



Statistical Learning II

Lecture 11 - PCA & Feature maps

Bruno Loureiro
@ CSD, DI-ENS & CNRS

brloureiro@gmail.com

Principal component analysis (PCA)

Variance reduction

- As we saw in Lecture 6, the OLS estimator suffers from high-variance in directions with small singular values.

$$\hat{\boldsymbol{\theta}}_{OLS}(X, \mathbf{y}) = \boldsymbol{\theta}_{\star} + \sum_{j=1}^d \frac{1}{\sigma_j} \langle \mathbf{u}_j, \boldsymbol{\varepsilon} \rangle \mathbf{v}_j$$

Variance reduction

- As we saw in Lecture 6, the OLS estimator suffers from high-variance in directions with small singular values.

$$\hat{\boldsymbol{\theta}}_{OLS}(X, \mathbf{y}) = \boldsymbol{\theta}_{\star} + \sum_{j=1}^d \frac{1}{\sigma_j} \langle \mathbf{u}_j, \boldsymbol{\varepsilon} \rangle \mathbf{v}_j$$

- In Lectures 7 to 10, we studied **regularisation** as a form to mitigate this problem. For instance, ridge regression:

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

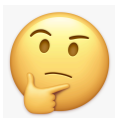
Variance reduction

- As we saw in Lecture 6, the OLS estimator suffers from high-variance in directions with small singular values.

$$\hat{\boldsymbol{\theta}}_{OLS}(X, \mathbf{y}) = \boldsymbol{\theta}_{\star} + \sum_{j=1}^d \frac{1}{\sigma_j} \langle \mathbf{u}_j, \boldsymbol{\varepsilon} \rangle \mathbf{v}_j$$

- In Lectures 7 to 10, we studied **regularisation** as a form to mitigate this problem. For instance, ridge regression:

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

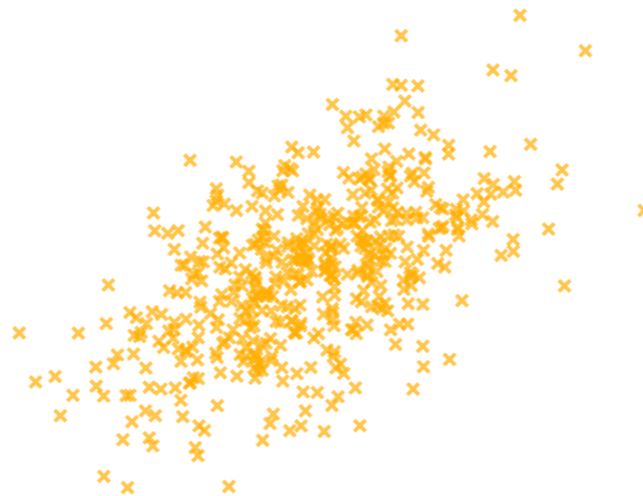


What about getting rid of these directions directly?

Principal component analysis

Let $\mathbf{x}_i \in \mathbb{R}^d$ denote $i = 1, \dots, n$ i.i.d. covariates. Define $\mathbf{X} \in \mathbb{R}^{n \times d}$.

Without loss of generality, assume data is centred.



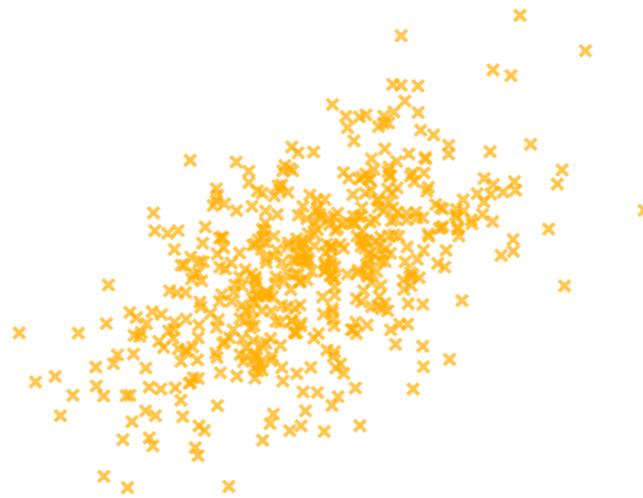
$n = 500$

$d = 2$

Principal component analysis

Let $\mathbf{x}_i \in \mathbb{R}^d$ denote $i = 1, \dots, n$ i.i.d. covariates. Define $\mathbf{X} \in \mathbb{R}^{n \times d}$.

Without loss of generality, assume data is centred.



$n = 500$

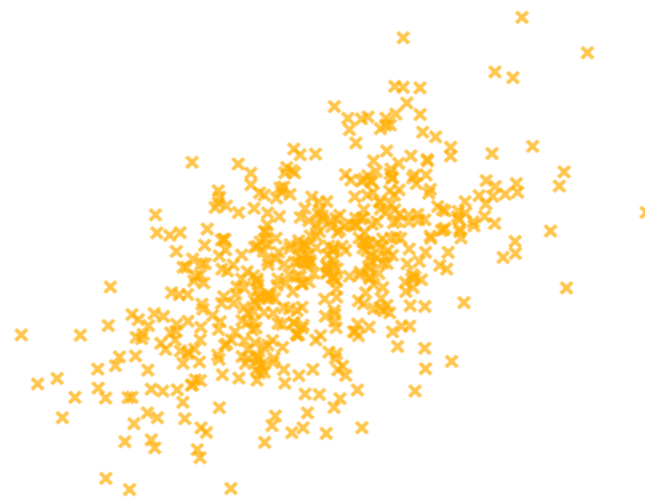
$d = 2$

Goal: find a **lower dimensional** approximation of \mathbf{X} .

Principal component analysis

Let $\mathbf{x}_i \in \mathbb{R}^d$ denote $i = 1, \dots, n$ i.i.d. covariates. Define $\mathbf{X} \in \mathbb{R}^{n \times d}$.

Without loss of generality, assume data is centred.



$n = 500$

$d = 2$

Goal: find a **lower dimensional** approximation of \mathbf{X} .

Simplest case: find best k -dimensional **linear approximation** of \mathbf{X} .

Principal component analysis

Mathematically:

- Let $z_i \in \mathbb{R}^d$, $i = 1, \dots, n$ such that $\text{span}(z_1, \dots, z_n) = \mathbb{R}^k$ with $k \leq d$

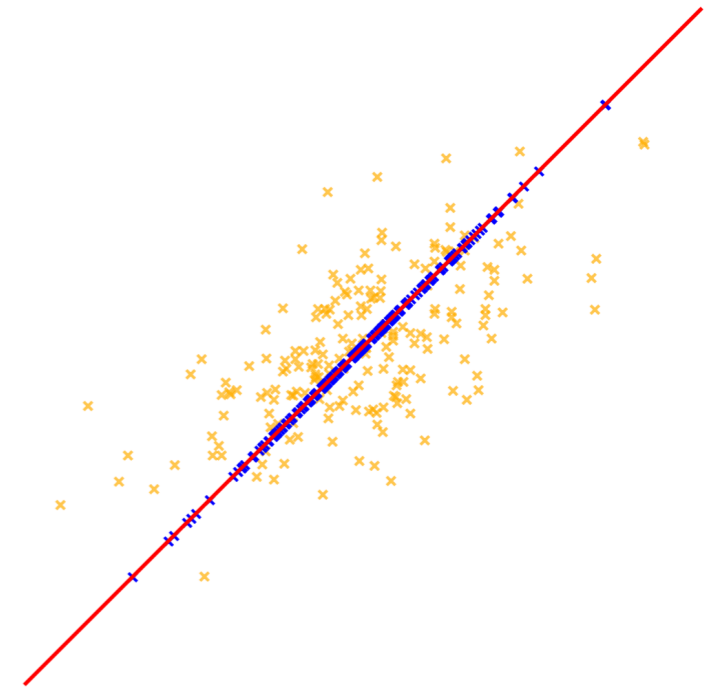
Principal component analysis

Mathematically:

- Let $\mathbf{z}_i \in \mathbb{R}^d$, $i = 1, \dots, n$ such that $\text{span}(\mathbf{z}_1, \dots, \mathbf{z}_n) = \mathbb{R}^k$ with $k \leq d$

The PCA problem consists of:

$$\min_{\mathbf{z}_1, \dots, \mathbf{z}_n} \sum_{i=1}^n \|\mathbf{z}_i - \mathbf{x}_i\|_2^2$$



Principal component analysis

Mathematically:

