

Réseaux de neurones  
IFT 603-712

Réseaux de neurones multicouches  
Par  
Pierre-Marc Jodoin

1

---

---

---

---

---

---

---

---

Rappel réseaux de neurones

(Perceptron, régression logistique, SVM)

2

---

---

---

---

---

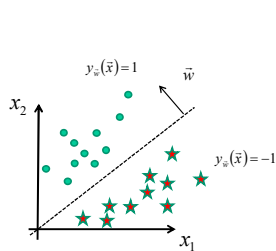
---

---

---

Séparation linéaire

(2D et 2 classes)



biais      poids

$$y_{\vec{w}}(\vec{x}) = w_0 + w_1x_1 + w_2x_2$$
$$= w_0 + \vec{w}^T \vec{x}$$
$$= \vec{w}'^T \vec{x}'$$

$$y_{\vec{w}}(\vec{x}) = \vec{w}^T \vec{x}$$

Par simplicité

- 2 grands **avantages**. Une fois l'entraînement terminé,
1. Plus besoin de données d'entraînement
  2. Classification est très rapide (**produit scalaire** entre 2 vecteurs)

3

---

---

---

---

---

---

---

---

### Perceptron

(2D et 2 classes)

Neurone

Produit scalaire + fonction d'activation

Fonction de coût Perceptron (loss) et gradient

$$E_D(\vec{w}) = \sum_{n=1}^N -t_n \vec{w}^T \vec{x}_n \quad \text{où } M \text{ est l'ensemble des données mal classées}$$

$$\nabla E_D(\vec{w}) = \sum_{\vec{x}_n \in M} -t_n \vec{x}_n$$

4

---

---

---

---

---

---

---

---

### Hinge Loss

(2D et 2 classes)

Neurone

Produit scalaire + fonction d'activation

Fonction de coût SVM (hinge loss) et gradient

$$E_D(\vec{w}) = \sum_{n=1}^N \max(0, 1 - t_n \vec{w}^T \vec{x}_n)$$

$$\nabla E_D(\vec{w}) = \sum_{\vec{x}_n \in M} -t_n \vec{x}_n$$

5

---

---

---

---

---

---

---

---

### Régression logistique

(2D, 2 classes)

Nouvelle fonction d'activation : **sigmoïde logistique**

Neurone

Fonction de coût (loss) et gradient

$$E_D(\vec{w}) = -\sum_{n=1}^N (t_n \ln(y_w(\vec{x}_n)) + (1-t_n) \ln(1-y_w(\vec{x}_n)))$$

$$\nabla E_D(\vec{w}) = \sum_{n=1}^N (y_w(\vec{x}_n) - t_n) \vec{x}_n$$

6

---

---

---

---

---

---

---

---

### Perceptron Multiclasse (2D et 3 classes)

The diagram shows a neural network with three input nodes:  $x_1$ ,  $x_2$ , and a bias node  $1$ . Each input node is connected to three hidden nodes representing weights  $\vec{w}_0^T \vec{x}$ ,  $\vec{w}_1^T \vec{x}$ , and  $\vec{w}_2^T \vec{x}$ . The output is  $\arg \max y_{w,j}(\vec{x})$ . To the right, a 2D plot shows three decision regions: blue for  $y_{w,0}(\vec{x})$  est max, green for  $y_{w,1}(\vec{x})$  est max, and red for  $y_{w,2}(\vec{x})$  est max.

$$y_W(\vec{x}) = W^T \vec{x}$$

|           |           |           |       |
|-----------|-----------|-----------|-------|
| $w_{0,0}$ | $w_{0,1}$ | $w_{0,2}$ | 1     |
| $w_{1,0}$ | $w_{1,1}$ | $w_{1,2}$ | $x_1$ |
| $w_{2,0}$ | $w_{2,1}$ | $w_{2,2}$ | $x_2$ |

7

---

---

---

---

---

---

---

---

---

---

### Perceptron Multiclasse

Exemple

A green star at coordinates (1.1, -2.0) is shown in the 2D plot, with a red arrow pointing to the red region labeled  $y_{w,2}(\vec{x})$  est max.

$$y_W(\vec{x}) = \begin{bmatrix} -2 & -3.6 & 0.5 \\ -4 & 2.4 & 4.1 \\ -6 & 4 & -4.9 \end{bmatrix} \begin{bmatrix} 1 \\ 1.1 \\ -2 \end{bmatrix} = \begin{bmatrix} -6.9 \\ -9.6 \\ 8.2 \end{bmatrix} \begin{matrix} \text{Classe 0} \\ \text{Classe 1} \\ \text{Classe 2} \end{matrix}$$

8

---

---

---

---

---

---

---

---

---

---

### Perceptron Multiclasse

Fonction de coût (**Perceptron loss**)

$$E_D(W) = \sum_{\vec{x}_n \in M} (\vec{w}_j^T \vec{x}_n - \vec{w}_l^T \vec{x}_n)$$

Annotations for the equation above:

- Summe sur l'ensemble des données mal classées (points to the summation symbol)
- Score de la mauvaise classe (points to  $\vec{w}_j^T \vec{x}_n$ )
- Score de la bonne classe (points to  $\vec{w}_l^T \vec{x}_n$ )

$$\nabla E_D(\vec{w}) = \sum_{\vec{x}_n \in M} \vec{x}_n$$

9

---

---

---

---

---

---

---

---

---

---

## Hinge Multiclasse

Fonction de coût (**Hinge loss**)

$$E_D(W) = \sum_{n=1}^N \sum_j \max(0, 1 + \vec{W}_j^T \vec{x}_n - \vec{W}_c^T \vec{x}_n)$$

Score de la mauvaise classe

Score de la bonne classe

$$\nabla E_D(\vec{W}) = \sum_{\vec{x}_n \in M} \vec{x}_n$$

10

---

---

---

---

---

---

---

---

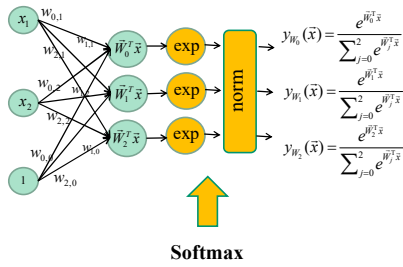
---

---

---

---

## Régression logistique multiclasse



11

---

---

---

---

---

---

---

---

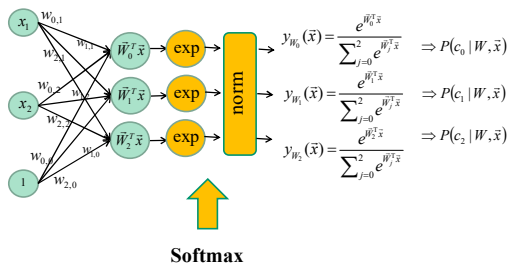
---

---

---

---

## Régression logistique multiclasse



12

---

---

---

---

---

---

---

---

---

---

---

---

### Régression logistique multiclasse

|            |  |   |
|------------|--|---|
| airplane   |  | 'airplane' $\Rightarrow t = [1000000000]$   |
| automobile |  | 'automobile' $\Rightarrow t = [0100000000]$ |
| bird       |  | 'bird' $\Rightarrow t = [0010000000]$       |
| cat        |  | 'cat' $\Rightarrow t = [0001000000]$        |
| deer       |  | 'deer' $\Rightarrow t = [0000100000]$       |
| dog        |  | 'dog' $\Rightarrow t = [0000010000]$        |
| frog       |  | 'frog' $\Rightarrow t = [0000001000]$       |
| horse      |  | 'horse' $\Rightarrow t = [0000000100]$      |
| ship       |  | 'ship' $\Rightarrow t = [0000000010]$       |
| truck      |  | 'truck' $\Rightarrow t = [0000000001]$      |

Étiquettes de classe : **one-hot vector**

13

---

---

---

---

---

---

---

---

---

---

---

---

### Régression logistique multiclasse

Fonction de coût est une **entropie croisée** (*cross entropy loss*)

$$E_D(W) = -\sum_{n=1}^N \sum_{k=1}^K t_{kn} \ln y_{W_k}(\bar{x}_n)$$

**Devoir 4**

$$\nabla E_D(W) = \frac{1}{N} \sum_{n=1}^N \bar{x}_n (y_W(\bar{x}_n) - t_{kn})$$

14

---

---

---

---

---

---

---

---

---

---

---

---

$\bar{x} = \begin{bmatrix} -15 \\ 22 \\ -44 \\ 56 \end{bmatrix}, t = 2$

|      |      |       |      |      |     |
|------|------|-------|------|------|-----|
| 0.0  | 0.01 | -0.05 | 0.1  | 0.05 | 1   |
| 0.2  | 0.7  | 0.2   | 0.05 | 0.16 | -15 |
| -0.3 | 0.0  | -0.45 | -0.2 | 0.03 | 22  |
|      |      |       |      |      | -44 |
|      |      |       |      |      | 56  |

$W$

**Score**

Hinge loss

$\max(0, -2.85 - 0.28 + 1) + \max(0, 0.86 - 0.28 + 1) = 1.58$

**Score**

Entropie croisée

Softmax:  $\frac{e^{-2.85}}{e^{-2.85} + e^{0.86} + e^{0.28}} = \frac{0.06}{2.36} = 0.025$

$-\ln(0.025) = 0.452$

15

---

---

---

---

---

---

---

---

---

---

---

---

## Maximum a posteriori

Régularisation

$$\arg \min_W = E_D(W) + \lambda R(W)$$

Fonction de perte      Constante  
Régularisation

En général L1 ou L2  $R(W) = \|W\|_1$  ou  $\|W\|_2$

16

---

---

---

---

---

---

---

---

## Optimisation

**Descente de gradient**

$$w^{[k+1]} = w^{[k]} - \eta^{[k]} \nabla_w E$$

↗ Gradient de la fonction de coût  
 ↘ Taux d'apprentissage ou "learning rate".

|   |  |
|---|--|
| <p><b>Descente de gradient stochastique</b></p> <p>Initialiser w<br/>k=0<br/>FAIRE k=k+1<br/>FOR n = 1 to N<br/>  <math>w = w - \eta^{[k]} \nabla E(\tilde{x}_n)</math></p> <p>JUSQU'À ce que toutes les données soient bien classées ou k=MAX_ITER</p> | <p><b>Optimisation par Batch</b></p> <p>Initialiser w<br/>k=0<br/>FAIRE k=k+1<br/>  <math>w = w - \eta^{[k]} \sum_i \nabla E(\tilde{x}_i)</math></p> <p>JUSQU'À ce que toutes les données soient bien classées ou k=MAX_ITER</p> |
|---|--|

Parfois  $\eta^{[k]} = cst / k$

17

---

---

---

---

---

---

---

---

Maintenant, rendons le réseau  
**profond**  
PROFOND  
Maintenant, rendons le réseau

18

---

---

---

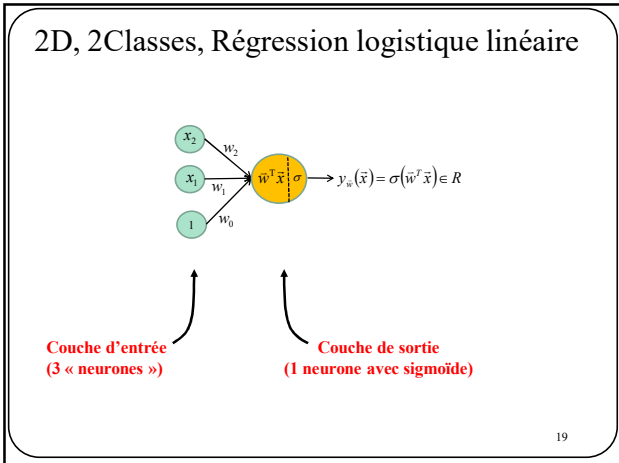
---

---

---

---

---




---

---

---

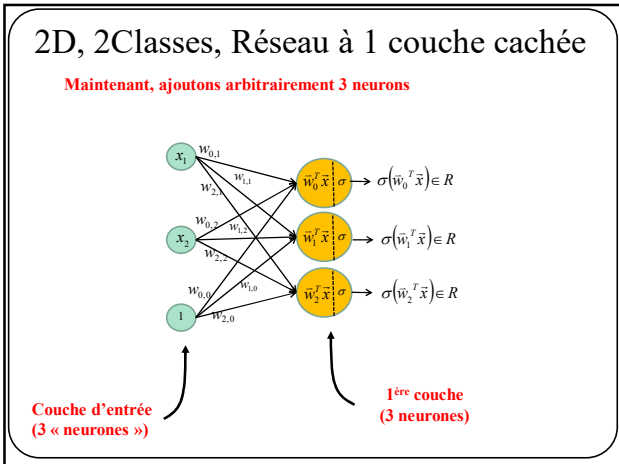
---

---

---

---

---




---

---

---

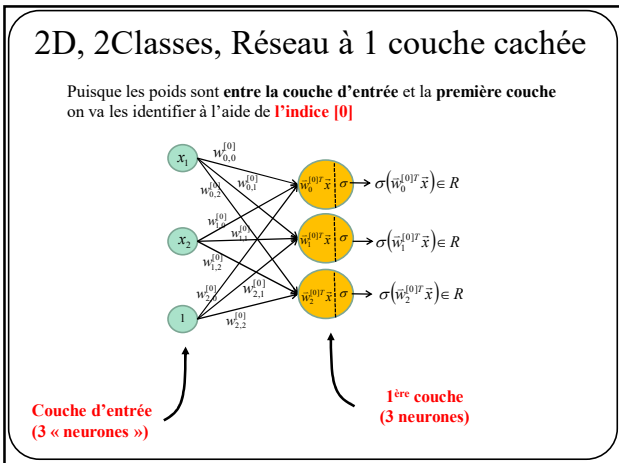
---

---

---

---

---




---

---

---

---

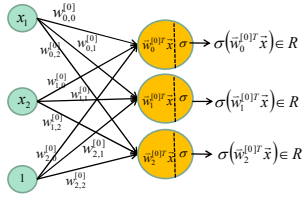
---

---

---

---

### 2D, 2Classes, Réseau à 1 couche cachée



**NOTE:** à la sortie de la première couche, on a **3 réels** calculés ainsi

$$\sigma \left( \begin{bmatrix} W_{0,0}^{[0]} & W_{0,1}^{[0]} & W_{0,2}^{[0]} \\ W_{1,0}^{[0]} & W_{1,1}^{[0]} & W_{1,2}^{[0]} \\ W_{2,0}^{[0]} & W_{2,1}^{[0]} & W_{2,2}^{[0]} \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ 1 \end{bmatrix} \right)$$

