbf{X}_S)^{-1} \mathbf{X}_S^\top \mathbf{y} - n\lambda \mathbf{s}_S^\top (\mathbf{X}_S^\top \mathbf{X}_S)^{-1} \mathbf{s}_S \\ &< ||(\mathbf{X}_S^\top \mathbf{X}_S)^{-1} \mathbf{X}_S^\top \mathbf{y}||_1 \end{aligned}$$

# LASSO: beyond orthogonal

Then, the by the optimality condition  $\hat{\boldsymbol{\theta}}_S \in \mathbb{R}^{|S|}$  satisfies:

$$\mathbf{X}_S^\top (\mathbf{y} - \mathbf{X}_S \hat{\boldsymbol{\theta}}_S) = n\lambda \mathbf{s}_S$$

Therefore, the LASSO solution satisfies:

$$\hat{\boldsymbol{\theta}}_S = \underbrace{(\mathbf{X}_S^\top \mathbf{X}_S)^{-1}}_{\text{OLS}} \underbrace{(\mathbf{X}_S^\top \mathbf{y} - n\lambda \mathbf{s}_S)}_{\text{Shrinkage}} \quad \hat{\boldsymbol{\theta}}_{-S} = \mathbf{0}_{d-|S|}$$



It can be shown uniqueness imply  $\mathbf{X}_S^\top \mathbf{X}_S \succ 0$

In particular, note that:

$$\begin{aligned} ||\hat{\boldsymbol{\theta}}_\lambda||_1 &= \mathbf{s}_S^\top \hat{\boldsymbol{\theta}}_S = \mathbf{s}_S^\top (\mathbf{X}_S^\top \mathbf{X}_S)^{-1} \mathbf{X}_S^\top \mathbf{y} - n\lambda \mathbf{s}_S^\top (\mathbf{X}_S^\top \mathbf{X}_S)^{-1} \mathbf{s}_S \\ &< ||(\mathbf{X}_S^\top \mathbf{X}_S)^{-1} \mathbf{X}_S^\top \mathbf{y}||_1 \quad ||\hat{\boldsymbol{\theta}}_{\text{LASSO}}||_1 \leq ||\hat{\boldsymbol{\theta}}_{\text{OLS}}||_1 !!! \end{aligned}$$

# LASSO in practice

Beyond the orthogonal case, the LASSO problem:

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

does not admit an explicit solution. How do we do in practice?

# LASSO in practice

Beyond the orthogonal case, the LASSO problem:

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

does not admit an explicit solution. How do we do in practice?



LASSO = OLS +  $\ell_1$  penalty

Idea: alternate between these two.

# LASSO in practice

Beyond the orthogonal case, the LASSO problem:

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

does not admit an explicit solution. How do we do in practice?



LASSO = OLS +  $\ell_1$  penalty

Idea: alternate between these two.

Iterative Shrinkage-Thresholding Algorithm (ISTA)

$$\boldsymbol{\theta}^{k+1} = S_{\eta\lambda} \left( \boldsymbol{\theta}^k + \frac{\eta}{n} \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}^k) \right)$$

# LASSO in practice

$$\lambda = 0.5$$

$$n = 10$$

$$d = 2$$

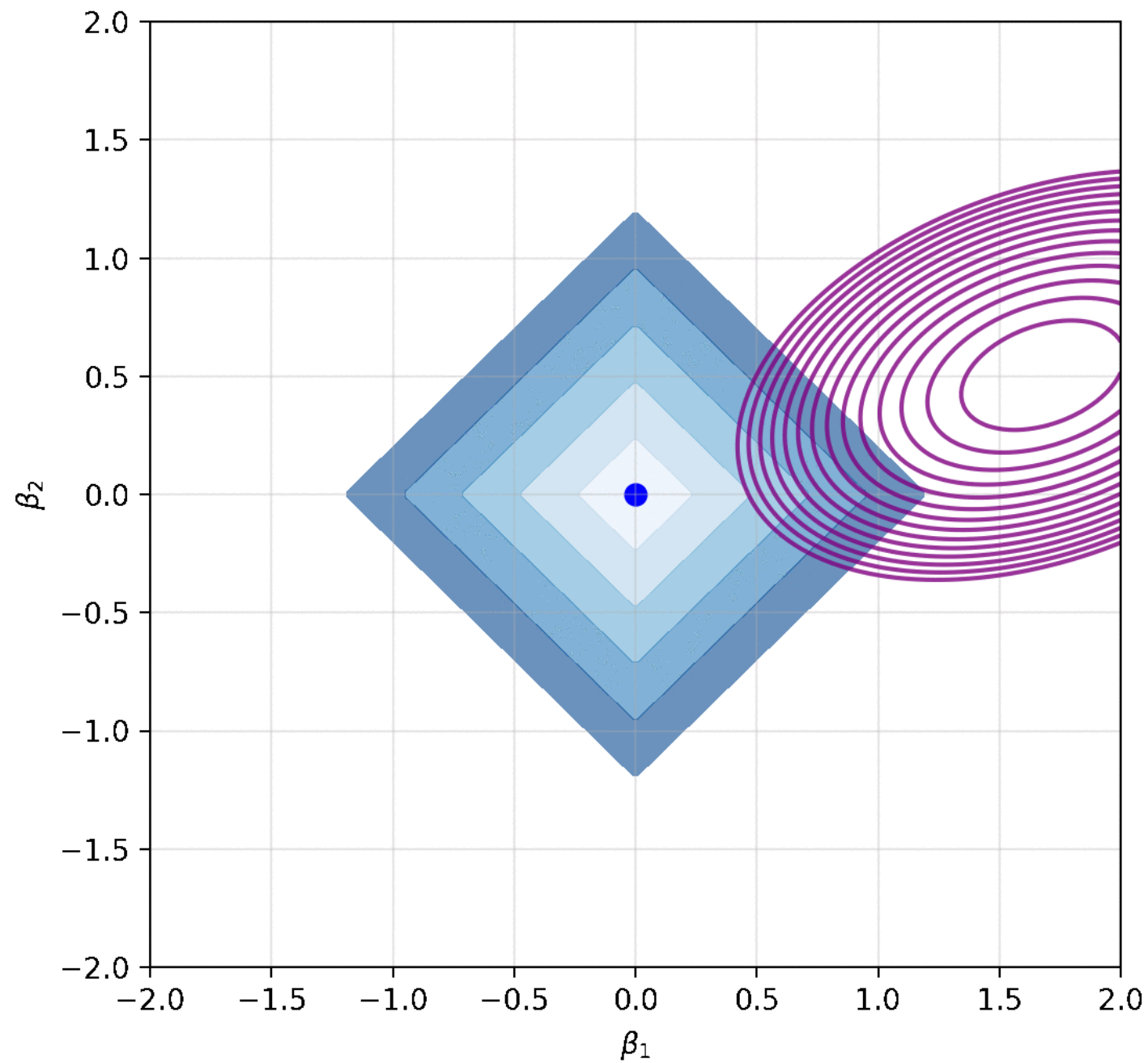
$$y_i = \langle \boldsymbol{\theta}_\star, \mathbf{x}_i \rangle + \varepsilon_i$$

