



IFSBM Module 11

**GUSTAVE/
ROUSSY**
CANCER CAMPUS
GRAND PARIS

université
PARIS-SACLAY

IFSBM
INSTITUT DE FORMATION
SUPÉRIEURE BIOMÉDICALE

Méthodes de classification supervisée

Yoann Pradat 1

Sommaire

1. Données et modélisation
 1. Le jeu de données
 2. L'estimation (model fitting)
2. Quelques modèles de classification supervisée
 1. Régression logistique binaire
 2. Régression logistique multinomiale
 3. Régression linéaire
 4. Régression pénalisée
 5. Réseau de neurones monocouche
 6. Réseau de neurones multicouches

1.1 Le jeu de données

Notations

$X : \Omega \mapsto \mathcal{X}$: variable aléatoire (scalaire ou vectorielle)
 $x \in \mathcal{X}$: une observation de la variable aléatoire

Exemples:

1. On lance un dé, $\mathcal{X} = \{1, 2, \dots, 6\}$

1.1 Le jeu de données

Notations

$X : \Omega \mapsto \mathcal{X}$: variable aléatoire (scalaire ou vectorielle)
 $x \in \mathcal{X}$: une observation de la variable aléatoire

Exemples:

1. On lance un dé, $\mathcal{X} = \{1, 2, \dots, 6\}$
2. On mesure la longueur et la largeur de fleurs $\mathcal{X} = \mathbb{R}^+ \times \mathbb{R}^+$

1.1 Le jeu de données

Notations

$X : \Omega \mapsto \mathcal{X}$: variable aléatoire (scalaire ou vectorielle)
 $x \in \mathcal{X}$: une observation de la variable aléatoire

Exemples:

1. On lance un dé, $\mathcal{X} = \{1, 2, \dots, 6\}$
2. On mesure la longueur et la largeur de fleurs $\mathcal{X} = \mathbb{R}^+ \times \mathbb{R}^+$
3. On mesure le niveau d'expression de gènes relatif à 1M (K gènes)
 $\mathcal{X} = [0, 1M]^K = [0, 1M] \times [0, 1M] \times \dots \times [0, 1M]$

1.1 Le jeu de données

Notations

$X : \Omega \mapsto \mathcal{X}$: variable aléatoire (scalaire ou vectorielle)
 $x \in \mathcal{X}$: une observation de la variable aléatoire

Exemples:

1. On lance un dé, $\mathcal{X} = \{1, 2, \dots, 6\}$

2. On mesure la longueur et la largeur de fleurs $\mathcal{X} = \mathbb{R}^+ \times \mathbb{R}^+$

3. On mesure le niveau d'expression de gènes relatif à 1M (K gènes)

$$\mathcal{X} = [0, 1M]^K = [0, 1M] \times [0, 1M] \times \dots \times [0, 1M]$$

$n \in \mathbb{N}^*$: nombre d'observations (=individus, échantillons)

$p \in \mathbb{N}^*$: nombre de variables (=covariables, prédicteurs, features)

$x_{1:n} = (x_1, \dots, x_n)$: ensemble d'observations (=échantillon, =dataset)

1.1 Le jeu de données

Exemple:

1. On réalise un essai clinique de réponse à un traitement sur 100 patients

$X =$ (Age, Poids, Nb Tx antérieurs,
Exp gene 1, Exp gene 2, Mutation gene 1,
Temps avant rechute, Meilleure réponse)

1.1 Le jeu de données

Exemple:

1. On réalise un essai clinique de réponse à un traitement sur 100 patients

$X =$ (Age, Poids, Nb Tx antérieurs,
Exp gene 1, Exp gene 2, Mutation gene 1,
Temps avant rechute, Meilleure réponse)

$x_1 =$ (45.2, 78.2, 3, 1032, 258, 1, 85, PR)

$x_2 =$ (81, 63, 6, 589, 903, 0, 390, SD)

...

1.1 Le jeu de données

Exemple:

1. On réalise un essai clinique de réponse à un traitement sur 100 patients

$X =$ (Age, Poids, Nb Tx antérieurs,
Exp gene 1, Exp gene 2, Mutation gene 1,
Temps avant rechute, Meilleure réponse)

$x_1 =$ (45.2, 78.2, 3, 1032, 258, 1, 85, PR)

$x_2 =$ (81, 63, 6, 589, 903, 0, 390, SD)

...

} $(x_1, x_2, \dots, x_{100})$
= jeu de données

Objectifs

1. **Proposer un modèle** mathématique pour modéliser une variable en fonction d'autres. Exemple Meilleure réponse vs (Exp gene 1, Exp gene 2, Age)
Modèle = Régression logistique multinomiale

1.1 Le jeu de données

Exemple:

1. On réalise un essai clinique de réponse à un traitement sur 100 patients

$X =$ (Age, Poids, Nb Tx antérieurs,
Exp gene 1, Exp gene 2, Mutation gene 1,
Temps avant rechute, Meilleure réponse)

$x_1 =$ (45.2, 78.2, 3, 1032, 258, 1, 85, PR)

$x_2 =$ (81, 63, 6, 589, 903, 0, 390, SD)

...

$(x_1, x_2, \dots, x_{100})$
= jeu de données

Objectifs

1. **Proposer un modèle** mathématique pour modéliser une variable en fonction d'autres. Exemple Meilleure réponse vs (Exp gene 1, Exp gene 2, Age)
Modèle = Régression logistique multinomiale
2. **Estimer les paramètres** du modèle à partir de $(x_1, x_2, \dots, x_{100})$

1.2 L'estimation

Exemple:

1. On veut predire le volume de la tumeur à partir de l'expression de 5000 gènes

-> (G^1, \dots, G^{5000}) profil d'expression et $V = \text{volume}$

$x_{1:n} = (g^1, \dots, g^{5000}, v)_{1:n}$ observations de $X = (G^1, \dots, G^{5000}, V)$

1.2 L'estimation

Exemple:

1. On veut prédire le volume de la tumeur à partir de l'expression de 5000 gènes

-> (G^1, \dots, G^{5000}) profil d'expression et $V = \text{volume}$

$x_{1:n} = (g^1, \dots, g^{5000}, v)_{1:n}$ observations de $X = (G^1, \dots, G^{5000}, V)$

-> on prend comme modèle un **modèle de regression linéaire**

$$\begin{aligned} V &= f_{\theta}(G^1, G^{5000}) \\ &= \theta_1 G^1 + \dots + \theta_{5000} G^{5000} \end{aligned}$$

Schématique, pas
mathématiquement exact

1.2 L'estimation

Exemple:

1. On veut prédire le volume de la tumeur à partir de l'expression de 5000 gènes

-> (G^1, \dots, G^{5000}) profil d'expression et $V = \text{volume}$

$x_{1:n} = (g^1, \dots, g^{5000}, v)_{1:n}$ observations de $X = (G^1, \dots, G^{5000}, V)$

-> on prend comme modèle un **modèle de regression linéaire**

$$V = f_{\theta}(G^1, G^{5000})$$
$$= \theta_1 G^1 + \dots + \theta_{5000} G^{5000}$$

Schématique, pas
mathématiquement exact

Fonction objectif (=coût) Ca veut dire quoi "coller au mieux aux données" ?

Mesure de distance entre l'observation v_i et la prédiction $\hat{v}_i = f_{\theta}(g_i^1, \dots, g_i^{5000})$?

1.2 L'estimation

Exemple:

1. On veut prédire le volume de la tumeur à partir de l'expression de 5000 gènes

-> (G^1, \dots, G^{5000}) profil d'expression et $V = \text{volume}$

$x_{1:n} = (g^1, \dots, g^{5000}, v)_{1:n}$ observations de $X = (G^1, \dots, G^{5000}, V)$

-> on prend comme modèle un **modèle de regression linéaire**

$$V = f_{\theta}(G^1, G^{5000})$$

Schématique, pas
mathématiquement exact

$$= \theta_1 G^1 + \dots + \theta_{5000} G^{5000}$$

Fonction objectif (=coût) Ca veut dire quoi "coller au mieux aux données" ?

Mesure de distance entre l'observation v_i et la prédiction $\hat{v}_i = f_{\theta}(g_i^1, \dots, g_i^{5000})$?

$$d(v_i, \hat{v}_i) = (v_i - \hat{v}_i) ?$$

1.2 L'estimation

Exemple:

1. On veut prédire le volume de la tumeur à partir de l'expression de 5000 gènes

-> (G^1, \dots, G^{5000}) profil d'expression et $V = \text{volume}$

$x_{1:n} = (g^1, \dots, g^{5000}, v)_{1:n}$ observations de $X = (G^1, \dots, G^{5000}, V)$

-> on prend comme modèle un **modèle de regression linéaire**

$$V = f_{\theta}(G^1, G^{5000})$$
$$= \theta_1 G^1 + \dots + \theta_{5000} G^{5000}$$

Schématique, pas
mathématiquement exact

Fonction objectif (=coût) Ca veut dire quoi "coller au mieux aux données" ?

Mesure de distance entre l'observation v_i et la prédiction $\hat{v}_i = f_{\theta}(g_i^1, \dots, g_i^{5000})$?

$$d(v_i, \hat{v}_i) = (v_i - \hat{v}_i) ? \quad (v_i - \hat{v}_i)^2 ?$$

1.2 L'estimation

Exemple:

1. On veut prédire le volume de la tumeur à partir de l'expression de 5000 gènes

-> (G^1, \dots, G^{5000}) profil d'expression et $V = \text{volume}$

$x_{1:n} = (g^1, \dots, g^{5000}, v)_{1:n}$ observations de $X = (G^1, \dots, G^{5000}, V)$

-> on prend comme modèle un **modèle de regression linéaire**

$$V = f_{\theta}(G^1, G^{5000})$$

Schématique, pas
mathématiquement exact

$$= \theta_1 G^1 + \dots + \theta_{5000} G^{5000}$$

Fonction objectif (=coût) Ca veut dire quoi "coller au mieux aux données" ?

Mesure de distance entre l'observation v_i et la prédiction $\hat{v}_i = f_{\theta}(g_i^1, \dots, g_i^{5000})$?

$$d(v_i, \hat{v}_i) = (v_i - \hat{v}_i) ? \quad (v_i - \hat{v}_i)^2 ? \quad (v_i - \hat{v}_i)^4 ?$$

1.2 L'estimation

Exemple:

1. On veut prédire le volume de la tumeur à partir de l'expression de 5000 gènes

-> (G^1, \dots, G^{5000}) profil d'expression et $V = \text{volume}$

$x_{1:n} = (g^1, \dots, g^{5000}, v)_{1:n}$ observations de $X = (G^1, \dots, G^{5000}, V)$

-> on prend comme modèle un **modèle de regression linéaire**

$$V = f_{\theta}(G^1, G^{5000})$$

Schématique, pas
mathématiquement exact

$$= \theta_1 G^1 + \dots + \theta_{5000} G^{5000}$$

Fonction objectif (=coût) Ca veut dire quoi "coller au mieux aux données" ?

Mesure de distance entre l'observation v_i et la prédiction $\hat{v}_i = f_{\theta}(g_i^1, \dots, g_i^{5000})$?

$$d(v_i, \hat{v}_i) = (v_i - \hat{v}_i) ? \quad (v_i - \hat{v}_i)^2 ? \quad (v_i - \hat{v}_i)^4 ? \quad e^{|v_i - \hat{v}_i|} - 1 ?$$

1.2 L'estimation

Exemple:

1. On veut prédire le volume de la tumeur à partir de l'expression de 5000 gènes

-> (G^1, \dots, G^{5000}) profil d'expression et $V = \text{volume}$

$x_{1:n} = (g^1, \dots, g^{5000}, v)_{1:n}$ observations de $X = (G^1, \dots, G^{5000}, V)$

-> on prend comme modèle un **modèle de regression linéaire**

$$V = f_{\theta}(G^1, G^{5000})$$
$$= \theta_1 G^1 + \dots + \theta_{5000} G^{5000}$$

Schématique, pas
mathématiquement exact

Fonction objectif (=coût) Ca veut dire quoi "coller au mieux aux données" ?

Mesure de distance entre l'observation v_i et la prédiction $\hat{v}_i = f_{\theta}(g_i^1, \dots, g_i^{5000})$?

$$d(v_i, \hat{v}_i) = (v_i - \hat{v}_i) ? \quad (v_i - \hat{v}_i)^2 ? \quad (v_i - \hat{v}_i)^4 ? \quad e^{|v_i - \hat{v}_i|} - 1 ?$$

Estimation $\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmin}} \sum_{i=1}^n d(v_i, \hat{v}_i)$

1.2 L'estimation

Quels modèles ? Le modèle peut avoir une interprétation probabiliste ou non.

-> si interprétation probabiliste

fonction objectif = - **maximum de vraisemblance**

Exemples: régression linéaire, régression logistique, classification bayésienne naïve

1.2 L'estimation

Quels modèles ? Le modèle peut avoir une interprétation probabiliste ou non.

-> si interprétation probabiliste

fonction objectif = - **maximum de vraisemblance**

Exemples: régression linéaire, régression logistique, classification bayésienne naïve

-> si pas d'interprétation probabiliste

fonction objectif = **a la main**

Exemples: machine à vecteur de support (SVM), forêt aléatoire, réseaux de neurones

1.2 L'estimation

Quels modèles ? Le modèle peut avoir une interprétation probabiliste ou non.

-> si interprétation probabiliste (=modèle statistique)

fonction objectif = - **maximum de vraisemblance**

Exemples: régression linéaire, régression logistique, classification bayésienne naïve

-> si pas d'interprétation probabiliste (=modèle statistique?)

fonction objectif = **a la main**

Exemples: machine à vecteur de support (SVM), forêt aléatoire, reseaux de neurones

Les reseaux de neurones sont-ils des modèles statistiques ?