22

---

---

---

---

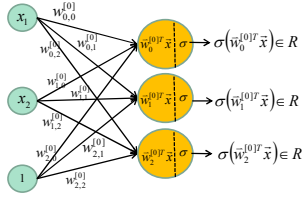
---

---

---

---

### 2D, 2Classes, Réseau à 1 couche cachée



**NOTE:** représentation plus simple de la sortie de la 1<sup>ère</sup> couche (**3 réels**)

$$\sigma(W^{[0]T} \vec{x})$$

23

---

---

---

---

---

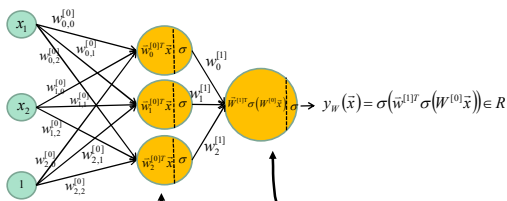
---

---

---

### 2D, 2Classes, Réseau à 1 couche cachée

Si on veut effectuer une **classification 2 classes** via une **régression logistique** (donc une **fonction coût par « entropie croisée »**) on doit ajouter un **neurone de sortie**.



Couche d'entrée  
(3 « neurones »)

1<sup>ère</sup> couche cachée  
(3 neurones)

Couche de sortie  
(1 neurone)

24

---

---

---

---

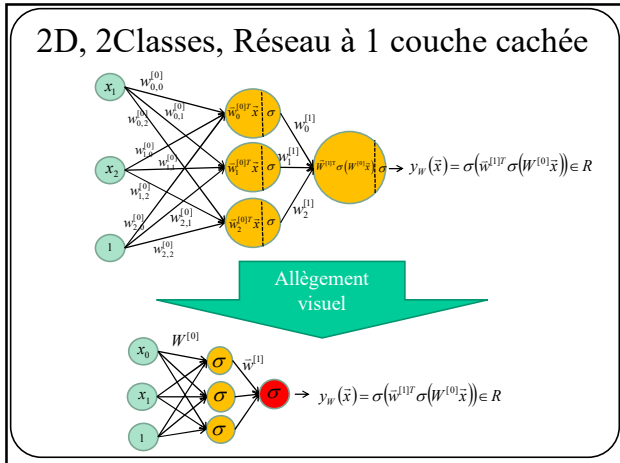
---

---

---

---





25

---

---

---

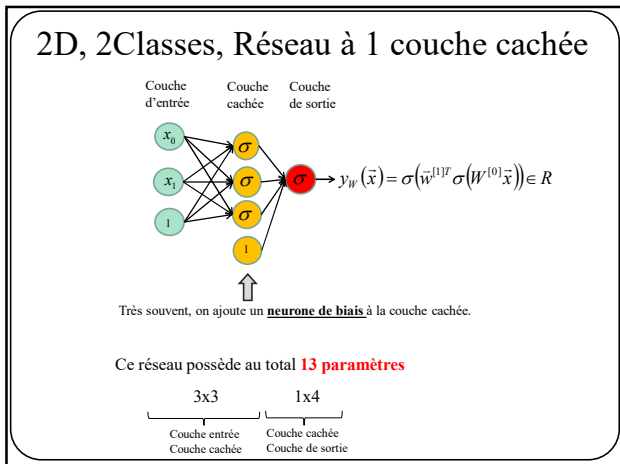
---

---

---

---

---



26

---

---

---

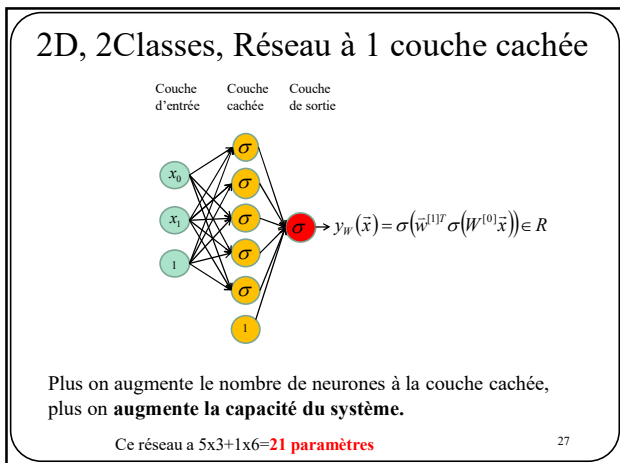
---

---

---

---

---



27

---

---

---

---

---

---

---

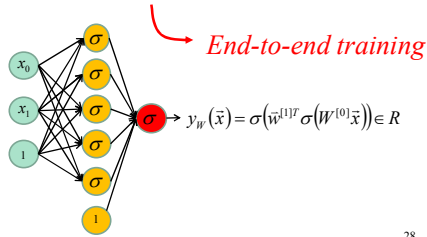
---

## NOTE Importante

Le but de la première couche est de **projeter les données d'entrée** (ici  $\vec{x} \in \mathbb{R}^2$ ) vers un espace dimensionnel plus grand (ici  $\sigma(W^{[0]}\vec{x}) \in \mathbb{R}^3$ ) là où les **classes sont linéairement séparables**.

Car il ne faut pas oublier que la **couche de sortie** est une **régression logistique linéaire**.

Par conséquent, au lieu de fixer nous même la fonction de base, on laisse le **réseau l'apprendre**.



28

28

---

---

---

---

---

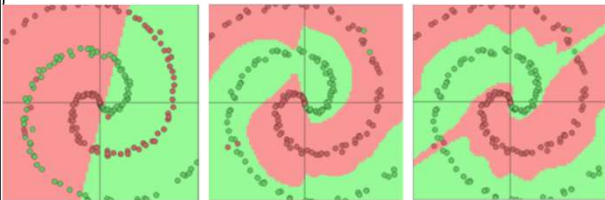
---

---

---

## Nombre de neurones VS Capacity

Aucun neurone caché    12 neurones cachés    60 neurones cachés



Classification linéaire **Underfitting** (pas assez de capacité)  
 Classification non linéaire **BON RESULTAT** (bonne capacité)  
 Classification non linéaire **Over fitting** (trop grande capacité)

<http://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html>

29

---

---

---

---

---

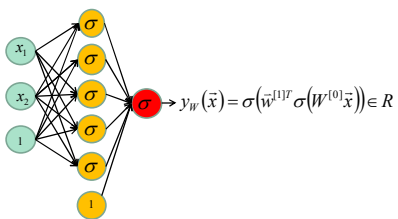
---

---

---

## 2D, 2Classes, Réseau à 1 couche cachée

Couche d'entrée    Couche cachée    Couche de sortie



30

30

---

---

---

---

---

---

---

---

### kD, 2Classes, Réseau à 1 couche cachée

Couche d'entrée    Couche cachée    Couche de sortie

$y_w(\vec{x}) = \sigma(\vec{w}^{[1]T} \sigma(W^{[0]}\vec{x})) \in R$

Au peut facilement augmenter la dimensionnalité des données d'entrée. Cela n'a pour effet que d'augmenter le nombre de colonnes dans  $W^{[0]}$

Ce réseau a  $5x(k+1)+1x6$  paramètres

31

31

---

---

---

---

---

---

---

---

### kD, 2Classes, Réseau à 2 couches cachées

Couche d'entrée    Couche cachée 1    Couche cachée 2    Couche de sortie

$W^{[0]} \in R^{5 \times k+1}$   
 $W^{[1]} \in R^{3 \times 6}$   
 $\vec{w}^{[2]} \in R^4$

$y_w(\vec{x}) = \sigma(\vec{w}^{[2]T} \sigma(W^{[1]} \sigma(W^{[0]}\vec{x})))$

En ajoutant une couche cachée, on ajoute une multiplication matricielle

Ce réseau a  $5x(k+1)+6x3 + 1x4$  paramètres

32

32

---

---

---

---

---

---

---

---

### kD, 2 Classes, Réseau à 4 couches cachées

Couche d'entrée    Couche cachée 1    Couche cachée 2    Couche cachée 3    Couche cachée 4    Couche de sortie

$W^{[0]} \in R^{5 \times k+1}$   
 $W^{[1]} \in R^{3 \times 6}$   
 $W^{[2]} \in R^{4 \times 4}$   
 $W^{[3]} \in R^{5 \times 7}$   
 $\vec{w}^{[4]} \in R^8$

$y_w(\vec{x}) = \sigma(\vec{w}^{[4]T} \sigma(W^{[3]} \sigma(W^{[2]} \sigma(W^{[1]} \sigma(W^{[0]}\vec{x}))))$