$$\mathbf{x}_i \sim \mathcal{N}(0, \mathbf{I}_2)$$

$$\varepsilon_i \sim \mathcal{N}(0, 1)$$

$$\boldsymbol{\theta}_\star = \begin{bmatrix} 1.5 \\ 0 \end{bmatrix}$$

$$\eta = 0.1$$



# LASSO in practice

$$\lambda = 0.1$$

$$n = 10$$

$$d = 2$$

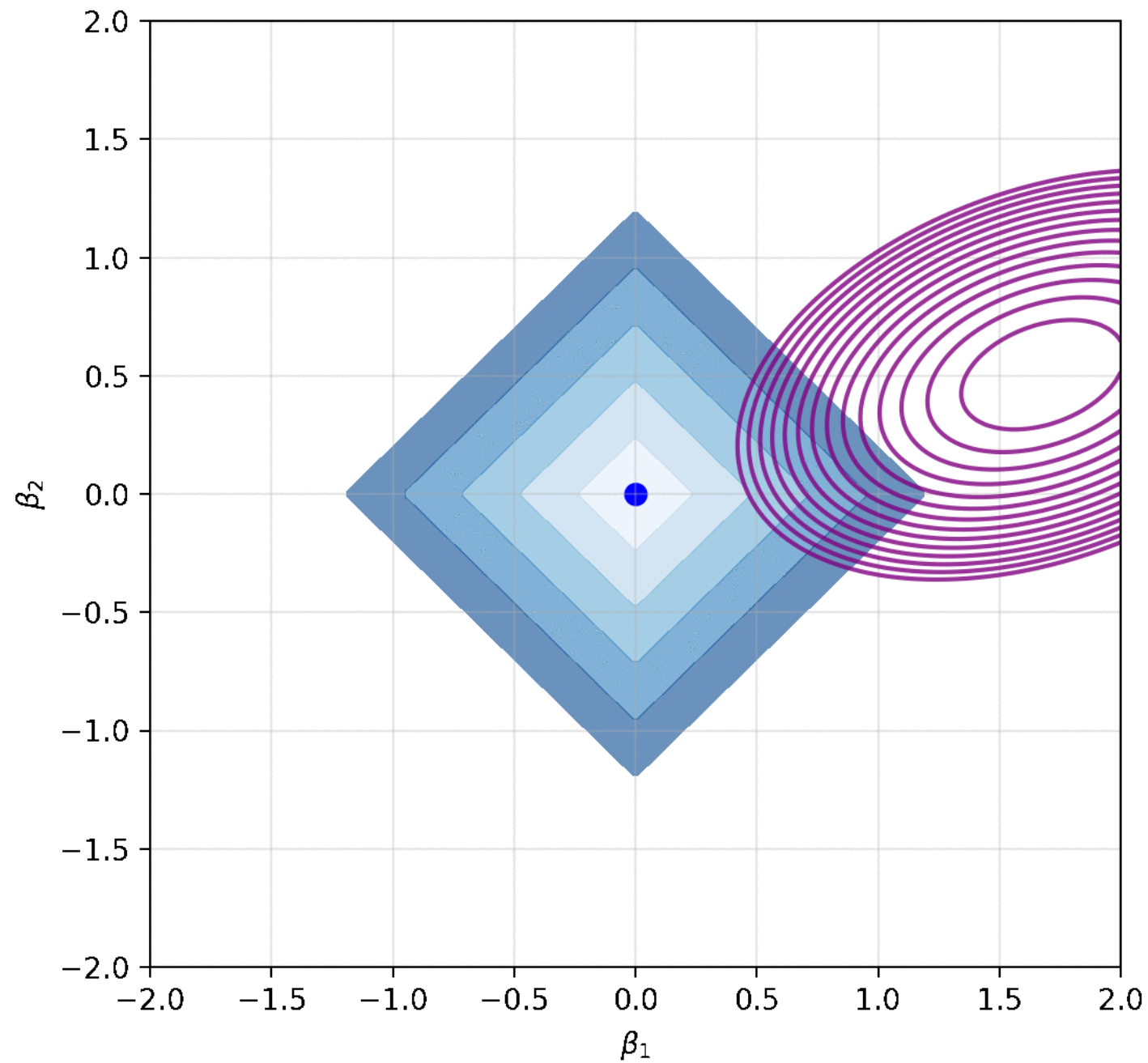
$$y_i = \langle \boldsymbol{\theta}_\star, \mathbf{x}_i \rangle + \varepsilon_i$$

$$\mathbf{x}_i \sim \mathcal{N}(0, \mathbf{I}_2)$$

$$\varepsilon_i \sim \mathcal{N}(0, 1)$$

$$\boldsymbol{\theta}_\star = \begin{bmatrix} 1.5 \\ 0 \end{bmatrix}$$

$$\eta = 0.1$$



# Elastic Net

---

The elastic net algorithm combines ridge with LASSO:

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \frac{1}{2n} \sum_{i=1}^n (y_i - \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle)^2 + \lambda_1 ||\boldsymbol{\theta}||_1 + \frac{\lambda_2}{2} ||\boldsymbol{\theta}||_2^2$$

And is particularly suited to the case where the covariate  $\mathbf{X}$  is badly conditioned.



# Feature maps

---

# Motivation

---

Up to know, our focus has been on parametric functions  $f_{\theta}(\mathbf{x})$  which are linear on both  $\theta \in \mathbb{R}^d$  and  $\mathbf{x} \in \mathbb{R}^d$ .

$$f_{\theta}(\mathbf{x}) = \langle \theta, \mathbf{x} \rangle$$

# Motivation

---

Up to know, our focus has been on parametric functions  $f_{\theta}(\mathbf{x})$  which are linear on both  $\theta \in \mathbb{R}^d$  and  $\mathbf{x} \in \mathbb{R}^d$ .

$$f_{\theta}(\mathbf{x}) = \langle \theta, \mathbf{x} \rangle$$

The main convenience of linear functions is that for convex loss functions, the ERM problem is convex:

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ell(y_i, f_{\theta}(\mathbf{x}))$$

# Motivation

---

Up to know, our focus has been on parametric functions  $f_{\theta}(\mathbf{x})$  which are linear on both  $\theta \in \mathbb{R}^d$  and  $\mathbf{x} \in \mathbb{R}^d$ .

$$f_{\theta}(\mathbf{x}) = \langle \theta, \mathbf{x} \rangle$$

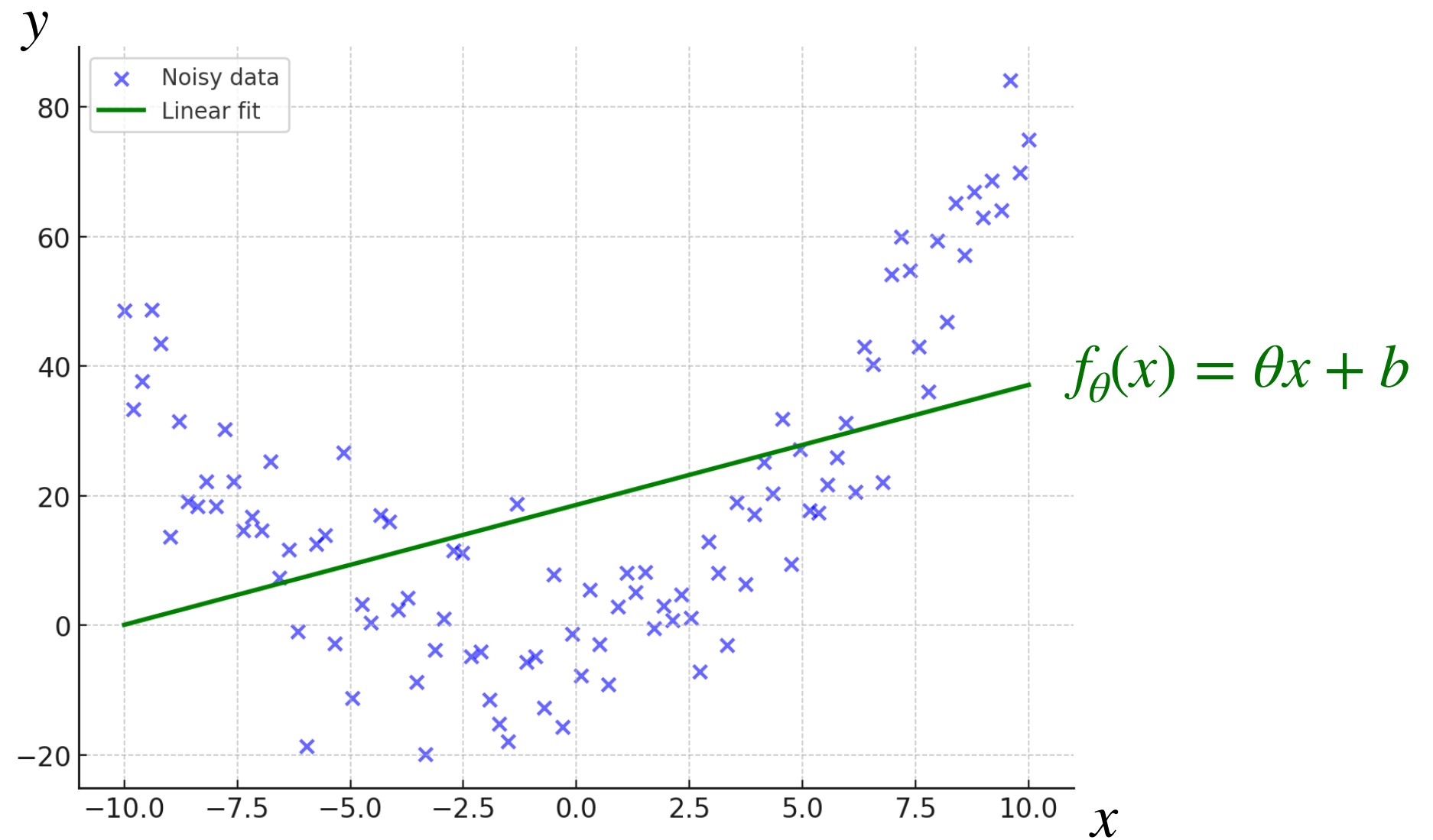
The main convenience of linear functions is that for convex loss functions, the ERM problem is convex:

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ell(y_i, f_{\theta}(\mathbf{x}))$$

But the main drawback is that we can only express linear relationships between the covariates and the labels...