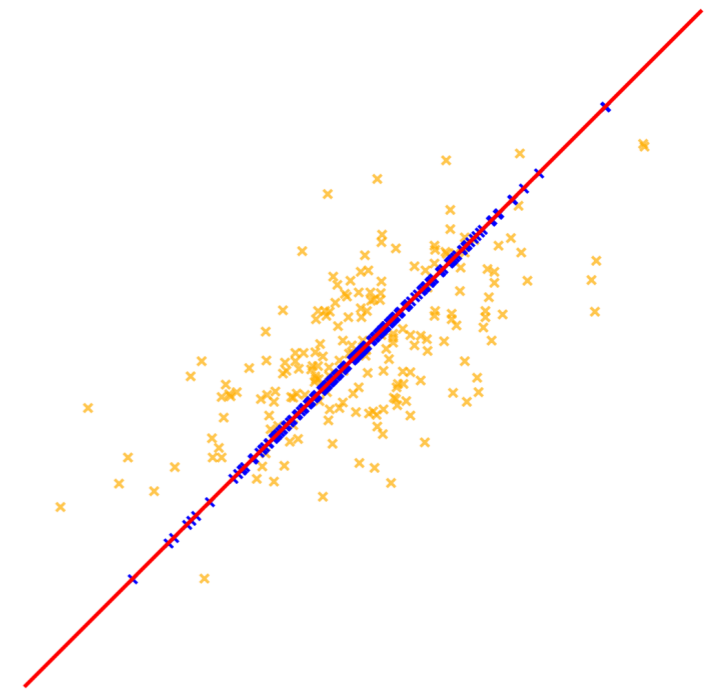
- Let $\mathbf{z}_i \in \mathbb{R}^d$, $i = 1, \dots, n$ such that $\text{span}(\mathbf{z}_1, \dots, \mathbf{z}_n) = \mathbb{R}^k$ with $k \leq d$

The PCA problem consists of:

$$\min_{\mathbf{z}_1, \dots, \mathbf{z}_n} \sum_{i=1}^n \|\mathbf{z}_i - \mathbf{x}_i\|_2^2$$

This can be equivalently written as:

$$\min_{\substack{\mathbf{Z} \in \mathbb{R}^{n \times d}, \\ \text{rank}(\mathbf{Z}) \leq k}} \|\mathbf{Z} - \mathbf{X}\|_F^2$$



Principal component analysis

As we saw in Lecture 1, the solution to this problem is the SVD!

$$\begin{array}{l} \operatorname{argmin} \\ \mathbf{Z} \in \mathbb{R}^{n \times d}, \\ \operatorname{rank}(\mathbf{Z}) \leq k \end{array} \|\mathbf{Z} - \mathbf{X}\|_F^2 = \sum_{j=1}^k \sigma_j \mathbf{u}_j \mathbf{v}_j$$

Principal component analysis

As we saw in Lecture 1, the solution to this problem is the SVD!

$$\underset{\substack{\mathbf{Z} \in \mathbb{R}^{n \times d}, \\ \text{rank}(\mathbf{Z}) \leq k}}{\text{argmin}} \|\mathbf{Z} - \mathbf{X}\|_F^2 = \sum_{j=1}^k \sigma_j \mathbf{u}_j \mathbf{v}_j$$

In other words, the **best k -dimensional linear approximation** to the data consists of retaining only the **top k singular values**.

Principal component analysis

As we saw in Lecture 1, the solution to this problem is the SVD!

$$\underset{\substack{\mathbf{Z} \in \mathbb{R}^{n \times d}, \\ \text{rank}(\mathbf{Z}) \leq k}}{\text{argmin}} \|\mathbf{Z} - \mathbf{X}\|_F^2 = \sum_{j=1}^k \sigma_j \mathbf{u}_j \mathbf{v}_j$$

In other words, the **best k -dimensional linear approximation** to the data consists of retaining only the **top k singular values**.



This result holds for more general norms, and is known as the **Eckart-Young-Minsky Theorem**

Principal component analysis

As we saw in Lecture 1, the solution to this problem is the SVD!

$$\begin{array}{l} \operatorname{argmin} \\ \mathbf{Z} \in \mathbb{R}^{n \times d}, \\ \operatorname{rank}(\mathbf{Z}) \leq k \end{array} \|\mathbf{Z} - \mathbf{X}\|_F^2 = \sum_{j=1}^k \sigma_j \mathbf{u}_j \mathbf{v}_j$$

In other words, the **best k -dimensional linear approximation** to the data consists of retaining only the **top k singular values**.



This result holds for more general norms, and is known as the **Eckart-Young-Minsky Theorem**

Remark: This is equivalent to keeping the k directions with largest variance in the data.

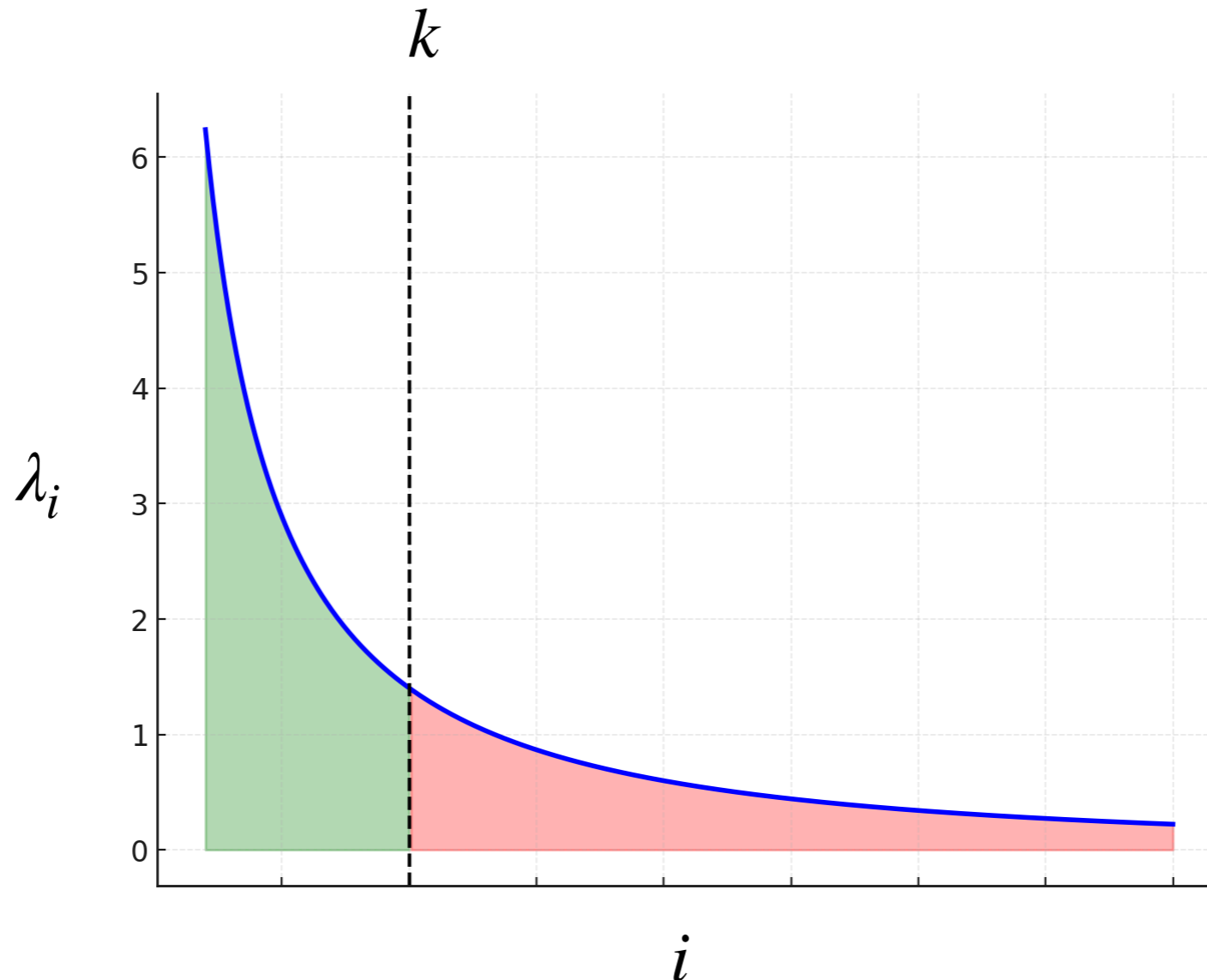
$$\operatorname{Tr}(\hat{\Sigma}_n) = \frac{1}{n} \sum_{i=1}^{\operatorname{rank}(\mathbf{X})} \lambda_i^2$$

PCA in practice

In practice, how to choose the k ?

