

# Reconnaissance de formes II

## Réduction des attributs – Réseaux de neurones

Yves Goussard

GBM6103A

26 novembre 2014

## 1 Réduction du nombre des attributs

- Position du problème
- Approche par minimum d'entropie
- Discriminant linéaire de Fisher
- Décomposition en composantes principales

## 2 Classification par réseaux de neurones

- Perceptron
- Réseaux à couches multiples

# Réduction du nombre d'attributs

## Position du problème

- Vecteurs d'attributs  $x$  de taille  $N$
- Si  $N$  est grand, lourdeur des calculs et possible détérioration des performances du classificateur

## Objectifs

- Réduction de la taille des attributs  $N \rightarrow M < N$
- Préservation, voire amélioration de la classification

## Cadre

- $l$  classes  $\omega_i$ ;  $1 \leq i \leq l$
- Pour chaque classe,  $m_i$  (moyenne) et  $R_i$  (matrice de covariance) connues
- Réduction :  $y = Tx$ ;  $y$  de taille  $M < N$ ;  $T$  : transformation linéaire

# Approche par minimum d'entropie (1)

## Approche

Choisir  $\mathcal{T}$  pour minimiser la dispersion dans chaque classe

## Entropie

- $f$  densité de probabilité ;  $H(f) = -E[\log f(x)] = -\int f(x) \log f(x) dx$
- Entropie  $H(f)$  : mesure la dispersion de  $f$
- $f$  uniforme de support de taille  $L$  :

$$H(f) = \log L$$

- $f$  gaussienne multivariée  $\mathcal{N}(\mathbf{m}, \mathbf{R})$  :

$$H(f) = \frac{N}{2} + \frac{1}{2} \log |\mathbf{R}| + \frac{N}{2} \log 2\pi$$

# Approche par minimum d'entropie (2)

## Hypothèses

- $f(\mathbf{x}|\omega_i)$  gaussienne  $\mathcal{N}(\mathbf{m}_i, \mathbf{R}_i)$
- Toutes les classes sont de covariance identique :  
 $\forall i \in \{1, 2 \dots I\} : \mathbf{R}_i = \mathbf{R}$

## Résultat

- Dans chaque classe :  $f(\mathbf{y}|\omega_i)$  gaussienne  $\mathcal{N}(\mathbf{T}\mathbf{m}_i, \mathbf{T}\mathbf{R}\mathbf{T}^t)$
- Après transformation :

$$H_i = \frac{M}{2}(1 + \log 2\pi) + \frac{1}{2} \sum_{j=1}^M \log \lambda_j \quad ; \quad \lambda_j : \text{val. propres de } \mathbf{T}\mathbf{R}\mathbf{T}^t$$

- Sous condition de normalisation,  $H_i$  minimisée pour  $\mathbf{T}$  composée des  $M$  vecteurs propres de  $\mathbf{R}$  associés aux valeurs propres les plus petites

# Discriminant linéaire de Fisher (1)

## Approche

Choisir  $T$  pour maximiser le rapport entre dispersion inter-classes et dispersion intra-classe

## Hypothèses

- Pour chaque classe,  $\{m_i, R_i\}$  connues

- $R = \sum_{i=1}^I R_i$        $m = \frac{1}{I} \sum_{i=1}^I m_i$

## Démarche

- Dispersion intra-classe :  $R$

- Dispersion inter-classes :  $B = \sum_{i=1}^I (m_i - m)(m_i - m)^t$

## Discriminant linéaire de Fisher (2)

### Démarche

- Rapport des dispersions avant transformation :  $|B|/|R|$
- Rapport des dispersions après transformation :  $|TBT^t|/|TRT^t|$
- Quantité maximisée pour  $T$  composée des  $M$  vecteurs propres de  $R^{-1}B$  associés aux valeurs propres les plus grandes

### Remarque

On doit choisir  $M < I$  pour que  $B$  soit inversible

# Décomposition en composantes principales

## Approche

- Meilleure approximation de  $\mathbf{x}$  par  $\mathbf{y}$  au sens de l'erreur quadratique moyenne
- Sélection de  $\mathbf{y}$  de manière à maximiser la dispersion *globale*

## Démarche

- Dispersion avant transformation :  $|\mathbf{R}|$
- Dispersion après transformation :  $|\mathbf{T}\mathbf{R}\mathbf{T}^t|$
- Sous condition de normalisation, quantité maximisée pour  $\mathbf{T}$  composée des  $M$  vecteurs propres de  $\mathbf{R}$  associés aux valeurs propres les plus grandes

## Remarques

- On utilise  $f(\mathbf{x})$  et non  $f(\mathbf{x}|\omega_i)$  !
- Pas d'effet garanti sur la séparation entre classes

