

Réseaux de neurones IFT 603-712

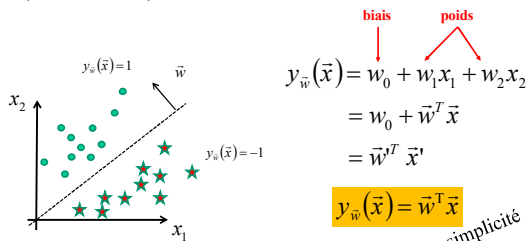
Réseaux de neurones multicouches
Par
Pierre-Marc Jodoin

Rappel réseaux de neurones

(Perceptron, régression logistique, SVM)

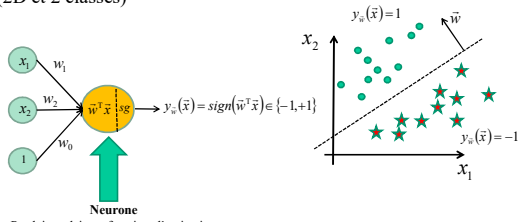
Séparation linéaire

(2D et 2 classes)



- 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)

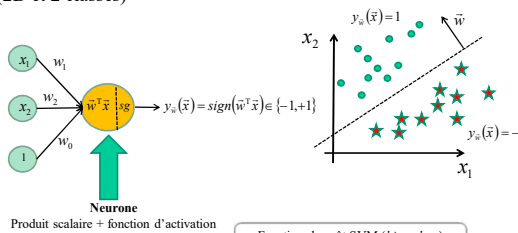
Perceptron (2D et 2 classes)



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

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

Hinge Loss (2D et 2 classes)

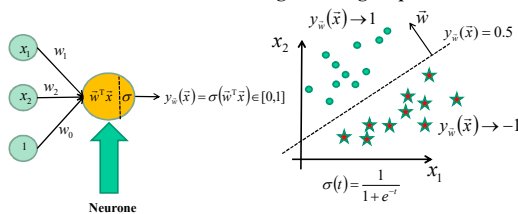


$$E_D(\vec{w}) = \sum_{n=1}^N \max(0, 1 - t_n \vec{w}^T \vec{x}_n)$$

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

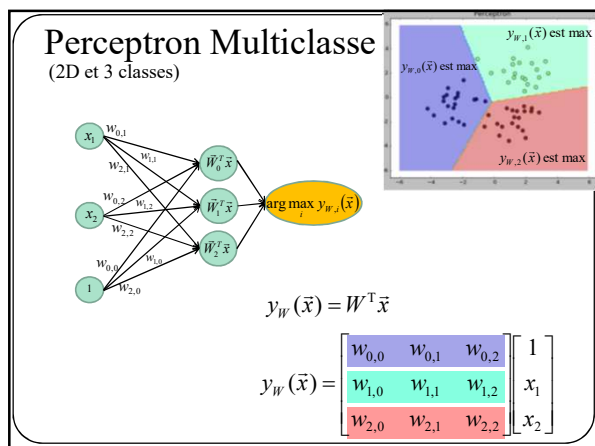
Régression logistique (2D, 2 classes)

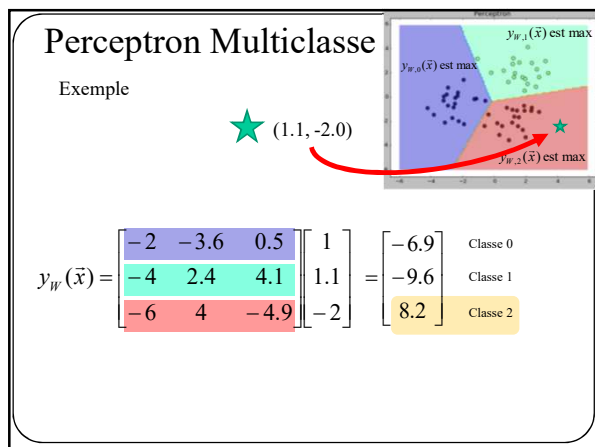
Nouvelle fonction d'activation : **sigmoïde logistique**

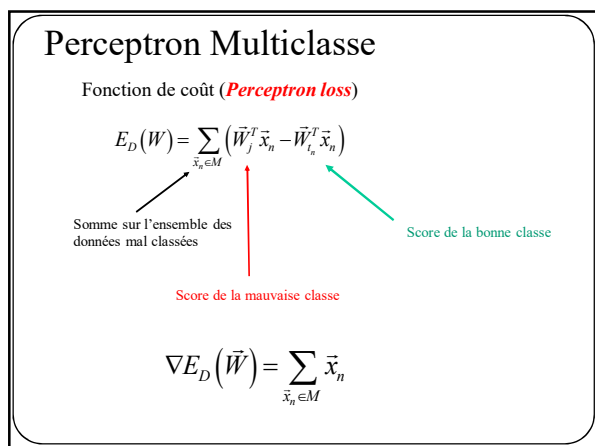


$$E_D(\vec{w}) = -\sum_{n=1}^N [t_n \ln(y_w(\vec{x}_n)) + (1 - t_n) \ln(1 - y_w(\vec{x}_n))]$$

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







Hinge Multiclasse

Fonction de coût (**Hinge loss**)

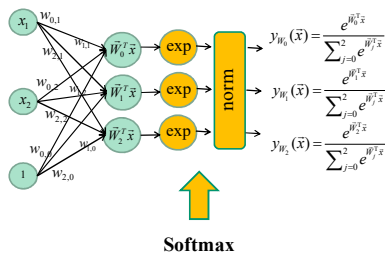
$$E_D(W) = \sum_{n=1}^N \sum_j \max(0, 1 + \vec{W}_j^T \vec{x}_n - \vec{W}_{c_n}^T \vec{x}_n)$$

Score de la mauvaise classe

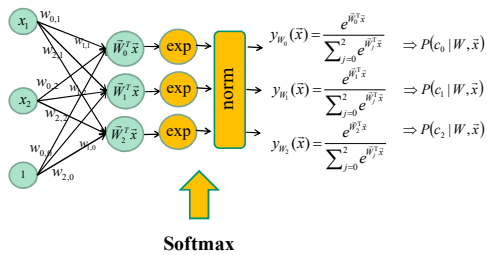
Score de la bonne classe

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

Régression logistique multiclass



Régression logistique multiclass



Régression logistique multiclasse

airplane		'airplane' $\Rightarrow t = [100000000]$
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