# Motivation

---



# Feature maps

---



Idea: Introduce a **feature map**:

$$\begin{aligned}\boldsymbol{\varphi} : \mathbb{R}^d &\rightarrow \mathbb{R}^p \\ \boldsymbol{x} &\mapsto \boldsymbol{\varphi}(\boldsymbol{x})\end{aligned}$$

And consider a linear predictor in **feature space**:

$$f_{\boldsymbol{\theta}}(\boldsymbol{x}) = \langle \boldsymbol{\theta}, \boldsymbol{\varphi}(\boldsymbol{x}) \rangle$$

# Feature maps

---



Idea: Introduce a **feature map**:

$$\begin{aligned}\boldsymbol{\varphi} : \mathbb{R}^d &\rightarrow \mathbb{R}^p \\ \boldsymbol{x} &\mapsto \boldsymbol{\varphi}(\boldsymbol{x})\end{aligned}$$

And consider a linear predictor in **feature space**:

$$f_{\boldsymbol{\theta}}(\boldsymbol{x}) = \langle \boldsymbol{\theta}, \boldsymbol{\varphi}(\boldsymbol{x}) \rangle$$



- Now we have  $\boldsymbol{\theta} \in \mathbb{R}^p$ .
- $f_{\boldsymbol{\theta}}$  still a linear function of  $\boldsymbol{\theta}$ .
- Typically  $p > d$ .
- More generally, we can consider  $\boldsymbol{\varphi} : \mathcal{X} \rightarrow \mathbb{R}^p$

# Feature maps

---



Idea: Introduce a **feature map**:

$$\begin{aligned}\boldsymbol{\varphi} : \mathbb{R}^d &\rightarrow \mathbb{R}^p \\ \boldsymbol{x} &\mapsto \boldsymbol{\varphi}(\boldsymbol{x})\end{aligned}$$

And consider a linear predictor in **feature space**:

$$f_{\boldsymbol{\theta}}(\boldsymbol{x}) = \langle \boldsymbol{\theta}, \boldsymbol{\varphi}(\boldsymbol{x}) \rangle$$



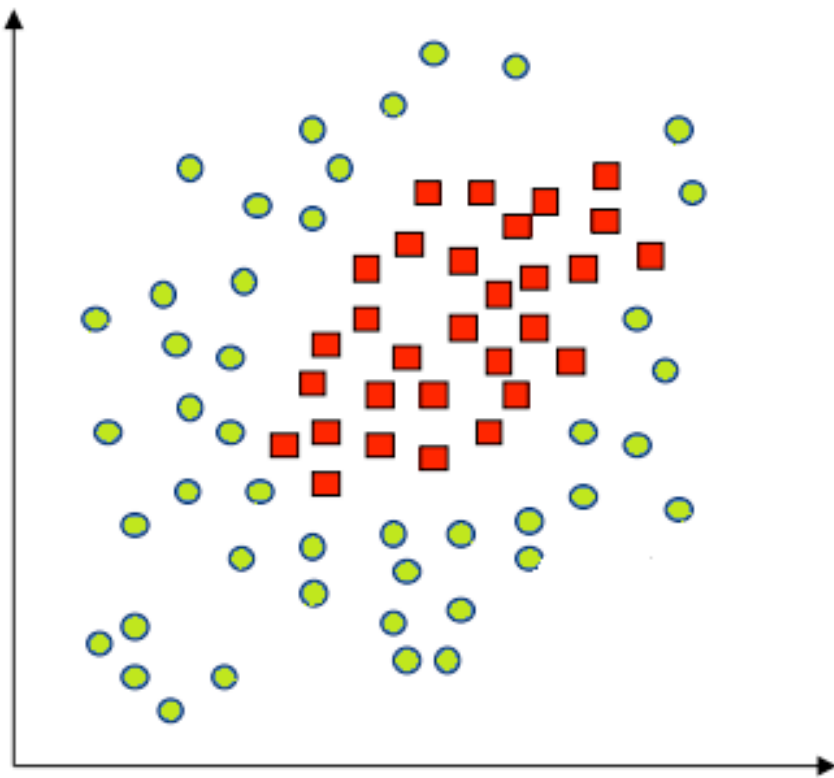
- Now we have  $\boldsymbol{\theta} \in \mathbb{R}^p$ .
- $f_{\boldsymbol{\theta}}$  still a linear function of  $\boldsymbol{\theta}$ .
- Typically  $p > d$ .
- More generally, we can consider  $\boldsymbol{\varphi} : \mathcal{X} \rightarrow \mathbb{R}^p$

Example:  $\mathcal{X}$  a collection of books.




# Feature maps

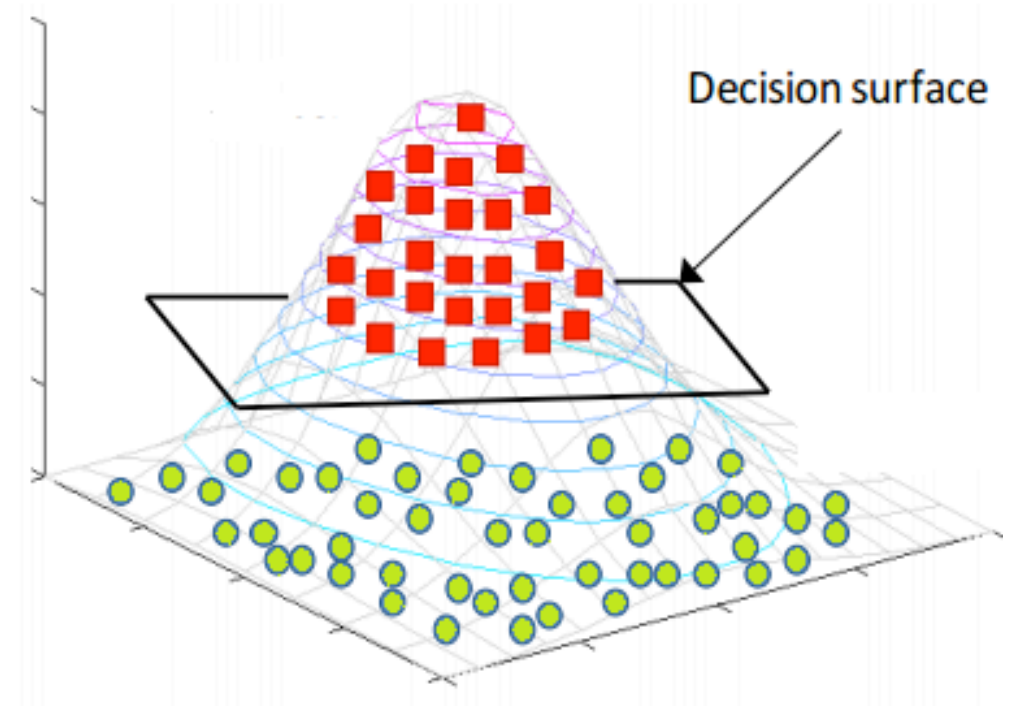
Intuition: Typically easier to linearly separate data in higher-dimensions



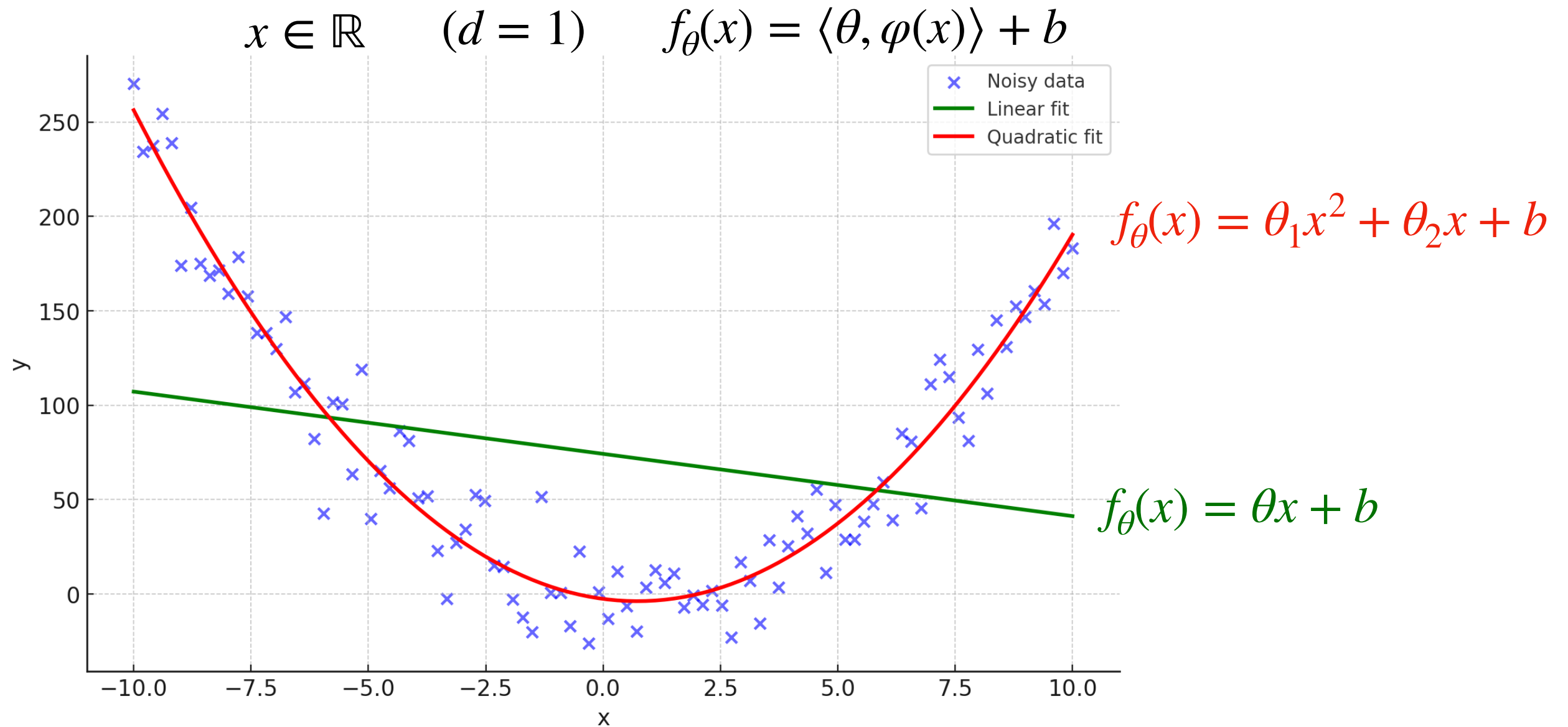
$\varphi$



A horizontal arrow pointing from the 2D scatter plot to the 3D plot, indicating a transformation  $\varphi$ .

