

Intelligence Artificielle

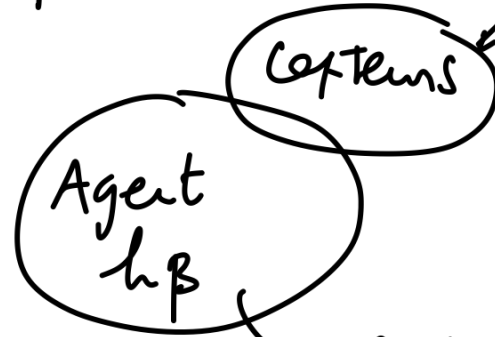
Apprentissage

Raisonnement
logique

Renforcement

Supervisé

Non-supervisé



sortie
 $y(x)$

Ensemble des
données
d'entraînement

Apprentissage Supervisé

agent a accès a un ensemble de données

$$\{x^{(i)}, t^{(i)}\}_{i=1}^N$$

observables $\vec{x}^{(i)} \in \mathbb{R}^D$

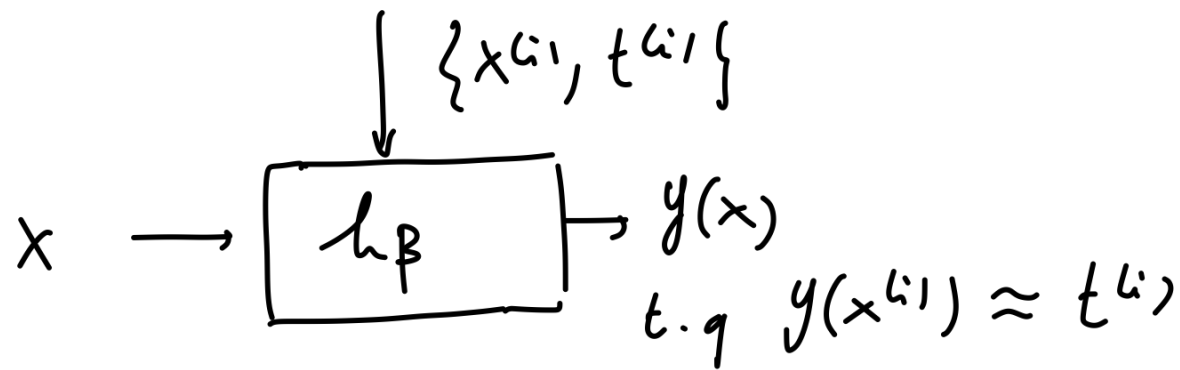
$$x^{(i)} = (x_1, x_2, x_3, \dots, x_D)$$

prédiction à produire à partir de observations

$$t^{(i)} \in \mathbb{R}$$

$\vec{x}^{(i)}$ = vecteur caractéristique $x_j^{(i)}$ → caractéristiques

Objectif : apprendre un modèle : h_{β}

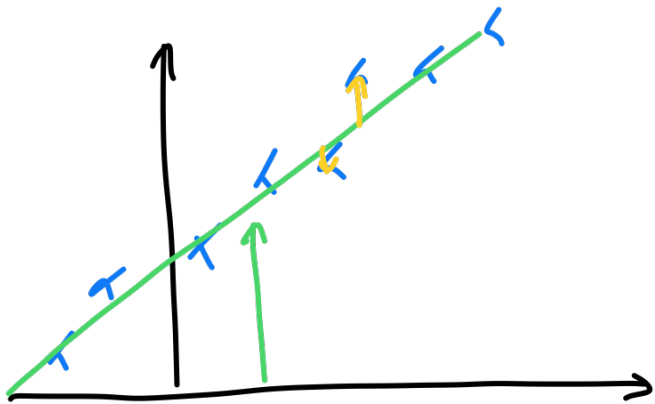


X : espace des données d'entrée

Y : espace de sortie

Une première approche: choisir $h_{\beta}(x)$ comme une combinaison linéaire

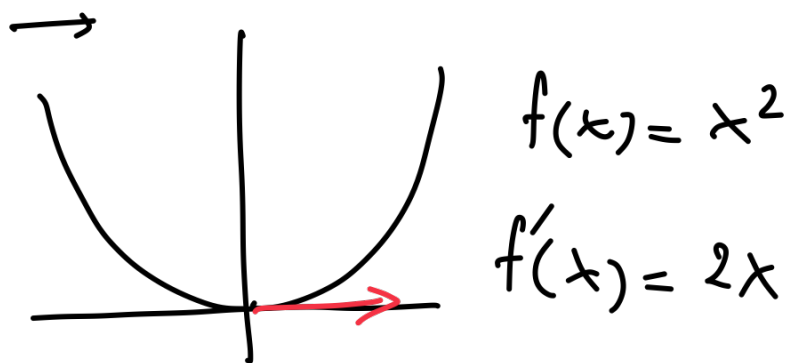
$$h_{\beta}(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_D x_D \quad \vec{\beta} \in \mathbb{R}^{D+1}$$



