# Apprentissage profond (RN-DEEP) Introduction à l'apprentissage profond

Olivier Pons et Arnaud Breloy \*
olivier.pons@lecnam.net
arnaud.breloy@lecnam.net
http://cedric.cnam.fr/vertigo/Cours/ml2/

Département Informatique Conservatoire National des Arts & Métiers, Paris, France

4 novembre 2024

\* Sur les versions précédentes de N. Thome, N. Audebert, C. Rambour

Introduction 1 / 80

#### Plan du cours

- 1 Introduction
- 2 Réseaux de neurones artificiel
- Rétropropagation du gradien
- 4 Optimisation des réseaux profonds
- 5 Implémentation de la rétropropagation
- 6 Activations, régularisations, initialisations

Introduction 2 / 80

#### Données massives

Vastes quantités de données non-structurées :

images, vidéos, sons, textes, signaux, etc.









BBC: 2,4m vidéos

Facebook: 350M images

100m caméras de surveillance

504m de tweets par jour

Besoin d'analyser, manipuler, ces données pour les reconnaître, classer, organiser, comprendre, interpréter, en générer ...

⇒ Sous quelle représentation?

Introduction 3 / 80

# Extraire du sens des signaux bas niveau

Problème : combler l'écart entre le signal et la sémantique.













Ce que l'humain perçoit vs.

ce que la machine perçoit

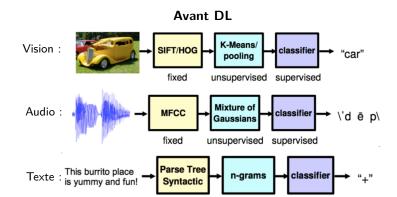
- Variations d'illumination
- Variations de points de vue
- Objets déformables
- Variance et diversité intra-classe
- **.**

⇒ Comment concevoir de "bonnes" représentations des données?

# Apprentissage statistique classique

#### Des données à la décison :

- Choix de caractéristiques descriptives hand-crafted
  - descripteurs images (SIFT/HOG), coefficients audio (MFCC), n-grammes/bag-of-words, etc.
  - nécessite des connaissances expertes du domaine
  - savoir a priori
  - Entraînement d'un modèle de décisionnel
  - Évaluation des performances



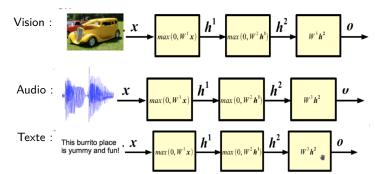
Introduction 4 / 80

# Apprentissage statistique classique

#### Des données à la décison :

- Choix de caractéristiques descriptives hand-crafted
  - descripteurs images (SIFT/HOG), coefficients audio (MFCC), n-grammes/bag-of-words, etc.
  - nécessite des connaissances expertes du domaine
  - savoir a priori
  - Entraînement d'un modèle de décisionnel
  - Évaluation des performances

Apprentissage profond ⇒ apprentissage de **représentations**Après DL



Introduction 5 / 80

## Historique des réseaux de neurones

- **1943**: Neurone artificiel McCulloch et Pitts 1943
- **1957**: Perceptron ROSENBLATT **1957**
- 1969: "Perceptrons An Introduction to Computational Geometry" MINSKY et PAPERT 2017(réedition)
- **1974**: Algorithme de rétropropagation WERBOS 1975
- 1980 : Neocognitron (premier réseau "profond") à poids partagés FUKUSHIMA 1980
- 1986 : Popularise la rétropropagation RUMELHART, HINTON et WILLIAMS 1986
- 1989 : LeNet-5 (premier réseau convolutif) LECUN et al. 1989
- 1989 : LSTM, résoud le problème de la disparition du gradient dans les RNN et permet de capturer des dépendances à long terme. HOCHREITER et SCHMIDHUBER 1997
- 2011/2012: DanNet/AlexNet (premiers modèles à l'état de l'art en reconnaissance d'image) CIREŞAN, MEIER et SCHMIDHUBER 2012; KRIZHEVSKY, SUTSKEVER et HINTON 2012

Introduction 6 / 80

## Historique récent des réseaux de neurones

- 2014 : Generative Adversarial Networks (GANs) GOODFELLOW et al. 2014 Entraînant deux réseaux de neurones (un générateur et un discriminateur) sont entraînés en opposition pour créer des données réalistes. Cela révolutionne la génération de données.
- 2016 : AlphaGo de DeepMind, premier programme à battre un champion du monde de Go SILVER et al. 2016 Il combine CNN, RNN et MCTS.
- 2017 : Transformer , "Attention Is All You Need" VASWANI et al. 2017 Popularise la notion d'attention qui révolutionne le traitement du langage naturel et sert de base à de nombreux modèles ultérieurs.
- 2018 : BERT (Bidirectional Encoder Representations from Transformers) DEVLIN et al. 2018 Modèle de langage basé sur des Transformers qui apprend les relations contextuelles entre les mots dans les deux directions
- 2019 : GPT-2 (Generative Pretrained Transformer 2) RADFORD et al. 2019 Modèle de langage auto-régressif basé sur les Transformers, entraîné pour prédire le mot suivant dans une séquence, capable de générer du texte cohérent et réaliste en se basant uniquement sur des données non supervisées.

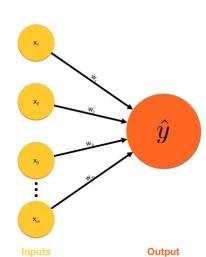
## Plan du cours

- 1 Introduction
- 2 Réseaux de neurones artificiels
- Rétropropagation du gradien
- 4 Optimisation des réseaux profonds
- 5 Implémentation de la rétropropagation
- 6 Activations, régularisations, initialisations

#### Neurone formel

- Modèle simplifié d'un neurone réel (1943) McCulloch et Pitts 1943
  - Entrée : vecteur  $\mathbf{x} \in \mathbb{R}^m$ ,  $\mathbf{x} = \{x_1, x_2, \dots, x_m\}$ ,
  - Poids w<sub>j</sub>, pondération de la connexion entre l'entrée j et la sortie,
  - Sortie  $\hat{y} \in \mathbb{R}$  (scalaire),  $\hat{y} = \sum_{j=1}^{m} w_j x_j$ .
- Expression matricielle :

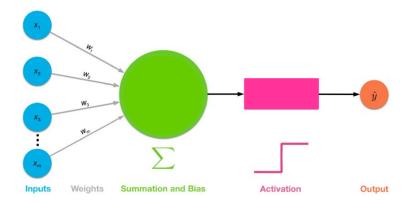
$$\hat{y} = \mathbf{w}^T \mathbf{x} = \begin{bmatrix} w_1 & w_2 & \cdots & w_m \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}$$



#### Neurone artificiel

- Fonction de  $\mathbf{x} \in \mathbb{R}^m$  vers  $\hat{\mathbf{y}} \in \mathbb{R}$  :
  - Fonction affine (linéaire plus translation par un biais b scalaire) :  $s = \mathbf{w}^T \mathbf{x} + b$
  - $oxed{2}$  Fonction d'activation non-linéaire  $\phi:\mathbb{R} 
    ightarrow \mathbb{R}: \hat{y} = \phi(s)$
- Expression du neurone artificiel "complet"

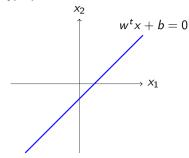
$$\hat{\mathbf{y}} = \phi \left( \mathbf{w}^\mathsf{T} \mathbf{x} + \mathbf{b} \right)$$



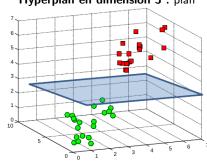
## Partie linéaire

- Projection affine :  $s = \mathbf{w}^{\top} \mathbf{x} + b = \sum_{i=1}^{m} w_i x_i + b$ 
  - ${f w}$  : vecteur à un hyperplan de  ${\Bbb R}^m\Rightarrow s=0$  définit une frontière linéaire
  - bias b définit un écart par rapport à la position de l'hyperplan

#### Hyperplan en dimension 2 : droite



#### Hyperplan en dimension 3 : plan

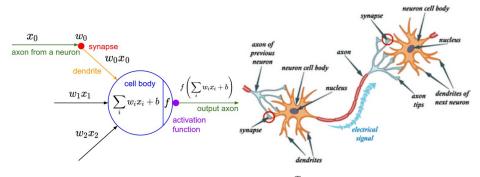


## Fonction d'activation non-linéaire

$$\hat{\mathbf{y}} = \sigma \left( \mathbf{w}^\mathsf{T} \mathbf{x} + \mathbf{b} \right)$$

- Fonction identité :  $\phi(x) = x$
- $\blacksquare \text{ Fonction de Heaviside (\'echelon)}: \phi(x) = \begin{cases} 1 & \text{six} \geq 0 \\ 0 & \text{sinon} \end{cases}$
- Sigmoïde :  $\sigma(x) = \frac{1}{1+e^{-x}}$
- Tangente hyperbolique :  $\tanh x = \frac{e^x e^{-x}}{e^x + e^{-x}}$
- Rectified Linear Unit (ReLU) :  $\phi(x) = \begin{cases} x & \text{six} \ge 0 \\ 0 & \text{sinon} \end{cases}$

# Lien avec les neurones biologiques

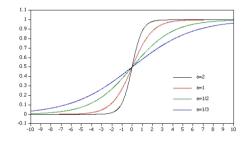


- Neurone formel, activation Heaviside  $H: \hat{y} = H(\mathbf{w}^{\top}\mathbf{x} + b)$ 
  - $\hat{\mathbf{y}} = 1 \text{ (neurone actif)} \Leftrightarrow \mathbf{w}^{\top} \mathbf{x} \ge -b$  $\hat{\mathbf{y}} = 0 \text{ (neurone inactif)} \Leftrightarrow \mathbf{w}^{\top} \mathbf{x} < -b$

■ Neurones biologiques : la sortie est active ⇔ la somme des entrées pondérées par le poids des connexions est synaptique est supérieure à un seuil.