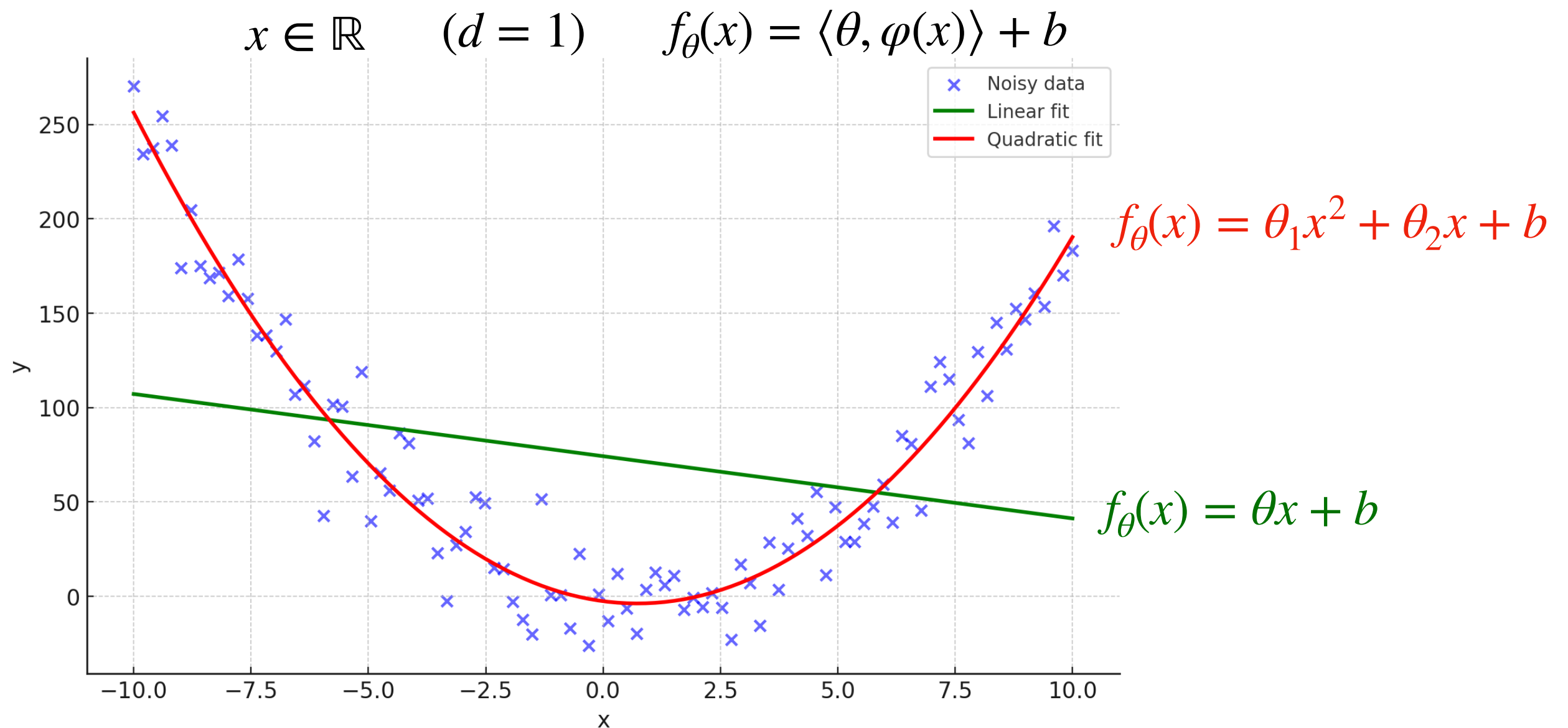


# Examples: quadratic function



Question: what is  $\varphi(x)$ ?

# Examples: quadratic function



Question: what is  $\varphi(x)$ ?

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$$

$$\varphi(x) = \begin{bmatrix} x^2 \\ x \end{bmatrix} \quad (p = 2)$$

# Polynomial regression

More generally, any polynomial of degree  $k \in \mathbb{N}$  over  $\mathbb{R}$

$$p(x) = \sum_{j=1}^k \theta_j x^j + b = \theta_k x^k + \theta_{k-1} x^{k-1} + \dots + \theta_1 x + b$$

# Polynomial regression

More generally, any polynomial of degree  $k \in \mathbb{N}$  over  $\mathbb{R}$

$$p(x) = \sum_{j=1}^k \theta_j x^j + b = \theta_k x^k + \theta_{k-1} x^{k-1} + \dots + \theta_1 x + b$$

Can be written as a linear function in  $\mathbb{R}^k$ :

$$p(x) = \langle \boldsymbol{\theta}, \boldsymbol{\varphi}(x) \rangle + b \qquad \boldsymbol{\varphi}(x) = \begin{bmatrix} x \\ x^2 \\ \vdots \\ x^k \end{bmatrix} \in \mathbb{R}^k$$

# Polynomial regression

More generally, any polynomial of degree  $k \in \mathbb{N}$  over  $\mathbb{R}$

$$p(x) = \sum_{j=1}^k \theta_j x^j + b = \theta_k x^k + \theta_{k-1} x^{k-1} + \dots + \theta_1 x + b$$

Can be written as a linear function in  $\mathbb{R}^k$ :

$$p(x) = \langle \boldsymbol{\theta}, \boldsymbol{\varphi}(x) \rangle + b \quad \boldsymbol{\varphi}(x) = \begin{bmatrix} x \\ x^2 \\ \vdots \\ x^k \end{bmatrix} \in \mathbb{R}^k$$

We can generalise this to degree  $k$  polynomials in  $\mathbb{R}^d$ :

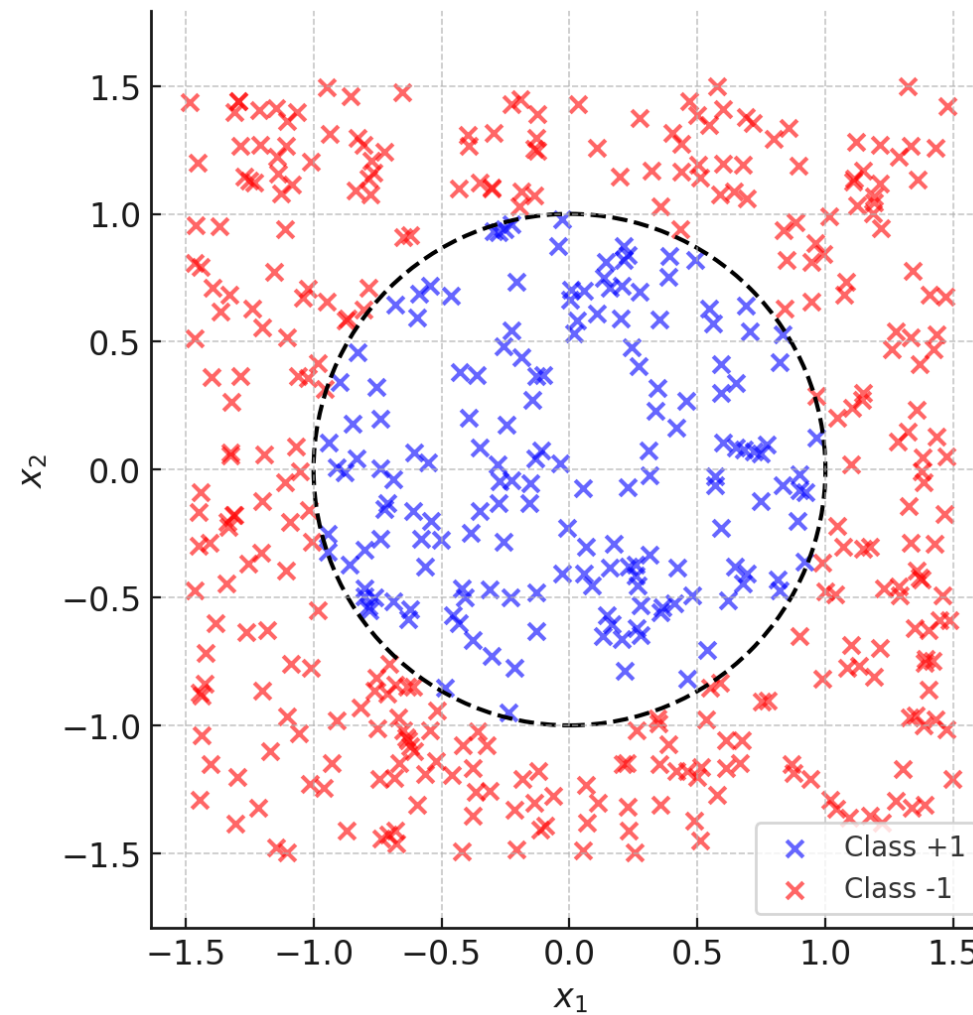
Example  $d = 2$ :

$$p(\mathbf{x}) = \langle \boldsymbol{\theta}, \boldsymbol{\varphi}(\mathbf{x}) \rangle + b \quad \boldsymbol{\varphi}(\mathbf{x}) = \begin{bmatrix} x_1 \\ x_2 \\ x_1^2 \\ x_1 x_2 \\ x_2^2 \end{bmatrix} \in \mathbb{R}^5$$

