

Sélection de Modèles et Méthodes

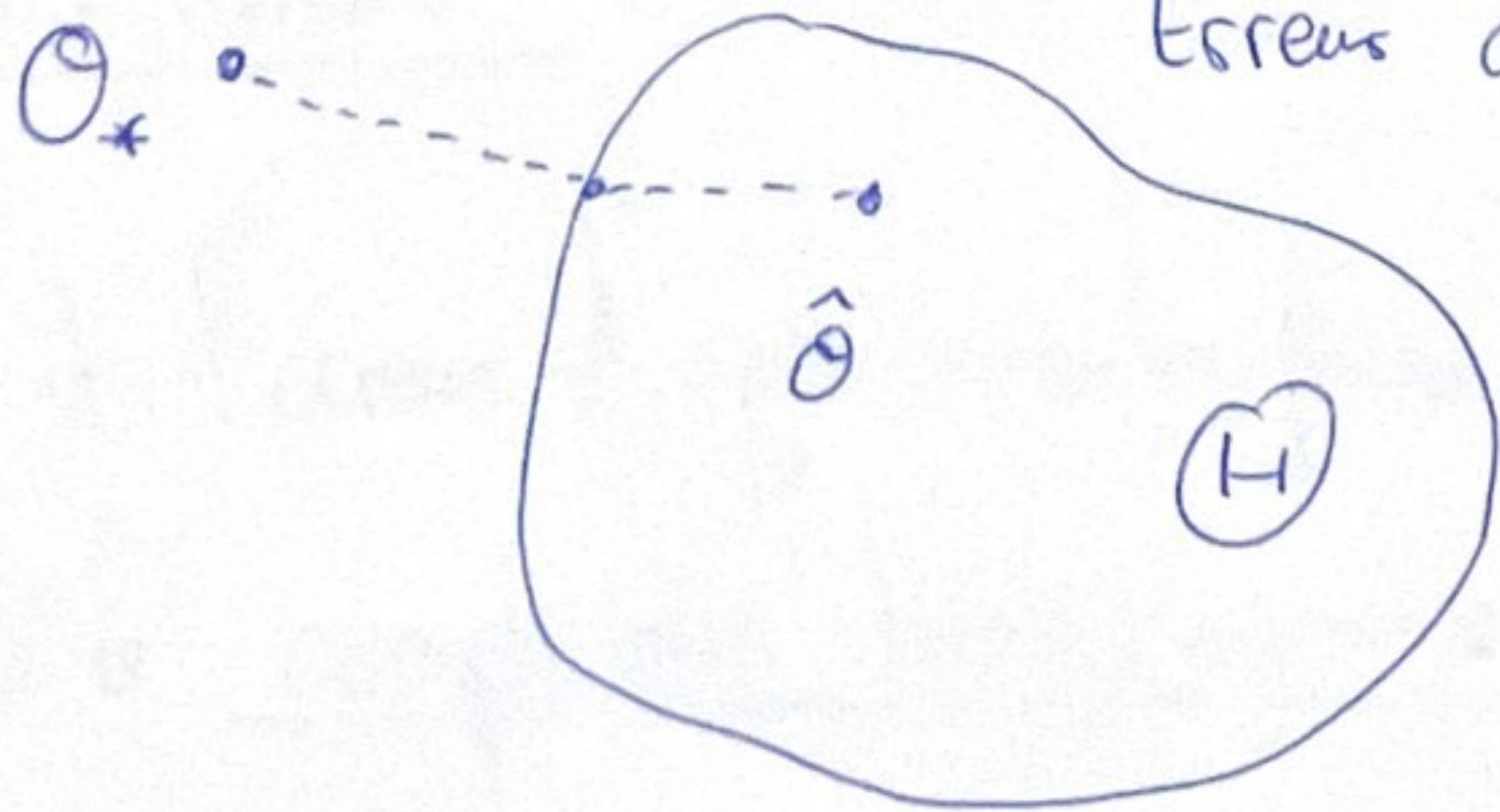
(1)