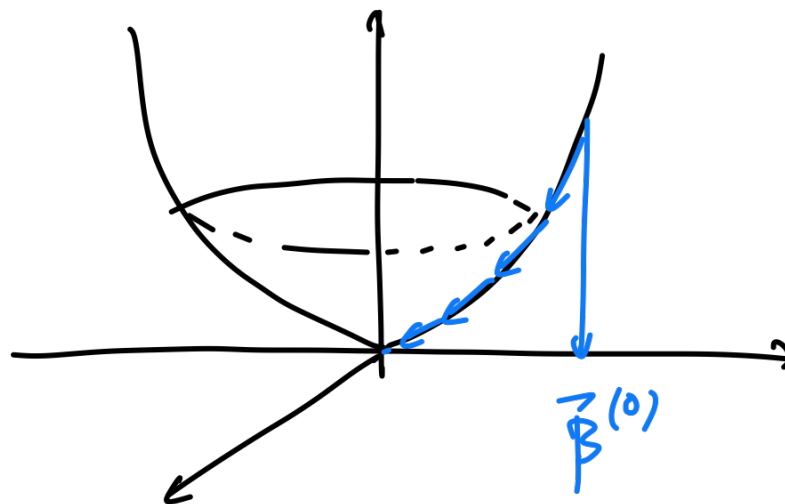
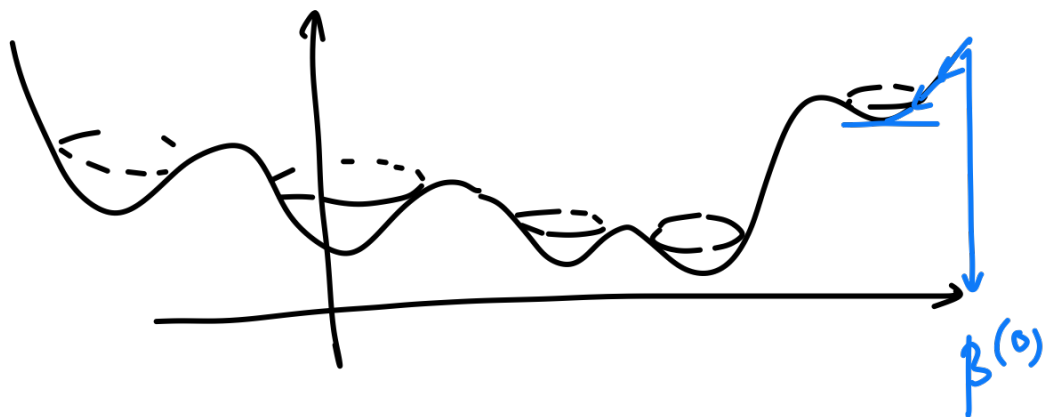
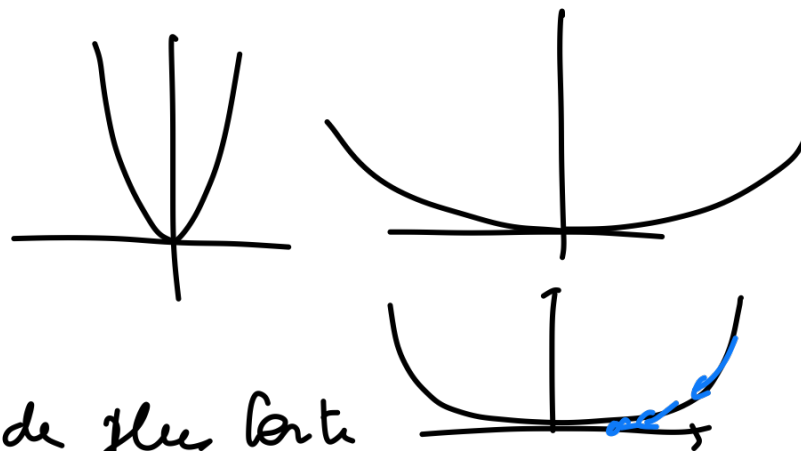
Question Comment déterminer les coefficients β_j de façon à ce que $h_{\beta}(\vec{x}^{(i)}) \approx t^{(i)}$

$$\begin{aligned} \rightarrow \mathcal{L}(\vec{\beta}) &= \frac{1}{N} \sum_{i=1}^N (t^{(i)} - h_{\beta}(x^{(i)}))^2 \quad (\text{fonction de coût / loss}) \\ &= \frac{1}{N} \sum_{i=1}^N (t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2 \end{aligned}$$