<https://ai.stackexchange.com/questions/10289/are-neural-networks-statistical-models/18580#18580>

Cf aussi sur la definition des modèles statistiques <https://www.stat.uchicago.edu/~pmcc/pubs/AOS023.pdf>

1.2 L'estimation

-> si interpretation probabiliste (=modèle statistique)

Modèle statistique = ensemble de lois (=mesures) de probabilité \mathbb{P} sur l'espace des observations \mathcal{X} . Si proba paramétriques, alors on parle de modèle paramétrique. Enfin si les proba ont une densité p alors le modèle s'écrit

$$\mathcal{M}_{\Theta} = \{p_{\theta} | \theta \in \Theta\}$$

1.2 L'estimation

-> si interprétation probabiliste (=modèle statistique)

Modèle statistique = ensemble de lois (=mesures) de probabilité \mathbb{P} sur l'espace des observations \mathcal{X} . Si proba paramétriques, alors on parle de modèle paramétrique. Enfin si les proba ont une densité p alors le modèle s'écrit

$$\mathcal{M}_{\Theta} = \{p_{\theta} | \theta \in \Theta\}$$

Quand on modélise, on fait l'hypothèse qu'il existe θ^* tel que $\mathbb{P}_X = \mathbb{P}_{\theta^*}$

Modéliser = calculer $\hat{\theta}$ en "espérant" que $\hat{\theta} \approx \theta^*$

1.2 L'estimation

Modèle statistique $\mathcal{M}_{\Theta} = \{p_{\theta} | \theta \in \Theta\}$ Il existe θ^* tel que $\mathbb{P}_X = \mathbb{P}_{\theta^*}$

Pour le modèle candidat \mathbb{P}_{θ}

$p_{\theta}(x_i)$ = vraisemblance échantillon i
 $\prod_{i=1}^n p_{\theta}(x_i)$ = vraisemblance jeu de données

1.2 L'estimation

Modèle statistique $\mathcal{M}_\Theta = \{p_\theta | \theta \in \Theta\}$ Il existe θ^* tel que $\mathbb{P}_X = \mathbb{P}_{\theta^*}$

Pour le modèle candidat \mathbb{P}_θ

$p_\theta(x_i)$ = vraisemblance échantillon i
 $\prod_{i=1}^n p_\theta(x_i)$ = vraisemblance jeu de données

$$L(\theta) = - \prod_{i=1}^n p_\theta(x_i)$$

Fonction de coût à minimiser

1.2 L'estimation

Modèle statistique $\mathcal{M}_\Theta = \{p_\theta | \theta \in \Theta\}$ Il existe θ^* tel que $\mathbb{P}_X = \mathbb{P}_{\theta^*}$

Pour le modèle candidat \mathbb{P}_θ

$p_\theta(x_i)$ = vraisemblance échantillon i
 $\prod_{i=1}^n p_\theta(x_i)$ = vraisemblance jeu de données

$$L(\theta) = - \prod_{i=1}^n p_\theta(x_i) \text{ équivalent à } \ell(\theta) = - \sum_{i=1}^n \log p_\theta(x_i)$$

Fonction de coût à minimiser = - log de la vraisemblance

Sommaire

1. Données et modélisation

1. Le jeu de données
2. L'estimation (model fitting)

2. Quelques modèles de classification supervisée

1. Régression logistique binaire
2. Régression logistique multinomiale
3. Régression linéaire
4. Régression pénalisée
5. Réseau de neurones monocouche
6. Réseau de neurones multicouches

2.1 Régression logistique binaire

Modèle statistique de regression Prédire Y à partir de X

Jeu de données $x_{1:n}, y_{1:n}$

Un modèle de regression est une famille de lois de probabilités *conditionnelles*

$$\mathcal{M}_{\Theta} = \{\mathbb{P}_{Y|X=x}^{\theta} | \theta \in \Theta, x \in \mathcal{X}\}$$

2.1 Régression logistique binaire

Modèle statistique de regression Prédire Y à partir de X

Jeu de données $x_{1:n}, y_{1:n}$

Un modèle de regression est une famille de lois de probabilités *conditionnelles*

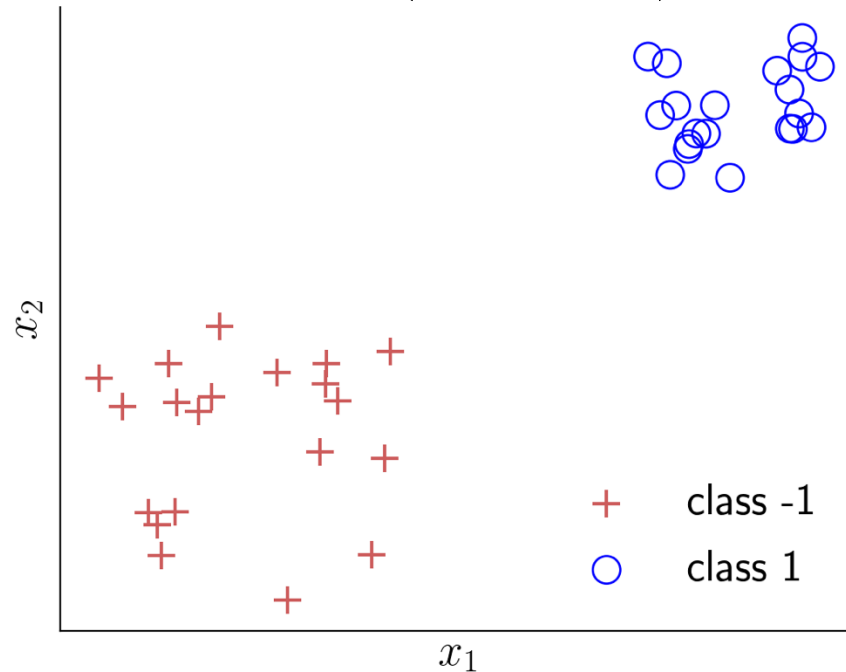
$$\mathcal{M}_{\Theta} = \{\mathbb{P}_{Y|X=x}^{\theta} | \theta \in \Theta, x \in \mathcal{X}\}$$

Regression logistique $Y \in \{0, 1\}, \quad X \in \mathbb{R}^p$

$$\mathbb{P}_{Y|X=x}^{\theta} = \text{Binomial}(\sigma(x^{\top} \theta)) \quad \sigma(u) = \frac{1}{1 + e^{-u}}$$

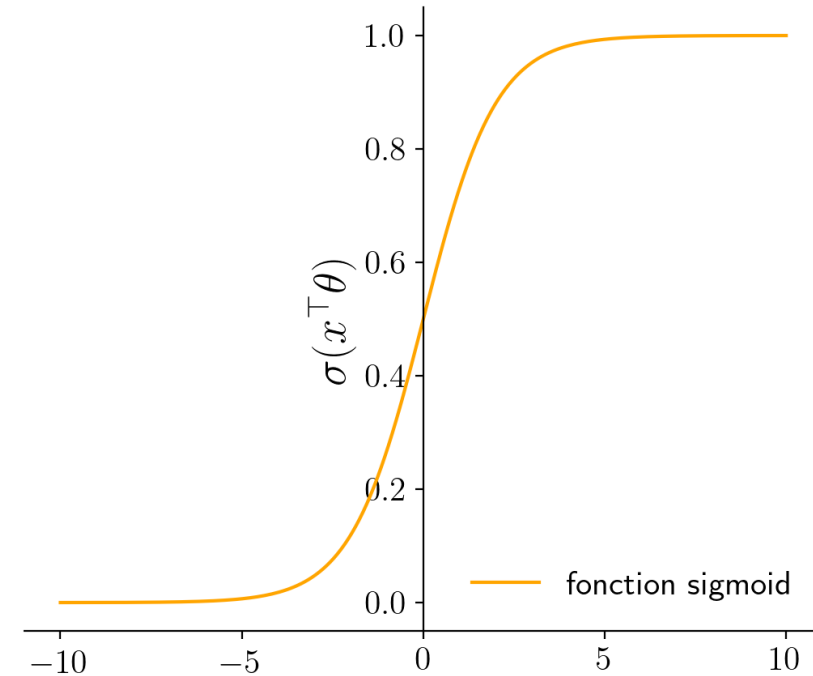
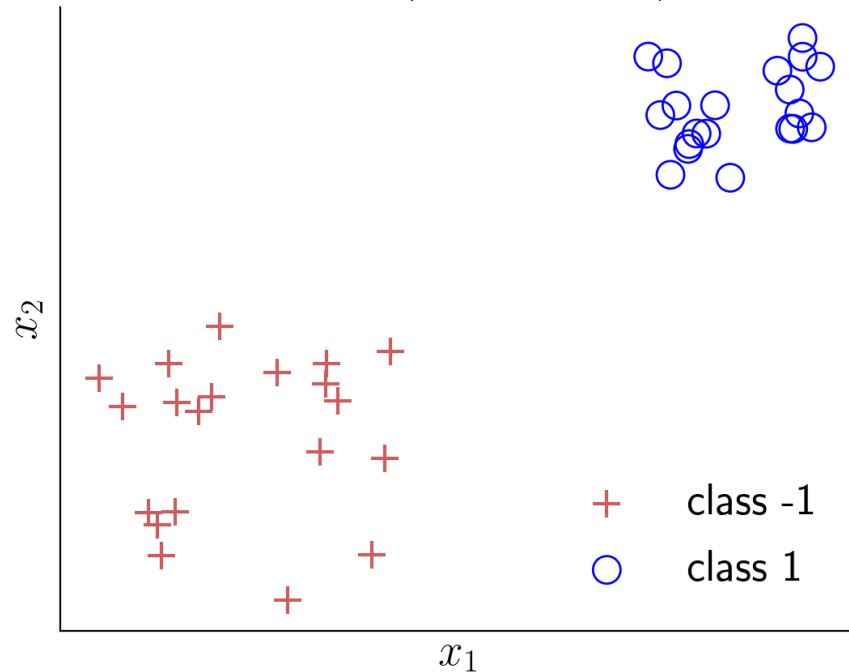
2.1 Régression logistique binaire

Regression logistique $Y \in \{0, 1\}$, $X \in \mathbb{R}^p$ $\mathbb{P}_{Y|X=x}^\theta = \text{Binomial}(\sigma(x^\top \theta))$
Exemple: $p=2$, $X = (X^1, X^2)$



2.1 Régression logistique binaire

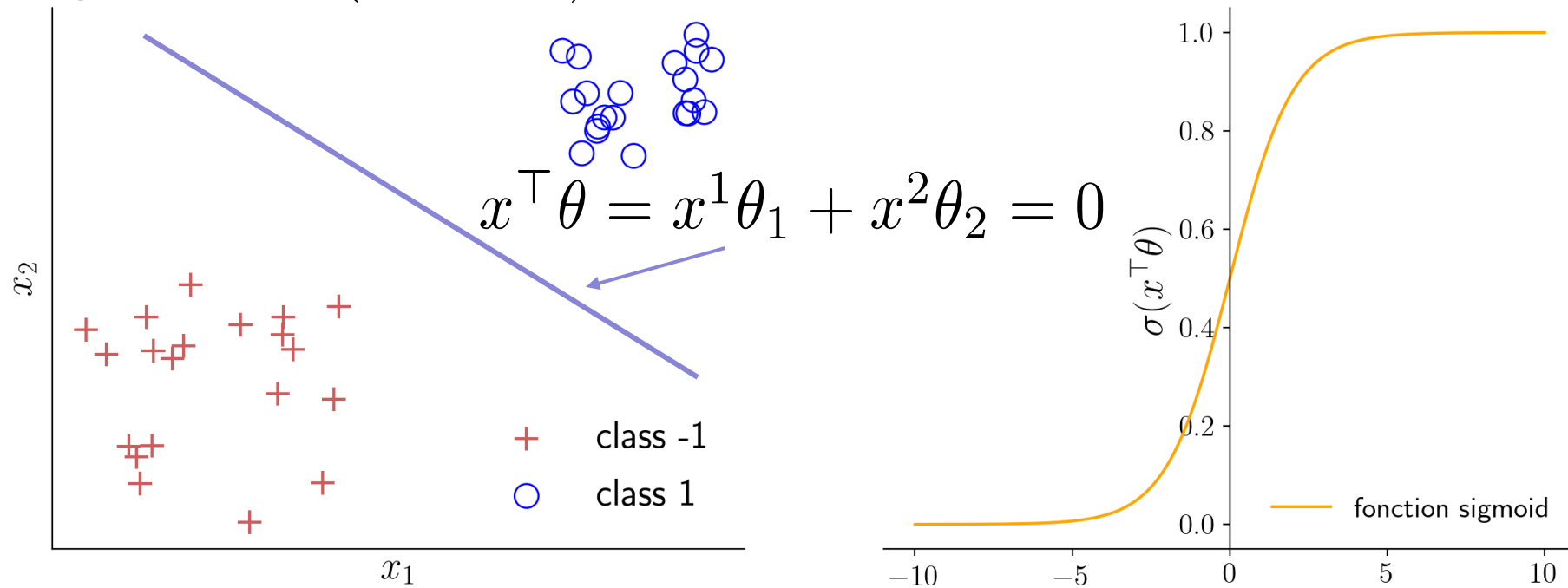
Regression logistique $Y \in \{0, 1\}$, $X \in \mathbb{R}^p$ $\mathbb{P}_{Y|X=x}^\theta = \text{Binomial}(\sigma(x^\top \theta))$
Exemple: $p=2$, $X = (X^1, X^2)$



Pour $X = x$, prédiction =
$$\begin{cases} 1 & \text{si } \sigma(x^\top \theta) > 0.5, \text{ i.e. } x^\top \theta > 0 \\ 0 & \text{si } x^\top \theta < 0 \end{cases}$$

2.1 Régression logistique binaire

Regression logistique $Y \in \{0, 1\}$, $X \in \mathbb{R}^p$ $\mathbb{P}_{Y|X=x}^\theta = \text{Binomial}(\sigma(x^\top \theta))$
Exemple: $p=2$, $X = (X^1, X^2)$



Pour $X = x$, prédiction =
$$\begin{cases} 1 & \text{si } \sigma(x^\top \theta) > 0.5, \text{ i.e. } x^\top \theta > 0 \\ 0 & \text{si } x^\top \theta < 0 \end{cases}$$

2.1 Régression logistique binaire

Regression logistique $Y \in \{0, 1\}, \quad X \in \mathbb{R}^p$

$$\mathbb{P}_{Y|X=x}^{\theta} = \text{Binomial}(\sigma(x^{\top} \theta)) \quad \sigma(u) = \frac{1}{1 + e^{-u}}$$

Comment s'écrit la vraisemblance de $y_i|x_i$?

$$p_{\theta}(\cdot|x_i) = \begin{cases} 1 & \text{avec probabilité } \sigma(x_i^{\top} \theta) \\ 0 & \text{avec probabilité } 1 - \sigma(x_i^{\top} \theta) \end{cases}$$

2.1 Régression logistique binaire

Regression logistique $Y \in \{0, 1\}, \quad X \in \mathbb{R}^p$

$$\mathbb{P}_{Y|X=x}^{\theta} = \text{Binomial}(\sigma(x^{\top} \theta)) \quad \sigma(u) = \frac{1}{1 + e^{-u}}$$

Comment s'écrit la vraisemblance de $y_i|x_i$?

$$p_{\theta}(\cdot|x_i) = \begin{cases} 1 & \text{avec probabilité } \sigma(x_i^{\top} \theta) \\ 0 & \text{avec probabilité } 1 - \sigma(x_i^{\top} \theta) \end{cases}$$

$$p_{\theta}(y_i|x_i) = \sigma(\theta^{\top} x_i)^{y_i} (1 - \sigma(\theta^{\top} x_i))^{1-y_i}$$

