

1. les différents types de réseaux de neurones

1.1 la dynamique neuronale

Les réseaux de neurones artificiels (ANN : artificial neural networks) constituent une représentation ^{stylisée} d'un système nerveux réel contenant un ensemble de cellules communiquant entre elles.

Le biologiste Hebb fut le premier à postuler que la connectivité dans le cerveau est continuellement en évolution au cours de la période d'apprentissage.

Des systèmes artificiels implémentent la règle de Hebb dans leurs algorithmes de modification des connexions intercellulaires pour tâches de classification et de traitement.

On peut estimer que le cerveau humain contient environ 10^9 neurones et 10^{14} synapses (synthèse de connexions inter-neuronales).
Sur chaque neurone, on peut estimer à plus de 1000 le nombre de synapses.

Le rapport temporel d'un potentiel d'action produit par un neurone est de l'ordre de quelques millisecondes, c'est-à-dire un temps 10^6 fois plus grand que le temps caractéristique de modification de bits des ordinateurs. La connectivité neuronale est cependant de milliards de fois plus grande que dans les ordinateurs.

Un neurone artificiel va collecter les signaux et les traiter. Si le résultat obtenu dépasse un certain seuil, il va envoyer un signal aux autres neurones du réseau.

5

Pour faire les calculs, on appelle généralement champ local la somme des valeurs des signaux x_j arrivant sur un neurone donné pondérée par les valeurs de efficacités des connexions synaptiques. Par exemple, pour un neurone i , dont les efficacités sont $w_{i1}, w_{i2}, \dots, w_{in}$



$$G_i = \sum_j w_{ij} x_j$$

Souvent noté net_i en
"Champ local net"

La valeur de sortie du neurone est obtenue à partir de cette somme par action de la fonction d'activation f

$$x_i = f(G_i)$$

f peut prendre plusieurs formes

— N —

1.2 les architectures : feedforward, récurrent

Concernant leurs structures, les réseaux neuronaux sont dans 2 types : réseaux feedforward (propagation vers l'avant) et réseaux récurrents (comportant des boucles de rétrocontrôle)*

Réseau Feedforward : les neurones sont groupés en couches.

Le signal se propage de la couche d'entrée vers la couche de sortie par des connexions unidirectionnelles, les neurones étant groupés entre couches mais pas à l'intérieur d'une couche. Les exemples sont

- le perceptron multicouche (MLP : multi-layer Perceptron Rumelhart et McClelland 1986)
- la quantification vectorielle (LVQ : linear Vector quantization Kohonen 1989)
- le réseau de groupement de données (GMDH : group method of data handling Hecht Nielsen 1990)

* De manière schématisée, on a la structure suivante



Les réseaux généralement réalisent des transformations statiques entre un espace d'entrée et un espace de sortie : la sortie à un instant est une fonction seulement de l'entrée à cet instant.

Réseaux récurrents : Dans ce cas, la sortie de certains neurones sont renvoyées vers ces mêmes neurones ou vers des neurones de couches inférieures - les signaux évoluent vers l'avant et vers l'arrière.

Les exemples sont :

- Réseau de Koffield (1982)
- le réseau de Elman (1990), Jordan (1986).

La mémoire de ce système peut être considérée comme dynamique : la sortie à un instant donné peuvent reflète la valeurs d'entrée à cet instant ou bien les valeurs à des instants antérieurs.

— N —

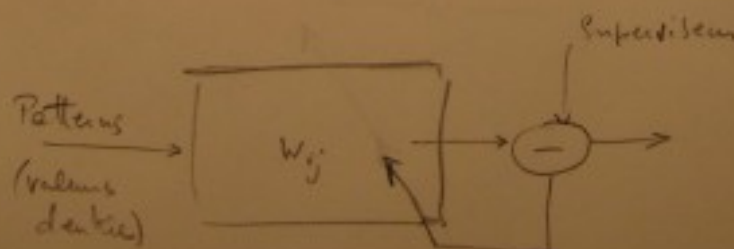
1.3 Les algorithmes d'apprentissage

Les réseaux sont entraînés par 2 types d'algorithmes : supervisé et non supervisé. Il existe un 3^e type, l'apprentissage par renforcement, qui peut être vu comme une forme spéciale d'apprentissage supervisé.

- Apprentissage supervisé. Le algorithme d'apprentissage supervisé ajuste le poids ou efficacités des connexions inter-neuronales suivant la différence entre des valeurs désirées pour les valeurs de sortie de réseau et les valeurs effectivement obtenues par le réseau, ceci pour une donnée des valeurs d'entrée.

L'apprentissage supervisé nécessite un superviseur pour indiquer les valeurs désirées en sortie.

Les exemples sont : la règle delta (Widrow et Hoff 1960)
l'algorithme de rétropropagation du gradient ou règle delta généralisée (Rumelhart, McClelland 1986)



Apprentissage non supervisé

Il ne s'occupe pas de la connaissance de valeurs de sortie. Durant l'apprentissage, seuls les patrons sont présentés au réseau qui adapte automatiquement les poids pour réaliser une tâche, celle de catégoriser les patrons en groupes ayant des caractéristiques.

Les exemples d'algorithmes sont ceux de :

- Kohonen (1984)
- Carpenter et Grossberg : ART Adaptive Resonance Theory (1988)

Ce sont des règles d'apprentissage dites compétitives.

Apprentissage par renforcement

Au lieu d'utiliser un superviseur, l'algorithme emploie une critère de sélection pour évaluer la qualité de la sortie.

L'exemple typique est

- l'algorithme génétique Holland 1975
Goldberg 1989.

— N —

1.4 le réseau de Hopfield :

1.4.1 Structure et fonction d'énergie



le réseau de Hopfield est composé de N neurones binaires
 c-à-d qui prennent les valeurs $x_i \in \{0, 1\}$ $i=1, 2, \dots, N$
 Affaire $G_i(x) = \sum_{j=1}^N w_{ij} x_j$, la somme pondérée par les
 poids des valeurs x_j de $X = (x_1, x_2, \dots, x_N)$. C'est le "champ local".
 le neurone i , à l'instant $t+1$, va changer son état x_i
 de la façon suivante

$$(H.1) \quad \begin{cases} x_i(t+1) = 1 & \text{si } G_i(x(t)) > \theta_i \\ x_i(t+1) = 0 & \text{si } G_i(x(t)) < \theta_i \\ x_i(t+1) = x_i(t) & \text{si } G_i(x(t)) = \theta_i \end{cases} \quad \begin{array}{l} \text{seuil} \\ \text{(seuil pour passer} \\ \text{ou la cellule,} \\ \text{à l'instant)} \end{array}$$

Il peut s'agir de plusieurs variantes.

Le mise à jour (updating) peut se faire de 2 façons, suivant le nombre de cellules considérées.

Dans le cas d'un schéma d'itération parallèle, toutes les cellules sont mises à jour à chaque cycle suivant la règle (A.1)

Dans le schéma d'itération asynchrone, seule une cellule est mise à jour à chaque pas d'itération. - Dans ce cas, la suite d'indices cellulaires est mise prédefinie, soit choisie aléatoirement. De plus, la mise à jour par bloc peut également être faite.

Soit une mise à jour asynchrone, avec la condition que toutes les cellules ont été mises à jour dans un temps fini à chaque valeur converger vers des points fixes (points dans l'espace de état $\text{cod} (X_1^0, X_2^0, \dots, X_N^0)$) quand les convergences ont certaines formes.

