



UNIVERSITÉ PARIS 13  
NORD




# Apprentissage par réseaux de neurones artificiels

Artificial Neural Networks Learning

Younès BENNANI  
Professeur

EPAT<sup>14</sup> : École de Printemps sur l'Apprentissage arTificiel 2014  
7-12 juin 2014 Carry-le-Rouet (France)



## Plan du cours

- **Éléments de base (le neurone, architectures, paramètres)**
- **Critères et algorithmes d'apprentissage**
- **Réseaux de neurones artificiels supervisés**
- **Réseaux de neurones artificiels non supervisés**
- **Réseaux de neurones artificiels profonds**
- **Propriétés et liens avec l'analyse factorielle**
- **Conclusion**



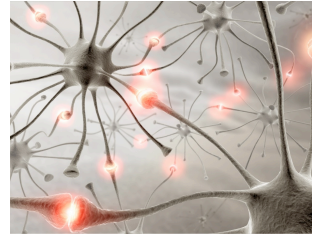
## L'idée d'origine (~ 1940)

• Les organismes vivants, même assez primitifs (ex. insectes), réalisent des tâches complexes de traitement de l'information :

- orientation
- communication
- comportement social
- ...

• La puissance de traitement de leur système nerveux vient de l'interconnexion (ex.  $10^{14}$  connexions chez l'homme)

- d'un grand nombre (ex.  $10^{11}$  chez l'homme)
- d'unités de traitement simples et similaires : les neurones



• La motivation initiale était de faire du neuro-mimétisme

- toutefois, la vision des années 1940 était assez simpliste;
- la réalité biologique c'est avéré plus complexe depuis.

• En revanche, cette idée c'est avérée très féconde en mathématiques et en ingénierie



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

3

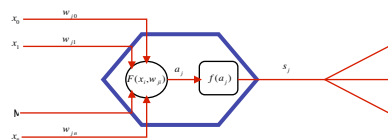
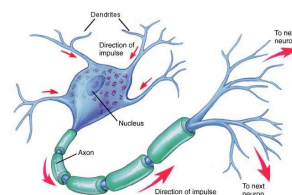


## Le Neurone Formel

1943, Mc Culloch & Pitts

Un modèle et non une copie du neurone biologique  
= un processeur élémentaire caractérisé par :

- signaux d'entrée  $x_0, x_1, \dots, x_n$
- poids des connexions  $w_{i0}, w_{i1}, \dots, w_{in}$
- fonction d'activation  $F(x, w)$
- état interne d'activation  $a = F(x, w)$
- fonction de transition  $f(a)$
- état de sortie  $s = f(a)$



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

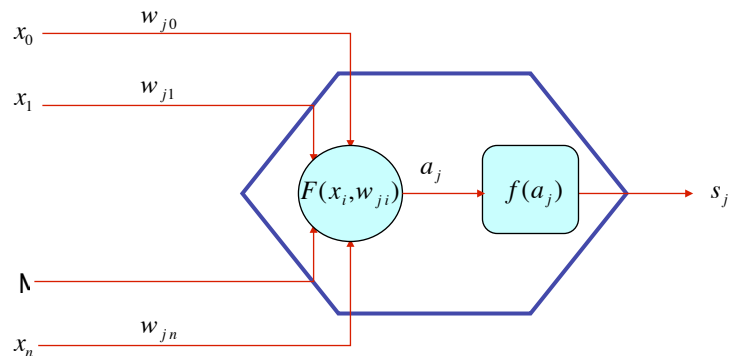
Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

4



## Le Neurone Formel

**Définition :** *Un neurone formel (artificiel) est une unité de traitement qui reçoit des données en entrée, sous la forme d'un vecteur, et produit une sortie réelle. Cette sortie est une fonction des entrées et des poids des connexions.*



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

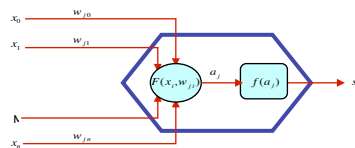
Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

5



## Poids et connexion

**Définition :** *Une connexion entre deux unités  $i$  et  $j$  indique la possibilité d'une relation physique entre ces deux unités.*



**Définition :** *La valeur numérique du poids associé à une connexion entre deux unités reflète la force de la relation entre ces deux unités. Si cette valeur est positive, la connexion est dite excitatrice, sinon elle est dite inhibitrice. La convention usuelle est de noter le poids de la connexion reliant le neurone  $i$  au neurone  $j$  :  $w_{ji}$ .*



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

6

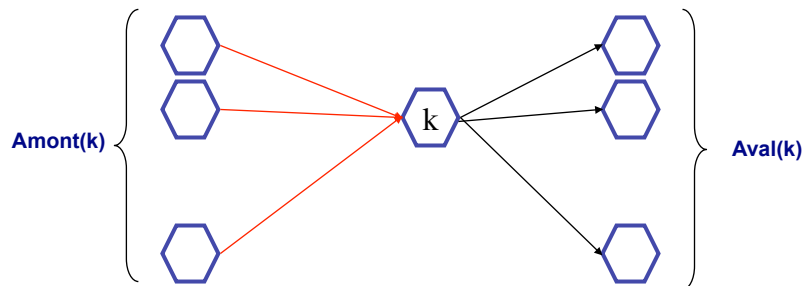




## Le Neurone Formel

**Définition :** Si l'on considère l'unité  $k$ , et que l'on appelle  $Amont(k)$  l'ensemble des neurones dont les sorties servent d'entrées au neurone  $k$ ,  $Aval(k)$  l'ensemble des neurones qui utilisent comme entrée la sortie du neurone  $k$ .

$E$  l'ensemble des neurones d'entrée, et  $S$  l'ensemble des neurones de sortie.  
Alors, par définition on a :  $\forall i \in E, Amont(i) = \emptyset$  et  $\forall i \in S, Aval(i) = \emptyset$ .



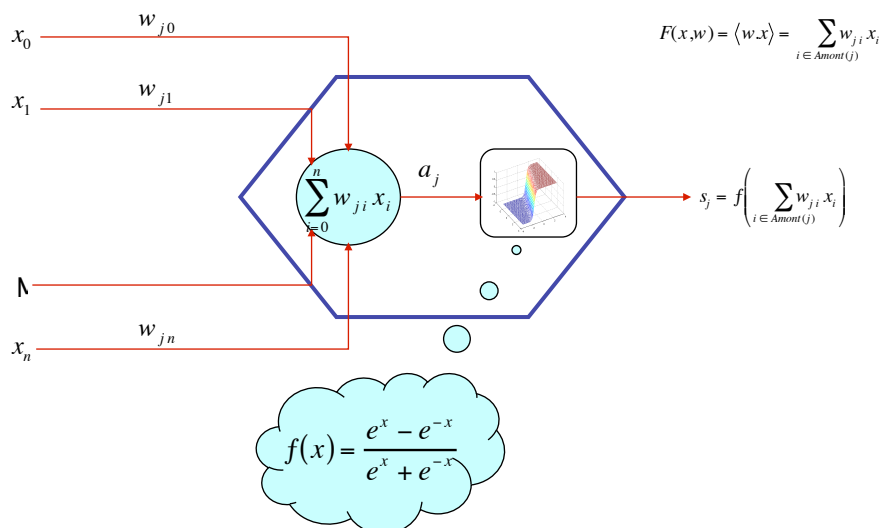
EPAT 14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

7



## Le Neurone Produit Scalaire



EPAT 14, Carry-Le-Rouet 7-12 juin 2014

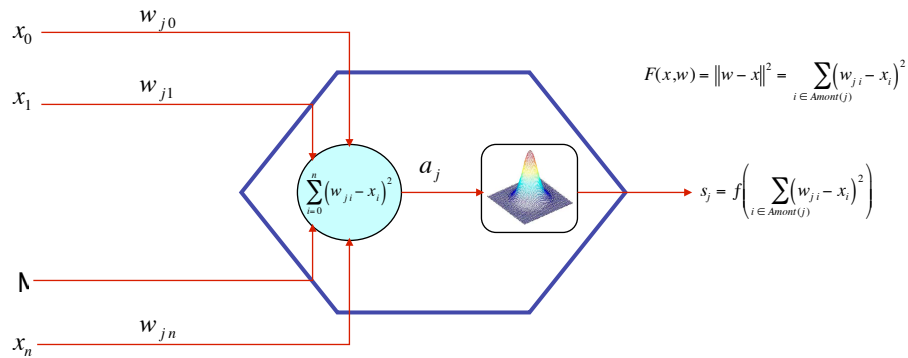
Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

8





## Le Neurone Distance



$$\|w - x\|^2 = \|w\|^2 - 2\langle w, x \rangle + \|x\|^2$$



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

9

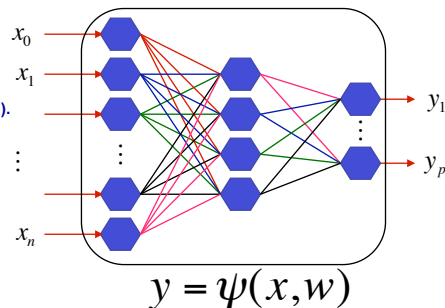


## Réseau Connexionniste / Réseau de neurones

**Définition :** *Un réseau de neurones est un graphe valué orienté, constitué d'un ensemble d'unités (ou automates), réalisant des calculs élémentaires, structurées en couches successives capables d'échanger des informations au moyen de connexions qui les relient.*

- Architecture massivement parallèles.  
- Système basé sur la coopération de plusieurs unités simples (neurones formels).

- Un réseau se caractérise par :  
- son architecture  
- les fonctions de ses éléments



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

10





## Fonctions d'un réseau connexionniste

### Le Mode Apprentissage :

Le réseau est modifié en fonction des exemples d'apprentissage

On peut modifier :

- les poids de ses connexions
- son architecture
- les fonctions de transition des unités

Il existe essentiellement deux sortes d'algorithmes d'apprentissage :

- l'apprentissage supervisé
- l'apprentissage non supervisé



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

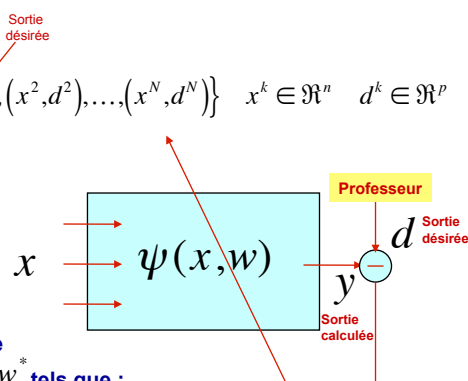
11



## Apprentissage Supervisé

On se donne :

- **N exemples étiquetés**  $E_N = \{(x^1, d^1), (x^2, d^2), \dots, (x^N, d^N)\}$   $x^k \in \mathfrak{R}^n$   $d^k \in \mathfrak{R}^p$
- **une architecture de réseau**  $A$
- **des poids initiaux**  $w_0$



On cherche, au moyen de l'algorithme

d'apprentissage à trouver des poids  $w^*$  tels que :

- les exemples sont reconnus :  $d^k = \psi(x^k, w)$
  - on obtient une bonne généralisation :  $d = \psi(x, w)$
- est une réponse raisonnable pour l'entrée  $x$



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

12

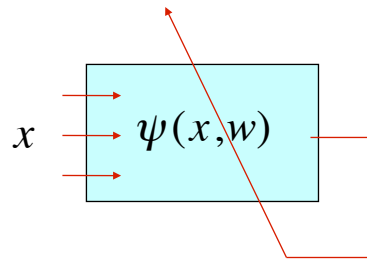




## Apprentissage Non Supervisé

**On se donne :**

- **N exemples non étiquetés**  $E_N = \{x^1, x^2, \dots, x^N\}$
- **une architecture de réseau**  $A$
- **des poids initiaux**  $w_0$



**On cherche à trouver des poids  $w^*$  tels que :**

- **les exemples sont correctement regroupés**
- **on obtient une bonne généralisation**



EPAT 14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

13



## Mode Reconnaissance

**Le réseau est utilisé pour traiter des données**

- **les poids et l'architecture du réseau sont fixés**
- **on lui présente un exemple  $x$  en entrée et le réseau**

**calcule une sortie  $y = \psi(x, w)$  qui dépend de l'architecture  $A$  et des poids  $w^*$  appris.**

**La sortie peut être :**

- **l'étiquette de la classe (identité)**
- **le N° du cluster**
- **la valeur à prévoir**
- **un vecteur du même type que l'entrée**

**Classement/Discrimination**

**Classification/Quantification**

**Prévision**

**Codage/Compaction**



EPAT 14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

14

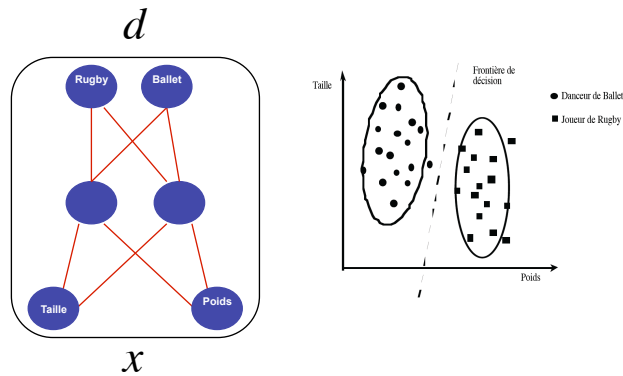




## Utilisation des modèles connexionnistes

Grand nombre de mesures + Loi sous-jacente inconnue

**Classement**



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

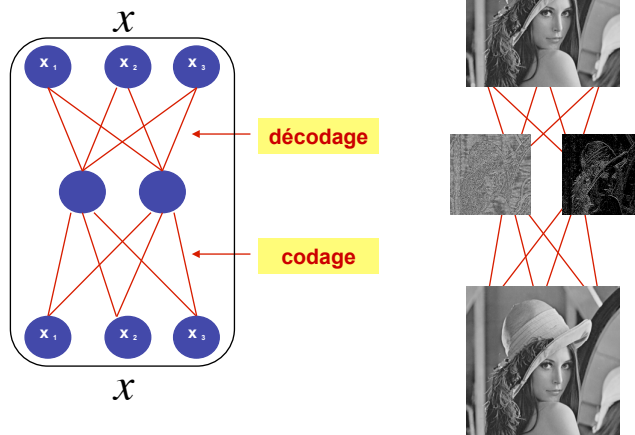
15



## Utilisation des modèles connexionnistes

Grand nombre de mesures + Loi sous-jacente inconnue

**Compression**



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

16

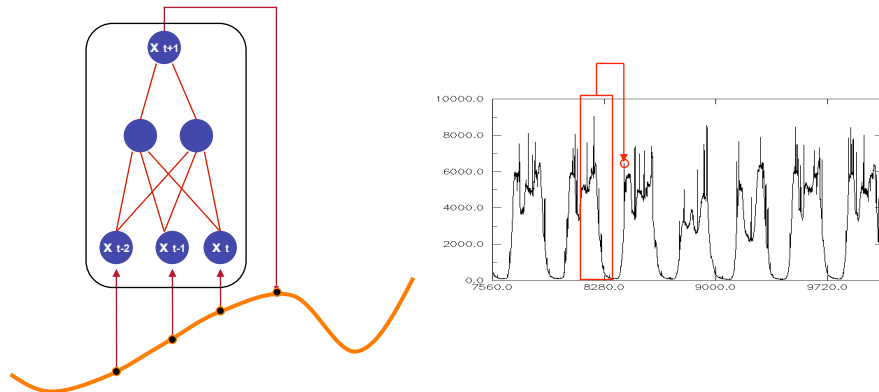




## Utilisation des modèles connexionnistes

Grand nombre de mesures + Loi sous-jacente inconnue

Régression



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

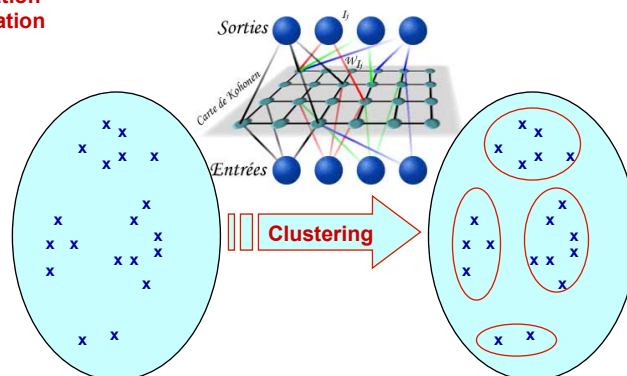
17



## Utilisation des modèles connexionnistes

Grand nombre de mesures + Loi sous-jacente inconnue

Classification  
Quantification



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

18





## Apprentissage à partir d'exemples

### But :

- Construire un système qui explique des données ou une relation pertinente liant ces données
- Estimer un modèle statistique qui a généré les données

### Connaissances :

- Données
- Domaine

### Démarche :

- Choisir un système possédant une bonne capacité de généralisation
- Agir sur la complexité effective des modèles



## Formalisme

La problématique de l'apprentissage se présente souvent comme la minimisation d'une fonction de risque :

$$R(w) = \int_{\mathcal{Z}} L(z, w) dp(z)$$

Risque théorique  
 L'erreur de généralisation

Fonction de perte  
 (Loss function)

Où  $z$ , à valeur dans un espace  $\mathcal{Z}$  (des conditions extérieures ou des exemples), représente les observations d'une v.a.  $Z$  de distribution de probabilité  $p(z)$ , fixe mais **inconnue**, et  $w$  les paramètres du réseau.