2.1 Régression logistique binaire

Regression logistique $Y \in \{0, 1\}, \quad X \in \mathbb{R}^p$

$$\mathbb{P}_{Y|X=x}^{\theta} = \text{Binomial}(\sigma(x^{\top} \theta)) \quad \sigma(u) = \frac{1}{1 + e^{-u}}$$

Comment s'écrit la vraisemblance de $y_i|x_i$?

$$p_{\theta}(\cdot|x_i) = \begin{cases} 1 & \text{avec probabilité } \sigma(x_i^{\top} \theta) \\ 0 & \text{avec probabilité } 1 - \sigma(x_i^{\top} \theta) \end{cases}$$

$$p_{\theta}(y_i|x_i) = \sigma(\theta^{\top} x_i)^{y_i} (1 - \sigma(\theta^{\top} x_i))^{1-y_i}$$

Fonction de coût

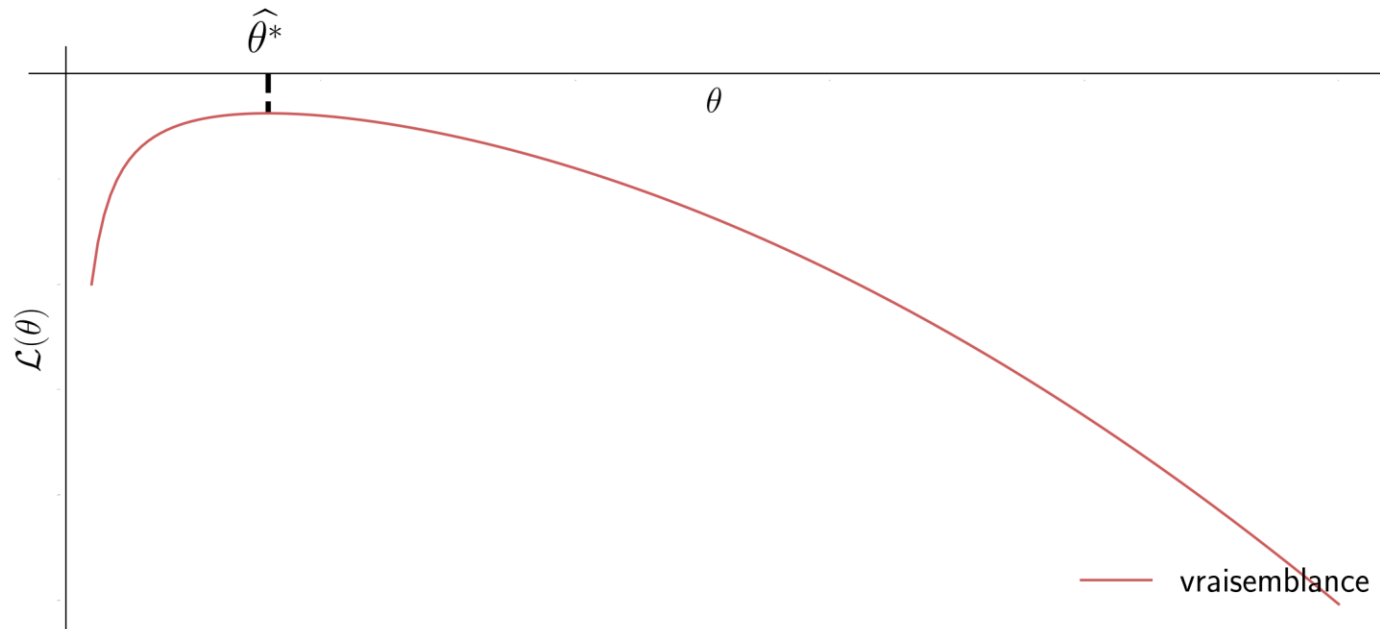
$$\ell(\theta; x_{1:n}, y_{1:n}) = - \sum_{i=1}^n y_i \log \sigma(\theta^{\top} x_i) + (1 - y_i) \log(1 - \sigma(\theta^{\top} x_i))$$

2.1 Régression logistique binaire

Regression logistique

Fonction de coût

$$\ell(\theta; x_{1:n}, y_{1:n}) = - \sum_{i=1}^n y_i \log \sigma(\theta^\top x_i) + (1 - y_i) \log(1 - \sigma(\theta^\top x_i))$$



Concave !

2.1 Régression logistique binaire

Régression logistique

Fonction de coût

$$\ell(\theta; x_{1:n}, y_{1:n}) = - \sum_{i=1}^n y_i \log \sigma(\theta^\top x_i) + (1 - y_i) \log(1 - \sigma(\theta^\top x_i))$$

Minimisation par descente de gradient

$$\ell(\theta^{t+1}) = \ell(\theta^t) + (\theta^{t+1} - \theta^t) \nabla \ell(\theta^t) + o(\theta^{t+1} - \theta^t)$$

(dvp de Taylor ordre 1)

2.1 Régression logistique binaire

Régression logistique

Fonction de coût

$$\ell(\theta; x_{1:n}, y_{1:n}) = - \sum_{i=1}^n y_i \log \sigma(\theta^\top x_i) + (1 - y_i) \log(1 - \sigma(\theta^\top x_i))$$

Minimisation par descente de gradient

$$\ell(\theta^{t+1}) = \ell(\theta^t) + (\theta^{t+1} - \theta^t) \nabla \ell(\theta^t) + o(\theta^{t+1} - \theta^t)$$

(dvp de Taylor ordre 1)

Idée. Choisir θ^{t+1} tel que $(\theta^{t+1} - \theta^t) = -\nabla \ell(\theta^t)$

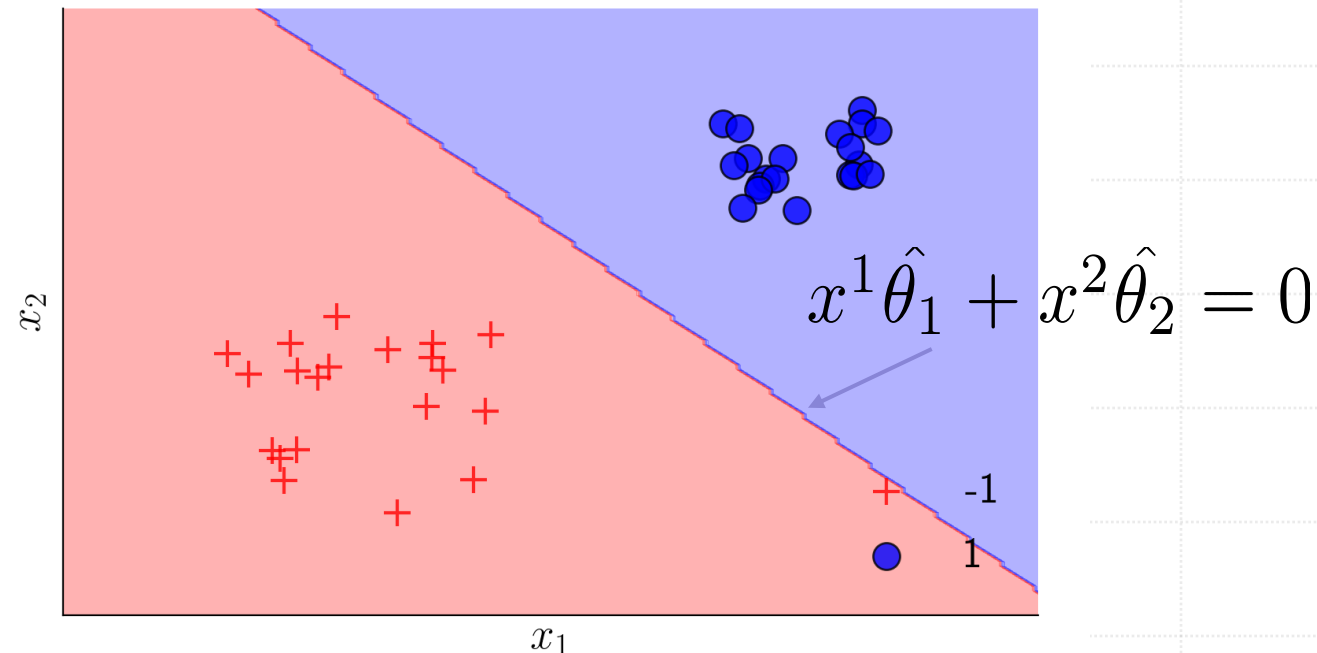
2.1 Régression logistique binaire

Regression logistique

Fonction de coût

$$\ell(\theta; x_{1:n}, y_{1:n}) = - \sum_{i=1}^n y_i \log \sigma(\theta^\top x_i) + (1 - y_i) \log(1 - \sigma(\theta^\top x_i))$$

Minimisation par descente de gradient



2.2 Régression logistique multinomiale

Regression logistique multinomiale $Y \in \{1, 2, \dots, K\}, \quad X \in \mathbb{R}^p$

$$\mathbb{P}_{Y|X=x}^{\Theta} = \text{Multinomial}(\sigma(\Theta^{\top} x))$$

$$\sigma(z) = \frac{1}{\sum_{k=1}^K e^{z_k}} \begin{bmatrix} e^{z_1} \\ \vdots \\ e^{z_K} \end{bmatrix}$$

$$\Theta = \begin{bmatrix} \theta_{1,1} & \cdots & \theta_{1,p} \\ \vdots & \ddots & \vdots \\ \theta_{K,1} & \cdots & \theta_{K,p} \end{bmatrix}$$

2.2 Régression logistique multinomiale

Regression logistique multinomiale $Y \in \{1, 2, \dots, K\}, \quad X \in \mathbb{R}^p$

$$\mathbb{P}_{Y|X=x}^{\Theta} = \text{Multinomial}(\sigma(\Theta^{\top} x))$$
$$\sigma(z) = \frac{1}{\sum_{k=1}^K e^{z_k}} \begin{bmatrix} e^{z_1} \\ \vdots \\ e^{z_K} \end{bmatrix}$$
$$\Theta = \begin{bmatrix} \theta_{1,1} & \cdots & \theta_{1,p} \\ \vdots & \ddots & \vdots \\ \theta_{K,1} & \cdots & \theta_{K,p} \end{bmatrix}$$
$$\begin{bmatrix} \mathbb{P}_{Y|X=x}^{\Theta}(Y=1) \\ \vdots \\ \mathbb{P}_{Y|X=x}^{\Theta}(Y=K) \end{bmatrix} = \frac{1}{\sum_{k=1}^K e^{(\Theta^{\top} x)_k}} \begin{bmatrix} e^{(\Theta^{\top} x)_1} \\ \vdots \\ e^{(\Theta^{\top} x)_K} \end{bmatrix}$$

2.2 Régression logistique multinomiale

Regression logistique multinomiale $Y \in \{1, 2, \dots, K\}, \quad X \in \mathbb{R}^p$

Comment s'écrit la vraisemblance de $y_i | x_i$?

$p_{\Theta}(\cdot | x_i) = k$ avec probabilité $\sigma(\Theta^{\top} x)_k$

$$p_{\Theta}(y_i | x_i) = \prod_{k=1}^K \sigma(\Theta^{\top} x_i)_k^{\mathbb{1}_{y_i=k}}$$

2.2 Régression logistique multinomiale

Regression logistique multinomiale $Y \in \{1, 2, \dots, K\}, \quad X \in \mathbb{R}^p$

Comment s'écrit la vraisemblance de $y_i | x_i$?

$p_{\Theta}(\cdot | x_i) = k$ avec probabilité $\sigma(\Theta^{\top} x)_k$

$$p_{\Theta}(y_i | x_i) = \prod_{k=1}^K \sigma(\Theta^{\top} x_i)_k^{\mathbb{1}_{y_i=k}}$$

Fonction de coût

$$\ell(\Theta; x_{1:n}, y_{1:n}) = - \sum_{i=1}^n \sum_{k=1}^K \mathbb{1}_{y_i=k} \log \sigma(\Theta^{\top} x_i)_k$$

2.2 Régression logistique multinomiale

Estimation par descente de gradient $(\Theta^{t+1} - \Theta^t) = -\nabla \ell(\Theta^t)$

2.2 Régression logistique multinomiale

Estimation par descente de gradient $(\Theta^{t+1} - \Theta^t) = -\nabla \ell(\Theta^t)$

NOTE 1: modèle surparamétré, car

$$\text{si } \Theta \leftarrow \Theta - \psi = \begin{bmatrix} \Theta_{1,:} - \psi \\ \vdots \\ \Theta_{K,:} - \psi \end{bmatrix}$$

2.2 Régression logistique multinomiale

Estimation par descente de gradient $(\Theta^{t+1} - \Theta^t) = -\nabla \ell(\Theta^t)$

NOTE 1: modèle surparamétré, car

$$\text{si } \Theta \leftarrow \Theta - \psi = \begin{bmatrix} \Theta_{1,:} - \psi \\ \vdots \\ \Theta_{K,:} - \psi \end{bmatrix} \text{ alors } \begin{bmatrix} \mathbb{P}_{Y|X=x}^{\Theta}(Y=1) \\ \vdots \\ \mathbb{P}_{Y|X=x}^{\Theta}(Y=K) \end{bmatrix} = \begin{bmatrix} \mathbb{P}_{Y|X=x}^{\Theta-\psi}(Y=1) \\ \vdots \\ \mathbb{P}_{Y|X=x}^{\Theta-\psi}(Y=K) \end{bmatrix}$$

2.2 Régression logistique multinomiale

Estimation par descente de gradient $(\Theta^{t+1} - \Theta^t) = -\nabla \ell(\Theta^t)$

NOTE 1: modèle surparamétré, car

$$\text{si } \Theta \leftarrow \Theta - \psi = \begin{bmatrix} \Theta_{1,:} - \psi \\ \vdots \\ \Theta_{K,:} - \psi \end{bmatrix} \text{ alors } \begin{bmatrix} \mathbb{P}_{Y|X=x}^{\Theta}(Y=1) \\ \vdots \\ \mathbb{P}_{Y|X=x}^{\Theta}(Y=K) \end{bmatrix} = \begin{bmatrix} \mathbb{P}_{Y|X=x}^{\Theta-\psi}(Y=1) \\ \vdots \\ \mathbb{P}_{Y|X=x}^{\Theta-\psi}(Y=K) \end{bmatrix}$$

On fixe donc $\Theta_{K,:} = \mathbf{1}$

2.2 Régression logistique multinomiale

Estimation par descente de gradient $(\Theta^{t+1} - \Theta^t) = -\nabla \ell(\Theta^t)$

