

# Réseaux de neurones

Machine learning

**Le Chenadec** Gilles (M103)

Quidu Isabelle & Thomas Hélène

# 1 Objectifs du cours

## 2 Introduction

### 3 Apprentissage supervisé

- Régression
  - Régression linéaire
- Classification supervisée
- Régression logistique
- Réseaux de neurones



# Objectifs du cours

# Introduction

# Apprentissage supervisé

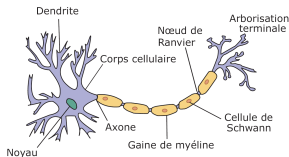
- Régression
- Classification supervisée
- Régression logistique
- Réseaux de neurones



# Réseaux de neurones

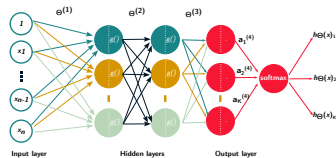
## Origines:

- Imitation le cerveau humain: inspiration du fonctionnement élémentaire du système nerveux
- Ont été très utilisés dans les années 80 et début 90 puis leur popularité a diminué.



## Avantages

- Récente résurgence: technique State-of-the-art pour beaucoup d'applications
- Pas de modélisation analytique
- Adaptabilité aux données
- Architecture parallèle



## Inconvénients

- Boîte noire: apprentissage d'une fonction complexe incompréhensible
- Sur-apprentissage

# Un neurone

Un neurone est une unité de calculs en 3 phases:

- lit l'information décrite par ses entrées  
(  $\mathbf{x} = [1, x_1 \dots, x_n]^T$  )

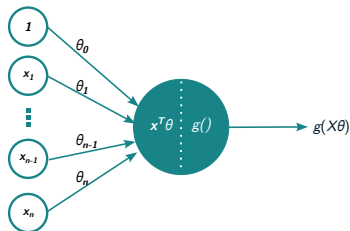
- combine linéairement les entrées

$$z(\mathbf{x}) = \mathbf{x}^T \boldsymbol{\theta} = \theta_0 + \sum_{i=1}^n \theta_i x_i$$

- paramètres:  $\boldsymbol{\theta} = [\theta_0, \dots, \theta_n]^T$

- passé le résultat dans une fonction linéaire

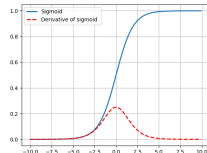
- $a(\mathbf{x}) = g(z(\mathbf{x})) = g(\mathbf{x}^T \boldsymbol{\theta})$
- $g(\cdot)$ : fonction d'activation (sigmoïde, ReLU)



# Différents choix de la fonction d'activation

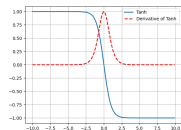
**Sigmoïde:**  $g(x) = \sigma(x) = \frac{1}{1+\exp^{-x}}$

- Régression logistique
- rend les valeurs entre 0 et 1
- toujours positif, borné, monotone croissante



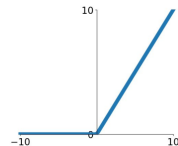
**Tangente hyperbolique**  $g(x) = \tanh(x) = \frac{\exp^{2x} - 1}{\exp^{2x} + 1}$

- rend les valeurs entre -1 et 1
- positif et négatif, borné, monotone croissante



**Rectified linear unite (ReLU):**  $g(x) = \text{ReLU}(x) = \max(0, x)$

- toujours positive ou nulle, borné à 0 pour les valeurs négatives, pas de borne supérieure
- monotone croissante, tend à donner des neurones avec des activations parcimonieuses (*sparse*) i.e une large gamme de valeurs d'entrée (toutes celles négatives) sera mise à 0.

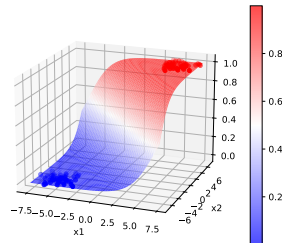
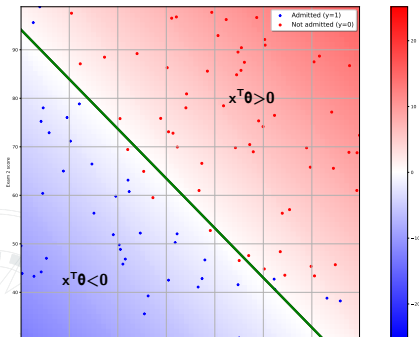