→ Pour trouver les β_j on minimise la fonction de coût $l(\beta)$
par rapport à $\vec{\beta}$



direction de plus forte
diminution de la fonction



$$l(\beta) = \frac{1}{N} \sum_{i=1}^N (t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2$$

$$\text{grad}_{\beta} l(\beta) = \left[\frac{\partial l}{\partial \beta_0}, \dots, \frac{\partial l}{\partial \beta_D} \right]$$

$$\frac{d f(x)^2}{dx} = 2 f(x) \cdot f'(x)$$

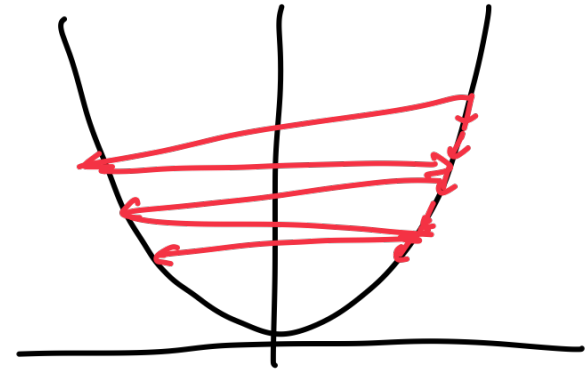
$$\frac{\partial l}{\partial \beta_0} = \frac{1}{N} \sum_{i=1}^N 2 (t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)})) \cdot (-1)$$

$$\frac{\partial l}{\partial \beta_j} = \frac{1}{N} \sum_{i=1}^N 2 (t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)})) \cdot (-x_j^{(i)})$$

$j \neq 0$

$x_j^{(i)}$

#2 Descent de gradient (BATCH)



Soit $\beta^{(0)}$ un itéré initial

for le im. max. arage (maxIter):

$$\left\{ \begin{array}{l} \beta_0 \leftarrow \beta_0 - \eta \frac{\partial \mathcal{L}}{\partial \beta_0} \\ \beta_j \leftarrow \beta_j - \eta \frac{\partial \mathcal{L}}{\partial \beta_j} \end{array} \right.$$

taux d'apprentissage
("learning rate")

Descent de gradient Stochastique (SGD)

→ for iter in maxIter

→ shuffle $\{x^{(i)}, t^{(i)}\}$

→ for i in $\{x^{(i)}, t^{(i)}\}$

$$\beta \leftarrow \beta - \eta \text{grad}_{\beta} \left(t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}) \right)^2$$

Descent de type MINI BATCH

(à chaque itération on échantillonne un sous ensemble de données de taille $m < N$)

#2 Annulation de la dérivée (Résolution des équations normales)

$$\rightarrow \vec{t} = \begin{bmatrix} t^{(1)} \\ \vdots \\ t^{(N)} \end{bmatrix} \quad \vec{\beta} = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_D \end{bmatrix} \rightarrow \text{paramètres du modèle / coefficients de régression}$$

$$\tilde{x} = [1, x_1, x_2, \dots, x_D] \\ = [1, \vec{x}]$$

$$h_{\beta}(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_D x_D = [\beta_0 \dots \beta_D] \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_D \end{bmatrix} \\ = \vec{\beta}^T \tilde{x}$$

$$\underset{\sim}{X} = \begin{bmatrix} 1 & \xrightarrow{x_1^{(1)}} & & & \\ \vdots & \vdots & & & \\ 1 & \xrightarrow{x_1^{(N)}} & & & \end{bmatrix} = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_D^{(1)} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_1^{(N)} & x_2^{(N)} & \dots & x_D^{(N)} \end{bmatrix}$$

vecteur d'erreur $\vec{\varepsilon} = \vec{t} - \underset{\sim}{X} \vec{\beta}$

$$\begin{bmatrix} \varepsilon^{(1)} \\ \varepsilon^{(2)} \\ \vdots \\ \varepsilon^{(N)} \end{bmatrix} = \begin{bmatrix} t^{(1)} \\ \vdots \\ t^{(N)} \end{bmatrix} - \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_D^{(1)} \\ \vdots & \vdots & \dots & \vdots \\ 1 & x_1^{(N)} & \dots & x_D^{(N)} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_D \end{bmatrix}$$

