

III - Réseaux Multicouches (Partie 1)

PMC: Perceptron Multi-couche (MLP: Multi Layer Perceptron)

3ème séance :

- 6) Estimation de l'Espérance
 - 6.1) Approche théorique
 - 6.2) Cas de la classification

- 7) Estimation de la Variance
 - 7.1) Maximum de vraisemblance (rappel)
 - 7.2) Application à la régression par PMC
 - 7.3) Calcul de la vraisemblance

Préambule

On se focalise dans ce cours sur les performances attendues lorsque l'on met au point un système de diagnostique automatique. On veut donc connaître les erreurs commises par celui-ci dans tous les cas de fonctionnement possibles. On s'intéresse donc à étudier l'approximation obtenue quand on remplace une réalité (fonction réelle sous jacente aux données) par une modélisation mathématique (fonction mathématique choisie par le modélisateur). Le but de ce cours est de présenter les résultats théoriques des réseaux de neurones en se plaçant dans ce contexte. Ces résultats établissent la capacité des PMC à approximer, à partir de l'ensemble d'apprentissage, espérances et variances conditionnelles qui permettent de caractériser la nature du phénomène étudié.

PMC (partie 1)

6 - Estimation de l'Espérance

6.1 - Approche théorique (1)

- Un système réel est décrit comme un vecteur aléatoire x régi par une densité de probabilité $P(x)$.
- A chaque vecteur x est associée une valeur d suivant la loi conditionnelle $P(d/x)$.
- Généralement, l'apprentissage consiste à déterminer les paramètres W^* pour une famille de fonctions $\{F(x, W)\}$ qui minimisent l'**erreur d'apprentissage** sur un ensemble de N exemples : $App = \{(x_1, d_1), (x_2, d_2), \dots, (x_N, d_N)\}$

$$E_{App} = \frac{1}{N} \sum_{i=1}^N (d_i - F(x_i, W))^2$$

Mais la fonction $F(x, w^*)$ ne minimise pas nécessairement l'**erreur en généralisation** :

$$E_{géné} = \iint (d - F(x, W^*))^2 p(x) \cdot p(d/x) dx dd$$

commentaire

PMC (partie 1)

6 - Estimation
Espérance

6.1 - Approche théorique (2)

Soit : $E(d/x) = \int d \cdot p(d/x) dd$ on démontre le résultat suivant :

$$\text{Min } E_{géné} \Leftrightarrow \text{Min } \int (E(d/x) - F(x, W))^2 p(x) dx$$

Remarque :

La meilleure fonction en généralisation est celle qui approxime au mieux $E(d/x)$

Démonstration : $E_{géné} = \iint (d - E(d/x) + E(d/x) - F(x, W))^2 p(x, d) dx dd$

$$= \iint (d - E(d/x))^2 p(x, d) dx dd \quad (1) + \iint (E(d/x) - F(x, W))^2 p(x) dx \quad (2)$$

$$+ 2 \iint (d - E(d/x)) (E(d/x) - F(x, W)) p(x, d) dx dd \quad (3)$$

$$(3) = 2 \int (E(d/x) - F(x, W)) \left[\underbrace{\int (d - E(d/x)) p(d/x) dd}_{\int d \cdot p(d/x) dd - E(d/x) = 0} \right] p(x) dx$$

Donc $E_{géné} = (1) + (2)$ et comme (1) est indépendant de W , on tire :

$$\text{Min } E_{géné} \Leftrightarrow \text{Min } (2)$$

commentaire

PMC (partie 1) 6 - Estimation
Espérance

6.1 - Approche théorique (3)

But de l'apprentissage : $F(\mathbf{x}, W^*) \simeq E(d/\mathbf{x})$... Mais $E_{\text{géné}}$ et $E(d/\mathbf{x})$ sont inconnus

Dans la pratique, on cherche à estimer $E(d/\mathbf{x})$ à partir d'un ensemble d'apprentissage : $App = \{(x_1, d_1), (x_2, d_2), \dots, (x_N, d_N)\}$

L'apprentissage consiste à minimiser le « risque empirique » :