Régression logistique multiclasse

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

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

Devoir 4

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

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

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

\vec{x}

Hinge loss

Score

$$\begin{aligned} & \max(0, -2.85 - 0.28 + 1) + \\ & \max(0, 0.86 - 0.28 + 1) \\ & = \\ & \max(0, -2.13) + \max(0, 1.58) \\ & = \\ & \mathbf{1.58} \end{aligned}$$

Entropie croisée

Score

$$\begin{aligned} & \begin{bmatrix} -2.85 \\ 0.86 \\ 0.28 \end{bmatrix} \xrightarrow{\text{exp}} \begin{bmatrix} 0.06 \\ 2.36 \\ 1.32 \end{bmatrix} \xrightarrow{\text{norm}} \begin{bmatrix} 0.02 \\ 0.63 \\ 0.35 \end{bmatrix} \\ & \text{(Softmax)} \end{aligned}$$

$-\ln(0.35) = 0.452$

Maximum *a posteriori*

Régularisation

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

Fonction de perte Régularisation

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

Optimisation

Descente de gradient

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

↗ 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
 $w = w - \eta^{(k)} \nabla E(\tilde{x}_n)$
 JUSQU'À ce que toutes les données
 soient bien classées ou $k=MAX_ITER$

Devoir 4

Optimisation par Batch

Initialiser w
 $k=0$
 FAIRE $k=k+1$
 $w = w - \eta^{(k)} \sum_i \nabla E(\tilde{x}_i)$
 JUSQU'À ce que toutes les données
 sont bien classées ou $k=MAX_ITER$

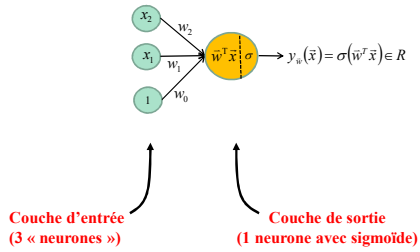
Parfois $\eta^{[k]} = cst / k$

Maintenant, rendons le réseau

profond

Maintenant, rendons le réseau

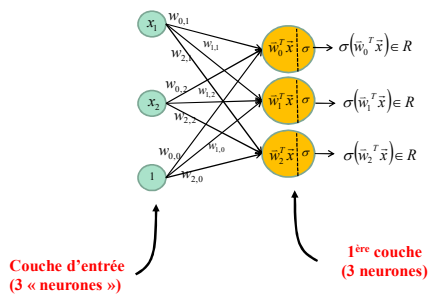
2D, 2Classes, Régression logistique linéaire



19

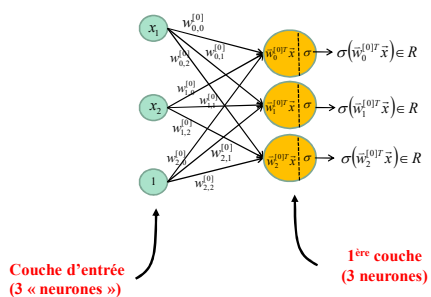
2D, 2Classes, Réseau à 1 couche cachée

Maintenant, ajoutons arbitrairement 3 neurones

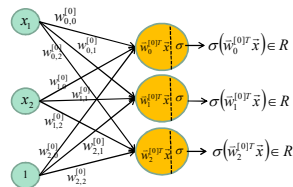


2D, 2Classes, Réseau à 1 couche cachée

Puisque les poids sont **entre la couche d'entrée et la première couche** on va les identifier à l'aide de l'indice **[0]**



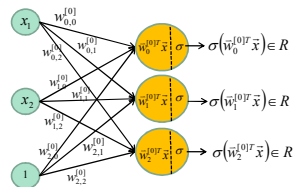
2D, 2Classes, Réseau à 1 couche cachée



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

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

2D, 2Classes, Réseau à 1 couche cachée

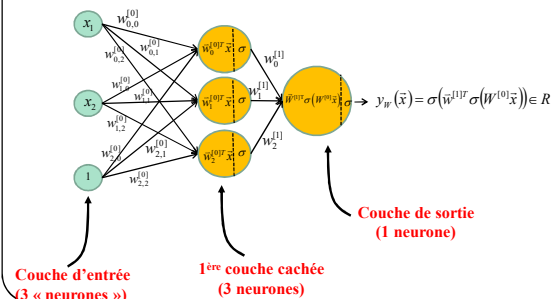


NOTE: représentation plus simple de la sortie de la 1^{ère} couche (3 réels)

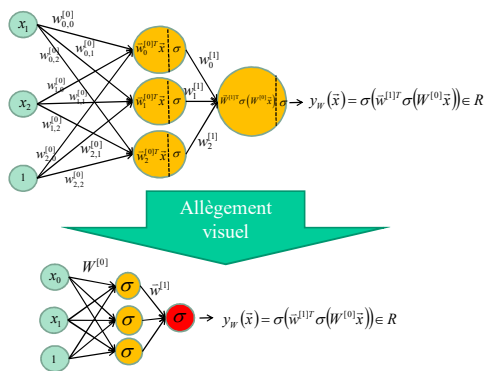
$$\sigma(W^{[0]}\vec{x})$$

2D, 2Classes, Réseau à 1 couche cachée

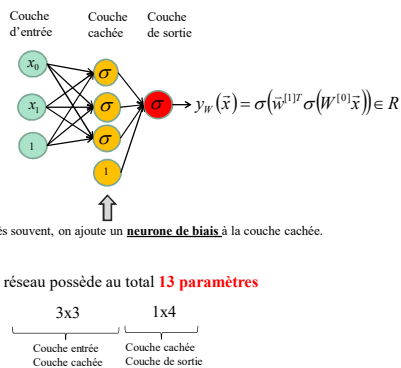
Si on veut effectuer une **classification 2 classes** via une **régression logistique** (donc **une fonction coût par « entropie croisée »**) on doit ajouter **un neurone de sortie**.



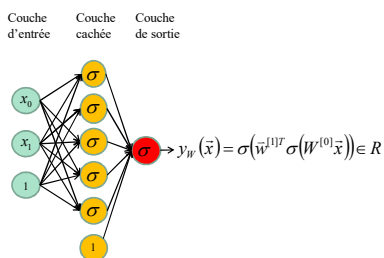
2D, 2Classes, Réseau à 1 couche cachée



2D, 2Classes, Réseau à 1 couche cachée



2D, 2Classes, Réseau à 1 couche cachée



Plus on augmente le nombre de neurones à la couche cachée, plus on **augmente la capacité du système**.

