

Modèles Linéaires

①

I. Rappels sur les applications linéaires - Paramétrisation

des modèles linéaires.

1) Applications Linéaires.

Définition: $f: M^d \rightarrow M^{d'}$ est linéaire si:

$$\forall x, f(x) = M_f x \quad \text{ou} \quad M_f \in M^{d' \times d}$$

est la matrice associée à f .

Remarque: On a alors, $M_f = \begin{pmatrix} | & | & & | \\ f(e_1) & f(e_2) & \dots & f(e_d) \\ | & | & & | \end{pmatrix} \begin{matrix} \leftarrow e_1' \\ \leftarrow e_2' \\ \vdots \\ \leftarrow e_{d'}' \end{matrix}$

où (e_1, \dots, e_d) est la base canonique de M^d et $(e_1', \dots, e_{d'}')$ est la base canonique de $M^{d'}$.

Définition: Le rang de f est la dimension du sous-espace vectoriel de $M^{d'}$ engendré par les colonnes de M_f , qui est la même que la dimension du sous-espace de M^d engendré par les lignes de M_f .

Théorème (du rang): $\dim(M^d) = \dim(\text{Ker}(f)) + \dim(\text{Im}(f))$.

2) Le cas des formes linéaires

Définition: Lorsque $d' = 1$, on parle de forme linéaire.

Comme toute application linéaire f peut être décomposée comme

$$f = f_1 e_1' + f_2 e_2' + \dots + f_{d'} e_{d'}'$$

on va s'intéresser au cas des formes linéaires.

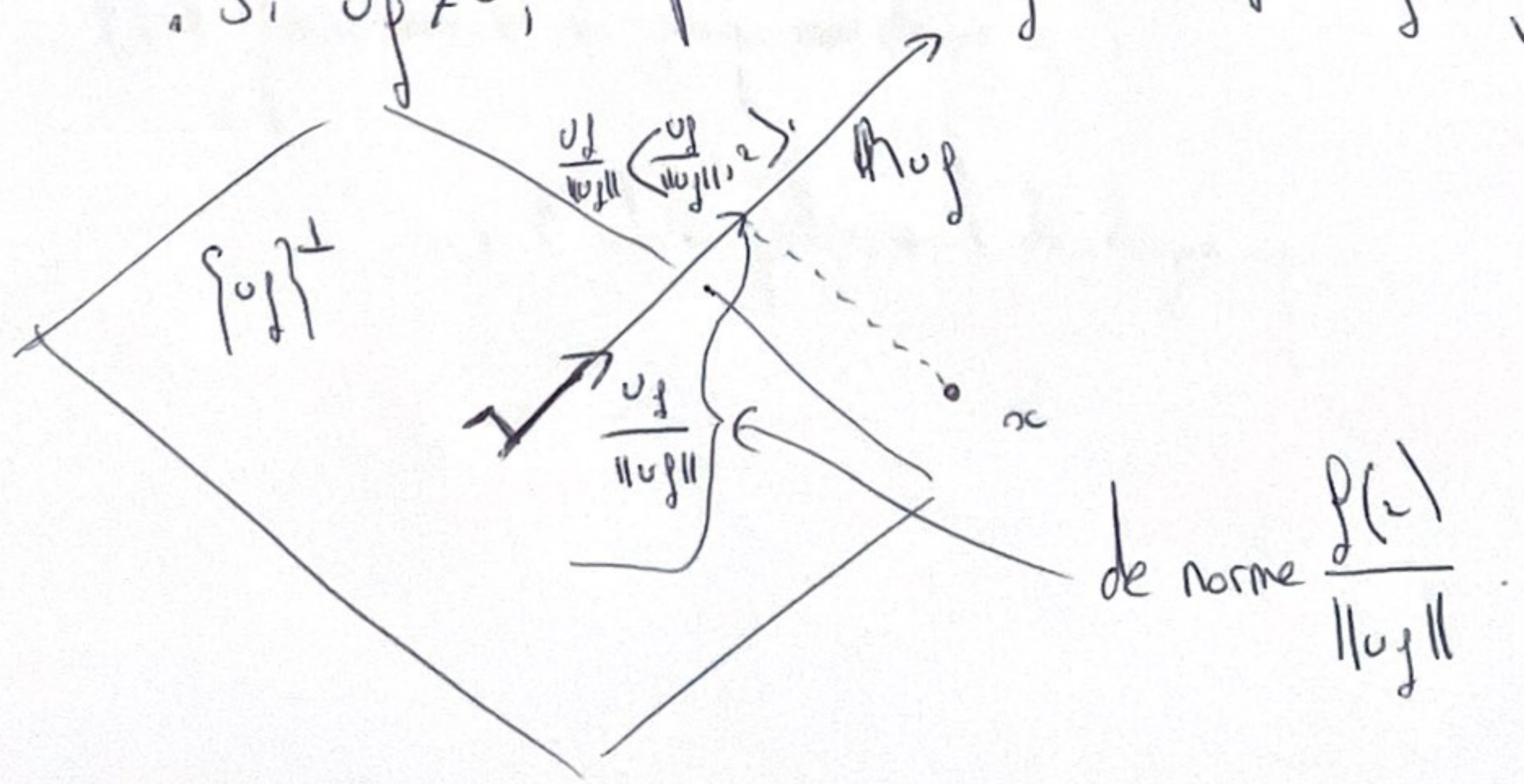
Proposition: Si f est une forme linéaire, il existe $u_f \in \mathbb{M}^d$

