

Apprentissage (supervisé) → Régression linéaire

↳ Dérivée de gradient

→ Résolution des équations des  
équations normales

→ Régularisation (normes  $l_p$ )

↳ Ridge  $(\sum_{j=2}^p \beta_j^2)$

↳ LASSO  $(\sum_{j=1}^p |\beta_j|)$

→ Sélection de caractéristiques

→ Estimateurs Statistiques

→ Maximum de vraisemblance  
(MLE)

⇔ minimisation somme des  
carrés des résidus (OLS)

→ Maximum a Posteriori (MAP)

↳ Dans le cas d'un a priori

→ A priori Gaussien ⇒ Ridge

→ A priori Laplace ⇒ LASSO

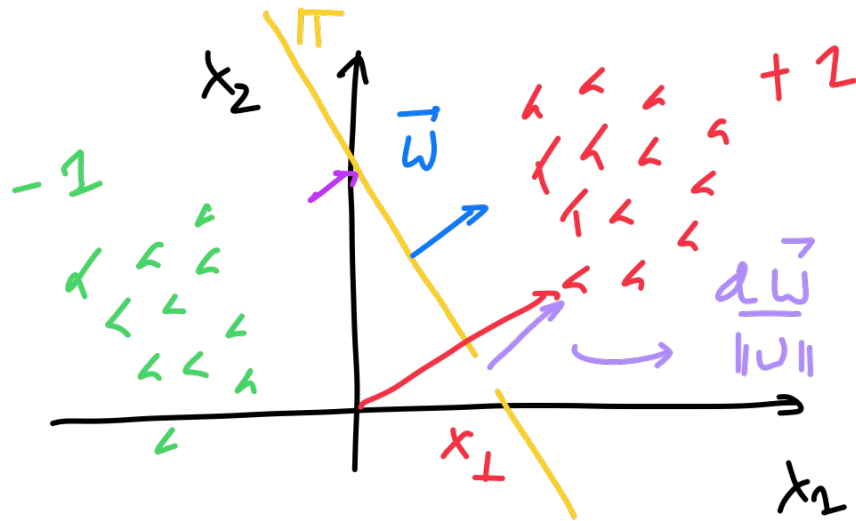
→ Equilibre Biais Variable

$$MSE = \text{biais}^2 + \text{variance}$$

# Aujourd'hui → Classification linéaire

$$\vec{w} = (\beta_1, \beta_2)$$

vecteur  
normal au  
plan



$$y(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

pour tout  $x \in \mathbb{R}^2$  on a

$$x = x_{\perp} + d \frac{\vec{w}}{\|\vec{w}\|}$$

$$x_{\perp} \in \pi \Rightarrow \beta_0 + (\beta_1, \beta_2)^T x_{\perp} = 0$$

$$\beta_0 + (\beta_1, \beta_2)^T x = \beta_0 + (\beta_1, \beta_2)^T x_{\perp} + (\beta_1, \beta_2)^T d \frac{\vec{w}}{\|\vec{w}\|}$$

$$= \begin{cases} d \frac{\|\vec{w}\|^2}{\|\vec{w}\|} > 0 & \text{si } x \text{ est situé au} \\ & \text{dessus de } \pi \\ < 0 & \text{si } x \text{ est situé en dessous} \end{cases}$$

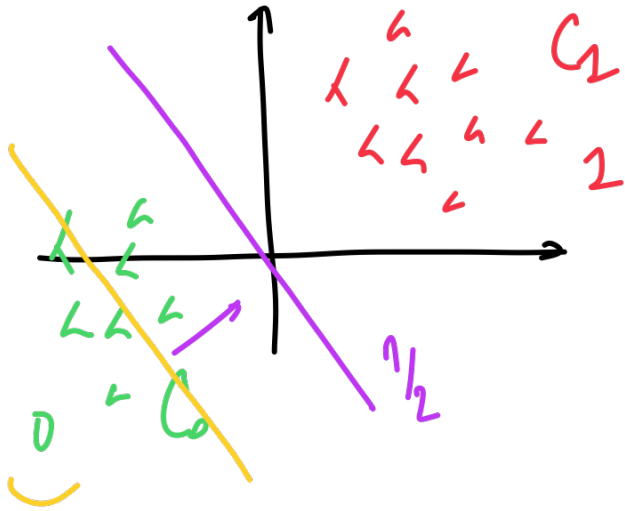
On peut donc utiliser le modèle linéaire pour séparer l'espace en deux classes  $C_1$  et  $C_2$  correspondant aux 2 demi-espaces (aux deux régions)