Ce réseau a $5 \times 3 + 1 \times 6 = 21$ paramètres

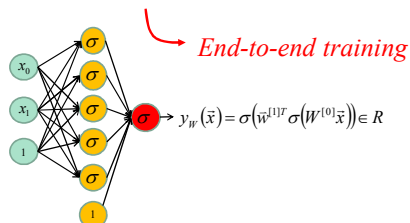
27

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 $\sigma(W^{[0]}\vec{x}) \in \mathbb{R}^4$) 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

Nombre de neurones VS Capacity

Aucun neurone caché



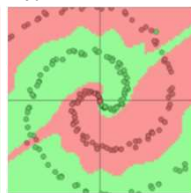
Classification linéaire
Underfitting
(pas assez de capacité)

12 neurones cachés



Classification non linéaire
BON RÉSULTAT
(bonne capacité)

60 neurones cachés

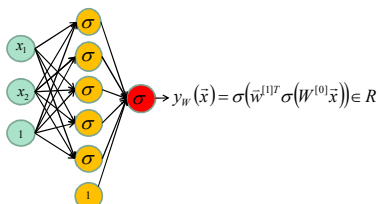


Classification non linéaire
Over fitting
(trop grande capacité)

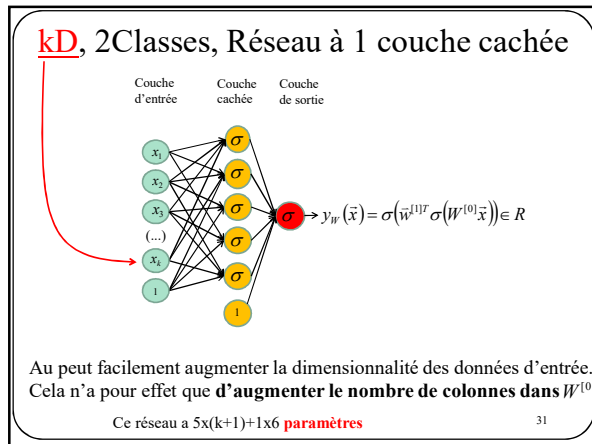
<http://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html>

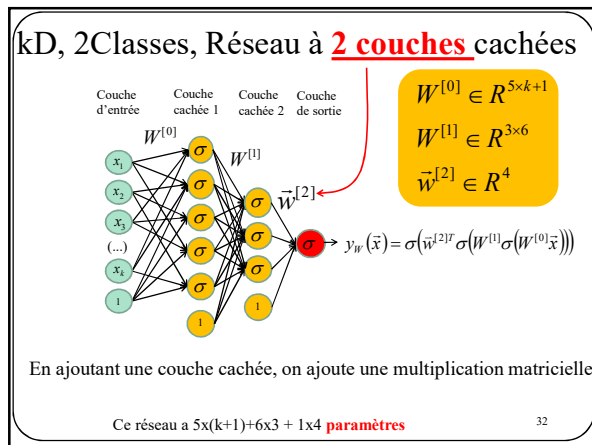
2D, 2Classes, Réseau à 1 couche cachée

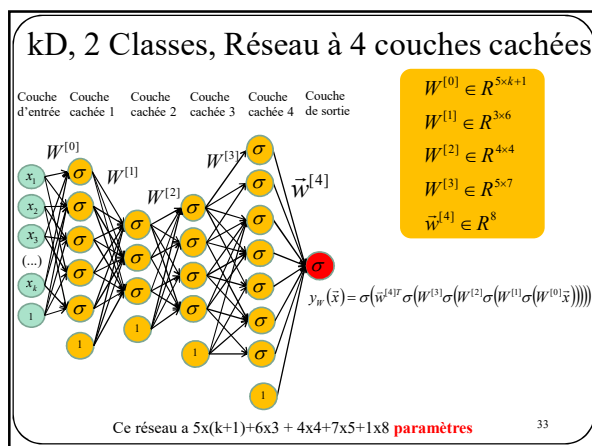
Couche d'entrée Couche cachée Couche de sortie



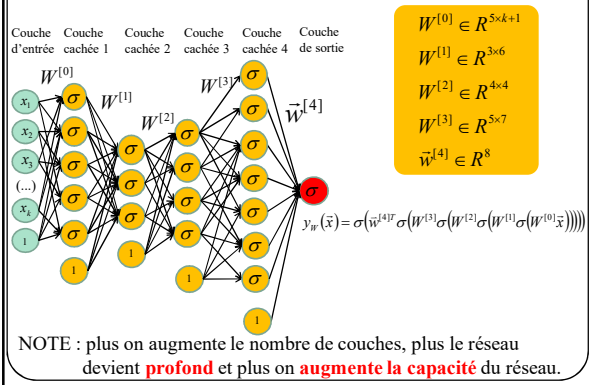
30



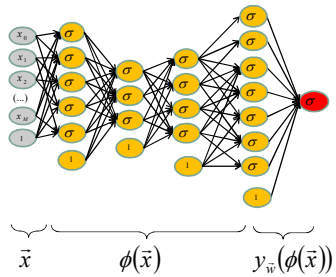




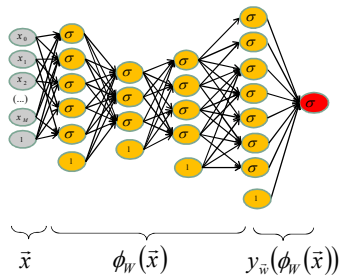
kD, 2 Classes, Réseau à 4 couches cachées