Fonction d'entropie

Supposons les convergences de forme asynchrone cod

$W_{ij} = W_{ji}$, le système admet une fonction de Lyapunov (ou fonction d'entropie) qui est une fonction bornée inférieurement et qui est non croissante lors de l'évolution du système $t \rightarrow t+1$

Cette fonction est

$$(H.2) \quad E(X) = -\frac{1}{2} \sum_{i,j=1}^N w_{ij} X_i X_j + \sum_{i=1}^N \theta_i X_i$$

(H.3) En fait, on a $E(X(t+1)) \leq E(X(t))$

De plus, comme $X(t+1)$ et $X(t)$ ne diffèrent que par un élément au plus, on a

$$E(X(t+1)) = E(X(t))$$

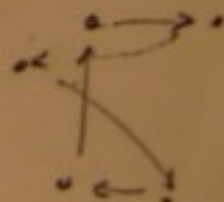
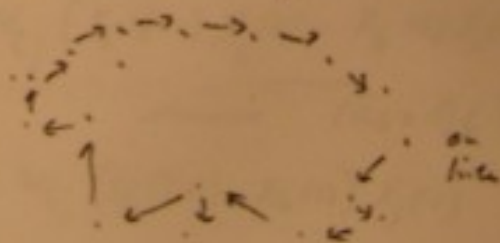
Ceci implique $X(t+1) = X(t)$. Or on déduit que dans le cas d'une dynamique asymptotique, les seuls attracteurs sont des points fixes. Les attracteurs sont le domaine de l'espace des états vers lesquels les états asymptotiquement les états.

Des attracteurs qui se seraient par des points fixes pourraient être des cycles, des attracteurs "étranges" (de structure chaotique)

$$X^1 = (x_1^1 \dots x_N^1)$$

$$X^2 = (x_1^2 \dots x_N^2)$$

cycle à 3 éléments



Par la (H.3) - l'expression pour la dérivée seconde a été
donnée pour la route à jour à l'instant t , et soit
 $\delta X_k(t)$ sa variation alors que $\delta E(t)$ denote la
variation de la fonction d'énergie. On a :

$$\begin{aligned} (H.4) \quad \delta E(t) &= E(t+1) - E(t) \\ &= -\frac{1}{2} \sum_{\substack{i,j \\ i \neq j}} W_{ij} [X_i(t+1)X_j(t+1) - X_i(t)X_j(t)] (\delta_{ik} + \delta_{jk}) \\ &\quad + \sum_i \theta_i (X_i(t+1) - X_i(t)) \delta_{ik} \end{aligned}$$

Pour le 1^{er} membre, on a

$$\sum_i \theta_i (X_i(t+1) - X_i(t))$$

$$\text{avec } X_i(t+1) - X_i(t) = \begin{cases} 0 & \text{si } i \neq k \\ \neq 0 & \text{si } i = k \end{cases}$$

$$\text{Donc on peut écrire } \sum_i \theta_i (X_i(t+1) - X_i(t)) \delta_{ik}$$

De la même façon

$$= -\frac{1}{2} \sum_{\substack{i,j \\ i \neq j}} W_{ij} [X_i(t+1)X_j(t+1) - X_i(t)X_j(t)]$$

soit $A = 0$ si $i \neq k$ et $j \neq k$

si $i \neq k$ et $j = k$

$$\begin{aligned} W_{ik} (X_i(t+1)X_k(t+1) - X_i(t)X_k(t)) &= W_{ik} (X_i(t)X_k(t+1) - X_i(t)X_k(t)) \\ &= W_{ik} (X_i(t) (X_k(t+1) - X_k(t))) \neq 0 \end{aligned}$$

Où nous pourrions aussi écrire

Après tout, si $i \neq k$ et $j \neq k$, ces deux termes sont
 en fait des termes en $i \neq j$ (du fait des deux termes
 $W_{ik} = 0$)

En fait, on a

$$(H.5) \quad \delta E(t) = -\delta X_k (G_k(X(t)) - \theta_k)$$

avec (rapel) $G_k(X(t)) = \sum_{\substack{j=1 \\ j \neq k}}^N W_{kj} X_j(t)$

le 2^e terme donne bien

$$(H.6) \quad \sum_i \theta_i (X_i(t+1) - X_i(t)) \delta_{ik} = \theta_k \delta X_k(t)$$

(il s'agit du 4^e).

Donc on peut

$$(H.7) \quad -\delta X_k G_k(X(t)) = (X_k(t+1) - X_k(t)) \sum_{\substack{j=1 \\ j \neq k}}^N W_{kj} X_j(t)$$

$$= - \sum_{\substack{j=1 \\ j \neq k}}^N W_{kj} (X_k(t+1) - X_k(t)) X_j(t)$$

Cela correspond bien au 2^e membre de (H.4). En fait

$$(H.4) \Rightarrow -\frac{1}{2} \sum_{\substack{j=1 \\ j \neq k}} W_{kj} [X_k(t+1) X_j(t) - X_k(t) X_j(t)]$$

(ici $X_k(t+1) X_j(t) - X_k(t) X_j(t)$ est pas min \rightarrow pour

(on suppose $(\delta_{ik} + \delta_{jk}) \rightarrow \delta_{ik}$)

$$(H.8) \quad = -\frac{1}{2} \sum_{\substack{j=1 \\ j \neq k}} W_{kj} (X_k(t+1) - X_k(t)) X_j(t)$$

$$2^{\text{ème}} \text{ terme } (\delta_i k + \delta_j k) \rightarrow \delta_j k$$

$$(44) \rightarrow -\frac{1}{2} \sum_{i \neq k} W_{ik} [X_i(t+1) X_k(t+1) - X_i(t) X_k(t)]$$

$$= -\frac{1}{2} \sum_{i \neq k} W_{ik} [X_k(t+1) - X_k(t)] X_i(t)$$

$$= -\frac{1}{2} \sum_{i \neq k} W_{ik} [X_k(t+1) - X_k(t)] X_i(t)$$

c'est bien (k.d)
en faisant la somme.

On a donc

$$(45) \quad \delta E(t) = -\delta X_k (G_k(x(t)) - \theta_k)$$

Si $\delta X_k(t) \neq 0$, $\delta X_k(t)$ a toujours le même
signe que $G_k(x(t)) - \theta_k$.

Exemple

$$i. \quad X_k(t) = 0 \rightarrow X_k(t+1) = 1 \text{ si } G_k(x(t)) - \theta_k > 0$$

$$\delta(X_k(t)) > 0 \quad \text{-----} > 0$$

$$ii. \quad X_k(t) = 1 \rightarrow X_k(t+1) = 0 \text{ si } G_k(x(t)) - \theta_k < 0$$

$$\delta(X_k(t)) < 0 \quad \text{-----} < 0$$

$$\text{Donc } \delta E(t) \leq 0$$

Le réseau de Hopfield est donc un système dynamique à plusieurs degrés de liberté, et est non linéaire. Sa dynamique lui permet de converger vers des points fixes qui sont des attracteurs.

Ainsi, une fois pour une large classe de conditions initiales, les états asymptotiques sont les points fixes. Cette insensibilité par rapport aux conditions initiales est utile à posit dans le but de reconnaître de patterns bruits qui sont convergus vers aux conditions initiales.



Il s'agit à présent d'identifier les points fixes de la dynamique avec des patterns que l'on désire stocker dans le système neuronal. Pour cela, Hopfield a proposé la prescription

suivante

1.4.2 Enregistrement et identification

Différents patterns sont stockés M patterns différents $A^1 \dots A^M$, chaque pattern étant un vecteur à N dimensions (N cellules en entrée)

le reste pour construire w_{ij} et alors le neurone

$$w_{ij} = (1 - \delta_{ij}) \sum_{m=1}^M (2 A_i^m - 1)(2 A_j^m - 1)$$

si $A_i^m = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$ - et a ainsi $w_{ii} = 0$ par $(1 - \delta_{ii}) = 0$

Comme $(2 A_i^m - 1) = \begin{cases} -1 & \text{si } A_i^m = 0 \\ +1 & \text{si } A_i^m = 1 \end{cases}$

on voit que w_{ij} est défini par le pattern A^m si A_i^m et A_j^m prennent la même valeur (0, ou 1)

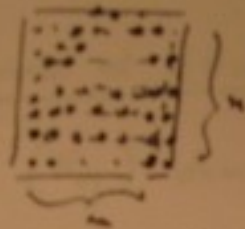
Soit $A_i^m = 0$ et $A_j^m = 0$ Soit $A_i^m = 1$ et $A_j^m = 1$

Cette concomitance des valeurs de patterns charge les connexions. Ainsi, si $A_i^m = 0$ et $A_j^m = 1$, la connexion w_{ij} est diminuée pour ce qui concerne

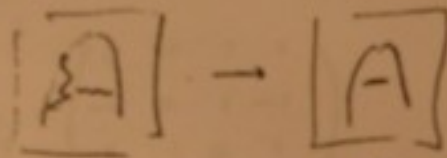
le pattern A^m - les patterns A^m sont des points clés pour la dynamique neuronale

l'algorithme est implémenté dans une procédure C qui est distribuée.

l'exemple traité portera sur des patterns ayant la forme d'images.



le nombre de cellules N en entrée (Partie) du réseau est
 $N = M \times m$ - les patterns ont certaines formes qui peuvent
 de lettres ou bien avoir d'autres structures.