# Examples: data in circle

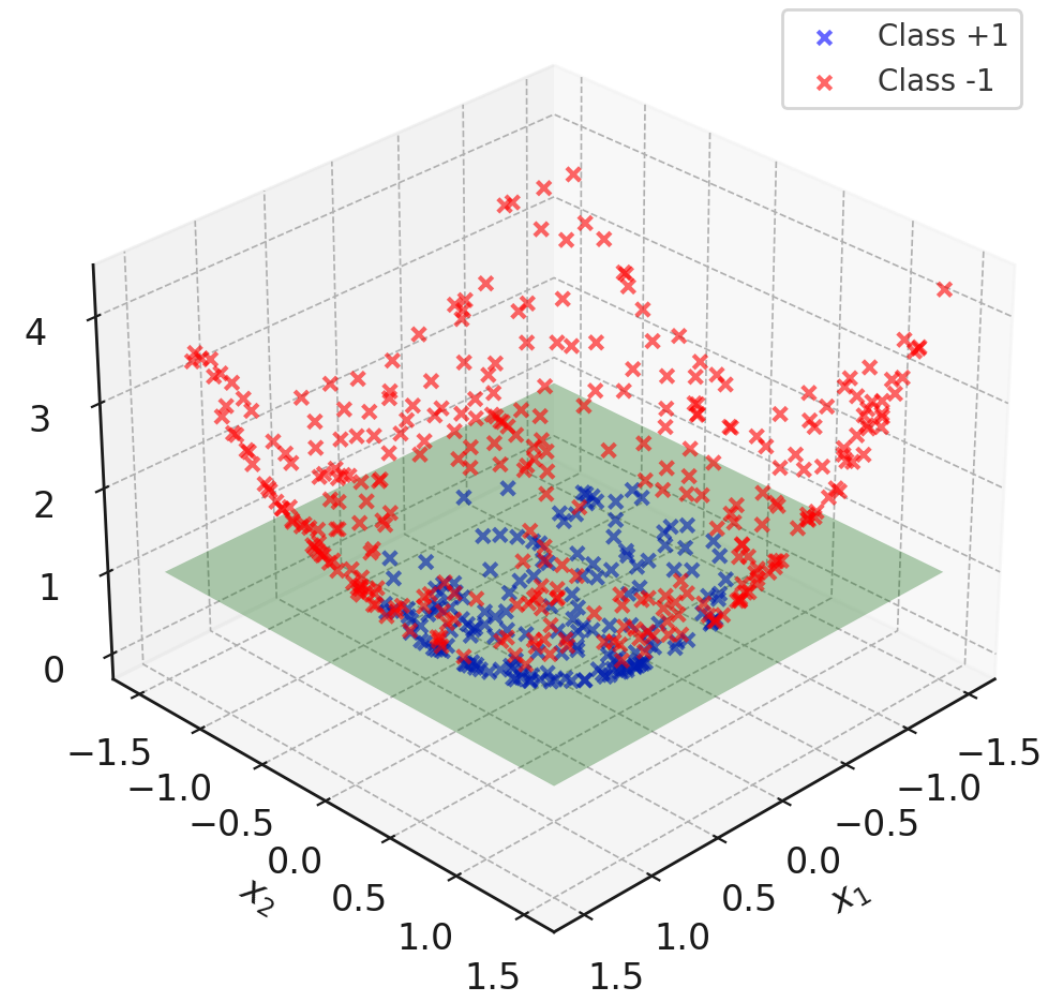
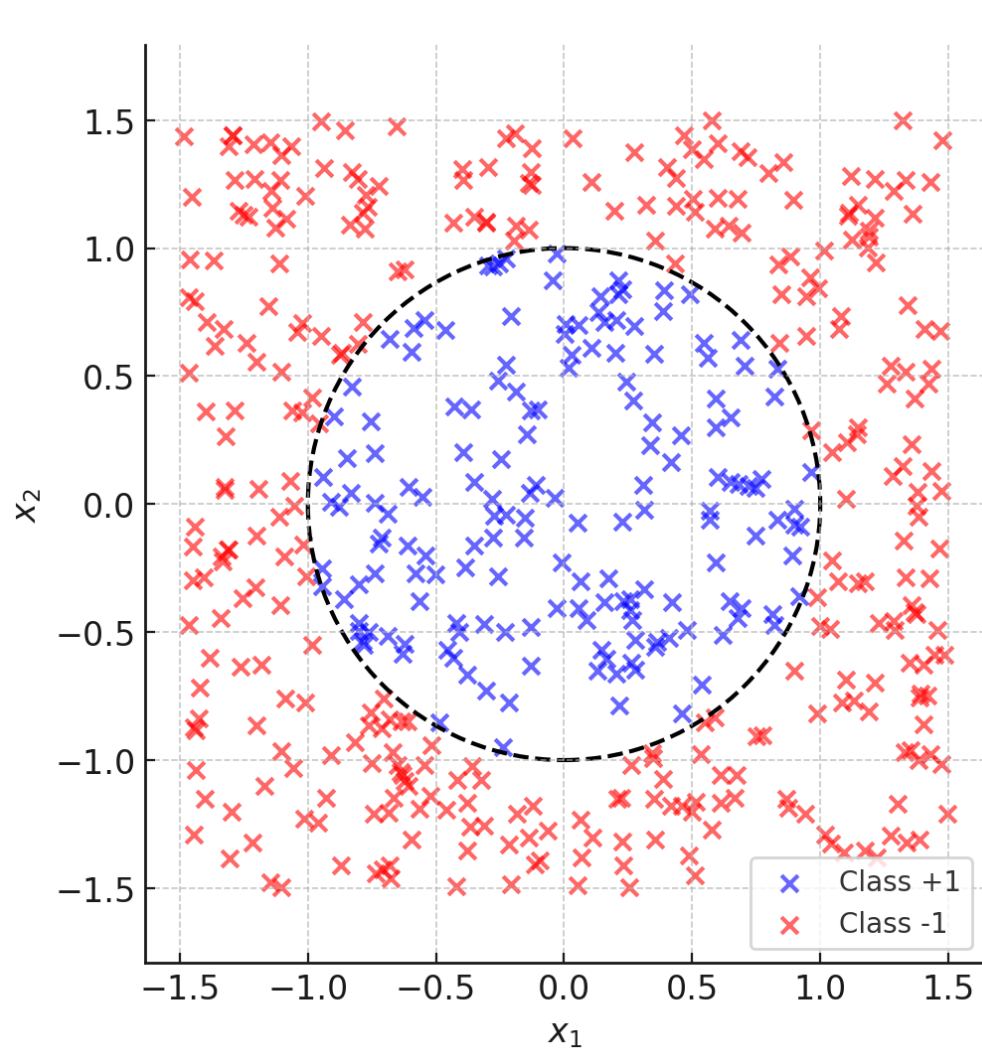
---

$$x \in \mathbb{R}^2 \quad (d = 2) \quad y = \begin{cases} +1 & \text{if } x_1^2 + x_2^2 \leq 1 \\ -1 & \text{if } x_1^2 + x_2^2 > 1 \end{cases}$$



# Examples: data in circle

$$x \in \mathbb{R}^2 \quad (d = 2) \quad y = \begin{cases} +1 & \text{if } x_1^2 + x_2^2 \leq 1 \\ -1 & \text{if } x_1^2 + x_2^2 > 1 \end{cases}$$

