

Réseaux Multicouches

Apprentissage et généralisation

PMC (partie 2)

1 - Le choix du modèle (résultats théoriques)

Plusieurs résultats théoriques généraux (biais-variance, dimension de Vapnik-Chernovenskis, ...) ont montré que, l'objectif de l'apprentissage est de trouver le modèle dont la « complexité » réalise un bon « compromis » qui tient compte de la complexité du problème à traiter et de la taille de la base d'apprentissage.

La méthode adoptée consiste, considérant une suite emboîtée de L espaces de fonctions $(F_l)_{l=1..L}$ de « richesse » croissante, dans le sens suivant :

$$F_1 \subset F_2 \subset \dots \subset F_{L-1} \subset F_L$$

à appliquer l'algorithme « *minimisation du risque structurel (MRS)* » ci-dessous :

Étant donnée une même base d'apprentissage A_{pp} de taille N :

Pour tout l de 1 à L

- Appliquer l'algorithme d'apprentissage sur la famille F_l ce qui permet d'obtenir une des fonctions f_l de F_l . On note f_{l*} cette fonction particulière de F_l qui minimise l'erreur en apprentissage.
- Calculer l'erreur en test de la fonction f_{l*} .

Choisir, la fonction f_{l*} pour laquelle l'erreur en test est la plus petite.

PMC (partie 2) 1-choix modèle

1.2 - choix du modèle pour les réseaux multicouches

Plan

Plusieurs méthodes sont envisageables :

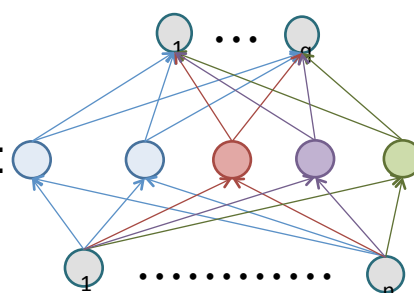
- Construction empirique de l'architecture
- Architecture à masques et poids partagés
- Régularisation de la fonction de coût
- Apprentissage avec bruit
- Suppression de poids (ou de cellules)

PMC (partie 2) 2 - Construction empirique de l'architecture

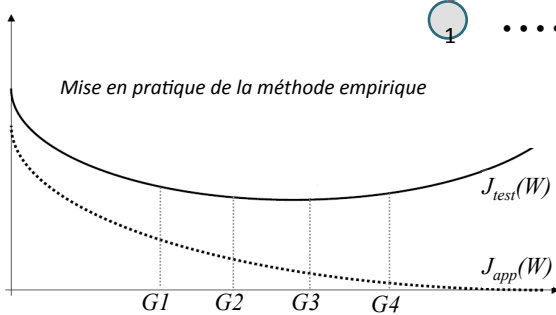
2.1- Familles de fonctions emboîtées de R^n dans R^q .

exemple ici avec une seule couche cachée et initialement 2 neurones :

$G_1 \subset G_2 \subset G_3 \subset G_4 :$



Mise en pratique de la méthode empirique



PMC (partie 2) 2-construction empirique

2.2 - Propriété d'approximation universelle des PMC

Théorème de Funahashi K.I.:
 Soit $f(x)$ une fonction non constante, bornée, monotone croissante et continue. Soit K un sous ensemble compact de R^n . (« par exemple une sphère de R^n ») pour toute fonction T continue, $T : K \rightarrow R^q$, définie par $x = (x_1, x_2, \dots, x_n) \rightarrow (T_1(x), \dots, T_q(x))$, et pour tout $\varepsilon > 0$, il existe un réseau A à une couche cachée, dont la fonction d'activation est $f(x)$ pour chacune des cellules cachées et linéaire pour les cellules de sortie, tel que :
 $\max_{x \in K} [d(T(x) - G(x, W))] < \varepsilon$
 où W est la matrice des poids de connexions associées, $d(\cdot)$ la métrique euclidienne de R^q et $G(\cdot, W)$ la fonction générée par le PMC.