Total variance of data given by:

$$\text{Tr}(\hat{\Sigma}_n) = \frac{1}{n} \sum_{i=1}^{\text{rank}(X)} \lambda_i^2$$



Feature maps

Motivation

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

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

Motivation

Up to now, 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:

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

Motivation

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

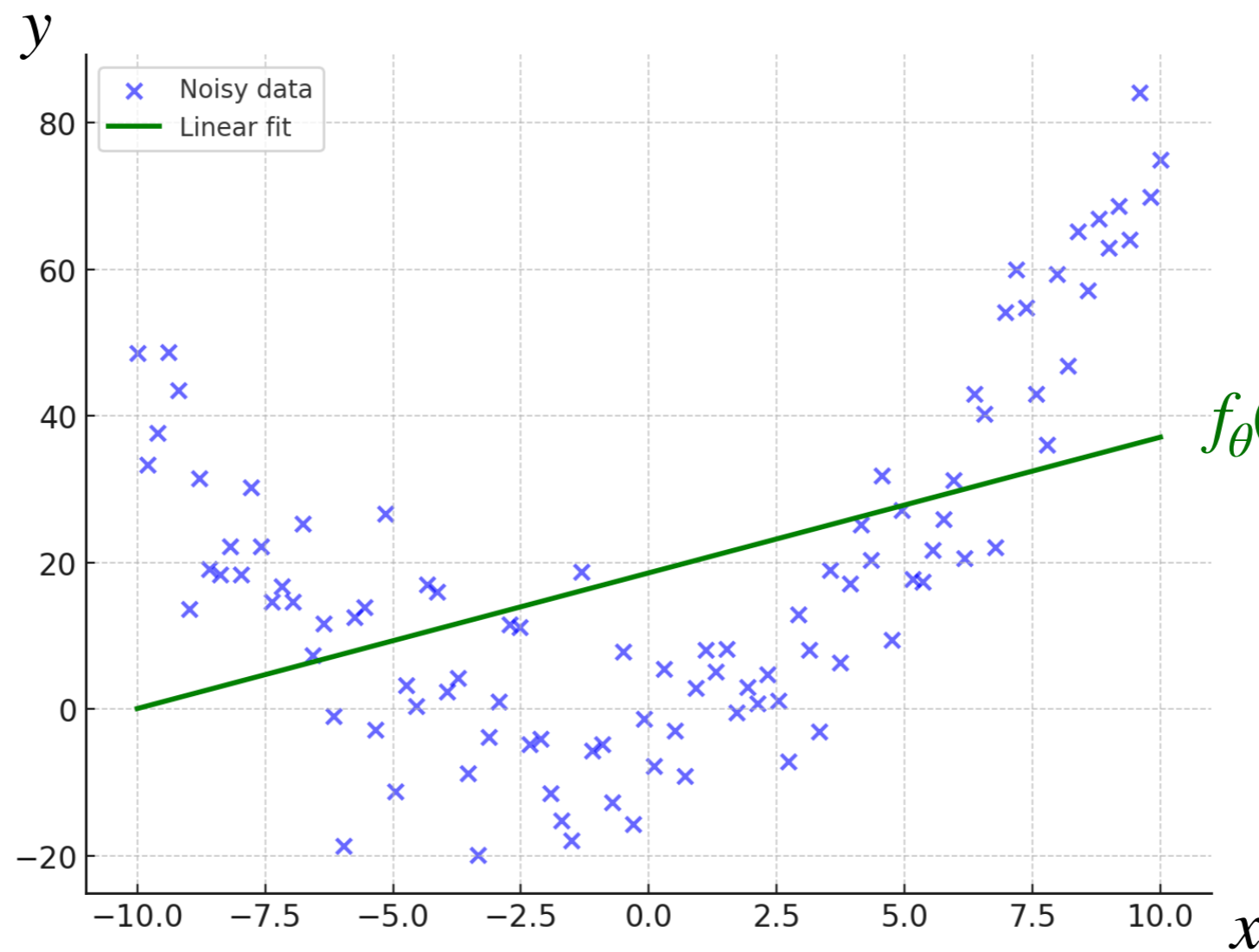
$$f_{\theta}(\mathbf{x}) = \langle \theta, \mathbf{x} \rangle$$

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

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

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

Motivation



$$f_{\theta}(x) = \theta x + b$$

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**:

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

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

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



- 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**:

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

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

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

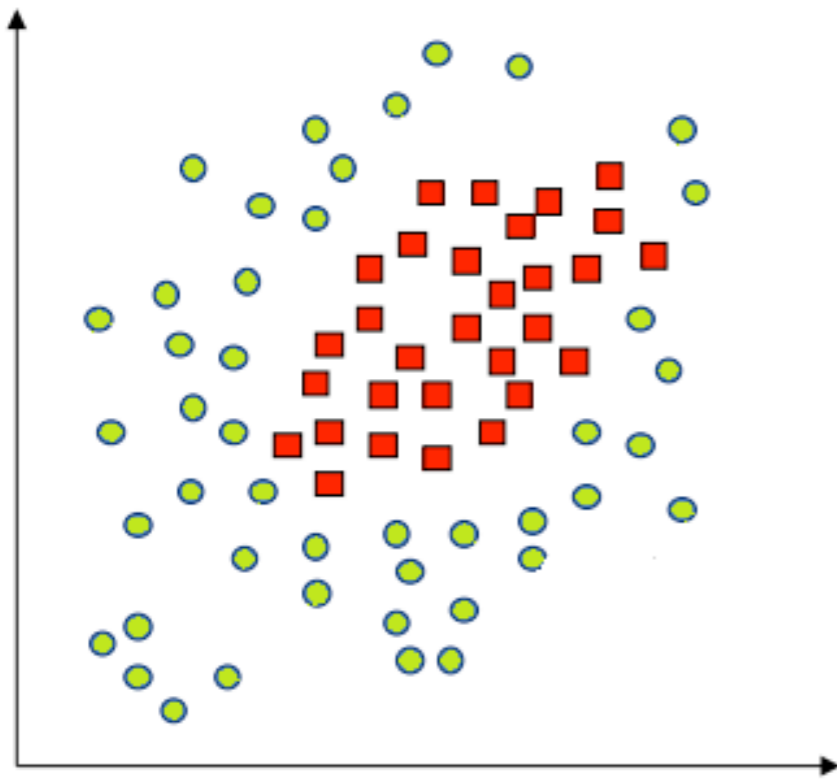


- Now we have $\boldsymbol{\theta} \in \mathbb{R}^p$.
- $f_{\boldsymbol{\theta}}$ still a linear function of $\boldsymbol{\theta}$.
- Typically $p > d$.
- More generally, we can consider $\boldsymbol{\varphi} : \mathcal{X} \rightarrow \mathbb{R}^p$

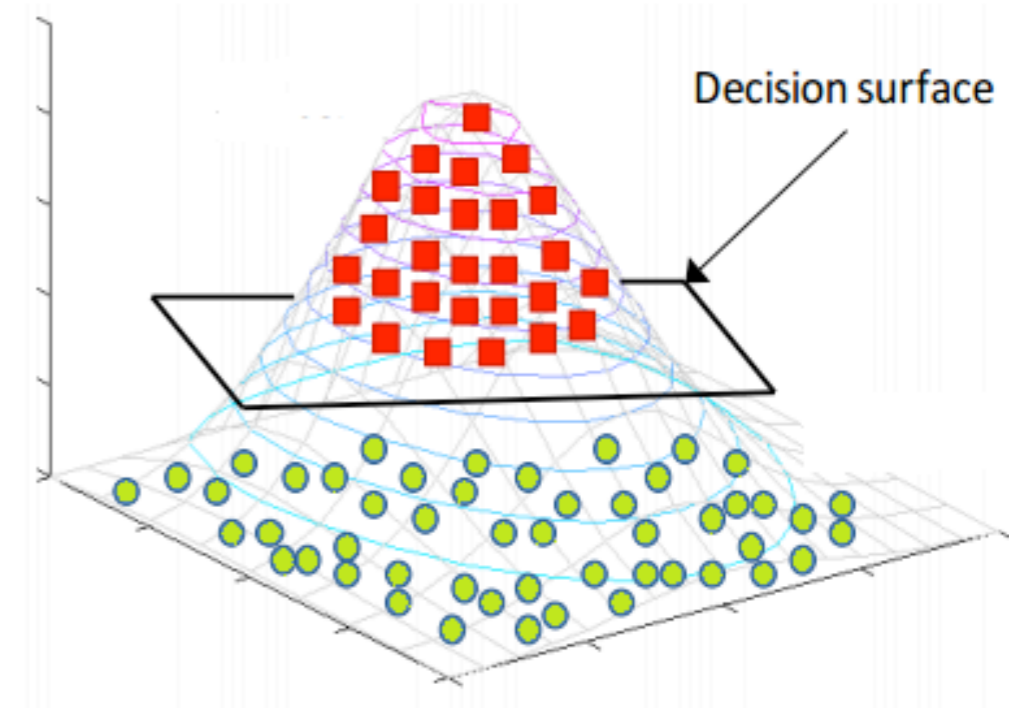

Example: \mathcal{X} a collection of books.