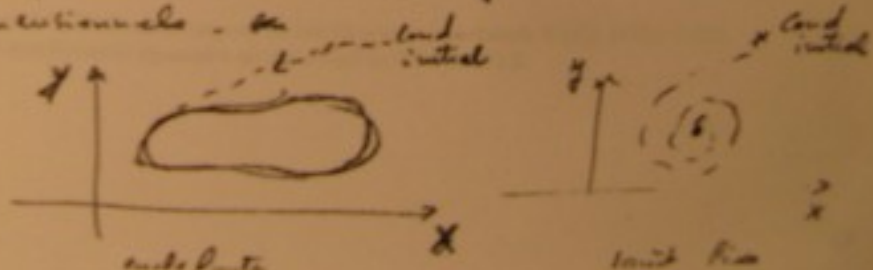


point fixe

L'existence d'une fonction d'énergie associée à l'existence
 d'un point fixe permet d'imaginer de présenter (comme
 conditions initiales de la dynamique) des patterns bruits
 introduits dans le bassin d'attraction d'un point fixe
 donné et d'observer un passage vers ce point fixe
 établissant ainsi une forme de reconnaissance.



Ceci peut être visualisé sur des systèmes plus simples
 bidimensionnels.



S'orties du réseau $y_i = z_i(\pi) \dots$

Trajectoires et Attracteurs.

(W_{ij}) donnée \rightarrow Trajectoires dans l'espace de tous les états

Exemple, 4 cellules, 2^4 états possibles

$(-1 -1 -1 -1) \quad (-1 -1 -1 1) \quad \dots \quad (1 1 1 1)$

0 1 2 3 4 5 6 7 8 9 A B C D E F

$$W_{ij} = \begin{pmatrix} 0 & -1 & 1 & 1 \\ -1 & 0 & 1 & -1 \\ 1 & 1 & 0 & 1 \\ 1 & -1 & 1 & 0 \end{pmatrix}$$

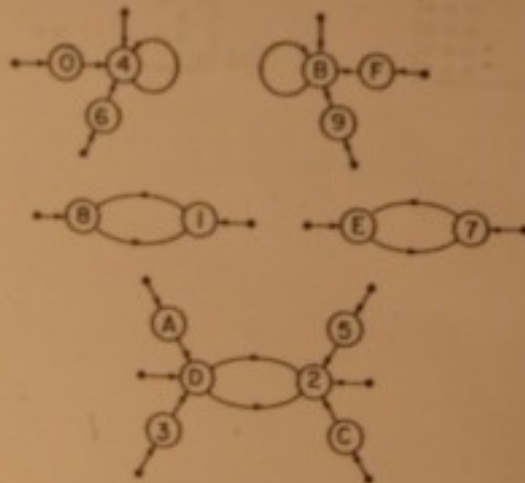


Figure 2.7: Trajectory map of 16 initial states in network 2 (J_{ij} of Eq. 2.25) with synchronous dynamics as described in Section 2.2.2.

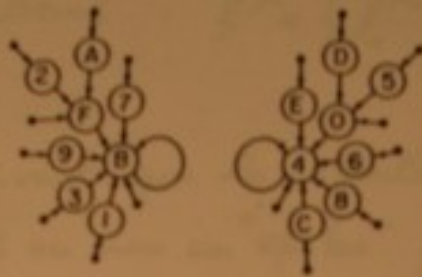


Figure 2.8: Trajectory map for network 2 with sequential, asynchronous dynamics. See e.g., Section 2.2.3.

variable d'ensemble : "Energie"

$$E(k) = -\frac{1}{2} \sum_{i \neq j} w_{ij} S_i(k) S_j(k)$$

$$E(k+1) \leq E(k)$$

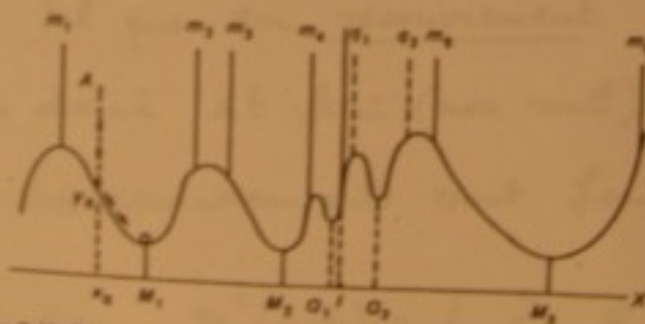
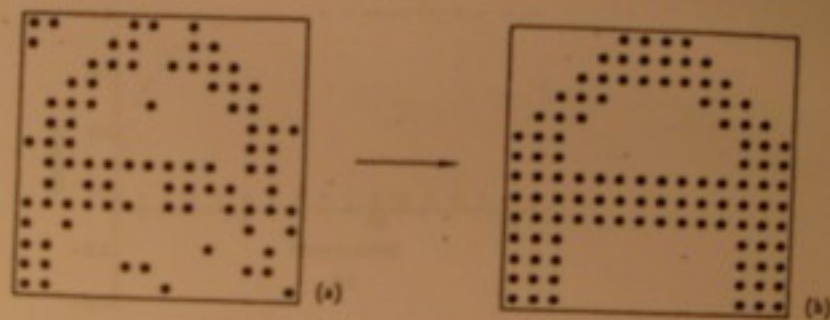


Figure 2.10: The one-dimensional landscape metaphor for association, random walkable memory. M_1, M_2 are memories, Q_1, Q_2 are spurious states, m_1, \dots, m_{10} are maxima delimiting basins of attraction.



Apprentissage

• Stocker p patterns : $\xi^\mu = \{ \xi_1^\mu \xi_2^\mu \dots \xi_N^\mu \}$

$\mu = 1, 2, \dots, p$

• $\xi_i^\mu = \pm 1$ suivant que la cellule d'input a émis ou non un "spike" pour le pattern μ

Prescription de Hopfield.

$$\Delta W_{ij} = \frac{1}{N} \sum_{\mu=1}^p \xi_i^\mu \xi_j^\mu \quad \sim \text{Règle de Hebb}$$

Reconnaissance

Expérience numérique.

Présentation au réseau d'un stimulus $\{S_i^{\text{in}}\}$ tel que son recouvrement avec un pattern donné est fini (non nul) alors que les autres recouvrements sont faibles

$$\underline{\text{Ex:}} \quad m^1(t=0) = \frac{1}{N} \sum_i \xi_i^1 S_i^{\text{in}} \neq 0$$

$$m^\mu(t=0) \approx 0 \quad \mu \neq 1$$

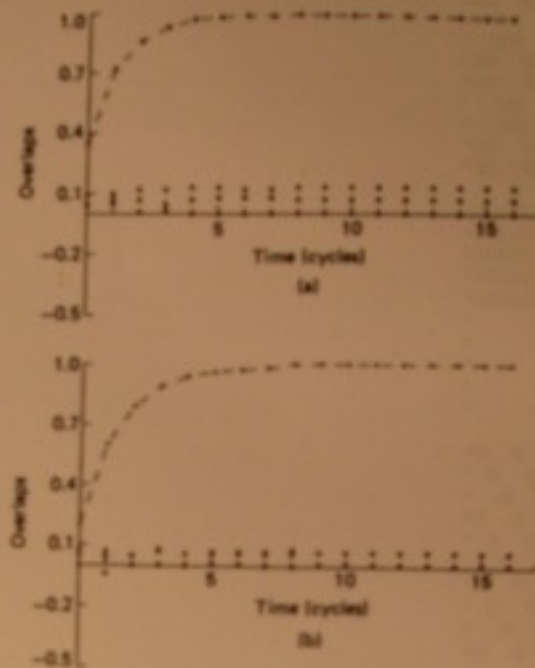


Figure 4.4: Time evolution of overlap of a network state with the stored memories for two initial configurations. (a) $m^1=0.35$; (b) $m^1=0.2$. Simulation with $N=400$, $p=7$, and noiseless asynchronous dynamics.