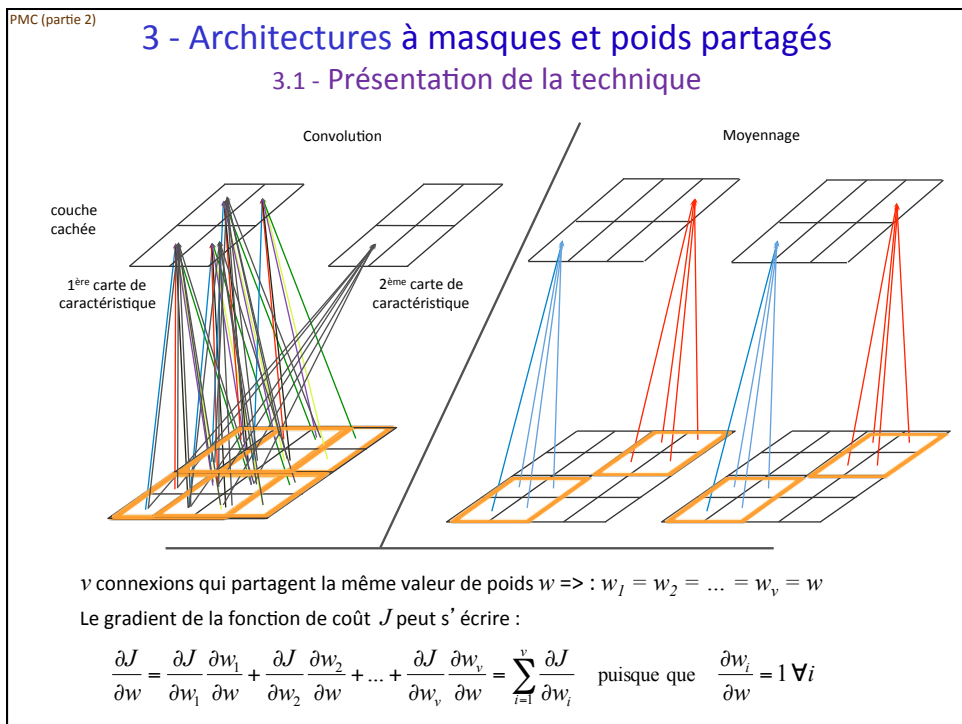
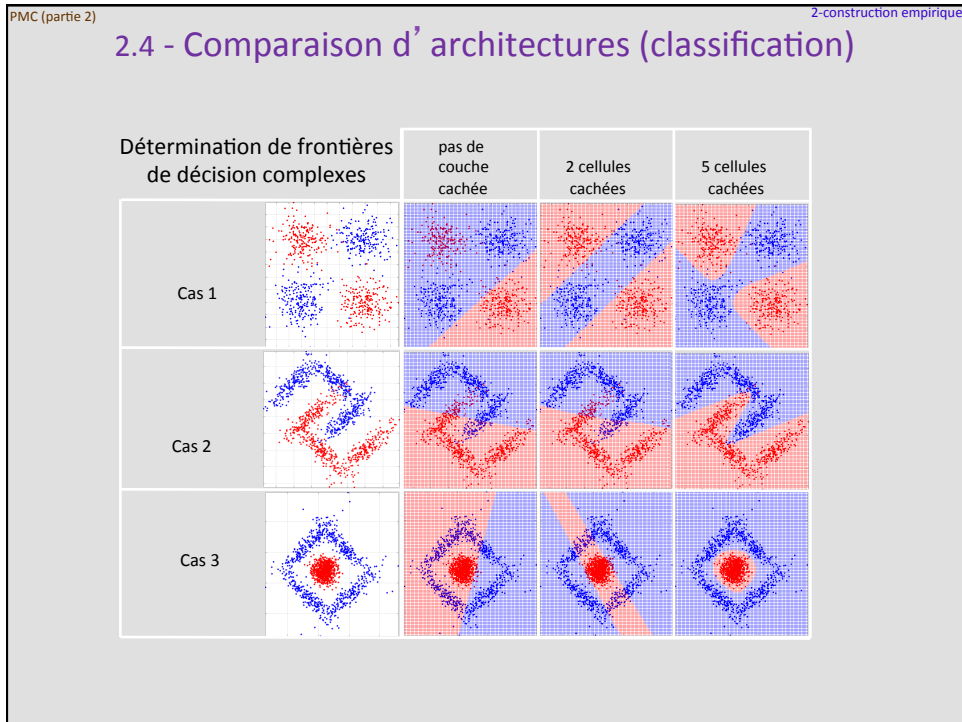
PMC (partie 2) 2-construction empirique