tel que $\forall x, f(x) = u_f^T x = \langle u_f, x \rangle$

Interprétation Géométrique:

• Si $u_f = 0$, la fonction est nulle.

• Si $u_f \neq 0$, on peut écrire $f(x) = \|u_f\| \langle \frac{u_f}{\|u_f\|}, x \rangle \quad \forall x$.



3) Reparamétrisation et feature maps

(3)

💡 IP est souvent possible, parce qu'une fonction $f: \mathbb{M}^d \rightarrow \mathbb{M}^{d'}$ n'est pas linéaire d'écrire

$$f = g \circ \varphi \quad \text{où} \quad \varphi: \mathbb{M}^d \rightarrow \mathbb{M}^{d''}$$

$g: \mathbb{M}^{d''} \rightarrow \mathbb{M}^{d'}$ est appelée une feature map, et où φ est linéaire.

Exemples:

- $\varphi(x) = x \rightarrow$ Modèle linéaire classique.

- $\varphi(x) = \begin{pmatrix} x \\ x^2 \\ \vdots \end{pmatrix} \rightarrow$ Modèle affine classique (souvent simplement appelé modèle linéaire).

- $\varphi(x) = (\text{monômes en } x \text{ jusqu'au degré } q)$

\rightarrow Modèle polynomial de degré q .

9) Séparabilité Linéaire.