# Capacité d'un neurone

Un neurone peut:

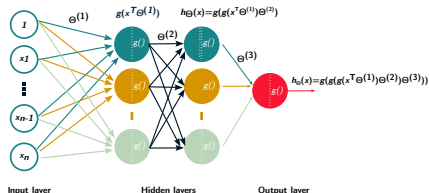
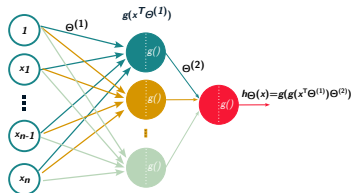
- résoudre uniquement de **classification binaire**
  - ▶ En utilisant la fonction d'activation sigmoïde (bornée entre 0 et 1), un neurone peut être interprété comme un estimateur de la probabilité de  $h_{\theta}(\mathbf{x}) = a(\mathbf{x}) = \mathbb{P}\{y = 1|\mathbf{x}\}, y \in \{0, 1\}$
  - ▶ connu comme le **classifieur de régression logistique**
  - ▶ si  $h_{\theta}(\mathbf{x}) > 0.5$  alors la classe 1 est prédite, sinon c'est la classe 0.
- La fonction d'activation  $\tanh$  peut mener à un classification binaire équivalente.



# Réseau de neurones avec des couches cachées

Le **réseau de neurones entièrement connectés** permet de résoudre des problèmes avec des frontières non-linéaires:

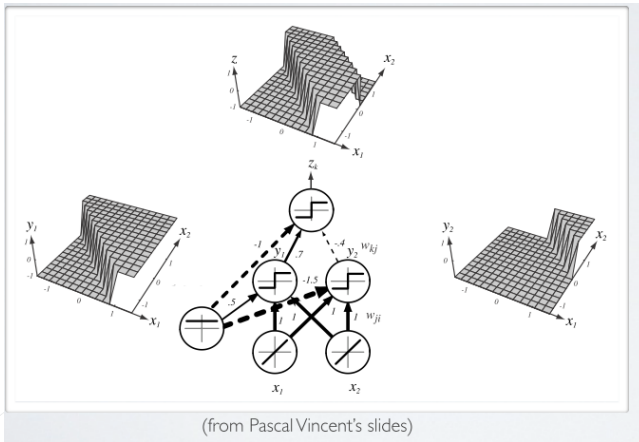
- par empilement des neurones entièrement connectés dans une couche cachée
- en rajoutant des couches cachées



Tous les neurones de la couche cachée sont connectés à un seul neurone de sortie (i.e: seule une **classification binaire** est possible)

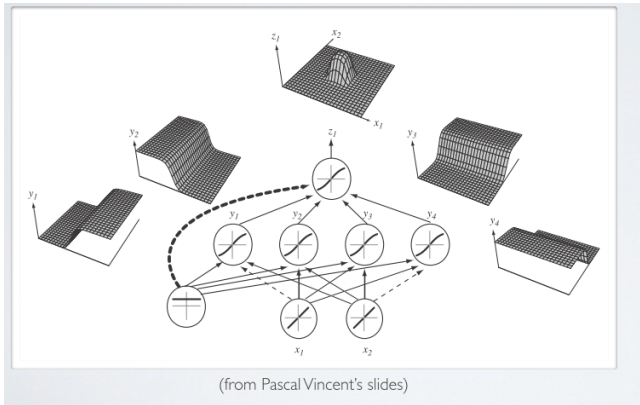
# Capacité d'un réseau de neurones à une couche cachée (deux neurones)

Réseaux de neurones



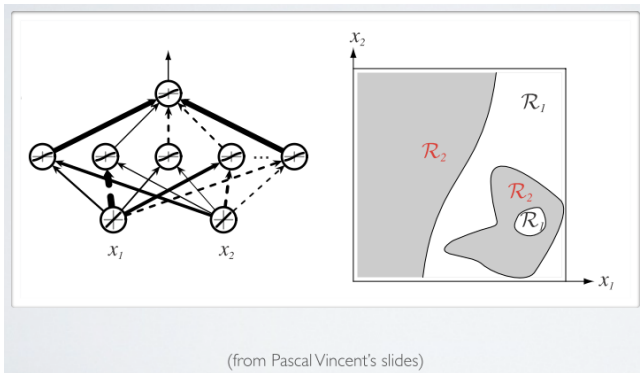
# Capacité d'un réseau de neurones à une couche cachée (plus de deux neurones)