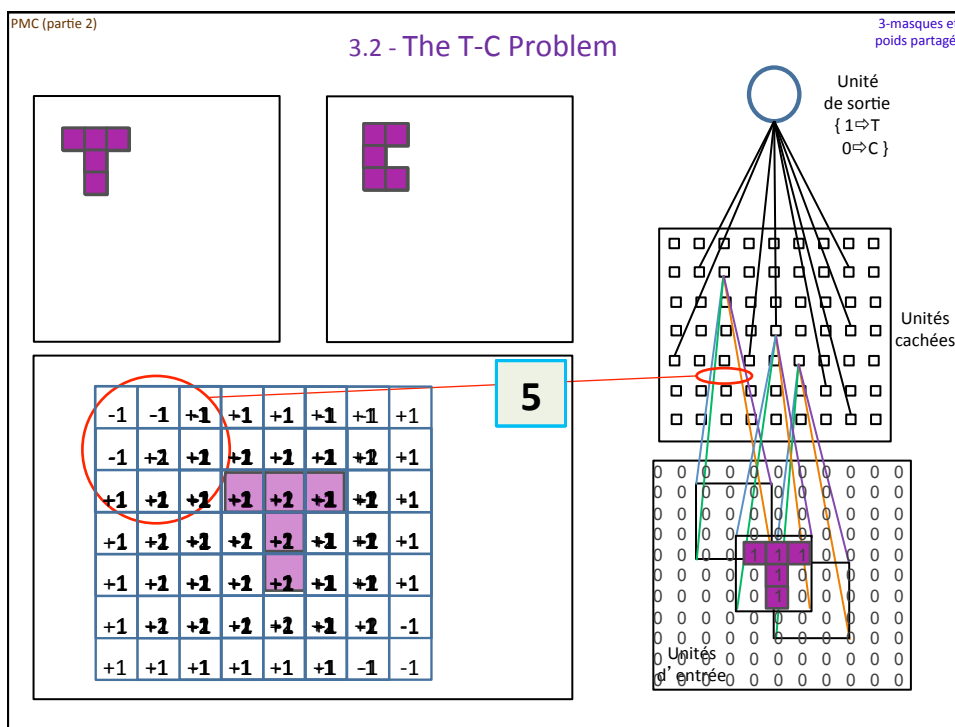
2.3 - Choix du meilleur réseau (régression)

ensemble d'Apprentissage

ensemble de Test

Unité cachées	E_{app}	E_{rest}
1	0,012	0,0204
2	0,0093	0,0147
3	0,0077	0,0382
4	0,0043	0,0784
5	0,0019	0,1569





PMC (partie 2) 3-masques et poids partagés

3.3 - Exemple: Reconnaissance de chiffres manuscrits

Figure 1: exemples de codes postaux originaux extraits de la base d'apprentissage provenant d'une poste newyorkaise (Buffalo).