L'apprentissage revient à trouver les paramètres :

$$w^* = \arg \min_w R(w)$$



## Formalisme

Le risque théorique n'est pas calculable,  $p(z)$  est inconnue.

Mais un échantillon  $E_N = \{z^k\}_{k=1}^N$  d'exemples indépendants, tirés de  $p(z)$  est connu.

En pratique, on ne peut pas minimiser  $R(w)$ , on utilise alors un principe d'induction.

Le plus courant = minimisation du risque empirique (MRE) :

Risque empirique  
(Erreur d'apprentissage)

$$\tilde{R}(w) = \frac{1}{N} \sum_{k=1}^N L(z^k, w)$$



Choix :  $w^+ = \arg \min_w \tilde{R}(w)$



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

21



## Algorithmes d'Optimisation

$$\frac{\partial \tilde{R}(w)}{\partial w} = \lim_{\delta w \rightarrow 0} \frac{\tilde{R}(w + \delta w) - \tilde{R}(w)}{\delta w}$$

La règle d'adaptation :

**Gradient stochastique**

$$w(t+1) = w(t) - \varepsilon(t) \nabla_w \tilde{R}^k(w)$$

$$\tilde{R}^k(w) = L(z^k, w)$$

$$\nabla_w \tilde{R}^k(w) = \left( \frac{\partial \tilde{R}^k(w)}{\partial w_1}, \dots, \frac{\partial \tilde{R}^k(w)}{\partial w_i}, \dots, \frac{\partial \tilde{R}^k(w)}{\partial w_n} \right)$$

**Gradient total**

$$w(t+1) = w(t) - \frac{1}{N} \varepsilon(t) \sum_{k=1}^N \nabla_w \tilde{R}^k(w)$$



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

22





## Systèmes d'Apprentissage Supervisé



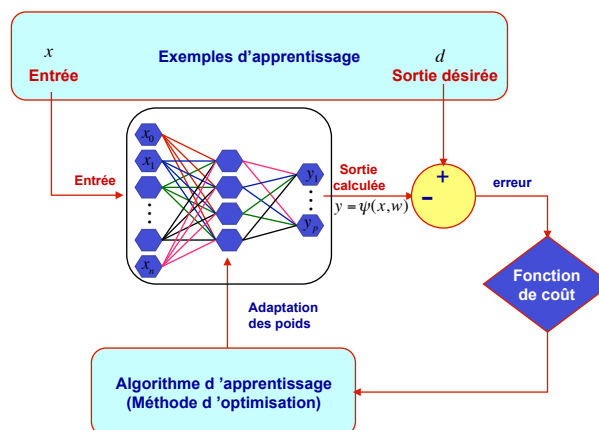
EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

23



## Système d'Apprentissage Supervisé



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

24





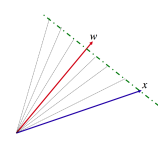
# Adaline : Adaptive Linear Element

Stanford, 1960, Bernard Widrow\*

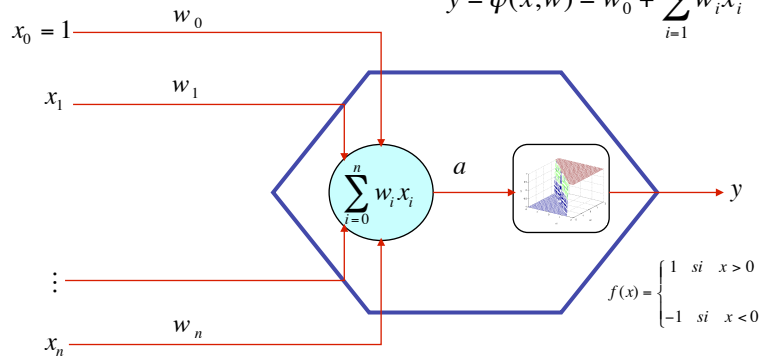
C'est un élément linéaire adaptatif :  $y = \sum_{i=0}^n w_i x_i$

L'unité  $x_0$ , dont l'activation fixée à 1, est dite unité de biais. Elle permet d'introduire les termes de biais dans le réseau.

$$y = \sum_{i=0}^n w_i x_i = \|w\| \|x\| \cos \phi$$



\* Widrow B., Hoff M.E. (1960) : « Adaptive switching circuits », IRE WESCON Conv. Record, part 4, pp. 96-104.



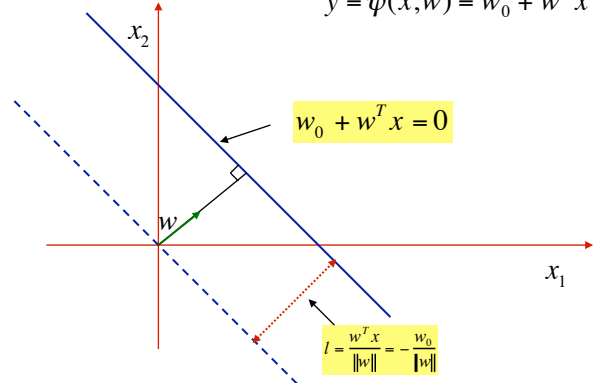
$$y = \psi(x, w) = w_0 + \sum_{i=1}^n w_i x_i$$

$$f(x) = \begin{cases} 1 & \text{si } x > 0 \\ -1 & \text{si } x < 0 \end{cases}$$

# Adaline : Adaptive Linear Element

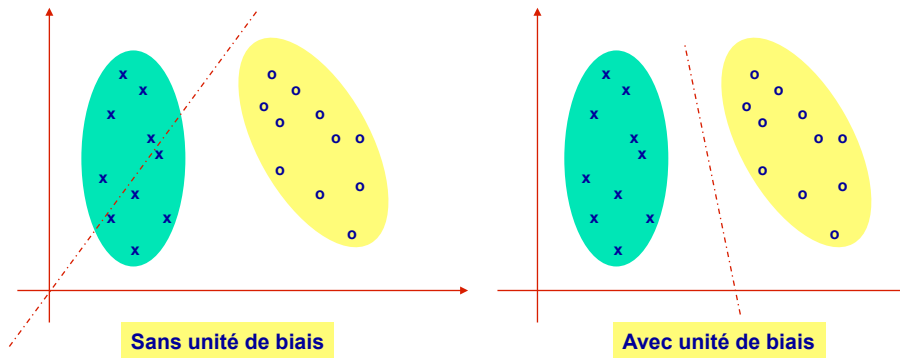
$$y = \psi(x, w) = w_0 + \sum_{i=1}^n w_i x_i$$

$$y = \psi(x, w) = w_0 + w^T x$$





## Adaline : Adaptive Linear Element



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

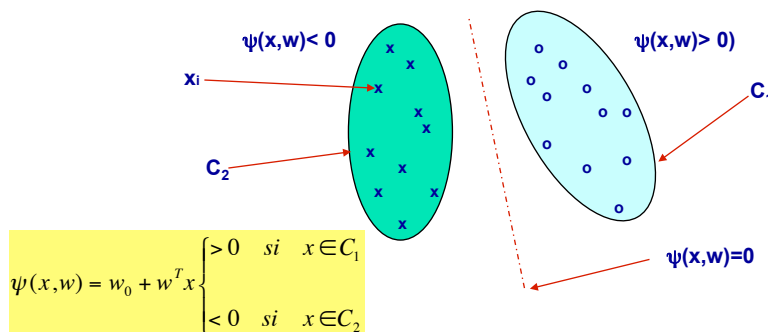
Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

27



## Adaline : Adaptive Linear Element

On peut l'utiliser en discrimination (classement) entre 2 classes :



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

28





## Adaline : Adaptive Linear Element

Si l'on appelle  $z^k = (x^k, d^k)$  la forme prise en compte à l'itération  $k$ ,  
On définit le carré de l'erreur instantanée associée à la forme  $(x^k, d^k)$  par :

$$\tilde{R}_{Adaline}^k(w) = (d^k - wx^k)^2$$

Sortie désirée    Sortie calculée :  $y$

L'erreur quadratique globale ou (MSE) est définie comme la moyenne observée des carrés des erreurs instantanées sur l'ensemble de toutes les formes :

$$\tilde{R}_{Adaline}(w) = \frac{1}{N} \sum_{k=1}^N \tilde{R}_{Adaline}^k(w)$$

Il existe plusieurs algorithmes d'apprentissage.



## Adaline : Adaptive Linear Element

Techniques de descente de gradient (la plus grande pente) :  
supposons qu'à l'instant  $t$ , les poids de l'Adaline soient  $w(t)$   
et qu'on présente la forme  $(x^k, d^k)$ , alors les poids seront modifiés par :

$$w(t+1) = w(t) - \epsilon(t) \nabla_w \tilde{R}_{Adaline}^k(w)$$

Le pas du gradient

Le gradient instantané

$$\nabla_w \tilde{R}_{Adaline}^k(w) = \frac{\partial \tilde{R}_{Adaline}^k(w)}{\partial w} = -2(d^k - wx^k)x^k$$

Cette règle est appelée règle du gradient stochastique  
ou règle de Widrow-Hoff  
ou règle du delta de Widrow-Hoff  
ou règle  $\mu$ -LMS (Least Mean Square)





## Adaline : Adaptive Linear Element

1- Tirer au hasard  $w_0$

2- Présenter une forme  $(x^k, d^k)$

3- Déterminer la valeur de l'écart

$$e^k(t) = (d^k - wx^k)$$

4- Calculer une approximation du gradient

$$\nabla_w \tilde{R}_{Adaline}^k(w) = -2e^k(t)x^k$$

5- Adapter les poids  $w(t)$

$$w(t+1) = w(t) - \varepsilon(t) \nabla_w \tilde{R}_{Adaline}^k(w)$$

Où  $\varepsilon(t)$  est le pas du gradient.

6- Répéter de 2 à 4 jusqu'à l'obtention d'une valeur acceptable de l'erreur



## Adaline : Exemples

**Données :** Table de vérité d'une fonction booléenne de 2 variables

$x_1$	$x_2$	$d$
1	1	1
1	-1	1
-1	1	1
-1	-1	-1

**Fonction :**  $\psi(x,w) = (x_1 \text{ ou } x_2)$

$$E = \left\{ \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}; \begin{pmatrix} 1 \\ -1 \\ 1 \end{pmatrix}; \begin{pmatrix} -1 \\ 1 \\ 1 \end{pmatrix}; \begin{pmatrix} -1 \\ -1 \\ -1 \end{pmatrix} \right\}$$

**Problème :**

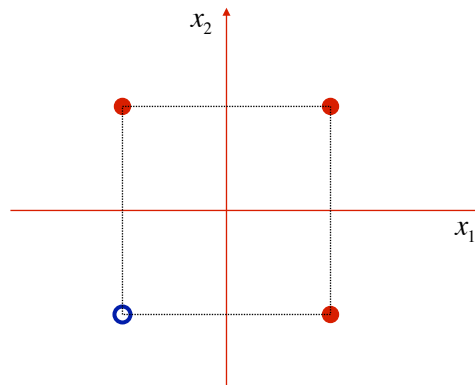
**Trouver un Adaline capable d'apprendre la table de vérité d'une fonction booléenne de 2 variables**



## Adaline : Exemples

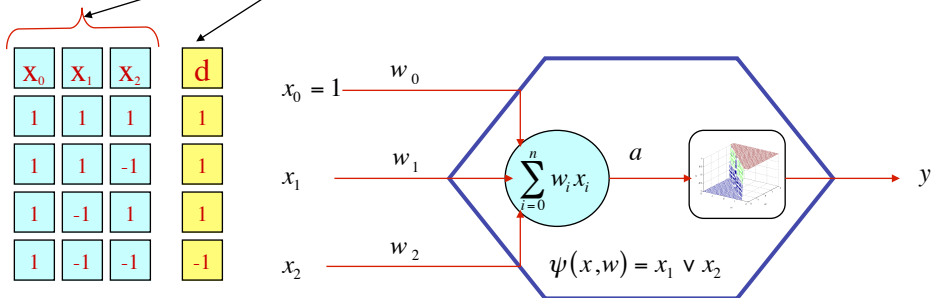
Fonction :  $\psi(x,w) = (x_1 \text{ ou } x_2)$

$x_1$	$x_2$	$d$
1	1	1
1	-1	1
-1	1	1
-1	-1	-1



## Adaline : Exemples

$$E_N = \{(x^1, d^1), (x^2, d^2), \dots, (x^N, d^N)\}$$



### Adaline : Exemples

Inputs:  $x_0 = 1$ ,  $x_1$ ,  $x_n$

Weights:  $w_0$ ,  $w_1$ ,  $w_2$

Summation:  $\sum_{i=0}^n w_i x_i$

Activation:  $a$

Transfer Function:  $\psi(x, w) = x_1 \vee x_2$

$w_1 x_1 + w_2 x_2 + w_0 = 0$

EPAT'14, Carry-Le-Rouet 7-12 juin 2014
Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)
35

### Adaline : Exemples

Inputs:  $x_0 = 1$ ,  $x_1$ ,  $x_2$

Weights:  $w_0$ ,  $w_1$ ,  $w_2$

Summation:  $\sum_{i=0}^n w_i x_i$

Activation:  $a$

Transfer Function:  $\psi(x, w) = x_1 \wedge x_2$

$w_1 x_1 + w_2 x_2 + w_0 = 0$

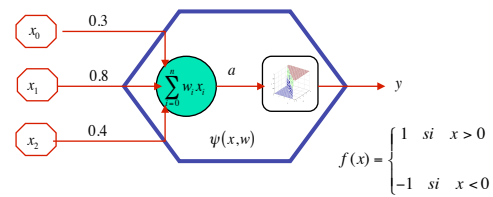
$x_0$	$x_1$	$x_2$	$d$	
1	1	1	1	●
1	1	-1	-1	○
1	-1	1	-1	○
1	-1	-1	-1	○