## Réseaux de neurones



# Capacité d'un réseau de neurones à une couche cachée (plus de deux neurones)

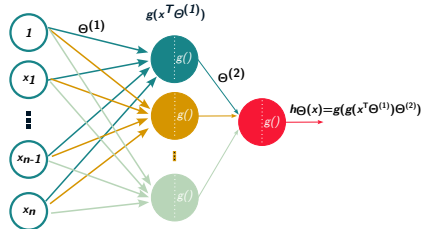
Réseaux de neurones



Il est prouvé<sup>1</sup> théoriquement que si on choisit un nombre suffisant de neurones, n'importe quelle fonction peut-être approximée (approximateur universel).

Kurt Hornik, "Approximation capabilities of multilayer feedforward networks", Neural Networks, vol.4, no.2, 1991, p.251–257 (DOI 10.1016/0893-6080(91)90009-T)

# Réseau de neurones avec une couche cachée



$\Theta_{i,j}^{(1)}$ : poids entre le neurone  $i$  de la couche cachée  $l = 1$  et l'entrée  $j$

- chaque neurone de la couche cachée a

- une préactivation  $\mathbf{z}(\mathbf{x}) = (z(\mathbf{x})_i)_{i=1,\dots}$  avec  $z(\mathbf{x})_i = \sum_j \Theta_{i,j}^{(1)} x_j$

- une activation  $\mathbf{a}(\mathbf{x}) = (a(\mathbf{x})_i)_{i=1,\dots} = (g(z(\mathbf{x})_i))_{i=1,\dots}$

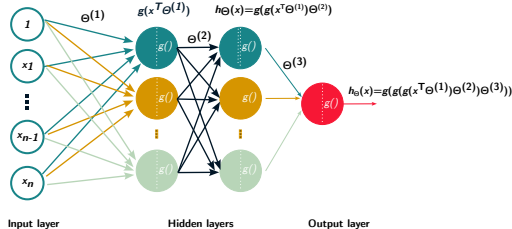
- Matriciellement, c'est plus facile à écrire

- la préactivation de la couche cachée:  $\mathbf{z}^{(1)}(\mathbf{x}) = \Theta^{(1)} \mathbf{x}^T$

- l'activation de la couche cachée:  $\mathbf{a}^{(1)}(\mathbf{x}) = g(\mathbf{z}^{(1)}(\mathbf{x}))$

- Pour la couche de sortie:  $h_{\Theta}(\mathbf{x}) = a^{(2)}(\mathbf{x}) = g(\mathbf{z}^{(2)}(\mathbf{x})) = g(\Theta^{(2)} \mathbf{a}^{(1)}(\mathbf{x}))$

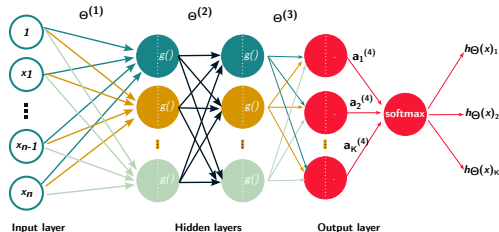
# Réseau de neurones avec $L$ couches cachées



- $l \in \{1, \dots, L\}$  est l'indice des couches cachées du réseau
- Les neurones de la couche  $l$  sont connectés à tous les neurones de la couche  $(l - 1)$
- $\Theta_{i,j}^{(l)}$ : poids entre neurones  $i$  de la couche  $l$  et des neurones  $j$  de la couche  $l - 1$
- Préactivation: fonctions d'activation de la couche  $l$  :  $z^{(l)}(x) = \Theta^{(l)} a^{(l-1)}(x)$
- Activation:  $h^{(l)}(x) = g(z^{(l)}(x))$

# Multiclass classification

Pour faire de la classification multiclasse



- Définir  $K$  neurones de sorties
- Choisir la fonction d'activation: **softmax()**

$$\text{softmax}(\mathbf{a}) = \left[ \frac{\exp(a_1)}{\sum_k \exp(a_k)}, \dots, \frac{\exp(a_k)}{\sum_k \exp(a_k)} \right]^T$$

- ▶ avec  $(\mathbf{a}_k^{(L+1)}) \forall k$  les préactivations de la dernière couche
- ▶ i.e. estimer la probabilité conditionnelle  $\mathbb{P}\{y = k | \mathbf{x}\}, k \in \{1, \dots, K\}$



# Apprentissage des réseaux de neurones

A partir d'un **ensemble d'apprentissage** ( $E$ ), déterminer les meilleurs paramètres du réseau de neurones qu'on cherche à entraîner pour une tâche  $T$ .