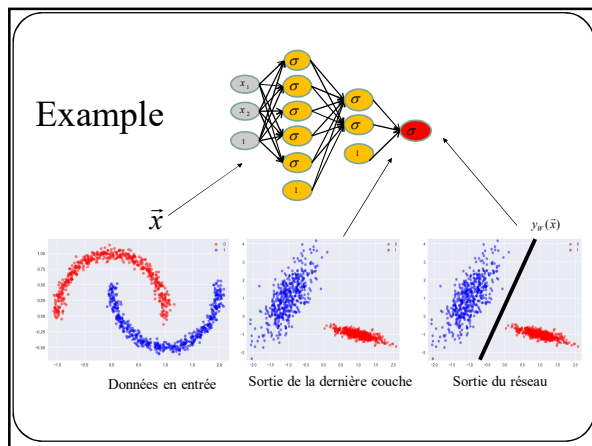


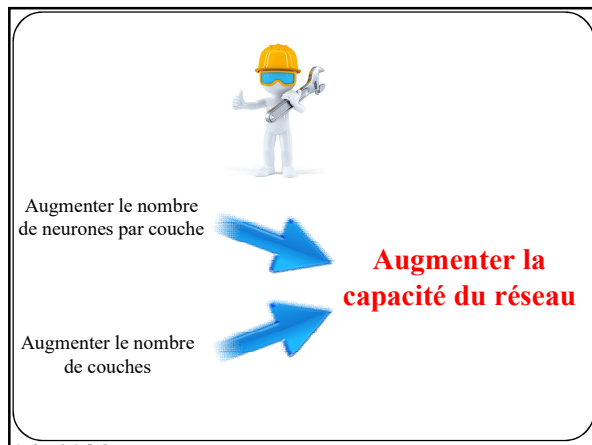
Réseau de neurones multicouches = apprendre la fonction de base



Réseau de neurones multicouches = apprendre la fonction de base









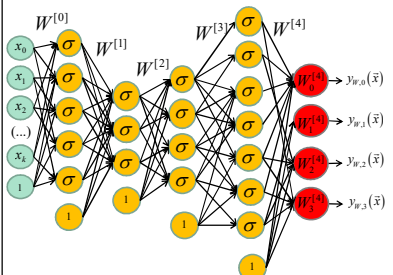


Lorsqu'un réseau doit prédire **plus de 2 classes**, on lui assigne **K neurones de sortie**, une par classe.

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
$x_0, x_1, x_2, \dots, x_i, 1$	$W^{[1]}$	$W^{[2]}$	$W^{[3]}$	$W^{[4]}$	$W^{[5]}$


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

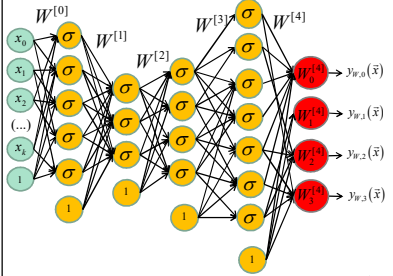


$y_w(\vec{x}) = W^{[4]} \sigma \left(W^{[3]} \sigma \left(W^{[2]} \sigma \left(W^{[1]} \sigma \left(W^{[0]} \vec{x} \right) \right) \right) \right)$

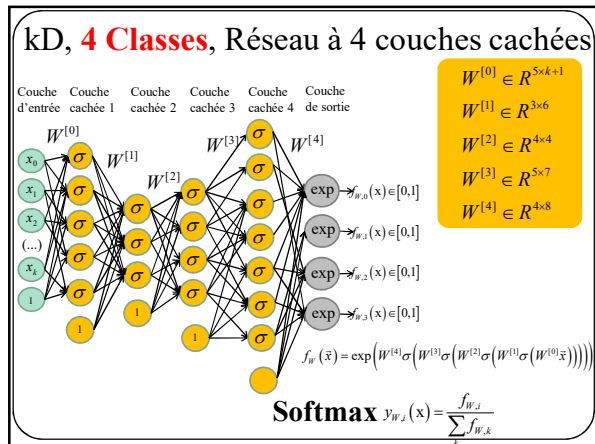
kD, 4 Classes, Réseau à 4 couches cachées

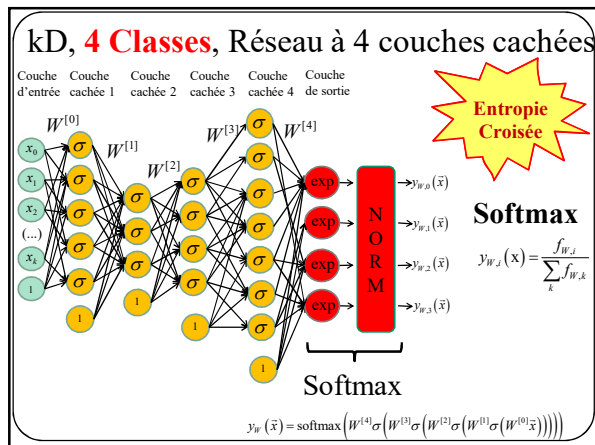
Couche d'entrée	Couche cachée 1	Couche cachée 2	Couche cachée 3	Couche cachée 4	Couche de sortie
$x_0, x_1, x_2, \dots, x_i, 1$	$W^{[1]}$	$W^{[2]}$	$W^{[3]}$	$W^{[4]}$	$W^{[5]}$





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





Simulation

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

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

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

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

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

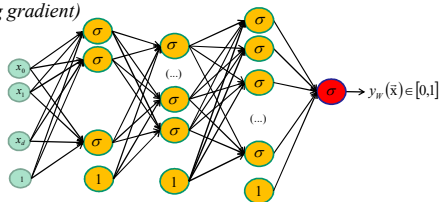
Comment optimiser les paramètres?

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

50

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

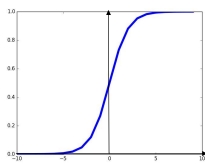
disparition du gradient

Démonstration au tableau

51

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

Fonction d'activation



Sigmoïde

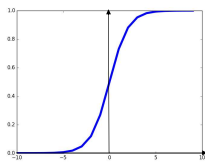
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

- Ramène les valeurs entre 0 et 1
- Historiquement populaire

3 Problèmes :

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

Fonction d'activation



Sigmoïde

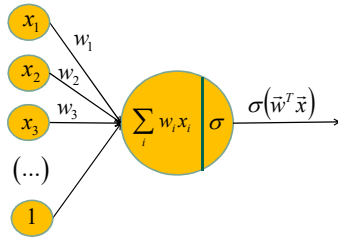
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

- Ramène les valeurs entre 0 et 1
- Historiquement populaire

3 Problèmes :

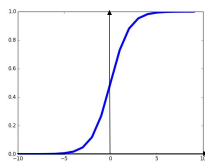
- Un neurone saturé a pour effet de « **tuer** » **les gradients**
- Sortie d'une sigmoïde n'est **pas** **centrée à zéro**.