EPAT'14, Carry-Le-Rouet 7-12 juin 2014
Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)
36



## Exercice à faire

$$E = \left\{ \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, -1 \right\}; \begin{pmatrix} 1 \\ 1 \\ -1 \end{pmatrix}, 1 \right\}; \begin{pmatrix} 1 \\ -1 \\ 1 \end{pmatrix}, 1 \right\}; \begin{pmatrix} 1 \\ -1 \\ -1 \end{pmatrix}, 1 \right\}$$



- Représenter dans un repère orthogonal l'ensemble des échantillons.
- Utiliser l'algorithme Adaline pour adapter les poids du modèle ( $\epsilon=0.1$ ).
- Donner l'équation de l'hyperplan séparant les deux classes.
- Représenter l'hyperplan dans le même repère orthogonal



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

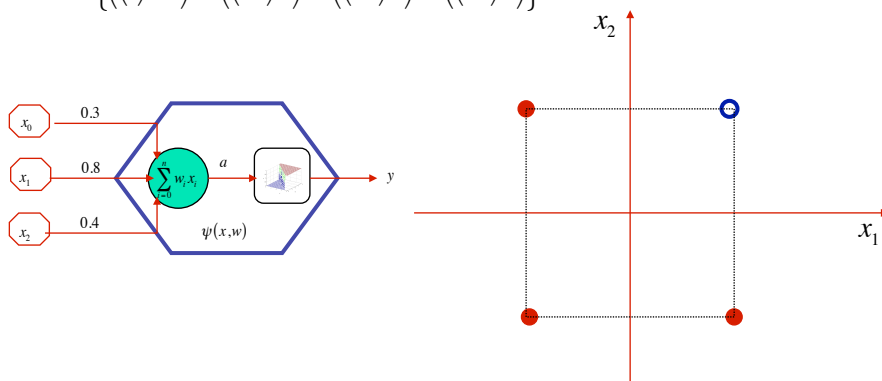
Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

37



## Représentation graphique

$$E = \left\{ \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, -1 \right\}; \begin{pmatrix} 1 \\ 1 \\ -1 \end{pmatrix}, 1 \right\}; \begin{pmatrix} 1 \\ -1 \\ 1 \end{pmatrix}, 1 \right\}; \begin{pmatrix} 1 \\ -1 \\ -1 \end{pmatrix}, 1 \right\}$$



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

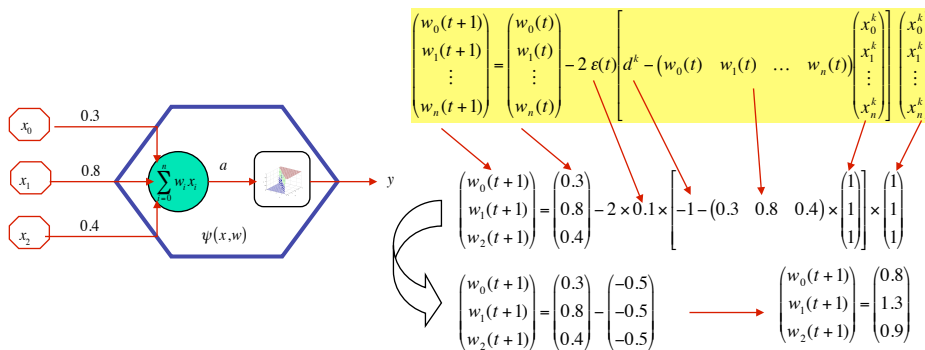
Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

38

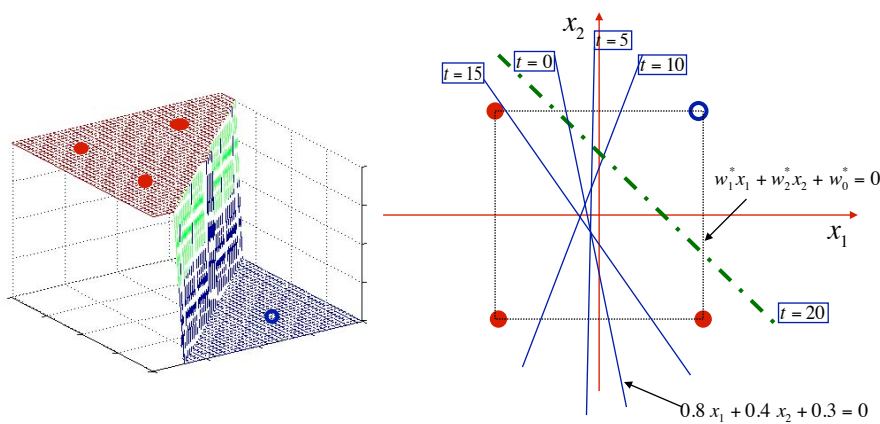


## Adaptation des poids

$$D = \left\{ \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, -1; \begin{pmatrix} 1 \\ 1 \\ -1 \end{pmatrix}, 1; \begin{pmatrix} 1 \\ -1 \\ 1 \end{pmatrix}, 1; \begin{pmatrix} 1 \\ -1 \\ -1 \end{pmatrix}, 1 \right\}$$



## Évolution de l'apprentissage

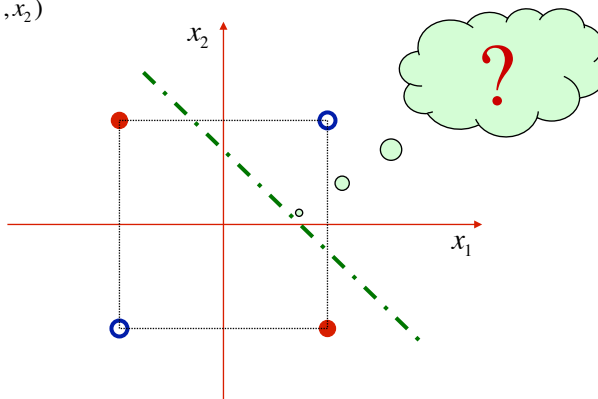




## Adaline : limites

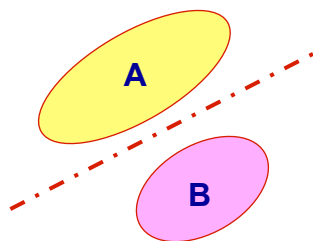
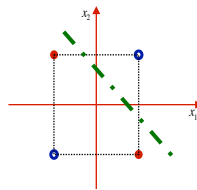
$$\psi(x, w) = x_1 \oplus x_2 = \text{XOR}(x_1, x_2)$$

	$x_0$	$x_1$	$x_2$	$d$	
$x^1$	1	1	1	-1	○
$x^2$	1	1	-1	1	●
$x^3$	1	-1	1	1	●
$x^4$	1	-1	-1	-1	○

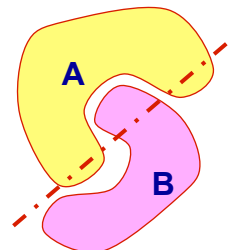


## Séparabilité linéaire

« Deux classes d'objets, décrits dans un espace de dimension  $n$ , sont dits « linéairement séparables » s'ils se trouvent de part et d'autre d'un hyperplan dans l'espace des descripteurs »



Linéairement séparable

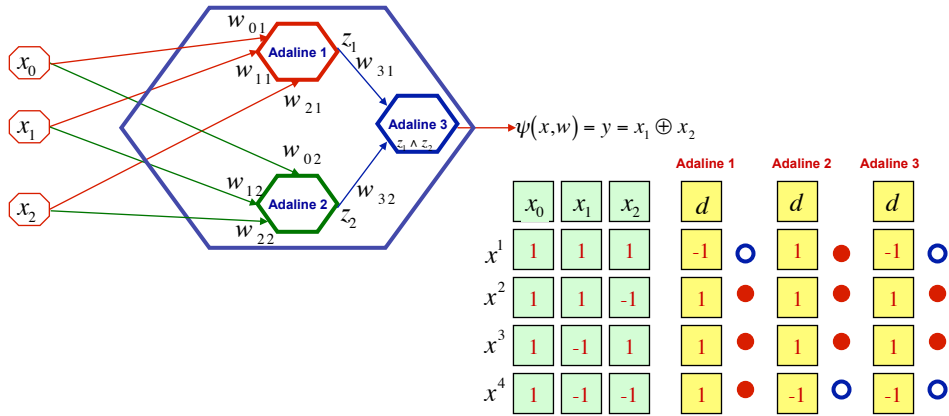


Non-linéairement séparable



# Madaline : Multi-Adaptive Linear Element

Madaline = un ensemble d'Adalines parallèles



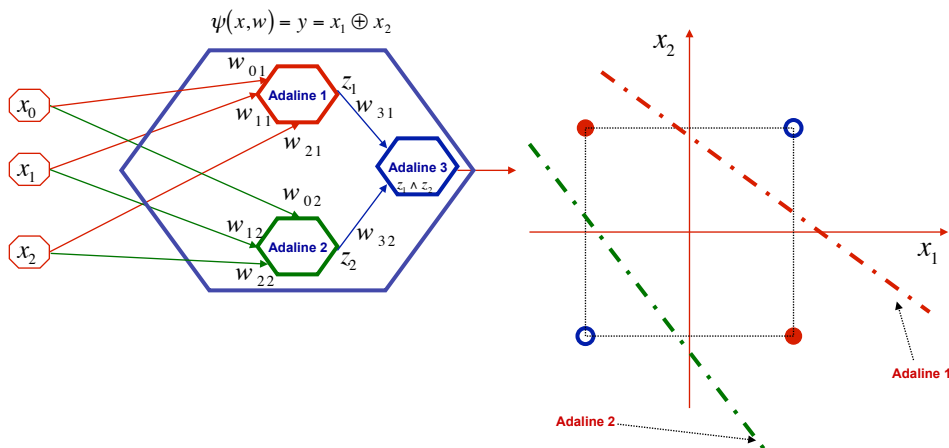
EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

43



# Madaline : Multi-Adaptive Linear Element



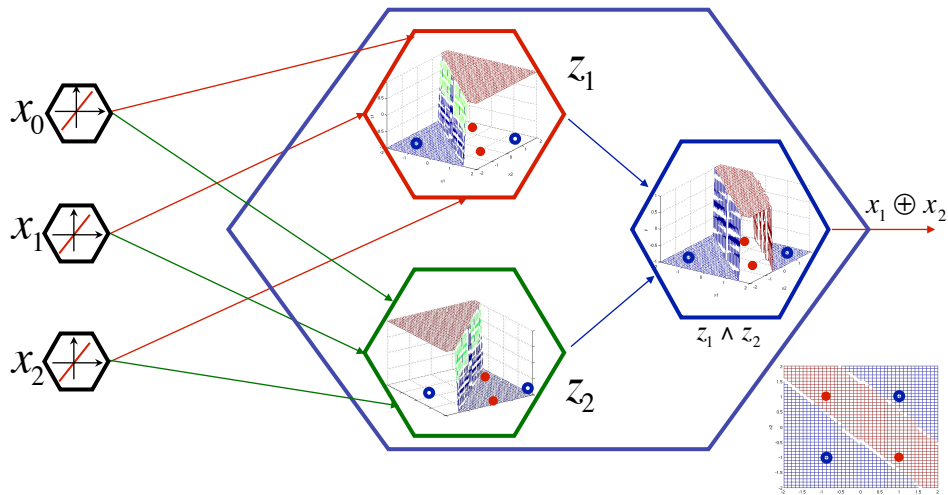
EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

44

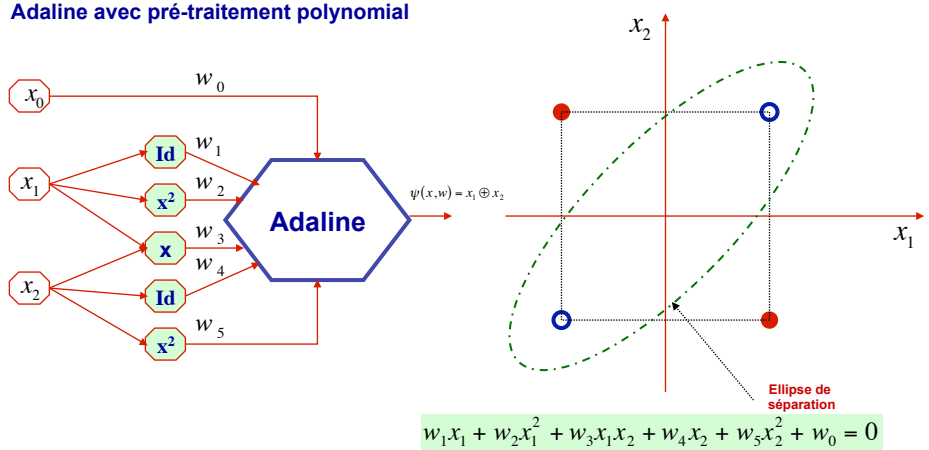


PARIS 13 Madaline : Multi-Adaptive Linear Element



PARIS 13 Madaline : Multi-Adaptive Linear Element

Adaline avec pré-traitement polynomial



# Perceptron

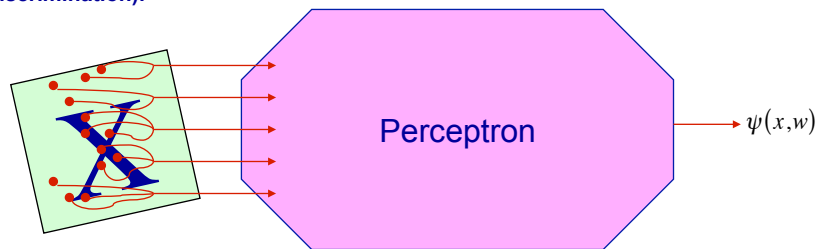
Rosenblatt F., 1957, 1962\*

1957, Frank Rosenblatt

Le perceptron ne désigne pas un seul modèle mais regroupe une importante famille d'algorithmes.



Le perceptron = machine adaptative employée pour résoudre des problèmes de **classement** (discrimination).



\* Rosenblatt F. (1957) : « The perceptron: a perceiving and recognizing automaton », Reports 85-460-1, Cornell Aeronautical Lab., Ithaca, N.Y.  
 \* Rosenblatt F. (1962) : « Principles of Neurodynamics: perceptrons and theory of brain mechanisms », Spartan Books, Washington.



