

Analyse Factorielle Discriminante

Machine Learning

Cédric RICHARD

Université Nice Sophia Antipolis

ANALYSE FACTORIELLE DISCRIMINANTE

Objectif

Contexte : Chaque individu \mathbf{x}_i du tableau \mathbf{X} est considéré comme un point d'un espace vectoriel \mathcal{E} de dimension p .

L'ensemble des n individus constitue un nuage de points dans \mathcal{E} .

On suppose que

— n_1 individus appartiennent à la classe ω_1

— n_2 individus appartiennent à la classe ω_2

tel que $n = n_1 + n_2$.

Objectif : On cherche à projeter les données sur un axe $\Delta(\mathbf{u})$ de vecteur directeur \mathbf{u} maximisant la *séparabilité* de ω_1 et ω_2 .

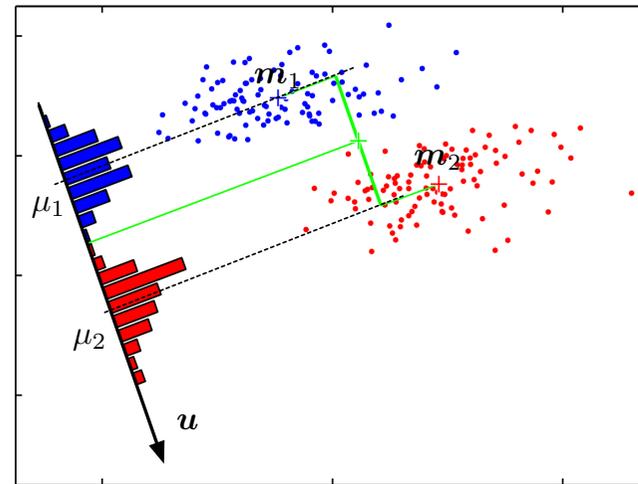
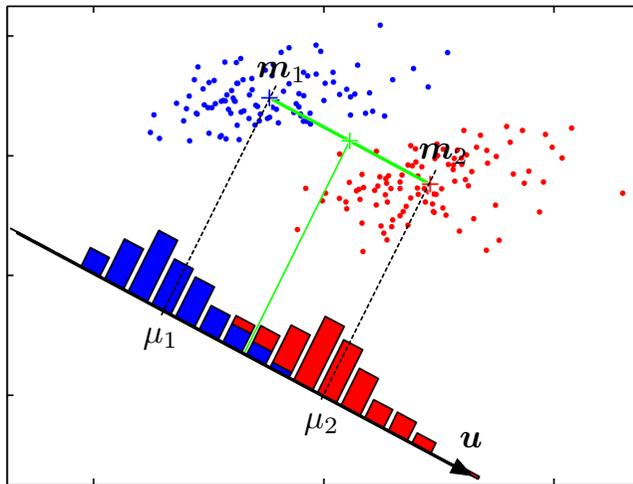
La règle de décision s'écrit

$$d(\mathbf{x}) = \begin{cases} \omega_2 & \text{si } \lambda(\mathbf{x}) = \mathbf{u}^\top \mathbf{x} > \lambda_0 \\ \omega_1 & \text{sinon} \end{cases}$$

ANALYSE FACTORIELLE DISCRIMINANTE

Objectif

Objectif : maximiser la séparabilité des données



ANALYSE FACTORIELLE DISCRIMINANTE

Mesure de séparabilité

Séparabilité : Afin de trouver un axe discriminant $\Delta(\mathbf{u})$, il convient de définir une mesure de séparabilité.

On note \mathbf{m}_1 et \mathbf{m}_2 les moyennes de chaque classe, soit

$$\mathbf{m}_1 = \frac{1}{n_1} \sum_{\mathbf{x} \in \omega_1} \mathbf{x} \quad \mathbf{m}_2 = \frac{1}{n_2} \sum_{\mathbf{x} \in \omega_2} \mathbf{x}$$

Solution (mauvaise) par l'exemple : On pourrait, par exemple, considérer l'écart quadratique entre les moyennes projetées de chaque classe

$$J(\mathbf{u}) = (\mathbf{u}^\top [\mathbf{m}_1 - \mathbf{m}_2])^2 = (\mu_1 - \mu_2)^2$$

Les figures précédentes montrent qu'il convient de prendre en compte la dispersion des classes.

ANALYSE FACTORIELLE DISCRIMINANTE

Discriminant de Fisher

Critère utilisé : Le critère de Fisher suggère de maximiser l'écart inter-classe, normalisé par une mesure de la dispersion intra-classe.