# Classification par réseaux de neurones

## Objectif général

Mise en œuvre simple d'une méthode de classification par minimum de distance quadratique

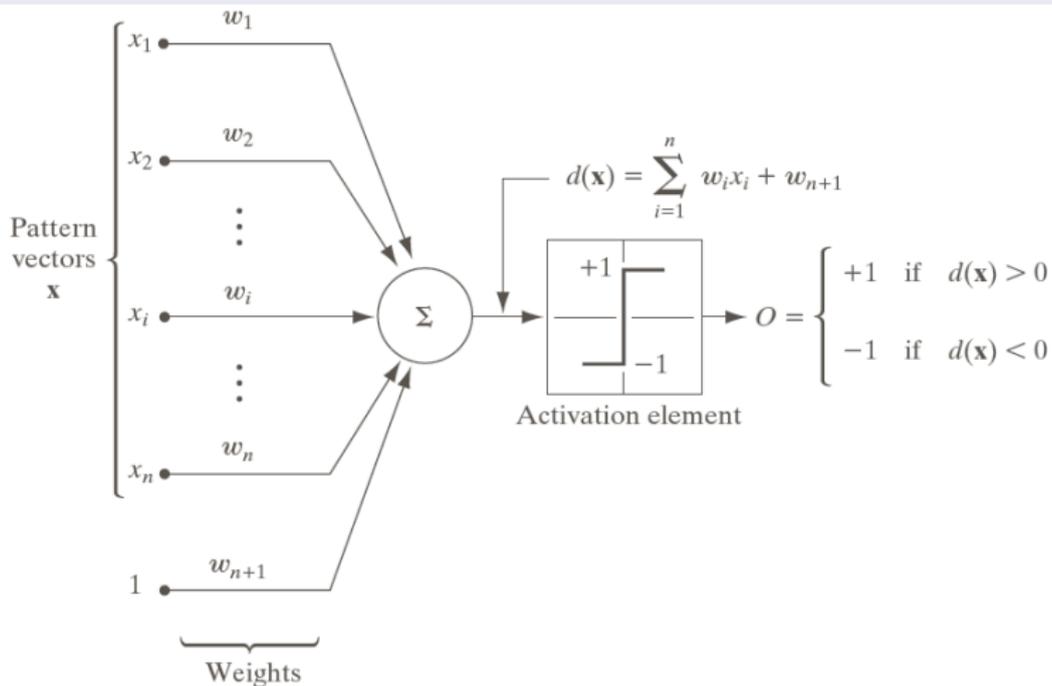
## Perceptron

- Deux classes
- Distance quadratique : surface de décision = hyperplan

Structure du classificateur  $w^t x \gtrsim t$

# Perceptron (1)

## Structure



©1992-2008 R. C. Gonzalez & R. E. Woods

# Perceptron (2)

## Formalisation et entraînement

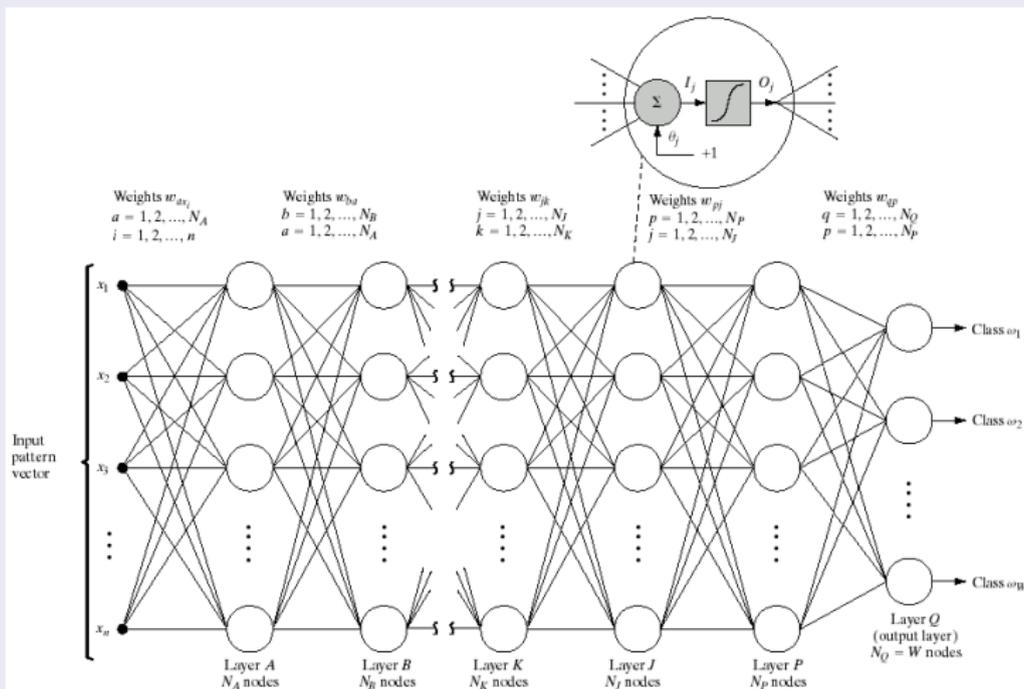
- Relation entrée-sortie :  $\mathbf{y}^t \mathbf{w} = r$
- Concaténation pour les éléments de l'ensemble d'apprentissage :  
 $\mathbf{Y}^t \mathbf{w} = \mathbf{r}$
- Estimation de  $\mathbf{w}$  par moindres carrés :  $\hat{\mathbf{w}} = (\mathbf{Y}^t \mathbf{Y})^{-1} \mathbf{Y}^t \mathbf{r}$
- Minimisation par algorithme du gradient :  
$$\hat{\mathbf{w}}^{(i+1)} = \hat{\mathbf{w}}^{(i)} + \alpha \mathbf{y}_{i+1} (r_{i+1} - \mathbf{y}_{i+1}^t \hat{\mathbf{w}}^{(i)})$$
- Divers choix possibles pour  $\alpha$ , selon la difficulté du problème