$$R_1 = \{x \in \mathbb{R}^2 \mid y(x) \geq 0\}$$

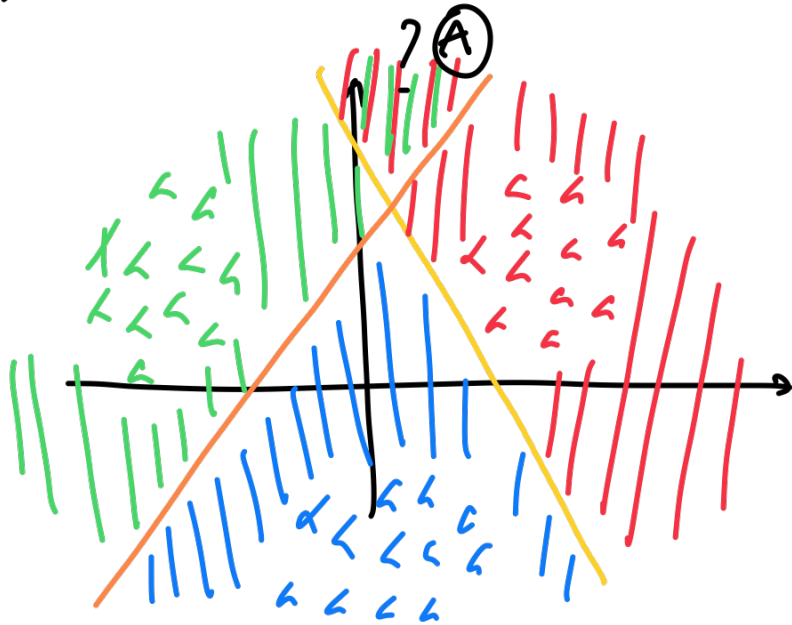
$$R_2 = \{x \in \mathbb{R}^2 \mid y(x) < 0\}$$

Comme pour la régression linéaire, on peut entraîner le modèle à l'aide de la méthode de moindres carrés OLS

$$l(\beta) = \frac{1}{N} \sum_{i=1}^N \left( \underbrace{t^{(i)}}_{-1} - \left( \beta_0 + \beta_1 \underbrace{x_1^{(i)}}_{-1} + \beta_2 \underbrace{x_2^{(i)}}_{-1} \right) \right)^2$$