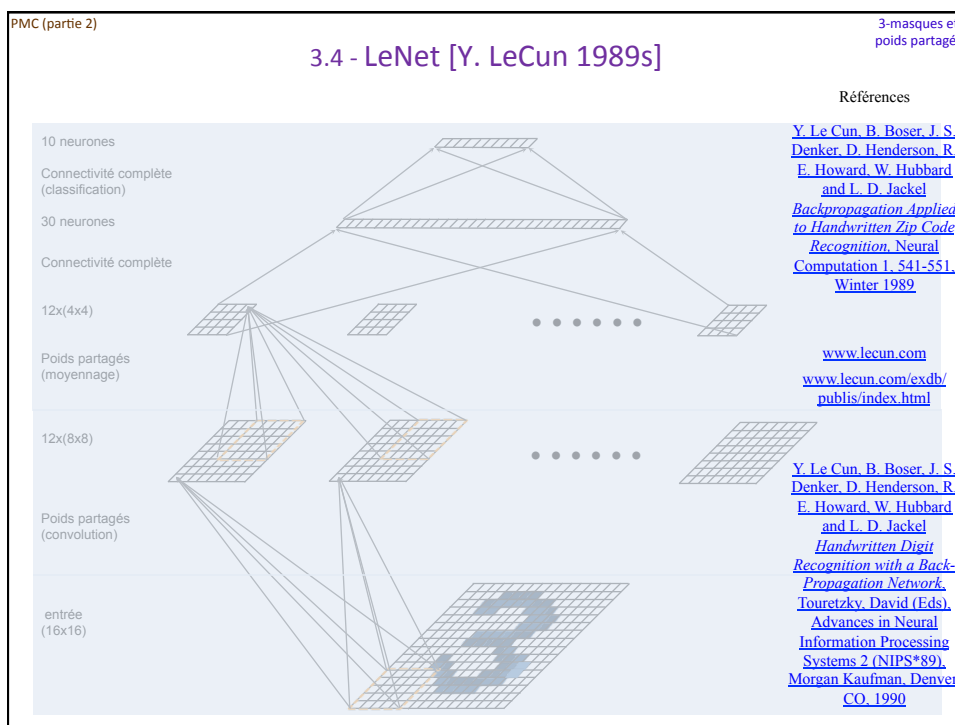
40004 75416
 14199-2087 23505
 96203 44310
 44151 05753

Figure 2: exemples de chiffres normalisés après localisation et segmentation [Wang & Srihari, 1988]

1 5 4 8 5 7 2 6 8 0
 2 0 2 9 9 2 9 9 7 2
 1 1 5 9 1 0 1 0 6 1
 1 1 0 3 0 4 7 5 2 6
 7 6 7 2 8 5 5 7 1 3
 3 0 1 2 7 1 1 2 9 9
 0 7 5 9 7 3 3 1 9 7
 1 2 2 5 5 1 8 2 8 1
 2 1 6 5 5 4 6 0 3 5
 1 0 8 5 0 3 0 4 7 5

Tiré de : Y. Le Cun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard and L. D. Jackel *Handwritten Digit Recognition with a Back-Propagation Network*, AT&T Bell Laboratories, Holmdel, N. J. 07733

commentaire



PMC (partie 2)

4 - Terme de pénalité dans la fonction de coût

4.1 - Principe

Soit la fonction $E_\lambda(\mathbf{W}) = J(\mathbf{W}) + \lambda R(\mathbf{W})$ à minimiser

où:

- $J(\mathbf{W})$ l'erreur quadratique obtenue sur un ensemble d'apprentissage D
- $R(\mathbf{W})$ une fonction qui dépend des paramètres \mathbf{W} et qui correspond à un terme de pénalité. On suppose $R(\mathbf{W}) \geq 0$
- λ est un paramètre à déterminer.

Résultat 1

On démontre qu'il existe $k > 0$ tel que minimiser $E_\lambda(\mathbf{W})$ est équivalent à :

minimiser $J(\mathbf{W})$ sous la contrainte $R(\mathbf{W}) \leq k$

et donc, minimiser $E_\lambda(\mathbf{W})$ est équivalent à minimiser $J(\mathbf{W})$ dans la famille de fonctions

$$G_\lambda = \{g \in G(\cdot, \mathbf{W}) \mid R(\mathbf{W}) \leq k\}$$

Résultat 2

Soit λ_1 et λ_2 tel que $\lambda_1 < \lambda_2$, alors on démontre que les 2 nombres k_1 et k_2 (qui découlent du résultat précédent) vérifient l'inégalité $k_2 \leq k_1$.