# Activation sigmoïde

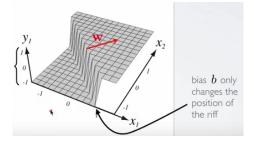
- Équation du neurone artificiel :  $\hat{y} = \phi(\mathbf{w}^{\top}\mathbf{x} + b)$
- Sigmoïde :  $\phi(z) = \sigma_a(z) = (1 + e^{-az})^{-1}$



- $\blacksquare$   $a \uparrow$ : se rapproche de l'échelon de Heaviside  $(a \to \infty)$
- Sigmoïde : deux régimes
  - $\mathbf{x} \approx 0 \rightarrow \text{régime linéaire } (\sigma(x) \approx ax)$
  - $\blacksquare \ x << 0 \ \text{ou} \ x >> 0 \ \rightarrow \ \text{r\'egime} \ \text{de saturation} \ \big(\sigma(x) \approx \pm 1\big)$

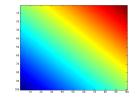
# Application à la classification binaire

- $\blacksquare$  Chaque entrée x appartient soit à la classe 0, soit à la classe 1.
- Sortie du neurone artificiel (sigmoïde) :  $\hat{y} = \frac{1}{1 + e^{-a(\mathbf{w}^{\top}\mathbf{x} + b)}}$
- Interprétation probabiliste  $\Rightarrow \hat{y} \sim P(1|\mathbf{x})$ 
  - lacksquare L'entrée  ${f x}$  est classée comme classe 1 si  $P(1|{f x})>0.5$   $\Leftrightarrow$   ${f w}_{oldsymbol{\bot}}^{oldsymbol{\top}}{f x}+b>0$
  - L'entrée x est classée comme classe 0 si  $P(1|\mathbf{x}) < 0.5 \Leftrightarrow \mathbf{w}^{\top}\mathbf{x} + b < 0$ 
    - $\Rightarrow \operatorname{signe}(\mathbf{w}^{\top}\mathbf{x} + b) : \mathsf{la} \text{ frontière de classification est linéaire dans l'espace d'entrée}\,!$



■ En 2D 
$$(m = 2)$$
,  $\mathbf{x} = \{x_1, x_2\} \in [-5; 5]^2$ 

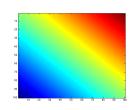
- Poids fixés  $\mathbf{w} = [1; 1]$  et b = -2
- Résultat de la fonction affine :  $s = \mathbf{w}^{\top}\mathbf{x} + b = x_1 + x_2 2$



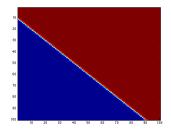
$$lacksquare$$
 Activation non-linéaire, sigmoïde  $a=10: \hat{y}=\left(1+e^{-a(\mathbf{w}^{ op}\mathbf{x}+b)}\right)^{-1}$ 

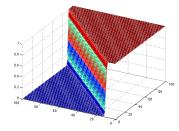
■ En 2D 
$$(m = 2)$$
,  $\mathbf{x} = \{x_1, x_2\} \in [-5; 5]^2$ 

- Poids fixés  $\mathbf{w} = [1; 1]$  et b = -2
- $\blacksquare$  Résultat de la fonction affine :  $s = \mathbf{w}^{\top} \mathbf{x} + b = x_1 + x_2 2$



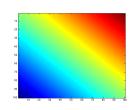
lacksquare Activation non-linéaire, sigmoïde  $a=10: \hat{y}=\left(1+e^{-a(\mathbf{w}^{ op}\mathbf{x}+b)}
ight)^{-1}$ 



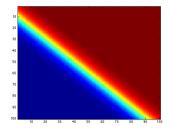


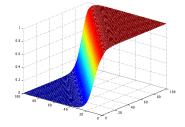
■ En 2D 
$$(m = 2)$$
,  $\mathbf{x} = \{x_1, x_2\} \in [-5; 5]^2$ 

- Poids fixés  $\mathbf{w} = [1; 1]$  et b = -2
- Résultat de la fonction affine :  $s = \mathbf{w}^{\top}\mathbf{x} + b = x_1 + x_2 2$



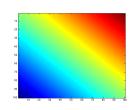
lacksquare Activation non-linéaire, sigmoïde a=1 :  $\hat{y}=\left(1+e^{-a(\mathbf{w}^{ op}\mathbf{x}+b)}
ight)^{-1}$ 



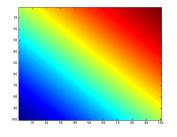


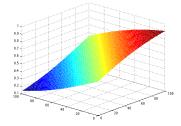
■ En 2D 
$$(m = 2)$$
,  $\mathbf{x} = \{x_1, x_2\} \in [-5; 5]^2$ 

- Poids fixés  $\mathbf{w} = [1; 1]$  et b = -2
- $\blacksquare$  Résultat de la fonction affine :  $\mathbf{s} = \mathbf{w}^{\top}\mathbf{x} + \mathbf{b} = \mathbf{x}_1 + \mathbf{x}_2 2$



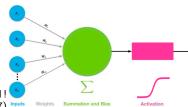
lacksquare Activation non-linéaire, sigmoïde  $a=0.1: \hat{y}=\left(1+e^{-a(\mathbf{w}^{ op}\mathbf{x}+b)}
ight)^{-1}$ 





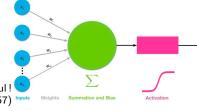
## Du neurone formel au réseau de neurones

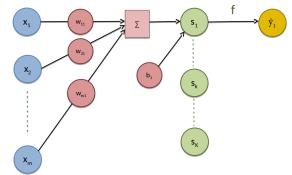
- Neurone artificiel :
  - $\blacksquare$  Une seule sortie scalaire  $\hat{y}$
  - Prontière linéaire pour la classification binaire
- Sortie scalaire unique : expressivité limitée pour la plupart des tâches
  - Comment gérer la classification multi-classe?
  - ⇒ utiliser plusieurs neurones de sortie plutôt qu'un seul!
  - $\Rightarrow$  modèle neuronal dit **Perceptron** (ROSENBLATT 1957)



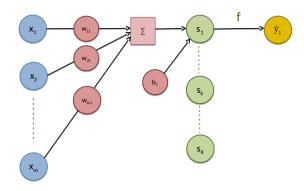
## Du neurone formel au réseau de neurones

- Neurone artificiel :
  - $\blacksquare$  Une seule sortie scalaire  $\hat{y}$
  - Prontière linéaire pour la classification binaire
- Sortie scalaire unique : expressivité limitée pour la plupart des tâches
  - Comment gérer la classification multi-classe?
  - ⇒ utiliser plusieurs neurones de sortie plutôt qu'un seul!
  - ⇒ modèle neuronal dit **Perceptron** (ROSENBLATT 1957)



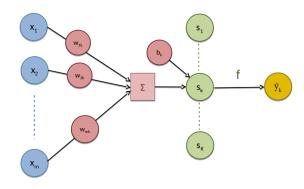


## Perceptron et classification multi-classe



- Entrée x in  $\mathbb{R}^m$  ("couche"d'entrée)
- Chaque sortie  $\hat{y_1}$  est un neurone artificiel ("couche"de sortie). Par exemple, pour la sortie  $\hat{y_1}$ :
  - Transformation affine :  $s_1 = \mathbf{w_1}^T \mathbf{x} + b_1$
  - lacksquare Activation non-linéaire  $\sigma: \hat{y_1} = \sigma(s_1)$
- $\blacksquare$  Paramètres de la transformation linéaire la sortie  $\hat{y_1}$ :
  - $\blacksquare$  poids  $\mathbf{w}_1 = \{w_{1,1}, w_{2,1}, \dots, w_{m,1}\} \in \mathbb{R}^m$
  - $\blacksquare$  biais  $b_1 \in \mathbb{R}$

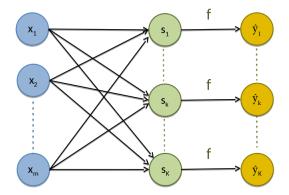
## Perceptron et classification multi-classe



- $\blacksquare$  Entrée x in  $\mathbb{R}^m$
- Chaque sortie  $\hat{y_k}$  est un neurone artificiel. Pour chaque sortie  $\hat{y_k}$ :
  - Transformation affine :  $s_k = \mathbf{w_k}^T \mathbf{x} + b_k$
  - Activation non-linéaire  $\sigma: \hat{y_k} = \sigma(s_k)$
- Paramètres de la transformation linéaire la sortie  $\hat{y_k}$ :
  - lacksquare poids  $\mathbf{w_k} = \left\{ w_{1,k}, w_{2,k}, \ldots, w_{m,k} 
    ight\} \in \mathbb{R}^m$
  - lacksquare biais  $b_k \in \mathbb{R}$

# Perceptron en résumé

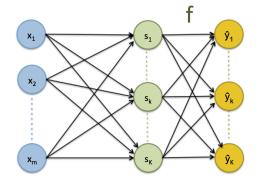
- Entrée  $\mathbf{x} \in \mathbb{R}^m$   $(1 \times m)$ , sortie  $\hat{y}$  : concaténation de k neurones
- $\blacksquare$  Transformation affine (  $\sim$  multiplication matricielle) : s = xW + b
  - lacksquare avec f W une matrice de poids de dimensions m imes k les colonnes sont les vecteurs  $f w_k$ ,
  - $\blacksquare$  et **b** le vecteur de biais dimension  $1 \times k$ .
- lacktriangle Activation non-linéaire élément par élément (pointwise) :  $\hat{\mathbf{y}} = \sigma(\mathbf{s})$



## Application à la classification multi-classe

- Comment généraliser la sigmoïde au cas à plusieurs classes?
  - lacksquare on cherche f telle que  $f(s_j) = P(j|\mathbf{x})$  (interprétation probabiliste)
  - $\blacksquare$  contrainte :  $\sum_{j=1}^{k} f(s_j) = 1$ .
- Activation softmax :

$$\hat{y}_i = f(s_i) = \frac{e^{s_i}}{\sum_{j=1}^k e^{s_j}}$$

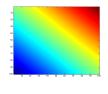


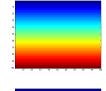
⇒ on appelle ce modèle la régression logistique.

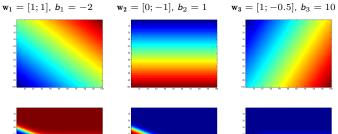
# Exemple jouet en 2D pour trois classes

■ 
$$\mathbf{x} = \{x_1, x_2\} \in [-5; 5] \times [-5; 5], \quad \hat{y} \in \{0, 1, 2\} \text{ (3 classes)}$$

Transformation affine:  $\mathbf{s}_k = \mathbf{w}_k^T \mathbf{x} + b_k$ 







Softmax:  $P(k|\mathbf{x},\mathbf{W})$ 





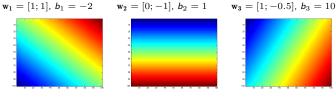


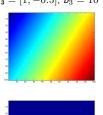
# Exemple jouet en 2D pour trois classes

$$\mathbf{x} = \{x_1, x_2\} \in [-5; 5] \times [-5; 5], \quad \hat{y} \in \{0, 1, 2\} \text{ (3 classes)}$$

Transformation affine:  $\mathbf{s}_k = \mathbf{w}_k^T \mathbf{x} + b_k$ 







Softmax:  $P(k|\mathbf{x},\mathbf{W})$ 

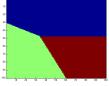






Classe prédite :

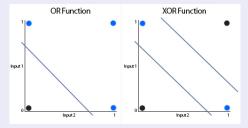
$$k^* = \underset{k}{\operatorname{arg max}} P(k|\mathbf{x}, \mathbf{W})$$



## Limites du perceptron

## Le problème du XOR

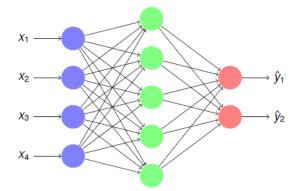
- Régression logistique (LR) : perceptron à une couche d'entrée et une couche de sortie
  - LR ne peut exprimer que des frontières de décisions linéaires
- XOR: (NON 1 ET 2) OU (NON 2 ET 1)
  - La frontière optimale pour XOR est non-linéaire.



 $\Rightarrow$  constat pessimiste : le perceptron échoue sur un problème en apparence très simple...