- l'ensemble des paramètres est:  $\Theta = \{\Theta^{(1)}, \dots, \Theta^{(L+1)}\}$
- le **modèle** du réseau de neurones multicouches est composé de **fonctions imbriquées** menant à une fonction à n'importe **quelle complexité**

$$h_{\theta}(x) = g \left[ \Theta^{(L)} \dots g \left[ \Theta^{(3)} g \left[ \Theta^{(2)} g \left[ \Theta^{(1)} \mathbf{x}^T \right] \right] \right] \dots \right]$$

**Trouver  $\Theta$ :**

- minimisant une mesure de performance  $P$  sous la forme d'une fonction de coût  $J(\Theta)$
- par descente de gradient
  - ▶  $\Theta_{ij}^{(l)} \leftarrow \Theta_{ij}^{(l)} - \alpha \frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$
  - ▶ estimation des gradients par rapport à tous les paramètres
  - ▶ **algorithme de rétropropagation** (*backpropagation*)

# Fonction de coût

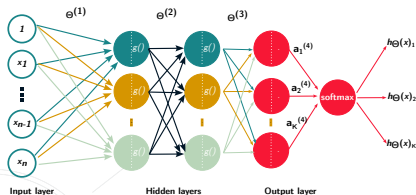
- L'apprentissage va consister à minimiser une fonction dépendante des paramètres et des données d'apprentissage

► minimisation de la fonction de coût  $\arg \min_{\theta} J(\theta)$

► Par exemple: cross-entropie pour la tâche  $T$  de classification

$$J(\theta) = \frac{1}{m} \sum_{i=0}^{m-1} \sum_{k=1}^K \left[ -y_k^i \log(h_{\theta}(\mathbf{x}^i))_k - (1 - y_k^i) \log(1 - h_{\theta}(\mathbf{x}^i)_k) \right]$$

- La fonction de coût est changée car passage à un problème multiclasse: recodage de  $y^i \in \{1, \dots, K\}$  en un codage "one-hot-encoding"  $y_{\text{one-hot}}^i$



$$y_{\text{one-hot}}^i = (y_k^i)_{k \in \{1, \dots, K\}}$$

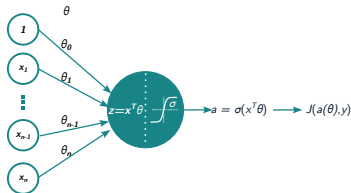
$$= \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \dots \text{ or } \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}.$$

# Backpropagation pour la régression logistique (1/2)

L'apprentissage va consister en deux passes:

## Feedforward pass

- calcul des (pré-)activations et du coût  $J(\theta)$  pour un couple  $(x, y)$

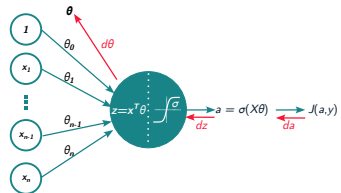


Objectif:

- calculer tous les  $\frac{\partial J(a, y, \theta)}{\partial \theta}$  pour la descente de gradient
- avec  $J(a, y, \theta) = -y \log(a) - (1 - y) \log(1 - a)$

## Backward pass

- calcul des gradients par rétropropagation des gradients pour un couple  $(x, y)$

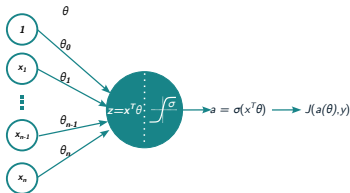


# Backpropagation pour la régression logistique (2/2)

L'apprentissage va consister en deux passes:

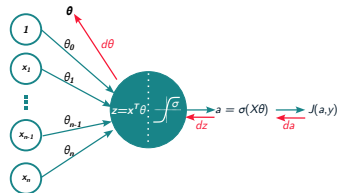
## Feedforward pass

- calcul des (pré-)activations et du coût  $J(\theta)$  pour un couple  $(x, y)$



## Backward pass

- calcul des gradients par rétropropagation des gradients pour un couple  $(x, y)$