Paramétriques

I. Décomposition du risque

Soit (\mathcal{H}) une famille de paramètres, et soit $\hat{\theta} \in (\mathcal{H})$ un estimateur.

$$\underbrace{M(\hat{\theta})}_{\text{Risque de l'estimateur}} - M^* = \underbrace{\left(M(\hat{\theta}) - \inf_{\theta \in \mathcal{H}} M(\theta) \right)}_{\text{Erreur d'estimation}} + \underbrace{\left(\inf_{\theta \in \mathcal{H}} M(\theta) - M^* \right)}_{\text{Erreur d'approximation}}$$



1) L'erreur d'estimation

- Elle caractérise à quel point le risque de notre estimateur est proche du meilleur estimateur de la classe.

- Elle est reliée à la variance de l'estimateur.

Plus (H) est compliquée, plus il faudra de points pour que cette erreur soit petite.

2) L'erreur d'approximation

- Modélise à quel point (H) modélise la réalité fidèlement.

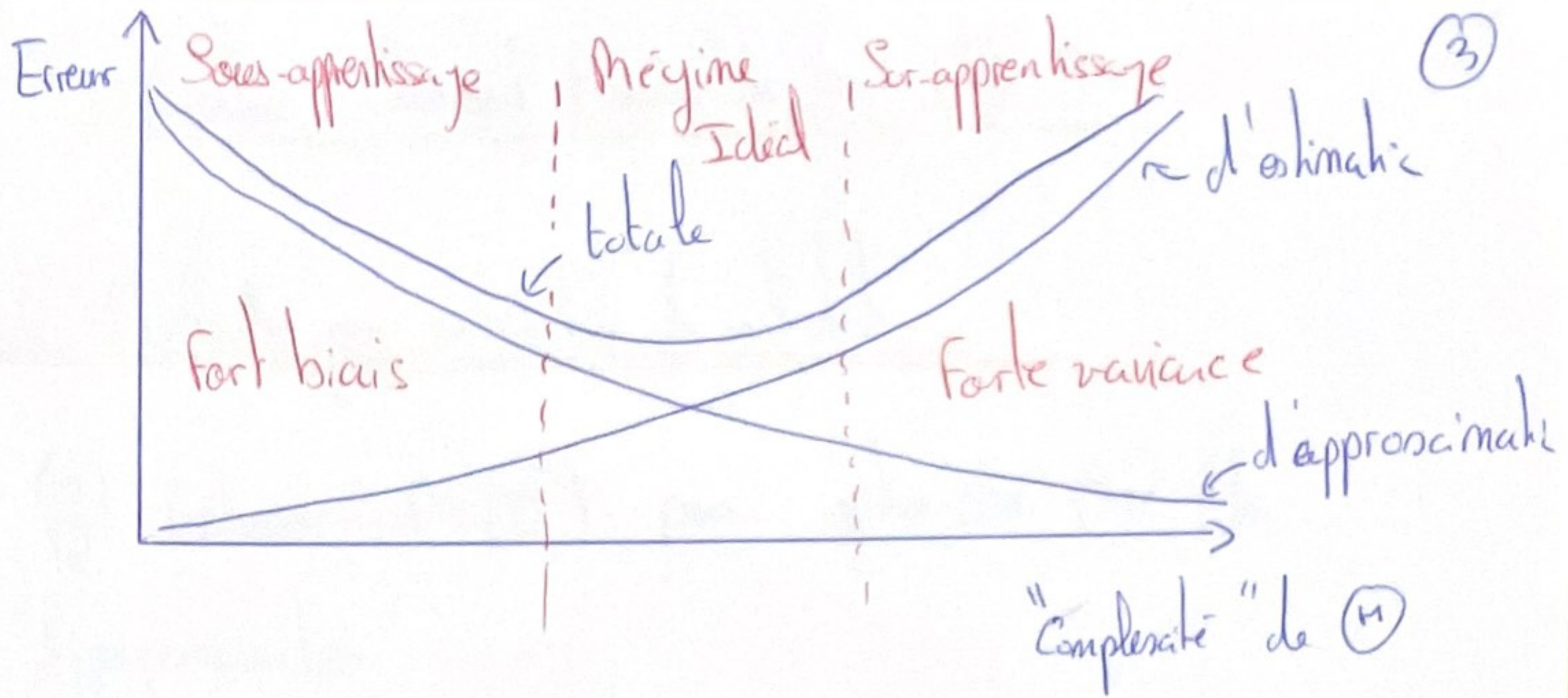
- Plus (H) est grande, plus l'erreur d'approximation sera faible.