Le discriminant linéaire de Fisher est défini par $\lambda(\mathbf{x}) = \mathbf{u}^\top \mathbf{x}$, où \mathbf{u} maximise le critère dit de Fisher donné par

$$J(\mathbf{u}) = \frac{(\mu_1 - \mu_2)^2}{p_1 \sigma_1^2 + p_2 \sigma_2^2}$$

avec

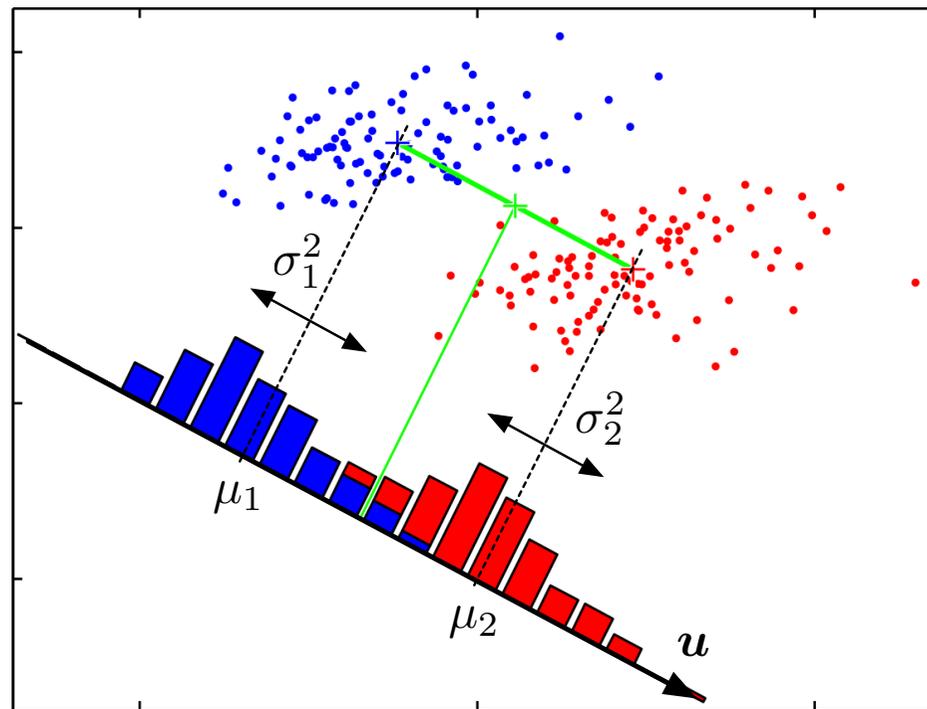
$$\mu_i = \mathbf{u}^\top \mathbf{m}_i$$

$$\sigma_i^2 = \frac{1}{n_i} \sum_{\mathbf{x} \in \omega_i} (\mathbf{u}^\top [\mathbf{x}_i - \mathbf{m}_i])^2$$

$$p_i = \frac{n_i}{n}$$

ANALYSE FACTORIELLE DISCRIMINANTE

Discriminant de Fisher



ANALYSE FACTORIELLE DISCRIMINANTE

Discriminant de Fisher

▷ On définit les quantités suivantes dans l'espace des individus

$$\Sigma_i = \frac{1}{n_i} \sum_{\mathbf{x} \in \omega_i} (\mathbf{x} - \mathbf{m}_i) (\mathbf{x} - \mathbf{m}_i)^\top$$

$$\Sigma_w = p_1 \Sigma_1 + p_2 \Sigma_2$$

La matrice Σ_w est appelée matrice de dispersion intra-classe

▷ On peut exprimer Σ_w en fonction de σ_1^2 et σ_2^2 . En effet

$$\begin{aligned} \sigma_i^2 &= \frac{1}{n_i} \sum_{\mathbf{x} \in \omega_i} \mathbf{u}^\top (\mathbf{x} - \mathbf{m}_i) (\mathbf{x} - \mathbf{m}_i)^\top \mathbf{u} \\ &= \mathbf{u}^\top \Sigma_i \mathbf{u} \end{aligned}$$

En conséquence, on a

$$p_1 \sigma_1^2 + p_2 \sigma_2^2 = \mathbf{u}^\top \Sigma_w \mathbf{u}$$

ANALYSE FACTORIELLE DISCRIMINANTE

Discriminant de Fisher

▷ On définit les quantités suivantes dans l'espace des individus

$$\Sigma_b = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^\top$$

La matrice Σ_b est appelée matrice de dispersion inter-classe

Remarque : la matrice Σ_b est de rang 1.