Feature maps

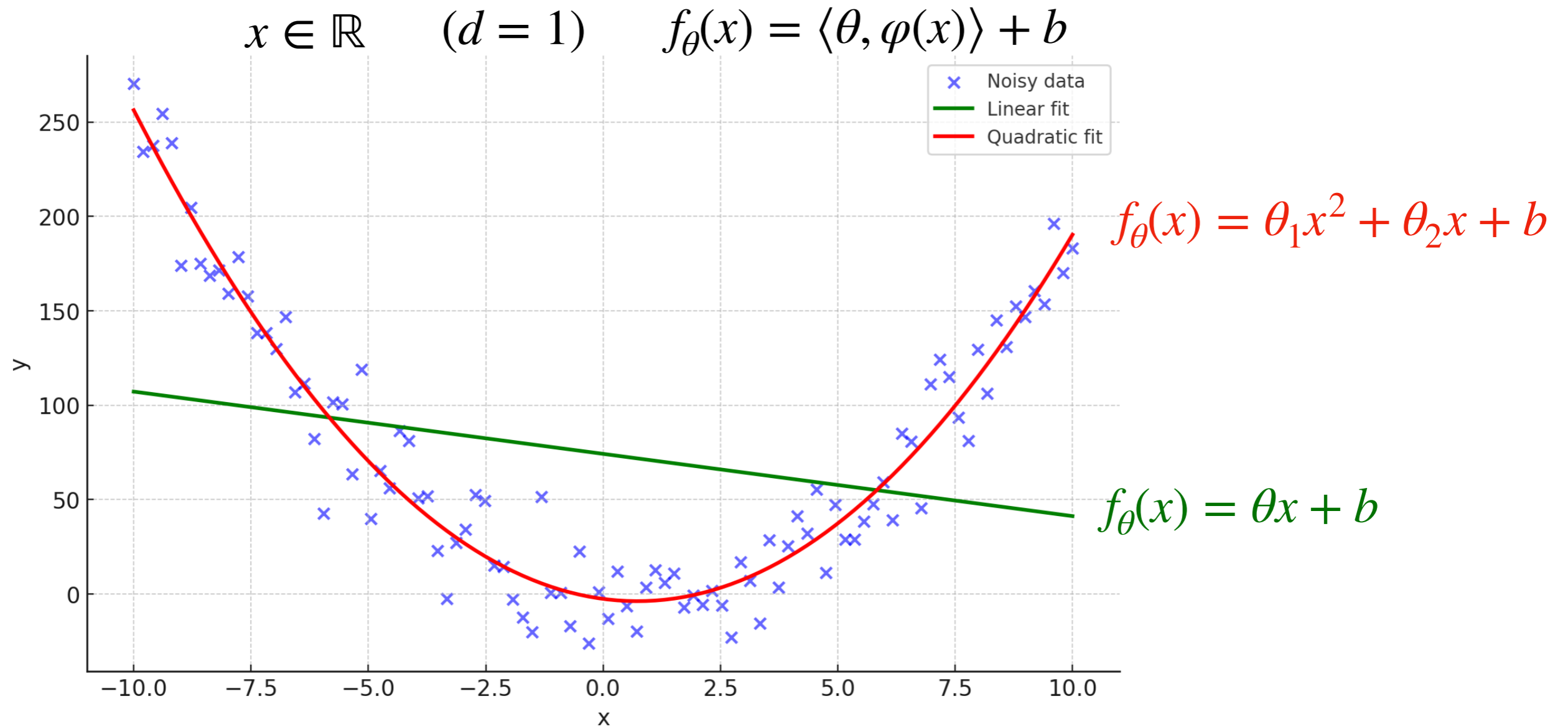
Intuition: Typically easier to linearly separate data in higher-dimensions



φ

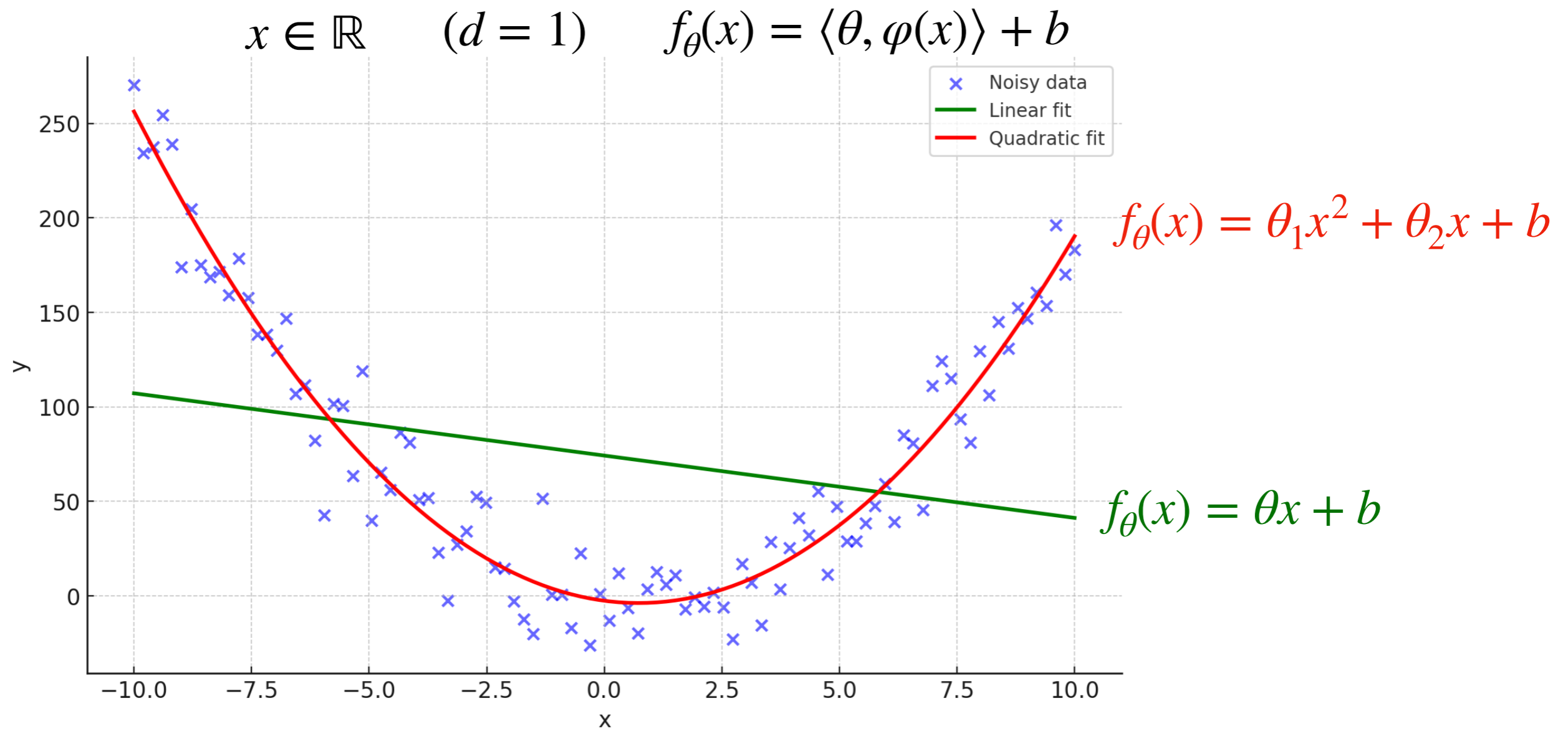


Examples: quadratic function



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 :

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

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 :

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

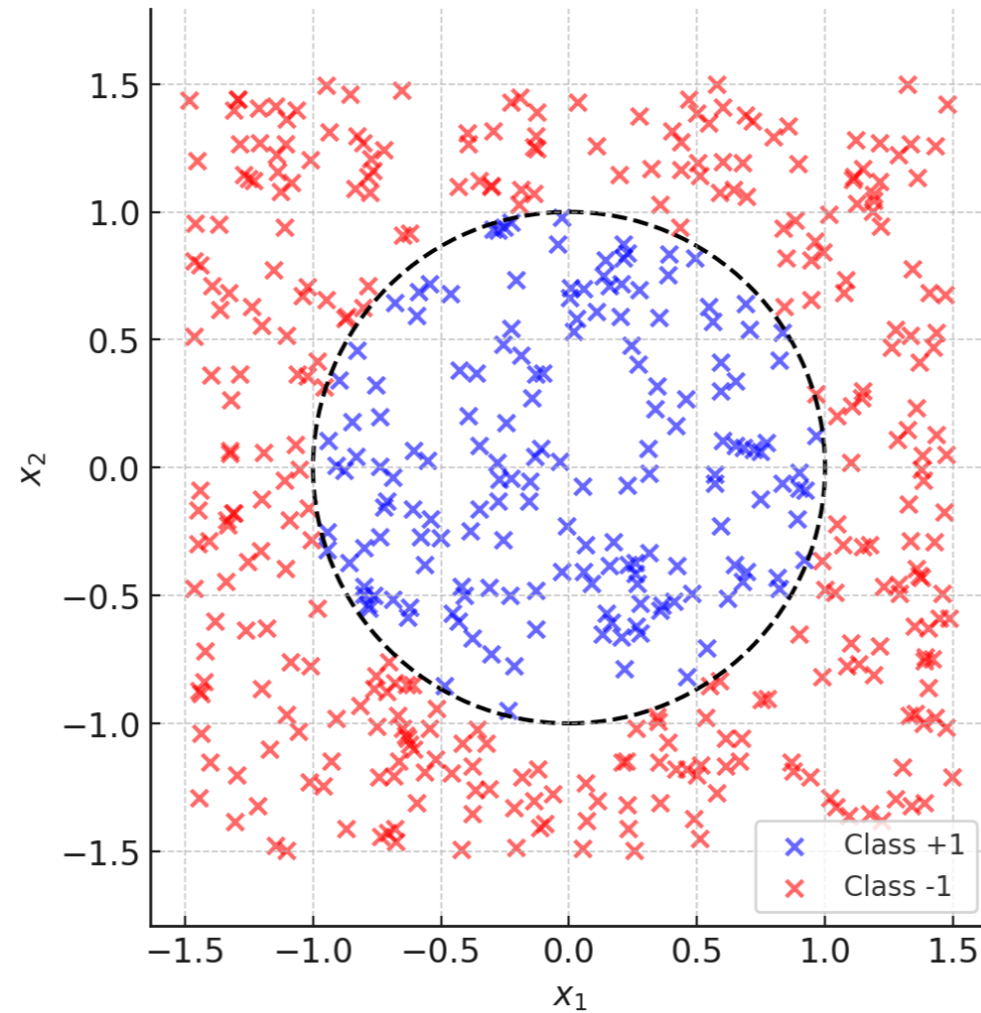
We can generalise this to degree k polynomials in \mathbb{R}^d :