Ce réseau a  $5x(k+1)+6x3 + 4x4+7x5+1x8$  paramètres

33

33

---

---

---

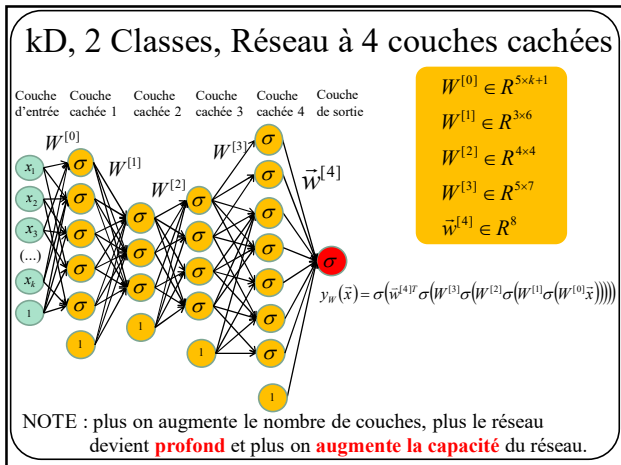
---

---

---

---

---



34

---

---

---

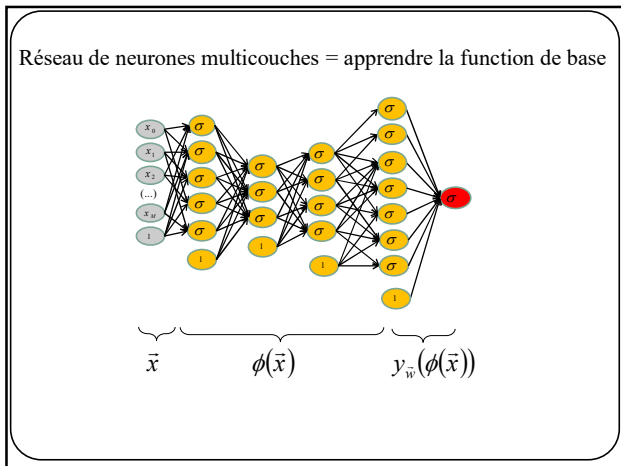
---

---

---

---

---



35

---

---

---

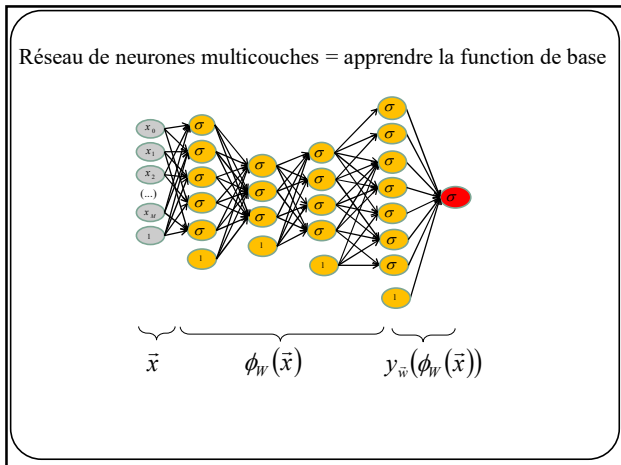
---

---

---

---

---



36

---

---

---

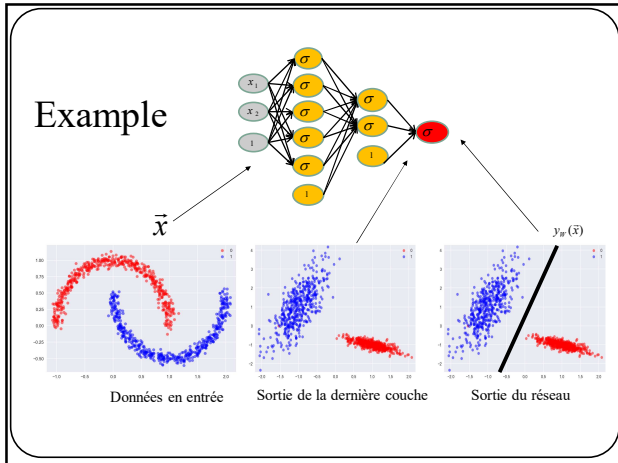
---

---

---

---

---



37

---

---

---

---

---

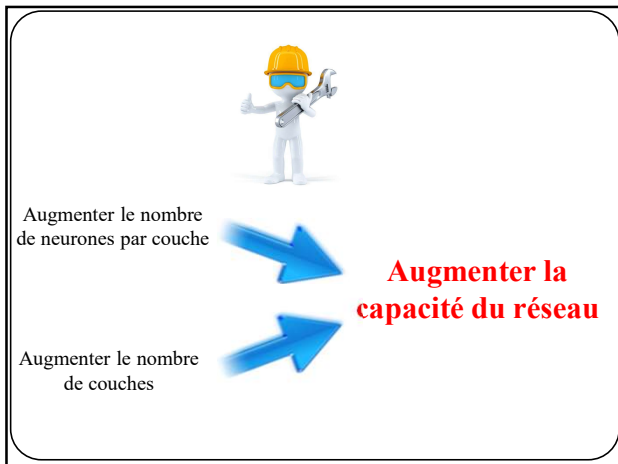
---

---

---

---

---



38

---

---

---

---

---

---

---

---

---

---



39

---

---

---

---

---

---

---

---

---

---



Lorsqu'un réseau doit prédire **plus de 2 classes**, on lui assigne **K neurones de sortie**, une par classe.

40

---

---

---

---

---

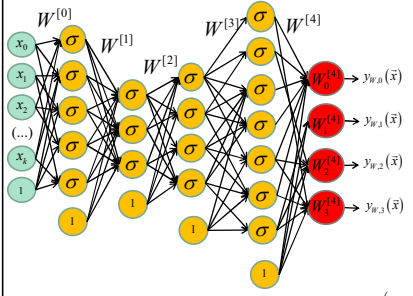
---

---

---

### kD, 4 Classes, Réseau à 4 couches cachées

Couche d'entrée    Couche cachée 1    Couche cachée 2    Couche cachée 3    Couche cachée 4    Couche de sortie



$W^{[0]} \in \mathbb{R}^{5 \times k+1}$   
 $W^{[1]} \in \mathbb{R}^{3 \times 6}$   
 $W^{[2]} \in \mathbb{R}^{4 \times 4}$   
 $W^{[3]} \in \mathbb{R}^{5 \times 7}$   
 $W^{[4]} \in \mathbb{R}^{4 \times 8}$

$$y_w(\vec{x}) = W^{[4]} \sigma(W^{[3]} \sigma(W^{[2]} \sigma(W^{[1]} \sigma(W^{[0]} \vec{x}))))$$

41

---

---

---

---

---

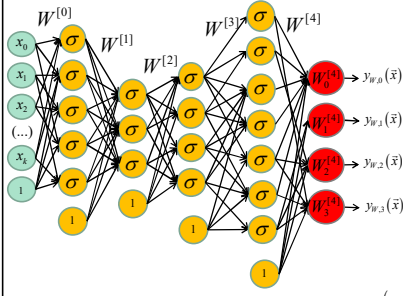
---

---

---

### kD, 4 Classes, Réseau à 4 couches cachées

Couche d'entrée    Couche cachée 1    Couche cachée 2    Couche cachée 3    Couche cachée 4    Couche de sortie



**Hinge loss**

$$y_w(\vec{x}) = W^{[4]} \sigma(W^{[3]} \sigma(W^{[2]} \sigma(W^{[1]} \sigma(W^{[0]} \vec{x}))))$$

42

---

---

---

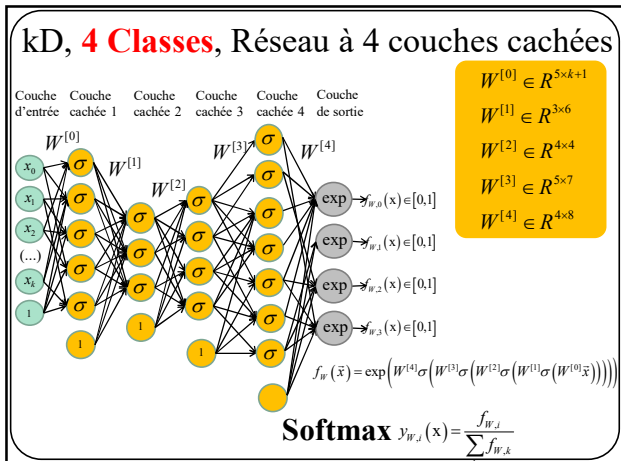
---

---

---

---

---



43

---

---

---

---

---

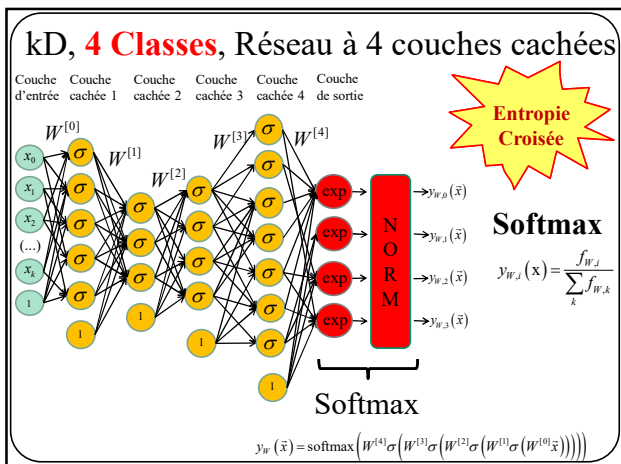
---

---

---

---

---



44

---

---

---

---

---

---

---

---

---

---

## Simulation

<http://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html>

45

---

---

---

---

---

---

---

---

---

---

## Comment faire une prédiction?

Ex.: faire transiter un signal de l'entrée à la sortie  
d'un réseau à 3 couches cachées

```
import numpy as np

def sigmoid(x):
    return 1.0 / (1.0+np.exp(-x))

x = np.insert(x,0,1) # Ajouter biais

H1 = sigmoid(np.dot(W0,x))
H1 = np.insert(H1,0,1) # Ajouter biais } Couche 1

H2 = sigmoid(np.dot(W1,H1))
H2 = np.insert(H2,0,1) # Ajouter biais } Couche 2

H2 = sigmoid(np.dot(W2,H1))
H2 = np.insert(H2,0,1) # Ajouter biais } Couche 3

y_pred = np.dot(W3,H2) } Couche sortie
```

Forward pass

46

---

---

---

---

---

---

---

---

## Comment optimiser les paramètres?

0- Partant de

$$W = \arg \min_W E_D(W) + \lambda R(W)$$

Trouver une fonction de régularisation. En général

$$R(W) = \|W\|_1 \text{ ou } \|W\|_2$$

47

47

---

---

---

---

---

---

---

---

## Comment optimiser les paramètres?

1- Trouver une loss  $E_D(W)$  comme par exemple

**Hinge loss**

**Entropie croisée (cross entropy)**



N'oubliez pas d'ajuster la sortie du réseau en fonction de la loss que vous aurez choisi.

*cross entropy => Softmax*

48

---

---

---

---

---

---

---

---



## Comment optimiser les paramètres?

2- Calculer le gradient de la loss par rapport à chaque paramètre

$$\frac{\partial (E_D(W) + \lambda R(W))}{\partial w_{a,b}^{[c]}}$$

et lancer un algorithme de descente de gradient pour mettre à jour les paramètres.

$$w_{a,b}^{[c]} = w_{a,b}^{[c]} - \eta \frac{\partial (E_D(W) + \lambda R(W))}{\partial w_{a,b}^{[c]}}$$

49

49

---

---

---

---

---

---

---

---

## Comment optimiser les paramètres?

$$\frac{\partial (E_D(W) + \lambda R(W))}{\partial w_{a,b}^{[c]}} \Rightarrow \text{calculé à l'aide d'une rétropropagation}$$

50

50

---

---

---

---

---

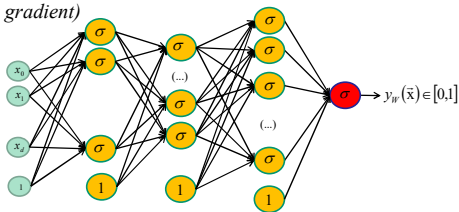
---

---

---

## Disparition du gradient

(vanishing gradient)



Malheureusement, l'entraînement d'un **réseau profond** avec **rétro-propagation** et des fonctions d'activations **sigmoïdales** entraîne des problèmes de

**disparition du gradient**

Démonstration au tableau

51

51

---

---

---

---

---

---

---

---

On résoud le problème de la disparition du gradient à l'aide **d'autres fonctions d'activations**

52

---

---

---

---

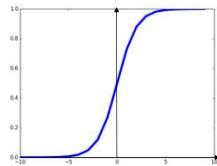
---

---

---

---

### Fonction d'activation



**Sigmoïde**

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

- Ramène les valeurs entre 0 et 1
- Historiquement populaire

#### 3 Problèmes :

- Un neurone saturé a pour effet de « **tuer** » les **gradients**

53

---

---

---

---

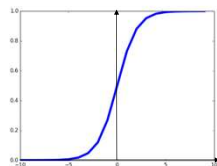
---

---

---

---

### Fonction d'activation



**Sigmoïde**

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

- Ramène les valeurs entre 0 et 1
- Historiquement populaire

#### 3 Problèmes :

- Un neurone saturé a pour effet de « **tuer** » les **gradients**
- Sortie d'une sigmoïde n'est **pas centrée à zéro**.

54

---

---

---

---

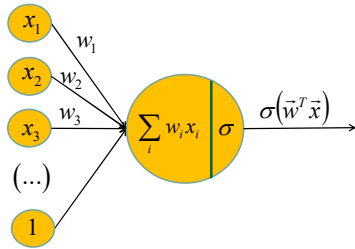
---

---

---

---

Qu'arrive-t-il lorsque le vecteur d'entrée  $\vec{x}$  d'un neurone est toujours positif?



Le gradient par rapport à  $\vec{w}$  est ... **Positif? Négatif?**

55

---

---

---

---

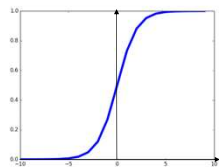
---

---

---

---

### Fonction d'activation



**Sigmoïde**

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

- Ramène les valeurs entre 0 et 1
- Historiquement populaire

**3 Problèmes :**

- Un neurone saturé a pour effet de « **tuer** » les **gradients**
- Sortie d'une sigmoïde n'est **pas centrée à zéro**.
- $\exp()$  est **coûteux** lorsque le nombre de neurones est élevé.