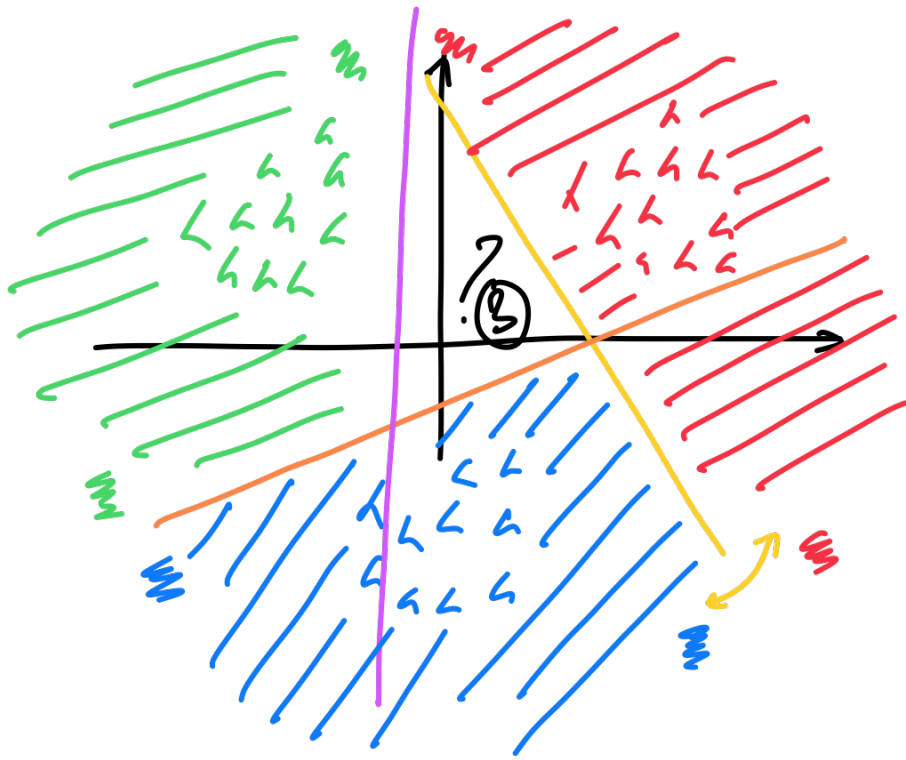


Question ? → Multi-classes ?



pour  $k$  classes, le modèle nécessite  $k-1$  discriminants

→ Ambiguïté au niveau de  $A$



Extension #2 : 1 centre 1

→ indétermination (B)

pour  $K$  classes : le module  
nécessite

$$\frac{K(K-2)}{2}$$

De façon à lever les ambiguïtés on peut considérer  $K$   
discriminants qu'on entraîne simultanément

pour chaque point  $x^{(i)}$  on introduit le vecteur cible  $\vec{t}^{(i)}$

$$\vec{t}^{(i)} = [0, 0, \underbrace{1, 0, \dots, 0}_{f: x^{(i)} \in C_3}]$$

On stocke les vecteurs cibles  $\vec{t}^{(i)}$  dans la matrice  $\underline{T}$

$$\underline{T} = \begin{bmatrix} \vec{t}^{(1)} \\ \vdots \\ \vec{t}^{(N)} \end{bmatrix} \in \mathbb{R}^{N \times K} \quad \underline{\tilde{X}} = \begin{bmatrix} 1 & \vec{x}^{(1)} \\ \vdots & \vdots \\ 1 & \vec{x}^{(N)} \end{bmatrix} \in \mathbb{R}^{N \times D+1}$$

On introduit une matrice de discriminants

$$\underline{W} = \begin{bmatrix} | & | & \dots & | \\ \vec{w}^{(1)} & \vec{w}^{(2)} & \dots & \vec{w}^{(K)} \\ | & | & & | \end{bmatrix} \in \mathbb{R}^{D+1 \times K} \quad \mathcal{L}(\beta) = \sum_i \xi_i^2$$

où  $\vec{w}^{(k)}$  représente le discriminant associé à la classe  $C_k$

$$E = (T - \underline{\tilde{X}} \underline{W}) \rightarrow \text{matrice des résidus}$$