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

47



# Perceptron

Rosenblatt F., 1957, 1962

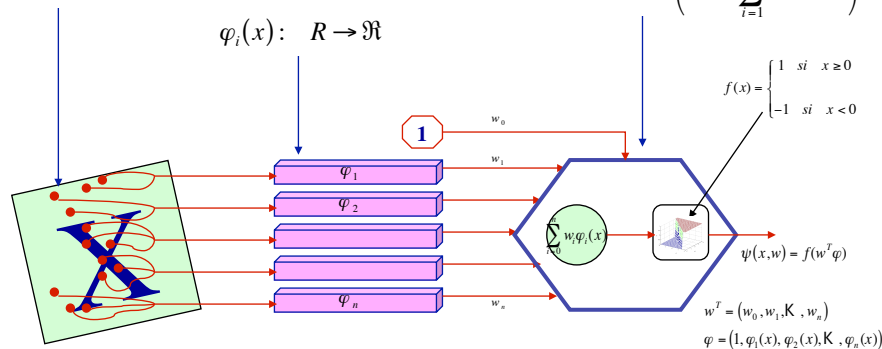
**la rétine R**  
qui reçoit les informations de l'extérieur

**les cellules d'association**  
chaque cellule possède une fonction de transition définie sur la rétine :

$$\varphi_i(x) : R \rightarrow \mathfrak{R}$$

**la cellule de décision**

$$\psi(x, w) = f\left(w_0 + \sum_{i=1}^n w_i \varphi_i(x)\right)$$



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

48



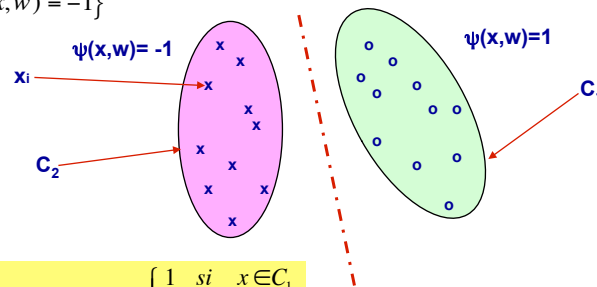
## Perceptron

Cas de deux classes

Un perceptron peut être vu comme un classificateur à 2 classes :

$$C_1 = \{x \in R : \psi(x, w) = 1\}$$

$$C_2 = \{x \in R : \psi(x, w) = -1\}$$



$$\psi(x, w) = f\left(w_0 + \sum_{i=1}^n w_i \varphi_i(x)\right) = \begin{cases} 1 & \text{si } x \in C_1 \\ -1 & \text{si } x \in C_2 \end{cases}$$



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

49



## Perceptron

Cas de deux classes

Si l'on appelle  $(x^k, d^k)$  la forme prise en compte à l'itération  $k$ ,

On définit le carré de l'erreur instantanée associée à la forme  $(x^k, d^k)$  par :

$$\begin{cases} d^k = 1 & \text{si } x^k \in C_1 \\ d^k = -1 & \text{si } x^k \in C_2 \end{cases} \Rightarrow \begin{cases} w^T \varphi^k > 0 & \text{pour } x^k \in C_1 \\ w^T \varphi^k < 0 & \text{pour } x^k \in C_2 \end{cases} \Rightarrow \forall x^k, w^T(\varphi^k d^k) > 0$$

$$\varphi^k = (1, \varphi_1(x^k), \varphi_2(x^k), \dots, \varphi_n(x^k))$$

$$\tilde{R}_{\text{Perceptron}}^k(w) = -w^T(\varphi^k d^k)$$

L'erreur quadratique globale ou (MSE) est :

$$\tilde{R}_{\text{Perceptron}}(w) = \frac{1}{N} \sum_{k=1}^N \tilde{R}_{\text{Perceptron}}^k(w) = - \sum_{\{k: x^k \text{ mal classé}\}} w^T(\varphi^k d^k)$$

Il existe plusieurs algorithmes d'apprentissage.



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

50



# Perceptron

Cas de deux classes

Techniques de descente de gradient (la plus grande pente) :  
 supposons qu'à l'instant  $t$ , les poids du Perceptron soient  $w(t)$   
 et qu'on présente la forme  $(x^k, d^k)$ , alors les poids seront  
 modifiés par :

$$w(t+1) = w(t) - \varepsilon(t) \nabla_w \tilde{R}_{\text{Perceptron}}^k(w)$$

Le pas du gradient

Le gradient instantané

$$\nabla_w \tilde{R}_{\text{Perceptron}}^k(w) = \frac{\partial \tilde{R}_{\text{Perceptron}}^k(w)}{\partial w} = -\varphi^k d^k$$

$$\varphi^k = (1, \varphi_1(x^k), \varphi_2(x^k), K, \varphi_n(x^k))$$

Critère d'arrêt = taux de classement satisfaisant.



# Perceptron

## Algorithme d'apprentissage : cas de 2 classes

1- A  $t=0$ , tirer au hasard  $w(0) = (w_0(0), w_1(0), K, w_n(0))^T$

2- Présenter une forme  $(x^k, d^k)$

3- Calculer la sortie du perceptron et la comparer à  $y^k$

$$f\left(w_0 + \sum_{i=1}^n w_i \varphi_i(x)\right) = f(w^T \varphi^k)$$

4- Si  $x^k$  est bien classé :  $f(w^T \varphi^k) = d^k$

$$w(t+1) = w(t)$$

Si  $x^k$  est mal classé :  $f(w^T \varphi^k) \neq d^k$

$$w(t+1) = w(t) + \varepsilon(t) \varphi^k d^k$$

$$\varphi^k = (1, \varphi_1(x^k), \varphi_2(x^k), K, \varphi_n(x^k))$$

$$\begin{pmatrix} w_0(t+1) \\ w_1(t+1) \\ \vdots \\ w_n(t+1) \end{pmatrix} = \begin{pmatrix} w_0(t) \\ w_1(t) \\ \vdots \\ w_n(t) \end{pmatrix} + \varepsilon(t) d^k \begin{pmatrix} 1 \\ \varphi_1(x^k) \\ \vdots \\ \varphi_n(x^k) \end{pmatrix}$$

Où  $\varepsilon(t)$  est le pas du gradient.

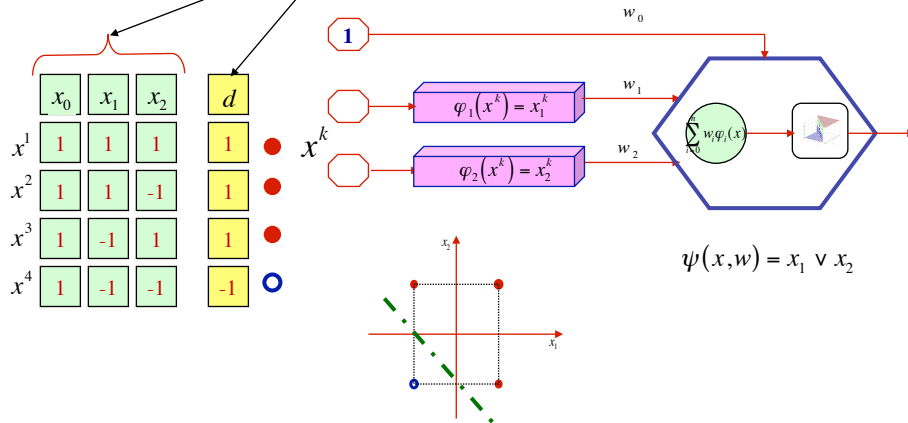
5- Répéter de 2 à 4 jusqu'à l'obtention d'une valeur acceptable de l'erreur



## Exemple de Perceptron

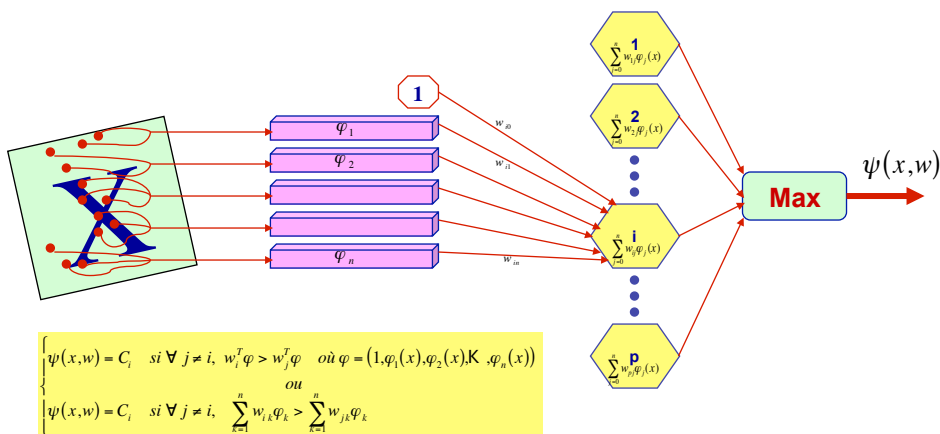
Applet : <http://lcn.epfl.ch/tutorial/french/perceptron/html/index.html>

$$E_N = \{(x^1, d^1), (x^2, d^2), \dots, (x^N, d^N)\}$$



## Perceptron

Cas de p classes :  $C_1, C_2, \dots, C_p$



$$\begin{cases} \psi(x, w) = C_i & \text{si } \forall j \neq i, w_j^i \varphi_j > w_j^i \varphi_j \text{ où } \varphi = (1, \varphi_1(x), \varphi_2(x), \dots, \varphi_n(x)) \\ \text{ou} \\ \psi(x, w) = C_i & \text{si } \forall j \neq i, \sum_{k=1}^n w_{ik} \varphi_k > \sum_{k=1}^n w_{jk} \varphi_k \end{cases}$$



## Perceptron

### Algorithme d'apprentissage : cas de p classes

- 1- A  $t=0$ , tirer au hasard la matrice des poids  $w(0)$
- 2- Présenter une forme  $x^k \in C_i$
- 3- Calculer la sortie  $\psi(x^k, w)$  du perceptron et la comparer à  $C_i$   

$$\psi(x^k, w) = C_j \Leftrightarrow j = \text{Arg max}_i (w_i^T \varphi^k)$$
- 4- Si  $x^k$  est bien classé :  $C_j = C_i$   
 $w(t+1) = w(t)$   
 Si  $x^k$  est mal classé :  $C_j \neq C_i$   
 $w_i(t+1) = w_i(t) + \varepsilon(t) \varphi^k$   
 $w_j(t+1) = w_j(t) - \varepsilon(t) \varphi^k$   
 $w_l(t+1) = w_l(t) \quad \forall l \neq i \text{ et } l \neq j$
- 5- Répéter de 2 à 4 jusqu'à l'obtention d'une valeur acceptable de l'erreur.  
 Où  $\varepsilon(t)$  est le pas du gradient.



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

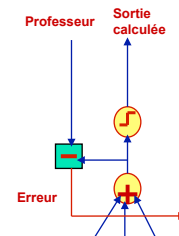
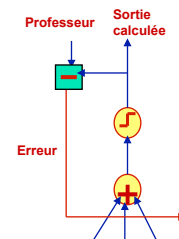
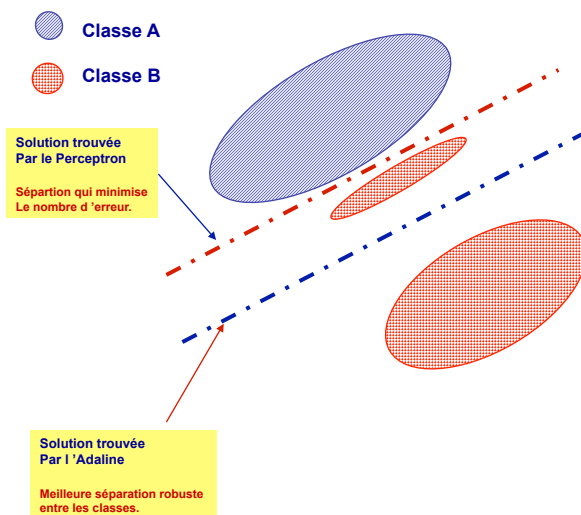
Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

55



## Perceptron et Adaline

Cas de 2 classes



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

56







## Théorème de convergence du Perceptron

« Si un ensemble de formes est **linéairement séparable**, alors l'algorithme d'apprentissage du Perceptron converge vers une solution correcte en un nombre fini d'itération. »

Arbib M.A. (1987) : « Brains, Machines, and Mathematics »  
Berlin, Springer-Verlag.

Rosenblatt F. (1962) : « Principles of Neurodynamics »  
N.Y., Spartan.

Block H.D. (1962) : « The Perceptron: A Model for Brain Functioning »  
Reviews of Modern Physics 34, 123-135.

Minsky M.L. & Papert S.A. (1969) : « Perceptrons »  
Cambridge, MIT Press.

Diederich S. & Opper M. (1987) : « Learning of Correlated Patterns  
Spin-Glass Networks by Local Learning Rules »  
Physical Review Letters 58, 949-952.



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

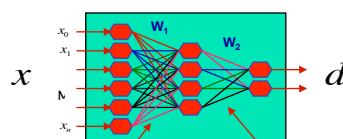
57



## Architecture multi-couches

Le « **credit assignment problem** »

On se donne un réseau en couches, et un ensemble d'exemples composés de paires entrées-sorties.



On ne connaît pas les sorties  
Désirées des unités cachées !

On ne peut pas appliquer  
l'algorithme d'apprentissage  
du Perceptron pour  
déterminer  $w_1$

Perceptron  
Appliquer l'algorithme  
d'apprentissage pour  
déterminer  $w_2$

L'algorithme de **rétro-propagation du gradient** apporte une solution d'une simplicité déroutante à ce problème.



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

58





## L'après pause connexionniste

### Problèmes de contrôle

1963, Bryson A., Denham W., Dreyfuss S.,  
"la rétro-propagation du gradient"

Les années 80, la portée exacte de l'ouvrage de Minsky & Papert sera correctement perçue.

### La "redécouverte" de la "rétro-propagation du gradient"

1986, LeCun Y.

1986, Rumelhart D., Hinton G.E., Williams R.

Bryson A., Denham W., Dreyfuss S. (1963) : « Optimal Programming Problem With Inequality Constraints. I: Necessary Conditions for Extremal Solutions », AIAA Journal, Vol. 1, pp. 25-44.

LeCun Y. (1986) : « Learning Processes in Asymmetric Threshold Network » Disordered Systems and Biological Organizations, Les Houches, France, Springer, pp. 223-240.

Rumelhart D., Hinton G.E., Williams R. (1986) : « Learning Internal Representations by Error Propagation » In Parallel Distributed Processing: exploring the microstructure of cognition, Vol 1, Badford Books, Cambridge, MA, pp. 318-362, MIT Press.