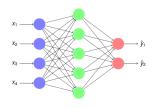
## Au-delà des frontières linéaires

- La régression logistique est limitée aux frontières linéaires.LR : limited to linear boundaries
  - ⇒ Solution : ajouter une couche de neurones!
- Entrée :  $\mathbf{x} \in \mathbb{R}^m$ , ici par exemple m = 4.
- Sortie :  $\hat{\mathbf{y}} \in \mathbb{R}^k$  (k classes), par ex. k = 2.
- **Couche "cachée"** de taille  $I: \mathbf{h} \in \mathbb{R}^I$ , par ex. I = 5



# Perceptron multi-couche

- Couche cachée h: projection de l'entrée x vers un nouvel espace  $\mathbb{R}^I$ 
  - h est une représentation intermédiaire de x qui sera ensuite utilisée pour la classification  $\hat{\mathbf{y}}:\mathbf{h}=\phi\left(\mathbf{W}_{s}\mathbf{h}+\mathbf{b}_{s}\right)$  par la régression logistique.
- Réseau de neurones avec au moins une couche cachée : PMC, perceptron multi-couche (multi-layer perceptron, MLP)

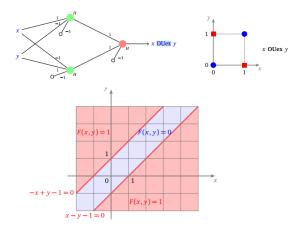


#### Attention!

La non-linéarité de la fonction d'activation est indispensable pour le perceptron multi-couche :

# Garanties théoriques

■ Le (non) problème XOR avec plusieurs couches



Note : Le neurone du haut réalise « $\times$  ET non y», le neurone du bas «non  $\times$  ET y» et celui de sortie l'opération «OU».

# Garanties théoriques

- Quelle est l'expressivité des perceptrons multi-couche?
  - Peut-on approcher n'importe quelle fonction à l'aide d'un PMC?

# Théorème d'approximation universelle

Soit f une fonction continue sur (des compacts de)  $\mathbb{R}^d$ . Alors, pour tout  $\epsilon>0$ , il existe un perceptron  $\hat{f}:x\to \mathbf{W}_\mathrm{s}\phi(\mathbf{W}\cdot x+b)$  multi-couche à une couche cachée de nombre de neurones fini qui est une approximation de f à  $\epsilon$  près, c'est-à-dire :

$$\max_{x \in \mathbb{R}^d} \lVert f(x) - \hat{f}(x) \rVert < \epsilon$$

si  $\phi$  est une fonction d'activation non-polynomiale.

- Démontré pour les fonctions d'activation sigmoïde en 1989 CYBENKO 1989
- Étendu à l'ensemble des perceptrons multi-couche en 1991 HORNIK 1991
- D'autres résultats existent : profondeur arbitraire, largeur arbitraire, etc.

#### Limitations

- Une seule couche cachée suffit mais le nombre de neurones n'est pas borné
- La démonstration prouve l'existence mais ne donne pas de méthode pour déterminer les paramètres  $W_s$ , W et b.

# Garanties théoriques

- Quelle est l'expressivité des perceptrons multi-couche?
  - Peut-on approcher n'importe quelle fonction à l'aide d'un PMC?

## Théorème d'approximation universelle

Soit f une fonction continue sur (des compacts de)  $\mathbb{R}^d$ . Alors, pour tout  $\epsilon > 0$ , il existe un perceptron  $\hat{f}: x \to \mathbf{W_s} \phi(\mathbf{W} \cdot x + b)$  multi-couche à une couche cachée de nombre de neurones fini qui est une approximation de f à  $\epsilon$  près, c'est-à-dire :

$$\max_{x \in \mathbb{R}^d} ||f(x) - \hat{f}(x)|| < \epsilon$$

si  $\phi$  est une fonction d'activation non-polynomiale.

- Démontré pour les fonctions d'activation sigmoïde en 1989 CYBENKO 1989
- Étendu à l'ensemble des perceptrons multi-couche en 1991 HORNIK 1991
- D'autres résultats existent : profondeur arbitraire, largeur arbitraire, etc.

#### Limitations

- Une seule couche cachée suffit mais le nombre de neurones n'est pas borné.
- La démonstration prouve l'existence mais ne donne pas de méthode pour déterminer les paramètres  $W_s$ , W et b.

# Idée de la preuve

#### Fonctions marches:









Heavmarche decalée à droiteiside

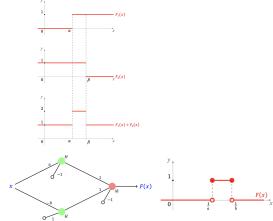
marche à l'envers décalée vers la droite

•••

# Idée de la preuve

#### Fonctions créneaux :

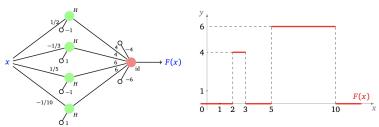
Réaliser un « créneau », en additionnant une marche à l'endroit et une marche à l'envers.



Pour une marche plus haute varier les poids du neurone de sortie d'un facteur k.

# Idée de la preuve

#### Fonctions en escaliers :

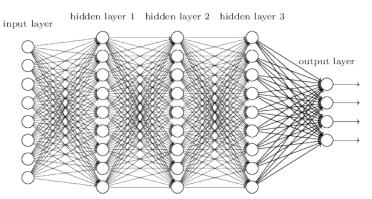


Thèoreme de heine (continuite uniforme) + approximation uniforme

#### Réseaux de neurones profonds

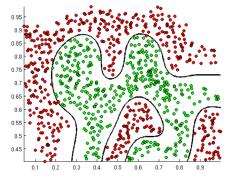
- Augmenter le nombre de couche cachées : réseaux de neurones profonds (deep neural networks)
- lacktriangle Chaque couche lacktriangle projette les activations de la couche précédente lacktriangle dans un nouvel espace intermédiaire
- ⇒ construction incrémentale d'un espace intermédiaire de représentation utile pour la tâche visée

#### l'apprentissage de représentation est la base de l'apprentissage profond!



#### En résumé

- Neurone artificiel : représentation grossière d'un neurone biologique
  - Poids des connexions synaptiques
- Fonction de transfert non-linéaire
- Perceptron : une matrice de poids projetant la couche d'entrée vers la couche de sortie
  - limité à des frontières de décision linéaires
- Réseaux de neurones profonds : perceptrons multi-couche applicables aux problèmes dont les frontières de décision sont non-linéaires



 $\Rightarrow$  comment trouver les paramètres du modèle (poids des connexions) ?  $\to$  apprentissage supervisé

#### Plan du cours

- 1 Introduction
- 2 Réseaux de neurones artificiels
- 3 Rétropropagation du gradient
- 4 Optimisation des réseaux profonds
- 5 Implémentation de la rétropropagation
- 6 Activations, régularisations, initialisations

# Algorithme du gradient

Soit  $f: \mathbb{R}^d \to \mathbb{R}$  une fonction réelle à d variables différentiable.

■ On note  $\nabla f(\mathbf{x})$  le gradient de f au point  $\mathbf{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d$ .

$$\nabla f = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \cdots, \frac{\partial f}{\partial x_d}\right)$$

- Si  $\nabla f(\mathbf{x}^*) = 0$ , alors  $\mathbf{x}^*$  est un extremum local de f.
  - $\blacksquare$  Si f est convexe ou concave, alors c'est un extremum global.
- lacktriangle On cherche à *minimiser f* , c'est-à-dire trouver  $\mathbf{x}^*$  tel que  $f(\mathbf{x}^*) < f(\mathbf{x})$  pour  $\mathbf{x} \in \mathbb{R}^d$ .

## Algorithme du gradient CAUCHY 1847

Soit  $\mathbf{x}^{(0)} \in \mathbb{R}^d$  un point de départ et  $\epsilon > 0$ . On définit la suite  $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \ldots$  telle que :

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \alpha_t \nabla f(\mathbf{x}^{(t+1)})$$

Si  $\|\nabla f(\mathbf{x}^{(t)})\| \leq \epsilon$ , on s'arrête.

lacktriangledown  $lpha_t$  peut être obtenu par divers moyens, comme une recherche linéaire. Dans notre cas, on fixera  $lpha_t=lpha>0$  constant.

# Optimisation du multi-layer perceptron (MLP)

- $\blacksquare$  Entrée  $\mathbf{x} \in \mathbb{R}^m$ , sortie  $\mathbf{y} \in \mathbb{R}^p$
- Paramètres w (poids et biais du modèle)
- lacksquare Le MLP est un modèle paramétrique  $\mathbf{x}\Rightarrow\mathbf{y}:\mathit{f}_{\mathbf{w}}(\mathbf{x}_{i})=\hat{\mathbf{y}_{i}}$
- lacksquare Apprentissage **supervisé** : jeu d'entraînement annoté  $\mathcal{A} = \{(\mathbf{x}_i, \mathbf{y}_i^*)\}_{i \in \{1, 2, ..., N\}}$ 
  - lacksquare Fonction de coût  $\mathcal{L} = \sum_{i=1}^N \ell(\hat{\mathbf{y}_i}, \mathbf{y}_i^*)$  entre la prédiction du modèle et l'annotation
- lacksquare Hypothèses : les paramètres  $\mathbf{w} \in \mathbb{R}^d$  sont des réels,  $\mathcal L$  est différentiable.
- Gradient  $\nabla_{\mathbf{w}} = \frac{\partial \mathcal{L}}{\partial \mathbf{w}}$ : direction de plus forte pente dans l'espace des paramètres permettant de faire décroître le coût  $\mathcal{L}$

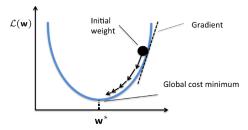


# Algorithme d'optimisation

Considérons le vecteur des paramètres  $\mathbf{w} \in \mathbb{R}^d = (w_1, w_2, \dots, w_d)$ . Le gradient du coût  $\mathcal{L}$  par rapport à  $\mathbf{w}$  est le vecteur des dérivées partielles du coût par rapport à chaque paramètre  $w_i$ :

$$\nabla_{\mathbf{w}} \mathcal{L} = \begin{pmatrix} \frac{\partial \mathcal{L}(\hat{y}, \mathbf{y}^*; \mathbf{w})}{\partial w_1} & \frac{\partial \mathcal{L}(\hat{y}, \mathbf{y}^*; \mathbf{w})}{\partial w_2} & \cdots & \frac{\partial \mathcal{L}(\hat{y}, \mathbf{y}^*; \mathbf{w})}{\partial w_d} \end{pmatrix}$$

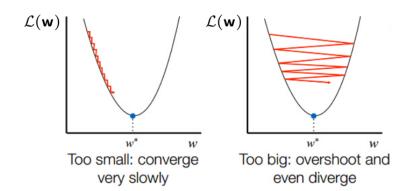
- Algorithme de la descente de gradient :
  - Initialisation aléatoire des paramètres w
  - $oxed{2}$  Mise à jour à l'itération  $t: \left| \mathbf{w}_i^{(t+1)} = \mathbf{w}_i^{(t)} \eta rac{\partial \mathcal{L}}{\partial \mathbf{w}_i} 
    ight|$
  - Répéter l'étape 2 jusqu'à convergence, par ex.  $\|\nabla_{\mathbf{w}}\mathcal{L}\|^2 \approx 0$  (le coût cesse de décroître).