56

---

---

---

---

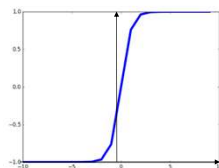
---

---

---

---

### Fonction d'activation



**Tanh(x)**

- Ramène les valeurs entre -1 et 1
- **Sortie centrée à zéro** 😊
- **Disparition du gradient** lorsque la fonction sature 😞

[LeCun et al., 1991]

57

---

---

---

---

---

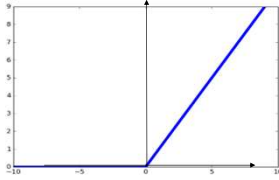
---

---

---

### Fonction d'activation

$\text{ReLU}(x) = \max(0, x)$



**ReLU(x)**  
(Rectified Linear Unit)

- Aucune **saturation** ☺
- Super **rapide** ☺
- Converge plus rapide** que sigmoïde/tanh (5 à 10x) ☺
- Sortie **non centrée à zéro** ☹
- Un inconvénient** : qu'arrive-t-il au gradient lorsque  $x < 0$ ? ☹

[Krizhevsky et al., 2012]

58

---

---

---

---

---

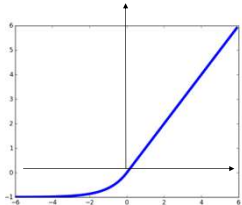
---

---

---

### Fonction d'activation

$\text{ELU}(x) = \begin{cases} x & \text{si } x > 0 \\ \alpha(e^x - 1) & \text{sinon} \end{cases}$



**ELU(x)**  
(Exponential Linear Unit)

- Tous les avantages de **ReLU** ☺
- Sortie plus « **centrée à zéro** » ☺
- Converge plus rapide** que sigmoïde/tanh (5 à 10x) ☺
- Gradients meurent plus lentement** ☺
- $\exp()$  est **coûteux** ☹

[Clevert et al., 2015]

59

---

---

---

---

---

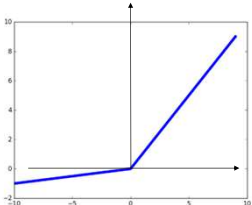
---

---

---

### Fonction d'activation

$\text{LReLU}(x) = \max(0.01x, x)$



**Leaky ReLU(x)**

- Aucune **saturation** ☺
- Super **rapide** ☺
- Converge plus rapide** que sigmoïde/tanh (5 à 10x) ☺
- Gradients ne meurent pas** ☺
- 0.01 est un **hyperparamètre** ☹

[Mass et al., 2013]  
[He et al., 2015]

60

---

---

---

---

---

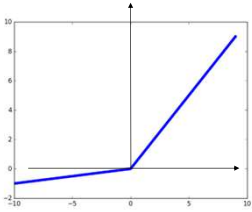
---

---

---

### Fonction d'activation

$\text{PReLU}(x) = \max(\alpha x, x)$



- Aucune **saturation** 😊
- Super **rapide** 😊
- **Converge plus rapide** que sigmoïde/tanh (5 à 10x) 😊
- **Gradients ne meurent pas** 😊
- $\alpha$  **appris** lors de la rétro-propagation 😊

**Parametric ReLU(x)**

[Mass et al., 2013]  
[He et al., 2015]

61

---

---

---

---

---

---

---

---

### En pratique

- Par défaut, le gens utilisent **ReLU**.
- Essayez **Leaky ReLU / PReLU / ELU**
- Essayez **tanh** mais n'attendez-vous pas à grand chose
- **Ne pas utiliser de sigmoïde** sauf à la sortie d'un réseau 2 classes.

62

---

---

---

---

---

---

---

---

### Les bonnes pratiques

63

---

---

---

---

---

---

---

---

## Optimisation

**Descente de gradient**

$$\mathbf{w}^{[k+1]} = \mathbf{w}^{[k]} - \eta^{[k]} \nabla E$$

↗ Gradient de la fonction de coût  
 ↘ Taux d'apprentissage ou "learning rate".

Descente de gradient stochastique

```

Initialiser w
k=0
FAIRE k=k+1
  FOR n = 1 to N
    w = w - η[k] ∇E(xn)
  
```

JUSQU'À ce que toutes les données sont bien classées ou k=MAX\_ITER

Optimisation par Batch

```

Initialiser w
k=0
FAIRE k=k+1
  w = w - η[k] ∑i ∇E(xi)
  
```

JUSQU'À ce que toutes les données sont bien classées ou k=MAX\_ITER

Parfois  $\eta^{[k]} = cst / k$

---

---

---

---

---

---

---

---

---

---

64

## Optimisation

**Descente de gradient**

$$\mathbf{w}^{[k+1]} = \mathbf{w}^{[k]} - \eta^{[k]} \nabla E$$

↗ Gradient de la fonction de coût  
 ↘ Taux d'apprentissage ou "learning rate".

Optimisation par **mini-batch**

```

Initialiser w
k=0
FAIRE k=k+1
  FAIRE n=0 à N par sauts de MBS /*Mini-batch size*/
    w = w - η[k] ∑i=nn+MBS ∇E(xi)
  
```

} Itération

JUSQU'À ce que toutes les données sont bien classées ou k=MAX\_ITER

---

---

---

---

---

---

---

---

---

---

65

## Optimisation

**Descente de gradient**

$$\mathbf{w}^{[k+1]} = \mathbf{w}^{[k]} - \eta^{[k]} \nabla E$$

↗ Gradient de la fonction de coût  
 ↘ Taux d'apprentissage ou "learning rate".

Optimisation par **mini-batch**

```

Initialiser w
k=0
FAIRE k=k+1
  FAIRE n=0 à N par sauts de MBS /*Mini-batch size*/
    w = w - η[k] ∑i=nn+MBS ∇E(xi)
  
```

} Epoch

JUSQU'À ce que toutes les données sont bien classées ou k=MAX\_ITER

---

---

---

---

---

---

---

---

---

---

66

Propagation avant pour une donnée (7 étapes)

|   |                        |
|---|------------------------|
| $\vec{x}$                                   | $\in \mathbb{R}^4$     |
| $W_1 \vec{x}$                               | $\in \mathbb{R}^5$     |
| $h(W_1 \vec{x})$                            | $\in \mathbb{R}^{5+1}$ |
| $W_2 h(W_1 \vec{x})$                        | $\in \mathbb{R}^3$     |
| $h(W_2 h(W_1 \vec{x}))$                     | $\in \mathbb{R}^{3+1}$ |
| $\vec{W}_3 (h(W_2 h(W_1 \vec{x})))$         | $\in \mathbb{R}$       |
| $\sigma(\vec{W}_3 (h(W_2 h(W_1 \vec{x}))))$ | $\in \mathbb{R}$       |

67

---

---

---

---

---

---

---

---

Propagation avant pour un batch de 3 données (7 étapes)

POUR  $i$  allant de 1 à 3

|  |                        |
|--|------------------------|
| $\vec{x}_i = X[i]$                                   | $\in \mathbb{R}^{M+1}$ |
| $W_1 \vec{x}_i$                                      | $\in \mathbb{R}^5$     |
| $h(W_1 \vec{x}_i)$                                   | $\in \mathbb{R}^{5+1}$ |
| $W_2 h(W_1 \vec{x}_i)$                               | $\in \mathbb{R}^3$     |
| $h(W_2 h(W_1 \vec{x}_i))$                            | $\in \mathbb{R}^{3+1}$ |
| $\vec{W}_3 (h(W_2 h(W_1 \vec{x}_i)))$                | $\in \mathbb{R}$       |
| $Y[i] = \sigma(\vec{W}_3 (h(W_2 h(W_1 \vec{x}_i))))$ | $\in \mathbb{R}$       |

Solution naïve et peu efficace

TP4

68

---

---

---

---

---

---

---

---

### Solution

Il est plus efficace d'effectuer UNE multiplication matricielle que PLUSIEURS produits scalaires (exemple de la 6<sup>e</sup> étape)

|   |   |
|---|---|
| $(w_1 \ w_2 \ w_3) \begin{pmatrix} a \\ b \\ c \end{pmatrix} = (w_1 a + w_2 b + w_3 c)$ | $\begin{pmatrix} w_1 a + w_2 b + w_3 c \\ w_1 d + w_2 e + w_3 f \\ w_1 g + w_2 h + w_3 i \end{pmatrix} = Y$ |
| $(w_1 \ w_2 \ w_3) \begin{pmatrix} d \\ e \\ f \end{pmatrix} = (w_1 d + w_2 e + w_3 f)$ |   |
| $(w_1 \ w_2 \ w_3) \begin{pmatrix} g \\ h \\ i \end{pmatrix} = (w_1 g + w_2 h + w_3 i)$ |   |

TROIS produits scalaires

69

---

---

---

---

---

---

---

---

## Solution

Il est plus efficace d'effectuer **UNE multiplication matricielle** que **PLUSIEURS** produits scalaires (exemple de la 6<sup>e</sup> étape)

$$\begin{bmatrix} w_1 & w_2 & w_3 \end{bmatrix} \begin{pmatrix} a & d & g \\ b & e & h \\ c & f & i \end{pmatrix} = \begin{pmatrix} w_1 a + w_2 b + w_3 c \\ w_1 d + w_2 e + w_3 f \\ w_1 g + w_2 h + w_3 i \end{pmatrix} = Y$$



70

---

---

---

---

---

---

---

---

## Solution

Il est plus efficace d'effectuer **UNE multiplication matricielle** que **PLUSIEURS** matrice-vecteur (exemple de la 1<sup>e</sup> étape, batch de 3)

$$\begin{matrix} W_1 \vec{x}_1 = \begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} \\ w_{21} & w_{22} & w_{23} & w_{24} \\ w_{31} & w_{32} & w_{33} & w_{34} \\ w_{41} & w_{42} & w_{43} & w_{44} \\ w_{51} & w_{52} & w_{53} & w_{54} \end{bmatrix} \begin{pmatrix} a \\ b \\ c \\ d \end{pmatrix} \\ W_2 \vec{x}_2 = \begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} \\ w_{21} & w_{22} & w_{23} & w_{24} \\ w_{31} & w_{32} & w_{33} & w_{34} \\ w_{41} & w_{42} & w_{43} & w_{44} \\ w_{51} & w_{52} & w_{53} & w_{54} \end{bmatrix} \begin{pmatrix} e \\ f \\ g \\ h \end{pmatrix} \\ W_3 \vec{x}_3 = \begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} \\ w_{21} & w_{22} & w_{23} & w_{24} \\ w_{31} & w_{32} & w_{33} & w_{34} \\ w_{41} & w_{42} & w_{43} & w_{44} \\ w_{51} & w_{52} & w_{53} & w_{54} \end{bmatrix} \begin{pmatrix} i \\ j \\ k \\ l \end{pmatrix} \end{matrix}$$



71

---

---

---

---

---

---

---

---

## Solution

Il est plus efficace d'effectuer **UNE multiplication matricielle** que **PLUSIEURS** matrice-vecteur (exemple de la 1<sup>e</sup> étape)

$$W_1 X = \begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} \\ w_{21} & w_{22} & w_{23} & w_{24} \\ w_{31} & w_{32} & w_{33} & w_{34} \\ w_{41} & w_{42} & w_{43} & w_{44} \\ w_{51} & w_{52} & w_{53} & w_{54} \end{bmatrix} \begin{pmatrix} a & d & h \\ b & e & i \\ c & f & j \\ d & g & k \end{pmatrix} = \begin{bmatrix} u_1 & v_1 & z_1 \\ u_2 & v_2 & z_2 \\ u_3 & v_3 & z_3 \\ u_4 & v_4 & z_4 \\ u_5 & v_5 & z_5 \end{bmatrix}$$



72

---

---

---

---

---

---

---

---



## Vectorisation de la propagation avant

En résumé, lorsqu'on propage une « batch »

|                    |                           |  |
|--------------------|---------------------------|--|
| Au niveau neuronal | Multi.<br>Vecteur-Matrice | $\vec{W}X = [w_1 \ w_2 \ w_3] \begin{pmatrix} a & d & g \\ b & e & h \\ c & f & i \end{pmatrix}$ |
|--------------------|---------------------------|--|

|                        |                           |   |
|------------------------|---------------------------|---|
| Au niveau de la couche | Multi.<br>Matrice-Matrice | $WX = \begin{pmatrix} w_{11} & w_{12} & w_{13} & w_{14} \\ w_{21} & w_{22} & w_{23} & w_{24} \\ w_{31} & w_{32} & w_{33} & w_{34} \\ w_{41} & w_{42} & w_{43} & w_{44} \\ w_{51} & w_{52} & w_{53} & w_{54} \end{pmatrix} \begin{pmatrix} a & d & h \\ b & e & i \\ c & f & j \\ d & g & k \end{pmatrix}$ |
|------------------------|---------------------------|---|

73

---

---

---

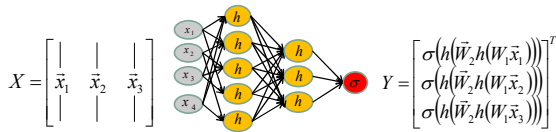
---

---

---

---

---



Vectoriser la rétropropagation

74

---

---

---

---

---

---

---

---

## Vectoriser la rétropropagation

Exemple simple pour un neurone et une batch de 3

$$\begin{matrix} [w_1 & w_2 & w_3] \\ \vec{W} & X & Y \end{matrix} \begin{pmatrix} a & d & g \\ b & e & h \\ c & f & i \end{pmatrix} = \begin{pmatrix} w_1 a + w_2 b + w_3 c \\ w_1 d + w_2 e + w_3 f \\ w_1 g + w_2 h + w_3 i \end{pmatrix}^T$$

En supposant qu'on connaît le gradient pour les 3 éléments de Y provenant de sortie du réseau, comment faire pour propager le gradient vers W?