- Il s'agit d'un terme de biais

3) Compromis Biais / Variance

L'erreur d'estimation et l'erreur d'approximation forment

un compromis appelé le compromis biais / variance

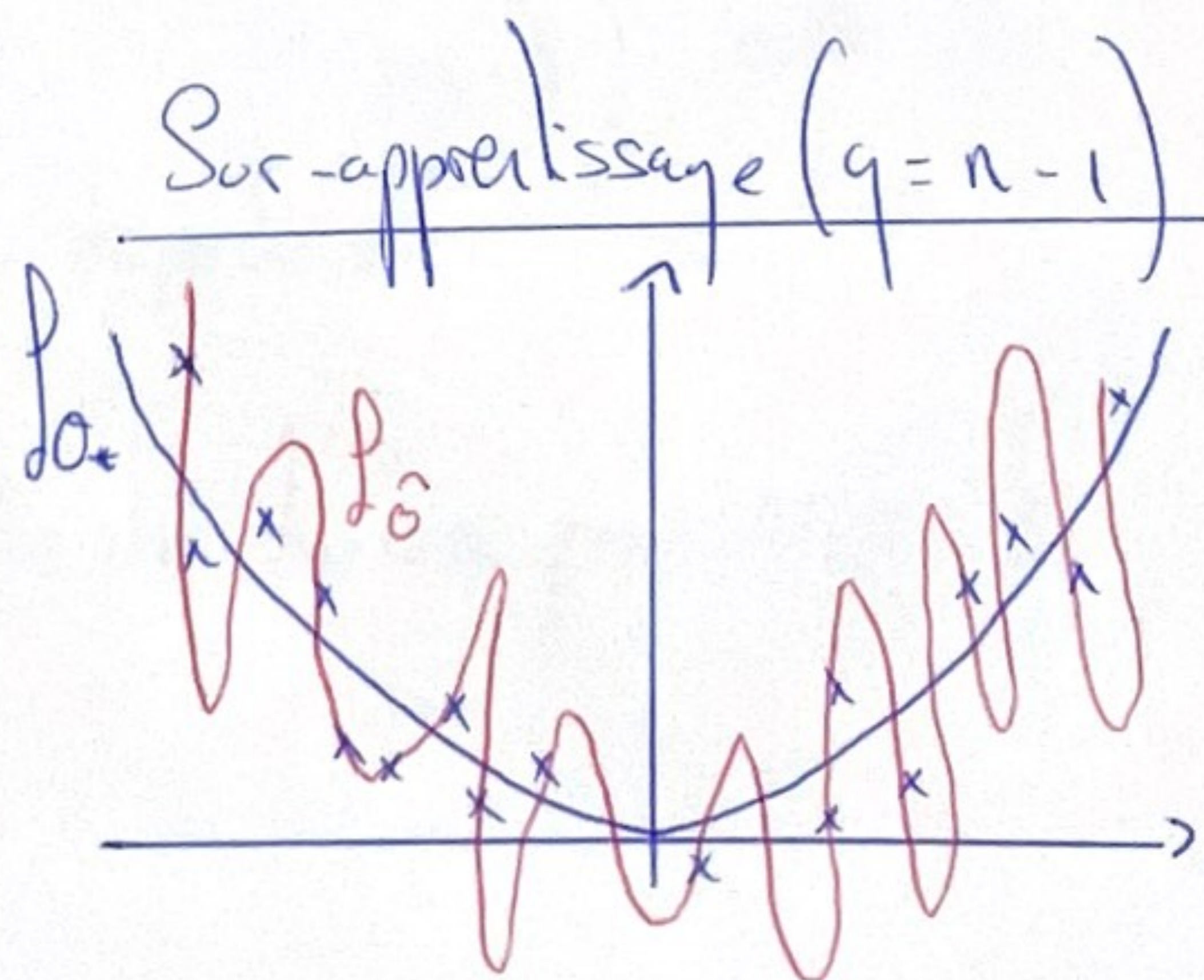
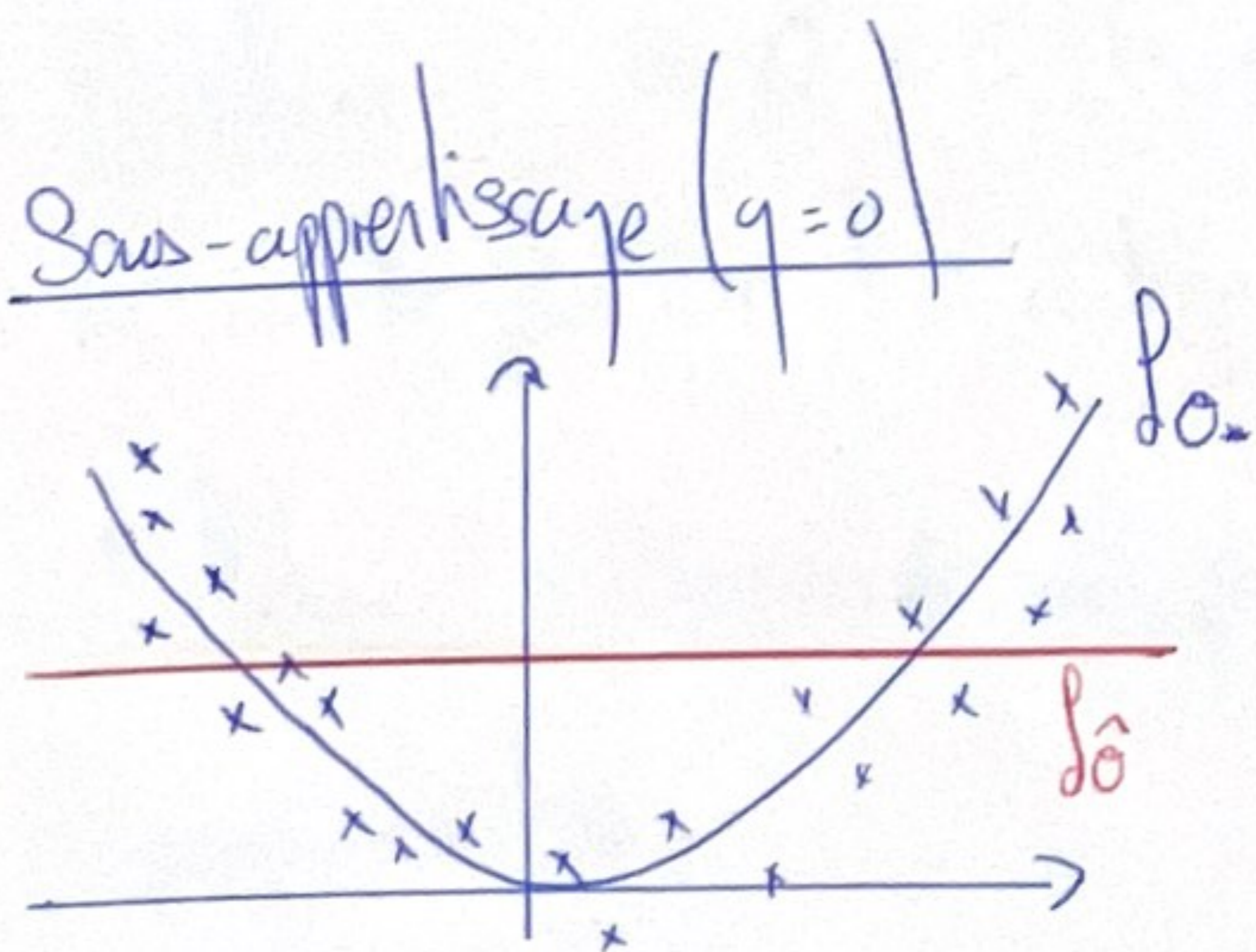