# Pas d'apprentissage

Équation de mise à jour :  $\left|\mathbf{w}_{i}^{(t+1)}=\mathbf{w}_{i}^{(t)}-\eta \frac{\partial \mathcal{L}}{\partial \mathbf{w}_{i}}\right|$   $\eta$  pas d'apprentissage (learning rate)

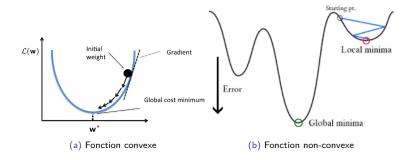
■ Peut-on garantir la convergence de l'algorithme vers un minimum?  $\Rightarrow$  seulement pour une valeur de  $\eta$  "bien choisie".



## Minimum local ou global

Équation de mise à jour : 
$$\mathbf{w}_i^{(t+1)} = \mathbf{w}_i^{(t)} - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{w}_i}$$

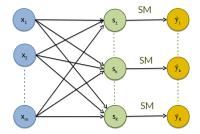
- Le minimum trouvé est-il le meilleur (minimum global)?
- $\Rightarrow$  garanti seulement si la fonction  $\mathcal{L}(w)$  à minimiser est convexe



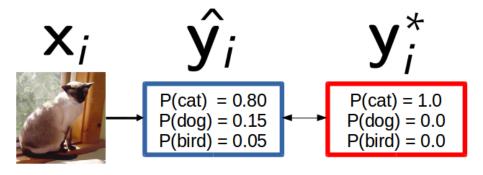
dans la quasi-totalité des cas, la fonction de coût utilisée n'est pas convexe par rapport aux paramètres du modèle...

# Fonctions de coût pour le perceptron

- Perceptron :  $\hat{\mathbf{y}} = \phi \left( \mathbf{x}_i \mathbf{W} + \mathbf{b} \right)$
- Régression :
  - Fonction de coût : L2 (MSE)  $\|\hat{\mathbf{y}} \mathbf{y}\|^2$ , L1 (MAE)  $|\hat{\mathbf{y}} \mathbf{y}|$ , etc.
- Classification multi-classe à K classes (régression logistique) :
  - $\mathbf{y} \in \{1; 2; \dots; K\}$
  - Activation softmax :  $\phi(\mathbf{s}_k) = P(k|\mathbf{x_i}) = \frac{e^{\mathbf{s}_k}}{\sum_{k'=1}^K e^{\mathbf{s}_{k'}}}$
  - $\mathbf{m} \ \hat{\mathbf{y}_i} = \underset{k}{\operatorname{arg\,max}} \ P(k|\mathbf{x}_i; \mathbf{W}, \mathbf{b})$
  - Fonction de coût :  $\ell_{0/1}(\hat{\mathbf{y}_i}, \mathbf{y}_i^*) = \begin{cases} 1 & \text{si } \hat{\mathbf{y}_i} \neq \mathbf{y}_i^* \\ 0 & \text{sinon} \end{cases}$  : **perte 0/1**
  - $\Rightarrow$  non-différentiable!



# Régression logistique : encodage des variables



- Vérité terrain (étiquettes de supervision)  $\mathbf{y}_i^* \in \{1, 2, \dots, K\}$
- Encodage *one hot* pour chaque étiquette :

$$y_{c,i}^* = egin{cases} 1 & ext{si } c ext{ est la classe de } \mathbf{x}_i \ 0 & ext{sinon} \end{cases}$$

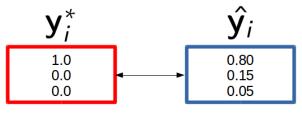
$$\mathbf{y}_{i}^{*} = \left(0, 0, \dots, \underbrace{1}_{c^{e} \text{composante}}, \dots, 0\right)$$

# Régression logistique : fonction de coût

- Fonction de coût : entropie croisée (*cross-entropy*, CE) :  $\ell_{CE}$
- $\blacksquare$   $\ell_{\textit{CE}}$  : divergence de Kullback-Leiber entre la distribution de probabilité  $y_i^*$  de la vérité terrain et la distribution prédite  $\widehat{y}_i$

$$\ell_{\textit{CE}}(\mathbf{y}_i^*, \widehat{\mathbf{y}}_i) = \mathrm{KL}(\mathbf{y}_i^*, \widehat{\mathbf{y}}_i) = -\sum_{c=1}^K \underbrace{y_{c,i}^*}_{\text{s suf pour } c^*} \log(\hat{y}_{c,i}) = -\log(\hat{y}_{c^*,i})$$

■ Attention! La divergence KL (et donc l'entropie croisée) est asymétrique :  $KL(\widehat{\mathbf{y}}_i,\mathbf{y}_i^*) \neq KL(\mathbf{y}_i^*,\widehat{\mathbf{y}}_i)$ 



$$KL(\mathbf{y_{i}^{*}}, \hat{\mathbf{y}_{i}}) = -log(\hat{y}_{c^{*},i}) = -log(0.8) \approx 0.22$$

# Entraînement de la régression logistique

- lacksquare  $\ell_{\it CE}$  est une borne supérieure de la perte  $\ell_{0/1}$
- lacksquare minimiser  $\ell_{\it CE} \Rightarrow$  minimiser la perte  $\ell_{0/1}$
- $\ell_{CE}$  est différentiable  $\Rightarrow$  descente de gradient!
  - $\blacksquare$   $\ell_{CE}$  est convexe mais seulement par rapport à y (pas par rapport à w)

#### Calcul du gradient

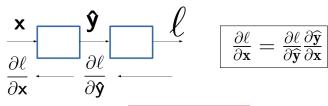
Descente de gradient :  $\mathbf{W}^{(t+1)} = \mathbf{W}^{(t)} - \eta \frac{\partial \mathcal{L}_{CE}}{\partial \mathbf{W}}$  et  $\mathbf{b}^{(t+1)} = \mathbf{b}^{(t)} - \eta \frac{\partial \mathcal{L}_{CE}}{\partial \mathbf{b}}$ 

- Principal problème : calculer  $\frac{\partial \mathcal{L}_{CE}}{\partial \mathbf{W}} = \frac{1}{N} \sum_{i=1}^{N} \frac{\partial \ell_{CE}}{\partial \mathbf{W}}$
- Propriété centrale : chain rule (dérivation des fonctions composées)

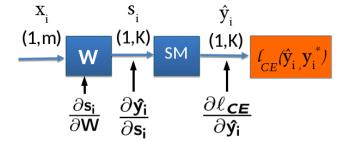
$$\frac{\partial \mathbf{x}}{\partial z} = \frac{\partial \mathbf{x}}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial z}$$

■ rétropropagation du gradient de l'erreur par rapport à y en erreur par rapport à w

#### Chain rule



Régression logistique : 
$$\frac{\partial \ell_{\textit{CE}}}{\partial \mathbf{W}} = \frac{\partial \ell_{\textit{CE}}}{\partial \hat{\mathbf{y_i}}} \frac{\partial \hat{\mathbf{y_i}}}{\partial \mathbf{s_i}} \frac{\partial \mathbf{s_i}}{\partial \mathbf{W}}$$



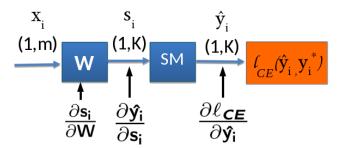
# Régression logistique : calcul des gradients (1/2)

$$\frac{\partial \ell_{CE}}{\partial \mathbf{W}} = \frac{\partial \ell_{CE}}{\partial \hat{\mathbf{y}_i}} \frac{\partial \hat{\mathbf{y}_i}}{\partial \mathbf{s}_i} \frac{\partial \mathbf{s}_i}{\partial \mathbf{W}}, \ \ell_{CE}(\hat{\mathbf{y}_i}, \mathbf{y}_i^*) = -log(\hat{y}_{c^*,i})$$

Mise à jour pour 1 example :

$$\begin{array}{c|c} & & & \\ \hline \bullet & \frac{\partial \ell_{\mathit{CE}}}{\partial \hat{\mathbf{y}}_{i}^{*}} = \frac{-1}{\hat{y}_{c^{*},i}} = \frac{1}{\hat{\mathbf{y}}_{i}} \odot \delta_{\mathbf{c},\mathbf{c}^{*}} \\ \hline \bullet & & \\ \hline \bullet & \frac{\partial \ell_{\mathit{CE}}}{\partial \mathbf{s}_{i}} = \hat{\mathbf{y}}_{i} - \mathbf{y}_{i}^{*} = \delta_{i}^{\mathbf{y}} \\ \end{array} \right] \text{ avec } \delta_{\mathbf{c},\mathbf{c}^{*}} = (0,\ldots,\underbrace{1}_{c^{*}},\ldots,0)$$

$$rac{\partial \ell_{\mathit{CE}}}{\partial \mathbf{W}} = \mathbf{x_i}^{\mathsf{T}} \delta_{\mathbf{i}}^{\mathbf{y}}$$



# Régression logistique : calcul des gradients (2/2)

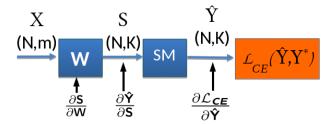
Pour l'ensemble du jeu de données X ( $N \times m$ ), étiquettes  $Y^*$  et prédictions  $\hat{Y}$  ( $N \times K$ ) :

$$\mathcal{L}_{CE}(\mathbf{W}, \mathbf{b}) = -\frac{1}{N} \sum_{i=1}^{N} log(\hat{y}_{c^*,i}), \ \frac{\partial \mathcal{L}_{CE}}{\partial \mathbf{W}} = \frac{\partial \mathcal{L}_{CE}}{\partial \hat{\mathbf{Y}}} \frac{\partial \hat{\mathbf{Y}}}{\partial \mathbf{S}} \frac{\partial \mathbf{S}}{\partial \mathbf{W}}$$

$$\frac{\partial \mathcal{L}_{CE}}{\partial \mathbf{S}} = \hat{\mathbf{Y}} - \mathbf{Y}^* = \Delta^y$$

$$\frac{\partial \mathcal{L}_{CE}}{\partial \mathbf{S}} - \mathbf{Y}^T \wedge y$$

$$\blacksquare \quad \frac{\partial \mathcal{L}_{CE}}{\partial \mathbf{W}} = \mathbf{X}^T \Delta^y$$



#### Entraînement du modèle

La fonction de coût est calculée sur tout le jeu de données d'entraînement :

$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^{N} \ell\left(\hat{\mathbf{y}}_{i}, \mathbf{y}_{i}^{*}; \mathbf{w}, \mathbf{b}\right)$$

Optimisation par descente de gradient :

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{w}} \left( \mathbf{w}^{(t)} \right)$$

- La complexité du calcul du gradient  $\nabla_{\mathbf{w}}^{(t)} = \frac{1}{N} \sum_{i=1}^{N} \frac{\partial \ell(\hat{\mathbf{y}}_i, \mathbf{y}_i^*)}{\partial \mathbf{w}} \left(\mathbf{w}^{(t)}\right)$  croît linéairement avec :
  - la dimensionalité de w (nombre de paramètres du modèle),
  - N, la taille du jeu de données (nombre d'exemples d'apprentissage).
- ⇒ coûteux, même pour des modèles de dimensionalité modérée et des jeux de données de taille moyenne!