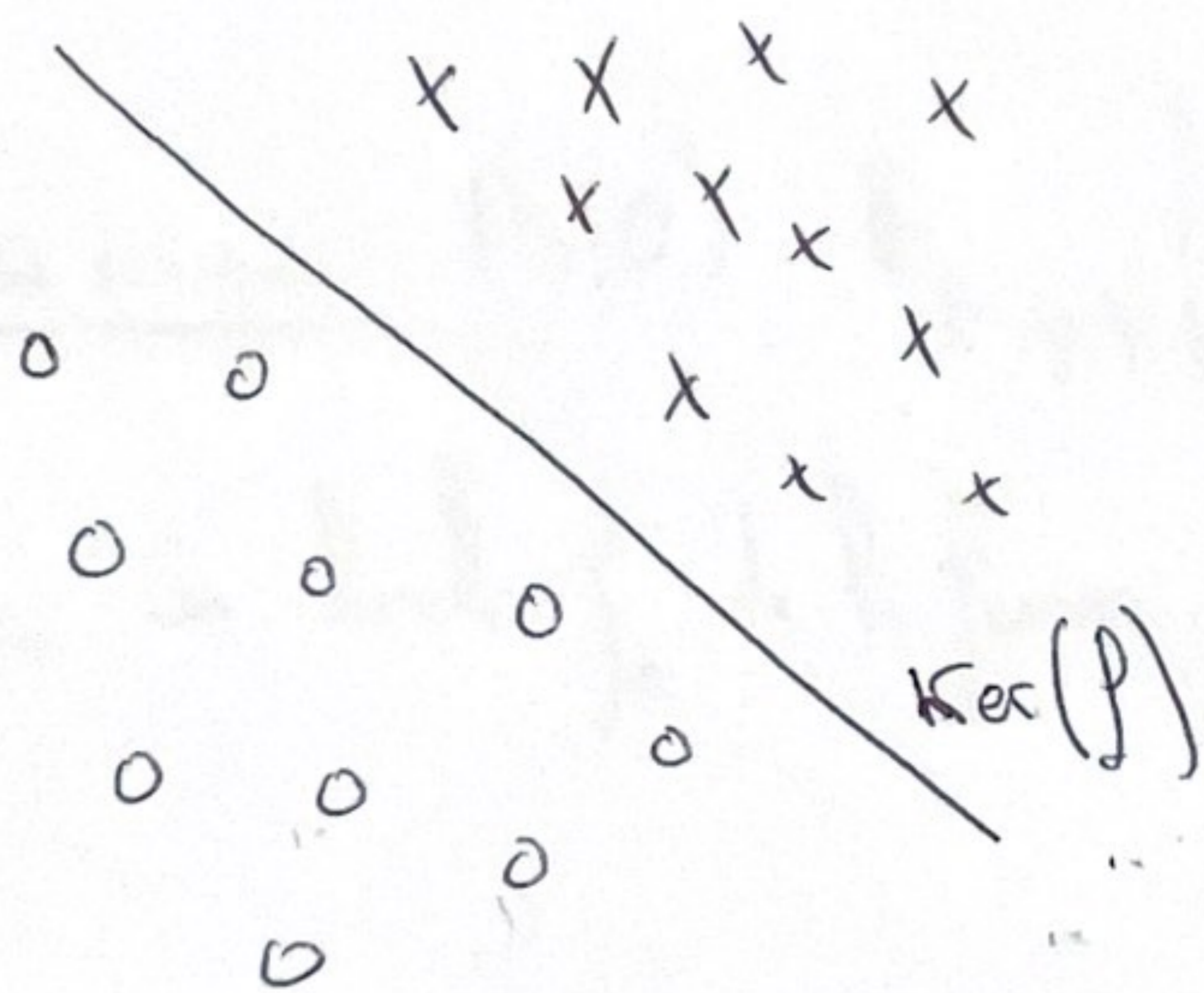
(9)

Soient $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^d \times \{-1, 1\}$

Identifié à "0" Identifié à "x"

On dit que ces points sont séparables linéairement si il existe une application linéaire $f: \mathbb{R}^d \rightarrow \mathbb{R}$ telle que