1.5 Systèmes autoorganisant

Von der Malsburg (73) Grossberg (76) Fukushima (77)
Kohonen (89)

Obermeyer (91) : Applications au système visuel.

Classifieurs de grande quantité de données.

Mécanisme de plasticité synaptique : Compétitif

Amplifier les réponses des cellules ayant initialement une forte activité

Kohonen : Cellules avec interactions latérales de portée finie

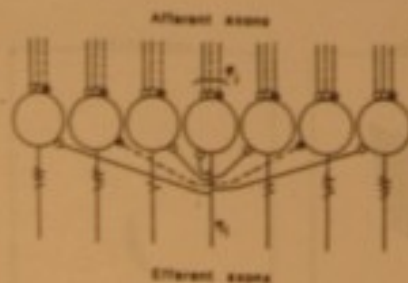


Fig. 5.1. Laterally interacting neurons

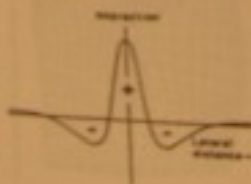


Fig. 5.2. The "Mexican-hat function" of lateral interaction

Conséquence : Amplification des réponses de certaines cellules dites "vainqueurs"

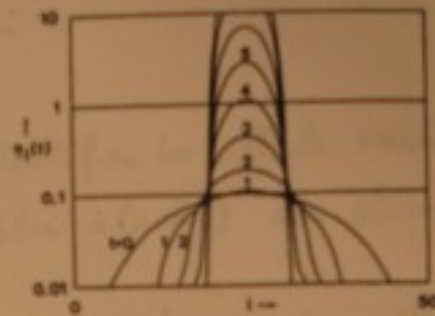
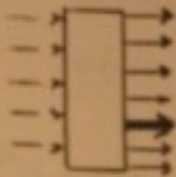


Fig. 5.4. Clustering of activity in a one-dimensional array



Affrentinage (non supervisé) pour l'autoorganisation

But : modifier les connexions de façon dynamique afin de clarifier les patrons (E_1, \dots, E_m)

~ Transformations somatosensorielles.

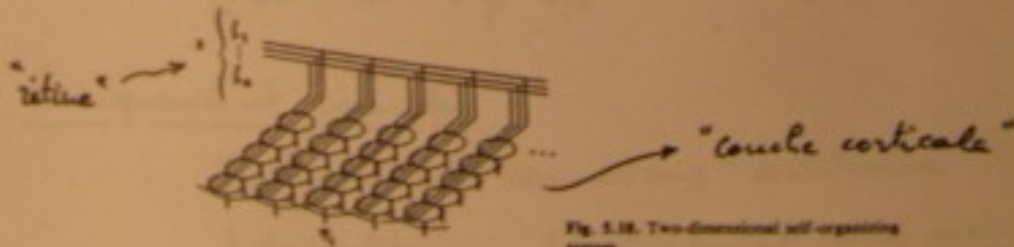
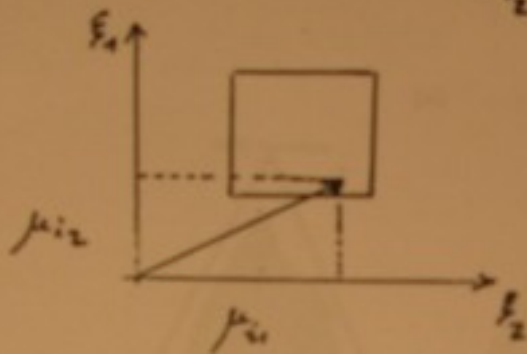


Fig. 5.18. Two-dimensional self-organizing system

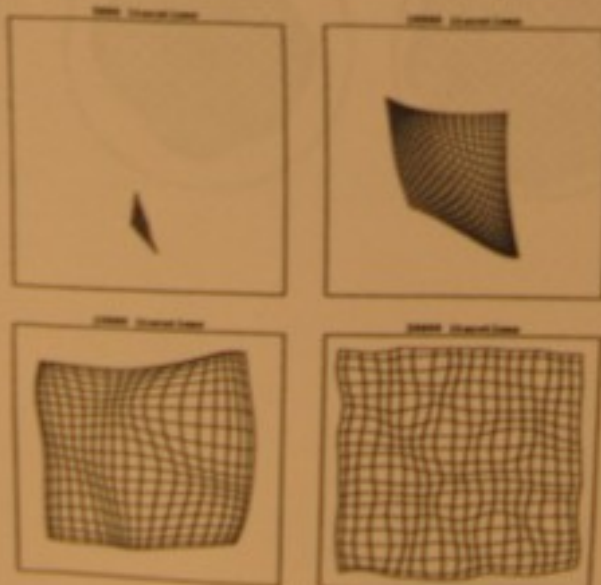
Cablage : Représentation des connexions dans le même espace que celui des patrons

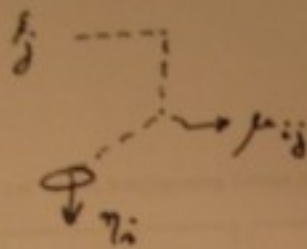
Exemple Rétine = Carré

$$\vec{x}_i = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$



On connecte 2 extrémités de vecteurs \vec{x}_i et \vec{x}_j si i et j labellent 2 cellules voisines sur la couche corticale





Algorithme

- 1) Identifier la cellule vainqueur i^* (η_{i^*} le plus élevé) au stimulus présenté
- 2) Adaptation (de type Hebbien) dans un voisinage de la cellule vainqueur, itérative après présentation de chaque pattern.

$$\frac{dw_{ij}}{dt} = \alpha (\eta_i x_j - \gamma(\eta_i) w_{ij}) \quad 0 \leq \alpha \leq 1$$

\nwarrow \searrow
 valeur de valeur d'entrée
 sortie de la en cellule j
 cellule i

$$\eta_i = 1 \quad \text{si } i \in \mathcal{V}_{i^*} \quad \gamma(0) = 0$$

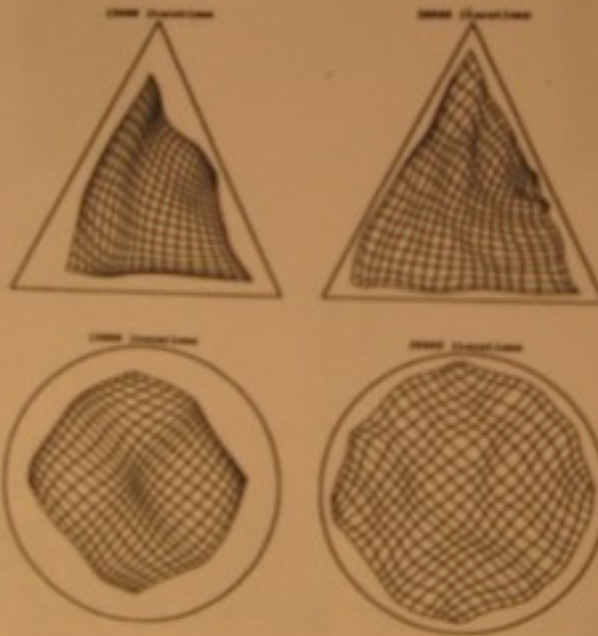
$$\eta_i = 0 \quad \text{si } i \notin \mathcal{V}_{i^*} \quad \gamma(1) = 1$$