Example $d = 2$:

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

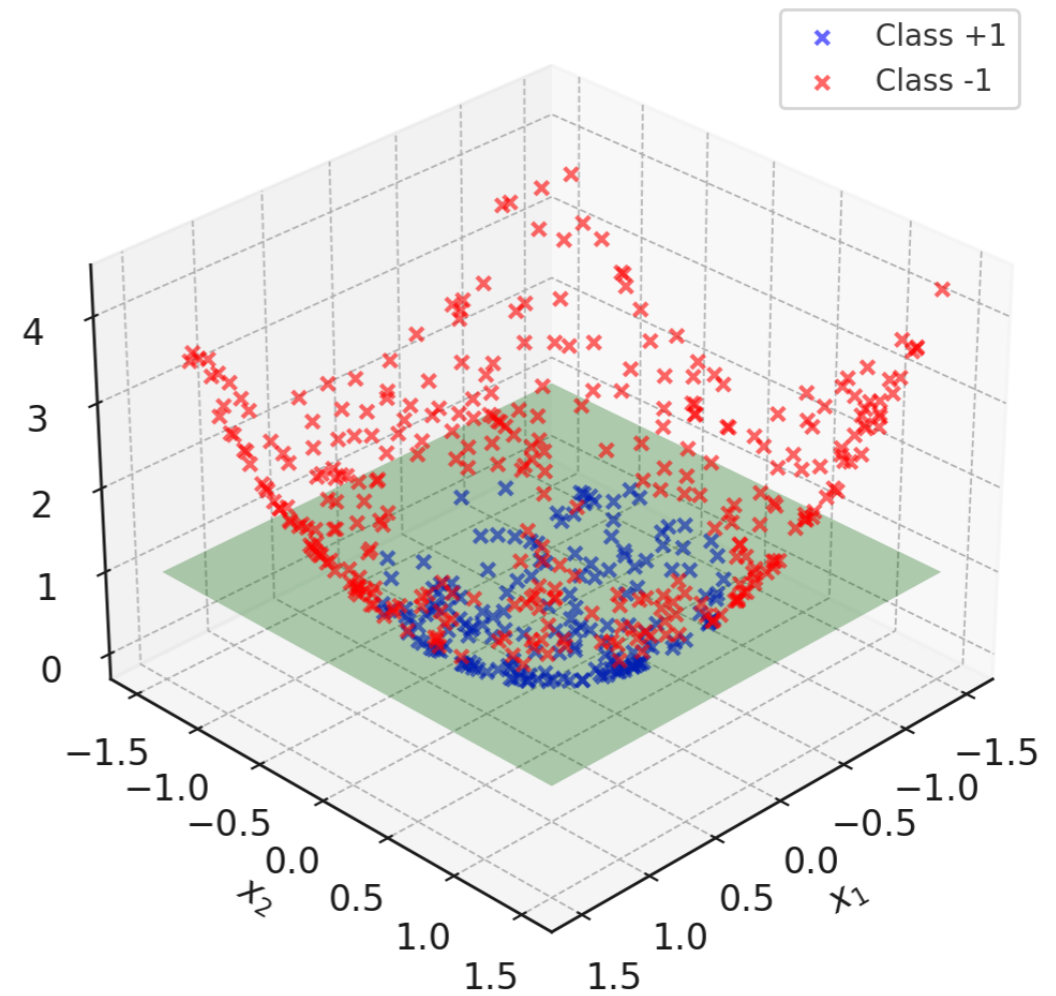
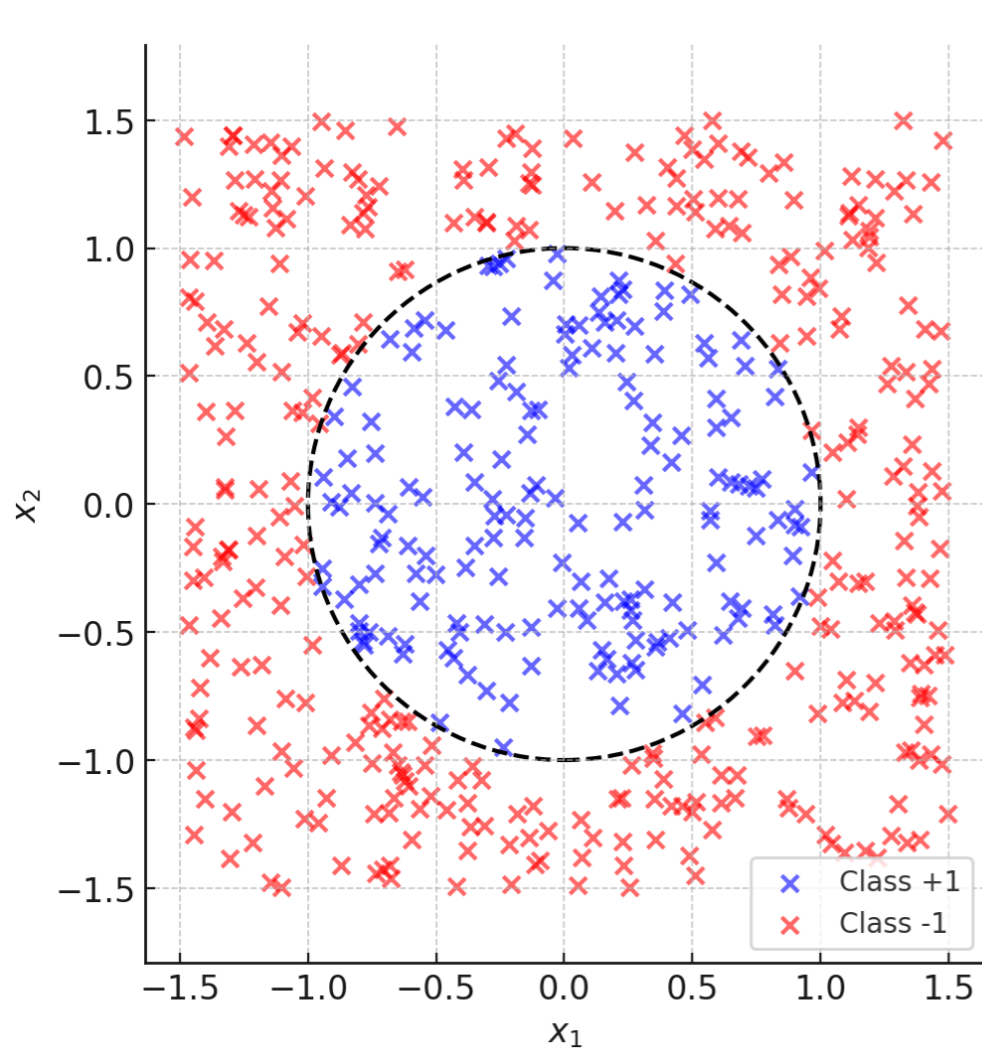
Examples: data in circle