Utilisation de la règle de dérivation en chaîne

$$da = \frac{\partial J(a, y)}{\partial a} = -\frac{y}{a} + \frac{1-y}{1-a} = \frac{-(1-a)y + a(1-y)}{a(1-a)} = \frac{a-y}{a(1-a)}$$

$$dz = \frac{\partial J(a, y)}{\partial z} = \frac{\partial J(a, y)}{\partial a} \frac{\partial a}{\partial z} = da \cdot \sigma'(z) = \frac{a-y}{a(1-a)} \sigma(z)(1-\sigma(z)) = a-y$$

$$d\theta = \frac{\partial J(a, y)}{\partial \theta} = dz \cdot \frac{\partial z}{\partial \theta} = (a-y) \cdot x$$

# Apprentissage des réseaux de neurones (1/1)

Exemple du tp3: on reprend le pb de reconnaissance des chiffres manuscrits

7	9	6	5	8	7	4	4	1	0
0	7	3	3	2	4	8	4	5	7
6	6	3	2	9	2	3	3	2	6
1	3	7	1	5	6	5	2	4	4
7	0	9	2	7	5	8	9	5	4
4	6	6	5	0	2	1	3	6	9
8	5	1	8	9	7	8	7	3	6
1	0	2	8	2	5	0	5	1	5
6	7	8	2	5	3	9	7	0	0
7	9	3	9	8	5	7	2	9	8

- $i$ : indice des  $m(= 5000)$  données d'apprentissage
- $X$ : matrice (ici à  $n = 20 \cdot 20 = 400$  colonnes) des descripteurs, input variable, feature
- $y$ : variable cible, target, output variable dont les valeurs sont entre 0 et 9

On fournit les poids d'un réseau de neurones:

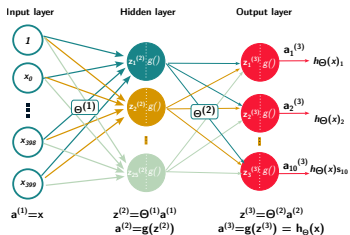
- prenant en 400 données d'entrée
- ayant une couche cachée de 25 neurones
- ayant 10 neurones de sorties
- $\Theta^{(1)}$  la matrice 25x401 liant les 25 neurones à toutes les entrées
- $\Theta^{(2)}$  la matrice 10x26 liant les 10 neurones de sortie à toutes les sorties des neurones de la couche cachée
- $y$  doit être recodé en "one-hot"

# Apprentissage des réseaux de neurones (1/2)

L'apprentissage va consister en deux passes:

## Feedforward pass

- calcul des (pré-)activations et du coût  $J(\theta)$

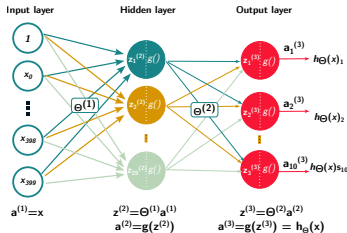


# Apprentissage des réseaux de neurones (2/2)

L'apprentissage va consister en deux passes:

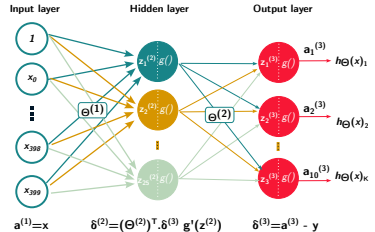
## Feedforward pass

- calcul des (pré-)activations et du coût  $J(\theta)$



## Backward pass

- calcul des gradients par rétropropagation des gradients



# Conclusions

- Un neurone
  - ▶ = régression logistique
  - ▶ classification binaire et linéaire
- Réseau de neurones
  - ▶ classification multiclasse et non-linéaire
  - ▶ apprentissage d'une fonction imbriquée (la sortie d'un neurone est l'entrée d'un neurone)
  - ▶ dont la complexité dépend du nombre de couches et du nombre de neurones (approximateur universel)

## Apprentissage

- à partir des données
- avec une fonction de coût  $J$  reflétant la tâche  $T$  à accomplir
- à minimiser par rétropropagation des gradients
- pour trouver les paramètres optimaux
- paramètres en grand nombre