On a donc :

$$\lambda_1 < \lambda_2 \Rightarrow G_{\lambda_2} \subset G_{\lambda_1}$$

PMC (partie 2)

4.2 - Lien avec l' algorithme de minimisation du risque structurel

- Étant donné :
- une suite de paramètres $\lambda_L > \lambda_{L-1} \dots > \lambda_2 > \lambda_1$,
 - une famille de fonctions $G(\cdot, W)$ (définie par un PMC)
 - une même base d' apprentissage A_{pp} de taille N :

l' algorithme « **minimisation du risque structurel** » devient :

Pour tout l de 1 à L

- Appliquer l' algorithme de minimisation de la fonction $E_M(W)$ dans la famille de fonctions G , ce qui permet d' obtenir une fonction g_{l^*} de G .
- Calculer l' erreur en test de la fonction g_{l^*} .

Choisir, la fonction g_{l^*} pour laquelle l' erreur en test est la plus petite.

PMC (partie 2)

4-terme de pénalité

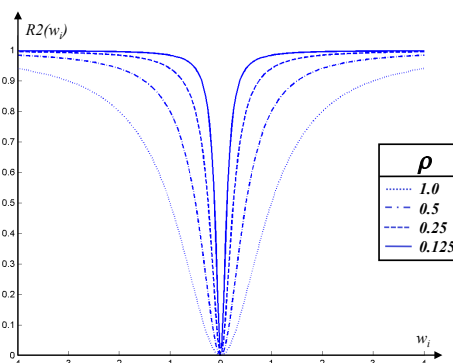
4.3 - Exemples de terme de pénalité

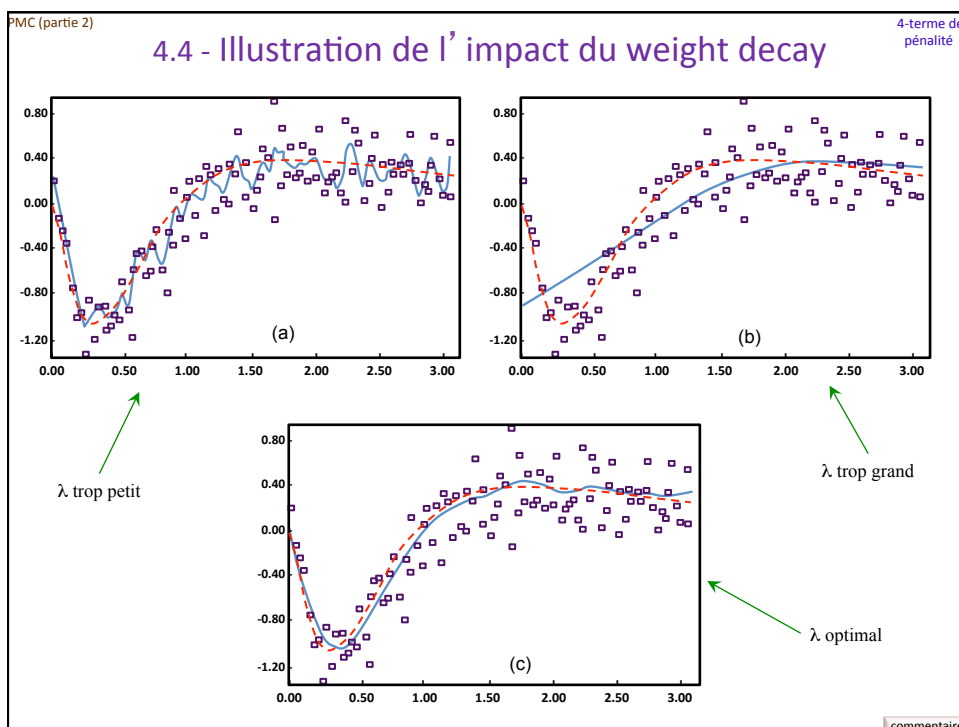
1) weight decay :

$$R1(W) = \sum_i w_i^2$$

2) :

$$R2(W) = \sum_i \left(\frac{w_i^2}{\rho^2} \right) \left(1 + \frac{w_i^2}{\rho^2} \right)$$





PMC (partie 2)

5 - Apprentissage avec bruit (sur les entrées)