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



Examples: data in circle

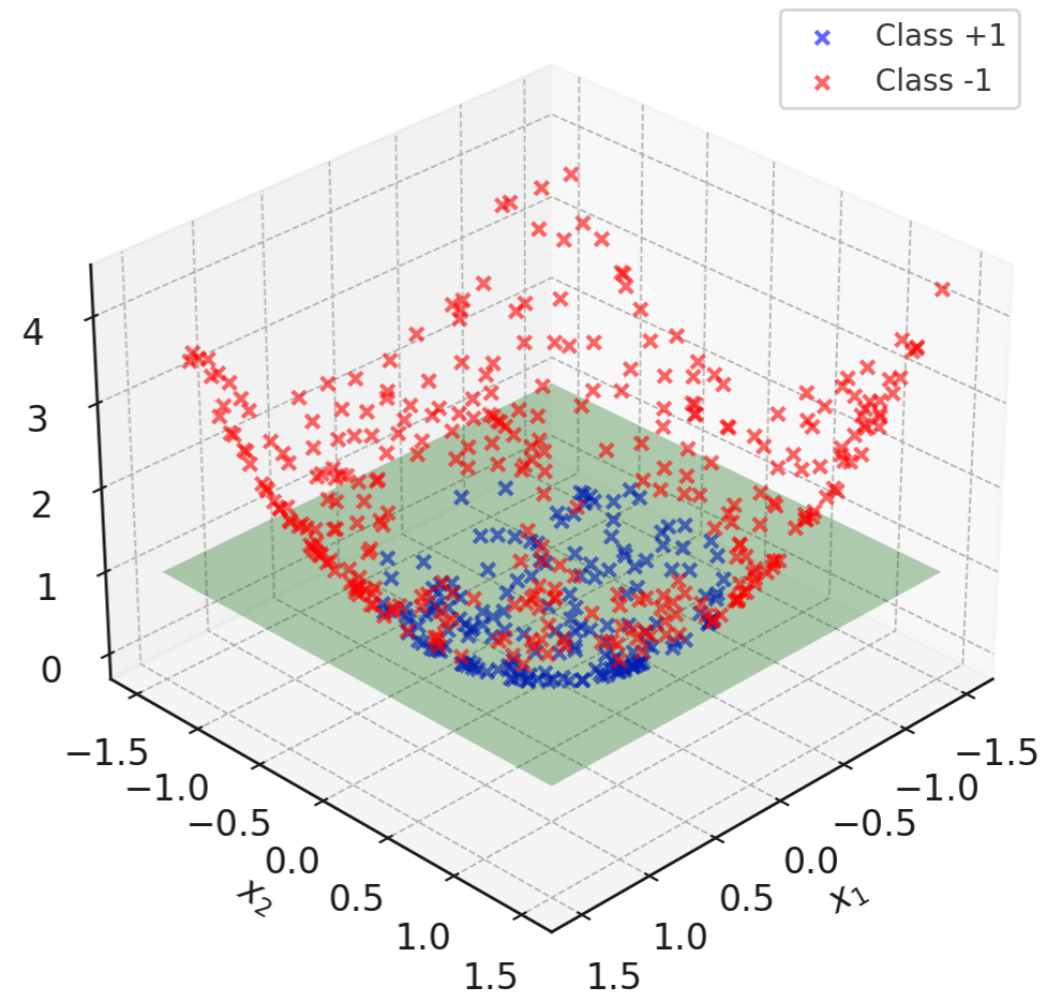
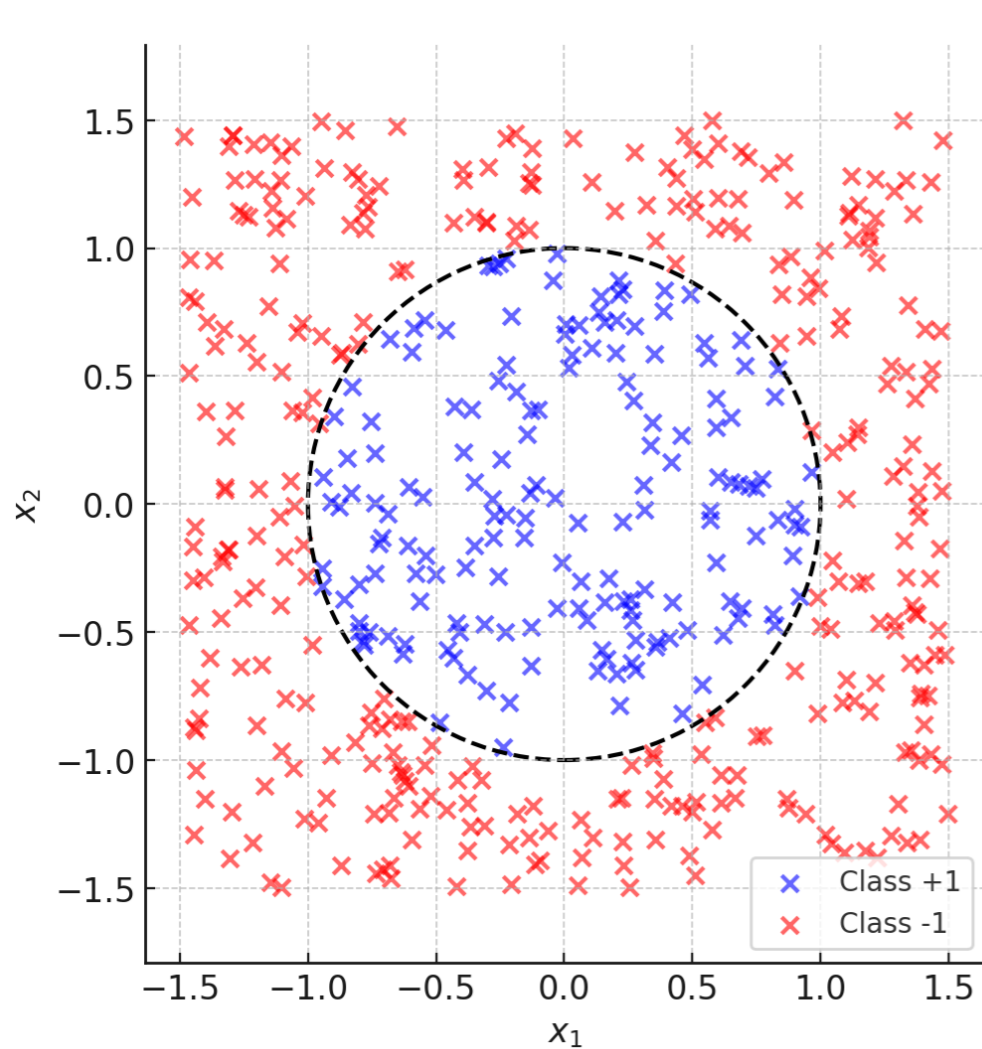
$$x \in \mathbb{R}^2 \quad (d = 2) \quad y = \begin{cases} +1 & \text{if } x_1^2 + x_2^2 \leq 1 \\ -1 & \text{if } x_1^2 + x_2^2 > 1 \end{cases}$$



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

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

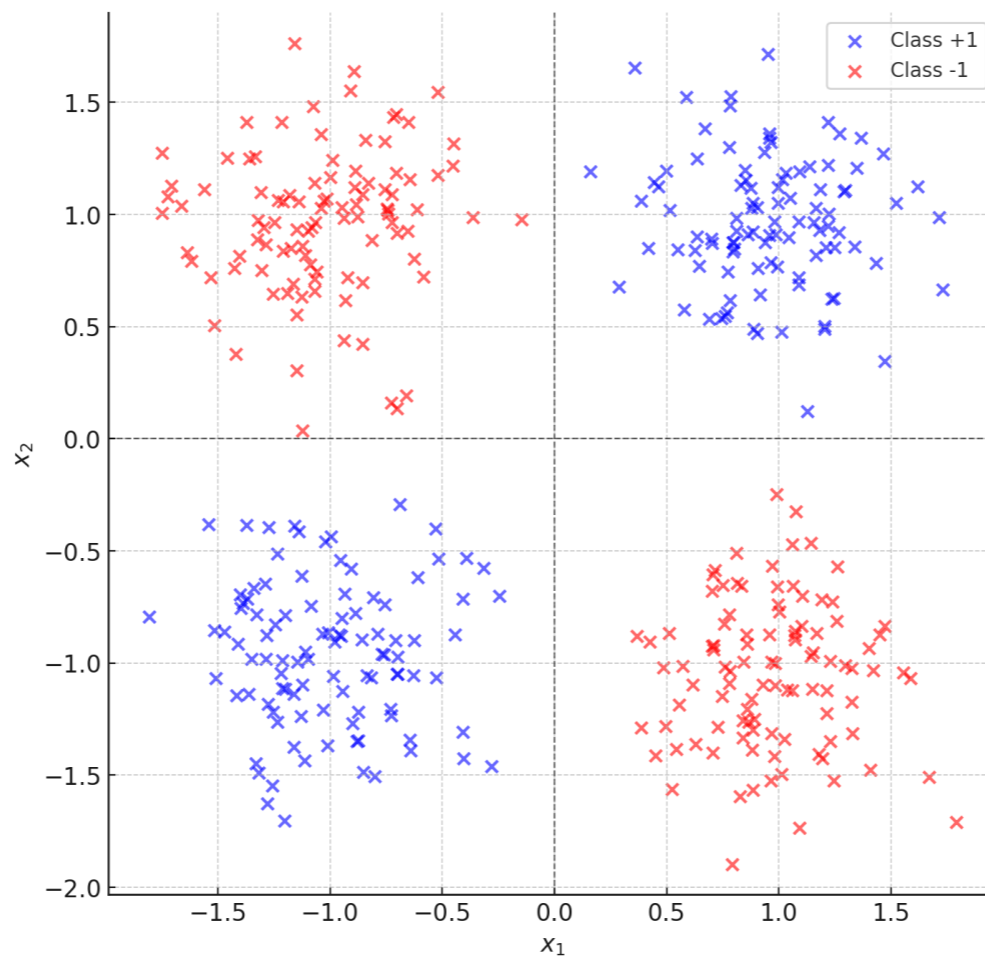
$$x \in \mathbb{R}^2 \quad (d = 2) \quad p(\mathbf{x}) = \frac{1}{4} \sum_{k=1}^4 \mathcal{N}(\boldsymbol{\mu}_k, \mathbf{I}_2)$$