$$\begin{bmatrix} \varepsilon^{(1)} \\ \vdots \\ \varepsilon^{(N)} \end{bmatrix} = \begin{bmatrix} t^{(1)} \\ \vdots \\ t^{(N)} \end{bmatrix} - \begin{bmatrix} \beta_0 + \beta_1 x_1^{(1)} + \beta_2 x_2^{(1)} + \dots + \beta_D x_D^{(1)} \\ \vdots \\ \beta_0 + \beta_1 x_1^{(N)} + \beta_2 x_2^{(N)} + \dots + \beta_D x_D^{(N)} \end{bmatrix}$$

$$\rightarrow \mathcal{L}(\beta) = \frac{1}{N} \langle \vec{\varepsilon}, \vec{\varepsilon} \rangle = \frac{1}{N} \vec{\varepsilon}^T \vec{\varepsilon} = \frac{1}{N} \sum_{i=1}^N (\varepsilon^{(i)})^2$$

$$\mathcal{L}(\beta) = \frac{1}{N} (\vec{t} - \underline{\tilde{X}} \vec{\beta})^T (\vec{t} - \underline{\tilde{X}} \vec{\beta})$$

$$\rightarrow ? \text{ gradient } \frac{\partial \mathcal{L}(\beta)}{\partial \beta} = 0$$

$$l(\beta) = \frac{1}{N} \left(\cancel{E^T E} - \beta^T \tilde{X}^T E - E^T \tilde{X} \beta + \beta^T \tilde{X}^T \tilde{X} \beta \right)$$

$$\frac{\partial l}{\partial \beta}$$

$$\beta^T \nu$$

$$-2 \tilde{X}^T E$$

$$\frac{\partial}{\partial \beta}$$

$$2 \tilde{X}^T \tilde{X} \beta$$

$$\frac{\partial}{\partial \beta} (\beta_0 \nu_0 + \beta_1 \nu_1 + \dots + \beta_D \nu_D) \rightarrow [\nu_0, \nu_1, \dots, \nu_D]$$

$$\frac{\partial}{\partial \beta} (\beta^T \nu) = \vec{\nu}$$

$$\beta^T \tilde{X}^T \tilde{X} \beta \Rightarrow \frac{\partial}{\partial \beta} (\beta^T A \beta)$$

$$\begin{aligned}
 [\beta_1 \quad \beta_2] \begin{bmatrix} A_{11} & A_{12} \\ A_{12} & A_{22} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} &\rightarrow \begin{bmatrix} \beta_1 A_{11} + \beta_2 A_{12} & \beta_1 A_{12} + \beta_2 A_{22} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} \\
 &= \beta_1^2 A_{11} + \beta_1 \beta_2 A_{12} \\
 &\quad + \beta_1 \beta_2 A_{12} + \beta_2^2 A_{22}
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial}{\partial \beta_1} &\rightarrow 2\beta_1 A_{11} + 2\beta_2 A_{12} \\
 \frac{\partial}{\partial \beta_2} &\rightarrow 2\beta_2 A_{22} + 2\beta_1 A_{12}
 \end{aligned}
 = 2 \begin{bmatrix} A_{11} & A_{12} \\ A_{12} & A_{22} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix}$$

En regroupant chacun des termes on obtient

$$\frac{\partial \mathcal{L}}{\partial \beta} = -2 \underline{\tilde{X}}^T \underline{\tilde{t}} + 2 \underline{\tilde{X}}^T \underline{\tilde{X}} \underline{\tilde{\beta}} = 0$$

$\in \mathbb{R}^{D+1 \times D+1}$

$$\underline{\tilde{X}}^T \underline{\tilde{X}} \underline{\tilde{\beta}} = \underline{\tilde{X}}^T \underline{\tilde{t}} \rightarrow \underline{\tilde{\beta}} = \left(\underline{\tilde{X}}^T \underline{\tilde{X}} \right)^{-1} \underline{\tilde{X}}^T \underline{\tilde{t}} + \delta$$

Equations normales

vecteur des coefficients
du modèle

$$A = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \dots & \\ & & & \lambda_{D+1} \end{bmatrix} \rightarrow A^{-1} = \begin{bmatrix} \lambda_1^{-2} & & & \\ & \dots & & \\ & & \dots & \\ & & & \lambda_{D+1}^{-2} \end{bmatrix}$$