NOTE 1: modèle surparamétré, car

NOTE 2: généralisation de la régression logistique binaire

2.2 Régression logistique multinomiale

Estimation par descente de gradient $(\Theta^{t+1} - \Theta^t) = -\nabla \ell(\Theta^t)$

NOTE 1: modèle surparamétré, car

NOTE 2: généralisation de la régression logistique binaire

NOTE 3: donne le même modèle de classification que K modèles de regression logistique binaire combinés en stratégie multiclasse one-vs-rest

$$f_{\theta^1}^1 = \text{Reg. log. binaire } Y = 1 \text{ vs } Y \neq 1$$

...

$$f_{\theta^K}^K = \text{Reg. log. binaire } Y = K \text{ vs } Y \neq K$$

2.3 Régression linéaire

Modèle statistique de regression Prédire Y à partir de X

Jeu de données $x_{1:n}, y_{1:n}$

Un modèle de regression est une famille de lois de probabilités *conditionnelles*

$$\mathcal{M}_{\Theta} = \{\mathbb{P}_{Y|X=x}^{\theta} | \theta \in \Theta, x \in \mathcal{X}\}$$

Regression linéaire $Y \in \mathbb{R}, \quad X \in \mathbb{R}^p$

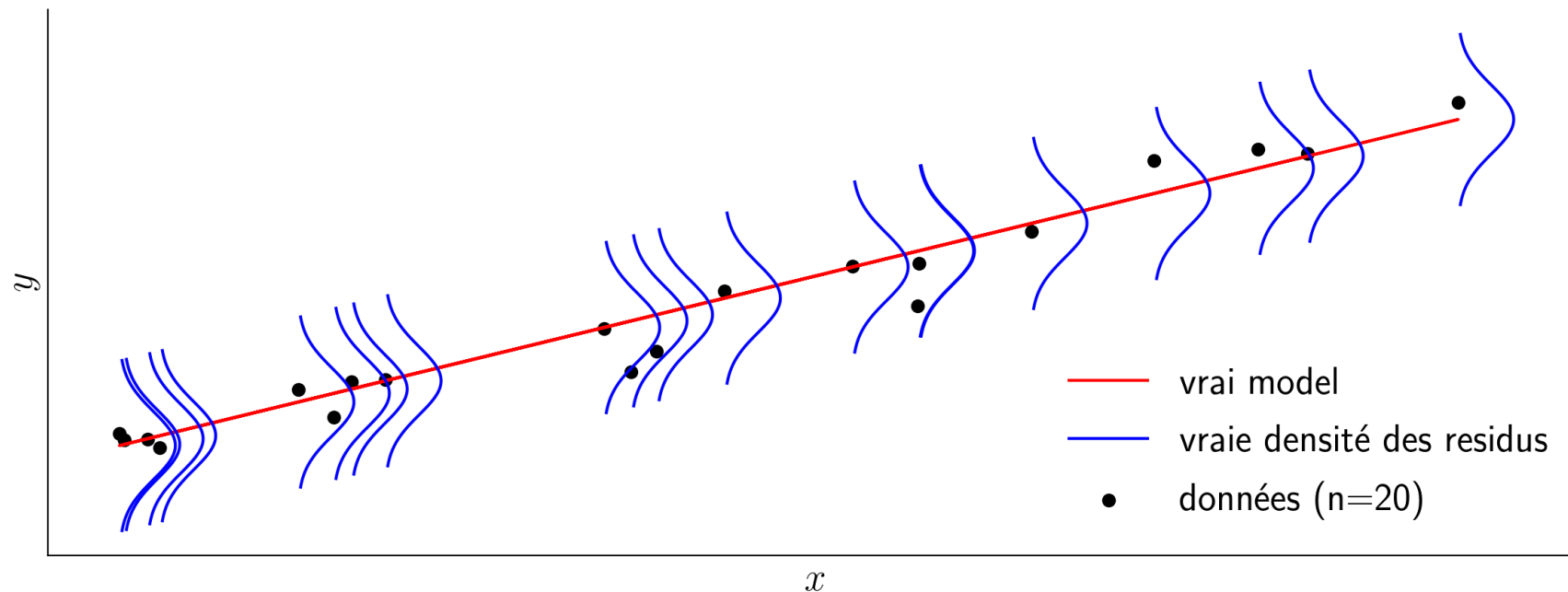
$$\mathbb{P}_{Y|X=x}^{\theta} = \mathcal{N}(\beta_0 + x^{\top} \beta, \sigma^2) \quad \theta = (\beta_0, \beta_1, \dots, \beta_p, \sigma^2)$$

2.3 Régression linéaire

Régression linéaire

$$Y \in \mathbb{R}, \quad X \in \mathbb{R}^p$$

$$\mathbb{P}_{Y|X=x}^{\theta} = \mathcal{N}(\beta_0 + x^{\top} \beta, \sigma^2) \quad \theta = (\beta_0, \beta_1, \dots, \beta_p, \sigma^2)$$



2.3 Régression linéaire

Estimation par maximum de vraisemblance

$$\begin{aligned}\mathcal{L}(\theta; \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) &= \prod_{i=1}^n p_{Y_i|X_i=\mathbf{x}_i}(y_i) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{\sigma^2 2\pi}} e^{-\frac{1}{2\sigma^2} (y_i - \beta_0 - \mathbf{x}_i^\top \beta)^2}\end{aligned}$$

2.3 Régression linéaire

Estimation par maximum de vraisemblance

$$\begin{aligned}\mathcal{L}(\theta; \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) &= \prod_{i=1}^n p_{Y_i|X_i=\mathbf{x}_i}(y_i) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{\sigma^2 2\pi}} e^{-\frac{1}{2\sigma^2} (y_i - \beta_0 - \mathbf{x}_i^\top \beta)^2}\end{aligned}$$

On peut montrer que $\theta \mapsto \mathcal{L}(\theta; \mathbf{x}_{1:n}, \mathbf{y}_{1:n})$ est concave.

$\theta \mapsto \ell(\theta) = -\log \mathcal{L}(\theta; \mathbf{x}_{1:n}, \mathbf{y}_{1:n})$ est convexe.

2.3 Régression linéaire

Estimation par maximum de vraisemblance

$$\begin{aligned}\mathcal{L}(\theta; \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) &= \prod_{i=1}^n p_{Y_i|X_i=\mathbf{x}_i}(y_i) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{\sigma^2 2\pi}} e^{-\frac{1}{2\sigma^2} (y_i - \beta_0 - \mathbf{x}_i^\top \beta)^2}\end{aligned}$$

On peut montrer que $\theta \mapsto \mathcal{L}(\theta; \mathbf{x}_{1:n}, \mathbf{y}_{1:n})$ est concave.

$\theta \mapsto \ell(\theta) = -\log \mathcal{L}(\theta; \mathbf{x}_{1:n}, \mathbf{y}_{1:n})$ est convexe.

$$\ell(\theta; \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) = \frac{n}{2} 2\pi + \frac{n}{2} \log(\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}_i^\top \beta)^2$$

2.3 Régression linéaire

Estimation par maximum de vraisemblance

$$\begin{aligned}\mathcal{L}(\theta; \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) &= \prod_{i=1}^n p_{Y_i|X_i=\mathbf{x}_i}(y_i) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{\sigma^2 2\pi}} e^{-\frac{1}{2\sigma^2} (y_i - \beta_0 - \mathbf{x}_i^\top \beta)^2}\end{aligned}$$

On peut montrer que $\theta \mapsto \mathcal{L}(\theta; \mathbf{x}_{1:n}, \mathbf{y}_{1:n})$ est concave.

$\theta \mapsto \ell(\theta) = -\log \mathcal{L}(\theta; \mathbf{x}_{1:n}, \mathbf{y}_{1:n})$ est convexe.

Moindres carrés !

$$\ell(\theta; \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) = \frac{n}{2} 2\pi + \frac{n}{2} \log(\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}_i^\top \beta)^2$$

2.3 Régression linéaire

Estimation par maximum de vraisemblance

Le minimum est réalisé aux dérivées nulles, i.e

$$\frac{\partial \ell}{\partial \beta_0} = 0, \frac{\partial \ell}{\partial \beta_1} = 0, \dots, \frac{\partial \ell}{\partial \beta_p} = 0, \text{ et } \frac{\partial \ell}{\partial \sigma^2} = 0$$

2.3 Régression linéaire

Estimation par maximum de vraisemblance

Le minimum est réalisé aux dérivées nulles, i.e

$$\frac{\partial \ell}{\partial \beta_0} = 0, \frac{\partial \ell}{\partial \beta_1} = 0, \dots, \frac{\partial \ell}{\partial \beta_p} = 0, \text{ et } \frac{\partial \ell}{\partial \sigma^2} = 0$$

Tous calculs faits

$$\hat{\beta}(\mathbf{X}_{1:n}, \mathbf{Y}_{1:n}) = \left(\tilde{\mathbf{X}}_{1:n}^\top \tilde{\mathbf{X}}_{1:n} \right)^{-1} \tilde{\mathbf{X}}_{1:n}^\top \mathbf{Y}_{1:n} \quad (1)$$

Avec

$$\tilde{\mathbf{X}}_{1:n} = \begin{bmatrix} 1 & \cdots & 1 \\ x_{1,1} & \cdots & x_{1,p} \\ \vdots & \ddots & \vdots \\ x_{n,1} & \cdots & x_{n,p} \end{bmatrix}$$

2.4 Pénalisation lasso

Observation: Les estimateurs précédents sans biais mais forte variance.

Idée: Diminuer la variance au prix de biais

$$\min_{\beta_0, \beta} = \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}_i^\top \beta)^2$$



$$\min_{\beta_0, \beta} = \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}_i^\top \beta)^2$$

$$\text{s.t } \sum_{i=0}^p |\beta_i| = \|\beta\|_1 \leq t$$

2.4 Pénalisation lasso

Observation: Les estimateurs précédents sans biais mais forte variance.

Idée: Diminuer la variance au prix de biais

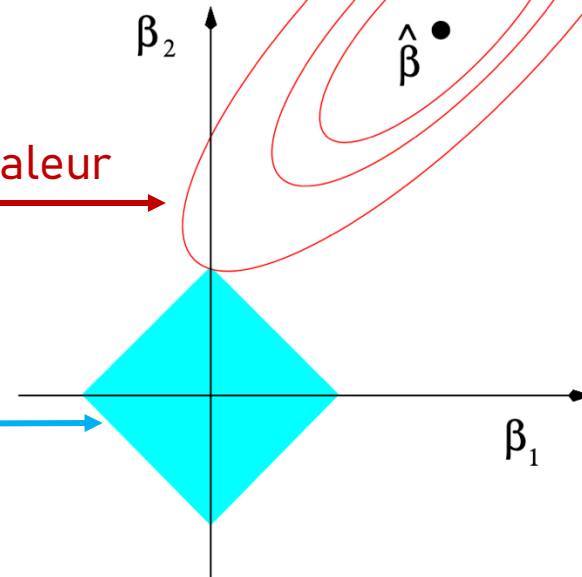
$\hat{\beta}$ Estimateur moindres carrés

$$\min_{\beta_0, \beta} = \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}_i^\top \beta)^2$$

$$\min_{\beta_0, \beta} = \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}_i^\top \beta)^2$$

$$\text{s.t } \sum_{i=0}^p |\beta_i| = \|\beta\|_1 \leq t$$

isovaleur



2.4 Pénalisation lasso

Observation: Les estimateurs précédents sans biais mais forte variance.

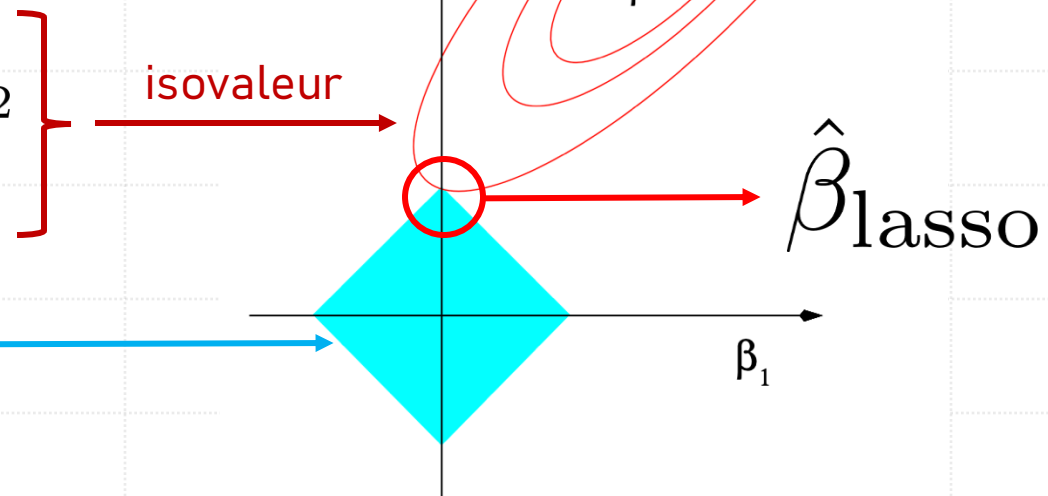
Idée: Diminuer la variance au prix de biais

$$\min_{\beta_0, \beta} = \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}_i^\top \beta)^2$$

$\hat{\beta}$ Estimateur moindres carrés

$$\min_{\beta_0, \beta} = \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}_i^\top \beta)^2$$

$$\text{s.t } \sum_{i=0}^p |\beta_i| = \|\beta\|_1 \leq t$$



2.4 Pénalisation lasso

$$\begin{aligned} \min_{\beta_0, \beta} &= \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}_i^\top \beta)^2 \\ \text{s.t. } &\sum_{i=0}^p |\beta_i| = \|\beta\|_1 \leq t \end{aligned} \quad \begin{array}{c} \text{Dualité de Lagrange} \\ \Longleftrightarrow \end{array} \quad \begin{aligned} \min_{\beta_0, \beta} &= \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}_i^\top \beta)^2 + \lambda \|\beta\|_1 \end{aligned}$$

Diagram illustrating the relationship between the constrained minimization problem and the unconstrained minimization problem with a penalty term. A red circle highlights the variable t in the constraint, and a red arrow points from it to the variable λ in the penalty term, indicating the relationship $\lambda = t/2$. Another red arrow points from the λ term to the ratio $1:1$.

Question: Comment choisir t/λ ?

Réponse: Par validation croisée!