## Descente de gradient stochastique

- **Solution**: approximation du vrai gradient  $\nabla_{\mathbf{w}}^{(t)} = \frac{1}{N} \sum_{i=1}^{N} \frac{\partial \ell(\hat{y_i}, \mathbf{y_i^*})}{\partial \mathbf{w}} \left(\mathbf{w}^{(t)}\right)$  sur un échantillon d'exemples (un *batch* ou "lot")
- ⇒ Stochastic Gradient Descent (SGD)
  - Version *online* : mise à jour à partir d'un seul exemple

$$abla_{\mathbf{w}}^{(t)} pprox rac{\partial \ell(\hat{\mathbf{y}_i}, \mathbf{y}_i^*)}{\partial \mathbf{w}} \left(\mathbf{w}^{(t)}\right)$$

■ Version par mini-batch : mise à jour sur B << N exemples :

$$\nabla_{\mathbf{w}}^{(t)} \approx \frac{1}{B} \sum_{i=1}^{B} \frac{\partial \ell(\hat{\mathbf{y}}_{i}, \mathbf{y}_{i}^{*})}{\partial \mathbf{w}} \left(\mathbf{w}^{(t)}\right)$$



### Avantages et inconvénients

La descente de gradient stochastique produit une approximation du véritable gradient  $\nabla_{\mathbf{w}}$ .

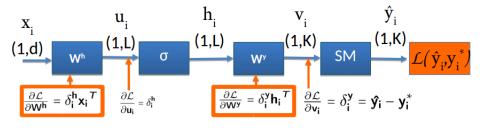
- estimation bruitée, pouvant envoyer une direction incorrecte (pente forte pour un exemple mais pas pour les autres),
- sensibilité aux outliers pour la version en ligne,
- + mises à jour des poids plus nombreuses car itérations moins coûteuses :  $\times N$  mises à jour (online),  $\times \frac{N}{R}$  mises à jour (mini-batch)
- + à nombre de calculs égal, convergence plus rapide

La SGD est incontournable pour la convergence des modèles profonds sur des jeux de données massifs



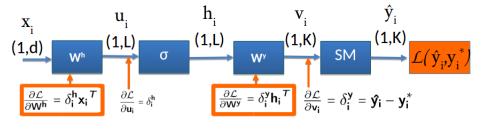
# Entraı̂nement du perceptron : rétropropagation

- $\blacksquare$  Passer de la régression logistique au perceptron multi-couche implique l'ajout d'une couche cachée (+ sigmoïde)
- Entraı̂nement : apprendre les paramètres  $\mathbf{W}^y$  et  $\mathbf{W}^h$  (+ bias) par **rétropropagation**  $\begin{bmatrix} \frac{\partial \mathcal{L}_{CE}}{\partial \mathbf{W}^y} = \frac{1}{N} \sum_{i=1}^N \frac{\partial \ell_{CE}}{\partial \mathbf{W}^h} \end{bmatrix}$  et  $\begin{bmatrix} \frac{\partial \mathcal{L}_{CE}}{\partial \mathbf{W}^h} = \frac{1}{N} \sum_{i=1}^N \frac{\partial \ell_{CE}}{\partial \mathbf{W}^h} \end{bmatrix}$



- La dernière couche est équivalente à une régression logistique.
- $\blacksquare \text{ Couche cachée}: \frac{\partial \ell_{\textit{CE}}}{\partial \mathbf{w}^h} = \mathbf{x_i}^{\textit{T}} \frac{\partial \ell_{\textit{CE}}}{\partial \mathbf{u_i}} \Rightarrow \text{calcul de } \frac{\partial \ell_{\textit{CE}}}{\partial \mathbf{u_i}} = \delta_i^h$

# Rétropropagation du gradient



$$lacksquare$$
 Coût pour un seul exemple :  $\mathcal{L}_W(\widehat{\mathbf{y}}_i, \mathbf{y}_i^*) = -\log(\hat{y}_{c^*,i})$ 

- Pour revenir en arrière : on calcule  $\frac{\partial \mathcal{L}}{\partial u_i} = \delta_i^h$
- Si on connaît  $\delta_i^h \Rightarrow \frac{\partial \mathcal{L}}{\partial W^h} = \delta_i^h {x_i}^T \sim \text{régression logistique}$

## Rétropropagation dans le perceptron

$$\begin{array}{c} X_{i} \\ (1,d) \\ \hline \frac{\partial \mathcal{L}}{\partial W^{h}} = \delta_{i}^{h} x_{i}^{T} \\ \hline \end{array} \begin{array}{c} U_{i} \\ (1,L) \\ \hline \frac{\partial \mathcal{L}}{\partial W^{y}} = \delta_{i}^{y} h_{i}^{T} \\ \hline \end{array} \begin{array}{c} V_{i} \\ (1,K) \\ \hline \frac{\partial \mathcal{L}}{\partial v_{i}} = \delta_{i}^{y} = \hat{\mathbf{y}}_{i} - \mathbf{y}_{i}^{*} \\ \hline \end{array}$$

■ Le calcul de  $\frac{\partial \ell_{\it CE}}{\partial u_i} = \delta^h_i \Rightarrow$  s'obtient par *chain rule* :

$$\frac{\partial \ell_{\textit{CE}}}{\partial u_i} = \frac{\partial \ell_{\textit{CE}}}{\partial v_i} \frac{\partial v_i}{\partial h_i} \frac{\partial h_i}{\partial u_i}$$

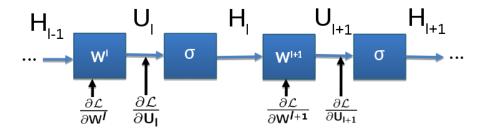
■ ce qui permet d'obtenir :

$$\frac{\partial \ell_{\textit{CE}}}{\partial \mathbf{u_i}} = \delta_i^h = \delta_i^{\textit{yT}} \mathbf{W}^{\textit{y}} \odot \sigma^{'}(\mathbf{h_i}) = \delta_i^{\textit{yT}} \mathbf{W}^{\textit{y}} \odot (\mathbf{h_i} \odot (1 - \mathbf{h_i}))$$

## Rétropropagation dans les réseaux profonds

- Perceptron multi-couche : ajout d'encore plus de couches
- Rétropropagation  $\sim$  Perceptron : **en supposant que l'on connaise**  $rac{\partial \mathcal{L}}{\partial \mathrm{U}_{1+1}} = \Delta^{l+1}$

$$\begin{array}{c|c} \hline & \frac{\partial \mathcal{L}}{\partial \mathbf{W}^{l+1}} = \mathbf{H_l}^T \Delta^{l+1} \\ \hline & \mathbf{Calcul} \ \mathrm{de} \ \frac{\partial \mathcal{L}}{\partial \mathbf{U_l}} = \Delta^l \\ \hline & \frac{\partial \mathcal{L}}{\partial \mathbf{W}^l} = \mathbf{H_{l-1}}^T \Delta^{h_l} \\ \hline \end{array} \right) \\ (= \Delta^{l+1}^T \mathbf{W}^{l+1} \odot \mathbf{H_l} \odot (1 - \mathbf{H_l}) \ \mathrm{pour} \ \mathrm{la} \ \mathrm{sigmo\"{ide}}) \\ \hline \end{array}$$



### Plan du cours

- 1 Introduction
- 2 Réseaux de neurones artificiels
- 3 Rétropropagation du gradient
- 4 Optimisation des réseaux profonds
- 5 Implémentation de la rétropropagation
- 6 Activations, régularisations, initialisations

## La rétropropagation en résumé

- Rétropropagation (backpropagation) : solution pour l'entraînement de bout en bout de tous les paramètres des réseaux profonds
  - Algorithme clé de l'apprentissage profond!

#### Historique de la rétropropagation

- Dérivation des fonctions composées (Leibniz, 1676)
- Optimisation par descente de gradient (CAUCHY 1847)
- Premières applications de la rétropropagation (programmation dynamique, contrôle Kelley 1960; Bryson et Ho 1969)
- Formalisation de la rétropropagation pour les réseaux de neurones (LINNAINMAA 1970; WERBOS 1975)
- Application aux réseaux de neurones profonds (RUMELHART, HINTON et WILLIAMS 1986; LECUN 1988)

Mais y a-t-il des inconvénients à utiliser la rétropropagation?

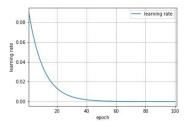
# Optimisation : politique du pas d'apprentissage

- lacksquare Mise à jour par descente de gradient :  $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} \eta \nabla_{\mathbf{w}}^{(t)}$
- Comment choisir le pas d'apprentissage  $\eta$ ?

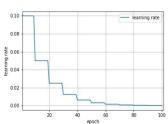
## Politiques d'adaptation du pas d'apprentissage

En général, l'heurisiste consiste à faire décroître  $\eta$  durant l'apprentissage (*learning rate "decay"*).

- Décroissance inversement proportionnelle à l'itération :  $\eta_t = \frac{\eta_0}{1+r \cdot t}$ , r vitesse de décroissance
- Décroissance exponentielle :  $\eta_t = \eta_0 \cdot e^{-\lambda t}$
- Décroissance en escalier :  $\eta_t = \eta_0 \cdot r^{\frac{t}{t_u}}$  ...



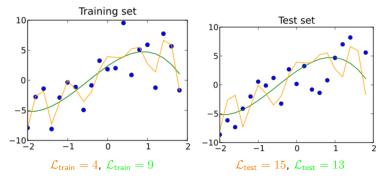
Décroissance exponentielle ( $\eta_0 = 0.1, \lambda = 0.1$ )



Décroissance step ( $\eta_0 = 0.1, r = 0.5, t_u = 10$ )

# Optimisation, généralisation et sur-apprentissage

- **Apprentissage** : minimisation d'une fonction de coût  $\mathcal L$  sur un jeu de données  $\implies$  risque empirique
  - Jeu d'apprentissage : ensemble d'échantillons représentatif de la distribution des données et des étiquettes
  - Objectif : apprendre une fonction de prédiction avec une erreur faible sur la distribution (inconnue) des données réelles ⇒ risque espéré



 $\textbf{Optimisation} \neq \textbf{apprentissage} \ ! \implies \textbf{sur-apprentissage} \neq \textbf{g\'en\'eralisation}$ 

#### Régularisation

#### Régularisation

Réduire la capacité du modèle pour réduire l'écart entre performances entraînement/test.

- Avec un modèle de suffisamment grande capacité, améliore la généralisation et donc les performances en test.
- Régularisation structurelle : ajout d'une contrainte sur les poids pour les forcer à suivre un a priori

$$\mathcal{L}_r(\mathbf{w}) = \mathcal{L}(\mathbf{y}; \mathbf{y}^*; \mathbf{w}) + \alpha R(\mathbf{w})$$

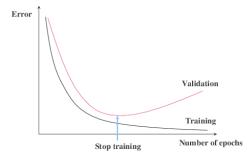
lacktriangle Weight decay : régularisation  $L_2$  pour pénaliser les poids de norme élevée :

$$R(\mathbf{w}) = ||\mathbf{w}||^2$$

- Utilisation courante pour les réseaux de neurones
- Justifications théoriques et bornes de généralisation dans le cas des SVM
- Autres possibilités pour  $R(\mathbf{w})$ : régularisation  $L_1$ , Dropout, etc.