1) Régression Polynomiale

Vérité : $y = x^2 + \text{bruit}$.

$$\hat{\theta} \in \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n (y_i - \theta^T \phi_q(x_i))^2$$

$$\text{où } \phi_q(x) = \begin{pmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^q \end{pmatrix}$$



Question: Comment trouver un bon \mathcal{H} ?

(4)

II Méthodes par ensemble de validation

💡 Estimer $M(\hat{\theta})$ puis optimiser \mathcal{H} (les "hyperparamètres").

1) Découpage Train / Val

$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) \stackrel{i.i.d}{\sim} P$.

On peut estimer

$$M_{\mathcal{H}}(f_0) = E_{(x, y)} \left(\frac{f_0(x, y)}{\theta} \right) \text{ par } \frac{1}{n} \sum_{i=1}^n \frac{f_0(x_i, y_i)}{\theta}$$

Problème: si $\hat{\theta} = \hat{\theta}(x_1, y_1, \dots, x_n, y_n)$, cette idée ne marche plus car il y a trop de dépendances.

Solution: Mésestimer une partie des données pour régler les hyperparamètres.

$(x_1, y_1), \dots, (x_{|\text{Entraînement}|}, y_{|\text{Entraînement}|}) +$
 Jeu de données d'entraînement

$(x_{|\text{Entraînement}|+1}, y_{|\text{Entraînement}|+1}), \dots, (x_{|\text{Entraînement}| + |\text{Validation}|}, y_{|\text{Entraînement}| + |\text{Validation}|})$
 Jeu de données de validation.

Maintenant si $\hat{\theta} = \hat{\theta}(\text{Jeu de données d'entraînement})$

Alors $\hat{\theta} \perp \text{Jeu de données de validation}$.

Gn a alors $R(\hat{\theta}) \approx \frac{1}{|\text{Validation}|} \sum_{(x,y) \in \text{Validation}} P(y, f_{\hat{\theta}}(x)) \quad (*)$