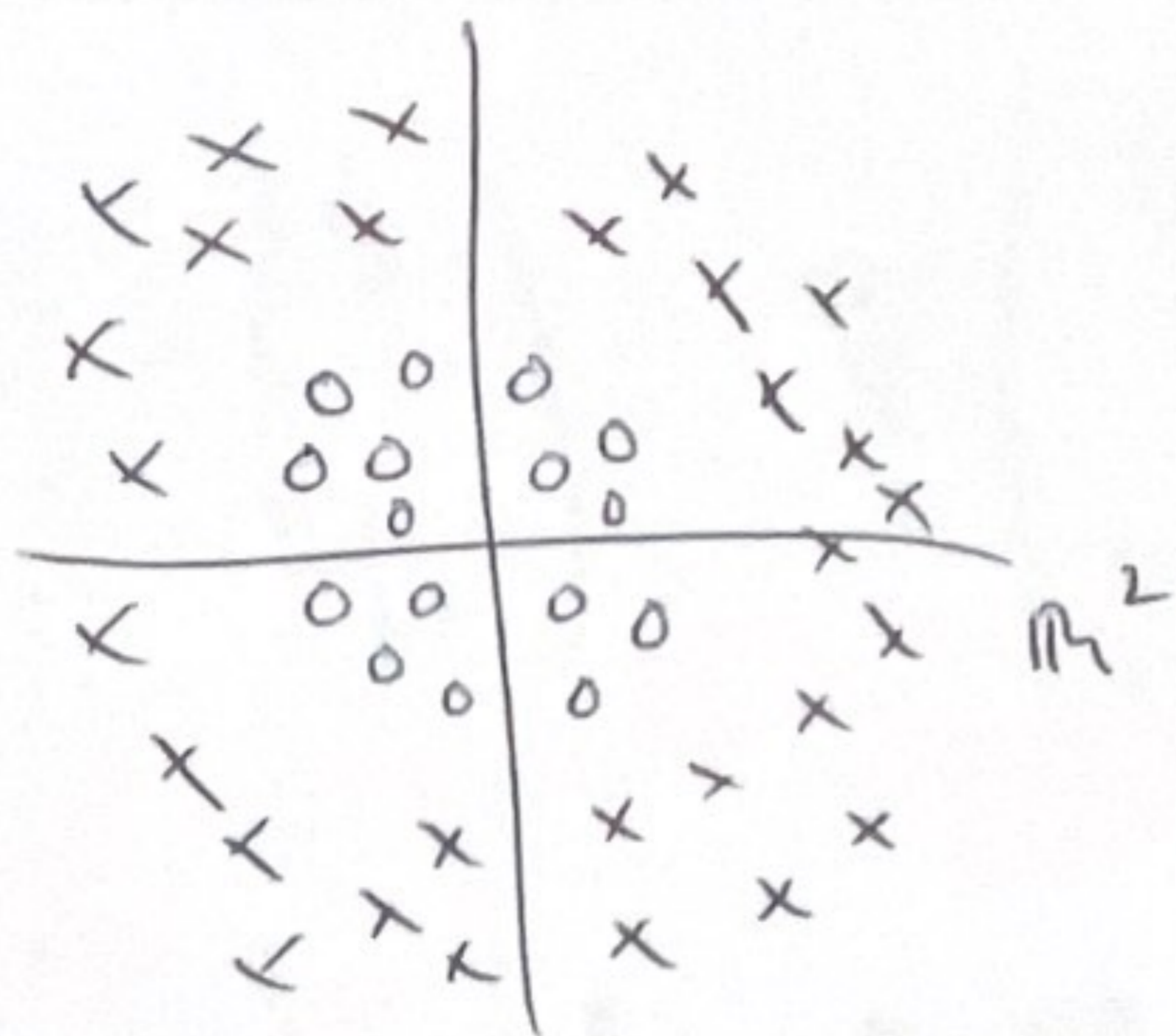
$$\forall i, y_i f(x_i) > 0.$$



Exercice: Trouver un ensemble de points qui ne sont pas séparables linéairement mais qui le deviennent en introduisant une "feature map"

Solution:

(5)



$\varphi(x) = \begin{pmatrix} \|x\|^2 \\ x \end{pmatrix} \in \mathbb{R}^2$
rend séparable les points
linéairement.

II. Régression Linéaire.

Contexte: $(x_1, y_1), \dots, (x_n, y_n) \in \underbrace{X}_{\subset \mathbb{R}^d} \times \underbrace{Y}_{\subset \mathbb{R}}$

Minimisation du risque $R(\theta) = \mathbb{E}_{(x,y)} \left((y - f_\theta(x))^2 \right), \theta \in \mathcal{H}$

Modèles linéaires: On cherche f_θ de la forme $f_\theta(x) = \theta^T \varphi(x)$

1) Moindres Carrés - Minimisation du Risque Empirique.

Empirical Risk Minimization (ERM)

$$\hat{\theta} \in \operatorname{argmin}_{\theta \in \mathbb{R}^d} \hat{R}(\theta) := \frac{1}{n} \sum_{i=1}^n (y_i - \theta^T \varphi(x_i))^2$$

Notation Matricielle

(6)

$$y := \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \quad \Phi = \begin{pmatrix} \varphi(x_1) \\ \vdots \\ \varphi(x_n) \end{pmatrix}$$

$$\hat{M}(\theta) = \frac{1}{n} \|y - \Phi\theta\|_2^2$$

Question: Comment trouver $\hat{\theta} \in \arg\min_{\theta} \hat{M}(\theta)$?

Hypothèse 1: Les colonnes de Φ sont indépendantes

(\Leftrightarrow) Il n'y a pas de dépendance linéaire entre les features.

En pratique, ce n'est souvent pas un problème si n est assez grand.

Définition - Proposition:

Sous l'hypothèse 1, $\hat{M}(\cdot)$ admet un unique minimiseur sur \mathbb{R}^d appelé estimateur des moindres carrés ordinaires ($\hat{\theta}$).