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?**

Fonction d'activation



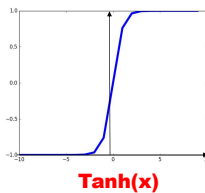
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

- 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.

Fonction d'activation

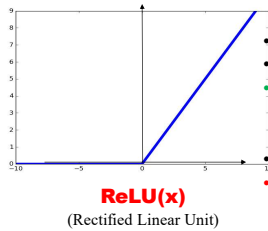


- Ramène les valeurs entre -1 et 1
- **Sortie centrée à zéro** 😊
- **Disparition du gradient** lorsque la fonction sature ☹️

[LeCun et al., 1991]

Fonction d'activation

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

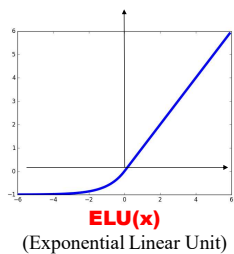


- Aucune **saturation** 😊
- Super **rapide** 😊
- Converge plus rapide** que sigmoïde/tanh (5 à 10x) 😊
- Sortie **non centrée à zéro** ☹️
- Un inconvénient** : qu'arrive-t-il au gradient lorsque $x < 0$? ☹️

[Krizhevsky et al., 2012]

Fonction d'activation

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

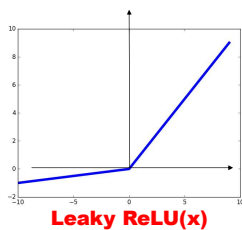


- Tous les avantages de **ReLU** 😊
- Sortie plus « **centrée à zéro** » 😊
- Converge plus rapide** que sigmoïde/tanh (5 à 10x) 😊
- Gradients meurent plus lentement** 😊
- $\exp()$ est **coûteux** ☹️

[Clevert et al., 2015]

Fonction d'activation

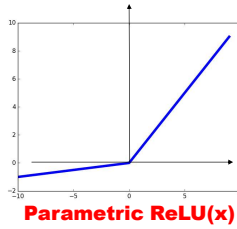
$$\text{LReLU}(x) = \max(0.01x, x)$$



- 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]

Fonction d'activation



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

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

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.

Les bonnes pratiques

Optimisation

Descente de gradient

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

∇E → Gradient de la fonction de coût
 η → 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 $\eta^{[k]} = cst / k$

Optimisation

Descente de gradient

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

∇E → Gradient de la fonction de coût
 η → 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)
  JUSQU'À ce que toutes les données sont bien classées ou
  k==MAX_ITER
    
```

} **Itération**

Optimisation

Descente de gradient

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

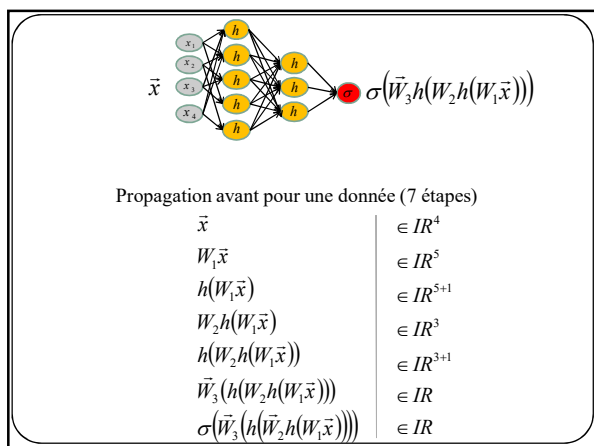
∇E → Gradient de la fonction de coût
 η → Taux d'apprentissage ou "learning rate".

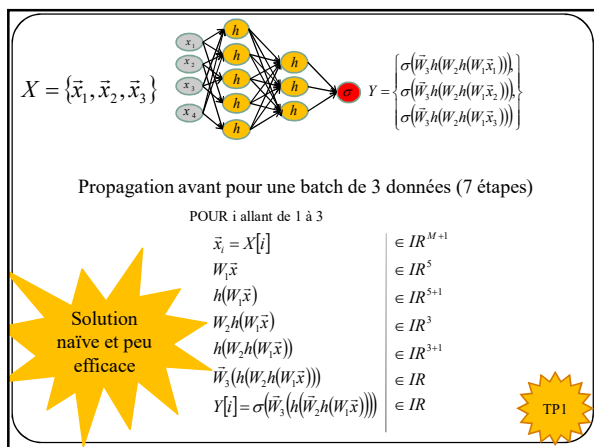
Optimisation par mini-batch

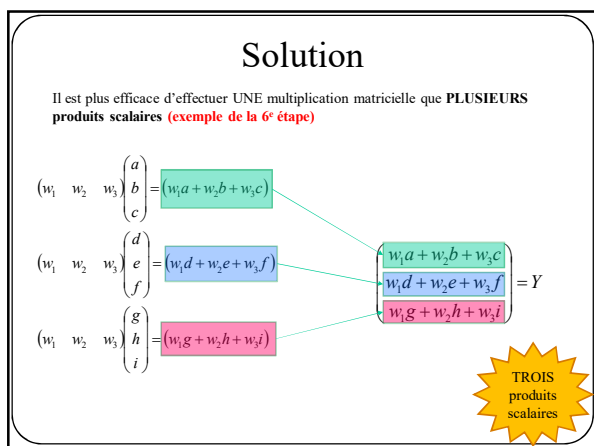
```

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)
  JUSQU'À ce que toutes les données sont bien classées ou
  k==MAX_ITER
    