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

59



## Perceptron Multi-Couche (PMC) Multi-Layer Perceptron (MLP)

### Architecture :

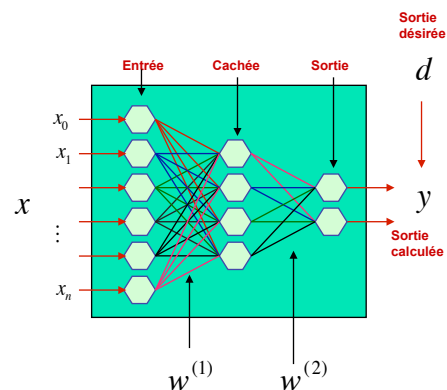
semblable à celle du Perceptron ou de Madaline  
+ des couches de traitement intermédiaire  
(couches cachées)

- Couches externes :  
Entrée (e unités),  
Sortie (s unités)
- Couches internes :  
Cachées (c unités)

Notation : < e | c | s > exemple : < 6 | 4 | 2 >

But :  $E_N = \{(x^1, d^1), (x^2, d^2), \dots, (x^N, d^N)\}$

$$\{x^k \rightarrow d^k\}^{k=1 \dots N}$$



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

60





## Perceptron Multi-Couche (PMC) Multi-Layer Perceptron (MLP)

**Notations :**

$E$  : l'ensemble des unités d'entrée

$S$  : l'ensemble des unités de sortie

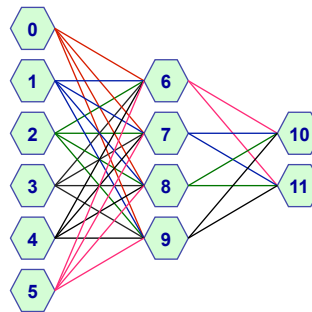
$Amont(k)$  : l'ensemble des unités dont les sorties servent d'entrées à l'unité  $k$

$Aval(k)$  : l'ensemble des unités qui utilisent comme entrée la sortie de  $k$

**Par définition, on a :**

$$\forall i \in E, Amont(i) = \emptyset$$

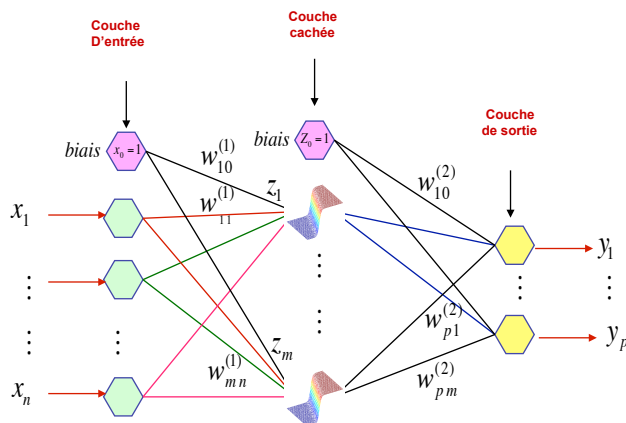
$$\forall i \in S, Aval(i) = \emptyset$$



$Amont(7) = \{0, 1, 2, 3, 4, 5\}$        $Aval(7) = \{10, 11\}$   
 $Amont(2) = \emptyset$        $Aval(2) = \{6, 7, 8, 9\}$   
 $Amont(11) = \{6, 7, 8, 9\}$        $Aval(11) = \emptyset$



## Perceptron Multi-Couche (PMC) Multi-Layer Perceptron (MLP)





## Perceptron Multi-Couche (PMC) Multi-Layer Perceptron (MLP)

L'activation de la  $j^{\text{ème}}$  cellule cachée est :

$$a_j = w_{0j}^{(1)} + \sum_{i \in \text{Amont}(j)} w_{ij}^{(1)} x_i$$

La sortie de cette  $j^{\text{ème}}$  cellule cachée s'obtient par une transformation non linéaire de l'activation :

$$z_j = f(a_j)$$

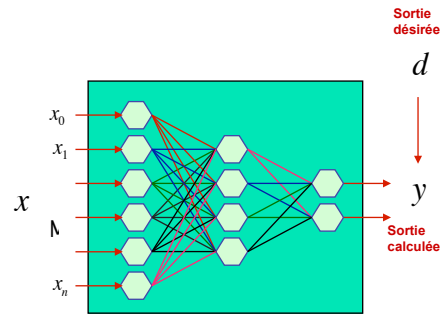
De la même façon, l'activation et la sortie de la  $k^{\text{ème}}$  unité de sortie peuvent s'obtenir comme suit :

$$a_k = w_{0k}^{(2)} + \sum_{j \in \text{Amont}(k)} w_{jk}^{(2)} z_j$$

$$y_k = f(a_k)$$

Si on combine le calcul des sorties des cellules cachées et celui des cellules de sortie on obtient pour la  $k^{\text{ème}}$  sortie du réseau l'expression suivante :

$$y_k = f \left( w_{0k}^{(2)} + \sum_{j \in \text{Amont}(k)} w_{jk}^{(2)} f \left( w_{0j}^{(1)} + \sum_{i \in \text{Amont}(j)} w_{ij}^{(1)} x_i \right) \right)$$



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

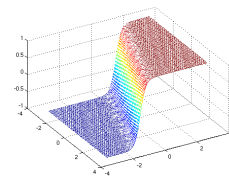
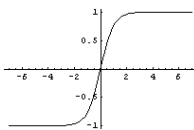
Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

63



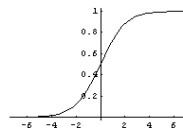
## Perceptron Multi-Couche (PMC) Types de fonction $f$

$$f(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$



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

$$f(x) = \frac{1}{1 + e^{-x}} \rightarrow f'(x) = f(x)(1 - f(x))$$



$$f(x) = \frac{x}{1 + |x|}$$



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

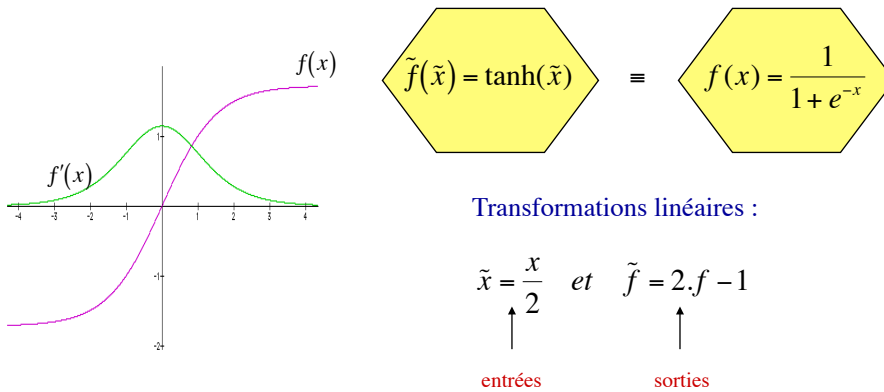
64





## Perceptron Multi-Couche (PMC)

Types de fonction  $f$



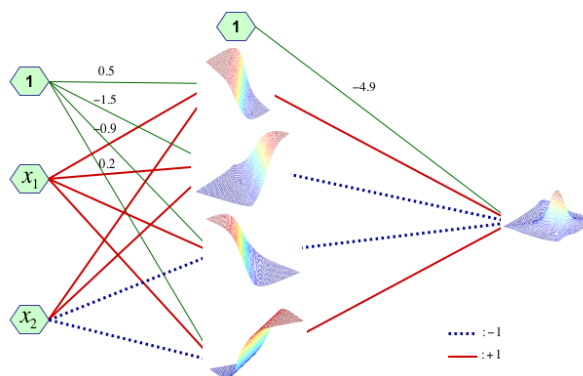
EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

65



## Surfaces de séparation et PMC



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

66





## Perceptron Multi-Couche (PMC)

### Critères d'apprentissage

$$R_{mse}(w) = \frac{1}{2} \sum_{k=1}^N (d^k - y^k)^2$$

$$R_{mse-pondéré}(w) = \sum_{k=1}^N \left( (d^k - y^k) \sum^{-1} (d^k - y^k) \right) \quad \text{avec} \quad \sum^{-1} = \frac{1}{N} \sum_{k=1}^N (d^k - y^k)(d^k - y^k)^T$$

$$R_{multiple-logistic}(w) = \sum_{k=1}^N \sum_{p=1}^n \left[ d_p^k \log \frac{d_p^k}{y_p^k} + (1 - d_p^k) \log \frac{1 - d_p^k}{1 - y_p^k} \right]$$

$$R_{log-likelihood}(w) = \sum_{k=1}^N d^k \log \frac{d^k}{p^k} \quad \text{avec} \quad p^k = \frac{e^{y^k}}{\sum_j e^{y^j}}$$



EPAT 14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

67



## Perceptron Multi-Couche (PMC)

### Règles d'adaptation

$$R(w) = \int_{\mathcal{Z}} L(z, w) dp(z) \quad \longrightarrow \quad \tilde{R}(w) = \frac{1}{N} \sum_{k=1}^N L(z^k, w)$$

$$w^+ = \underset{w}{\operatorname{argmin}} \tilde{R}(w)$$

$$z = (x, d) \quad \text{v.a. } (X, D) \quad p(x, d) = p(x)p(d/x)$$

$$\tilde{R}_{mlp}^k(w) = \frac{1}{2} \sum_{i=1}^n (y_i^k - d_i^k)^2$$

$$\tilde{R}_{mlp}^k(w) = \frac{1}{2} \sum_{m=1}^n \left( f \left( w_{m0}^{(2)} + \sum_{j \in \operatorname{Aumont}(m)} w_{mj}^{(2)} f \left( w_{j0}^{(1)} + \sum_{i \in \operatorname{Aumont}(j)} w_{ji}^{(1)} x_i \right) \right) - d_m^k \right)^2$$



EPAT 14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

68

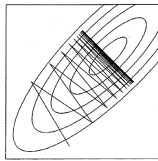




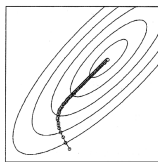
## Perceptron Multi-Couche (PMC)

### Règles d'adaptation

$$\tilde{R}_{mlp}^k(w) = \frac{1}{2} \sum_{m=1}^n \left( f \left( w_{m0}^{(2)} + \sum_{j \in \text{Amont}(m)} w_{mj}^{(2)} f \left( w_{j0}^{(1)} + \sum_{i \in \text{Amont}(j)} w_{ji}^{(1)} x_i \right) \right) - d_m^k \right)^2$$



$\epsilon$  grand

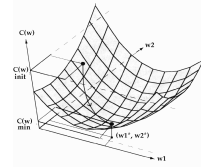


$\epsilon$  petit

$$w_{ji}(t+1) = w_{ji}(t) - \epsilon(t) \nabla_{w_{ji}} \tilde{R}_{mlp}^k(w)$$

$$\nabla_{w_{ji}} \tilde{R}_{mlp}^k(w) = \delta_j^k s_i$$

$$\begin{cases} \delta_j^k = f'(a_j)(y_j^k - d_j^k) & \text{si } j \in \text{Sortie} \\ \delta_j^k = f'(a_j) \sum_{h \in \text{Aval}(j)} w_{hj} \delta_h^k & \text{si } i \notin \text{Sortie} \end{cases}$$



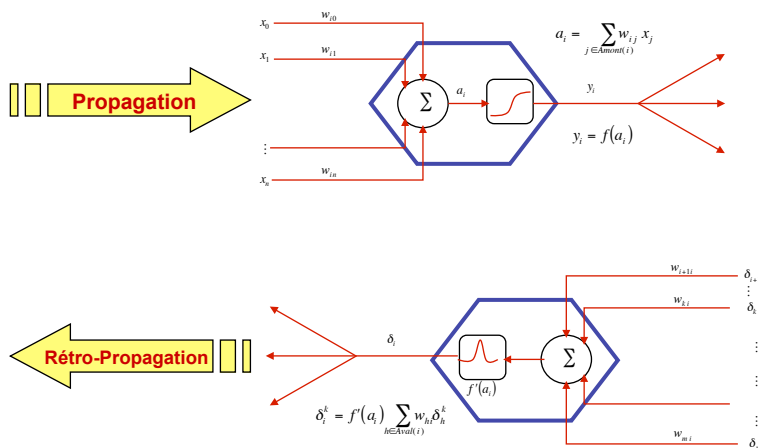
EPAT 14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

69



## Rétro-propagation du gradient



EPAT 14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

70





## Perceptron Multi-Couche

### Algorithme d'apprentissage

- 1- Tirer au hasard  $w_0$
  - 2- Présentation d'une forme  $(x^k, y^k)$
  - 3- Calcul de l'état du réseau par propagation
 
$$x_j = f(a_j) \quad a_j = w_{j0} + \sum_{i \in \text{Amont}(j)} w_{ji} x_i$$
  - 4- calcul des signaux d'erreur
 
$$\begin{cases} \delta_j^k = f'(a_j)(y_j^k - d_j^k) & \text{si } j \in \text{Sortie} \\ \delta_j^k = f'(a_j) \sum_{h \in \text{Aval}(j)} w_{hj} \delta_h^k & \text{si } i \notin \text{Sortie} \end{cases}$$
  - 5- Adaptation les poids
 
$$w_{ji}(t+1) = w_{ji}(t) - \varepsilon(t) \delta_j^k s_i$$
- Où  $\varepsilon(t)$  est le pas du gradient.
- 6- Répéter de 2 à 5 jusqu'à l'obtention d'une valeur acceptable de l'erreur



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

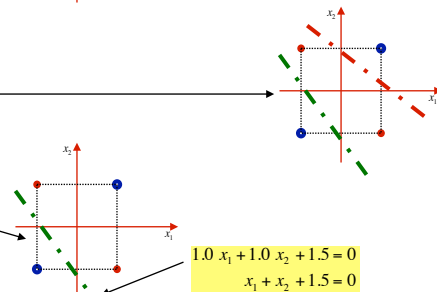
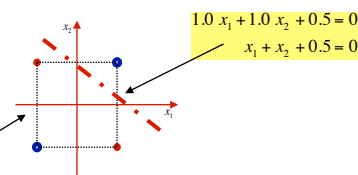
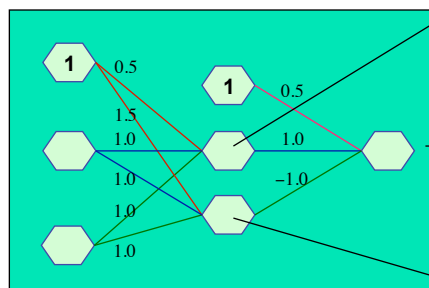
71



### Exemple de PMC (MLP)

$$D = \left\{ \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, -1; \begin{pmatrix} 1 \\ 1 \\ -1 \end{pmatrix}, 1; \begin{pmatrix} -1 \\ 1 \\ 1 \end{pmatrix}, 1; \begin{pmatrix} -1 \\ -1 \\ -1 \end{pmatrix}, -1 \right\}$$

$$\psi(x, w) = x_1 \oplus x_2$$



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

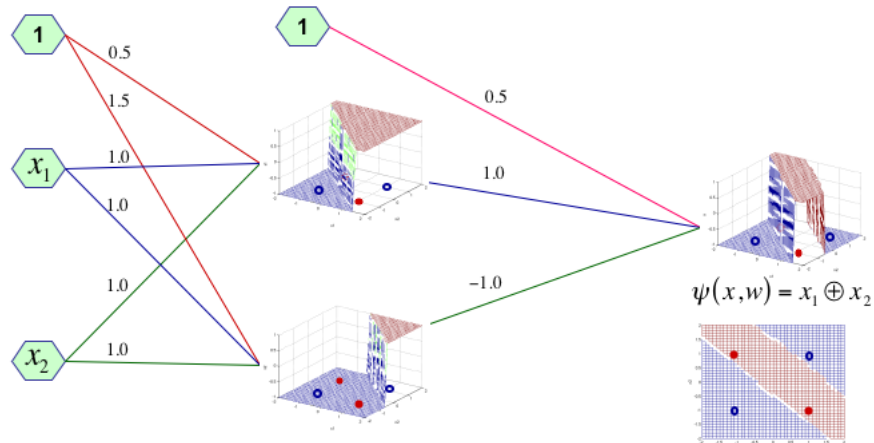
72







## Exemple de PMC (MLP)



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

73



## Approximation universelle

### Théorème [Hornik 1989\*] :

*Les PMC à trois couches avec des sorties linéaires, des cellules cachées saturantes croissantes "en nombre suffisant" peuvent approcher avec une précision arbitraire toute fonction bornée mesurable d'un espace de dimension finie dans un autre.*

Un PMC est un **approximateur universel**.

### Remarques :

– En théorie, il n'y aurait besoin d'aucune autre structure de réseau. Toutefois, dans les applications, il peut s'avérer plus pratique d'utiliser plusieurs couches, des sorties non linéaires, etc.

– C'est un théorème purement existentiel. Il ne dit pas comment déterminer un nombre approprié de neurones dans la couche cachée et les valeurs des poids pour approximer une fonction donnée avec une précision donnée !



