

Prérequis :

- proba, calcul d'espérance
- loi gaussienne
- Algèbre linéaire, calcul matriciel.

Surapprendre : est-ce une si mauvaise idée ?

Cadre d'étude de l'apprentissage supervisé :

(X, Y) variables aléatoires -

X : entrée / descripteur / caractéristiques $\in \mathbb{R}^d$

Y : sortie / cible / étiquette / label

Modèle : on suppose que l'entrée et la sortie sont liées de la manière suivante.

$$Y = f^*(X) + \varepsilon$$

↳ ε est un terme de bruit \perp de tt tq $E(\varepsilon) = 0$.

↳ f^* caractérise le lien entre X et Y : en particulier si on connaissait f^* on pourrait prédire Y à la simple connaissance de X .

• Qd Y est une variable à valeurs réelles, on parlera de régression

ex) $Y =$ nb de clics sur un lien
 $X =$ description de l'utilisateur, âge, dernier type de page consulté, nb de fs, ...

ex) $Y =$ prix d'une action en bourse

ex) $Y =$ rendement d'un plant de céiches

• Qd Y est une variable à valeurs discrètes et finies, on parlera de classification

ex) $Y =$ spam / non-spam
 $X =$ contenu de l'email

$Y = 1$ si le tableau a été peint par Picasso.
 $X =$ image du tableau.

$Y = 1$ si une transaction est frauduleuse.

} classif^o binaire

(ex) $Y =$ chiffre manuscrit
 $Y =$ expression du visage (joie, tristesse, colère, ...)
 $X =$ image faciale

$\left. \begin{array}{l} \text{classif}^e \\ \text{multi-classe} \end{array} \right\}$

But: faire la prédiction de Y sachant X

ou encore estimer f par une fct \hat{f} tq $\hat{f} \simeq f^*$

ou encore apprendre f .

Comment construire \hat{f} ?

En apprentissage supervisé, on apprend à faire des prédictions
 à partir d'exemples étiquetés, c'est-à-dire accompagnés de la valeur que l'on
 voudrait prédire. On a donc accès à un échantillon / des observations

$(X_1, Y_1), \dots, (X_n, Y_n)$ iid

parmi lesquelles les étiquettes $(Y_i)_{i \leq n}$ vont jouer le rôle de "professeur" et
 vont venir "superviser" l'apprentissage de f^* par \hat{f} .

On pourra qualifier les ds° $(X_1, Y_1), \dots, (X_n, Y_n)$ de exemples d'entraînement
 de jeu d'apprentissage.

Nous allons donc construire \hat{f} à partir du jeu d'apprentissage tq

$$Y_i \simeq \hat{f}(X_i) \quad i = 1, \dots, n.$$

c'est-à-dire tq l'erreur d'entraînement soit faible

mais le véritable enjeu est que $Y_{\text{new}} \simeq \hat{f}(X_{\text{new}})$ sc $(X_{\text{new}}, Y_{\text{new}})$ des

données que \hat{f} n'a pas "vu" lors de l'entraînement, c'est-à-dire
 que \hat{f} généralise bien la prédiction à des nouvelles données.

En régression ($Y \in \mathbb{R}$), la "vraie" fonction f peut être interprétée

comme $f^* \in \operatorname{argmin}_{f: \mathbb{R}^d \rightarrow \mathbb{R}} \underbrace{\mathbb{E}_{(X,Y)} [(Y - f(X))^2]}_{=: \mathcal{R}(f)}$

On appelle cette quantité le risque du prédicteur f , c'est-à-dire l'erreur moyenne quadratique commise par f sur des observations (Y, X) .

Or la loi conjointe de X et Y est inconnue.

Donc le risque d'un prédicteur n'est pas calculable. On peut encore \ominus le minimiser pour trouver f^* !

À la place, on va chercher à minimiser le risque empirique

calculable à partir du jeu d'entraînement: étant donné $(x_1, y_1) \dots (x_n, y_n)$

$$\hat{\mathcal{R}}(f) := \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2.$$

⚠ On ne va pas chercher n'importe quelle fct qui minimise le risque empirique, on va se donner une classe de prédicteurs possibles. Pourquoi? On pourrait construire la

$$\text{fct } f: \mathbb{R}^d \rightarrow \mathbb{R}$$

↳ elle a un risque empirique nul!

$$x \mapsto \begin{cases} y_i & \text{qd } x = x_i \\ 10^9 & \text{sinon.} \end{cases}$$

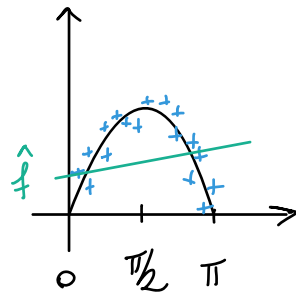
↳ elle fera certainement de très mauvaises prédictions pour la nouvelle donnée.