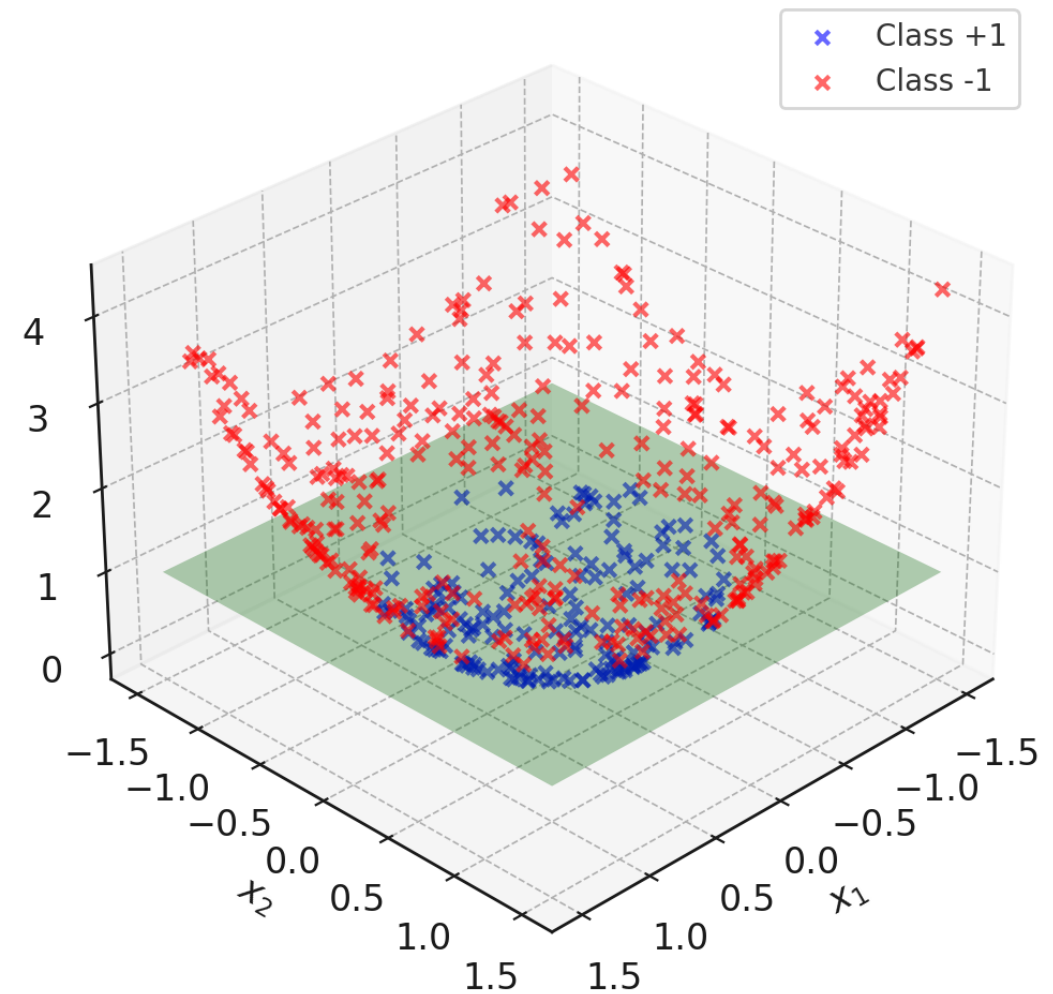
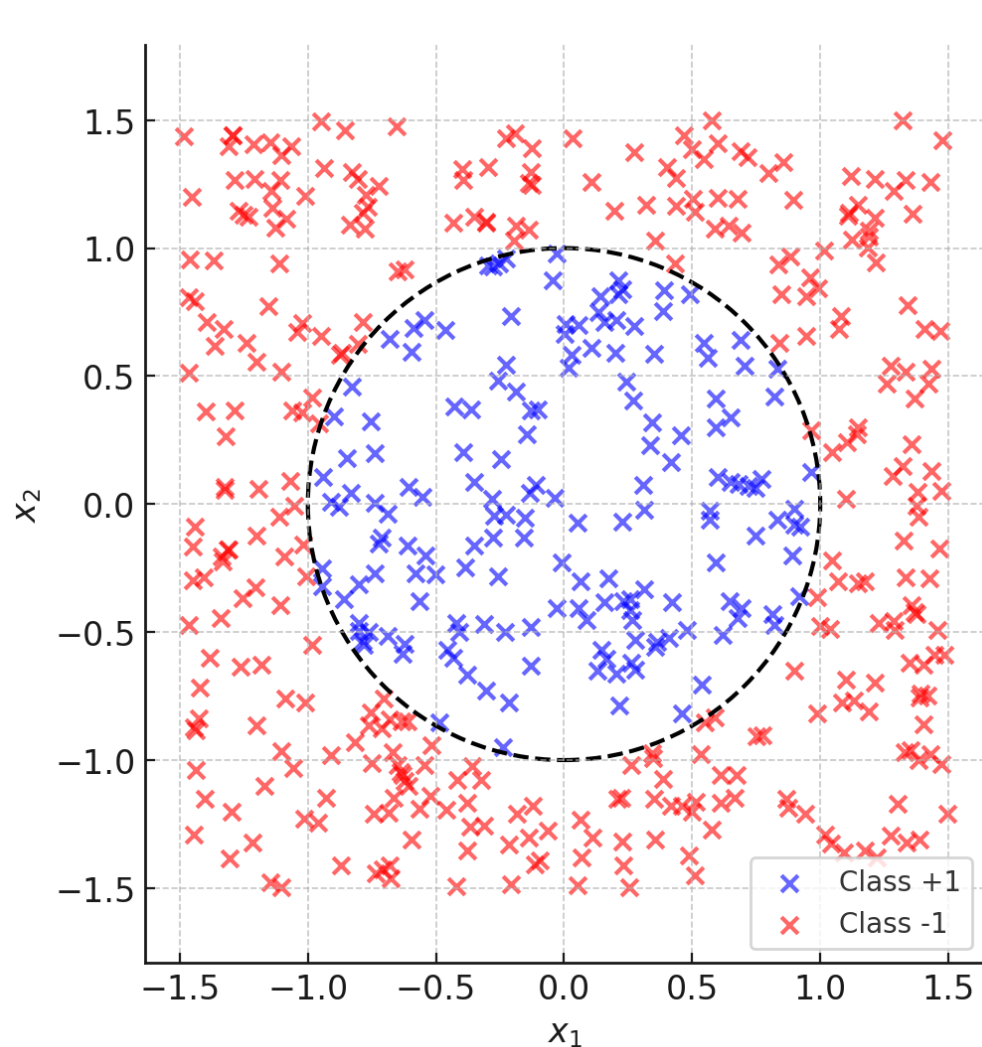


$$\varphi(x) = \begin{bmatrix} x_1 \\ x_2 \\ x_1^2 + x_2^2 \end{bmatrix} \quad (p = 3)$$



# Examples: data in circle

$$x \in \mathbb{R}^2 \quad (d = 2) \quad y = \begin{cases} +1 & \text{if } x_1^2 + x_2^2 \leq 1 \\ -1 & \text{if } x_1^2 + x_2^2 > 1 \end{cases}$$



$$\varphi(x) = \begin{bmatrix} x_1 \\ x_2 \\ x_1^2 + x_2^2 \end{bmatrix}$$

( $p = 3$ )



Not unique!

# Examples: XOR Gaussian mixture

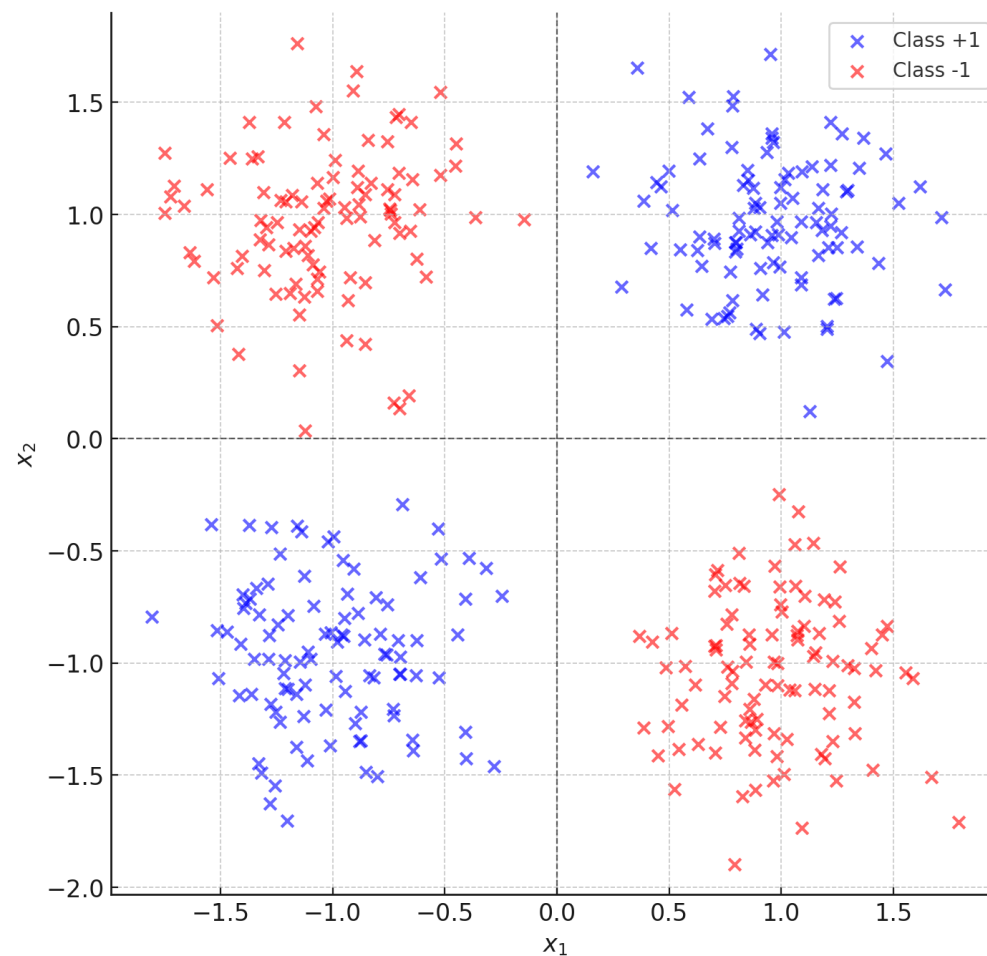
$$x \in \mathbb{R}^2 \quad (d = 2) \quad p(\mathbf{x}) = \frac{1}{4} \sum_{k=1}^4 \mathcal{N}(\boldsymbol{\mu}_k, \mathbf{I}_2)$$