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

74



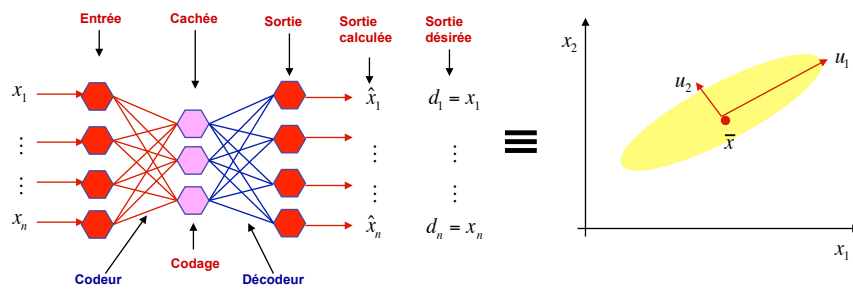


## Liens avec l'analyse factorielle

PMC & ACP \*\*\*

Un PMC auto-associatif avec des sorties linéaires, peut réaliser une transformation de Karhunen-Loève.

Un PMC auto-associatif est équivalent à une Analyse en Composantes Principales.



\* Bourlard H. & Kamp Y. (1988) : «Auto-association by multilayer perceptrons and singular value decomposition»  
Biological Cybernetics, Vol. 59, pp. 291-294.  
\*\* Baldi P. & Hornik K. (1989) : «Neural networks and principal component analysis: Learning from examples without local minima »  
Neural Networks, Vol. 2, N° 1, pp. 53-58.



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

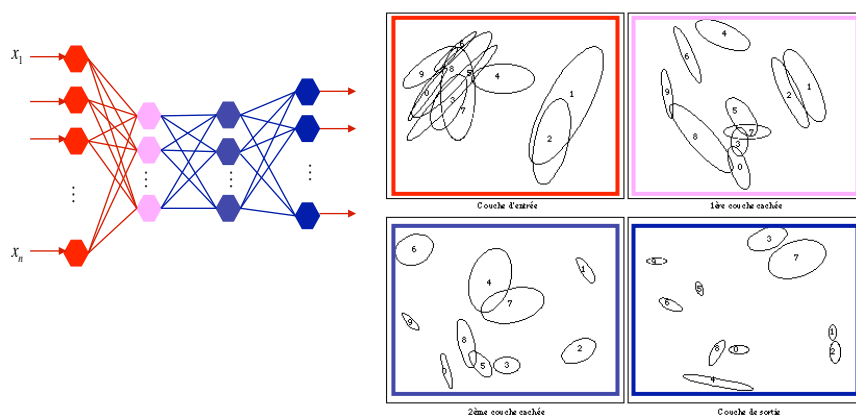
Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

75



## Liens avec l'analyse factorielle

PMC & AD \* : analyse et interprétation dans le cas non-linéaire



\* Bennani Y. (1992) : «Approches Connexionnistes pour la Modélisation et l'Identification»  
Thèse de Doctorat, LRI-Université Paris 11, Orsay.

\* Gallinari P., Thiria S., Badran F., Fogelman-Soulie F. (1991) : « On the relations between discriminant analysis and multilayer perceptrons »  
Neural Networks, Vol. 4, N° 3, pp. 349-360.



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

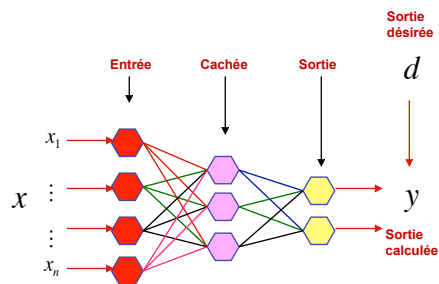
76





## Réseaux structurés

Connexions complètes



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

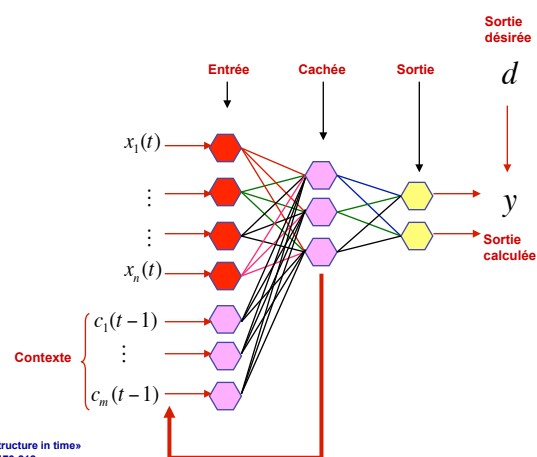
Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

77



## Réseaux structurés

Connexions complètes avec contexte [Elman\*]



\* Elman J.L. (1990) : «Finding structure in time»  
Cognitive Science, Vol. 14, pp. 179-212.



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

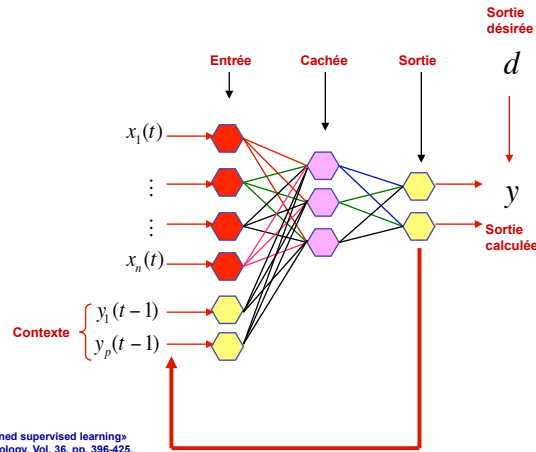
78





## Réseaux structurés

Connexions complètes avec contexte [Jordan\*]



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

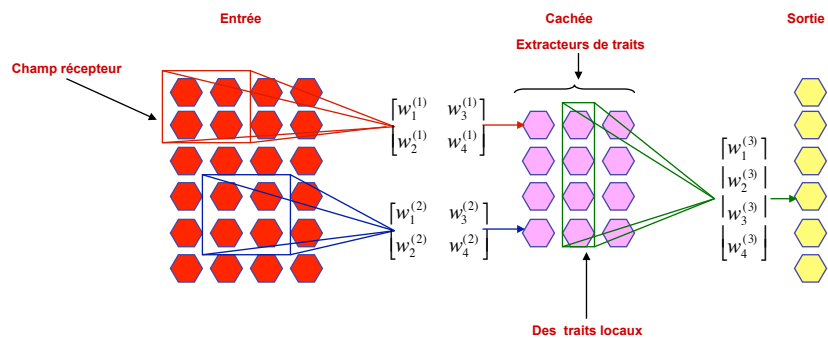
79



## Réseaux structurés

Connexions locales

L'utilisation de connexions locales diminue très fortement le nombre de poids d'un réseau.



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

80

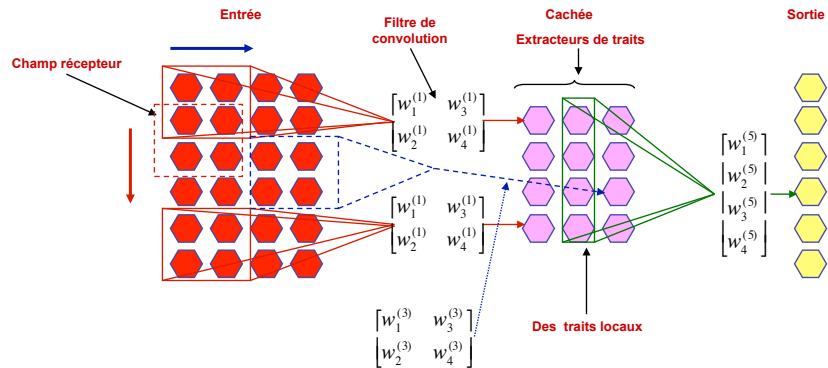




## Réseaux structurés

Connexions contraintes ou à poids partagés

Une propriété intéressante du mécanisme de partage des poids tient au nombre très faible de paramètres libres.



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

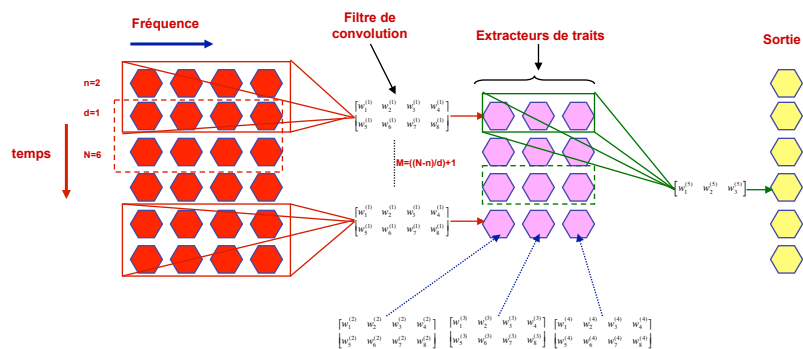
Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

81



## Réseaux structurés

TDNN (Time Delay Neural Network)



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

82

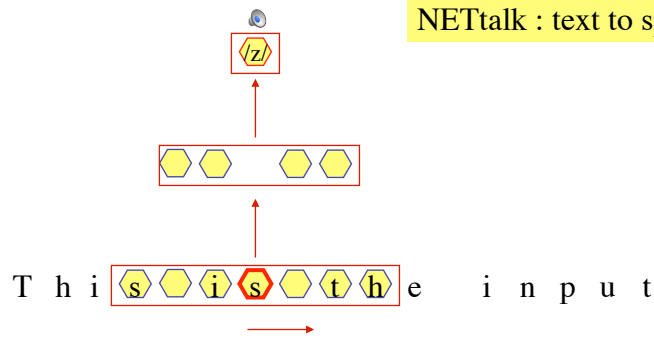




## Réseaux structurés

Synthèse de la parole [Sejnowski & Rosenberg\*]

### NETtalk : text to speech



\* Sejnowski T.J. & Rosenberg C.R. (1987) :  
«Parallel Networks that learn to pronounce English text»  
Complex systems, Vol 1, pp. 145-168.



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

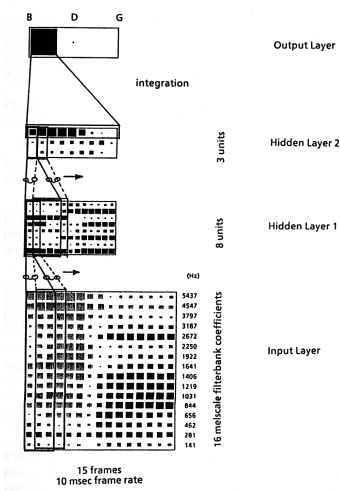
Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

83



## Réseaux structurés

Reconnaissance de la parole [Alex Waibel\*]



\* Waibel A., Hanazawa T., Hinton G., Shikano K., Lang K. (1987) :  
«Phoneme recognition using Time-Delay Neural Networks»  
Tech. Rep. ATR, TR-1-006.



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

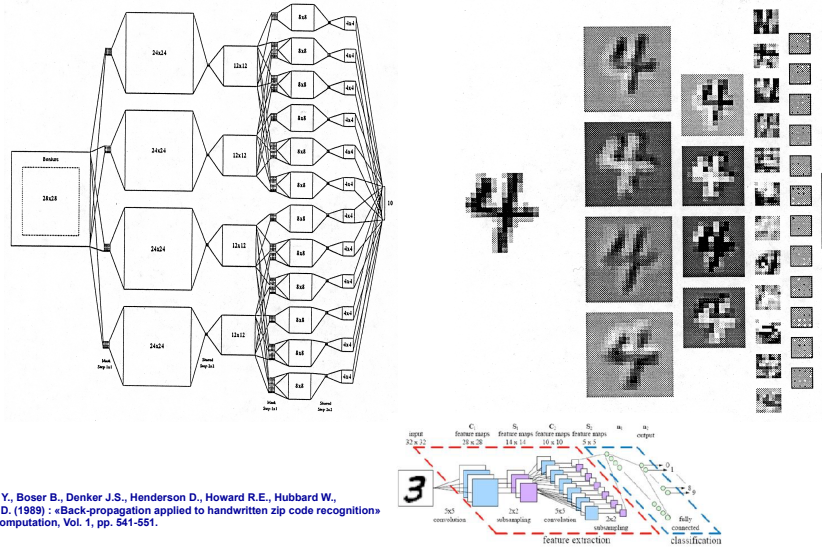
84





# Réseaux structurés : Réseaux à convolution

LeNet pour la reconnaissance de chiffres  
[Yann LeCun\*]



\* Le Cun Y., Boser B., Denker J.S., Henderson D., Howard R.E., Hubbard W., Jackel L.D. (1989) : «Back-propagation applied to handwritten zip code recognitions Neural Computation, Vol. 1, pp. 541-551.



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

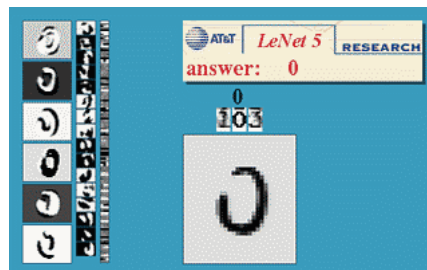
Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

85



# Réseaux structurés

LeNet pour la reconnaissance de chiffres  
[Yann LeCun\*]



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

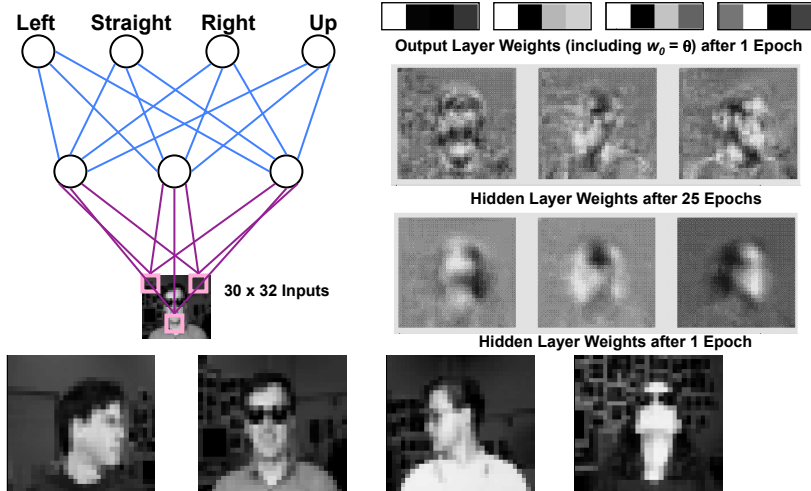
86





## Réseaux structurés

reconnaissance de la position du visage



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

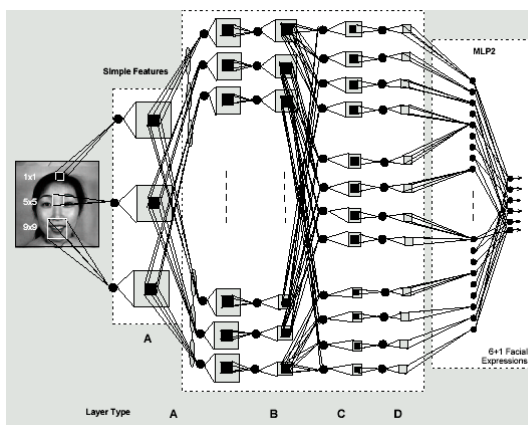
Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

87



## Réseaux structurés

reconnaissance d'expressions faciales



210 images (246x256 : TIFF)  
 10 femmes japonaises  
 6 expressions + 1 position neutre



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

88

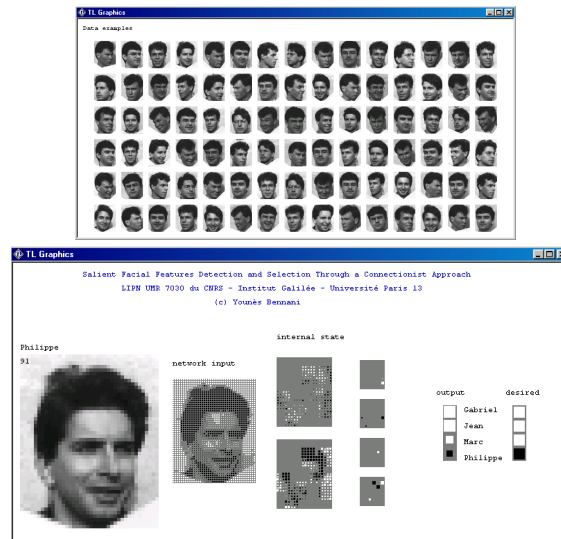






## Réseaux structurés

### identification de visages



## Systèmes d'Apprentissage Non Supervisé



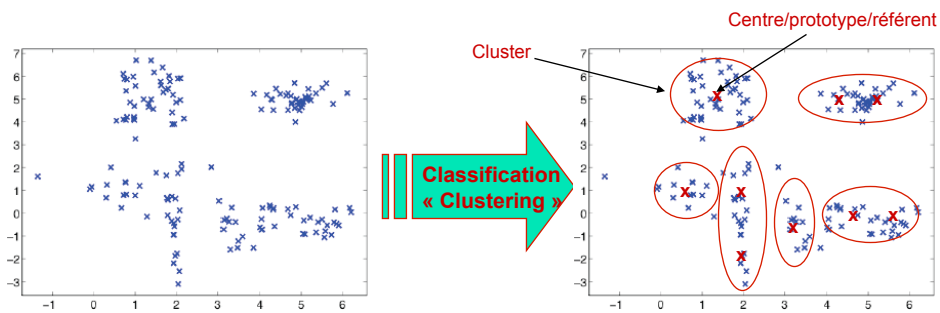


## Classification non supervisée « Clustering »

**Objectif :**

Définir sur un ensemble d'objets 2 à 2 comparables (ou une matrice de (dis)similarités), une partition (en groupes disjoints et complémentaires) qui respecte au mieux les ressemblances entre objets.