On parle de sur-apprentissage qd le prédicteur ne sait pas séparer l'information du bruit: l'erreur en test explose!

Un modèle qui sur-apprend est généralement un modèle trop complexe qui "calle" trop aux données et capture le bruit.

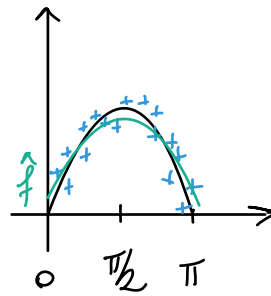
A l'inverse on parlera de **sous-apprentissage** qd le prédicteur est trop simple par rapport aux données.

(ex) $X_i \in [0; \pi]$ $Y_i \in \mathbb{R}$ $Y_i = f^*(X_i) + \epsilon_i$ $f^* = \sin$



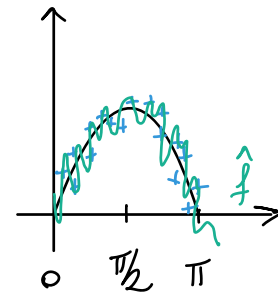
Sous-apprentissage

On contraint \hat{f} à être une fct affine. On trouve la "meilleure fct affine"



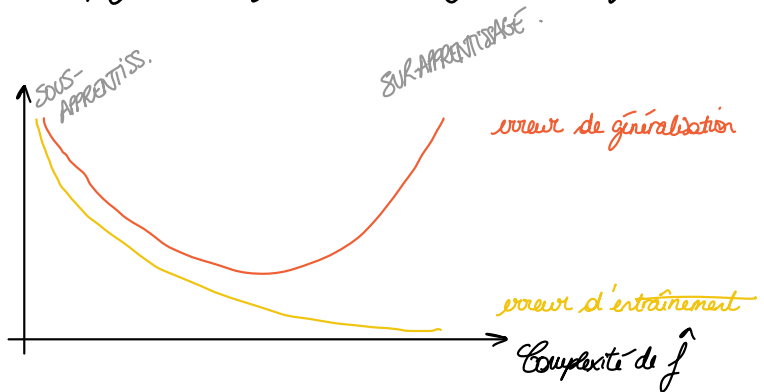
Bon compromis

polynôme de deg 4.



Sur-apprentissage

polynôme de deg 15.



C'est un schéma classique d'un 1^{er} cours de ML. Pourtant, aujourd'hui des prédicteurs de + en + complexes et entraînés (réseaux de neurones) et leur performance s'améliore tj. On va creuser plus en détail ce phénomène par un modèle en particulier : le **modèle linéaire**, c'ad f est une fct affine des entrées X et ce qui jouera le rôle de la "complexité" du prédicteur, ce sera la **dimension d** de l'espace des entrées. Nous venons finalement que le 1^{er} graphique n'est qu'un bout de l'histoire.

HELP ON GAUSSIAN?

Hyp: $Y_i = X_i^T \theta^* + \varepsilon_i$ pour $\theta^* \in \mathbb{R}^d$

$X_i \in \mathbb{R}^d$ chaque composante est $\mathcal{N}(0, 1)$ iid.

$\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$

↳ pouvez-vous me dire $\mathbb{E}[XX^T]$? Id.

On sait que $\theta^* = \underset{\theta}{\operatorname{argmin}} \mathbb{E}[(Y - X^T \theta)^2]$ pour (X, Y) copie de (X_i, Y_i)

$$\forall \theta \quad \mathbb{E}[(Y - X^T \theta)^2] \geq \mathbb{E}[(Y - X^T \theta^*)^2]$$

en effet
$$\mathbb{E}[(Y - X^T \theta)^2] = \mathbb{E}[(Y - X^T \theta^* + X^T \theta^* - X^T \theta)^2]$$
$$= \mathbb{E}[(Y - X^T \theta^*)^2] + \underbrace{\mathbb{E}[(X^T \theta^* - X^T \theta)^2]}_{\geq 0} + 2 \mathbb{E}[\dots]$$

Pour estimer θ^* , on va chercher à minimiser le risque empirique

$$\hat{\theta} \in \underset{\theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (Y_i - X_i^T \theta)^2 = \frac{1}{n} \|Y - X\theta\|_2^2$$

Fct conv en θ : si min local \Rightarrow min global.