```

} **Epoch**







Solution

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

$$\begin{bmatrix} w_1 & w_2 & w_3 \end{bmatrix} \begin{bmatrix} a & d & g \\ b & e & h \\ c & f & i \end{bmatrix} = \begin{bmatrix} w_1a + w_2b + w_3c \\ w_1d + w_2e + w_3f \\ w_1g + w_2h + w_3i \end{bmatrix} = Y$$

UNE
multiplication
matricielle

Solution

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

$$\begin{aligned} W_1 \vec{x}_1 &= \begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} \\ w_{21} & w_{22} & w_{23} & w_{24} \\ w_{31} & w_{32} & w_{33} & w_{34} \\ w_{41} & w_{42} & w_{43} & w_{44} \\ w_{51} & w_{52} & w_{53} & w_{54} \end{bmatrix} \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix} = \begin{bmatrix} u_1 \\ u_2 \\ u_3 \\ u_4 \\ u_5 \end{bmatrix} \\ W_1 \vec{x}_2 &= \begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} \\ w_{21} & w_{22} & w_{23} & w_{24} \\ w_{31} & w_{32} & w_{33} & w_{34} \\ w_{41} & w_{42} & w_{43} & w_{44} \\ w_{51} & w_{52} & w_{53} & w_{54} \end{bmatrix} \begin{bmatrix} e \\ f \\ g \\ h \end{bmatrix} = \begin{bmatrix} v_1 \\ v_2 \\ v_3 \\ v_4 \\ v_5 \end{bmatrix} \\ W_1 \vec{x}_3 &= \begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} \\ w_{21} & w_{22} & w_{23} & w_{24} \\ w_{31} & w_{32} & w_{33} & w_{34} \\ w_{41} & w_{42} & w_{43} & w_{44} \\ w_{51} & w_{52} & w_{53} & w_{54} \end{bmatrix} \begin{bmatrix} i \\ j \\ k \\ l \end{bmatrix} = \begin{bmatrix} z_1 \\ z_2 \\ z_3 \\ z_4 \\ z_5 \end{bmatrix} \end{aligned}$$

TROIS
multi
matrice-
vecteur

Solution

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

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

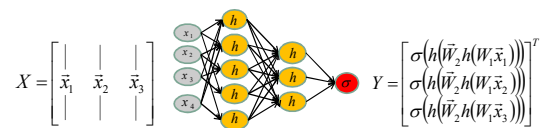
UNE
multiplication
matricielle

Vectorisation de la propagation avant

En résumé, lorsqu'on propage une « batch »

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

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



Vectoriser la rétropropagation

Vectoriser la rétropropagation

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

$$\begin{matrix} \begin{bmatrix} w_1 & w_2 & w_3 \end{bmatrix} & \begin{pmatrix} a & d & g \\ b & e & h \\ c & f & i \end{pmatrix} & = & \begin{pmatrix} w_1a + w_2b + w_3c \\ w_1d + w_2e + w_3f \\ w_1g + w_2h + w_3i \end{pmatrix}^T \\ \vec{W} & X & & Y \end{matrix}$$

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?

Vectoriser la rétropropagation

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

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

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

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

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

Concentrons-nous sur w_1

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

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

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

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

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

Et pour tous les poids

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

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

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

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

W
 X
 Y

Même chose pour 1 couche et une batch de 3

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

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

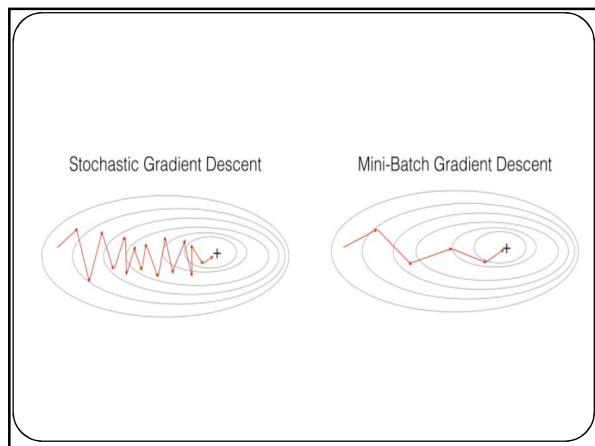
Vectorisation de la rétro-propagation

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

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

Pour plus de détails:

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



Comment initialiser un réseau de neurones?

$W = ?$


Initialisation

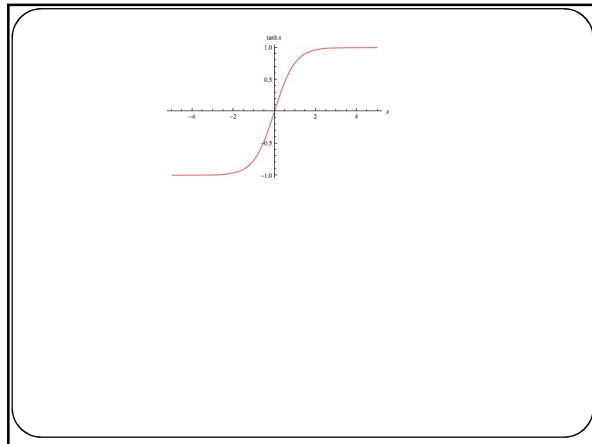
Première idée: faibles valeurs aléatoires
(Gaussienne $\mu = 0, \sigma = 0.01$)

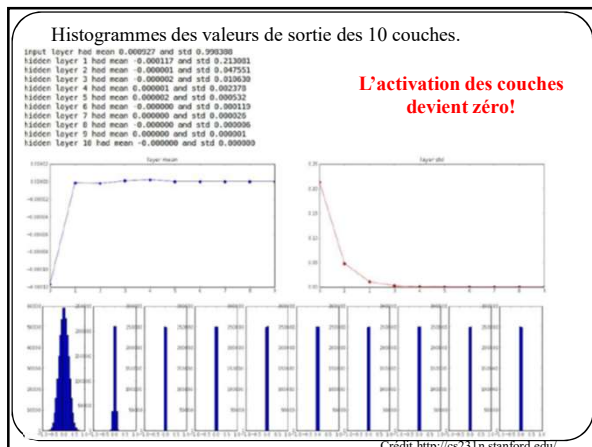
```
W_i = 0.01 * np.random.randn(H_i, H_im1)
```

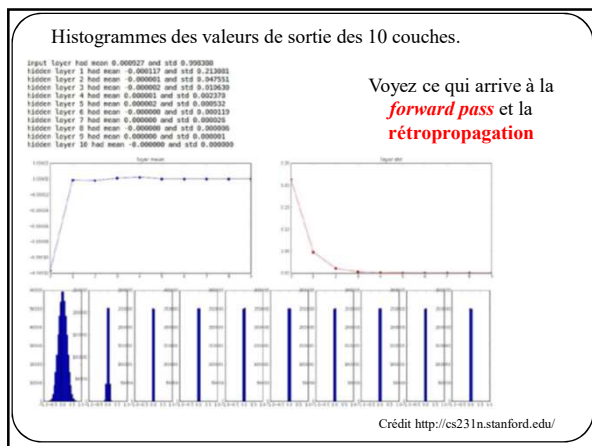
Fonctionne bien pour de petits réseaux mais
pas pour des réseaux profonds.

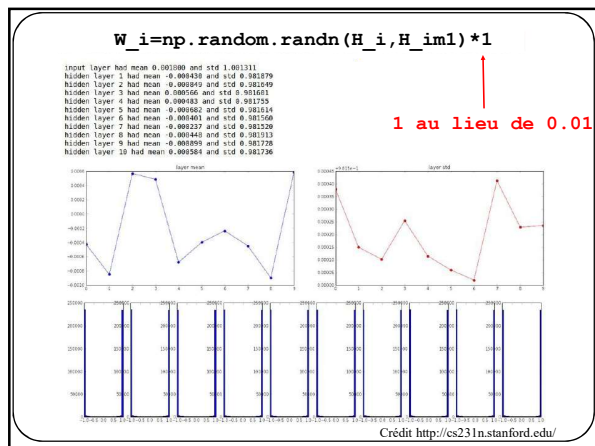
E.g. réseau à 10 couches avec 500 neurones par couche et des **tanh** comme fonctions d'activation.

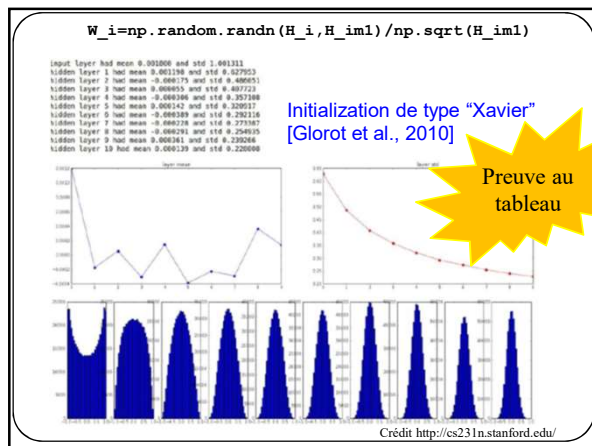






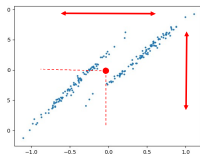
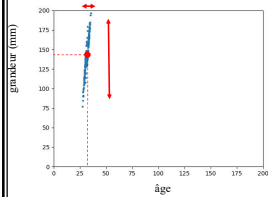






Prétraitement des données

Centrer et normaliser les données d'entrée