La ressemblance de 2 objets est grande lorsqu'ils figurent dans le même groupe et petite dans le cas contraire.



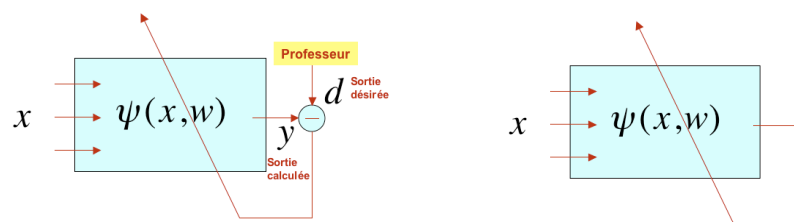
EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

91



## Supervisée vs Non-Supervisé



<b>Coût</b>	:	<b>Supervisé</b>	<b>Erreur de classification</b>	<b>Non supervisé</b>	<b>Inertie, ...</b>
<b>Données</b>	:	<b>Supervisé</b>	<b>Étiquetées</b> $X_w = \{(x^1, d^1), (x^2, d^2), \dots, (x^n, d^n)\}$	<b>Non supervisé</b>	<b>Non étiquetées</b> $X_w = \{x^1, x^2, \dots, x^n\}$
<b>Objectif</b>	:	<b>Supervisé</b>	<b>Prédiction</b>	<b>Non supervisé</b>	<b>Exploration</b>
<b># Classes</b>	:	<b>Supervisé</b>	<b>Connu</b>	<b>Non supervisé</b>	<b>Inconnu</b>
<b>Généralisation</b>	:	<b>Supervisé</b>	<b>Perf. N<sup>l</sup>es données</b>	<b>Non supervisé</b>	<b>Qualité sur les mêmes données</b>
<b>Théorie</b>	:	<b>Supervisé</b>	<b>Mature</b>	<b>Non supervisé</b>	<b>Jeune, très active</b>



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

92





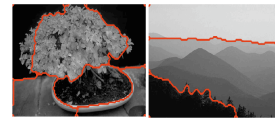
## Classification Non-Supervisée

### Domaines d'applications

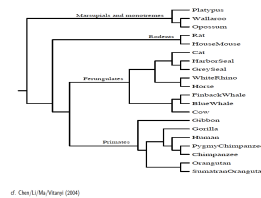
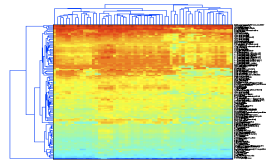
Outil très utilisé en **analyse exploratoire de données** dans de **nombreux domaines d'applications** :

- **Sciences sociales**
- **Biologie**
- **Médecine**
- **Astronomie**
- **Diagnostic**
- **Images**
- **Marketing**
- **Web mining**
- **Texte mining**
- **Graphe mining**
- **Data mining**
- ...

Corporate email communication (Adamic and Adar, 2005)



(from Zehbi-Masri, Perona, 2005)



d. Chen, Li, Wu, Wang (2006)



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

93



## SOM : Self Organizing Map

### Cartes auto-organisatrices/Cartes topologiques/Cartes de Kohonen

Développé par **Kohonen**, à partir des travaux de **Willshaw & Von Der Malsburg**

Cherche à transformer des signaux de **dimension quelconque**, en signaux à **une ou deux dimensions**.

Projeter les **données initiales** sur un espace discret et régulier de **faible dimension**.

Les espaces utilisés sont des **treillis réguliers** dont chacun des nœuds est occupé par un automate + **notion de voisinage** entre automates.

On a, dans ces cartes, la propriété suivante :

**des données similaires auront des projections proches sur la carte.**



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

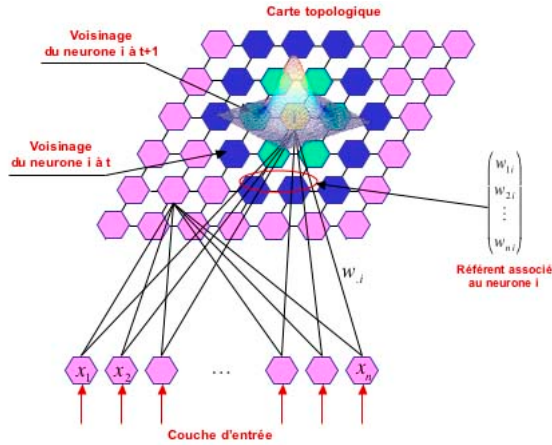
94





# SOM : Self Organizing Map

## Architecture du réseau



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

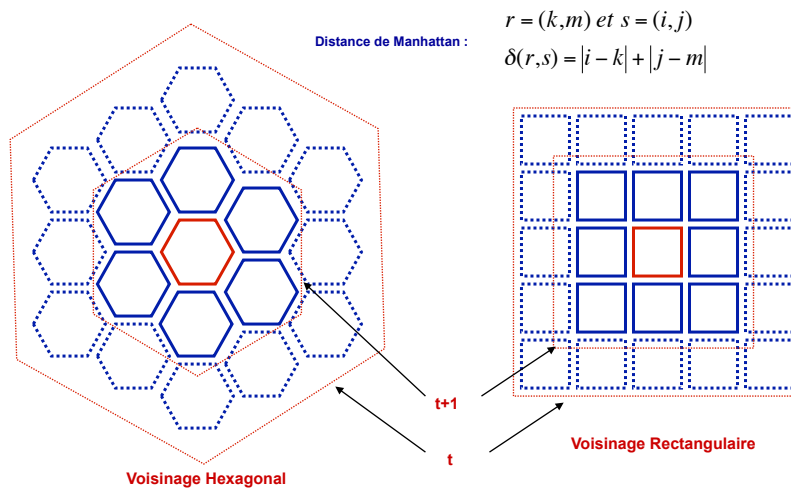
Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

95



# SOM : Self Organizing Map

## Types de voisinage



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

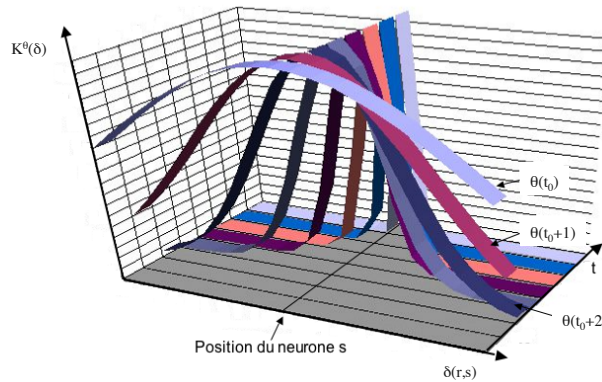
96





## SOM : Self Organizing Map

### Fonction de voisinage



$$K^\theta(\delta) = \frac{1}{\theta(t)} e^{-\frac{\delta^2}{\theta^2(t)}}$$

$$\theta(t) = \theta_i \left( \frac{\theta_f}{\theta_i} \right)^{\frac{t}{N_{iter}-1}}$$

$$r = (k, m) \text{ et } s = (i, j)$$

$$\delta(r, s) = |i - k| + |j - m|$$



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

97



## SOM : Self Organizing Map

### Formulation

#### Minimiser la fonction de coût :

$$R_{SOM}(w, \chi) = \sum_{i=1}^N \sum_{c=1}^K K^\theta(\delta(c, \chi(\bar{x}_i))) \|\bar{x}_i - \bar{w}_c\|^2$$

#### Calcul des gradients :

$$\nabla_{\bar{w}}(R_{SOM}) = \frac{\partial R_{SOM}}{\partial \bar{w}}$$

#### Règle d'adaptation :

$$\bar{w}_j^t = \bar{w}_j^{t-1} - \epsilon^t \nabla_{\bar{w}}(R_{SOM})$$

$$\bar{w}_j^t = \bar{w}_j^{t-1} - \epsilon^t K^\theta(\delta(j, \chi(\bar{x}_i))) (\bar{w}_j^{t-1} - \bar{x}_i)$$

#### Fonction d'affectation :

$$\chi(\bar{x}_i) = \arg \min_{1 \leq c \leq K} \|\bar{x}_i - \bar{w}_c\|^2$$

#### Fonction de voisinage :

$$K^\theta(\delta) = \frac{1}{\theta(t)} e^{-\frac{\delta^2}{\theta^2(t)}}$$

#### Distance entre unités :

$$r = (k, m) \text{ et } s = (i, j)$$

$$\delta(r, s) = |i - k| + |j - m|$$

#### Température :

$$\theta(t) = \theta_i \left( \frac{\theta_f}{\theta_i} \right)^{\frac{t}{N_{iter}-1}}$$



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

98





## SOM : Self Organizing Map

### Algorithme d'apprentissage

Phase  
d'initialisation

#### 1. Initialisation :

$$t = 0$$

$$w^0 = \{\bar{w}_1^0, \bar{w}_2^0, \dots, \bar{w}_K^0\} \quad \bar{w}_i \in \mathfrak{R}^n$$

Phase  
de compétition

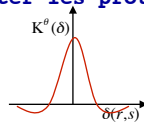
#### 2. Présenter un exemple $\bar{x}$ d'apprentissage

#### 3. Déterminer le gagnant :

$$\chi(\bar{x}_i) = \arg \min_{1 \leq c \leq K} \|\bar{x}_i - \bar{w}_c\|^2$$

Phase  
d'adaptation

#### 4. Adapter les prototypes :



$$\bar{w}_j^t = \bar{w}_j^{t-1} - \varepsilon^t K^{\theta}(\delta(j, \chi(\bar{x}_i))) (\bar{w}_j^{t-1} - \bar{x}_i)$$

$$K^{\theta}(\delta) = \frac{1}{\theta^t} e^{-\frac{\delta^2}{\theta^2(t)}} \quad \theta(t) = \theta \left( \frac{\theta}{\theta} \right)^{\frac{t}{N_{iter}-1}} \quad \varepsilon(t) = \varepsilon_i \left( \frac{\varepsilon_i}{\varepsilon_j} \right)^{\frac{t}{N_{iter}-1}}$$

#### 5. Incrémenter le nombre d'itération :

$$t = t + 1$$

#### 6. Test d'arrêt : si $t < N_{iter}$ aller en 2.

$$\theta_j = 2, \quad \theta_j = 0.5$$

$$\varepsilon_i = 0.5, \quad \varepsilon_j = 0.005$$

$$N_{iter} = 10000$$

$O(kN)$ , avec  $N$  le nombre d'objets,  
 $t$  le nombre d'itérations et en général  $t$  et  $k \ll N$



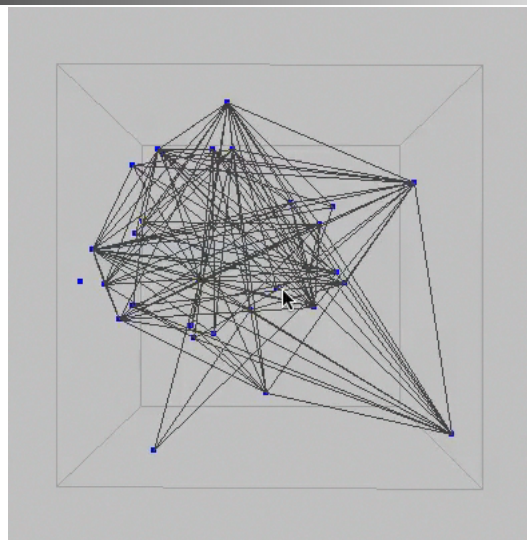
EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

99



## SOM : Self Organizing Map



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

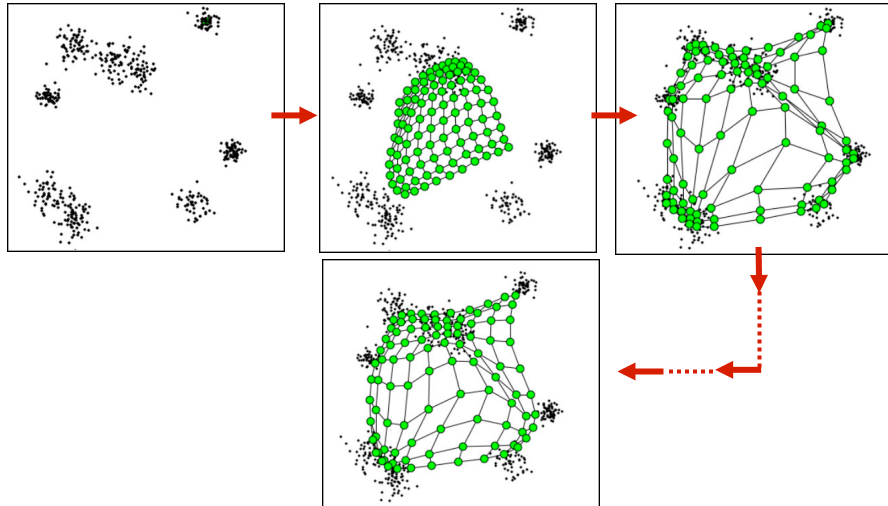
100





# SOM : Self Organizing Map

Phase d'auto-organisation



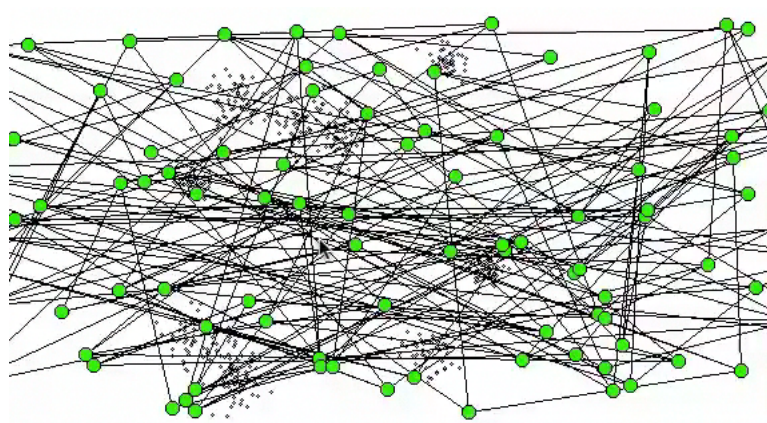
EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