Fct diff en θ : pt critique \Rightarrow min.

$$\nabla \hat{R}(\theta) = \begin{pmatrix} \frac{\partial \hat{R}(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial \hat{R}(\theta)}{\partial \theta_d} \end{pmatrix} = \frac{1}{n} 2 X^T (X\theta - Y)$$
$$= 0 \Leftrightarrow X^T X \theta = X^T Y$$

$d < n$

Régime sous-paramétrisé: plus d'éqs que d'inconnus.

$X^T X$ est inversible.

$$\hat{\theta} = (X^T X)^{-1} X^T Y.$$

$$\begin{aligned} \mathcal{R}(\hat{\theta}) &= \mathbb{E} \left[(X^T \hat{\theta} - y)^2 \right] \\ &= \mathbb{E} \left[(X^T \hat{\theta} - X^T \theta^*)^2 \right] + \mathbb{E} \left[(X^T \theta^* - y)^2 \right]. \end{aligned}$$

$$\begin{aligned} \Rightarrow \mathcal{R}(\hat{\theta}) - \mathcal{R}(\theta^*) &= \mathbb{E} \left[(X^T \hat{\theta} - X^T \theta^*)^2 \right]. \\ &\stackrel{\text{min}}{=} \mathbb{E} \left[\langle X, \hat{\theta} - \theta^* \rangle^2 \right] \\ &= \mathbb{E} \left[(\hat{\theta} - \theta^*)^T X X^T (\hat{\theta} - \theta^*) \right] \\ &= (\hat{\theta} - \theta^*)^T \underbrace{\mathbb{E} [X X^T]}_{\text{Id.}} (\hat{\theta} - \theta^*) \\ &= \|\hat{\theta} - \theta^*\|_2^2 \end{aligned}$$

$$\begin{aligned} \mathbb{E}_{\mathcal{D}_n} \left[\mathcal{R}(\hat{\theta}) - \mathcal{R}(\theta^*) \right] &= \mathbb{E}_{\mathcal{D}_n} \left[\left\| (X^T X)^{-1} X^T y - \theta^* \right\|_2^2 \right] \\ &= \mathbb{E}_{\mathcal{D}_n} \left[\left\| (X^T X)^{-1} X^T X \theta^* + (X^T X)^{-1} X^T \varepsilon - \theta^* \right\|_2^2 \right] \\ &= \mathbb{E}_{\mathcal{D}_n} \left[\left\| (X^T X)^{-1} X^T \varepsilon \right\|_2^2 \right] \\ &= \mathbb{E}_{\mathcal{D}_n} \left[\text{Tr} \left(\varepsilon^T X (X^T X)^{-1} (X^T X)^{-1} X^T \varepsilon \right) \right] \\ &= \sigma^2 \mathbb{E} \left[\text{Tr} \left((X^T X)^{-1} \right) \right] \end{aligned}$$

$X \in \mathbb{R}^{n \times d}$ est une matrice gaussienne

On sait que $X^T X \in \mathbb{R}^{d \times d}$ a une distrib^o de Wishart ac n ddl :

- presque sûrement inv.

théorème

$$\mathbb{E}_{\mathcal{D}_n} \left(\mathcal{R}(\hat{\theta}) - \mathcal{R}(\theta^*) \right) = \begin{cases} \frac{\sigma^2 d}{n-d-1} & \text{si } n \geq d+2 \\ \infty & \text{si } n = d \text{ ou } n = d+1. \end{cases}$$

$n \leq d$

"Régime surparamétré"

Il existe une infinité d'ERTL, de ce cas le risque empirique est nul dès lors que $Y = X\theta$,
i.e. $\forall i, Y_i = X_i^T \theta$.

On va se concentrer sur une solution en particulier,
celle qui est de norme ℓ^2 min:

$$\hat{\theta} \in \operatorname{argmin} \|\theta\|_2^2 \text{ s.t. } Y = X\theta \quad (P)$$

$$\text{i.e. } \hat{\theta} = X^T (XX^T)^{-1} Y$$

On peut tjrs réécrire (P) $\hat{\theta} \in \operatorname{argmin} \|\hat{\theta} + h\|_2^2$ s.t. $h \in \operatorname{Ker} X$
Or $\hat{\theta} \in \operatorname{Im} X^T = (\operatorname{Ker} X)^\perp$ de $\operatorname{argmin} \|\hat{\theta}\|_2^2 + \|h\|_2^2$ s.t. $h \in \operatorname{Ker} X$.
 $\Rightarrow h = \underline{0}!!$