Conséquences et exemples

Initialement connexions arbitraires

Adaptation \rightarrow Etablissement d'un role

(a)



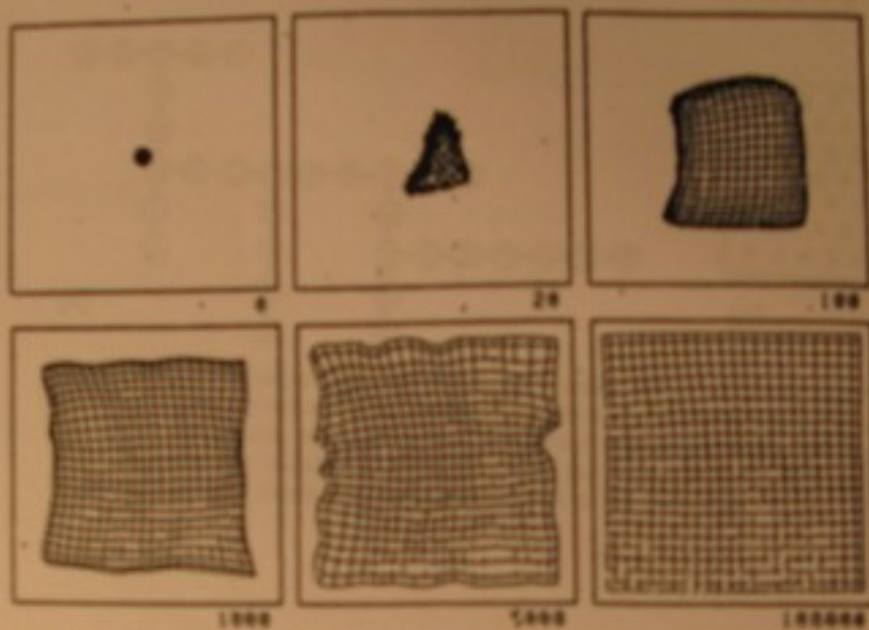


Fig. 3.16. Weight vectors during the ordering process

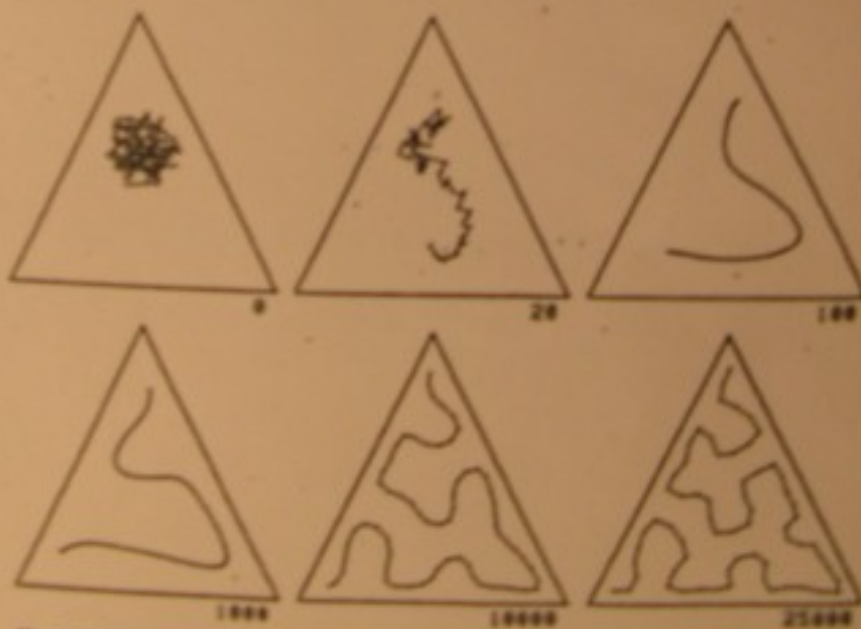


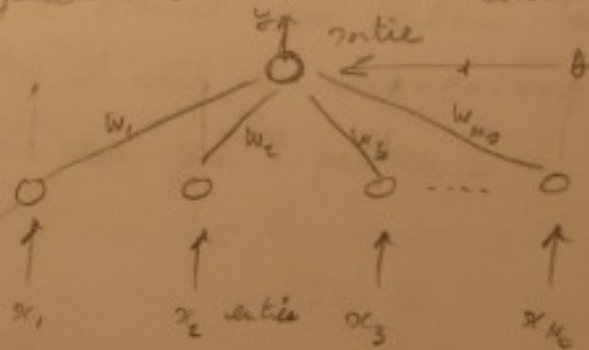
Fig. 3.17. Weight vectors during the ordering process

2. le rôle de Widrow-Hoff et de généralisation (apprentissage supervisé)

2.1 le réseau Adaline (adaptateur linéaire)

le réseau Adaline est le système le plus simple incluant un algorithme supervisé. Le but est de minimiser une fonction d'erreur calculée sur le écart entre valeurs de sortie et valeurs calculées par un réseau constitué d'une couche d'entrée et d'un neurone de sortie.

la fonction d'activation est linéaire



la valeur de sortie est donc

$$y = \sum_{j=1}^{N_0} w_j \cdot x_j + \theta$$

$$\text{Soit } \theta = w_{N_0+1}$$

on peut écrire

$$y = W^T z \quad \text{où}$$

$$W^T = (w_1 \ w_2 \ \dots \ w_{N_0} \ \theta)^T = \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_{N_0} \\ \theta \end{pmatrix}$$

$$z = (x_1 \ x_2 \ \dots \ x_{N_0} \ 1)^T = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_{N_0} \\ 1 \end{pmatrix}$$

Il s'agit d'obtenir la meilleure valeur de W telle que elle minimise l'erreur

$$Erreur = E = \frac{1}{2} \sum_{m=1}^M (y_j^{(m)} - t_j^{(m)})^2$$

entre la valeur d'entrée (la sortie) $y_j^{(m)}$ pour 14 patterns (ce sont des valeurs réelles - y_j = seul neurone en sortie) et la valeur $t_j^{(m)}$ calculée par le réseau pour chaque pattern qui nous est donné.

Il s'agit de trouver la valeur des poids W_j telle que E admet une valeur minimale.

Soit $W_j^{(m)}$ l'état des connexions après la présentation du $m^{ième}$ pattern.

qui représente l'état des connexions pour le pattern m .

L'erreur E est une fonction de $\{W_j^{(m)}\}_m$ par $\{y_j^{(m)}\}$ et la fonction (linéaire) d'activation

$$y_j^{(m)} = \sum_{i=1}^{N_0} W_{ij}^{(m)} x_i^{(m)} + \theta_j^{(m)}$$

Donc la dérivée partielle

$$\frac{\partial E}{\partial W_{ij}^{(m)}} = \frac{\partial E}{\partial y_j^{(m)}} \frac{\partial y_j^{(m)}}{\partial W_{ij}^{(m)}}$$

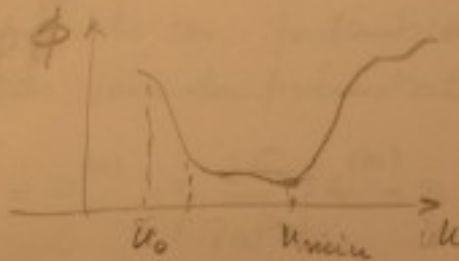
$$= (y^{(n)} - t^{(n)}) z_j^{(n)}$$

form vectorielle

$$\frac{\partial E}{\partial w^{(n)}} = (y^{(n)} - t^{(n)}) z^{(n)}$$

$$z = (x_1, x_2, \dots, x_{N_0+1})$$

A présent, comment modifier $w^{(n)}$ pour se diriger vers le minimum. Prenons l'exemple de la recherche du minimum d'une fonction



ϕ de u - nous aurons
c'est à dire $E \in E(w^{(n)})$
pour chaque n

On doit trouver u_{min} , de manière itérative. On part de u_0 et on fait

$$u_1 = u_0 - \eta \left(\frac{d\phi}{du} \right)_{u_0} \quad \eta > 0$$

$$\text{Comme } \left(\frac{d\phi}{du} \right)_{u_0} < 0 \text{ alors } u_1 > u_0$$

Arrivés en u_{min} , l'itération

$$\text{donne } u_k = u_{k-1} - \eta \left(\frac{d\phi}{du} \right)_k$$

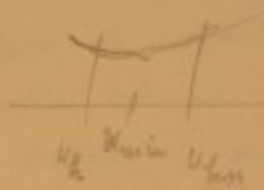
$$\text{si } \left(\frac{d\phi}{du} \right)_k = 0 \text{ exactement}$$

$$\text{alors } u_{k+1} = u_k$$

et le processus s'arrête

w_3

Ensuite, au fur de l'autre côté de u_{min}



$$\text{Alors } u_{k+2} = u_{k+1} - \eta \left(\frac{du}{dx} \right)_{k+1} > 0$$

$$\text{et } u_{k+2} < u_{k+1}$$

de nouveau vers u_{min} .