$$\boldsymbol{\mu}_2 = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

$$\boldsymbol{\mu}_1 = \begin{bmatrix} -1 \\ -1 \end{bmatrix}$$

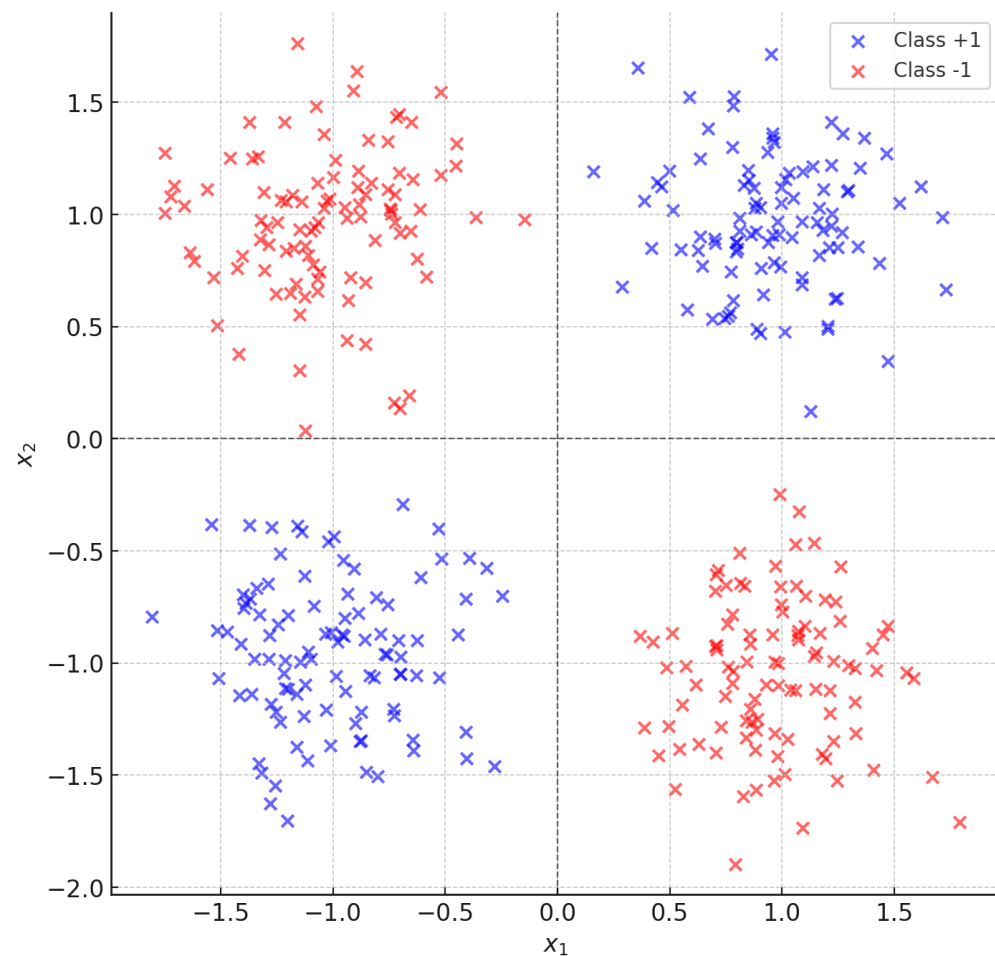
$$\boldsymbol{\mu}_3 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\boldsymbol{\mu}_4 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$



# Examples: XOR Gaussian mixture

$$x \in \mathbb{R}^2 \quad (d = 2) \quad p(x) = \frac{1}{4} \sum_{k=1}^4 \mathcal{N}(\mu_k, I_2)$$



Note that:

$$y = +1$$

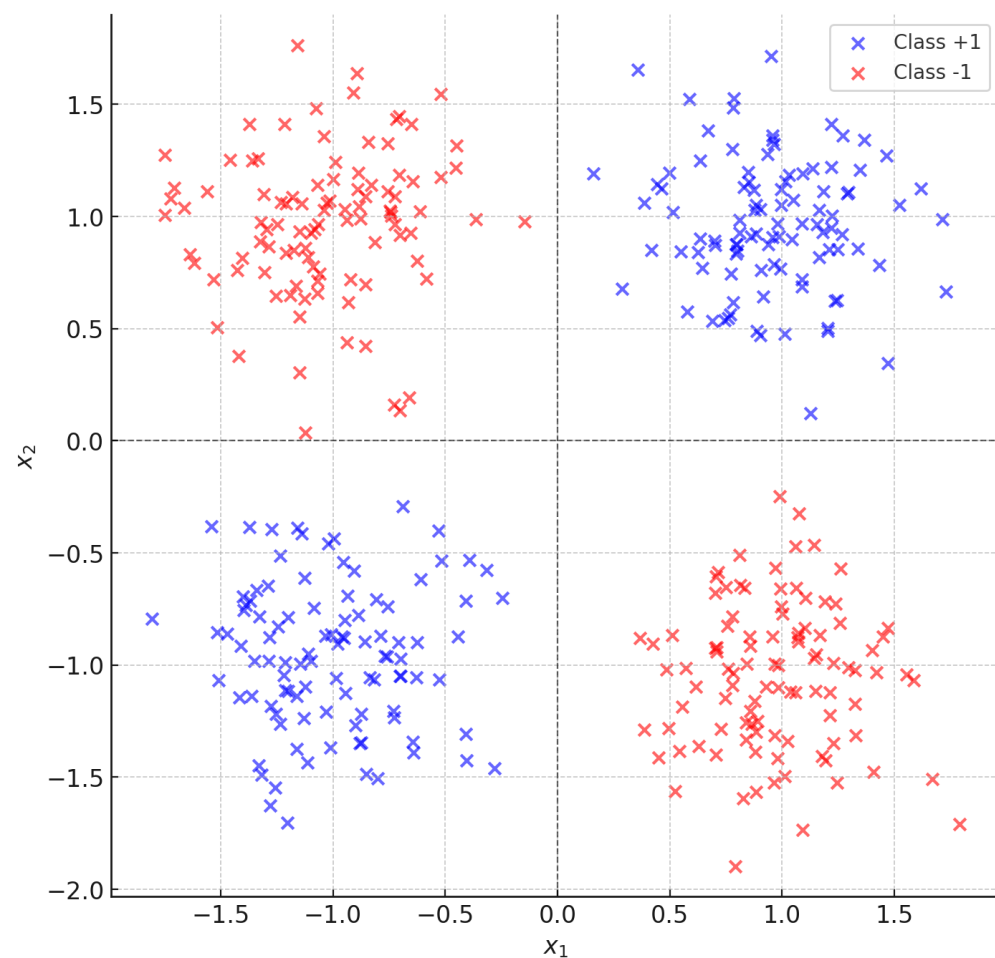
$$x_1, x_2 > 0 \text{ or } x_1, x_2 < 0$$

$$y = -1$$

$$x_1 > 0 \text{ and } x_2 < 0 \text{ or } x_1 < 0 \text{ and } x_2 > 0$$

# Examples: XOR Gaussian mixture

$$x \in \mathbb{R}^2 \quad (d = 2) \quad p(x) = \frac{1}{4} \sum_{k=1}^4 \mathcal{N}(\mu_k, I_2)$$



Note that:

$$y = +1 \quad x_1, x_2 > 0 \text{ or } x_1, x_2 < 0$$

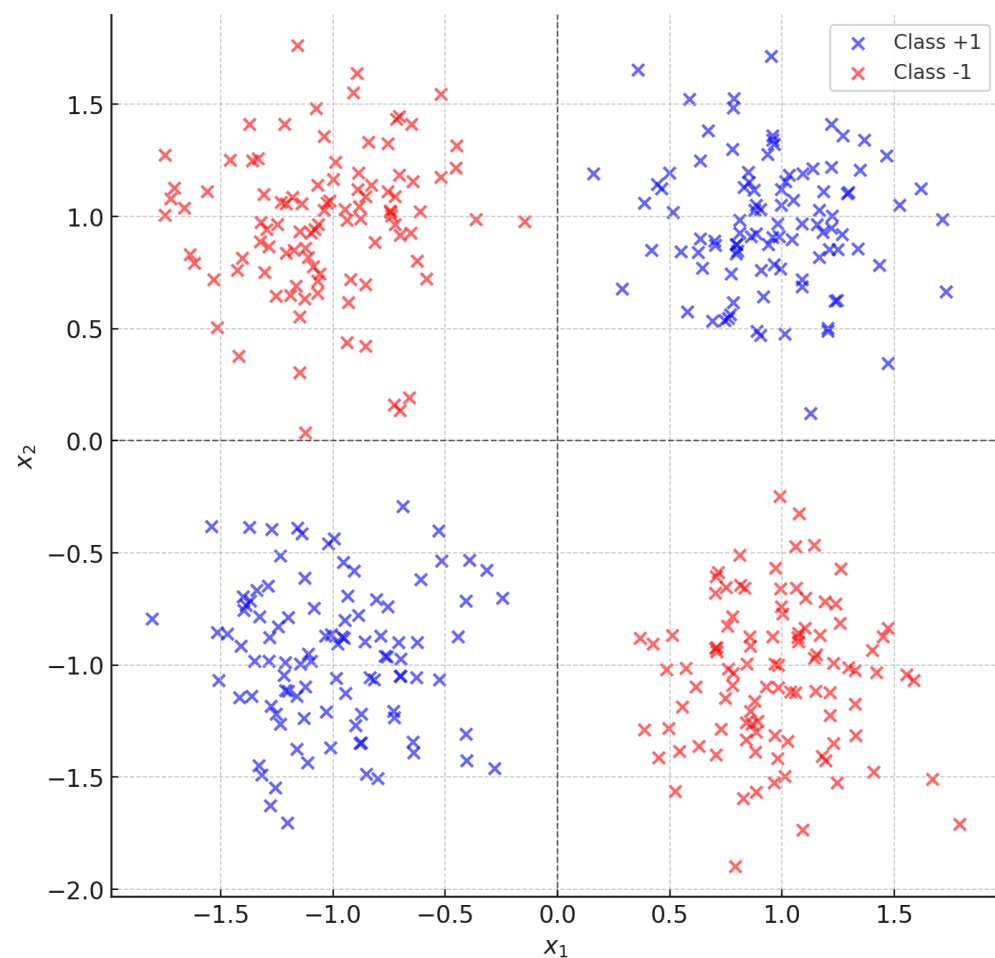
or...  $x_1 x_2 > 0$

$$y = -1 \quad x_1 > 0 \text{ and } x_2 < 0 \text{ or } x_1 < 0 \text{ and } x_2 > 0$$