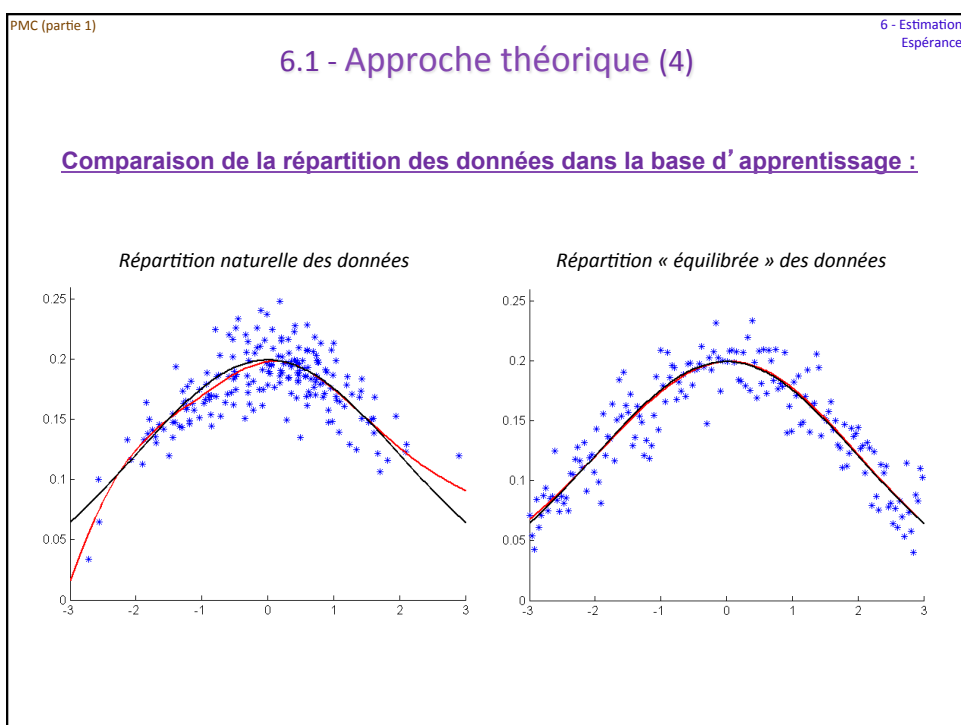
$$E_{App} = \frac{1}{N} \sum_{i=1}^N (d_i - F(\mathbf{x}_i, W))^2$$

La solution obtenue dépend alors de l'ensemble d'apprentissage : $F(\mathbf{x}, W^*/App)$

Remarque : Ce résultat reste valable dans le cas de l'erreur quadratique pondérée :

$$E_{\text{géné}} = \iint \frac{(d - F(\mathbf{x}, W))^2}{\sigma^2(\mathbf{x})} p(\mathbf{x}) p(d/\mathbf{x}) d\mathbf{x} dd$$

Ainsi,

$$\text{Min } E_{\text{géné}} \Leftrightarrow \text{Min } \int (E(d/\mathbf{x}) - F(\mathbf{x}, W))^2 \frac{p(\mathbf{x})}{\sigma^2(\mathbf{x})} d\mathbf{x}$$


PMC (partie 1) 6 - Estimation
Espérance

6.1 - Approche théorique (5)

Conditions d'approximation :

- Pour les PMC, la précision de l'approximation de $E(\mathbf{d}/\mathbf{x})$ dépend de l'architecture choisie
- Réalisation d'un codage des données approprié (par exemple une normalisation)
- L'échantillon d'apprentissage doit être de taille N suffisante et représentatif

Remarque :

- La précision de cette approximation pour un point \mathbf{x} fixé diffère selon la répartition des exemples $p(\mathbf{x})$.

$$\begin{array}{l} p(\mathbf{x}) \text{ grand} \quad \Leftrightarrow \quad F(\mathbf{x}, W) \text{ proche de } E(\mathbf{d}/\mathbf{x}) \\ p(\mathbf{x}) \text{ petit} \quad \Leftrightarrow \quad F(\mathbf{x}, W) ? \end{array}$$

En pratique les régions sous représentées dans l'ensemble d'apprentissage sont plus ou moins ignorées.

commentaire

PMC (partie 1) 6 - Estimation
Espérance

6.2 - Cas de la classification (1)