Ceci est le principe de la méthode.

Dans notre cas, partant de $w^{(m)}$ qui modifie les poids pour la présentation du pattern suivant

$$w_{(m+1)}^{(p)} = w^{(m)} - \eta \frac{\partial E}{\partial w^{(m)}} = w^{(m)} - \eta (y^{(m)} - t^{(m)}) z^{(m)}$$

L'algorithme est donc le suivant :

- 1) - on initialise les poids à des valeurs aléatoires petits
- 2) - on présente un pattern $(x_1, \dots, x_{N_0}, y) = z^{(m)}$ et on désigne une sortie désirée $t^{(m)}$ pour ce pattern.
- 3) - on calcule la sortie $y^{(m)}$ pour ce pattern
- 4) - les nouvelles valeurs de poids pour la présentation du pattern suivant est

$$w_{(m+1)}^{(p)} = w_{(m)}^{(p)} - \eta (y^{(m)} - t^{(m)}) z^{(m)}$$
 où η est un paramètre $0 < \eta < 1$
- 5) on retourne à 2) si l'erreur E est trop grande.

2.2 la propagation du gradient

En fait, la caractéristique linéaire des fonctions d'activation est une limitation très importante. Le réseau Adalines ne pourra pas fonctionner si les patterns sont bruités par exemple, comme dans le cas du réseau de Hopfield. Donc ~~pas~~ passer au non linéaire.

D'autre part, le passage à des structures multi-couche va augmenter les potentialités. Le prototype

la rétro-propagation du gradient

1 Architecture du réseau

Dans un réseau en couches, toute cellule d'entrée est suivie d'une succession de couches dites "cachées", pour finalement aboutir sur une couche de sortie.

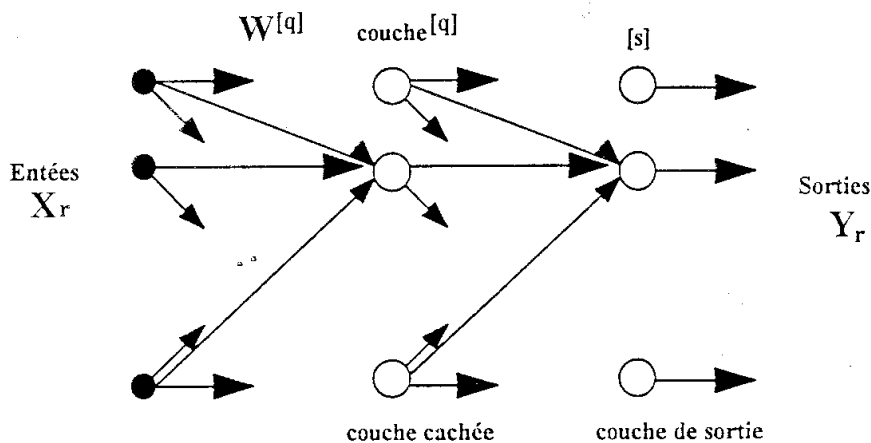


Fig. 1 - Architecture générale d'un réseau multi-couches

La présence de neurones des couches intermédiaires apporte une richesse à la structure qui doit accroître les capacités du réseau. Il faut insister sur le fait que les cellules internes n'ont aucune connexion prédéfinie, elles ne servent qu'à contribuer à l'obtention du résultat souhaité en sortie. Une autre particularité réside dans la fonction de transfert des unités. Ce sont des éléments semi-linéaires, qui compressent le potentiel du neurone dans l'intervalle $[-1, +1]$. Les fonctions qui réalisent cette tâche doivent être continues et dérivables.

Pour l'unité i de la couche $[q]$, en considérant que les $x_j[q]$ sont les entrées de cette couche, le potentiel $p_i[q]$ vaudra :

$$p_i[q] = \sum_j W_{ij} x_j[q]$$

et la sortie $y_i[q]$ sera déterminée, à partir du potentiel, comme :

$$y_i[q] = \sigma(p_i[q])$$

la fonction de transfert σ peut être, par exemple, égale à :

$$\sigma(v) = \text{th}(v) = \frac{e^v - e^{-v}}{e^v + e^{-v}}$$

où th dénote la fonction tangente hyperbolique dont la dérivée s'exprime par :

$$\text{th}(v)' = 1 - \text{th}^2(v)$$

Sur cette architecture, nous allons voir comment modifier les poids entre unités, c'est à dire comment faire l'apprentissage. L'architecture que nous venons de décrire comporte plusieurs paramètres, tels que la forme de la fonction de transfert, le nombre de couches, le nombre d'unités par couche. Actuellement, il n'existe pas de méthode analytique pour déterminer ces paramètres de structure. Il convient donc d'acquérir un savoir faire pour être capable d'exploiter au mieux ce modèle. C'est d'ailleurs ce qui le rend physiologiquement peu plausible et difficilement exploitable pour de gros problèmes.

2 Algorithme de retro-propagation du gradient

Comme nous venons de le voir , le réseau est structuré en couches successives, traversées par le signal d'entrée qui se propage vers la couche de sortie. Il n'est pas difficile de trouver un terme d'erreur à minimiser. Il suffira, en considérant qu'un superviseur est capable d'engendrer un objectif à atteindre, de calculer l'écart entre cet objectif et la réponse fournie par le réseau fig. 2 .

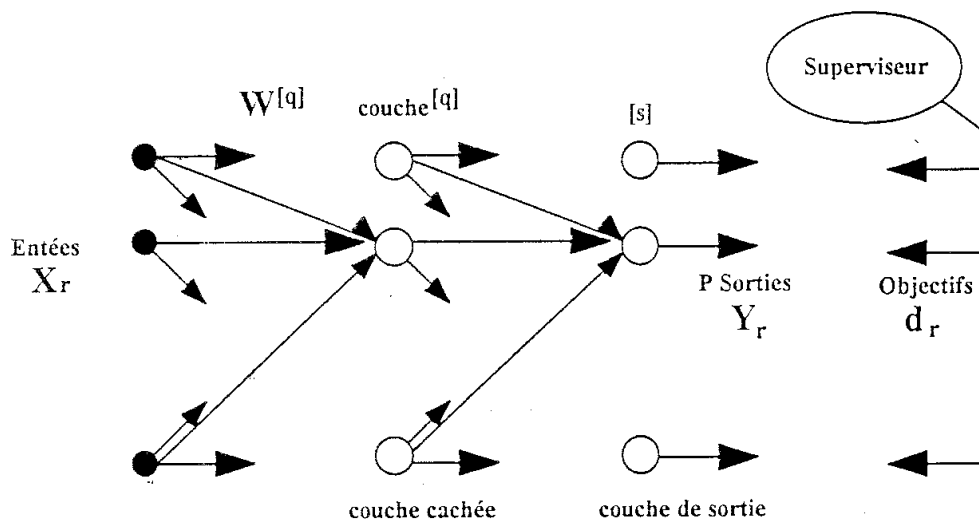


Fig. 2 -Calcul du terme d'erreur

L'algorithme d'ajustement des poids considère le terme d'erreur comme une grandeur à minimiser afin d'obtenir une configuration qui copie au mieux la fonction de transfert définie par les présentations successives des couples (x_r, d_r) . Or si nous voulons utiliser une technique de minimisation à base de gradient stochastique, seule la couche d'interconnexions reliée aux neurones de sortie a la connaissance de la grandeur à minimiser. Les couches internes, reliées aux neurones cachés, n'ont pas de terme à minimiser. Il n'y a ainsi aucune façon de connaître directement la variation des poids à appliquer. Le principe de l'algorithme de rétro-propagation du gradient réside dans la possibilité de donner un objectif aux neurones des couches internes et, par conséquent, une mesure de la variation des poids à appliquer pour le satisfaire. Un neurone placé dans une couche interne peut avoir l'erreur calculée en sortie du réseau à travers les poids qui le relient à la sortie. L'algorithme pourra ainsi moduler les poids qui arrivent sur un neurone afin de réduire sa participation à l'erreur calculée en sortie du réseau.