## Caractéristiques de l'approche

- Limitée en tant que telle
- « Brique de base » pour les réseaux de neurones

# Réseaux à couches multiples (1)

## Structure

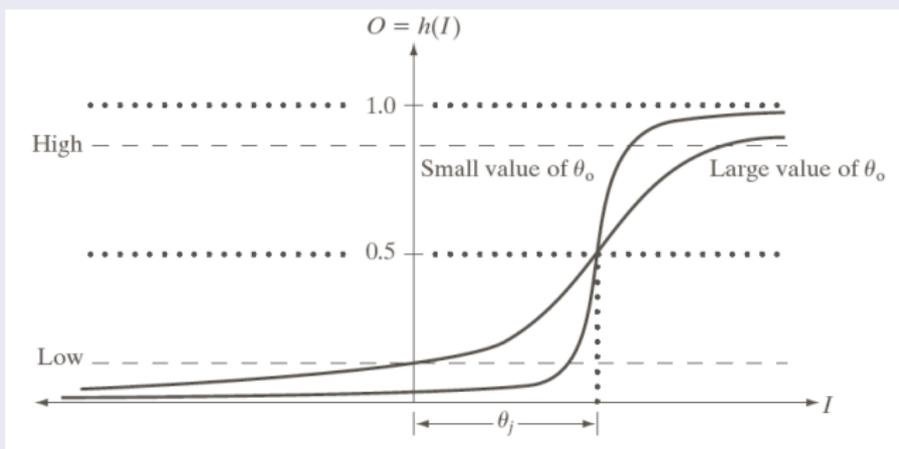


©1992-2008 R. C. Gonzalez & R. E. Woods

# Réseaux à couches multiples (2)

## Quantités à fixer

- Nombre de couches
- Nombre de cellules par couche
- Fonction seuil  $O = h(I)$



©1992-2008 R. C. Gonzalez & R. E. Woods

# Réseaux à couches multiples (3)

## Entraînement : algorithme de rétropropagation

- Entraînement du perceptron :

$$\Delta w_n = -\alpha \frac{\partial E}{\partial w_n} \quad E = \frac{1}{2} \left( r - \sum_{n=1}^{N+1} y_n w_n \right)^2$$

- Rétropropagation : extension par parcours du réseau de droite à gauche
- Dernière couche :

$$\Delta w_{qp} = \alpha \underbrace{(r_q - O_q) h'(I_q)}_{\delta_q} O_p \quad E_q = \frac{1}{2} \sum (r_q - O_q)^2$$

- Avant-dernière couche : on ne connaît pas la sortie désirée  $r_p$  !

$$E_p = \frac{1}{2} \sum_p (r_p - O_p)^2$$

$$\text{mais } \delta_p = \frac{\partial E_p}{\partial I_p} = h'(I_p) \sum_q \delta_q w_{qp} \quad \text{et} \quad \Delta w_{pj} = \alpha \delta_p O_j$$

# Réseaux à couches multiples (4)

## Avantages de l'approche

- Technique relativement simple
- Structure souple
- Fonctionnement de type boîte noire

## Limitations de l'approche

- Structure parfois difficile à spécifier
- Parfois difficile à entraîner (mauvaise sensibilité des paramètres internes, non-linéarités, algorithme lent)
- Fonctions de décision : hyperplans par morceaux