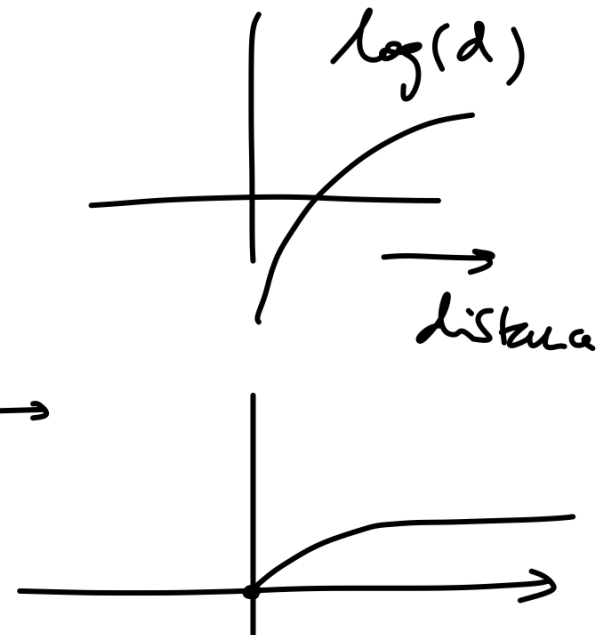
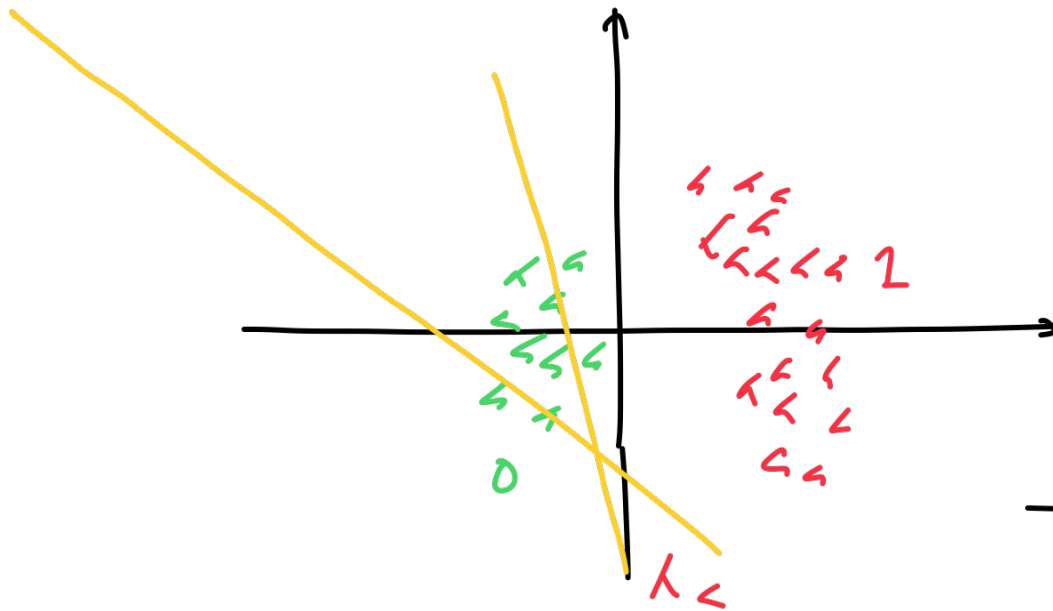
On peut définir la fonction de coût

$$L(\underline{W}) = \frac{1}{KN} \sum_{i=1}^N \sum_{j=1}^K (T - \underline{\tilde{X}} \underline{W})_{ij}^2$$

→

la minimisation de  $L(\underline{W})$  donne le modèle à  $K$  discriminants qui permet de lever les ambiguïtés des modèles 1 contre 1 et 1 contre tous.

↑

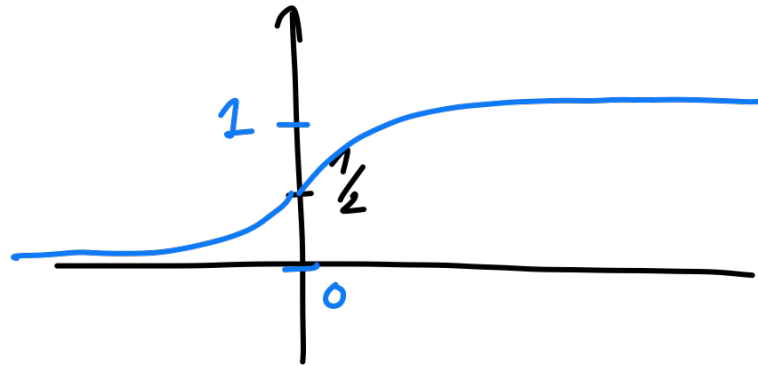




→ Idée: On utilise une fonction d'activation

$\sigma(x)$  et on va définir le modèle comme  $\sigma(\beta_0 + \beta_1 x_1 + \beta_2 x_2)$

choix #1  $\sigma(x) = \frac{1}{1 + e^{-x}}$



→ Modèle correspondant est appelé Régression logistique

$$y(x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \beta_2 x_2)}}$$

$$\left. \begin{aligned} p(t^{(i)}=1 | x^{(i)}, \beta) &= \sigma(\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}) \\ p(t^{(i)}=0 | x^{(i)}, \beta) &= 1 - \sigma(\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}) \end{aligned} \right\}$$