101



# SOM : Self Organizing Map



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

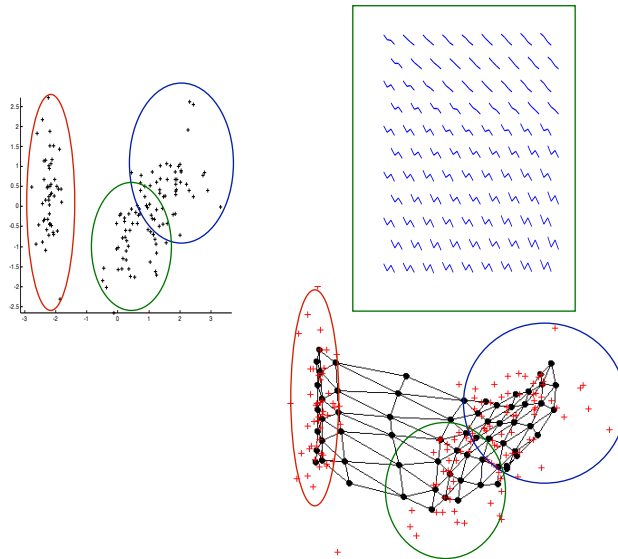
102





# Apprentissage Non-Supervisé

SOM : Self Organizing Map : Iris



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

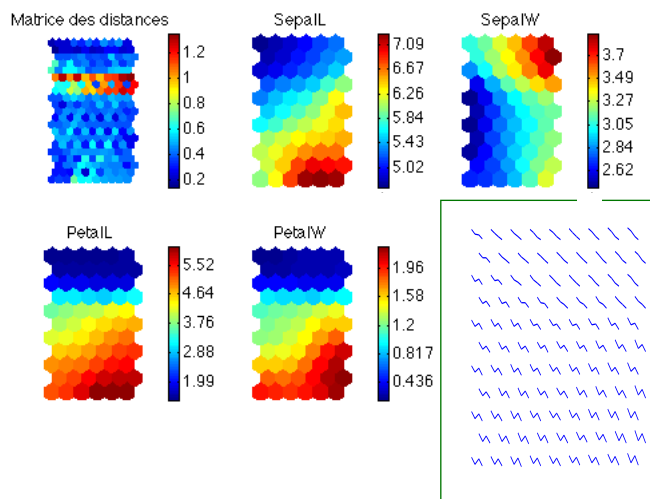
Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

103



# Apprentissage Non-Supervisé

SOM : Self Organizing Map : Iris



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

104







## Systèmes d'Apprentissage profond



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

105



## Deep Learning (Travaux pionniers)

- Fukushima (1980) – Neo-Cognitron
- LeCun (1998) – Convolutional Neural Networks
- Many layered MLP with backpropagation
  - Tried early but without much success
    - Lent
    - Diffusion du gradient
- Présentation du deep networks avec apprentissage non supervisé



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

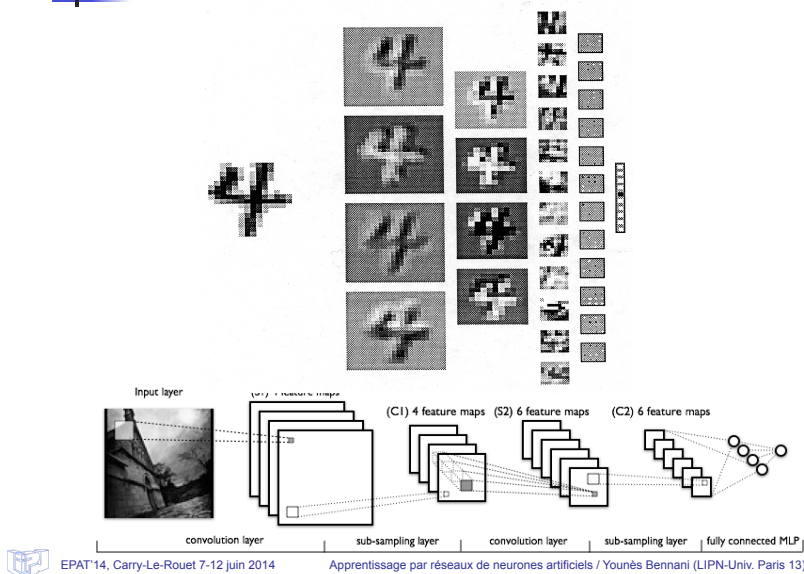
Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

106





## Convolutional Neural Networks



## Apprentissage des Deep Networks

- Construire un espace de représentation (feature space)
  - Notez que c'est ce que nous faisons avec les noyaux SVM, ou les couches cachées dans MLP, etc, mais maintenant, nous allons construire l'espace de représentation en utilisant les architectures profondes.
  - Apprentissage non supervisé entre les couches peut décomposer le problème en sous-problèmes distribués (avec des niveaux d'abstraction plus élevés) à être encore décomposé à des couches successives





## Le défi d'entraîner des réseaux de neurones profonds

- Difficultés d'apprentissage supervisé des réseaux profonds
  - Les premières couches du MLP ne sont pas bien apprises
    - Diffusion du Gradient - erreur s'atténue à mesure qu'elle se propage aux couches précédentes : le gradient se propage « mal » de la sortie vers l'entrée.
    - Conduit à un apprentissage très lent.
    - Plus le réseau est profond plus le degré de non-linéarité du réseau augmente, ce qui augmenterait les chances de trouver ces obstacles à l'optimisation.
    - Les couches inférieures restent avec des transformations pas très utiles de l'entrée.
    - Besoin d'un moyen pour aider les premières couches à faire un travail efficace
  - Souvent pas suffisamment de données étiquetées disponibles
    - Pouvons-nous utiliser des approches non supervisées / semi-supervisées pour profiter des données non étiquetées
  - Réseaux profonds ont tendance à avoir des problèmes de minima locaux plus que les réseaux peu profonds pendant l'apprentissage supervisé



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

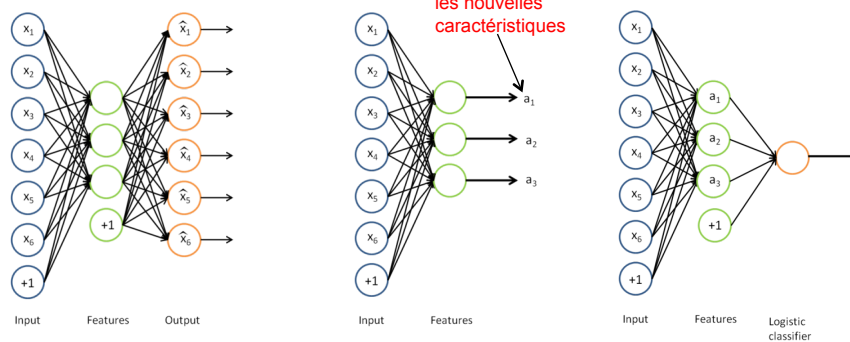
Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

109



## Auto-Encoders

- Un type d'apprentissage non supervisé qui tente de découvrir les caractéristiques génériques des données
  - Apprendre la fonction identité par apprentissage des sous-caractéristiques importantes des données
  - Compression, etc
  - Peut utiliser seulement les nouvelles caractéristiques dans la nouvelle série de d'apprentissage ou concaténer les deux



EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

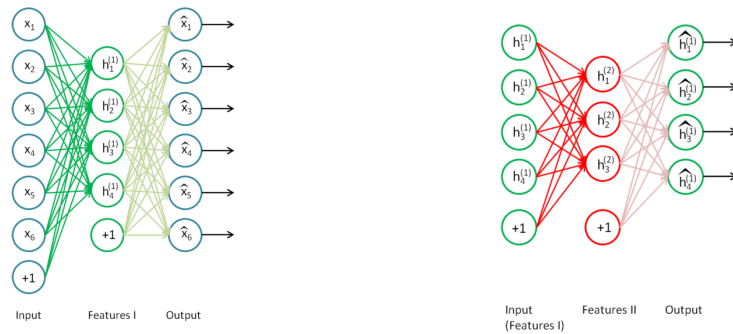
110





## Stacked Auto-Encoders

- Bengio (2007)  
 Empilez nombreux auto-encodeurs en succession  
 Déposer la couche de sortie de décodage à chaque fois



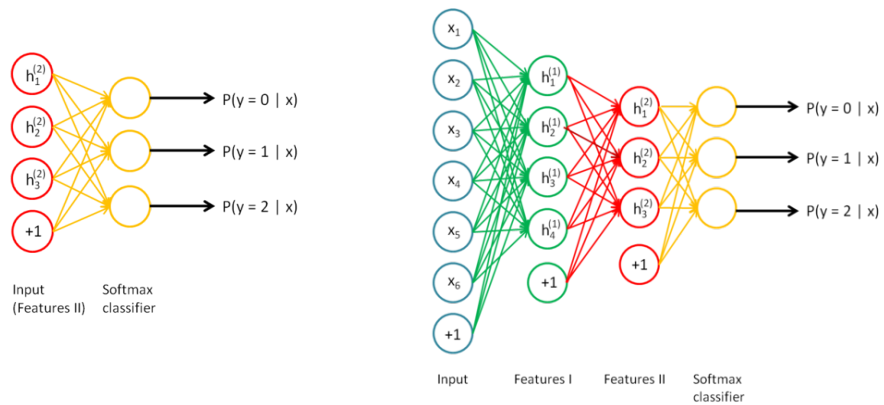
EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

111



## Stacked Auto-Encoders



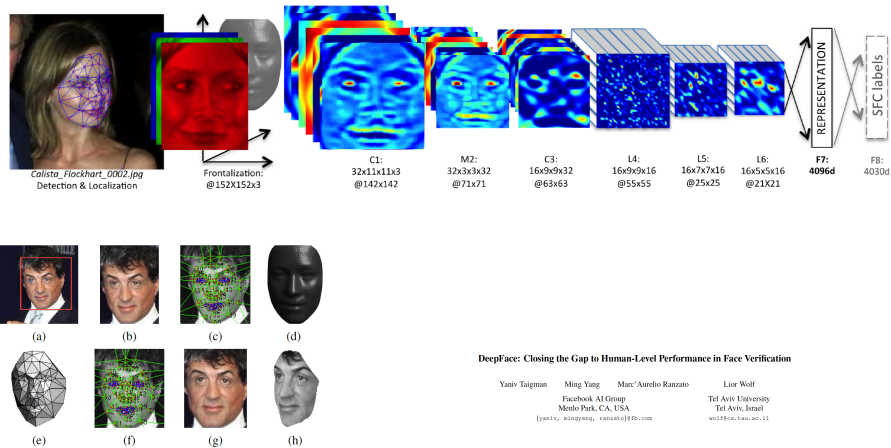
EPAT\*14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

112



## Un exemple : Deep Face



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

113



## Deep Learning

### References

- Bengio, Y. (2008). Learning deep architectures for AI. *Foundations and Trends in Machine Learning*, to appear.
- Bengio, Y., Delalleau, O., & Le Roux, N. (2006). The curse of highly variable functions for local kernel machines. In Weiss, Y., Schölkopf, B., & Platt, J. (Eds.), *Advances in Neural Information Processing Systems 18*, pp. 107–114. MIT Press, Cambridge, MA.
- Bengio, Y., Lamblin, P., Popovici, D., & Larochelle, H. (2007). Greedy layer-wise training of deep networks. In Schölkopf, B., Platt, J., & Hoffman, T. (Eds.), *Advances in Neural Information Processing Systems 19*, pp. 153–160. MIT Press.
- Bengio, Y., & LeCun, Y. (2007). Scaling learning algorithms towards AI. In Bottou, L., Chapelle, O., DeCoste, D., & Weston, J. (Eds.), *Large Scale Kernel Machines*. MIT Press.
- Hinton, G. E., Osindero, S., & Teh, Y.-W. (2006). A fast learning algorithm for deep belief nets. *Neural Computation*, 18, 1527–1554.
- Ranzato, M., Poultney, C., Chopra, S., & LeCun, Y. (2007). Efficient learning of sparse representations with an energy-based model. In Schölkopf, B., Platt, J., & Hoffman, T. (Eds.), *Advances in Neural Information Processing Systems 19*. MIT Press.



EPAT'14, Carry-Le-Rouet 7-12 juin 2014

Apprentissage par réseaux de neurones artificiels / Younés Bennani (LIPN-Univ. Paris 13)

114



## Conclusion

- Aujourd'hui, on observe un retour en force des réseaux de neurones artificiels avec des résultats impressionnants.
- Les réseaux de neurones artificiels (RN) forment un ensemble de techniques
  - matures
  - efficaces
  - avec une base théorique solide
  - largement utilisées dans de nombreux domaines
    - Médecine (diagnostic, prothèses, conseils d'urgence)
    - Prospection minière et pétrolière
    - Reconnaissance vocale, écriture manuscrite
    - Télécommunication (compression des données)
    - Finance (estimation immobilier, détection fausses déclarations, prédiction des cours)
    - Industrie (mesure, prédiction contrôle)
    - Transports (pilotage automatique, détection de risques, détection de pannes)



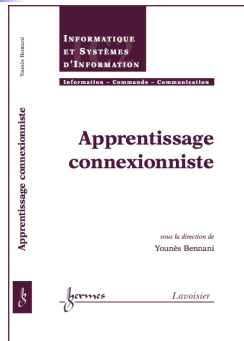
## De quoi n'a-t-on pas parlé ?

- Réseaux de neurones récurrents
- Autres modèles supervisés et non supervisés (LVQ, RBF, ART, ...)
- Architectures hybrides, modulaires
- Analyse du comportement (Dilemme Biais-Variance, VC-dim, ...)
- Problème de la généralisation
- Ajustement de la capacité de généralisation
  - Régularisation Formelle
  - Régularisation Structurelle
- ...





## Bibliographie



• **Apprentissage Connexionniste**  
Younès Bennani  
Editions Hermès Science (2006)  
ISBN: 2-7462-1337-0. (100 €)

• **Réseaux de neurones**  
G. Drefus, JM Martinez, M. Samuelides, MB Gordon, F.  
Badran, S. Thiria, L. Héroult  
Editions Eyrolles(2002) ISBN: 2-212-11019-7. (50 €)

• **Neural Networks for Pattern Recognition**

Christopher M. Bishop  
Clarendon Press - Oxford (1995)

• **Neural Smithing**

Supervised Learning in Feedforward Artificial Neural Networks  
Russell D. Reed & Robert J. Marks  
Massachusetts Institute of Technology Press (1999)

• **Pattern Recognition and Neural Networks**

B.D. Ripley  
Cambridge University Press (1996)

• **Neural Networks**

James A. Freeman & David M. Skapura  
Addison-Wesley Publishing Compagny (1991)

• **Adaptive Pattern Recognition and Neural Networks**

Yoh-Han Pao  
Addison-Wesley Publishing Compagny (1989)

• **Neural Networks in Computer Intelligence**

LiMin Fu  
Massachusetts Institute of Technology Press (1994)