2.4 Pénalisation lasso

$$\min_{\beta_0, \beta} = \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}_i^\top \beta)^2$$

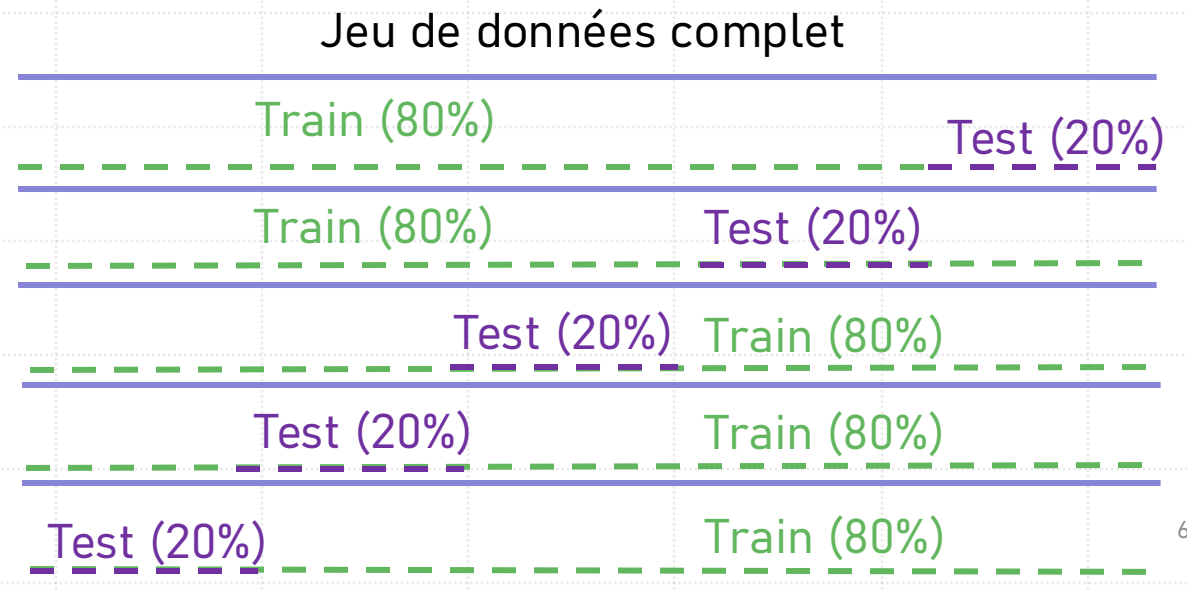
Dualité de Lagrange \iff

$$\min_{\beta_0, \beta} = \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}_i^\top \beta)^2 + \lambda \|\beta\|_1$$

s.t $\sum_{i=0}^p |\beta_i| = \|\beta\|_1 \leq t$ → 1:1 ←

Question: Comment choisir t/λ ?

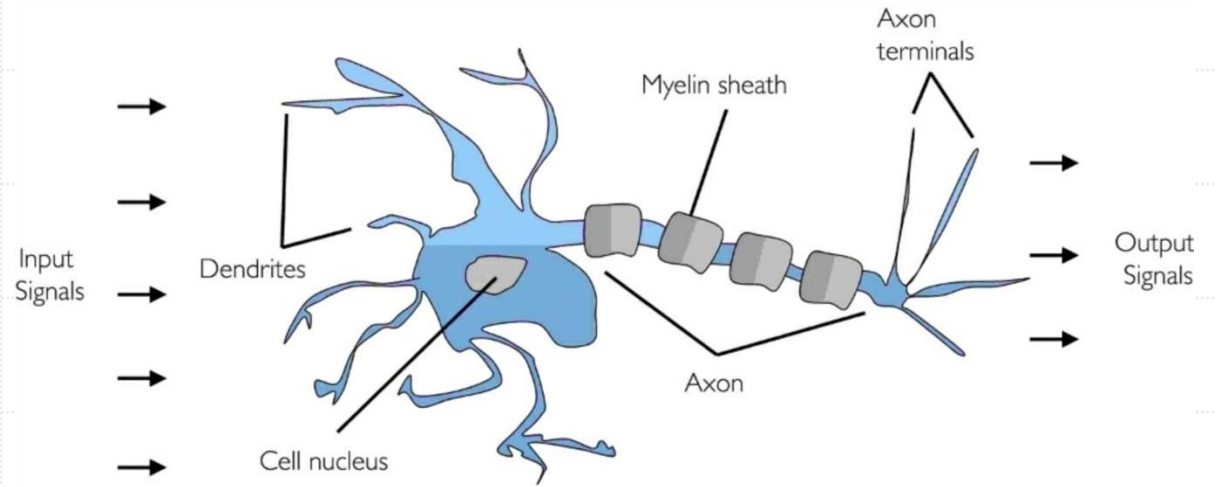
Réponse: Par validation croisée!



2.5 Réseau monocouche

Un neurone biologique:

si somme des signaux entrées > seuil
-> potentiel d'action généré
sinon inactif.



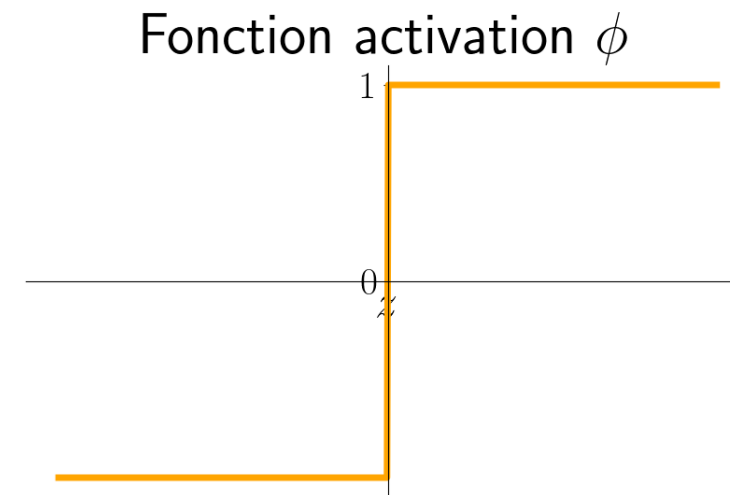
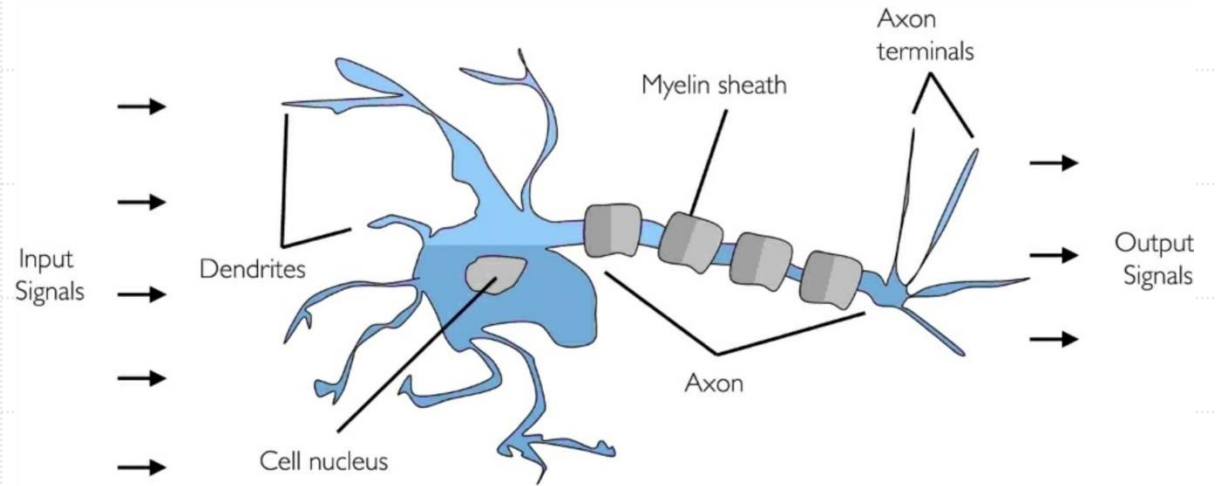
2.5 Réseau monocouche

Un neurone biologique:

si somme des signaux entrées > seuil
→ potentiel d'action généré
sinon inactif.

Un neurone artificiel Rosenblatt 1957

- Entrées $x = \begin{bmatrix} x_1 \\ \vdots \\ x_p \end{bmatrix}$, poids $w = \begin{bmatrix} w_1 \\ \vdots \\ w_p \end{bmatrix}$
- Entrées agrégées $z = x^\top w$
- Activation $\phi(z)$.



2.5 Réseau monocouche

Règle du Perceptron: jeu de données $x_{1:n}, y_{1:n}$

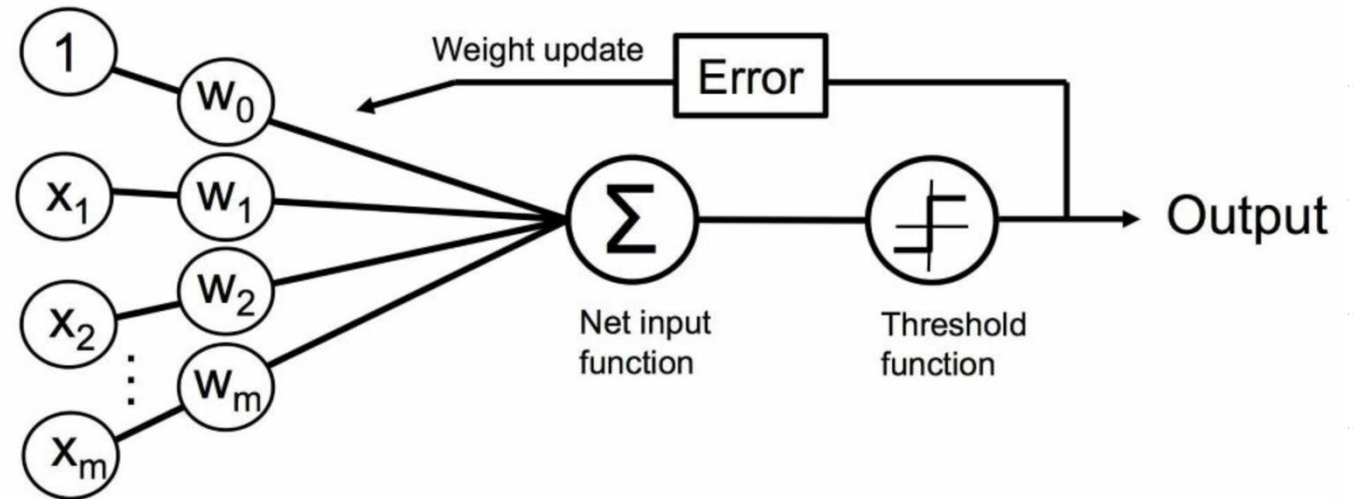
1. initialisation w aléatoire;

2. pour chaque $i = 1, \dots, n$,

(a) $\hat{y}_i = \phi(x_i^\top w)$

(b) $\delta = \eta(y_i - \hat{y}_i)$.

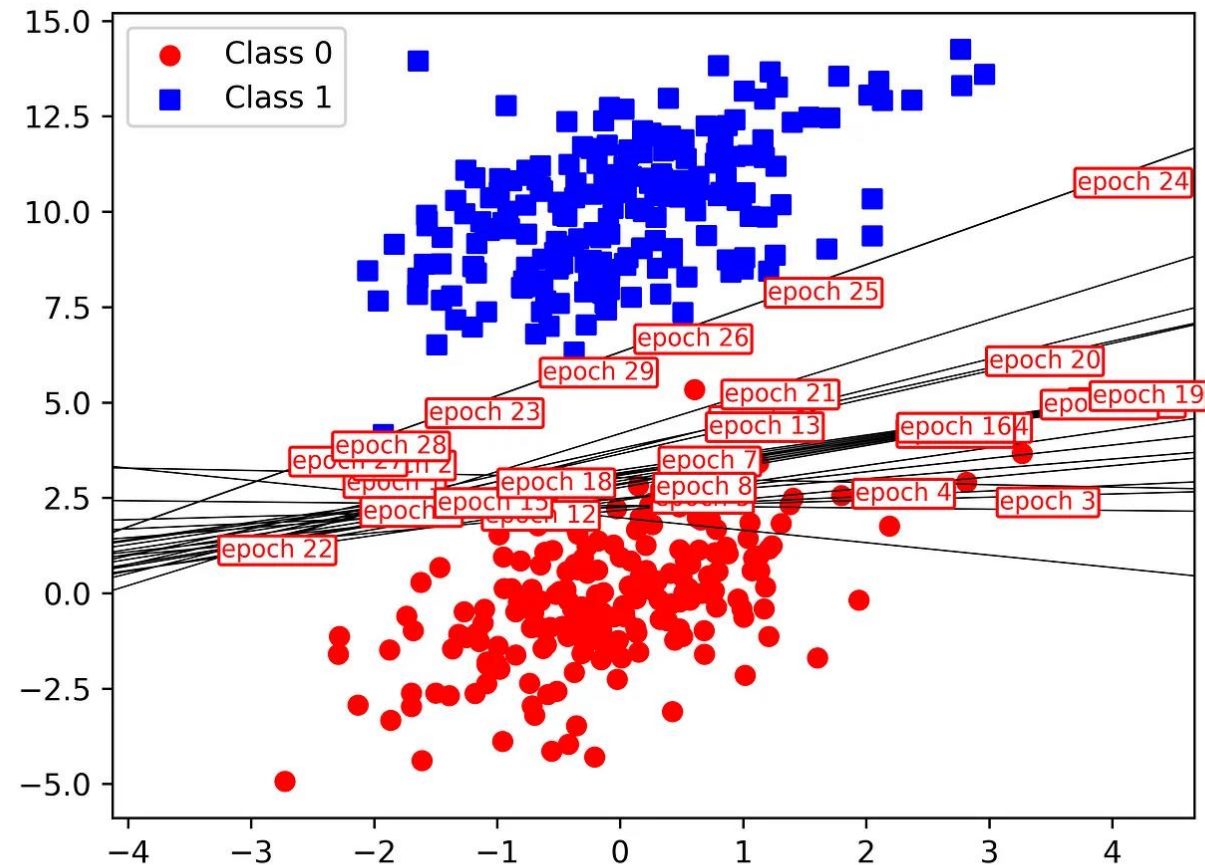
(c) $w = w + \delta x_i$



2.5 Réseau monocouche

Règle du Perceptron: jeu de données $x_{1:n}, y_{1:n}$

1. initialisation w aléatoire;
2. pour chaque $i = 1, \dots, n$,
 - (a) $\hat{y}_i = \phi(x_i^\top w)$
 - (b) $\delta = \eta(y_i - \hat{y}_i)$.
 - (c) $w = w + \delta x_i$