Cas de 2 classes

L'ensemble d'apprentissage est réparti en 2 classes. On code les réponses désirées :

$$\begin{cases} d = a & \text{si } x \in C_1 \\ d = b & \text{si } x \in C_2 \end{cases}$$

$F(\mathbf{x}, W)$ « approxime » : $E(\mathbf{d}/\mathbf{x}) = a \times p(C_1/\mathbf{x}) + b \times p(C_2/\mathbf{x})$

Interprétation :

Si $a = 1$ et $b = 0$ (alors) $E(\mathbf{d}/\mathbf{x}) = p(C_1/\mathbf{x})$

Si $a = 1$ et $b = -1$ (alors) $E(\mathbf{d}/\mathbf{x}) = p(C_1/\mathbf{x}) - p(C_2/\mathbf{x})$

Fonction discriminante de Bayes

Généralisation à Q classes :

Le vecteur qui code la classe k (C_k) a toutes ses composantes égales à b exceptée la k ème qui prend la valeur a .

$$E(\mathbf{d} / \mathbf{x}) = a p(C_k / \mathbf{x}) + \sum_{q \neq k} b p(C_q / \mathbf{x})$$

commentaire

PMC (partie 1) 6 - Estimation
Espérance

6.2 - Cas de la classification (2)

Approximation des probabilités : cas du codage a=1 b=0 :

- En général les sorties du réseau ne donnent pas de vraies probabilités

$$\sum_{k=1}^Q F_k(\mathbf{x}, W) \neq 1$$

- Il est cependant possible d'obtenir des probabilités en utilisant la fonction de transfert « softmax » pour les cellules de sortie :

$$y_k = f_k(A_k) = \frac{e^{A_k}}{\sum_{q \in \text{couche de sortie}} e^{A_q}}$$

on a bien : $\sum_k y_k = 1$
et : $y_k \geq 0 \quad \forall k$

commentaire

PMC (partie 1) 6 - Estimation
Espérance

6.2 - Cas de la classification (3)

Sortie 1

Sortie 2

commentaire

7 - Estimation de la Variance

7.1 - Maximum de Vraisemblance (rappel)

Soit un ensemble d'observations $\{\mathbf{x}_1, \dots, \mathbf{x}_p, \dots, \mathbf{x}_N\}$ réalisations indépendantes d'une v.a. X identiquement distribuée selon une loi de probabilité admettant θ comme paramètre de fonction de densité : $f(\mathbf{x}, \theta)$

On appelle « vraisemblance de θ » selon l'échantillon $\{\mathbf{x}_1, \dots, \mathbf{x}_p, \dots, \mathbf{x}_n\}$ le nombre :

$$L(\mathbf{x}_1, \dots, \mathbf{x}_i, \dots, \mathbf{x}_N, \theta) = f(\mathbf{x}_1, \theta) \times f(\mathbf{x}_2, \theta) \times \dots \times f(\mathbf{x}_N, \theta) = \prod_{i=1}^N f(\mathbf{x}_i, \theta)$$

- => Trouver la valeur du paramètre θ qui rend maximum cette vraisemblance.

$$\text{Max}_{\theta} \prod_{i=1}^N f(\mathbf{x}_i, \theta)$$

- **Estimateur par Maximum de Vraisemblance** noté EMV. L'EMV peut ne pas exister, et peut ne pas être unique.
- => Choisir analytiquement la fonction f .

PMC (partie 1)

7 - Estimation
Variance

7.2 - Application à la régression par PMC

$App = \{(\mathbf{x}_i, \mathbf{d}_i); i = 1, 2, \dots, N\}$ ($\mathbf{x}_i \in \mathbb{R}^p, \mathbf{d}_i \in \mathbb{R}^q$) de probabilités $p(\mathbf{x}, \mathbf{d})$.