3 Définition du terme d'erreur

Nous définissons un terme d'erreur qu'il faut minimiser afin de trouver un ensemble de poids susceptibles de satisfaire le problème posé.

Pour un prototype r ($r = 1 \dots R$), avec K_S le nombre d'éléments de la couche de sortie le terme d'erreur instantanée, qui est un scalaire, s'écrit:

$$\xi_r = \frac{1}{2} \sum_{k=1}^{K_S} (d_{kr} - y_{kr})^2$$

Le terme d'erreur doit être défini sur l'ensemble des erreurs instantanées mesurées à chaque présentation d'un couple (sortie désirée d_r , sortie obtenue y_r).

$$\xi = \sum_{r=1}^R \xi_r$$

Mais, dans la suite de la démonstration, nous ne considérerons que le terme d'erreur instantanée, sans nuire en cela à la généralité du calcul.

4 Détermination du gradient de l'erreur par rapport aux poids

Comme nous l'avons constaté, la fonction d'un réseau est basée sur trois paramètres: les poids, le potentiel du neurone et sa fonction de transfert non linéaire. Nous décomposerons les calcul de façon à faire apparaître toutes les opérations élémentaires ce qui nous permettra de simplifier la présentation de la démonstration. En premier lieu, nous examinons la variation de l'erreur par rapport au poids $W_{ij}^{[q]}$ comme la composition de la variation de l'erreur par rapport à une variation du potentiel du neurone i de la couche $[q]$ et de la variation de ce potentiel en fonction des poids incidents $W_{ij}^{[q]}$. Le potentiel p_i est choisi car il dépend de $W_{ij}^{[q]}$.

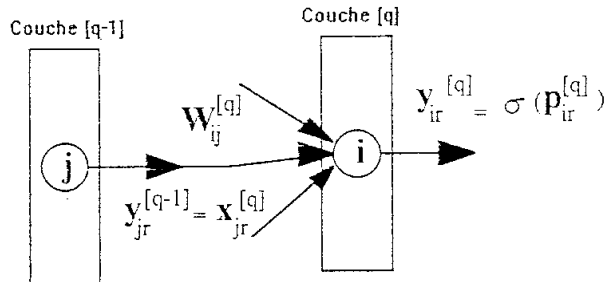
$$\frac{\partial \xi_r}{\partial W_{ij}^{[q]}} = \frac{\partial \xi_r}{\partial p_{ir}^{[q]}} \frac{\partial p_{ir}^{[q]}}{\partial W_{ij}^{[q]}}$$


Fig 3 -Propagation à travers les couches

Nous appelons signal d'erreur instantanée le premier terme de cette décomposition, et nous le notons de la façon suivante:

$$\delta_{ir}^{[q]} = - \frac{\partial \xi_r}{\partial p_{ir}^{[q]}}$$

Ce terme sera développé plus loin, en fonction de la couche $[q]$, et nous permettra de faire apparaître une forme récurrente pour son évaluation sur une couche interne à partir de la

sortie désirée. Considérons maintenant un neurone i placé sur une couche interne $[q]$. Son potentiel est calculé à partir des sorties de la couche précédente pondérées par les poids incidents. La variation de ce potentiel par rapport aux poids s'écrit:

si K_q est le nombre de neurones sur la couche q

$$\frac{\partial p_{ir}^{[q]}}{\partial W_{ij}^{[q]}} = \frac{\partial}{\partial W_{ij}^{[q]}} \left(\sum_{k=1}^{K_q} W_{ik}^{[q]} x_{kr}^{[q]} \right)$$

Le calcul de ce terme est immédiat et montre naturellement que le potentiel ne dépend que de l'entrée des neurones.

$$\frac{\partial p_{ir}^{[q]}}{\partial W_{ij}^{[q]}} = x_{jr}^{[q]}$$

En regroupant les deux membres de l'équation, nous pouvons dégager une forme simple de l'expression du gradient de l'erreur par rapport aux poids:

$$\frac{\partial \xi_r}{\partial W_{ij}^{[q]}} = -\delta_{ir}^{[q]} x_{jr}^{[q]}$$

5 Expression du signal d'erreur

Il faut maintenant expliciter le terme d'erreur δ_{ir} qui dépendra de la couche sur laquelle il est calculé. Nous décomposons le signal d'erreur comme la variation de l'erreur en fonction de la sortie multipliée par la variation de la sortie par rapport au potentiel:

$$\frac{\partial \xi_r}{\partial p_{ir}^{[q]}} = \frac{\partial \xi_r}{\partial y_{ir}^{[q]}} \frac{\partial y_{ir}^{[q]}}{\partial p_{ir}^{[q]}}$$

Ici encore, nous calculerons facilement le second membre de l'expression:

$$\frac{\partial y_{ir}^{[q]}}{\partial p_{ir}^{[q]}} = \sigma'(p_{ir}^{[q]})$$

c'est à dire que la variation de sortie n'est que la dérivée de la fonction de transfert du neurone appliquée à son potentiel. Il reste enfin à déterminer l'expression de la variation de l'erreur en fonction de celles de la sortie des neurones. Pour cela, nous allons considérer deux cas selon que nous nous trouvons sur la couche de sortie ou sur une couche cachée.

6.1. Couche de sortie: $q = s$

Pour la couche de sortie (s), qui contient K_s neurones, nous connaissons directement le terme d'erreur, comme il a été défini précédemment. L'expression du gradient de l'erreur par rapport aux sorties des neurones vaut donc:

$$\frac{\partial \xi_r}{\partial y_{ir}^{[s]}} = \frac{\partial}{\partial y_{ir}^{[s]}} \frac{1}{2} \sum_{k=1}^{K_s} (d_{kr} - y_{kr}^{[s]})^2$$

$$\frac{\partial \xi_r}{\partial y_{ir}^{[s]}} = -(d_{ir} - y_{ir}^{[s]})$$

Le signal d'erreur δ_{ir} , pour la couche de sortie vaut donc:

$$\delta_{ir}^{[s]} = (d_{ir} - y_{ir}^{[s]}) \sigma'(p_{ir}^{[s]})$$

6.2. Couches cachées: $q < s$

Pour les couches cachées, l'expression du terme d'erreur n'est pas directement connue: c'est la difficulté principale de cet algorithme. Nous allons considérer que le terme d'erreur est vu *au travers des interconnexions qui relient la couche interne à la couche de sortie*. Nous évaluerons l'erreur de la façon suivante:

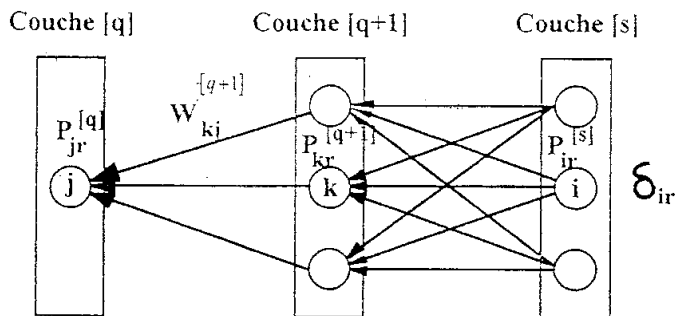


Fig Rétro-propagation du terme d'erreur

$$\frac{\partial \xi_r}{\partial y_{ir}^{[q]}} = \sum_{k=1}^{K^{[q+1]}} \frac{\partial \xi_r}{\partial p_{kr}^{[q+1]}} \frac{\partial p_{kr}^{[q+1]}}{\partial y_{ir}^{[q]}}$$