2.5 Réseau monocouche

Règle du Perceptron: jeu de données $x_{1:n}, y_{1:n}$

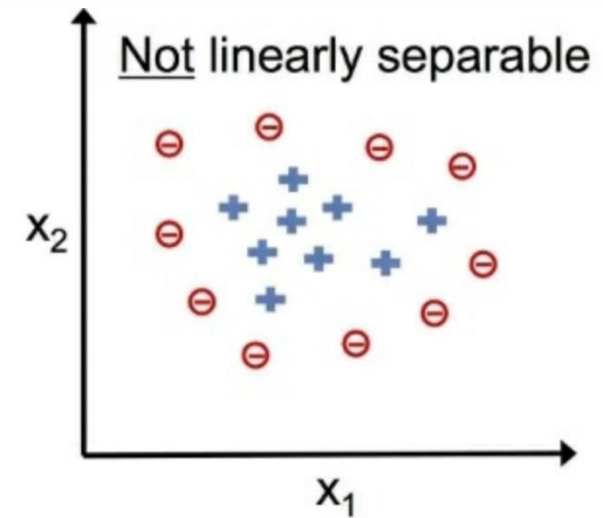
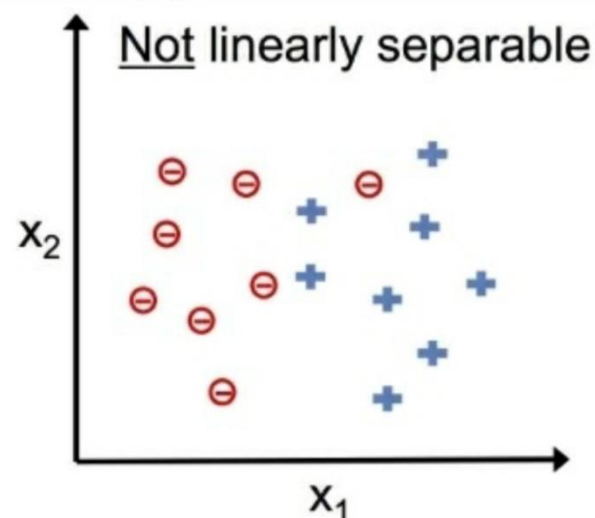
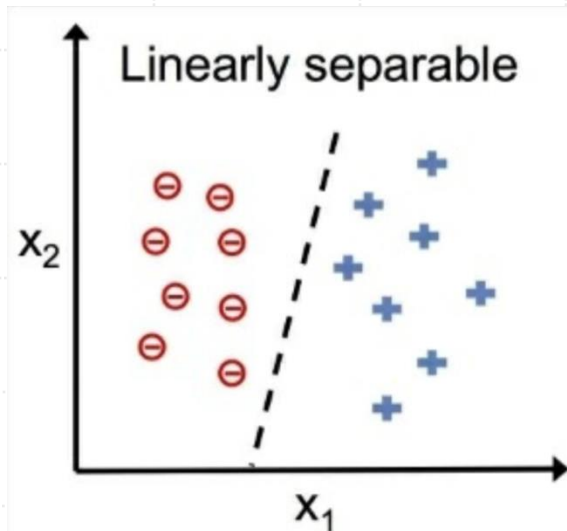
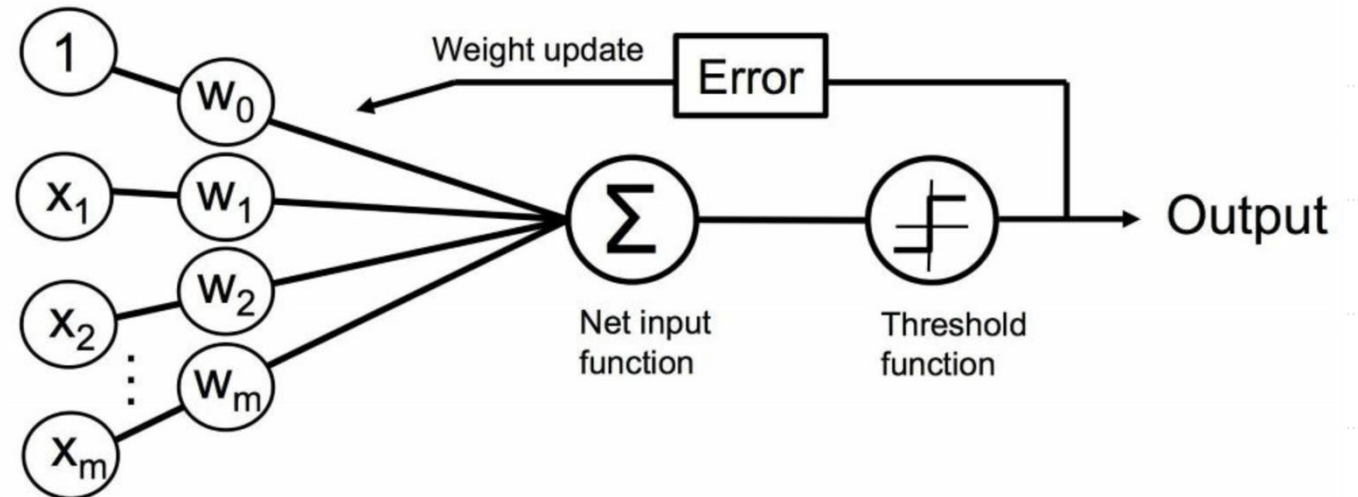
1. initialisation w aléatoire;

2. pour chaque $i = 1, \dots, n$,

(a) $\hat{y}_i = \phi(x_i^\top w)$

(b) $\delta = \eta(y_i - \hat{y}_i)$.

(c) $w = w + \delta x_i$



2.5 Réseau monocouche

Reseau Adaline: jeu de données $x_{1:n}, y_{1:n}$

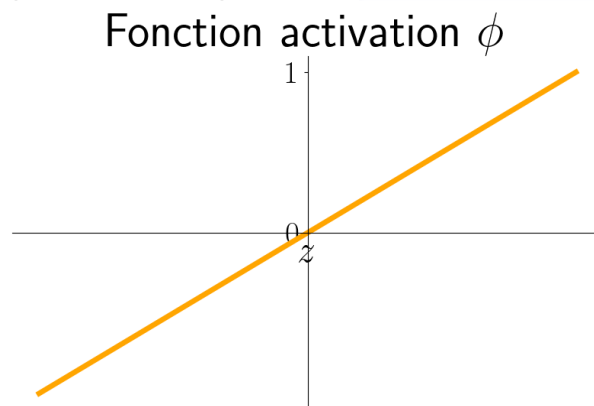
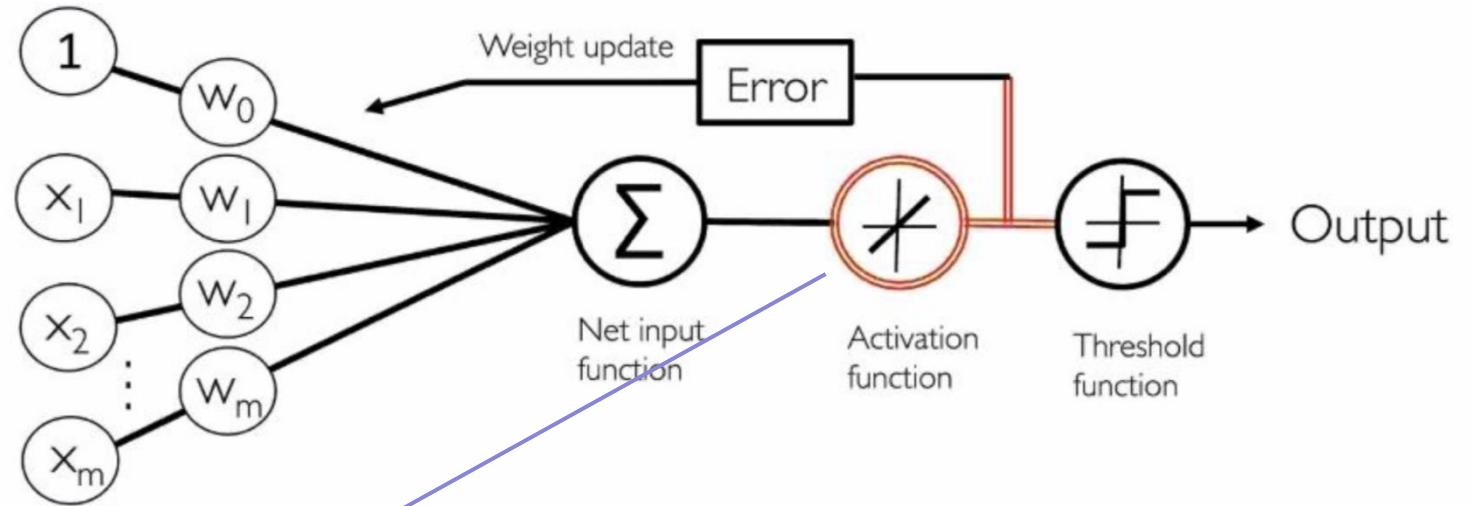
1. initialisation w aléatoire;

2. pour chaque $i = 1, \dots, n$,

(a) $\hat{y}_i = \phi(x_i^\top w)$

(b) $\delta = \eta(y_i - x_i^\top w)$.

(c) $w = w + \delta x_i$

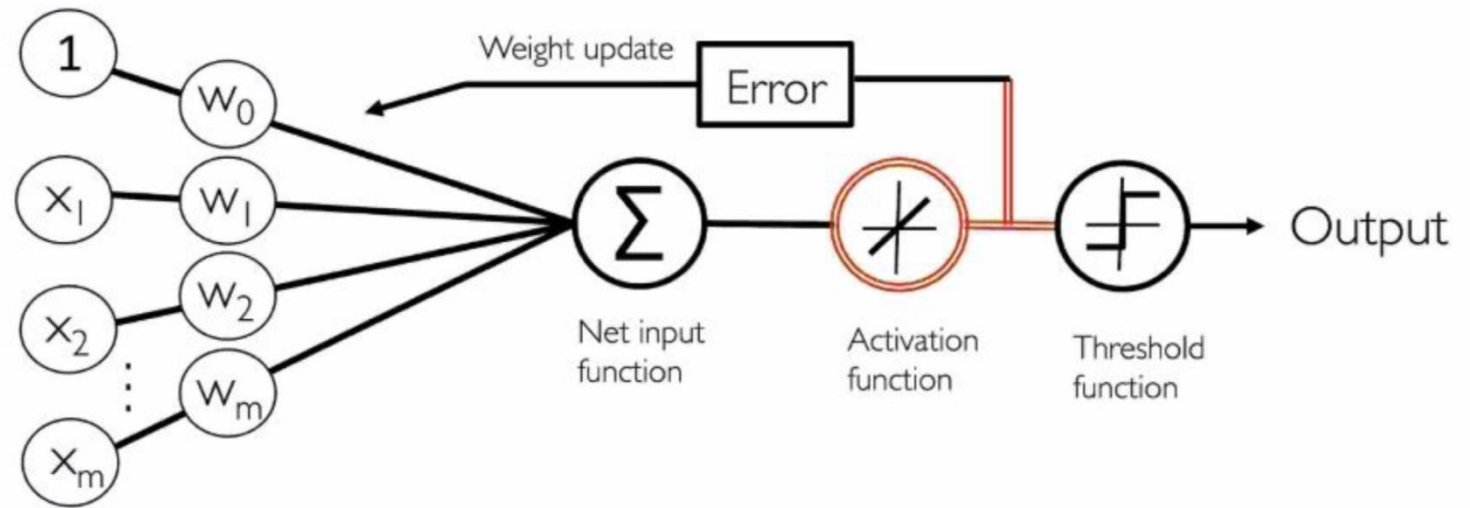


Adaptive Linear Neuron (Adaline)

2.5 Réseau monocouche

Reseau Adaline: jeu de données $x_{1:n}, y_{1:n}$

1. initialisation w aléatoire;
2. pour chaque $i = 1, \dots, n$,
 - (a) $\hat{y}_i = \phi(x_i^\top w)$
 - (b) $\delta = \eta(y_i - x_i^\top w)$.
 - (c) $w = w + \delta x_i$



Adaptive Linear Neuron (Adaline)

Pourquoi: $y_i - x_i^\top w$

$$\hat{w}(x_{1:n}, y_{1:n}) \in \operatorname{argmin}_{w \in \mathbb{R}^p} J(w) = \frac{1}{2} \sum_{i=1}^n (y_i - x_i^\top w)^2$$

descente de gradient sur un échantillon pas = $-\nabla J_i(w) = x_i(y_i - x_i^\top w)$

2.5 Réseau monocouche

Reseau Adaline: jeu de données $x_{1:n}, y_{1:n}$

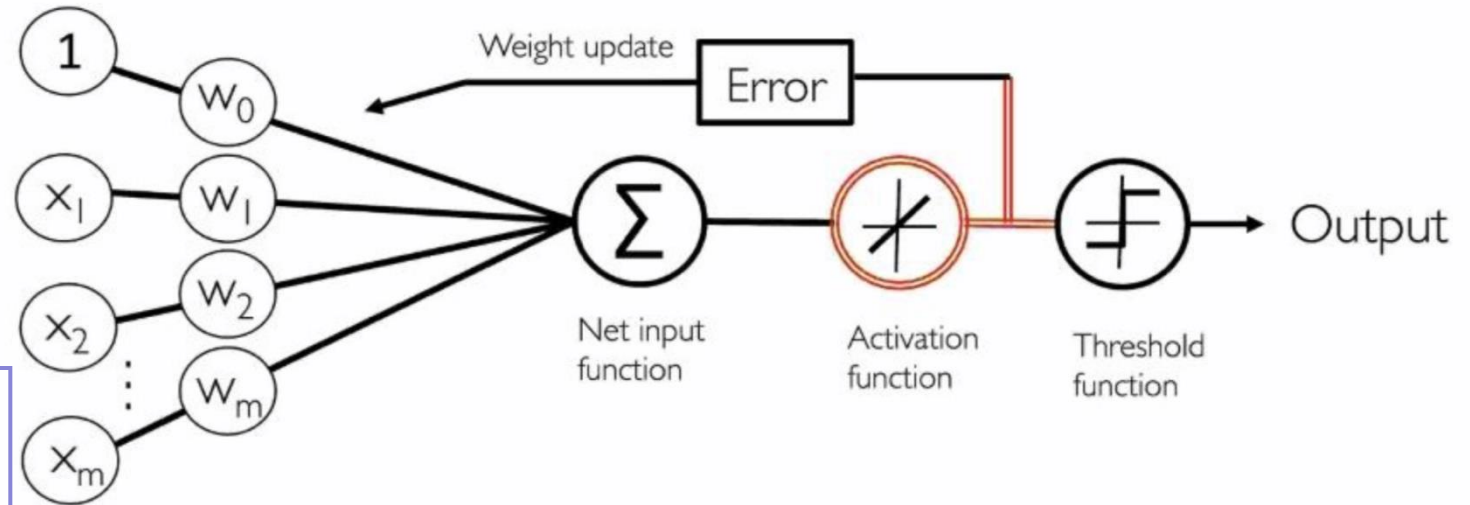
1. initialisation w aléatoire;