Il satisfait $\Phi^T \Phi \hat{\theta} = \Phi^T y$ (Équation Normale)

Donc

$$\hat{\theta} = (\Phi^T \Phi)^{-1} \Phi^T y.$$

Preuve Algébrique:

⑦

La fonction $\sigma \mapsto \hat{m}(\sigma)$ est différentiable sur M^d qui est ouvert.

Alors $\nabla \hat{m}(\hat{\sigma}) = 0$ est une condition nécessaire à l'optimalité de $\hat{\sigma}$.

$$\text{or, } \nabla \hat{m}(\sigma) = \frac{2}{A} \Phi^T (\Phi \sigma - y) \quad \forall \sigma$$

ce qui donne $\Phi^T \Phi \hat{\sigma} = \Phi^T y$ (équation normale).

Cette équation admet une unique solution, donc il existe un unique σ qui satisfait les conditions nécessaires d'optimalité.

Pour prouver que ces conditions sont suffisantes, il y a trois possibilités (au moins).

$$\bullet \nabla^2 \hat{m}(\sigma) = \frac{2}{A} \Phi^T \Phi \succ 0$$

$\Rightarrow \hat{\sigma}$ est un minimum local donc global car c'est le seul à vérifier la condition nécessaire.

$\bullet \hat{m}(\cdot)$ est convexe donc les conditions nécessaires deviennent suffisantes

\bullet Coercivité \Rightarrow on peut se ramener à l'étude sur un compact.

2) Le cas où l'hypothèse 1 n'est pas satisfaisable

(8)

Si l'hypothèse 1 n'est pas satisfaisable, $\text{Ker } \Phi \neq \emptyset$
et l'équation normale peut ne pas avoir de solution.

Par exemple, lorsqu'on peut écrire $\hat{\theta} = (\Phi^T \Phi)^{-1} \Phi^T y$
et que $\lambda_{\min}(\Phi^T \Phi) \rightarrow 0$, $\hat{\theta}$ "part à l'infini".

Pour résoudre ce problème, on "force" $\hat{\theta}$ à rester proche de 0 en pénalisant par sa norme.

$$\hat{\theta}_\lambda \in \underset{\theta}{\text{argmin}} \frac{1}{n} \|y - \Phi \theta\|_2^2 + \lambda \|\theta\|_2^2$$

$$\Rightarrow \hat{\theta}_\lambda = \frac{1}{n} \left(\frac{1}{n} \Phi^T \Phi + \lambda I \right)^{-1} \Phi^T y$$

Estimateur "Midge".