→ Maximum de vraisemblance

$$p(t^{(i)}=t | x^{(i)}, \beta) = \sigma(\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)})^t \times (1 - \sigma(\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^{1-t}$$

$$p(\{t^{(i)}\}_{i=1}^N | \{x^{(i)}\}_{i=1}^N, \vec{\beta}) = \prod_{i=1}^N \sigma(\beta_0 + \beta_1 x_1 + \dots + \beta_D x_D)^t \times (1 - \sigma(\beta_0 + \beta_1 x_1 + \dots + \beta_D x_D))^{1-t}$$

Or supprime l'indépendance

On considère la fonction de log vraisemblance

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

Entropie binaire croisée

$$h_{\beta} = \beta_0 + \beta_1 x_1 + \dots + \beta_D x_D$$

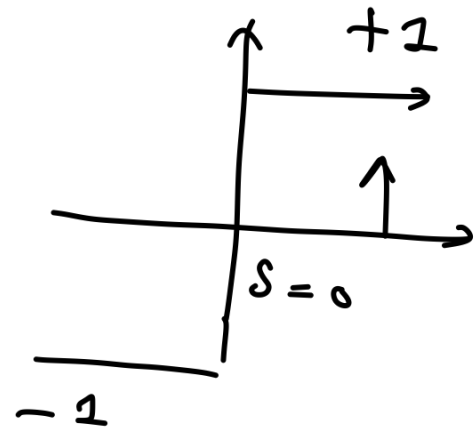
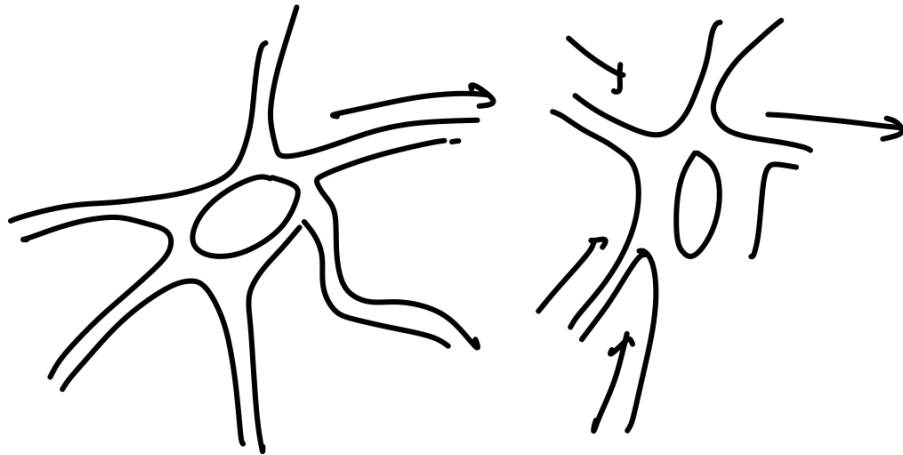
$$\rightarrow t \log(\sigma(\underline{z})) + (1-t) \log(1 - \sigma(\underline{z}))$$

$$\sigma'(x) = \frac{d}{dx} \left( \frac{1}{1+e^{-x}} \right) = \frac{e^{-x}}{(1+e^{-x})^2} = \frac{\sigma(x)}{1+e^{-x}} \left( 1 - \frac{1}{1+e^{-x}} \right)$$

$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \beta_j} &= \frac{\partial \mathcal{L}}{\partial h_j} \cdot \frac{\partial h_j}{\partial \beta_j} = \left( t \frac{\sigma'}{\sigma} + (1-t) \frac{(-\sigma')}{1-\sigma} \right) \cdot \tilde{x}_j \\ &= (t(1-\sigma) - (1-t) \cdot \sigma) \tilde{x}_j \\ &= (t - \sigma(\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)})) \cdot \tilde{x}_j^{(i)} \end{aligned}$$

Option #2



Heaviside  $H(x) = \begin{cases} +1 & x \geq 0 \\ -1 & x < 0 \end{cases}$

→ Algorithme du perceptron

$$\min - \sum_{i \in \text{misclassified}} \overbrace{t^{(i)}}^{>0} \underbrace{\left( \beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)} \right)}^{<0}$$