▷ On peut exprimer Σ_b en fonction de μ_1 et μ_2 . En effet

$$\begin{aligned}(\mu_1 - \mu_2)^2 &= [\mathbf{u}^\top (\mathbf{m}_1 - \mathbf{m}_2)]^2 \\ &= \mathbf{u}^\top \Sigma_b \mathbf{u}\end{aligned}$$

ANALYSE FACTORIELLE DISCRIMINANTE

Discriminant de Fisher

Définition : Le critère de Fisher s'exprime ainsi :

$$J(\mathbf{u}) = \frac{\mathbf{u}^\top \boldsymbol{\Sigma}_b \mathbf{u}}{\mathbf{u}^\top \boldsymbol{\Sigma}_w \mathbf{u}}$$

Maximisation : On calcule le gradient de $J(\mathbf{u})$, qu'on annule ensuite

$$\begin{aligned}\nabla J(\mathbf{u}) &= (\mathbf{u}^\top \boldsymbol{\Sigma}_w \mathbf{u}) \nabla \{\mathbf{u}^\top \boldsymbol{\Sigma}_b \mathbf{u}\} - (\mathbf{u}^\top \boldsymbol{\Sigma}_b \mathbf{u}) \nabla \{\mathbf{u}^\top \boldsymbol{\Sigma}_w \mathbf{u}\} \\ &= 2(\mathbf{u}^\top \boldsymbol{\Sigma}_w \mathbf{u}) \boldsymbol{\Sigma}_b \mathbf{u} - 2(\mathbf{u}^\top \boldsymbol{\Sigma}_b \mathbf{u}) \boldsymbol{\Sigma}_w \mathbf{u}\end{aligned}$$

$$\nabla J(\mathbf{u}) = 0 \quad \Rightarrow \quad (\mathbf{u}^\top \boldsymbol{\Sigma}_w \mathbf{u}) \boldsymbol{\Sigma}_b \mathbf{u} = (\mathbf{u}^\top \boldsymbol{\Sigma}_b \mathbf{u}) \boldsymbol{\Sigma}_w \mathbf{u}$$

En divisant chaque membre par $(\mathbf{u}^\top \boldsymbol{\Sigma}_w \mathbf{u})$, on aboutit à

$$\boxed{\boldsymbol{\Sigma}_w^{-1} \boldsymbol{\Sigma}_b \mathbf{u} = J(\mathbf{u}) \mathbf{u}}$$

ANALYSE FACTORIELLE DISCRIMINANTE

Discriminant de Fisher

Problème : Pour maximiser $J(\mathbf{u})$, il convient de résoudre

$$\Sigma_w^{-1} \Sigma_b \mathbf{u} = J(\mathbf{u}) \mathbf{u}$$

▷ On sait que $\Sigma_b = (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^\top$. En remplaçant ci-dessus, on obtient

$$J(\mathbf{u}) \mathbf{u} = \Sigma_w^{-1} (\mathbf{m}_1 - \mathbf{m}_2) \underbrace{(\mathbf{m}_1 - \mathbf{m}_2)^\top \mathbf{u}}_{\text{constante}}$$

▷ A une constante près, ceci implique que :

$$\boxed{\mathbf{u} = \Sigma_w^{-1} (\mathbf{m}_1 - \mathbf{m}_2)}$$

Le discriminant $\lambda(\mathbf{x}) = \mathbf{u}^\top \mathbf{x}$ correspondant est dit *discriminant de Fisher* (1936).

ANALYSE FACTORIELLE DISCRIMINANTE

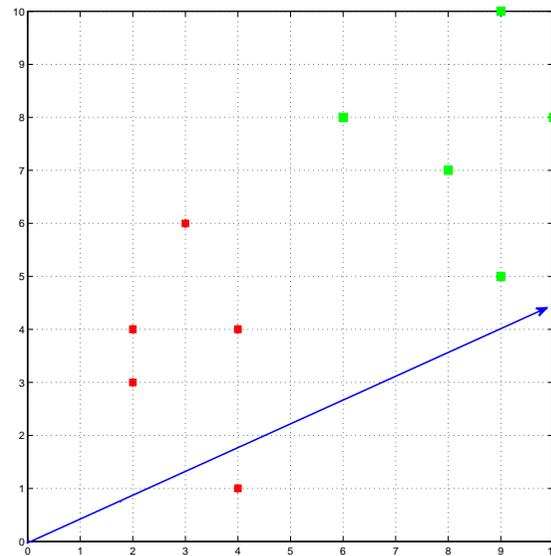
Exemple

Données : On considère les deux classes d'individus dans \mathbb{R}^2 suivantes :