- Soit une base d'apprentissage : $D = \{ (\mathbf{x}^1, \mathbf{d}^1), (\mathbf{x}^2, \mathbf{d}^2), \dots, (\mathbf{x}^N, \mathbf{d}^N) \}$
- A chaque exemple $(\mathbf{x}^k, \mathbf{d}^k)$ on construit l autres exemples : $(\mathbf{x}^k + \mathbf{e}^i, \mathbf{d}^k)$ $i=1, 2, \dots, l$ où les \mathbf{e}^i correspondent à des tirages indépendants suivant la loi normale $N(0, \sigma^2 I)$.
- On forme ainsi une base d'apprentissage D' ayant $N \times l$ éléments.

⇒ On peut montrer que l'apprentissage du PMC sur D' correspond à minimiser

$$E(\mathbf{w}) = J(\mathbf{w}) + \sigma^2 \Omega(\mathbf{x})$$

où $\Omega(\mathbf{x})$ est le terme de régularisation, et où σ intervient comme un paramètre de régularisation à déterminer.

σ_1^2

$g(\mathbf{x}, \mathbf{w})$: avec bruit de paramètre σ_1

$g(\mathbf{x}, \mathbf{w})$: sans bruit

σ_2^2

$g(\mathbf{x}, \mathbf{w})$: avec bruit de paramètre σ_2

$g(\mathbf{x}, \mathbf{w})$: sans bruit

2 illustrations avec $\sigma_1 < \sigma_2$: En terme de généralisation, la courbe orange (σ_2) est préférable à la verte (σ_1) toutes deux étant préférables à la courbe bleue pour laquelle on a pas ajouté de bruit.

PMC (partie 2)

6 - Méthode de suppression de poids

6.1 - Exemple : OBD : Optimal Brain Damage

La pertinence d'un poids est estimée à partir de la formule de Taylor :

$$\delta J = J(w + \delta w) - J(w) = \sum_i \frac{\partial J}{\partial w_i} \Big|_w \delta w_i + \frac{1}{2} \sum_i \sum_j \frac{\partial^2 J}{\partial w_i \partial w_j} \Big|_w \delta w_i \delta w_j + o(\delta w_i^3) \quad (1)$$

En un minimum local ce premier terme est nul.

Sous l'hypothèse que la Hessienne est diagonale : $\delta J \approx \frac{1}{2} \sum_i \frac{\partial^2 J}{\partial w_i^2} \delta w_i^2$ (2)

Mesure de pertinence :

Partant d'un vecteur poids

$$w = (w_1, \dots, w_k, \dots, w_N)$$

et d'une variation

$$\delta w = (0, \dots, 0, -w_k, 0, \dots, 0)$$

on obtient :

$$OBD(w_k) = \frac{1}{2} w_k^2 \frac{\partial^2 J}{\partial w_k^2} \quad (3)$$

référence: Y. Lecun, J. S. Denker, S. A. Solla (1990). *Optimal Brain Damage*. In D. S. Touretsky (Ed), *Advances in Neural Information Processing Systems*, Volum 2, pp.598-605. San Mateo, CA: Morgan Kaufmann.

PMC (partie 2)

6-suppression de poids

6.2 - Algorithme d'élimination de poids

Algorithme général

- 1) Faire une phase d'apprentissage complète, et calculer l'erreur en test.
- 2) Calculer la pertinence de chaque poids avec OBD.
- 3) Supprimer un (ou des) poids dont la pertinence est la plus petite.
- 4) Recommencer en 1) tant que l'erreur en test ne se dégrade pas.

Conclusion

Dans cette séance nous avons présenté l'algorithme de minimisation du risque structurel (MRS) qui consiste à construire des familles de fonctions emboîtées, et à retenir celle qui a la meilleure performance. Jusqu'à présent, nous avons exposé deux techniques qui permettent de mettre en œuvre cette méthode, il s'agit de :

- La construction empirique d'architectures, et des
- Architecture à masques et poids partagés.
- La régularisation de la fonction de coût
- L'apprentissage avec bruit, et
- La suppression de poids.

Ces sujets sont sous-tendus par des approches mathématiques plus complexes.