2. pour chaque $i = 1, \dots, n$,

(a) $\hat{y}_i = \phi(x_i^\top w)$

(b) $\delta = \eta(y_i - x_i^\top w)$

(c) $w = w + \delta x_i$



Adaptive Linear Neuron (Adaline)

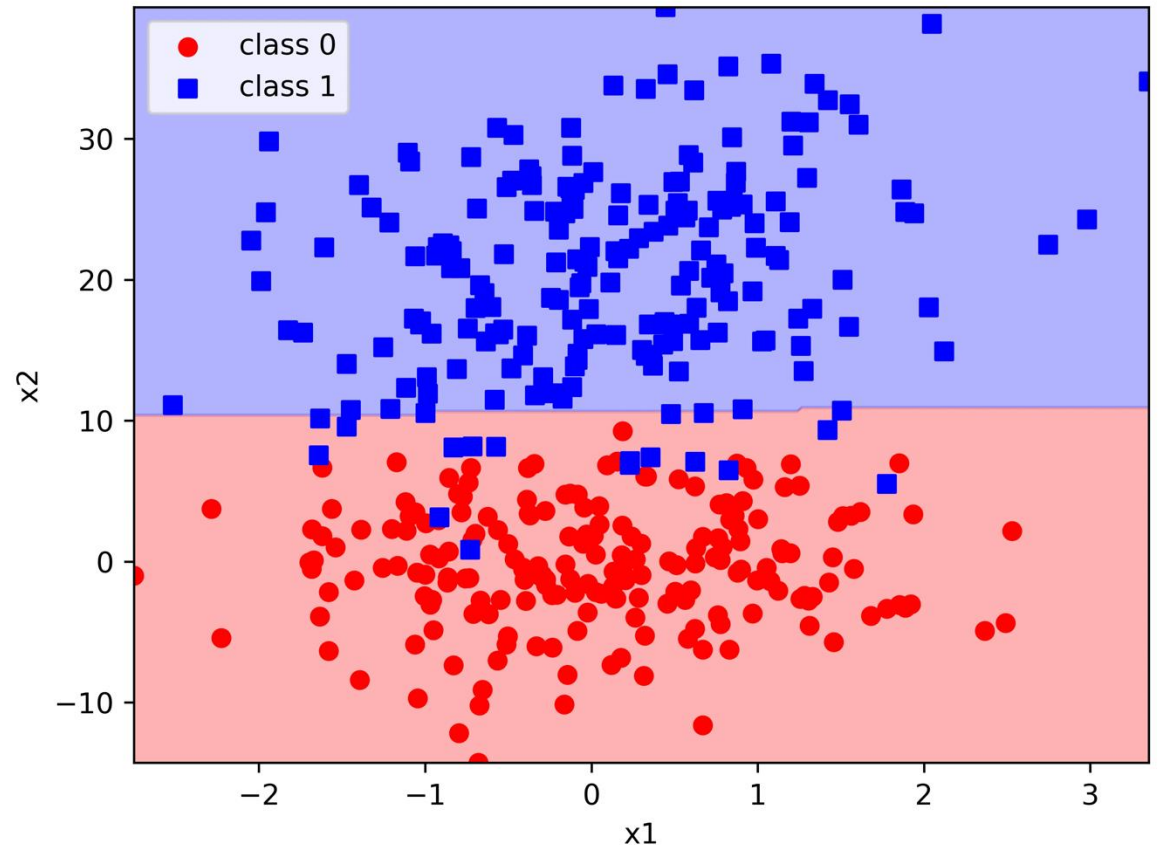
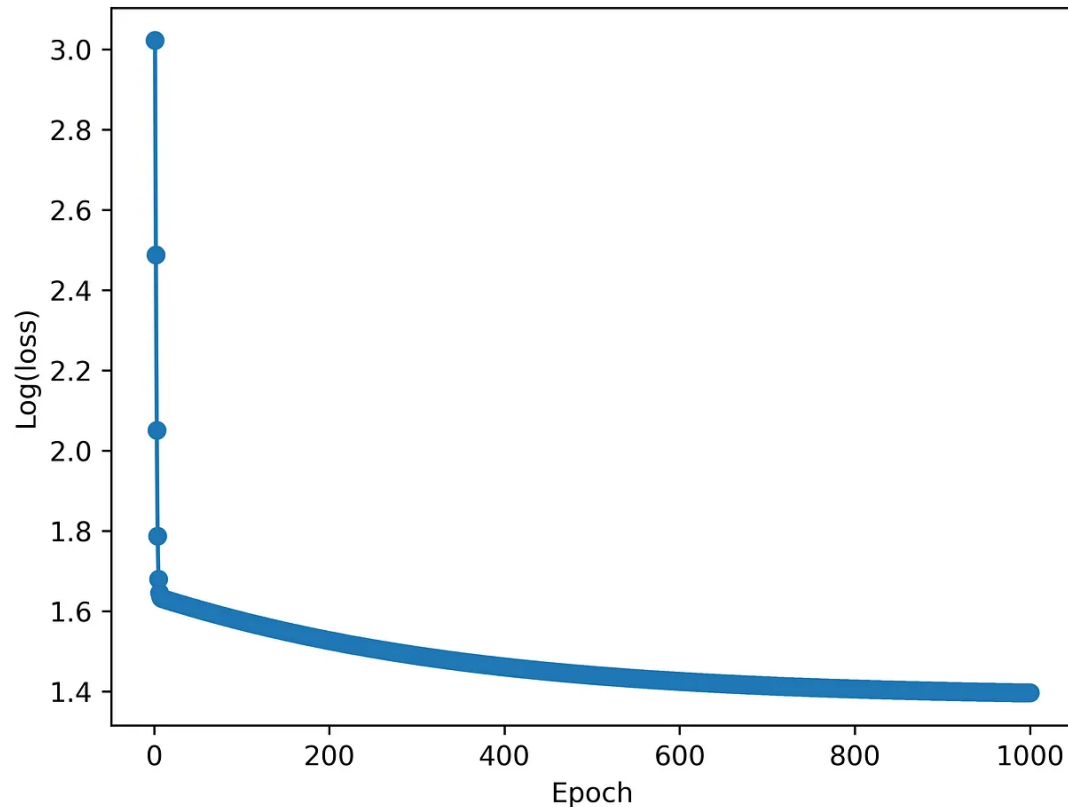
Pourquoi: $y_i - x_i^\top w$

$$\hat{w}(x_{1:n}, y_{1:n}) \in \operatorname{argmin}_{w \in \mathbb{R}^p} J(w) = \frac{1}{2} \sum_{i=1}^n (y_i - x_i^\top w)^2$$

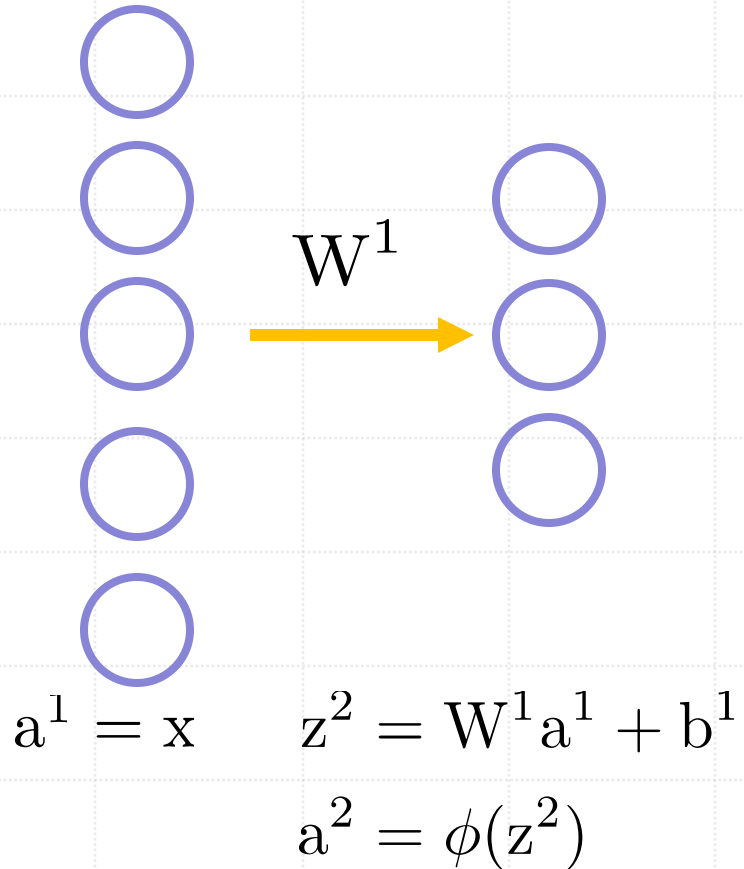
descente de gradient sur un échantillon pas = $-\nabla J_i(w) = x_i(y_i - x_i^\top w)$

2.5 Réseau monocouche

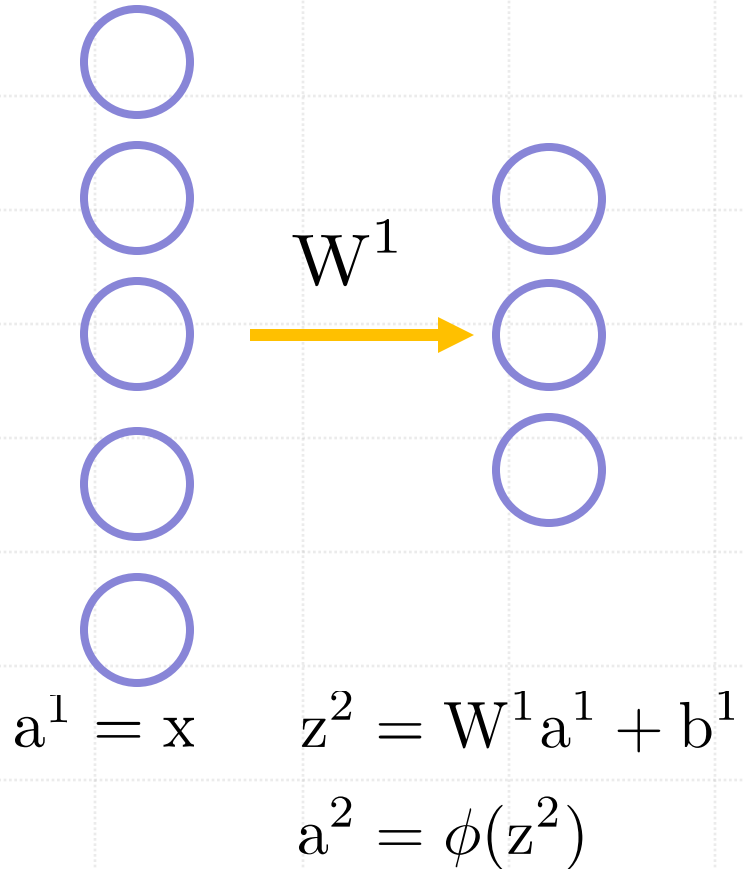
Reseau Adaline: jeu de données $x_{1:n}, y_{1:n}$



2.6 Réseau multicouches

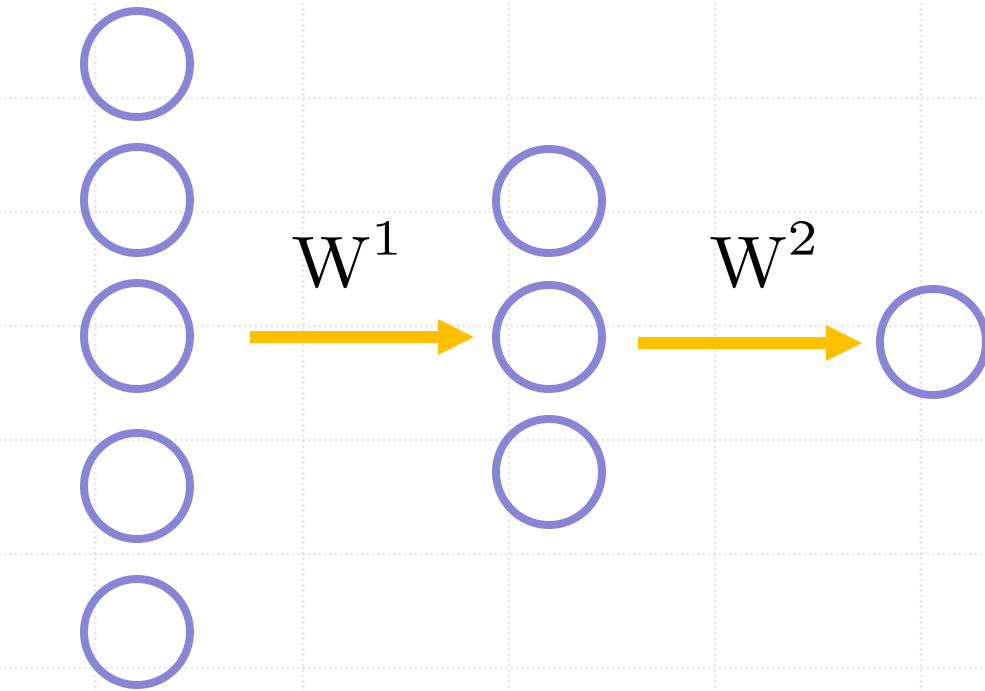


2.6 Réseau multicouches



$$W^1 = \begin{bmatrix} w_{11}^1 & \cdots & w_{15}^1 \\ \vdots & \ddots & \vdots \\ w_{31}^1 & \cdots & w_{35}^1 \end{bmatrix} \quad b^1 = \begin{bmatrix} b_1^1 \\ \vdots \\ b_3^1 \end{bmatrix}$$