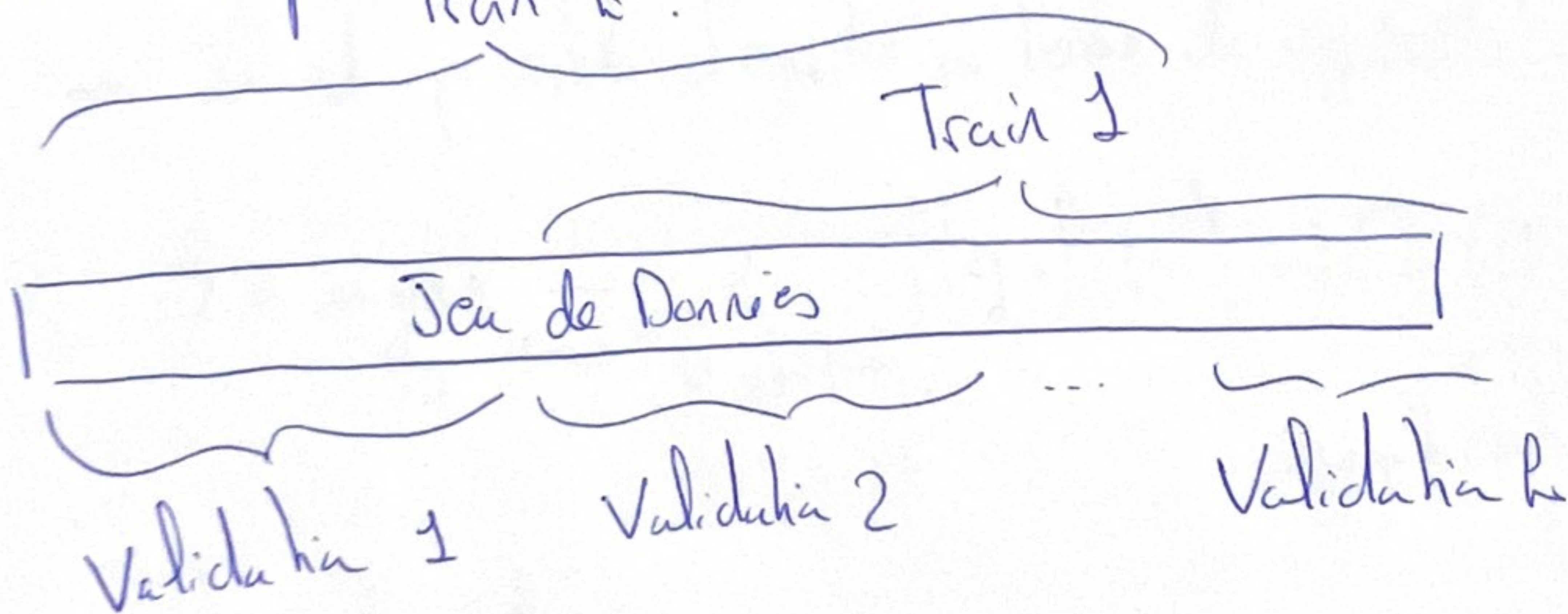
Pour régler les hyperparamètres, on peut donc

- (i) Estimer $\hat{\theta}$ pour plusieurs valeurs d'hyperparamètres
- (ii) Sélectionner les hyperparamètres qui minimisent (*).

Remarque: Il est aussi bon de réserver une partie des données, appelée ensemble de test, pour tester ~~les données~~ à la toute fin. (6)

2) Validation Croisée

⚡ L'idée de la validation croisée est de stabiliser l'étape de splitting du jeu de données en moyennant sur plusieurs découpages Train & Test.



III. Méthodes par pénalisation (Régularisation)

⑦

⚡ L'erreur d'entraînement se comporte à peu près comme l'erreur d'approximation, elle tend vers 0 quand la "complexité" de Θ tend vers $+\infty$.

\Rightarrow L'erreur d'entraînement capture mal l'explosion de la variance liée à la complexification

\Rightarrow Il est possible d'ajouter un terme de régularisation:

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmin}} \frac{1}{|\text{Train}|} \sum_{(x,y) \in \text{Train}} P(y, f_{\Theta}(x)) + \underbrace{\lambda \Omega(\Theta)}_{\text{Régularisation}}$$