$$\begin{aligned} \mathbb{E}_{2n} [R(\hat{\theta}) - R(\theta^*)] &= \mathbb{E}_{2n} [\|\hat{\theta} - \theta^*\|_2^2] \\ &= \mathbb{E}_{2n} [\|X^T (XX^T)^{-1} Y - \theta^*\|_2^2] \\ &= \mathbb{E}_{2n} [\|X^T (XX^T)^{-1} X \theta^* + X^T (XX^T)^{-1} \varepsilon - \theta^*\|_2^2] \\ &= \mathbb{E}_{2n} [\underbrace{\|X^T (XX^T)^{-1} X - I\|_2^2}_{\substack{\text{Matrice de} \\ \text{projection sur} \\ \operatorname{Ker} X \text{ au signe près.}}} \|\theta^*\|_2^2] + \mathbb{E}_{2n} [\underbrace{\|X^T (XX^T)^{-1} \varepsilon\|_2^2}_{\text{}}] \end{aligned}$$

$$\begin{aligned}
\mathbb{E} \left[\|\mathbf{X}^T (\mathbf{X}\mathbf{X}^T)^{-1} \boldsymbol{\varepsilon}\|_2^2 \right] &= \mathbb{E} \left[\boldsymbol{\varepsilon}^T (\mathbf{X}\mathbf{X}^T)^{-1} \mathbf{X}\mathbf{X}^T (\mathbf{X}\mathbf{X}^T)^{-1} \boldsymbol{\varepsilon} \right] \\
&= \mathbb{E} \left[\text{Tr} \left(\begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \end{array} \right) \right] \\
&= \sigma^2 \mathbb{E} \left[\text{Tr} \left((\mathbf{X}\mathbf{X}^T)^{-1} \right) \right] \\
&= \frac{\sigma^2 m}{d-m+1} \quad \text{si } d \geq m+2.
\end{aligned}$$

$$\begin{aligned}
\mathbb{E} \left[(\boldsymbol{\theta}^*)^T \left(\mathbf{I} - \mathbf{X}^T (\mathbf{X}\mathbf{X}^T)^{-1} \mathbf{X} \right) \boldsymbol{\theta}^* \right] &= \mathbb{E} \left[\left\| \text{Proj}_{(\text{Im} \mathbf{X}^T)^\perp} (\boldsymbol{\theta}^*) \right\|_2^2 \right] \\
&= \mathbb{E} \left[\left\| \text{Proj}_{\text{Ker} \mathbf{X}} (\boldsymbol{\theta}^*) \right\|_2^2 \right]
\end{aligned}$$

\nearrow
 matrice de projection
 de idempotente !

Focus sur $\mathbb{E} \left[(\boldsymbol{\theta}^*)^T \mathbf{X}^T (\mathbf{X}\mathbf{X}^T)^{-1} \mathbf{X} \boldsymbol{\theta}^* \right]$:

On pose la rotation $R^{(e)}$ $\boldsymbol{\theta}^* = \|\boldsymbol{\theta}^*\|_2 R^{(e)} e_e$, alors

$$\begin{aligned}
&= \|\boldsymbol{\theta}^*\|_2^2 \mathbb{E} \left[e_e^T R^{(e)T} \mathbf{X}^T (\mathbf{X}\mathbf{X}^T)^{-1} \mathbf{X} R^{(e)} e_e \right] \\
&= \|\boldsymbol{\theta}^*\|_2^2 e_e^T \mathbb{E} \left[(\mathbf{X}R^{(e)})^T (\mathbf{X}R^{(e)} R^{(e)T} \mathbf{X}^T)^{-1} (\mathbf{X}R^{(e)}) \right] e_e.
\end{aligned}$$

par orthogonalité de $R^{(e)}$

Or la matrice aléatoire $\mathbf{X}R^{(e)}$ est de même loi que \mathbf{X} , on peut s'en convaincre :

$$\begin{aligned}
(x_1, \dots, x_d) R^{(e)} &= (u_1, \dots, u_d) \\
\text{avec } u_j &= (x_1, \dots, x_d) R_{\cdot j}^{(e)} = \sum_{k=1}^d x_k R_{kj}^{(e)}
\end{aligned}$$

C'est une combinaison linéaire de gaussiennes de U_j et de loi gaussienne $\mathbb{E}[U_j] = \mathbb{E}\left[\sum_k X_k R_{kj}^{(e)}\right] = 0$

$$\text{Var}[U_j] = \sum_{k=1}^d R_{kj}^{2(e)} \underset{1}{\text{Var}(X_k)} = \|R_{\cdot j}^{(e)}\|_2^2 = 1.$$