La variation de l'erreur en sortie du réseau par rapport à la sortie d'une couche interne est décomposée comme la somme de toutes les variations de l'erreur par rapport aux potentiels de la couche suivante multipliée par la variation de ces potentiels par rapport à la sortie du neurone considéré. Le second terme vaut, en écrivant $p_{kr}^{[q+1]}$ en fonction de la sortie de la couche précédente:

$$\frac{\partial \xi_r}{\partial y_{ir}^{[q]}} = \sum_{k=1}^{K^{[q+1]}} \frac{\partial \xi_r}{\partial p_{kr}^{[q+1]}} \frac{\partial}{\partial y_{ir}^{[q]}} \sum_{i=1}^{K^{[q]}} w_{ki}^{[q+1]} y_{ir}^{[q]}$$

donc:

$$\frac{\partial \xi_r}{\partial y_{ir}^{[q]}} = \sum_{k=1}^{K_{[q+1]}} \delta_{kr}^{[q+1]} W_{ki}^{[q+1]}$$

En reprenant la définition du signal d'erreur:

$$\frac{\partial \xi_r}{\partial p_{ir}^{[q]}} = \frac{\partial \xi_r}{\partial y_{ir}^{[q]}} \frac{\partial y_{ir}^{[q]}}{\partial p_{ir}^{[q]}}$$

nous pouvons extraire une formulation récurrente de δ_i . Le signal d'erreur pour une couche interne vaut:

$$\delta_{ir}^{[q]} = \sum_{k=1}^{K_{q+1}} \delta_{kr}^{[q+1]} W_{ki}^{[q+1]} \cdot \sigma'(p_{ir}^{[q]})$$

Cette expression permet de répercuter le signal d'erreur depuis la couche externe, où il est directement calculé à partir des sorties désirées, vers les couches internes où il est estimé.

7. Expression de la variation des poids.

Enfin, en reprenant la règle générale de modification des poids, on modifie les $W_{ij}^{(q)}$ suivant:

$$\Delta W_{ij}^{[q]} = \alpha \left(-\frac{\partial \xi_r}{\partial W_{ij}^{[q]}} \right) = \alpha (\delta_i^{[q]} x_{jr}^{[q]})$$

nous trouvons les formules suivantes:

1- Pour les poids entre l'avant dernière couche et la couche de sortie:

$$\Delta W_{ij}^{[s]} = \alpha (d_{ir} - y_{ir}) x_{jr}^{[s]} \sigma'(p_{ir}^{[s]})$$

2- Pour tous les poids entre la couche (q) et la couche(q+1)

$$\Delta W_{ij} = \alpha \left(\sum_{k=1}^{K_{[q+1]}} \delta_{kr}^{[q+1]} W_{ki}^{[q+1]} \right) x_{jr}^{[q]} \sigma'(p_{ir}^{[q]})$$

Cette expression signifie que l'erreur en sortie du réseau est propagée en sens rétrograde avec les mêmes poids de connexion que ceux employés en sens direct.

La procédure de modification des poids se résume donc de la manière suivante:

-(1) Initialisation des poids $W(q)$ entre les neurones des couches consécutives $q=1,...,s$ à des petites valeurs aléatoires

-(2) Présentation d'une entrée x_r et de la sortie désirée d_r , r est le numéro de la présentation d'un couple (entrée, sortie désirée), $r = 1,..., R$

-(3) Calcul de la sortie actuelle par propagation directe à travers toutes les couches:

$$y_{jr}^{[q]} = \sigma\left(\sum_{i=1}^{K_q} W_{ji}^{[q]} y_{ir}^{[q-1]}\right) = \sigma\left(\sum_{i=1}^{K_q} W_{ji}^{[q]} x_{ir}^{[q]}\right)$$

où σ est la fonction neurone.

Pour la couche d'entrée $q = 1$, $y^{[q-1]} = x$

-(4) Calcul de l'erreur en sortie (pour le tracé de l'évolution du terme d'erreur, mais n'est pas nécessaire au calcul d'adaptation).

$$\xi_r = \frac{1}{2} \sum_{K=1}^{K_s} (d_{kr} - y_{kr}^{[s]})^2$$

où

d_r est la sortie désirée associée au vecteur x_r

$y_r^{[s]}$ est la sortie obtenue sur la dernière couche lors de la présentation du vecteur d'entrée x_r

ξ_r est l'erreur calculée pour le couple (x_r, d_r)

-(5) Rétro-propagation du signal d'erreur (δ_{ir}) depuis la dernière couche vers la première couche :

-pour chaque cellule de sortie:

$$\delta_{ir}^{[s]} = \sigma'(p_{ir}^{[s]})(d_{ir} - y_{ir}^{[s]})$$

- pour les cellules cachées:

$$\delta_{ir}^{[q]} = \sum_{k=1}^{K_{q+1}} \delta_{kr}^{[q+1]} W_{ki}^{[q+1]} \sigma'(p_{ir}^{[q]})$$

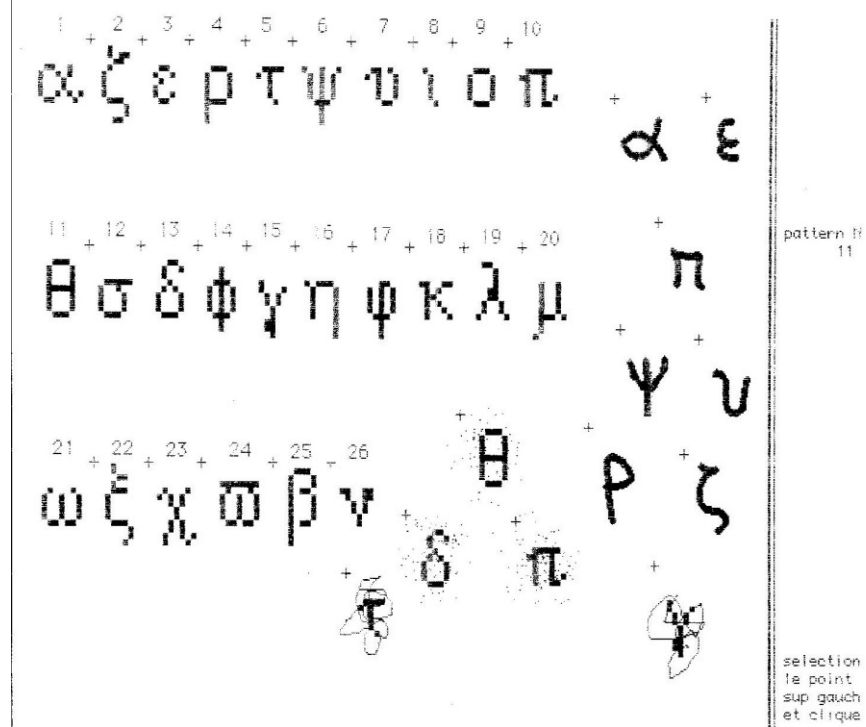
-(6) Mise à jour des poids selon la règle:

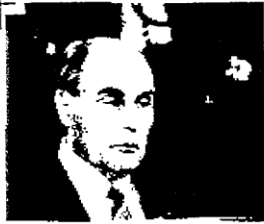
$$\Delta W_{ij}^{[q]} = \alpha (\delta_{ir}^{[q]} x_{jr}^{[q]})$$

où α est un paramètre d'apprentissage compris entre 0 et 1.

-(7) retour à (2) tant qu'il y a des couples à présenter

Applications : Identification de patterns





Prédiction de séries temporelles