- On suppose que toutes les observations $(\mathbf{x}_i, \mathbf{d}_i)$ sont indépendantes. Ainsi :

$$\begin{aligned} P(App) &= P((\mathbf{x}_1, \mathbf{d}_1), (\mathbf{x}_2, \mathbf{d}_2), \dots, (\mathbf{x}_N, \mathbf{d}_N)) \\ &= p(\mathbf{x}_1, \mathbf{d}_1) \times p(\mathbf{x}_2, \mathbf{d}_2) \times \dots \times p(\mathbf{x}_N, \mathbf{d}_N) \end{aligned}$$

$$P(App) = \prod_{i=1}^N p(\mathbf{x}_i, \mathbf{d}_i) = \prod_{i=1}^N p(\mathbf{d}_i / \mathbf{x}_i) p(\mathbf{x}_i)$$

- D'autre part, si on fait l'hypothèse : $\mathbf{x} \longrightarrow \mathbf{d} = \mu(\mathbf{x}) + \varepsilon(\mathbf{x})$ où $\mu(\cdot)$ est une fonction bien déterminée et $\varepsilon(\mathbf{x})$ une variable aléatoire qui suit une loi normale $N(0, \sigma(\mathbf{x})\mathbf{I})$. Ainsi :

$$p(\mathbf{d} / \mathbf{x}) = \frac{1}{\sigma(\mathbf{x}) \sqrt{2\pi}} \exp \left[-\frac{1}{2} \frac{\|\mathbf{d} - \mu(\mathbf{x})\|^2}{\sigma^2(\mathbf{x})} \right]$$

$$P(\mathbf{x}_i, \mathbf{d}_i) = \frac{1}{(\sqrt{2\pi})^N} \prod_{i=1}^N \frac{1}{\sigma(\mathbf{x}_i)} \exp \left[-\frac{1}{2} \frac{\|\mathbf{d}_i - \mu(\mathbf{x}_i)\|^2}{\sigma^2(\mathbf{x}_i)} \right] p(\mathbf{x}_i) \quad \text{Loi Gaussienne } N(\mu(\mathbf{x}), \sigma^2(\mathbf{x}))$$

PMC (partie 1) 7 - Estimation
Variance

7.3 - Calcul de la vraisemblance (1)

- Soit une famille de fonctions : $F = \{ F(x, W) \}$
- On désire approximer μ par une fonction $F(x, W)$ particulière de la famille F . Ainsi :

$$P(App, W) = \frac{1}{(\sqrt{2\pi})^N} \prod_{i=1}^N \frac{1}{\sigma(x_i)} \exp \left[-\frac{1}{2} \frac{\|d_i - F(x_i, W)\|^2}{\sigma^2(x_i)} \right] p(x_i)$$

Posons : $Err(W) = -\ln P(App, W)$

$$Err(W) = \frac{1}{2} \sum_{i=1}^N \frac{\|d_i - F(x_i, W)\|^2}{\sigma^2(x_i)} + \frac{1}{2} \sum_{i=1}^N \ln \sigma^2(x_i) + \frac{N}{2} \ln(2\pi) - \sum_{i=1}^N \ln p(x_i)$$

(4) (5) (6) (Const)

Méthode de maximum de vraisemblance.

Déterminer : $W^* = Arg \ Max_W P(App, W)$
 $= Arg \ Min_W Err(W)$

PMC (partie 1) 7 - Estimation
Variance

7.3 - Calcul de la vraisemblance (2)

7.3.1 - Cas où $\sigma^2(x) = \sigma^2 = \text{constante}$:

Dans ce cas, le terme (5) ne dépend pas de W .

On peut prendre : $Err(W) = \frac{1}{2} \sum_{i=1}^N \frac{\|d_i - F(x_i, W)\|^2}{\sigma^2}$

Ainsi : $Min_W Err(W) = Min_W \left(\sum_{i=1}^N \|d_i - F(x_i, W)\|^2 \right)$

On retrouve la somme des erreurs quadratiques.
 Soit W^* un système de poids qui minimise cette erreur.