75

---

---

---

---

---

---

---

---

# Vectoriser la rétropropagation

Exemple simple pour **1 neurone et une batch de 3**

$$\begin{matrix} [w_1 & w_2 & w_3] \\ W \end{matrix} \begin{matrix} \begin{pmatrix} a & d & g \\ b & e & h \\ c & f & i \end{pmatrix} \\ X \end{matrix} = \begin{matrix} \begin{pmatrix} w_1 a + w_2 b + w_3 c \\ w_1 d + w_2 e + w_3 f \\ w_1 g + w_2 h + w_3 i \end{pmatrix} \\ Y \end{matrix}^T$$

Rappelons que l'objectif est de faire une **descente de gradient**, i.e.

$$w_1 \leftarrow w_1 - \eta \frac{\partial E}{\partial w_1} \quad w_2 \leftarrow w_2 - \eta \frac{\partial E}{\partial w_2} \quad w_3 \leftarrow w_3 - \eta \frac{\partial E}{\partial w_3}$$

76

---

---

---

---

---

---

---

---

---

---

$$\begin{matrix} [w_1 & w_2 & w_3] \\ W \end{matrix} \begin{matrix} \begin{pmatrix} a & d & g \\ b & e & h \\ c & f & i \end{pmatrix} \\ X \end{matrix} = \begin{matrix} \begin{pmatrix} w_1 a + w_2 b + w_3 c \\ w_1 d + w_2 e + w_3 f \\ w_1 g + w_2 h + w_3 i \end{pmatrix} \\ Y \end{matrix}^T$$

Concentrons-nous sur  $w_1$

$$w_1 \leftarrow w_1 - \eta \frac{\partial E}{\partial w_1}$$

$$w_1 \leftarrow w_1 - \eta \frac{\partial E}{\partial Y} \frac{\partial Y}{\partial w_1} \quad (\text{par propriété de la dérivée en chaîne})$$

$$w_1 \leftarrow w_1 - \eta \left[ \frac{\partial E_1}{\partial Y} \quad \frac{\partial E_2}{\partial Y} \quad \frac{\partial E_3}{\partial Y} \right] \begin{bmatrix} a \\ d \\ g \end{bmatrix} \quad (\text{provient de la rétro-propagation})$$

$$w_1 \leftarrow w_1 - \eta \left( \frac{\partial E_1}{\partial Y} a + \frac{\partial E_2}{\partial Y} d + \frac{\partial E_3}{\partial Y} g \right)$$

77

---

---

---

---

---

---

---

---

---

---

$$\begin{matrix} [w_1 & w_2 & w_3] \\ W \end{matrix} \begin{matrix} \begin{pmatrix} a & d & g \\ b & e & h \\ c & f & i \end{pmatrix} \\ X \end{matrix} = \begin{matrix} \begin{pmatrix} w_1 a + w_2 b + w_3 c \\ w_1 d + w_2 e + w_3 f \\ w_1 g + w_2 h + w_3 i \end{pmatrix} \\ Y \end{matrix}^T$$

Et pour tous les poids

$$\begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}^T \leftarrow \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}^T - \eta \left[ \frac{\partial E_1}{\partial Y} \quad \frac{\partial E_2}{\partial Y} \quad \frac{\partial E_3}{\partial Y} \right] \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix}$$

$$\vec{W}^T \leftarrow \vec{W}^T - \eta \frac{\partial E}{\partial Y} \begin{bmatrix} \partial Y_1 / \partial w_1 & \partial Y_1 / \partial w_2 & \partial Y_1 / \partial w_3 \\ \partial Y_2 / \partial w_1 & \partial Y_2 / \partial w_2 & \partial Y_2 / \partial w_3 \\ \partial Y_3 / \partial w_1 & \partial Y_3 / \partial w_2 & \partial Y_3 / \partial w_3 \end{bmatrix}$$

$$\vec{W}^T \leftarrow \vec{W}^T - \eta \frac{\partial E}{\partial Y} \frac{\partial Y}{\partial \vec{W}} \quad \text{Matrice jacobienne}$$

78

---

---

---

---

---

---

---

---

---

---

$$\begin{matrix}
 W_{11} & W_{12} & W_{13} & W_{14} \\
 W_{21} & W_{22} & W_{23} & W_{24} \\
 W_{31} & W_{32} & W_{33} & W_{34} \\
 W_{41} & W_{42} & W_{43} & W_{44} \\
 W_{51} & W_{52} & W_{53} & W_{54}
 \end{matrix}
 \begin{matrix}
 a d h \\
 b e i \\
 c f j \\
 d g k
 \end{matrix}
 =
 \begin{matrix}
 u_1 v_1 z_1 \\
 u_2 v_2 z_2 \\
 u_3 v_3 z_3 \\
 u_4 v_4 z_4 \\
 u_5 v_5 z_5
 \end{matrix}$$

Même chose pour 1 couche et une batch de 3

$$W \leftarrow W^T - \eta \frac{\partial E}{\partial Y} \begin{matrix} \frac{\partial E}{\partial Y_{11}} & \frac{\partial E}{\partial Y_{12}} & \frac{\partial E}{\partial Y_{13}} \\ \frac{\partial E}{\partial Y_{21}} & \frac{\partial E}{\partial Y_{22}} & \frac{\partial E}{\partial Y_{23}} \\ \frac{\partial E}{\partial Y_{31}} & \frac{\partial E}{\partial Y_{32}} & \frac{\partial E}{\partial Y_{33}} \\ \frac{\partial E}{\partial Y_{41}} & \frac{\partial E}{\partial Y_{42}} & \frac{\partial E}{\partial Y_{43}} \\ \frac{\partial E}{\partial Y_{51}} & \frac{\partial E}{\partial Y_{52}} & \frac{\partial E}{\partial Y_{53}} \end{matrix} \begin{matrix} a & b & c & d \\ d & e & f & g \\ h & i & j & k \end{matrix}$$

$$W^T \leftarrow W^T - \eta \frac{\partial E}{\partial Y} \frac{\partial Y}{\partial W}$$

79

---

---

---

---

---

---

---

---

---

---

## Vectorisation de la rétro-propagation

En résumé, lorsqu'on rétro-propage un gradient d'une batch

|                    |                        |   |
|--------------------|------------------------|---|
| Au niveau neuronal | Multi. Vecteur-Matrice | $\bar{W}^T \leftarrow \bar{W}^T - \eta \frac{\partial \bar{E}}{\partial Y} \frac{\partial Y}{\partial \bar{W}}$ $\bar{W}^T \leftarrow \bar{W}^T - \eta \frac{\partial \bar{E}}{\partial Y} X^T$ |
|--------------------|------------------------|---|

|                        |                        |   |
|------------------------|------------------------|---|
| Au niveau de la couche | Multi. Matrice-Matrice | $W^T \leftarrow W^T - \eta \frac{\partial E}{\partial Y} \frac{\partial Y}{\partial W}$ $W^T \leftarrow W^T - \eta \frac{\partial E}{\partial Y} X^T$ |
|------------------------|------------------------|---|

80

---

---

---

---

---

---

---

---

---

---

## Pour plus de détails:

<https://medium.com/datathings/vectorized-implementation-of-back-propagation-1011884df84>  
<https://peterroelants.github.io/posts/neural-network-implementation-part04/>

81

---

---

---

---

---

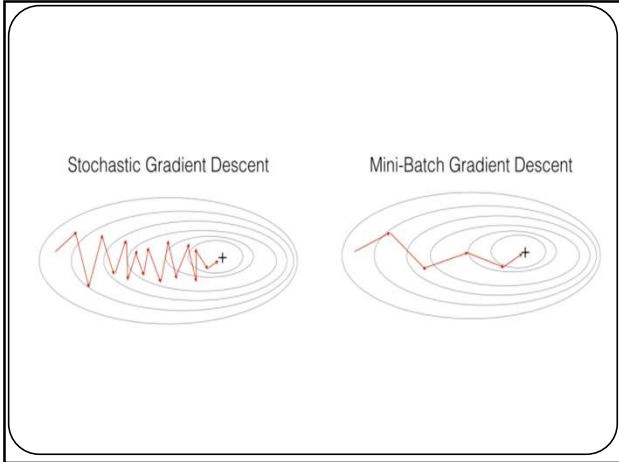
---

---

---

---

---



82

---

---

---

---

---

---

---

---

Comment initialiser un réseau de neurones?

$W = ?$

83

---

---

---

---

---

---

---


---

**Initialisation**

**Première idée:** faibles valeurs aléatoires  
(Gaussienne  $\mu = 0, \sigma = 0.01$ )

`W_i=0.01*np.random.randn(H_i,H_im1)`

Fonctionne bien pour de petits réseaux mais  
**pas pour des réseaux profonds.**

 E.g. réseau à 10 couches avec 500 neurones par couche et des **tanh** comme fonctions d'activation.

84

---

---

---

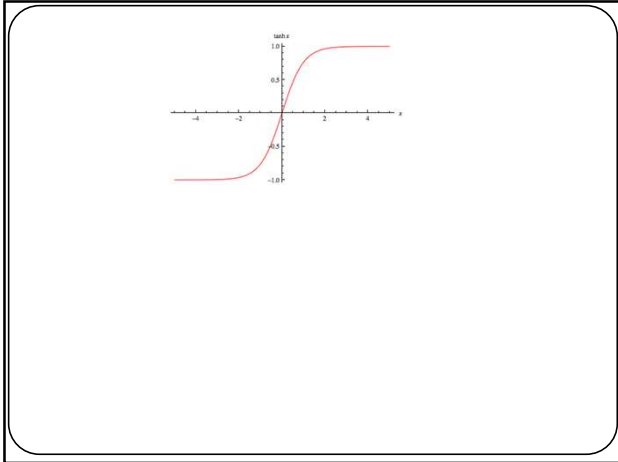
---

---

---

---

---



85

---

---

---

---

---

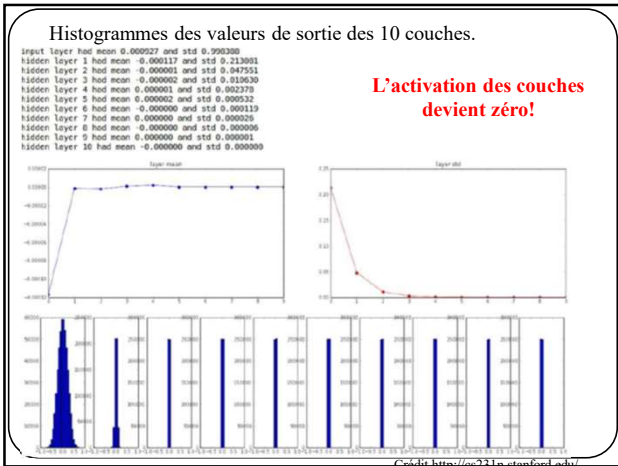
---

---

---

---

---



86

---

---

---

---

---

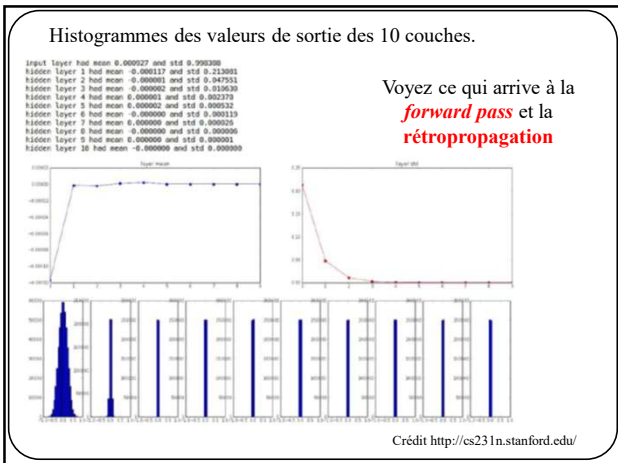
---

---

---

---

---



87

---

---

---

---

---

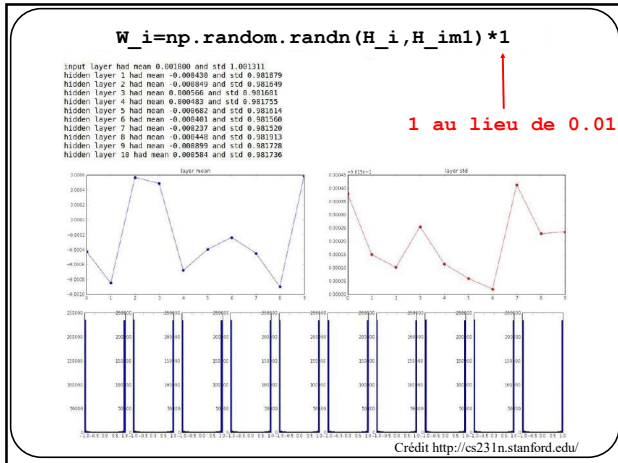
---

---

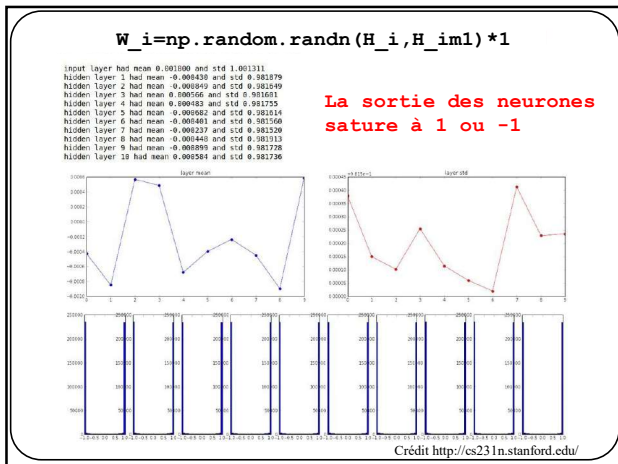
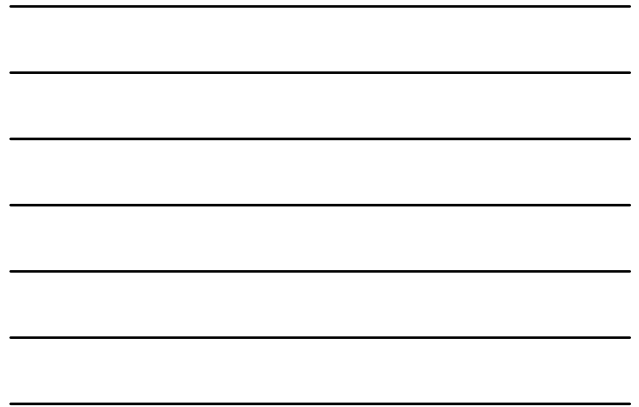
---