$$\omega_1 = \{(4, 1); (2, 4); (2, 3); (3, 6); (4, 4)\}$$

$$\omega_2 = \{(9, 10); (6, 8); (9, 5); (8, 7); (10, 8)\}$$

$$\Rightarrow \quad \mathbf{u} = (1.76, 0.75)^\top$$



CRITÈRES DU SECOND ORDRE

Taxonomie

Les critères *du second ordre*, tels que le critère de Fisher se contentent d'une quantité d'information modeste, limitée aux 2 premiers moments de $\lambda(\mathbf{x})$:

$$\begin{aligned}\mu_i &= \frac{1}{n_i} \sum_{\mathbf{x} \in \omega_i} \lambda(\mathbf{x}) = \mathbf{u}^\top \mathbf{m}_i - \lambda_0 \\ \sigma_i^2 &= \frac{1}{n_i} \sum_{\mathbf{x} \in \omega_i} (\lambda(\mathbf{x}) - \mu_i)^2 = \mathbf{u}^\top \Sigma_i \mathbf{u}\end{aligned}$$

Exemple : le rapport signal-sur-bruit généralisé

$$J_{\text{rsbg}}(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2) = \frac{(\mu_1 - \mu_2)^2}{\rho \sigma_1^2 + (1 - \rho) \sigma_2^2}.$$

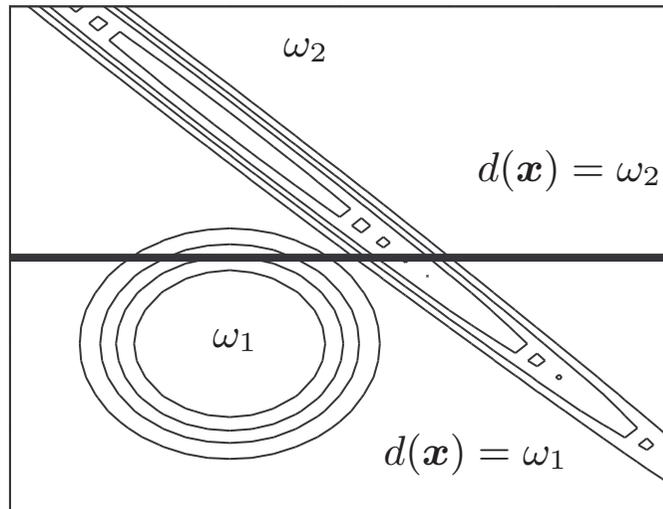
Il se décline suivant le critère de Fisher ($\rho = p_1$), la déflexion ($\rho = 1/2$) et le rapport signal-sur-bruit ($\rho = 1$).

CRITÈRES DU SECOND ORDRE

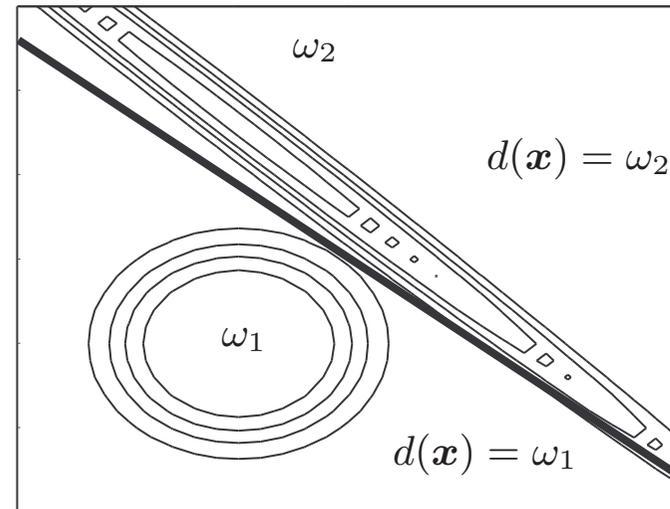
Choix optimum du critère

Le choix d'un critère du second ordre constitue un problème crucial car le paramètre ρ influe largement sur les performances :

rapport signal-sur-bruit ($\rho = 1$)



déflexion ($\rho = 1/2$)