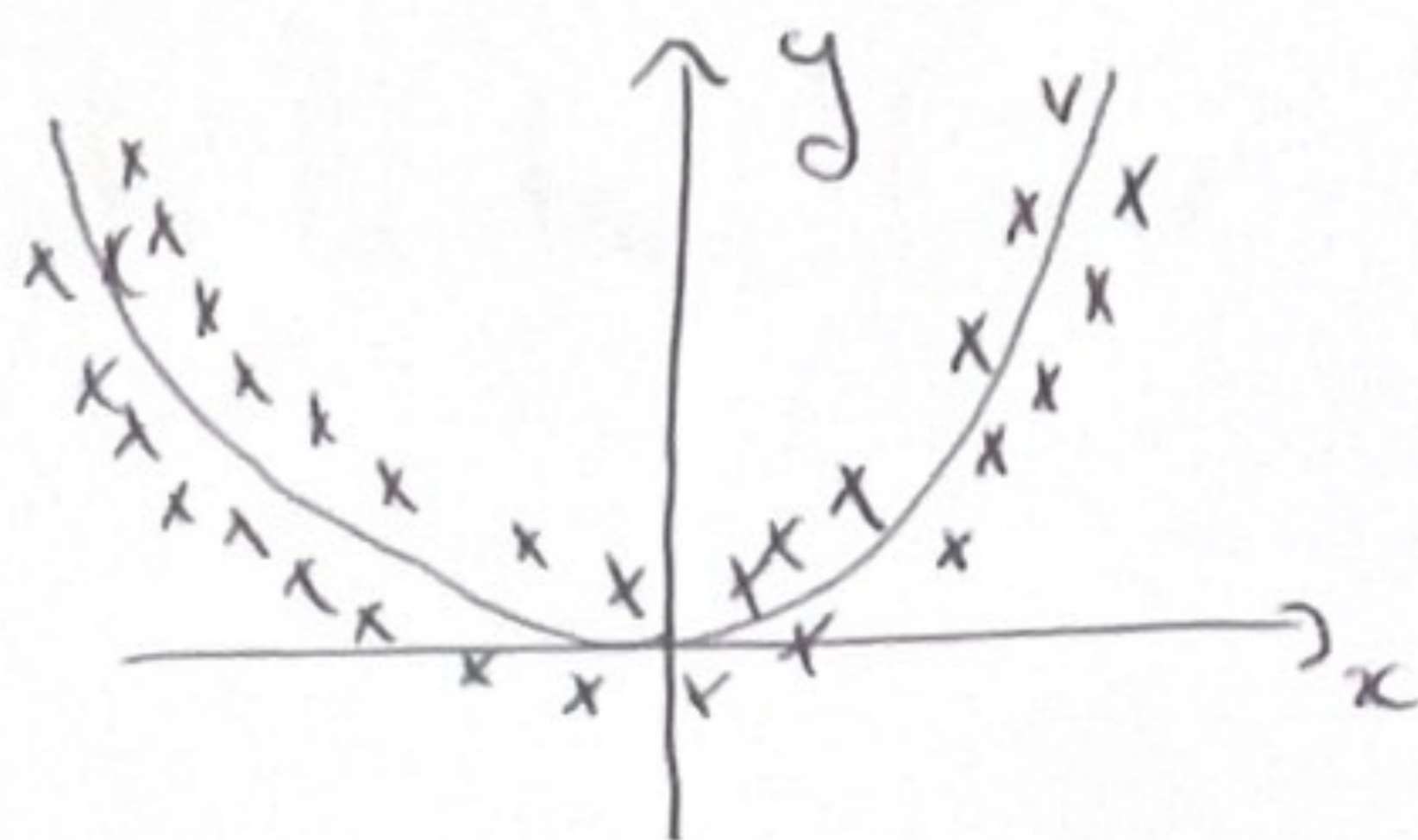
Remarque : En pratique c'est très utile lorsque :

- d est "grand" par rapport à n
- les features sont proches d'une situation de dépendance linéaire

3) Compromis sur/sous apprentissage.

9

Exemple: $y = x^2 + \text{bruit}$
 ↑
 centré.



Considérons: $\varphi_q(x) = \begin{pmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^q \end{pmatrix}$.

Cas extrême 1: $q = 0$

Le modèle ne peut apprendre que des droites horizontales

\Rightarrow Il n'est pas assez expressif pour comprendre le modèle

\Rightarrow On parle de sous-apprentissage.

Cas extrême 2: $q = n - 1$

Il existe un polynôme interpolateur (de Lagrange)

qui passe par tous les points

\Rightarrow Le modèle est trop expressif et apprend tout le bruit des données

=> On parle de sur-apprentissage.

(10)

En pratique, il faut prendre un q entre les deux (ici 2 marcesca).

III Modèles Linéaires pour la Classification - Mégression Logistique.

3) Maximum de vraisemblance - Régression Logistique

Fonction sigmoïde : $\sigma(z) := \frac{1}{1 + e^{-z}}$

Modèle Logistique : $p(y=1|x) = \frac{1}{1 + e^{-\theta^T x}} \quad y \in \{-1, 1\}$
 $= \sigma(\theta^T x) = \sigma(y \theta^T x)$

remarque

(10)

$$p(y = -1) = 1 - \sigma(\theta^T x)$$

$$= 1 - \frac{1}{1 + e^{-\theta^T x}} = \frac{e^{-\theta^T x}}{1 + e^{-\theta^T x}} = \frac{1}{1 + e^{\theta^T x}}$$

$$= \sigma(y \theta^T x)$$

donc $p(y|x) = \sigma(y \theta^T x)$.

Maximum de vraisemblance:

$$\hat{\theta} \text{ est argmax}_{\theta} \sum_{i=1}^n \ln \left(\frac{1}{1 + e^{-y_i \theta^T x_i}} \right)$$

$$\stackrel{\text{argmin}}{\theta} \frac{1}{n} \sum_{i=1}^n \ln(1 + e^{-y_i \theta^T x_i})$$

Remarque: Il s'agit de l'estimateur du minimum du risque empirique pour la perte

$$p_{\text{logistique}}(y, \beta) = \ln(1 + e^{-y\beta})$$