```
X=np.mean(X,axis=0)  
X/=np.std(X,axis=0)
```

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**.

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

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

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

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(32*32*3, 50, 10) # input size, hidden size, number of classes  
loss, grad = two_layer_net(x_train, model, y_train, 0.0) # disable regularization  
print loss
```

2.30261216167

loss ~2.3.

"correct" for

10 classes

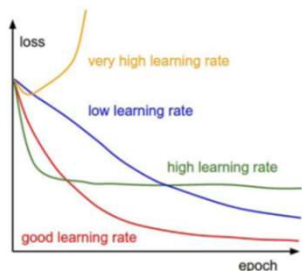
returns the loss and the

gradient for all parameters

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

Sanity checks

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



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 toujours 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

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

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

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

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

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

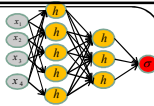
Vérification du gradient

(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

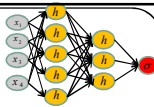
Vérification du gradient (exemple)



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

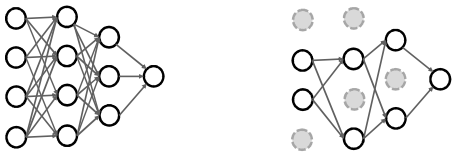
Vérification du gradient (exemple)



gradient W (numérique)		gradient W (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

Autre bonne pratique : *Dropout*

Forcer à zéro certains neurones de façon aléatoire à chaque itération



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/>

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/>

Autre bonne pratique : *Dropout*

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

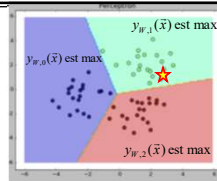
car **dropout ajoute du bruit** à la prédiction

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

↑
masque aléatoire

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

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

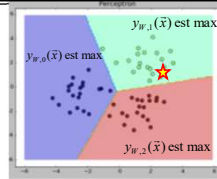


Si on lance le modèle 10 fois, on aura 10 réponses différentes

```
[ 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]
```

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]
```

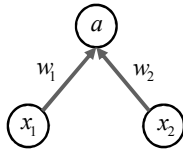
Exécuter le modèle un grand nombre de fois et **prendre la moyenne** revient à calculer **l'espérance mathématique**

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

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

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*.

```

*** 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)
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
    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/>

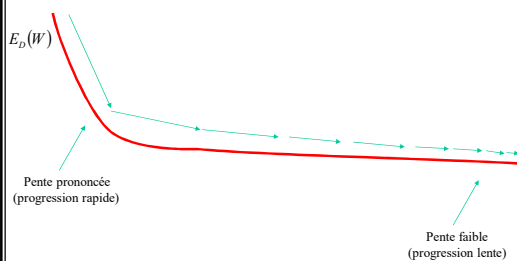
Descente de gradient version améliorée

Descente de gradient

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

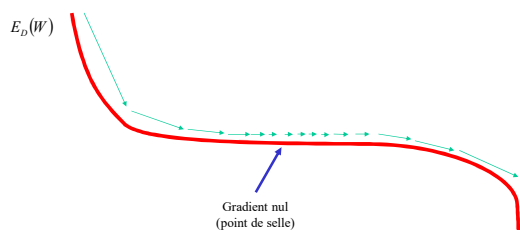
Descente de gradient : **problème**

Progrès quasi nul lorsque la pente est très faible



Descente de gradient : **problème**

Les points de selles sont fréquents en haute dimension



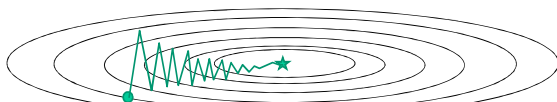
Descente de gradient : **problème**

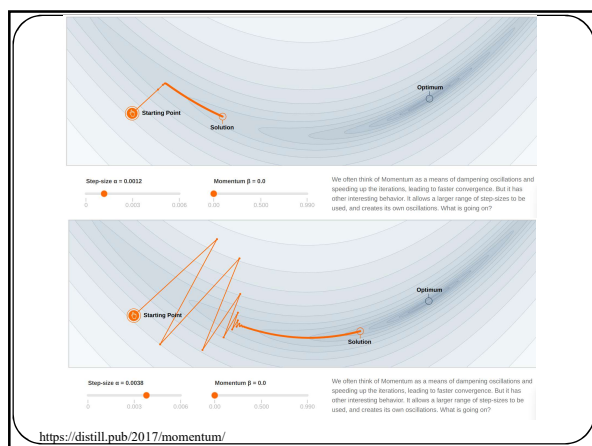
Qu'arrive-t-il si la fonction de coût (loss) a une pente prononcée dans une direction et moins prononcée dans une autre direction?