or...  $x_1 x_2 < 0$

# Examples: XOR Gaussian mixture

$$x \in \mathbb{R}^2 \quad (d = 2) \quad p(x) = \frac{1}{4} \sum_{k=1}^4 \mathcal{N}(\mu_k, I_2)$$



Note that:

$$y = +1 \quad x_1, x_2 > 0 \text{ or } x_1, x_2 < 0 \\ \text{or...} \quad x_1 x_2 > 0$$

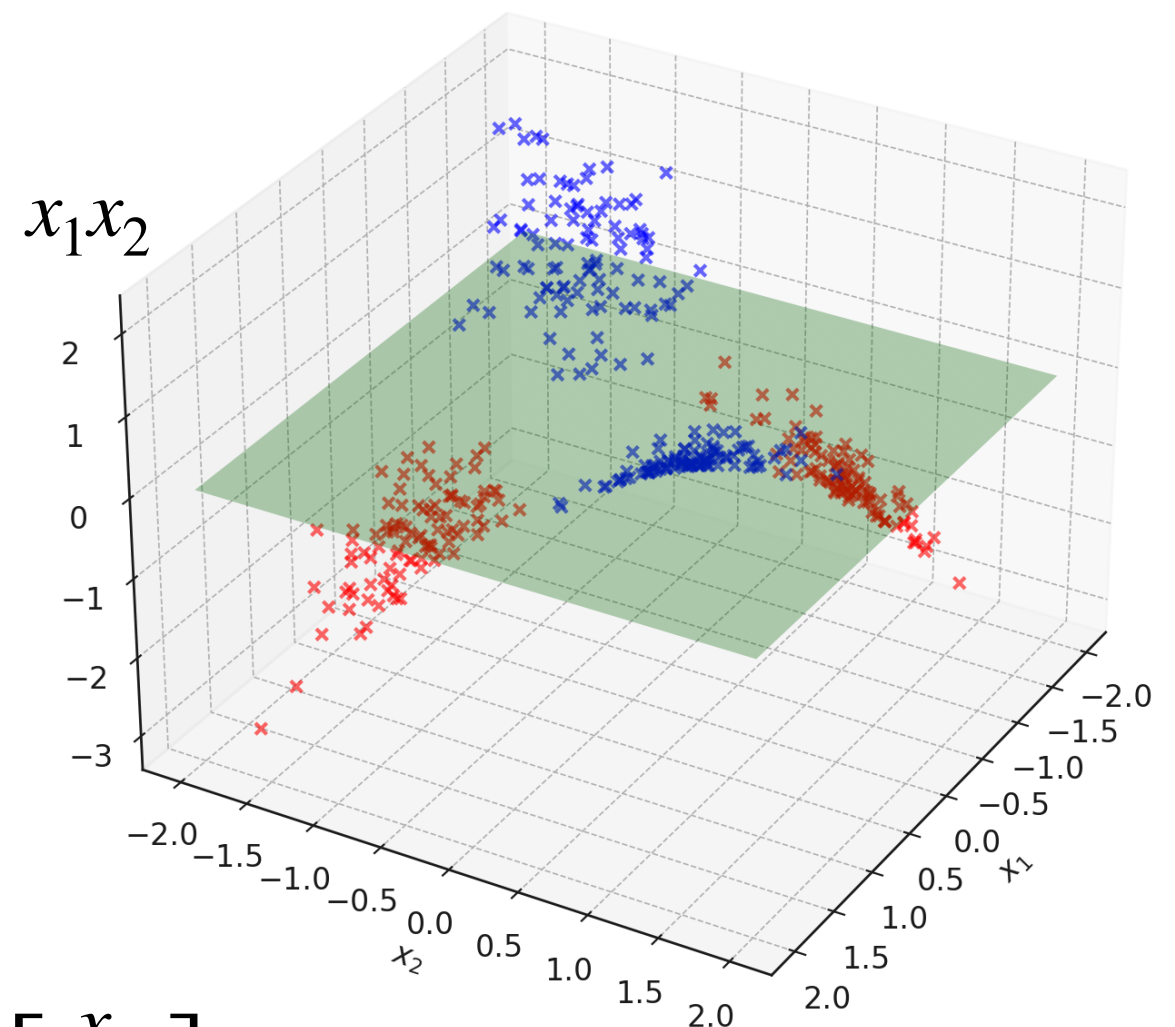
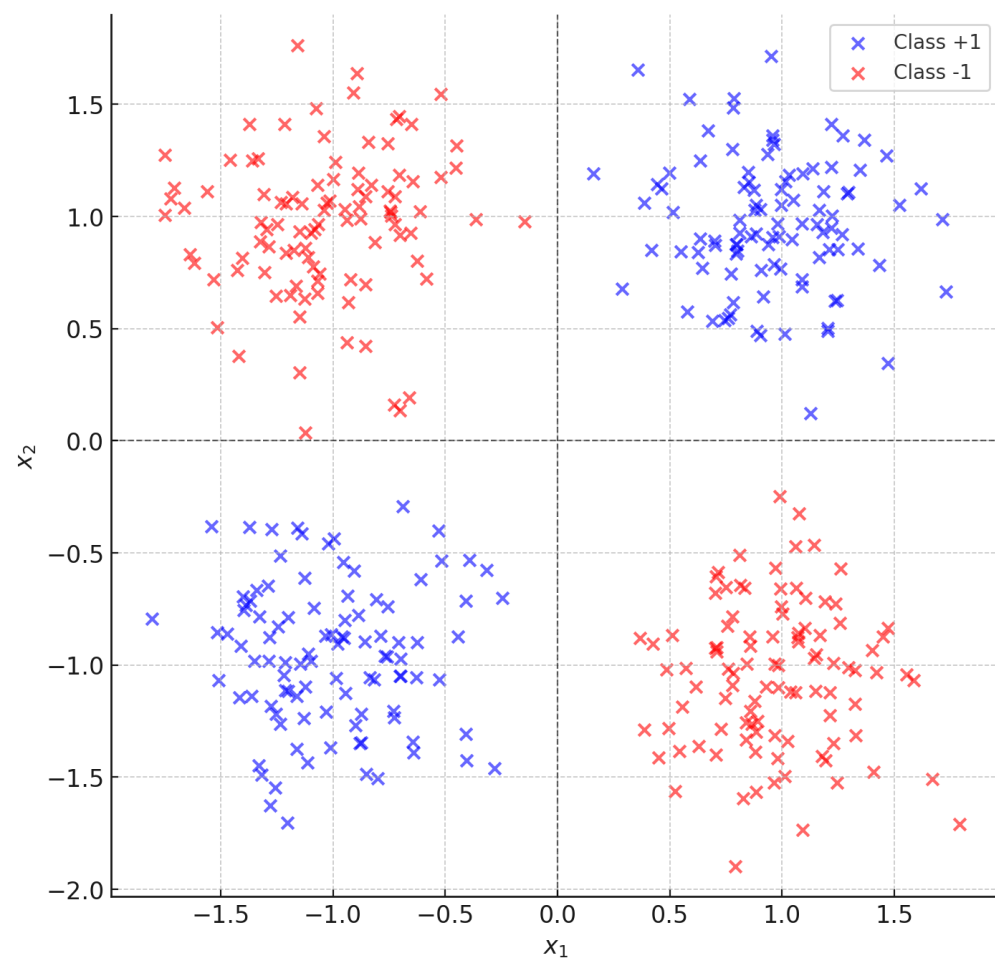
$$y = -1 \quad x_1 > 0 \text{ and } x_2 < 0 \text{ or} \\ x_1 < 0 \text{ and } x_2 > 0 \\ \text{or...} \quad x_1 x_2 < 0$$

This motivates a choice:

$$\varphi(x) = \begin{bmatrix} x_1 \\ x_2 \\ x_1 x_2 \end{bmatrix} \quad (p = 3)$$

# Examples: XOR Gaussian mixture

$$x \in \mathbb{R}^2 \quad (d = 2) \quad p(x) = \frac{1}{4} \sum_{k=1}^4 \mathcal{N}(\mu_k, I_2)$$



$$\varphi(x) = \begin{bmatrix} x_1 \\ x_2 \\ x_1x_2 \end{bmatrix}$$