---

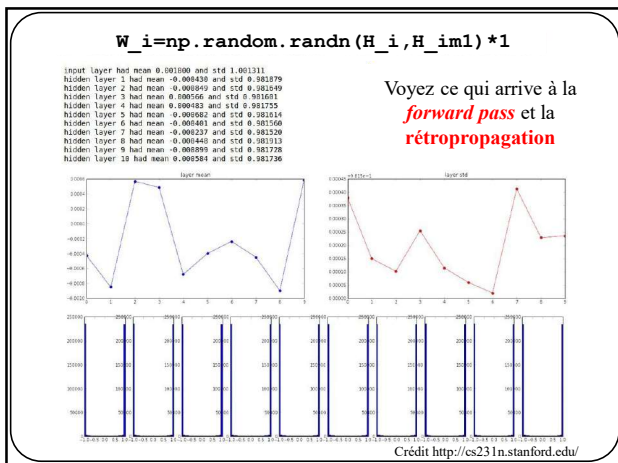
---



88

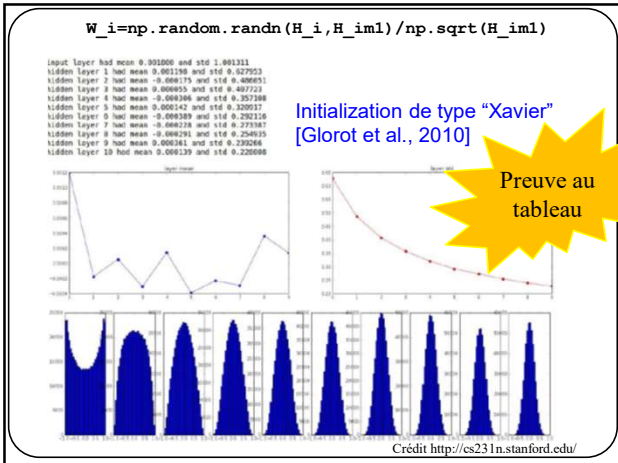


89



90






---

---

---

---

---

---

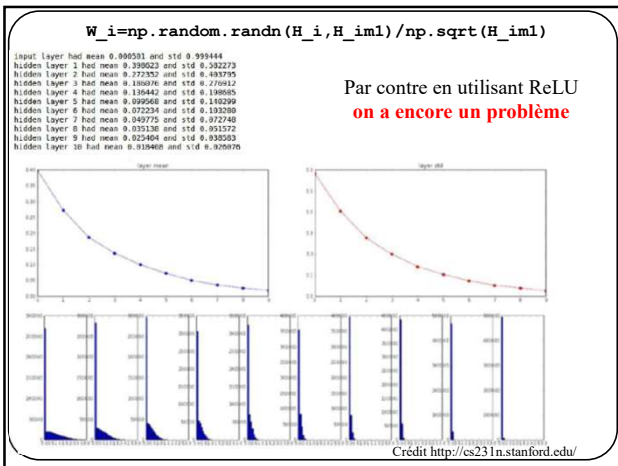
---

---

---

---

91




---

---

---

---

---

---

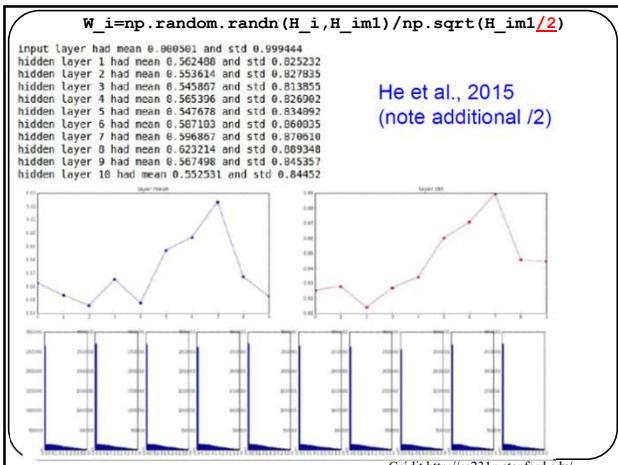
---

---

---

---

92




---

---

---

---

---

---

---

---

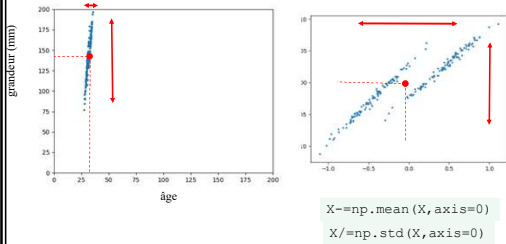
---

---

93

## Prétraitement des données

### Centrer et normaliser les données d'entrée



94

---

---

---

---

---

---

---

---

---

---

## Sanity checks

1. Toujours s'assurer qu'une initialization aléatoire donne une **perte (loss) maximale**

Exemple : pour le cas **10 classes**, une **régularisation à 0** et une **entropie croisée**.

$$E_D(\mathbf{W}) = -\frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K t_{kn} \ln y_{W,k}(\bar{x}_n)$$

Si l'initialisation est aléatoire, alors la probabilité sera égale pour chaque classe

$$\begin{aligned} E_D(\mathbf{W}) &= -\frac{1}{N} \sum_{n=1}^N \ln \frac{1}{10} \\ &= \ln(10) \\ &= 2.30 \end{aligned}$$

95

---

---

---

---

---

---

---

---

---

---

## Sanity checks

1. Toujours s'assurer qu'une initialization aléatoire donne une **perte (loss) maximale**

Exemple : pour le cas **10 classes**, une **régularisation à 0** et une **entropie croisée**.

```
def init_two_layer_model(input_size, hidden_size, output_size):  
    # initialize a model  
    model = {}  
    model['w1'] = 0.0001 * np.random.randn(input_size, hidden_size)  
    model['b1'] = np.zeros(hidden_size)  
    model['w2'] = 0.0001 * np.random.randn(hidden_size, output_size)  
    model['b2'] = np.zeros(output_size)  
    return model
```

```
model = init_two_layer_model(27*32*3, 50, 10) # input_size, hidden_size, number of classes  
loss, grad = two_layer_net(X_train, model, y_train) # disable regularization  
print loss  
2.30261216167
```

loss ~2.3  
"correct" for  
10 classes

returns the loss and the  
gradient for all parameters

Credit: <http://cs231n.stanford.edu>

96

---

---

---

---

---

---

---

---

---

---



## Sanity checks

- Et lorsqu'on **augmente la régularisation**, la perte augmente aussi

```
def init_two_layer_model(input_size, hidden_size, output_size):  
    # hidden_size = 20  
    model = {}  
    model['W1'] = 0.0001 * np.random.randn(input_size, hidden_size)  
    model['b1'] = np.zeros(hidden_size)  
    model['W2'] = 0.0001 * np.random.randn(hidden_size, output_size)  
    model['b2'] = np.zeros(output_size)  
    return model  
  
model = init_two_layer_model(12*12*3, 50, 10) # input_size, hidden_size, number of classes  
loss, grad = two_layer_net(x_train, model, y_train, 1e3) # crank up regularization  
print loss  
3.06859716482
```

loss went up, good. (sanity check)

Credit <http://cs231n.stanford.edu/>

97

---

---

---

---

---

---

---

---

---

---

## Sanity checks

- Toujours s'assurer qu'on peut « **over-fitter** » sur un **petit nombre de données**.

Lets try to train now...

Tip: Make sure that you can overfit very small portion of the training data

Very small loss, train accuracy 1.00, nice!

```
model = init_two_layer_model(12*12*3, 50, 10) # input_size, hidden_size, number of classes  
trainer = ClassifierTrainer()  
x_train = X_train[:100] # take 10% examples  
y_train = Y_train[:100]  
best_model, stats = trainer.train(x_train, y_train, X_val, Y_val, 0.100,  
                                num_epochs=100, lr=0.0001,  
                                optimizer='sgd', learning_rate_decay=0.9,  
                                num_hiddens=1, verbose=True)  
  
Finished epoch 1 / 100: cost 1.36353, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 2 / 100: cost 2.38238, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 3 / 100: cost 1.26149, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 4 / 100: cost 2.38136, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 5 / 100: cost 1.36084, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 6 / 100: cost 2.37984, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 7 / 100: cost 1.35935, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 8 / 100: cost 2.37785, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 9 / 100: cost 1.35786, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 10 / 100: cost 2.37587, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 11 / 100: cost 1.35637, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 12 / 100: cost 2.37388, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 13 / 100: cost 1.35488, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 14 / 100: cost 2.37189, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 15 / 100: cost 1.35339, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 16 / 100: cost 2.36990, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 17 / 100: cost 1.35190, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 18 / 100: cost 2.36791, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 19 / 100: cost 1.35041, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 20 / 100: cost 2.36592, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 21 / 100: cost 1.34892, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 22 / 100: cost 2.36393, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 23 / 100: cost 1.34743, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 24 / 100: cost 2.36194, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 25 / 100: cost 1.34594, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 26 / 100: cost 2.35995, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 27 / 100: cost 1.34445, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 28 / 100: cost 2.35796, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 29 / 100: cost 1.34296, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 30 / 100: cost 2.35597, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 31 / 100: cost 1.34147, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 32 / 100: cost 2.35398, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 33 / 100: cost 1.33998, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 34 / 100: cost 2.35199, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 35 / 100: cost 1.33849, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 36 / 100: cost 2.34999, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 37 / 100: cost 1.33700, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 38 / 100: cost 2.34800, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 39 / 100: cost 1.33551, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 40 / 100: cost 2.34601, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 41 / 100: cost 1.33402, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 42 / 100: cost 2.34402, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 43 / 100: cost 1.33253, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 44 / 100: cost 2.34203, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 45 / 100: cost 1.33104, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 46 / 100: cost 2.34004, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 47 / 100: cost 1.32955, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 48 / 100: cost 2.33805, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 49 / 100: cost 1.32806, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 50 / 100: cost 2.33606, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 51 / 100: cost 1.32657, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 52 / 100: cost 2.33407, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 53 / 100: cost 1.32508, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 54 / 100: cost 2.33208, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 55 / 100: cost 1.32359, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 56 / 100: cost 2.33009, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 57 / 100: cost 1.32210, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 58 / 100: cost 2.32810, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 59 / 100: cost 1.32061, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 60 / 100: cost 2.32611, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 61 / 100: cost 1.31912, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 62 / 100: cost 2.32412, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 63 / 100: cost 1.31763, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 64 / 100: cost 2.32213, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 65 / 100: cost 1.31614, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 66 / 100: cost 2.32014, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 67 / 100: cost 1.31465, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 68 / 100: cost 2.31815, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 69 / 100: cost 1.31316, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 70 / 100: cost 2.31616, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 71 / 100: cost 1.31167, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 72 / 100: cost 2.31417, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 73 / 100: cost 1.31018, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 74 / 100: cost 2.31218, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 75 / 100: cost 1.30869, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 76 / 100: cost 2.31019, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 77 / 100: cost 1.30720, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 78 / 100: cost 2.30820, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 79 / 100: cost 1.30571, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 80 / 100: cost 2.30621, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 81 / 100: cost 1.30422, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 82 / 100: cost 2.30422, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 83 / 100: cost 1.30273, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 84 / 100: cost 2.30223, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 85 / 100: cost 1.30124, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 86 / 100: cost 2.30024, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 87 / 100: cost 1.29975, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 88 / 100: cost 2.29825, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 89 / 100: cost 1.29825, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 90 / 100: cost 2.29626, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 91 / 100: cost 1.29676, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 92 / 100: cost 2.29427, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 93 / 100: cost 1.29527, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 94 / 100: cost 2.29228, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 95 / 100: cost 1.29378, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 96 / 100: cost 2.29029, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 97 / 100: cost 1.29229, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 98 / 100: cost 2.28830, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 99 / 100: cost 1.29080, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished epoch 100 / 100: cost 2.28631, train 0.00000, val 0.00000, lr 1.00000e-03  
Finished optimization. Best validation accuracy: 1.00000
```