# Régularisation et hyperparamètres

- Hyperparamètres pour les réseaux de neurones :
  - Hyperparamètres d'optimisation : pas d'apprentissage (et évolution), nombre d'itérations, régularisation, etc.
  - Hyperparamètres d'architecture : nombre de couches, nombre de neurones, choix de la non-linéarité, etc.
- Le réglage des hyperparamètres permet d'ajuster la capacité du modèle et d'améliorer la généralisation :
  - L'utilisation d'un jeu de validation est cruciale!



#### Plan du cours

- 1 Introduction
- 2 Réseaux de neurones artificiels
- Rétropropagation du gradien
- 4 Optimisation des réseaux profonds
- 5 Implémentation de la rétropropagation
- 6 Activations, régularisations, initialisations

## Optimisation des réseaux profonds

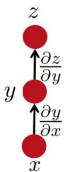
 $\blacksquare$  Optimisation des réseaux profonds : descente de gradient sur la fonction de coût  $\mathcal L$ 

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta \nabla \mathcal{L} \left( \mathbf{w}^t \right)$$

- Rétropropagation de l'erreur : algorithme pour calculer  $\nabla \mathcal{L}\left(\mathbf{w}^{t}\right)$  dans un réseau de neurones
  - Expression analytique du gradient par chain rule
  - mais le calcul du gradient peut être difficile à réaliser efficacement

#### Deux sources de difficultés

- Problèmes d'optimisation numérique
- Différenciation automatique (automatic differenciation)



- L'accumulation des approximations (arrondis) pose problème.
- Underflow :  $x \approx 0$  ou x = 0 provoque des comportements différents ■ Division par  $0 \Rightarrow$  indéterminé (not a number)
- *Overflow* : grandes valeurs de x > 0 ou  $x < 0 \Rightarrow not$  a number

Exemple : softmax sur 
$$\mathbf{x} = \{x_1, x_2, ..., x_K\}$$

$$s(\mathbf{x})_i = \frac{e^{x_i}}{\sum_{j=1}^K e^{x_j}}$$

- Que se passe-t-il si  $x_i = C$  pour tout i?
- En théorie,  $s(\mathbf{x})_i = \frac{1}{K}$  quel que soit i, mais en pratique

- L'accumulation des approximations (arrondis) pose problème.
- Underflow :  $x \approx 0$  ou x = 0 provoque des comportements différents ■ Division par  $0 \Rightarrow$  indéterminé (not a number)
- Overflow : grandes valeurs de x > 0 ou  $x < 0 \Rightarrow not$  a number

Exemple : softmax sur 
$$\mathbf{x} = \{x_1, x_2, ..., x_{\mathcal{K}}\}$$

$$s(\mathbf{x})_i = \frac{e^{x_i}}{\sum_{j=1}^K e^{x_j}}$$

- Que se passe-t-il si  $x_i = C$  pour tout i?
- $\blacksquare$  En théorie,  $s(\mathbf{x})_i = \frac{1}{K}$  quel que soit i, mais en pratique

- L'accumulation des approximations (arrondis) pose problème.
- Underflow:  $x \approx 0$  ou x = 0 provoque des comportements différents ■ Division par 0 ⇒ indéterminé (not a number)
- Overflow: grandes valeurs de x > 0 ou  $x < 0 \Rightarrow not$  a number

# Exemple : softmax sur $\mathbf{x} = \{x_1, x_2, ..., x_K\}$

$$s(\mathbf{x})_i = \frac{e^{x_i}}{\sum_{j=1}^K e^{x_j}}$$

- Que se passe-t-il si  $x_i = C$  pour tout i?
- En théorie,  $s(\mathbf{x})_i = \frac{1}{\kappa}$  quel que soit i, mais en pratique :

- L'accumulation des approximations (arrondis) pose problème.
- Underflow :  $x \approx 0$  ou x = 0 provoque des comportements différents Division par  $0 \Rightarrow$  indéterminé (not a number)
- Overflow : grandes valeurs de x > 0 ou  $x < 0 \Rightarrow not$  a number

# Exemple : softmax sur $\mathbf{x} = \{x_1, x_2, ..., x_K\}$

$$s(\mathbf{x})_i = \frac{e^{x_i}}{\sum_{j=1}^K e^{x_j}}$$

- Que se passe-t-il si  $x_i = C$  pour tout i?
- En théorie,  $s(\mathbf{x})_i = \frac{1}{K}$  quel que soit i, mais en pratique :
  - $C \to -\infty, e^C \to 0$ , division par 0, not a number! (underflow)
  - $C \to +\infty, e^C \to +\infty$ , not a number! (overflow)

- L'accumulation des approximations (arrondis) pose problème.
- Underflow :  $x \approx 0$  ou x = 0 provoque des comportements différents ■ Division par  $0 \Rightarrow$  indéterminé (not a number)
- Overflow : grandes valeurs de x > 0 ou  $x < 0 \Rightarrow not$  a number

# Exemple : softmax sur $\mathbf{x} = \{x_1, x_2, ..., x_K\}$

$$s(\mathbf{x})_i = \frac{e^{x_i}}{\sum_{j=1}^K e^{x_j}}$$

- Que se passe-t-il si  $x_i = C$  pour tout i?
- En théorie,  $s(\mathbf{x})_i = \frac{1}{K}$  quel que soit i, mais en pratique :
  - $C \to -\infty$ ,  $e^C \to 0$ , division par 0, not a number! (underflow)
  - lacksquare  $C o +\infty, e^C o +\infty$ , not a number! (overflow)

# Problèmes d'optimisation numérique

- L'accumulation des approximations (arrondis) pose problème.
- Underflow :  $x \approx 0$  ou x = 0 provoque des comportements différents ■ Division par  $0 \Rightarrow$  indéterminé (not a number)
- Overflow: grandes valeurs de x > 0 ou  $x < 0 \Rightarrow$  not a number

# Exemple : softmax sur $\mathbf{x} = \{x_1, x_2, ..., x_K\}$

$$s(\mathbf{x})_i = \frac{e^{x_i}}{\sum_{j=1}^K e^{x_j}}$$

- Que se passe-t-il si  $x_i = C$  pour tout i?
- En théorie,  $s(\mathbf{x})_i = \frac{1}{\kappa}$  quel que soit i, mais en pratique :
  - $lacksquare C 
    ightarrow -\infty, e^C 
    ightarrow 0$ , division par 0, not a number! (underflow)
  - $lacksquare C o +\infty, e^C o +\infty$ , not a number! (overflow)

### Une solution

Ajouter une stabilisation numérique au dénominateur :

$$\mathbf{z} = \mathbf{x} - \max_{i}(x_i) \Rightarrow s(\mathbf{z}) = s(\mathbf{x})$$

- $= \max_i(e^{z_i}) = 1 \Rightarrow \mathsf{pas} \mathsf{d'overflow}$
- $\max_i(e^{z_i}) = 1 \Rightarrow \mathsf{pas} \; \mathsf{d'underflow}$

### Cas de l'entropie croisée

$$\mathbf{z} = \mathbf{x} - \max_i(x_i) \Rightarrow s(\mathbf{z}) = s(\mathbf{x})$$

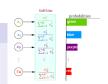
Que se passe-t-il pour le calcul de  $\log [s(\mathbf{z})]$ , par exemple dans l'entropie croisée ?

■ 
$$s(\mathbf{z}) = 0$$
 (underflow)  $\Rightarrow \log[s(\mathbf{z})] \to -\infty$ : not a number!

### Solution : stabiliser le $\log$

$$log[s(\mathbf{x})_i] = x_i - log\left[\sum_{j=1}^K e^{x_j}\right]$$

- $\blacksquare$  À nouveau, on substitue :  $\mathbf{z} = \mathbf{x} \max_i(x_i)$
- ⇒ underflow, overflow?



### Calcul des dérivées partielles

Plusieurs méthodes permettent de calculer les dérivées de façon automatique :

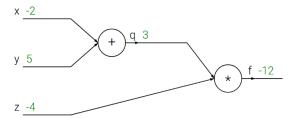
- Approximation numérique :  $f'(x) \approx \frac{f(x+h)*f(x)}{h}$  (méthode des éléments finis)
  - $\oplus$  Ne nécessite aucune connaissance sur f
  - Approximation, stabilité numérique (underflow/overflow)
  - → Requiert d'évaluer plusieurs fois f
- lacksquare Dérivation symbolique ( $\sim$  Mathematica)
  - ⊖ Expressions qui deviennent rapidement complexes, termes redondants ⇒ faible efficacité

$$\begin{array}{lll} f(x) & f'(x) & f'(x) \\ 64x(1-x)(1-2x)^2 & 128x(1-x)(-8+16x)(1-2x)^2(1-64(1-42x+504x^2-2640x^3+61(1-8x+8x^2)^2) \\ 8x^2+8x^2) + 64(1-x)(1-2x)^2(1-8x+7040x^4-9984x^3+7168x^6-2048x^7) \\ 8x^2)^2 - 64x(1-2x)^2(1-8x+8x^2)^2 \\ -266x(1-x)(1-2x)(1-8x+8x^2)^2 \end{array}$$

- Différenciation automatique
  - Alterne entre dérivation symboliques et étapes de simplification
  - Différenciation symbolique au niveau des opérations élémentaires
  - Simplification en conservant les résultats numériques intermédiaires

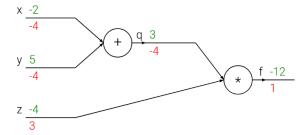
# Automatic Differentiation (AD)

- Définition d'un graphe de calcul : abstraction pour représenter les séquences d'opérations à inverser pour rétropropager le gradient
- Exemple : f(x, y, z) = (x + y) \* z
- Objectif : calculer les valeurs numériques de  $\frac{\partial f}{\partial x}$ ,  $\frac{\partial f}{\partial y}$  et  $\frac{\partial f}{\partial z}$
- Supposons que x = -2, y = 5, z = -4, propagation avant  $\Rightarrow q = 3$ , f = -12



# Rétropropagation et graphe de calcul

$$\blacksquare$$
  $\frac{\partial f}{\partial f}=1$  , rétropropagation pour déterminer  $\Rightarrow$   $\frac{\partial f}{\partial q}$ ,  $\frac{\partial f}{\partial z}$ 

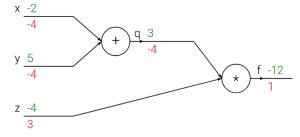


### Différenciation automatique : forward pass, backward pass

- Différenciation automatique dans le graphe de calcul : passe avant (forward) puis différenciation automatique en passe arrière (backward)
  - Backward pass : calcul récursif des dérivées de haut en bas
- Programmation dynamique : gain d'efficacité important par rapport à la différenciation

Une implémentation différent possible est l'autodifférenciation en passe avant :

En pratique, pour un graphe à N nœuds, il faut N passes avant pour la forward AD, contre 1 passe avant + 1 passe arrière pour la backward AD (backprop).