Donc finalement

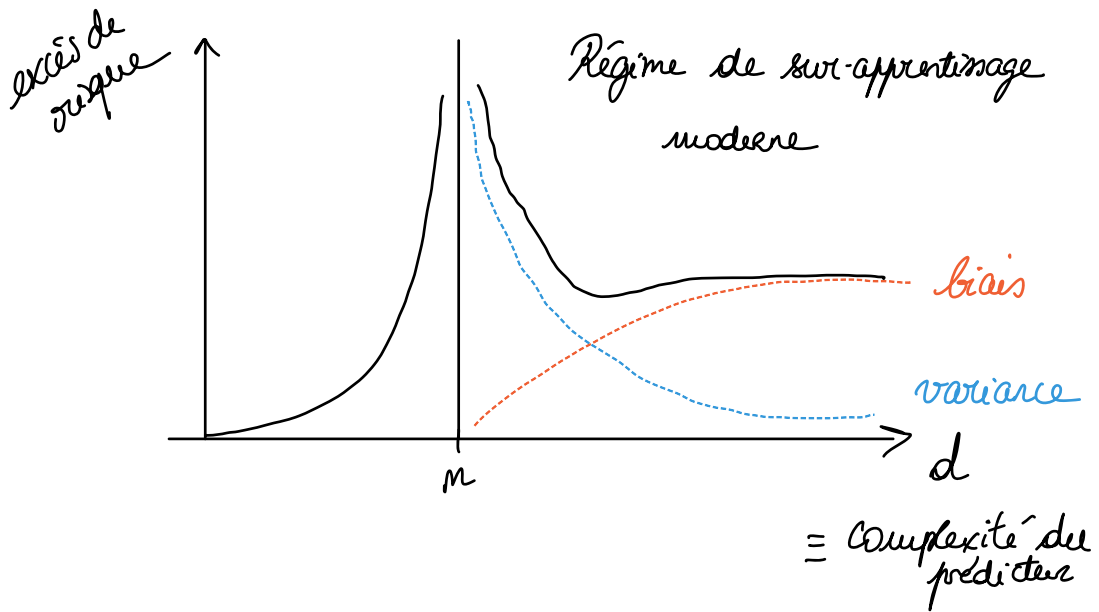
$$= \|\theta^*\|_2^2 e^T \mathbb{E}\left[X^T (XX^T)^{-1} X\right] e. \quad \underline{\underline{\theta e}}$$

$$\begin{aligned} \mathbb{E}\left[(\theta^*)^T X^T (XX^T)^{-1} X \theta^*\right] &= \frac{\|\theta^*\|_2^2}{d} \sum_{l=1}^d e^T \mathbb{E}\left[X^T (XX^T)^{-1} X\right] e \\ &= \frac{\|\theta^*\|_2^2}{d} \mathbb{E}\left[\text{Tr}\left(X^T (XX^T)^{-1} X\right)\right] \\ &= \frac{\|\theta^*\|_2^2}{d} \mathbb{E}\left[\text{Tr}(\text{Id}_n)\right] \\ &= \frac{\|\theta^*\|_2^2}{d} n. \end{aligned}$$

Finalement, $\mathbb{E}\left[(\theta^*)^T (I - X^T (XX^T)^{-1} X) \theta^*\right] = \|\theta^*\|_2^2 \frac{d-n}{d}$

$$\mathbb{E}_n\left[\mathcal{R}(\hat{\theta}) - \mathcal{R}(\theta^*)\right] = \int \sigma^2 \frac{n}{d-n+1} + \|\theta^*\|_2^2 \frac{d-n}{d} \quad \text{si } d \geq n+2$$

$\left. \begin{array}{l} +\infty \\ \text{si } \begin{cases} d = n+1 \\ d = n \end{cases} \end{array} \right\}$



Ds le régime de sur-apprentissage moderne, on voit que le risque n'explor pas qd $d \nearrow$, finalement on ne fait pas n'importe quoi qd on est à droite du graphique.

Evidemment ceci est vrai pour un minimiseur du risque empirique particulière, et pas n'importe lequel -

↳ le biais augmente en d/n : qd $d > n$, $\hat{\theta}$ vit dans le "raw space" de X , espace de dim n .

↳ La variance décroît en d/n ds le cas surparamétré

Pour des prédicteurs plus complexes tj les réseaux de neurones, l'entraînement revient à minimiser une fct non-convexe, donc il faut encore comprendre vers quel minimiseur du risque empirique l'entraînement nous amène et si celui-ci a de bonnes propriétés de généralisation.