Credit <http://cs231n.stanford.edu/>

98

---

---

---

---

---

---

---

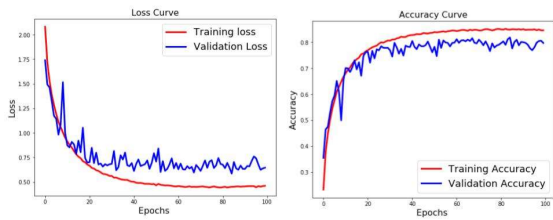
---

---

---

## Sanity checks

- Toujours visualiser les courbes d'apprentissage et de validation



99

---

---

---

---

---

---

---

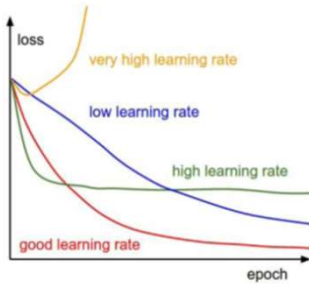
---

---

---

## Sanity checks

3. Toujours visualiser les courbes d'apprentissage et de validation



Credit: <http://cs231n.stanford.edu>

100

---

---

---

---

---

---

---

---

## Sanity checks

3. Toujours vérifier la validité d'un gradient

Comme on l'a vu, calculer un gradient est sujet à erreur. Il faut donc s'assurer que nos gradients sont bons au fur et à mesure qu'on rédige notre code. En voici la meilleure façon

Rappel

Approximation numérique de la dérivée

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

101

---

---

---

---

---

---

---

---

## Sanity checks

3. Toujours vérifier la validité d'un gradient

On peut facilement calculer un gradient à l'aide d'une approximation numérique.

Rappel

Approximation numérique du gradient

$$\nabla E(W) \approx \frac{E(W+H) - E(W)}{H}$$

En calculant

$$\frac{\partial E(W)}{\partial w_i} \approx \frac{E(w_i+h) - E(w_i)}{h} \quad \forall i$$

102

---

---

---

---

---

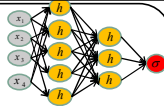
---

---

---

### Vérification du gradient

(exemple)



| W                | W+h                    | gradient W                      |
|------------------|------------------------|---------------------------------|
| $W_{00} = 0.34$  | $W_{00} = 0.34+0.0001$ | $-2.5=(1.25322-1.25347)/0.0001$ |
| $W_{01} = -1.11$ | $W_{01} = -1.11$       |                                 |
| $W_{02} = 0.78$  | $W_{02} = 0.78$        |                                 |
| ...              | ...                    |                                 |
| $W_{20} = -3.1$  | $W_{20} = -3.1$        |                                 |
| $W_{21} = -1.5,$ | $W_{21} = -1.5,$       |                                 |
| $W_{22} = 0.33$  | $W_{22} = 0.33$        |                                 |

$E(W)=1.25347$   $E(W+h)=1.25322$

103

---

---

---

---

---

---

---

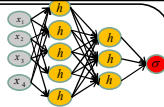
---

---

---

### Vérification du gradient

(exemple)



| W                | W+h                     | gradient W                     |
|------------------|-------------------------|--------------------------------|
| $W_{00} = 0.34$  | $W_{00} = 0.34$         | $-2.5$                         |
| $W_{01} = -1.11$ | $W_{01} = -1.11+0.0001$ | $0.6=(1.25353-1.25347)/0.0001$ |
| $W_{02} = 0.78$  | $W_{02} = 0.78$         |                                |
| ...              | ...                     |                                |
| $W_{20} = -3.1$  | $W_{20} = -3.1$         |                                |
| $W_{21} = -1.5,$ | $W_{21} = -1.5,$        |                                |
| $W_{22} = 0.33$  | $W_{22} = 0.33$         |                                |

$E(W)=1.25347$   $E(W+h)=1.25353$

104

---

---

---

---

---

---

---

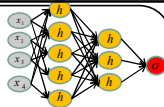
---

---

---

### Vérification du gradient

(exemple)



| W                | W+h                    | gradient W                     |
|------------------|------------------------|--------------------------------|
| $W_{00} = 0.34$  | $W_{00} = 0.34$        | $-2.5$                         |
| $W_{01} = -1.11$ | $W_{01} = -1.11$       | $0.6$                          |
| $W_{02} = 0.78$  | $W_{02} = 0.78+0.0001$ | $0.0=(1.25347-1.25347)/0.0001$ |
| ...              | ...                    |                                |
| $W_{20} = -3.1$  | $W_{20} = -3.1$        |                                |
| $W_{21} = -1.5,$ | $W_{21} = -1.5,$       |                                |
| $W_{22} = 0.33$  | $W_{22} = 0.33$        |                                |

$E(W)=1.25347$   $E(W+h)=1.25347$

105

---

---

---

---

---

---

---

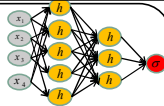
---

---

---

### Vérification du gradient

(exemple)



| W                | W+h              | gradient W |
|------------------|------------------|------------|
| $W_{00} = 0.34$  | $W_{00} = 0.34$  | -2.5       |
| $W_{01} = -1.11$ | $W_{01} = -1.11$ | 0.6        |
| $W_{02} = 0.78$  | $W_{02} = 0.78$  | 0.0        |
| ...              | ...              | ...        |
| $W_{20} = -3.1$  | $W_{20} = -3.1$  | 1.1        |
| $W_{21} = -1.5$  | $W_{21} = -1.5$  | 1.3        |
| $W_{22} = 0.33$  | $W_{22} = 0.33$  | -2.1       |

$E(W) = 1.25347$

106

---

---

---

---

---

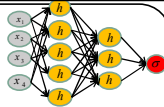
---

---

---

### Vérification du gradient

(exemple)



| gradient W<br>(numérique) |  | gradient W<br>(retro-propagation) |
|---------------------------|--|-----------------------------------|
| -2.5                      |  | -2.5                              |
| 0.6                       |  | 0.6                               |
| 0.0                       |  | 0.0                               |
| ...                       |  | ...                               |
| 1.1                       |  | 1.1                               |
| 1.3                       |  | 1.3                               |
| -2.1                      |  | -2.1                              |

107

---

---

---

---

---

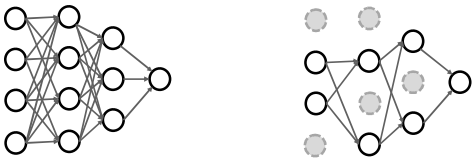
---

---

---

### Autre bonne pratique : Dropout

Forcer à zéro certains neurones de façon aléatoire à chaque itération



108

---

---

---

---

---

---

---

---

## Autre bonne pratique : *Dropout*

Idée : s'assurer que **chaque neurone apprend pas lui-même** en brisant au hasard des chemins.

Crédit <http://cs231n.stanford.edu/>

109

---

---

---

---

---

---

---

---

## Autre bonne pratique : *Dropout*

```
p = 0.5 # probability of keeping a unit active. higher = less dropout

def train_step(X):
    """ X contains the data """

    # forward pass for example 3-layer neural network
    H1 = np.maximum(0, np.dot(W1, X) + b1)
    U1 = np.random.rand(*H1.shape) < p # first dropout mask
    H1 *= U1 # drop!
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
    U2 = np.random.rand(*H2.shape) < p # second dropout mask
    H2 *= U2 # drop!
    out = np.dot(W3, H2) + b3

    # backward pass: compute gradients... (not shown)
    # perform parameter update... (not shown)
```

Crédit <http://cs231n.stanford.edu/>

110

---

---

---

---

---

---

---

---

## Autre bonne pratique : *Dropout*

Le problème avec *Dropout* est en **prédiction** (« test time »)

car *dropout* **ajoute du bruit** à la prédiction

$$pred = y_W(\vec{x}, Z)$$

↑  
masque aléatoire

111

---

---

---

---

---

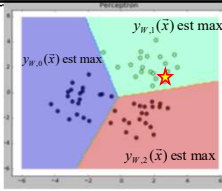
---

---

---

dropout **ajoute du bruit** à la prédiction.

Exemple simple :  $\vec{x} = \begin{pmatrix} 2.2 \\ 1.3 \end{pmatrix}, t=1$



**Si on lance le modèle 10 fois, on aura 10 réponses différentes**

|              |            |             |
|--------------|------------|-------------|
| [ 0.09378555 | 0.76511644 | 0.141098 ]  |
| [ 0.13982909 | 0.62885327 | 0.23131764] |
| [ 0.23658253 | 0.61960162 | 0.14381585] |
| [ 0.23779425 | 0.51357115 | 0.24863461] |
| [ 0.16005442 | 0.68060227 | 0.1593433 ] |
| [ 0.16303195 | 0.50583392 | 0.33113413] |
| [ 0.24183069 | 0.51319834 | 0.24497097] |
| [ 0.14521815 | 0.52006858 | 0.33471327] |
| [ 0.09952161 | 0.66276146 | 0.23771692] |
| [ 0.16172851 | 0.6044877  | 0.23378379] |

112

---

---

---

---

---

---

---

---

---

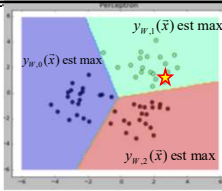
---

---

---

dropout **ajoute du bruit** à la prédiction.

Exemple simple :  $\vec{x} = \begin{pmatrix} 2.2 \\ 1.3 \end{pmatrix}, t=1$



**Solution**, exécuter le modèle un grand nombre de fois et **prendre la moyenne**.

|               |             |             |
|---------------|-------------|-------------|
| [ 0.09378555  | 0.76511644  | 0.141098 ]  |
| [ 0.13982909  | 0.62885327  | 0.23131764] |
| [ 0.23658253  | 0.61960162  | 0.14381585] |
| [ 0.23779425  | 0.51357115  | 0.24863461] |
| [ 0.16005442  | 0.68060227  | 0.1593433 ] |
| [ 0.16303195  | 0.50583392  | 0.33113413] |
| [ 0.24183069  | 0.51319834  | 0.24497097] |
| [ 0.14521815  | 0.52006858  | 0.33471327] |
| [ 0.09952161  | 0.66276146  | 0.23771692] |
| [ 0.16172851  | 0.6044877   | 0.23378379] |
| (...)         |             |             |
| [ 0.15933813, | 0.65957005, | 0.18109183] |

113

---

---

---

---

---

---

---

---

---

---

---

---

Exécuter le modèle un grand nombre de fois et **prendre la moyenne** revient à calculer **l'espérance mathématique**

$$pred = E_z [y_w(\vec{x}, \vec{z})] = \sum_i P(\vec{z}) y_w(\vec{x}, \vec{z})$$

Bonne nouvelle, on peut faire plus simple en approximant cette l'espérance mathématique!

114

---

---

---

---

---

---

---

---

---

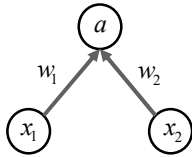
---

---

---

## Regardons pour un neurone

Avec une probabilité de *dropout* de 50%, en prédiction  $w_1$  et  $w_2$  seront **nuls 1 fois sur 2**



$$\begin{aligned} E[a] &= \frac{1}{4}(w_1x_1 + w_2x_2) + \frac{1}{4}(w_1x_1 + 0x_2) \\ &\quad + \frac{1}{4}(0x_1 + w_2x_2) + \frac{1}{4}(0x_1 + 0x_2) \\ &= \frac{1}{2}(w_1x_1 + w_2x_2) \end{aligned}$$

En prédiction, on a qu'à multiplier par la prob. de *dropout*.