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.





Descente de gradient + **Momentum**

Descente de gradient
stochastique

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

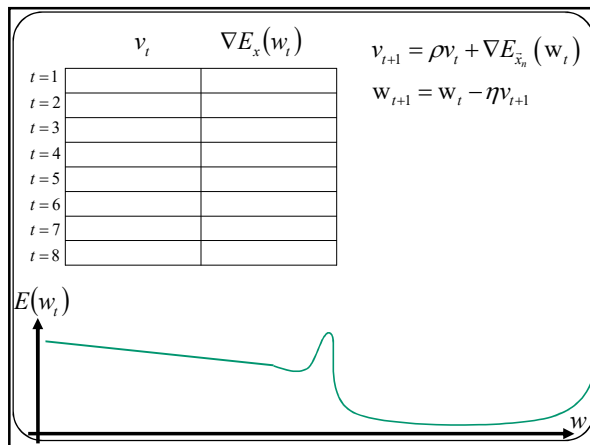
Descente de gradient
stochastique + **Momentum**

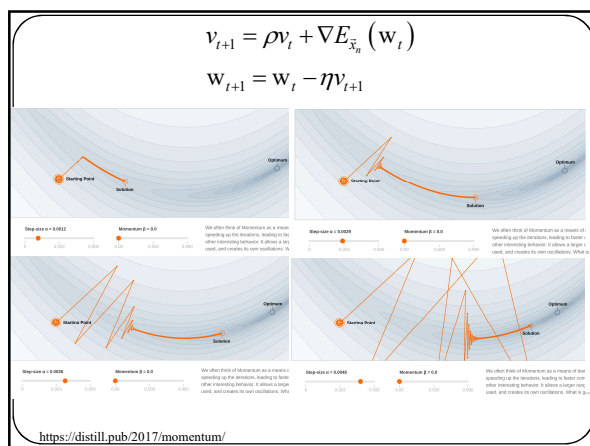
$$v_{t+1} = \rho v_t + \nabla E_{\tilde{x}_n}(w_t)$$

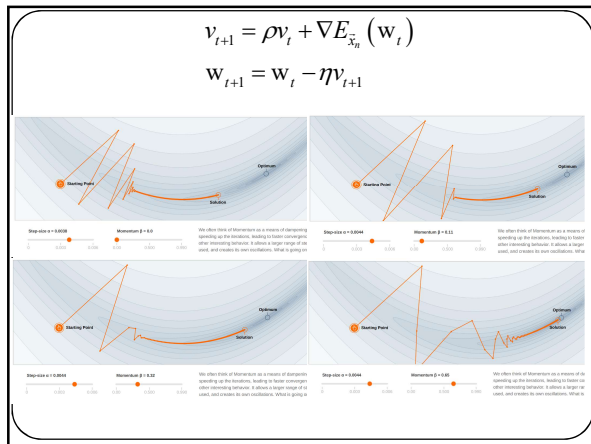
$$w_{t+1} = w_t - \eta v_{t+1}$$

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

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







AdaGrad (décroissance automatique de η)

Descente de gradient stochastique **AdaGrad**

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

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

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

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

AdaGrad (décroissance automatique de η)

Descente de gradient stochastique **AdaGrad**

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

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

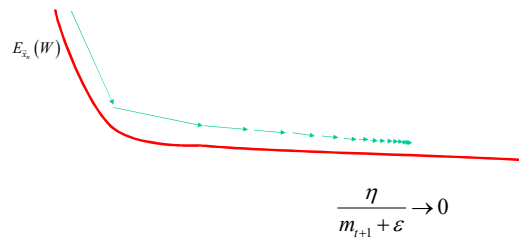
$$m_{t+1} = m_t + |dE_t|$$

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

η décroît sans cesse au fur et à mesure de l'optimisation

AdaGrad (décroissance automatique de η)

Qu'arrive-t-il à long terme?



RMSProp (AdaGrad amélioré)

AdaGrad

$$dE_t = \nabla E_{x_t}(w_t)$$

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

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

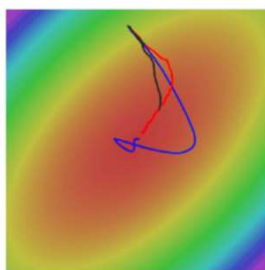
RMSProp

$$dE_t = \nabla E_{x_t}(w_t)$$

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

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

η décroît lorsque le gradient est élevé
 η augmente lorsque le gradient est faible



— SGD
 — SGD+Momentum
 — RMSProp

Adam (Combo entre Momentum et RMSProp)

Momentum

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

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

Adam

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

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

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

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

Adam (Combo entre Momentum et RMSProp)

Momentum

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

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

Adam

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

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

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

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

Momentum

Adam (Combo entre Momentum et RMSProp)

RMSProp

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

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

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

Adam

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

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

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

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

RMSProp

Adam (Version complète)

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

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

for t=1 à num_iterations

for n=0 à N

$$dE_t = \nabla E_{x_t}(w_t)$$

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

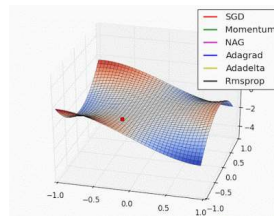
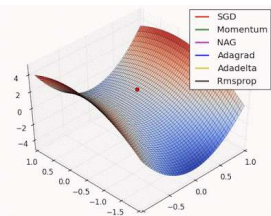
$$m_{t+1} = \gamma m_t + (1 - \gamma) |dE_t|$$

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

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

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

Illustrations



À voir sur :

www.denizvuret.com/2015/03/alec-radfords-animations-for.html

Autre excellent survol

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