2.6 Réseau multicouches

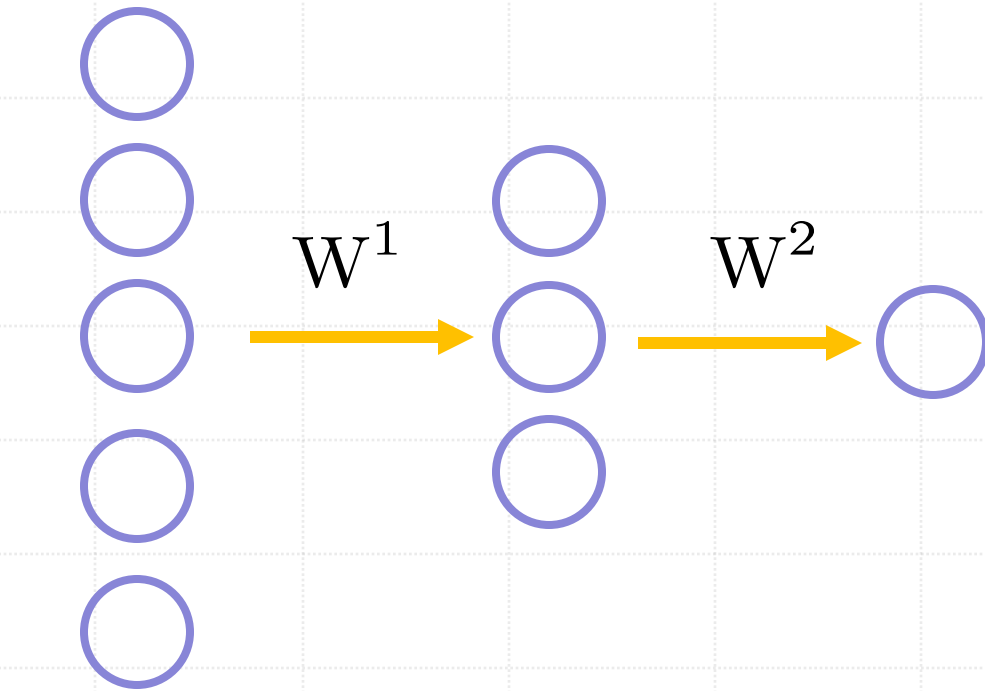


$$W^1 = \begin{bmatrix} w_{11}^1 & \cdots & w_{15}^1 \\ \vdots & \ddots & \vdots \\ w_{31}^1 & \cdots & w_{35}^1 \end{bmatrix} \quad b^1 = \begin{bmatrix} b_1^1 \\ \vdots \\ b_3^1 \end{bmatrix}$$

$$W^2 = [w_{11}^2 \quad \cdots \quad w_{13}^2] \quad b^2 = [b_1^2]$$

$$\begin{aligned} a^1 &= x & z^2 &= W^1 a^1 + b^1 & z^3 &= W^2 a^2 + b^2 \\ a^2 &= \phi(z^2) & a^3 &= \phi(z^3) \end{aligned}$$

2.6 Réseau multicouches



$$W^1 = \begin{bmatrix} w_{11}^1 & \cdots & w_{15}^1 \\ \vdots & \ddots & \vdots \\ w_{31}^1 & \cdots & w_{35}^1 \end{bmatrix} \quad b^1 = \begin{bmatrix} b_1^1 \\ \vdots \\ b_3^1 \end{bmatrix}$$

$$W^2 = [w_{11}^2 \quad \cdots \quad w_{13}^2] \quad b^2 = [b_1^2]$$

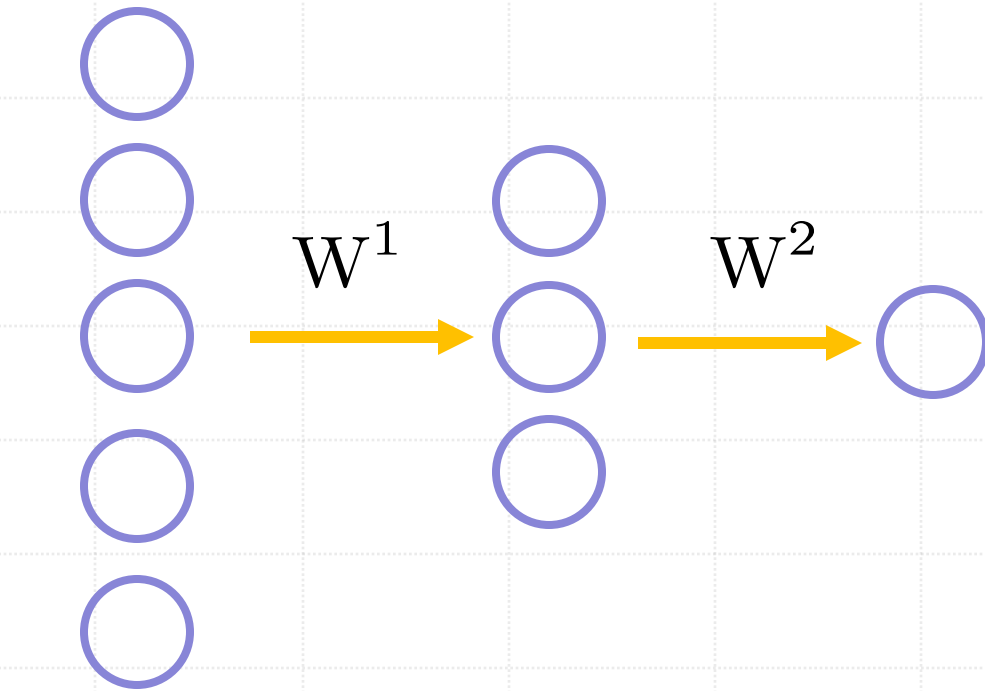
$$a^1 = x \quad z^2 = W^1 a^1 + b^1 \quad z^3 = W^2 a^2 + b^2$$

$$a^2 = \phi(z^2) \quad a^3 = \phi(z^3)$$

Q1: combien de paramètres ?

Q2: que se passe-t-il si $\phi = \text{Id}$?

2.6 Réseau multicouches



$$W^1 = \begin{bmatrix} w_{11}^1 & \cdots & w_{15}^1 \\ \vdots & \ddots & \vdots \\ w_{31}^1 & \cdots & w_{35}^1 \end{bmatrix} \quad b^1 = \begin{bmatrix} b_1^1 \\ \vdots \\ b_3^1 \end{bmatrix}$$

$$W^2 = [w_{11}^2 \quad \cdots \quad w_{13}^2] \quad b^2 = [b_1^2]$$

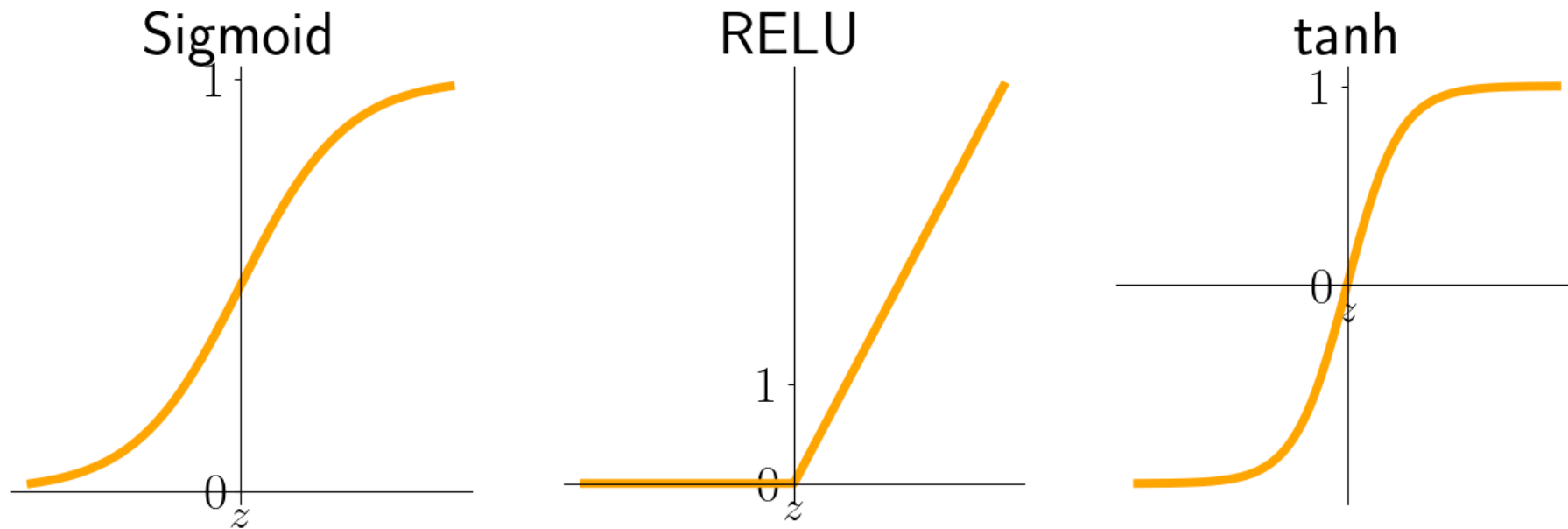
$$a^1 = x \quad z^2 = W^1 a^1 + b^1 \quad z^3 = W^2 a^2 + b^2$$

$$a^2 = \phi(z^2) \quad a^3 = \phi(z^3)$$

Q1: combien de paramètres ? $3 \times 5 + 3 + 1 \times 3 + 1 = 22$

Q2: que se passe-t-il si $\phi = \text{Id}$? $a^3 = W^2 W^1 x + W^2 b^1 + b^2$
 $= mx + p$

2.6 Réseau multicouches



Exemples de fonctions d'activation non linéaires

2.6 Réseau multicouches

Comment apprendre les paramètres $\Theta = [W^1, b^1, W^2, b^2]$?

Jeu de données $x_{1:n}, y_{1:n}$ Modèle $f_{\Theta}(x) = a^3(x)$

2.6 Réseau multicouches

Comment apprendre les paramètres $\Theta = [W^1, b^1, W^2, b^2]$?

Jeu de données $x_{1:n}, y_{1:n}$

Modèle $f_{\Theta}(x) = a^3(x)$

-> **Choix** d'une fonction de coût

Exemples: $y_i \in \mathbb{R}$ $J(\Theta) = \frac{1}{2n} \sum_{i=1}^n \|y_i - f_{\Theta}(x)\|^2$

`tf.keras.losses.MeanSquaredError`

2.6 Réseau multicouches

Comment apprendre les paramètres $\Theta = [W^1, b^1, W^2, b^2]$?

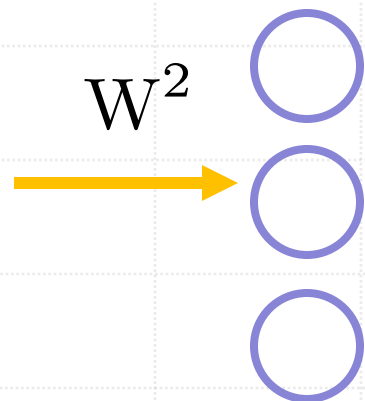
Jeu de données $x_{1:n}, y_{1:n}$

Modèle $f_{\Theta}(x) = a^3(x)$

-> **Choix** d'une fonction de coût

Exemples: $y_i \in \mathbb{R}$ $J(\Theta) = \frac{1}{2n} \sum_{i=1}^n \|y_i - f_{\Theta}(x_i)\|^2$ `tf.keras.losses.MeanSquaredError`

$y_i \in \{1, \dots, K\}$ $J(\Theta) = \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K \mathbb{1}_{y_i=k} \log((f_{\Theta}(x_i))_k)$ `tf.keras.losses.CategoricalCrossEntropy`



$$z^3 = W^2 a^2 + b^2$$

$$a^3 = \phi(z^3)$$

$$\phi(z) = \begin{bmatrix} \frac{e^{z_1}}{\sum_{k=1}^K e^{z_k}} \\ \vdots \\ \frac{e^{z_K}}{\sum_{k=1}^K e^{z_k}} \end{bmatrix} \quad \text{"probabilities" (!)}$$

Softmax activation

2.6 Réseau multicouches

Comment apprendre les paramètres $\Theta = [W^1, b^1, W^2, b^2]$?

Jeu de données $x_{1:n}, y_{1:n}$

Modèle $f_{\Theta}(x) = a^3(x)$

-> **Choix** d'une fonction de coût

-> Apprentissage par descente de gradient

$$\Theta^{t+1} = \Theta^t - \eta_t \nabla J(\Theta^t)$$

2.6 Réseau multicouches

Comment apprendre les paramètres $\Theta = [W^1, b^1, W^2, b^2]$?

Jeu de données $x_{1:n}, y_{1:n}$

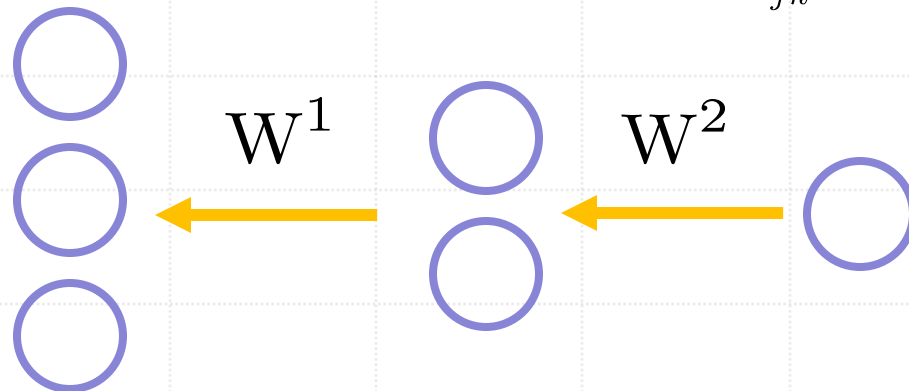
Modèle $f_{\Theta}(x) = a^3(x)$

-> **Choix** d'une fonction de coût

-> Apprentissage par descente de gradient

$$\Theta^{t+1} = \Theta^t - \eta_t \nabla J(\Theta^t)$$

Difficulté 1: comment calculer $\frac{\partial J}{\partial w_{jk}^l}(\Theta^t), \frac{\partial J}{\partial b_j^l}(\Theta^t)$ pour tout j, k, l ?



$\epsilon_i = y_i - f_{\Theta^t}(x_i)$ → Stochastic descent

$\epsilon_1 = y_i - f_{\Theta^t}(x_i)$

...

$\epsilon_m = y_i - f_{\Theta^t}(x_n)$

Batch descent

2.6 Réseau multicouches

Comment apprendre les paramètres $\Theta = [W^1, b^1, W^2, b^2]$?

Jeu de données $x_{1:n}, y_{1:n}$

Modèle $f_{\Theta}(x) = a^3(x)$

-> **Choix** d'une fonction de coût

-> Apprentissage par descente de gradient

$$\Theta^{t+1} = \Theta^t - \eta_t \nabla J(\Theta^t)$$

Difficulté 1: comment calculer $\frac{\partial J}{\partial w_{jk}^l}(\Theta^t), \frac{\partial J}{\partial b_j^l}(\Theta^t)$ pour tout j, k, l ?

Difficulté 2: comment choisir η_t ?

Constant, RMSProp, Adam, Adamax, Adadelata, Nadam, Adafactor, ...