115

---

---

---

---

---

---

---

---

```
*** Vanilla Dropout: Not recommended implementation (see notes below) ***
p = 0.5 # probability of keeping a unit active. higher = less dropout

def train_step(X):
    """ X contains the data """
    # forward pass for example 3-layer neural network
    H1 = np.maximum(0, np.dot(W1, X) + b1)
    U1 = np.random.rand(*H1.shape) < p # first dropout mask
    H1 *= U1 # drop!
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
    U2 = np.random.rand(*H2.shape) < p # second dropout mask
    H2 *= U2 # drop!
    out = np.dot(W3, H2) + b3

    # backward pass: compute gradients... (not shown)
    # perform parameter update... (not shown)

def predict(X):
    # ensemble forward pass
    H1 = np.maximum(0, np.dot(W1, X) + b1) * p # NOTE: scale the activations
    H2 = np.maximum(0, np.dot(W2, H1) + b2) * p # NOTE: scale the activations
    out = np.dot(W3, H2) + b3
```

En prédiction, tous les neurones sont actifs  
→ tout ce qu'il faut faire est de multiplier la sortie de chaque couche par la probabilité de dropout

Crédit <http://cs231n.stanford.edu/>

116

---

---

---

---

---

---

---

---

## Descente de gradient version améliorée

117

---

---

---

---

---

---

---

---

## Descente de gradient

$$W^{[t+1]} = W^{[t]} - \eta \nabla E_D(W^{[t]})$$

118

---

---

---

---

---

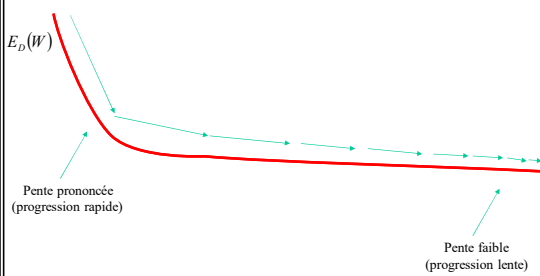
---

---

---

## Descente de gradient : **problème**

Progrès quasi nul lorsque la pente est très faible



119

---

---

---

---

---

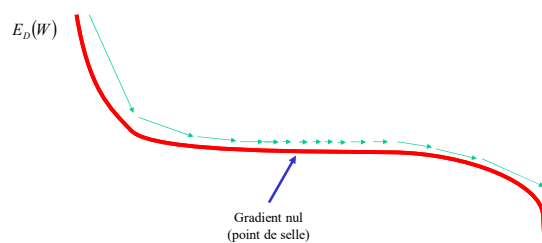
---

---

---

## Descente de gradient : **problème**

Les points de selles sont fréquents en haute dimension



120

---

---

---

---

---

---

---

---



## Descente de gradient : **problème**

Qu'arrive-t-il si la fonction de coût (loss) a une pente prononcée dans une direction et moins prononcée dans une autre direction?

121

---

---

---

---

---

---

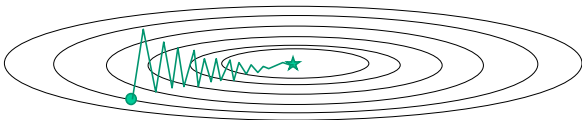
---

---

## Descente de gradient : **problème**

Qu'arrive-t-il si la fonction de coût (loss) a une pente prononcée dans une direction et moins prononcée dans une autre direction?

**Progrès très lent le long de la pente la plus faible et oscillation le long de l'autre direction.**



122

---

---

---

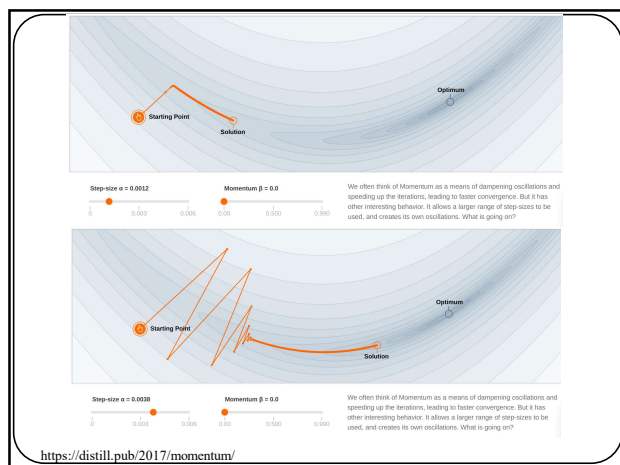
---

---

---

---

---



123

---

---

---

---

---

---

---

---

# Descente de gradient + Momentum

Descente de gradient  
stochastique

$$w_{t+1} = w_t - \eta \nabla E_{\bar{x}_n}(w_t)$$

Descente de gradient  
stochastique + **Momentum**

$$v_{t+1} = \rho v_t + \nabla E_{\bar{x}_n}(w_t)$$

$$w_{t+1} = w_t - \eta v_{t+1}$$

Provient de l'équation de la vitesse  
(à démontrer en devoir ou en exercice)

$\rho$  exprime la « friction », en général  $\in [0.5, 1[$

124

---

---

---

---

---

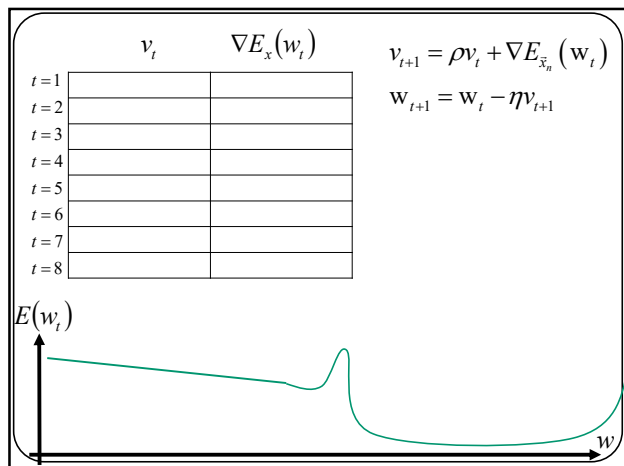
---

---

---

---

---



125

---

---

---

---

---

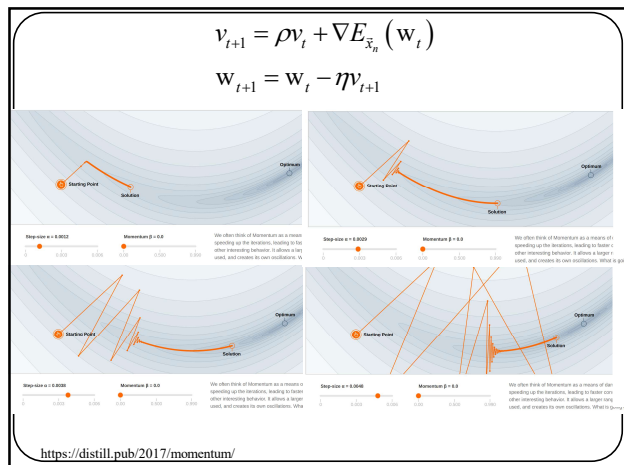
---

---

---

---

---



126

---

---

---

---

---

---

---

---

---

---

$$v_{t+1} = \rho v_t + \nabla E_{\bar{x}_n} (w_t)$$

$$w_{t+1} = w_t - \eta v_{t+1}$$

127

---

---

---

---

---

---

---

---

---

---

### AdaGrad (décroissance automatique de $\eta$ )

Descente de gradient stochastique AdaGrad

$$w_{t+1} = w_t - \eta \nabla E_{\bar{x}_n} (w_t)$$

$$dE_t = \nabla E_{\bar{x}_n} (w_t)$$

$$m_{t+1} = m_t + |dE_t|$$

$$w_{t+1} = w_t - \frac{\eta}{m_{t+1} + \varepsilon} dE_t$$

128

---

---

---

---

---

---

---

---

---

---

### AdaGrad (décroissance automatique de $\eta$ )

Descente de gradient stochastique AdaGrad

$$w_{t+1} = w_t - \eta \nabla E_{\bar{x}_n} (w_t)$$

$$dE_t = \nabla E_{\bar{x}_n} (w_t)$$

$$m_{t+1} = m_t + |dE_t|$$

$$w_{t+1} = w_t - \frac{\eta}{\underbrace{m_{t+1} + \varepsilon}} dE_t$$

$\eta$  décroît sans cesse au fur et à mesure de l'optimisation

129

---

---

---

---

---

---

---

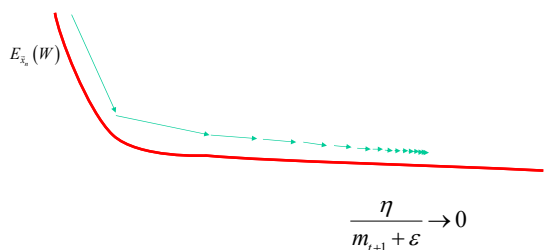
---

---

---

## AdaGrad (décroissance automatique de $\eta$ )

Qu'arrive-t-il à long terme?



130

---

---

---

---

---

---

---

---

## RMSProp (AdaGrad amélioré)

**AdaGrad**

$$dE_t = \nabla E_{\tilde{x}_t}(w_t)$$

$$m_{t+1} = m_t + |dE_t|$$

$$w_{t+1} = w_t - \frac{\eta}{m_{t+1} + \epsilon} dE_t$$

**RMSProp**

$$dE_t = \nabla E_{\tilde{x}_t}(w_t)$$

$$m_{t+1} = \gamma m_t + (1 - \gamma) |dE_t|$$

$$w_{t+1} = w_t - \frac{\eta}{m_{t+1} + \epsilon} dE_t$$

$\eta$  décroît lorsque le gradient est élevé  
 $\eta$  augmente lorsque le gradient est faible

131

---

---

---

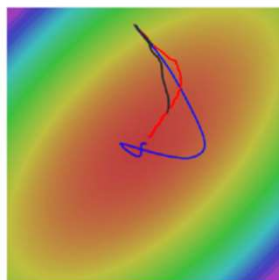
---

---

---

---

---



- SGD
- SGD+Momentum
- RMSProp

132

---

---

---

---

---

---

---

---

### Adam (Combo entre Momentum et RMSProp)

#### Momentum

$$v_{t+1} = \rho v_t + \nabla E_{\bar{x}_n}(\mathbf{w}_t)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta v_{t+1}$$

#### Adam

$$dE_t = \nabla E_{\bar{x}_n}(\mathbf{w}_t)$$

$$v_{t+1} = \alpha v_t + (1 - \alpha) dE_t$$

$$m_{t+1} = \gamma m_t + (1 - \gamma) |dE_t|$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\eta}{m_{t+1} + \epsilon} v_{t+1}$$

133

---

---

---

---

---

---

---

---

### Adam (Combo entre Momentum et RMSProp)

#### Momentum

$$v_{t+1} = \rho v_t + \nabla E_{\bar{x}_n}(\mathbf{w}_t)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta v_{t+1}$$

#### Adam

$$dE_t = \nabla E_{\bar{x}_n}(\mathbf{w}_t)$$

$$v_{t+1} = \alpha v_t + (1 - \alpha) dE_t$$

$$m_{t+1} = \gamma m_t + (1 - \gamma) |dE_t|$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\eta}{m_{t+1} + \epsilon} v_{t+1}$$

*Momentum*

134

---

---

---

---

---

---

---

---

### Adam (Combo entre Momentum et RMSProp)

#### RMSProp

$$dE_t = \nabla E_{\bar{x}_n}(\mathbf{w}_t)$$

$$m_{t+1} = \gamma m_t + (1 - \gamma) |dE_t|$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\eta}{m_{t+1} + \epsilon} dE_t$$

#### Adam

$$dE_t = \nabla E_{\bar{x}_n}(\mathbf{w}_t)$$

$$v_{t+1} = \alpha v_t + (1 - \alpha) dE_t$$

$$m_{t+1} = \gamma m_t + (1 - \gamma) |dE_t|$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\eta}{m_{t+1} + \epsilon} v_{t+1}$$

*RMSProp*

135

---

---

---

---

---

---

---

---

## Adam (Version complète)

$$v_{t=0} = 0$$

$$m_{t=0} = 0$$

for  $t=1$  à num\_iterations

for  $n=0$  à  $N$

$$dE_t = \nabla_{\vec{x}_n} (w_t)$$

$$v_{t+1} = \alpha v_t + (1 - \alpha) dE_t$$

$$m_{t+1} = \gamma m_t + (1 - \gamma) |dE_t|$$

$$v_{t+1} = \frac{v_{t+1}}{1 - \beta_1^t}, m_{t+1} = \frac{m_{t+1}}{1 - \beta_2^t}$$

$$\beta_1 = 0.9, \beta_2 = 0.99$$

$$w_{t+1} = w_t - \frac{\eta}{m_{t+1} + \epsilon} v_{t+1}$$

136

---

---

---

---

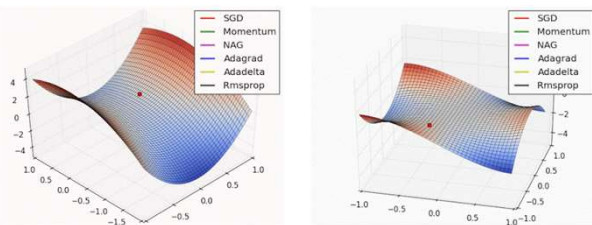
---

---

---

---

## Illustrations



À voir sur :

[www.denizyuret.com/2015/03/alec-radfords-animations-for.html](http://www.denizyuret.com/2015/03/alec-radfords-animations-for.html)

137

---

---

---

---

---

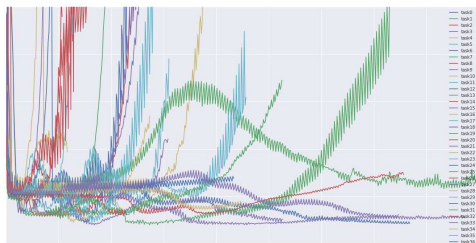
---

---

---

## Autre excellent survol

<http://ruder.io/optimizing-gradient-descent/>



138

---

---

---

---

---

---

---

---