### Quelles opérations implémenter?

- La différenciation symbolique est implémentée seulement au niveau des opérations "atomiques":
  - $\blacksquare$  Les opérateurs arithmétiques binaires :  $+,-,\times,/$  , etc.
  - $\blacksquare$  Les fonctions "élémentaires" :  $\exp, \log, \cos, \sin,$  etc.
- Les fonctions dont les dérivées sont connues, avec une expression analytique dont le comportement numérique est stable :
  - Sigmoïde :  $\sigma(x) = \frac{1}{1 + e^{-x}} \implies \sigma'(x) = \sigma(x) [1 \sigma(x)]$

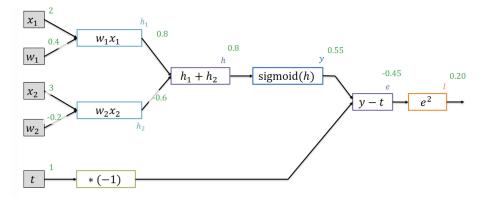
$$x = 1.0 \Rightarrow \sigma(x) = 0.73, \sigma'(x) = 0.2$$



Le graphe de calcul est remplaçable par un seul "bloc"dont l'inverse est  $\sigma'(x)$ .

# Mise en œuvre pour un réseau de neurones





 $\Rightarrow$  Backward : stockage des valeurs numériques des gradients  $\frac{\partial \mathbf{h}}{\partial \mathbf{W}}$ ,  $\frac{\partial \mathbf{h}}{\partial \mathbf{x}}$ 

# Ressources logicielles pour l'apprentissage profond

- Fonctionnalités indispensables :
  - Construction du graphe de calcul et auto-différenciation (autograd)
  - Exécution transparente sur CPU ou GPU (pas de programmation bas niveau)
- De nombreuses bibliothèques disponibles : MatConvNet (Matlab), Caffe/Caffe2 (C++/Python), Torch (Lua/Python/C++), Theano (Python), TensorFlow (Python), JAX (Python), MXNet (C++/Java), etc.
- Actuellement, les deux principales bibliothèques logicielles sont :
  - TensorFlow/Keras (soutenu par Google)
  - PyTorch (soutenu par Facebook)















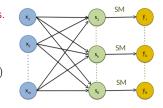


### Keras : implémentation de la régression logistique

- Entrée :  $\mathbf{x}_i \in \mathbb{R}^d$ , d=784, étiquette  $\mathbf{y}_i$  à K=10 classes.
- $\blacksquare$  Activations :  $\mathbf{s}_i = \mathbf{x}_i \mathbf{W} + \mathbf{b}$

from tensorflow import keras

- Sortie (softmax) :  $\hat{\mathbf{y}_k} = e^{\mathbf{s}_k} / (\sum_{k'=1}^K e^{\mathbf{s}_{k'}})$
- $\blacksquare$  Fonction de coût : entropie croisée =  $\frac{1}{N}\sum_{i=1}^N \ell_{CE}(\hat{\mathbf{y}}_i,\mathbf{y}_i^*)$
- Nombre de paramètres :  $784 \times 10 + 10 = 7850$



```
from Keras.datasets import mnist
(X train, y train), (X test, y.test) = keras.datasets.mnist.load_data()

# Pré-traitement et normalisation
X_train, X_test = X_train.reshape(-1, 784), X_test.reshape(-1, 784)
X_train, X_test = X_train / 255. X_test / 255.

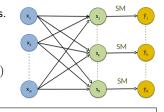
# Encodage one hot des vecteurs de classe
Y train, Y test = keras.utils.to categorical(v train, 10), keras.utils.to categorical(v test, 10)
```

### Keras : implémentation de la régression logistique

- Entrée :  $\mathbf{x}_i \in \mathbb{R}^d$ , d = 784, étiquette  $\mathbf{y}_i$  à K = 10 classes.
- Activations :  $\mathbf{s}_i = \mathbf{x}_i \mathbf{W} + \mathbf{b}$

from tensorflow import keras

- lacksquare Sortie (softmax) :  $\hat{\mathbf{y}_k} = e^{\mathbf{s}_k}/(\sum_{k'=1}^K e^{\mathbf{s}_{k'}})$
- Fonction de coût : entropie croisée =  $\frac{1}{N} \sum_{i=1}^{N} \ell_{CE}(\hat{\mathbf{y}}_i, \mathbf{y}_i^*)$
- Nombre de paramètres :  $784 \times 10 + 10 = 7850$



```
from keras.datasets import mnist
(X train, y train). (X, test, y test) = keras.datasets.mnist.load_data()

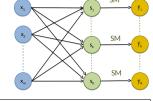
# Pré-traitement et normalisation
X train, X_test = X_train.reshape(-1, 784), X_test.reshape(-1, 784)
X_train, X_test = X_train. / 255., X_test / 255.

# Encodage one hot des vecteurs de classe
Y_train, Y_test = keras.utils.to_categorical(y_train, 10), keras.utils.to_categorical(y_test, 10)
```

```
# Instancie un modèle vide
from tensorflow.keras import Sequential
model = Sequential()
# Ajoute les couches au modèle
from tensorflow.keras.layers import Dense, Activation
model.add(Dense(10, input_dim=784, name='fc1'))
model.add(Activation('softmax'))
```

# Keras : implémentation de la régression logistique

- Entrée :  $\mathbf{x}_i \in \mathbb{R}^d$ , d = 784, étiquette  $\mathbf{y}_i$  à K = 10 classes.
- $\blacksquare$  Activations :  $\mathbf{s}_i = \mathbf{x}_i \mathbf{W} + \mathbf{b}$
- lacksquare Sortie (softmax) :  $\hat{\mathbf{y}_k} = e^{\mathbf{s}_k}/(\sum_{k'=1}^K e^{\mathbf{s}_{k'}})$
- Fonction de coût : entropie croisée  $= \frac{1}{N} \sum_{i=1}^{N} \ell_{CE}(\hat{\mathbf{y}}_i, \mathbf{y}_i^*)$
- Nombre de paramètres :  $784 \times 10 + 10 = 7850$



```
from tensorflow import keras
from keras.datasets import mnist
(X.train, y.train), (X.test, y.test) = keras.datasets.mnist.load_data()
# Pré-traitement et normalisation
X.train, X.test = X.train.reshape(-1, 784), X.test.reshape(-1, 784)
X.train, X.test = X.train./ 255. X.test / 255.
# Encodage one hot des vecteurs de classe
Y.train, Y.test = keras.utils.to.categorical(y.train, 10), keras.utils.to.categorical(y.test, 10)
```

```
# Instancie um modèle vide
from tensorflow.keras import Sequential
model = Sequential()
# Ajoute les couches au modèle
from tensorflow.keras.layers import Dense, Activation
model.add(Dense(10, input_dim=784, name='fc1'))
model.add(Activation('softmax'))
```

# Keras : apprentissage/inférence

■ Visualisation d'un résumé du modèle :

model.summary()

Layer (type)	Output Shape	Param #
fc1 (Dense)	(None, 10)	7850
activation_1 (Activation)	(None, 10)	0
Total params: 7,850.0 Trainable params: 7,850.0 Non-trainable params: 0.0		

■ Entraı̂nement des paramètres sur le jeu de données d'apprentissage :

```
# Apprentissage, méthode .fit() "à la scikit-learn"
model.fit(X_train, y_train,batch_size=128, epochs=20,verbose=1)
```

#### Keras : évaluation du modèle

Évaluation des performances sur le jeu de test :

```
scores = model.evaluate(X_test, Y_test, verbose=0)
print("%s TEST: %.2f%%" % (model.metrics_names[0], scores[0]*100))
print("%s TEST: %.2f%%" % (model.metrics_names[1], scores[1]*100))
```

```
Epoch 1/10
Epoch 2/10
Fpoch 3/10
Epoch 4/10
Epoch 5/10
Epoch 6/10
60000/60000 [============] - 1s - loss: 0.2883 - acc: 0.9196
Epoch 7/10
Epoch 8/10
Epoch 9/10
Epoch 10/10
loss TEST: 27.20%
acc TEST: 92.20%
```

⇒ mise en pratique dans la séance de travaux pratiques "Deep Learning avec Keras"

### Plan du cours

- 1 Introduction
- 2 Réseaux de neurones artificiels
- Rétropropagation du gradient
- 4 Optimisation des réseaux profonds
- 5 Implémentation de la rétropropagation
- 6 Activations, régularisations, initialisations

### Non-linéarités modernes

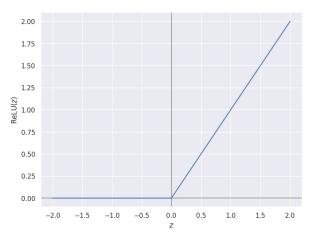
 $\blacksquare$  Activations non-linéaires classiques : sigmoïde, tanh

$$\sigma(z) = \frac{1}{1 + e^{-z}} \qquad \tanh(z) = \frac{e^z - e^{-z}}{e^z - e^{-z}}$$

- $\ominus$  Régime de saturation  $\Longrightarrow$  dérivée nulle hors de  $x \approx 0$ 
  - ⇒ Gradient "évanescent" (vanishing gradient) : rétropropagation limitée
  - ⇒ Convergence lente car la mise à jour des paramètres est de faible amplitude

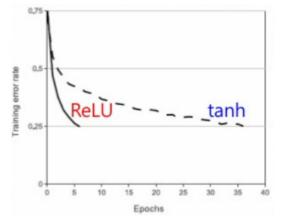
# Rectified Linear Unit (ReLU) NAIR et HINTON 2010

$$\mathit{ReLU}(z) = \begin{cases} z \text{ si } z \ge 0 \\ 0 \text{ sinon} \end{cases} = \max\{0, z\}$$



### Rectified Linear Unit (ReLU)

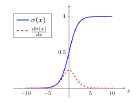
- Réduire l'occurrence des gradients évanescents permet :
  - ⇒ d'entraîner des modèles plus profonds (la décroissance de la norme du gradient est plus lente),
  - ⇒ d'accélérer la convergence du modèle.



Exemple : CNN à 4 couches entraîné sur CIFAR-10. Remplacer  $\tanh$  par ReLU permet de converger en  $6\times$  moins d'itérations. Extrait de KRIZHEVSKY, SUTSKEVER et HINTON 2012

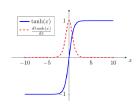
### Panel de fonctions d'activation

### Sigmoïde



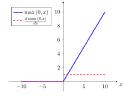
- Saturation
- Calcul coûteux
- Non centrée sur zéro

#### Tanh



- Saturation
- Calcul coûteux
- Centrée sur zéro

#### ReLU

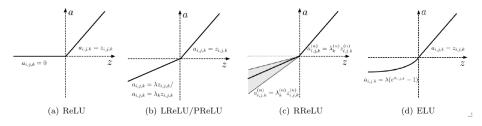


- Pas de saturation
- Efficace à calculer
- Non centrée sur zéro
- Les activations négatives sont ignorées

# Fonctions d'activation non-linéaire paramétriques

Deux observations concernant la fonction d'activation ReLU :

- Pour z < 0, l'activation ReLU vaut 0 avec un gradient nul  $\implies$  pas de gradient propagé.
- Peut-on déterminer automatiquement certains paramètres de la fonction d'activation (pente, seuil, etc.)?