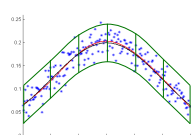
Estimation de σ^2 :

En retenant W^* , l'erreur Err devient :

$$Err(W^*) = \frac{1}{2} \sum_{i=1}^N \frac{\|d_i - F(x_i, W^*)\|^2}{\sigma^2} + N \ln \sigma$$

La valeur de σ^2 qui minimise $Err(W^*)$ est :

Erreur moyenne quadratique $\sigma^2 = \frac{1}{N} \sum_{i=1}^N \|d_i - F(x_i, W^*)\|^2$



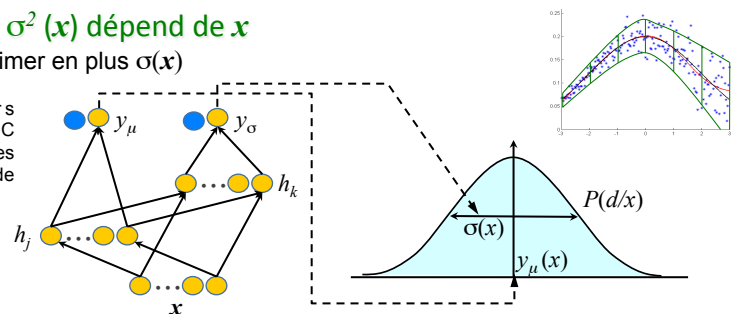
PMC (partie 1) 7 - Estimation
Variance

7.3 - Calcul de la vraisemblance (3)

7.3.2 - Cas où $\sigma^2(x)$ dépend de x

⇒ Estimer en plus $\sigma(x)$

On peut alors considérer un PMC avec deux neurones sur la couche de sortie :



- **1^{er} neurone**: entrée : A_μ fonction d'activation : Id sortie : $y_\mu = F_\mu(x, W)$
 $y_\mu = A_\mu \Rightarrow$ estime $\mu(x)$
- **2^{eme} neurone**: entrée : A_σ fonction d'activation : exp sortie : $y_\sigma = F_\sigma(x, W)$
 $y_\sigma = e^{A_\sigma} \Rightarrow$ estime $\sigma^2(x)$

• La fonction d'erreur à minimiser par rapport à W devient :

$$Err(W) = \frac{1}{2} \sum_{i=1}^N \frac{\|d_i - F_\mu(x_i, W)\|^2}{\sigma^2(x_i)} + \frac{1}{2} \sum_{i=1}^N \ln \sigma^2(x_i) \Rightarrow Err(W) = \frac{1}{2} \sum_{i=1}^N \frac{(d_i - A_\mu)^2}{e^{A_\sigma}} + \frac{1}{2} \sum_{i=1}^N A_\sigma$$

PMC (partie 1) 7 - Estimation
Variance

7.3 - Calcul de la vraisemblance (4)

7.3.3 - Algorithme neuronal

On pose :

$$Err_i(W) = \frac{1}{2} \frac{(d_i - A_\mu)^2}{e^{A_\sigma}} + \frac{1}{2} A_\sigma$$

Initialisation de la rétro-propagation du gradient :

- pour le neurone **1** : $\frac{\partial Err_i}{\partial A_\mu} = -\frac{(d_i - A_\mu)}{e^{A_\sigma}}$ (7)
- pour le neurone **2** : $\frac{\partial Err_i}{\partial A_\sigma} = -\frac{(d_i - A_\mu)^2}{e^{A_\sigma}} + \frac{1}{2}$ (8)

L' algorithme se présente alors en 3 phases ...

PMC (partie 1) 7 - Estimation
Variance

7.3 - Calcul de la vraisemblance (5)

7.3.3 - Algorithme neuronal (suite)

Phase 1 :
(Initialisation de la moyenne conditionnelle).

On suppose que $\sigma^2(x)$ est constante et on apprend l'erreur quadratique classique. On rétropropage alors uniquement à partir du neurone **1**.

(7)
$$\frac{\partial Err_i}{\partial A_\mu} = -\frac{(d_i - A_\mu)}{e^{A_\mu}}$$
 dans ce cas.

Soit W_1 les poids déterminés à la fin de cette phase.

Phase 2 :

On gèle la sortie du neurone **1** à :

$y_\mu = F_\mu(x; W_1)$ pour tout x_i .

On rétropropage alors uniquement à partir du neurone **2** en prenant :

(8)
$$\frac{\partial Err_i}{\partial A_\sigma} = -\frac{(d_i - A_\sigma)^2}{e^{A_\sigma}} + \frac{1}{2}$$

Soit W_2 les poids déterminés à la fin de cette phase.

Phase 3 :

Minimiser la fonction d'erreur complète en rétropropageant à partir des deux neurones **1** et **2**, en appliquant les formules (7) et (8).