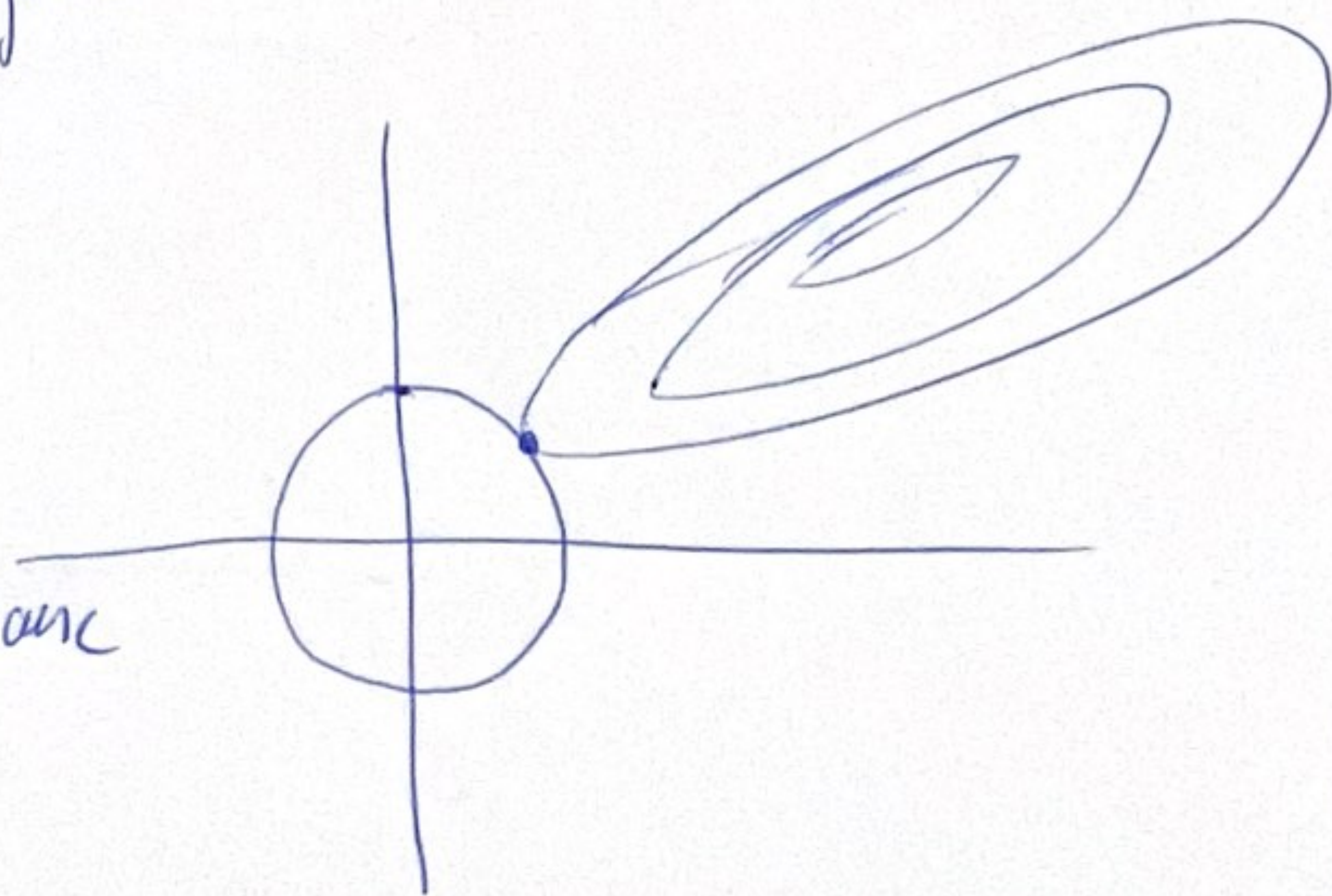
1) Régularisation L_2

$$\Omega(\Theta) = \frac{1}{2} \|\Theta\|^2 \Rightarrow \text{force } \Theta \text{ à être sans des bords}$$

de petite norme

ou: Régression Ridge

Weight Decay des réseaux de neurones.



2) Régularisation l_0 :

$\Omega(\theta) = \|\theta\|_0 = \text{nb de composantes de } \theta \text{ non nuls.}$

\Rightarrow force θ à avoir peu de coefficients non nuls.

\Rightarrow réduit le nombre de variables explicatives.

⚠ Problèmes difficiles à résoudre numériquement

3) Régularisation l_1 :

💡 $\|\cdot\|_1 =$ meilleure approximation convexe de $\|\cdot\|_0$

$=$ plus grande fonction g convexe tq $\forall x, g(x) \leq \|x\|_0$.

\Rightarrow On remplace $\|\cdot\|_0$ par $\|\cdot\|_1$, et on s'attend à avoir des effets similaires.

Le problème est maintenant "facile" à résoudre.