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

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

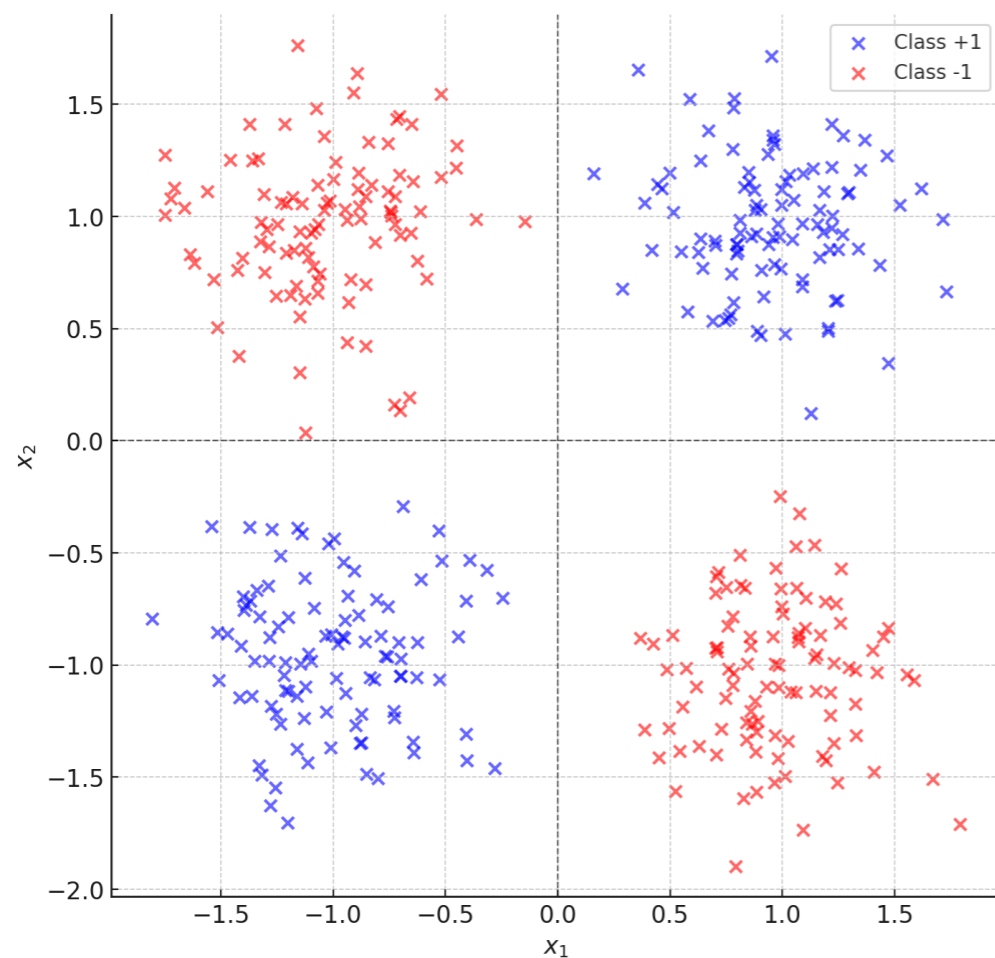
$$\boldsymbol{\mu}_3 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

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



Examples: XOR Gaussian mixture

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



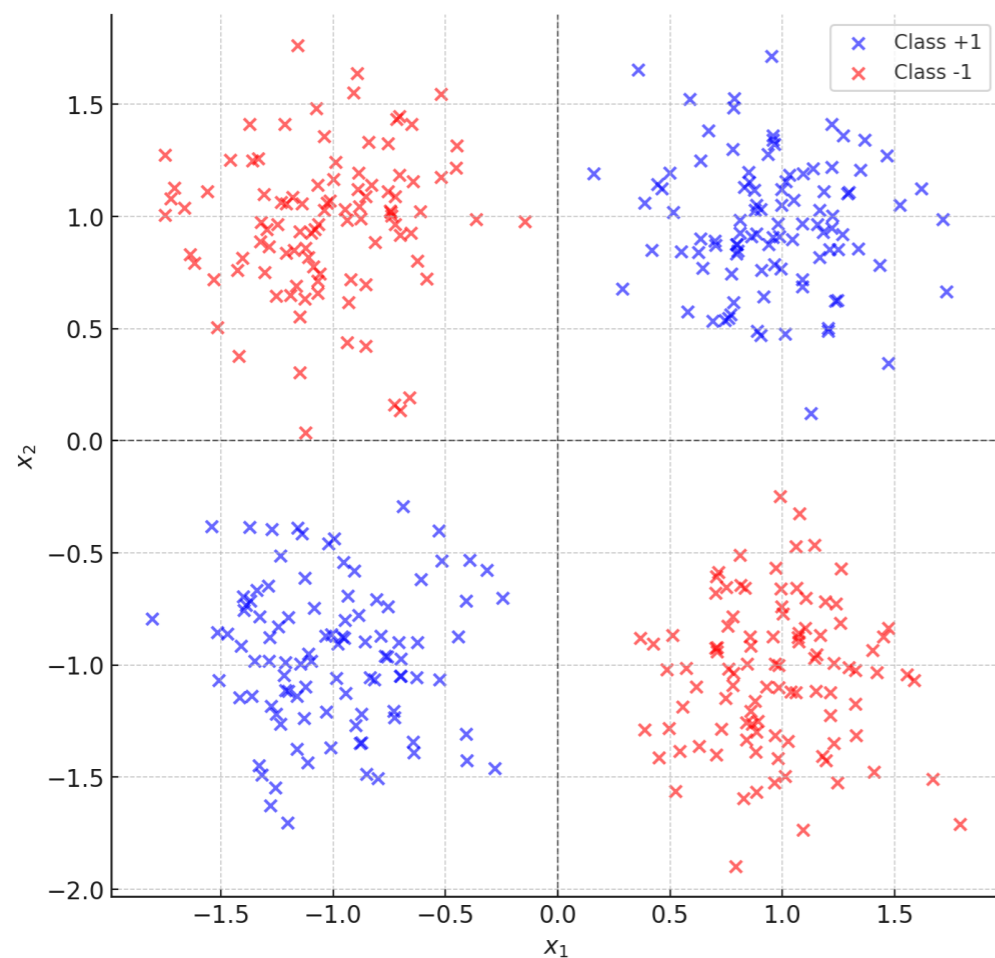
Note that:

$$y = +1 \quad (x_1, x_2) \in \{(-1, -1), (1, 1)\}$$

$$y = -1 \quad (x_1, x_2) \in \{(1, -1), (-1, 1)\}$$

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:

$$y = +1 \quad (x_1, x_2) \in \{(-1, -1), (1, 1)\}$$

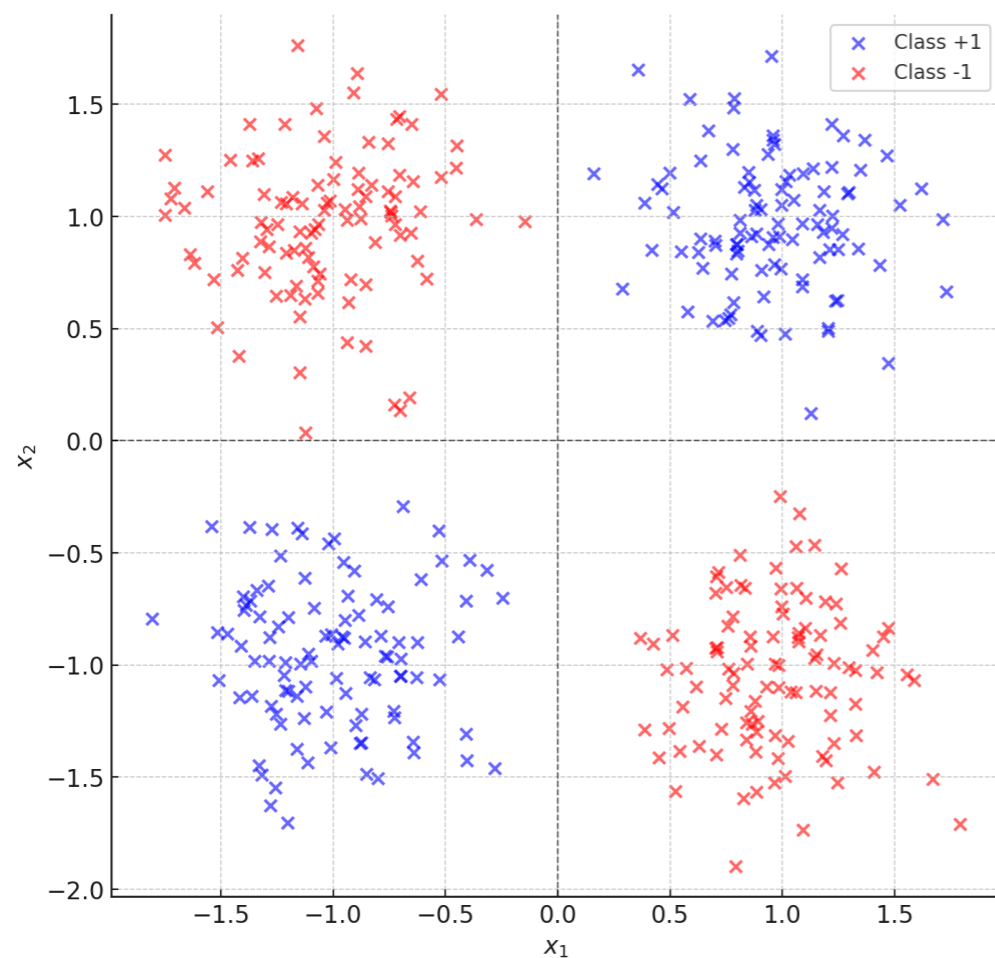
$$\text{or...} \quad x_1 x_2 = +1$$

$$y = -1 \quad (x_1, x_2) \in \{(1, -1), (-1, 1)\}$$

$$\text{or...} \quad x_1 x_2 = -1$$

Examples: XOR Gaussian mixture

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



Note that:

$$y = +1 \quad (x_1, x_2) \in \{(-1, -1), (1, 1)\}$$

$$\text{or...} \quad x_1 x_2 = +1$$

$$y = -1 \quad (x_1, x_2) \in \{(1, -1), (-1, 1)\}$$

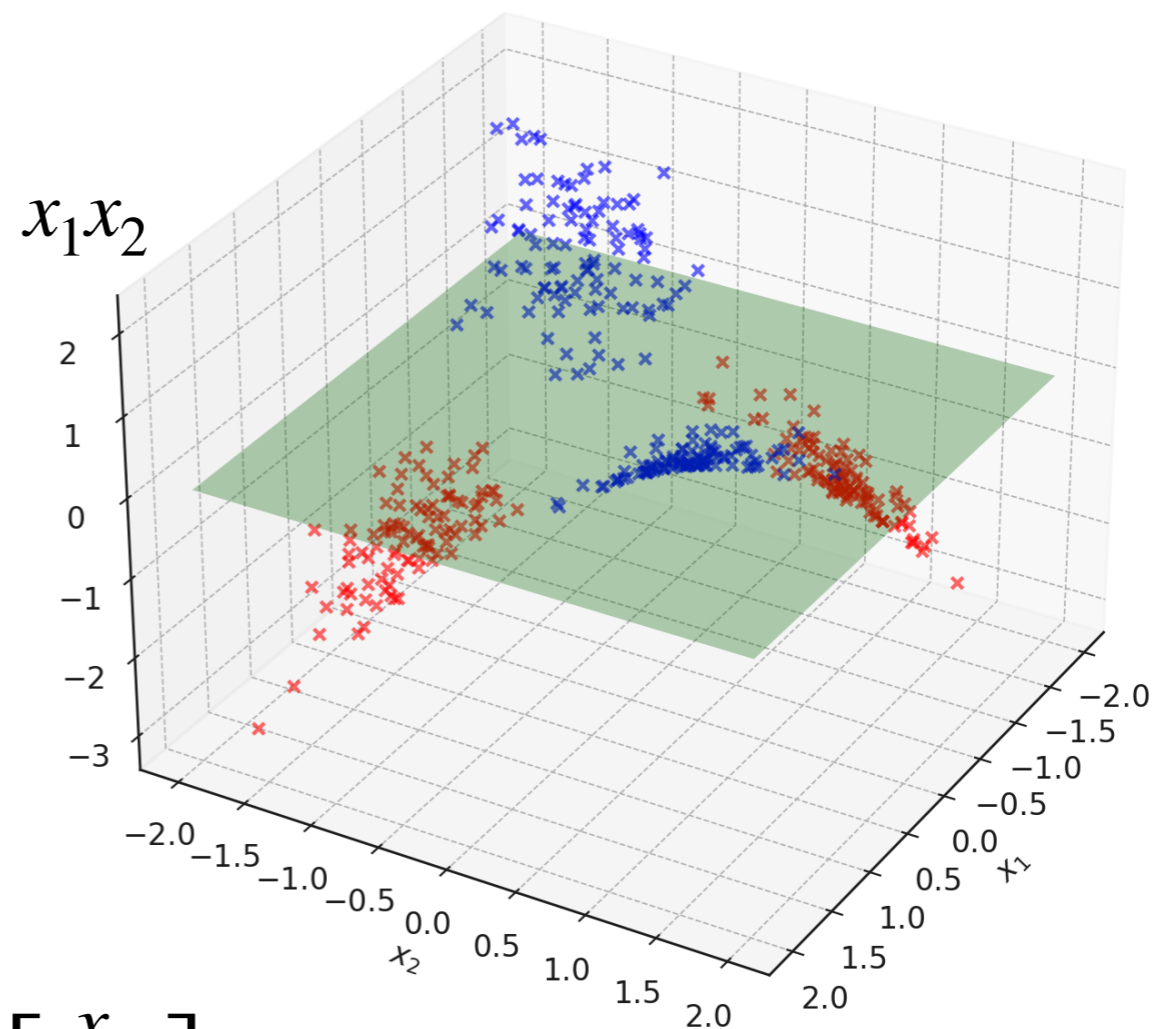
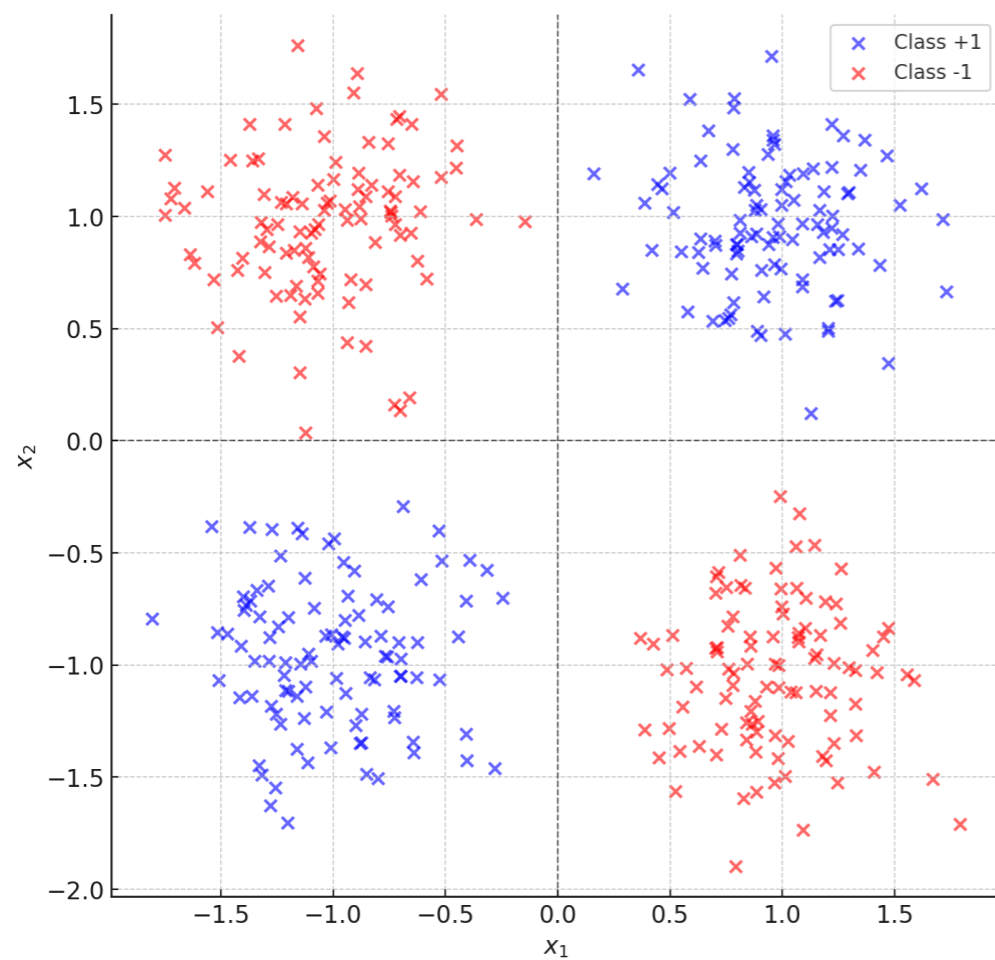
$$\text{or...} \quad x_1 x_2 = -1$$

This motivates a choice:

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

Examples: XOR Gaussian mixture

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



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