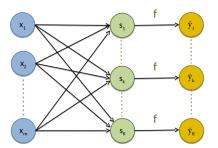
- Leaky ReLU (LReLU) : pente négative  $\lambda$  hyperparamètre à régler
- Parametric ReLU (PReLU) : pente négative  $\lambda_k$  apprise pendant l'entraînement
- Randomized ReLU (RReLU) : pente négative  $\lambda_k^n$  tirée selon une loi uniforme
- lacktriangle Exponential Linear Unit (ELU) : paramètre de lissage  $\lambda$  hyperparamètre à régler

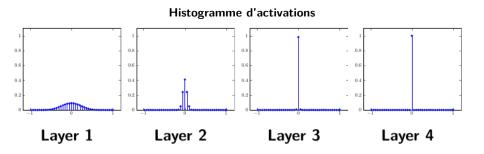
### Initialisation des poids

### Comment initialiser les paramètres du modèle?

La fonction objectif à minimiser non-convexe par rapport aux paramètres des modèles profonds  $\implies$  l'initialisation des poids est importante!

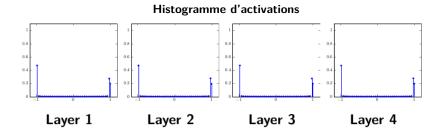
- Zero-init : toutes les connexions à zéro ⇒ tous les neurones ont la même sortie, et donc le même gradient!
- Initialisation aléatoire faible, par exemple  $\mathbf{W} \sim \mathcal{U}(-0.1, +0.1)$  ou  $\mathbf{W} \sim \mathcal{N}(0, \sigma^i)$
- Initialisation "Glorot":
  - lacksquare Pour une entrée x de dimension m et une sortie s, on a :  $\mathrm{Var}[s] = m\,\mathrm{Var}[\mathbf{w}]\,\mathrm{Var}[\mathbf{x}]$
  - Idée : initialiser les poids par  $\mathbf{W} \sim \frac{1}{\sqrt{m}} \mathcal{N}(0, \sigma^i)$  GLOROT et BENGIO 2010





- Perceptron multi-couche : 10 couches, 500 neurones par couche.
- Activation : tanh
- Initialisation :  $\mathbf{w} \sim \mathcal{N}(0, \sigma^i)$

 $\implies$  dans le cas  $\sigma^i$  faible, les activations sont toutes proches de 0, mauvaise initialisation.

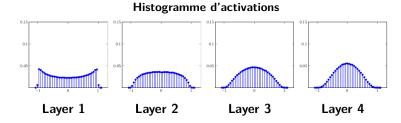


■ Perceptron multi-couche : 10 couches, 500 neurones par couche.

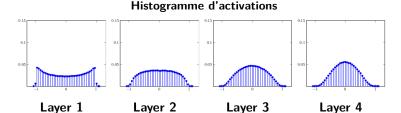
■ Activation : tanh

■ Initialisation :  $\mathbf{w} \sim \mathcal{N}(0, \sigma^i)$ 

 $\Rightarrow$  dans le cas  $\sigma^i$  grand, les activations sont toutes saturées ( $\pm 1$ ), mauvaise initialisation (gradient évanescent).



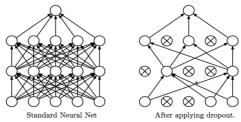
- Perceptron multi-couche : 10 couches, 500 neurones par couche.
- Activation : tanh
- Initialisation : "Glorot",  $\mathbf{W} \sim \frac{1}{\sqrt{500}} \mathcal{N}(0,1)$
- ⇒ variance adaptative et contrôlée pour chaque couche.



- Perceptron multi-couche : 10 couches, 500 neurones par couche.
- Activation : ReLU
- Dans le cas ReLU :  $Var[s] = 2m Var[\mathbf{w}] Var[\mathbf{x}] \implies \mathbf{W} \sim \frac{1}{\sqrt{2m}} \mathcal{N}(0, \sigma^i)$
- Initialisation : "Glorot",  $\mathbf{W} \sim \frac{1}{\sqrt{1000}} \mathcal{N}(0,1)$
- ⇒ variance adaptative et contrôlée pour chaque couche.

### Régularisation : Dropout SRIVASTAVA et al. 2014

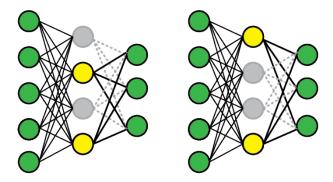
- Supprime aléatoirement un neurone caché avec une probabilité p, par exemple  $p = \frac{1}{2}$ .
- Méthode de régularisation pour limiter le sur-apprentissage
  - Éviter la co-adaptation
  - Forcer la redondance entre les poids
- Interprétation alternative : combinaison de nombreux réseaux de neurones qui diffèrent légèrement



Credits: Geoffrey E. Hinton, NIPS 2012

### Dropout : implémentation

- Optimisation : le dropout est différentiable
  - Les activations des neurones dropped sont mises à zéro
  - Les mises à jour de ces neurones sont ignorées
- À l'inférence, plusieurs possibilités :
  - Appliquer le dropout en inférence et moyenner sur les prédictions de sortie
  - $\blacksquare$  Alternative plus rapide : inférence normale mais multiplier par p toutes les activations
    - Équivalent à la moyenne géométrique pour un perceptron
    - Approximation assez bonne pour un perceptron multi-couche

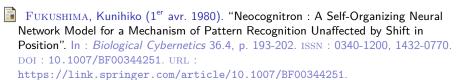


# Bibliographie I

- BRYSON, Arthur E. et Yu-Chi Ho (1969). Applied Optimal Control: Optimization, Estimation, and Control. A Blaisdell Book in the Pure and Applied Sciences.

  Waltham, Mass: Blaisdell Pub. Co. 481 p.
- CAUCHY, Augustin Louis (juill. 1847). Comptes rendus hebdomadaires des séances de l'Académie des sciences. Paris : Gauthier-Villars. URL : http://gallica.bnf.fr/ark:/12148/bpt6k2982c.
  - CIREŞAN, Dan C., Ueli MEIER et Jürgen SCHMIDHUBER (2012). "Multi-Column Deep Neural Networks for Image Classification". In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). Washington, DC, United States, p. 3642-3649. ISBN: 978-1-4673-1226-4. URL: http://dl.acm.org/citation.cfm?id=2354409.2354694.
  - CYBENKO, George (1er déc. 1989). "Approximation by Superpositions of a Sigmoidal Function". In: *Mathematics of Control, Signals and Systems* 2.4, p. 303-314. ISSN: 0932-4194, 1435-568X. DOI: 10.1007/BF02551274. URL: https://link.springer.com/article/10.1007/BF02551274.
  - DEVLIN, Jacob et al. (2018). "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding". In: Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, p. 4171-4186.

### Bibliographie II



GLOROT, Xavier et Yoshua BENGIO (31 mars 2010). "Understanding the Difficulty of Training Deep Feedforward Neural Networks". In: Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics. Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics, p. 249-256. URL: http://proceedings.mlr.press/v9/glorot10a.html (visité le 13/04/2018).

GOODFELLOW, lan et al. (2014). "Generative Adversarial Networks". In: Advances in Neural Information Processing Systems, p. 2672-2680. URL: https://proceedings.neurips.cc/paper/2014/file/5ca3e9b122f61f8f06494c97b1afccf3-Paper.pdf.

HOCHREITER, Sepp et Jürgen Schmidhuber (1997). "Long short-term memory". In: Neural computation 9.8, p. 1735-1780.

# Bibliographie III

- HORNIK, Kurt (1er jan. 1991). "Approximation Capabilities of Multilayer Feedforward Networks". In: Neural Networks 4.2, p. 251-257. ISSN: 0893-6080. DOI: 10.1016/0893-6080(91)90009-T. URL: http://www.sciencedirect.com/science/article/pii/089360809190009T.
- KELLEY, Henry J. (oct. 1960). "Gradient Theory of Optimal Flight Paths". In: ARS Journal 30.10, p. 947-954. DOI: 10.2514/8.5282. URL: https://arc.aiaa.org/doi/10.2514/8.5282 (visité le 07/04/2023).
  - KRIZHEVSKY, Alex, Ilya SUTSKEVER et Geoffrey E. HINTON (2012). "ImageNet Classification with Deep Convolutional Neural Networks". In: Proceedings of the Neural Information Processing Systems (NIPS). NeurIPS, p. 1097-1105. URL: http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf.
- LECUN, Yann (1988). "A Theoretical Framework for Back-Propagation". In: Proceedings of the 1988 Connectionist Models Summer School, CMU, Pittsburg, PA. Sous la dir. de D. Touretzky, G. Hinton et T. Sejnowski, p. 21-28.
- LECUN, Yann et al. (déc. 1989). "Backpropagation Applied to Handwritten Zip Code Recognition". In: Neural Computation 1.4, p. 541-551. ISSN: 0899-7667. DOI: 10.1162/neco.1989.1.4.541.

### Bibliographie IV

- LINNAINMAA, Seppo (1970). "The representation of the cumulative rounding error of an algorithm as a Taylor expansion of the local rounding errors". Helsinki University.
- MCCULLOCH, Warren S. et Walter H. PITTS (1943). "A Logical Calculus of the Ideas Immanent in Nervous Activity". In: Bulletin of Mathematical Biophysics 5, p. 115-133. URL: http://www.cse.chalmers.se/~coquand/AUTOMATA/mcp.pdf.
- MINSKY, Marvin et Seymour A. PAPERT (2017). Perceptrons: An Introduction to Computational Geometry. Réédition augmenté de l'édition originale de 1969. The MIT Press. ISBN: 0262534770.
- NAIR, Vinod et Geoffrey E. HINTON (2010). "Rectified Linear Units Improve Restricted Boltzmann Machines". In: Proceedings of the 27th International Conference on Machine Learning (ICML-10), p. 807-814.
- RADFORD, Alec et al. (2019). "Language models are unsupervised multitask learners". In: *OpenAl blog* 1.8, p. 9.
- ROSENBLATT, Frank (1957). The Perceptron : A Probabilistic Model for Information Storage and Organization In The Brain.

# Bibliographie V

- RUMELHART, D. E., G. E. HINTON et R. J. WILLIAMS (1986). "Learning Internal Representations by Error Propagation". In: sous la dir. de David E. RUMELHART, James L. McClelland et given-i=CORPORATE family=PDP RESEARCH GROUP given=CORPORATE. Cambridge, MA, USA: MIT Press, p. 318-362. ISBN: 978-0-262-68053-0. URL: http://dl.acm.org/citation.cfm?id=104279.104293 (visité le 13/04/2018).
- SILVER, David et al. (2016). "Mastering the Game of Go with Deep Neural Networks and Tree Search". In: t. 529. 7587, p. 484-489.
- SRIVASTAVA, Nitish et al. (2014). "Dropout: A Simple Way to Prevent Neural Networks from Overfitting". In: Journal of Machine Learning Research 15, p. 1929-1958. URL: http://jmlr.org/papers/v15/srivastava14a.html (visité le 19/01/2016).
- VASWANI, Ashish et al. (2017). "Attention Is All You Need". In: Advances in Neural Information Processing Systems, p. 5998-6008.
- WERBOS, Paul John (1975). "Beyond Regression: New Tools for Prediction and Analysis in the Behavioral Sciences". Harvard University. 906 p. Google Books: z81XmgEACAAJ.