CRITÈRES DU SECOND ORDRE

Choix optimum du critère

Soit $\lambda(\mathbf{x})$ la statistique linéaire définie par $\lambda(\mathbf{x}) = \mathbf{u}^\top \mathbf{x} - \lambda_0$. Soit $J(\eta_1, \eta_2, \sigma_1^2, \sigma_2^2)$ un critère du second ordre, ne dépendant donc que des variables suivantes :

$$\begin{aligned}\mu_i &= \frac{1}{n_i} \sum_{\mathbf{x} \in \omega_i} \lambda(\mathbf{x}) = \mathbf{u}^\top \mathbf{m}_i - \lambda_0 \\ \sigma_i^2 &= \frac{1}{n_i} \sum_{\mathbf{x} \in \omega_i} (\lambda(\mathbf{x}) - \mu_i)^2 = \mathbf{u}^\top \Sigma_i \mathbf{u}\end{aligned}$$

avec $\mathbf{m}_i = \frac{1}{n_i} \sum_{\mathbf{x} \in \omega_i} \mathbf{x}$ et $\Sigma_i = \frac{1}{n_i} \sum_{\mathbf{x} \in \omega_i} (\mathbf{x} - \mathbf{m}_i) (\mathbf{x} - \mathbf{m}_i)^\top$.

Objectif :

Rechercher \mathbf{u} et λ_0 maximisant $J(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2)$ sans expression du critère

CRITÈRES DU SECOND ORDRE

Choix optimum du critère

Les dérivés partielles de J par rapport à \mathbf{u} et λ_0 doivent être nulles, c'est-à-dire

$$\begin{cases} \frac{\partial J}{\partial \mathbf{u}} = \frac{\partial J}{\partial \sigma_1^2} \cdot \frac{\partial \sigma_1^2}{\partial \mathbf{u}} + \frac{\partial J}{\partial \sigma_2^2} \cdot \frac{\partial \sigma_2^2}{\partial \mathbf{u}} + \frac{\partial J}{\partial \mu_1} \cdot \frac{\partial \mu_1}{\partial \mathbf{u}} + \frac{\partial J}{\partial \mu_2} \cdot \frac{\partial \mu_2}{\partial \mathbf{u}} = 0 \\ \frac{\partial J}{\partial \lambda_0} = \frac{\partial J}{\partial \sigma_1^2} \cdot \frac{\partial \sigma_1^2}{\partial \lambda_0} + \frac{\partial J}{\partial \sigma_2^2} \cdot \frac{\partial \sigma_2^2}{\partial \lambda_0} + \frac{\partial J}{\partial \mu_1} \cdot \frac{\partial \mu_1}{\partial \lambda_0} + \frac{\partial J}{\partial \mu_2} \cdot \frac{\partial \mu_2}{\partial \lambda_0} = 0, \end{cases}$$

où les dérivées partielles de μ_i et σ_i^2 par rapport à \mathbf{u} et λ_0 sont données par

$$\frac{\partial \sigma_i^2}{\partial \mathbf{u}} = 2 \Sigma_i \mathbf{u}, \quad \frac{\partial \mu_i}{\partial \mathbf{u}} = \mathbf{m}_i, \quad \frac{\partial \sigma_i^2}{\partial \lambda_0} = 0, \quad \frac{\partial \mu_i}{\partial \lambda_0} = -1.$$

Ces résultats permettent de réécrire le système précédent ainsi

$$\begin{cases} 2 \left[\frac{\partial J}{\partial \sigma_1^2} \Sigma_1 + \frac{\partial J}{\partial \sigma_2^2} \Sigma_2 \right] \mathbf{u} = - \left[\frac{\partial J}{\partial \mu_1} \mathbf{m}_1 + \frac{\partial J}{\partial \mu_2} \mathbf{m}_2 \right] \\ \frac{\partial J}{\partial \mu_1} + \frac{\partial J}{\partial \mu_2} = 0. \end{cases}$$

CRITÈRES DU SECOND ORDRE

Choix optimum du critère

En introduisant la deuxième équation dans la première, et en notant que \mathbf{u} peut être défini à un coefficient multiplicatif près, on aboutit au système linéaire

$$[\rho \Sigma_1 + (1 - \rho) \Sigma_2] \mathbf{u} = (\mathbf{m}_2 - \mathbf{m}_1),$$

avec

$$\rho = \frac{\frac{\partial J}{\partial \sigma_1^2}}{\frac{\partial J}{\partial \sigma_1^2} + \frac{\partial J}{\partial \sigma_2^2}}.$$

Remarque : On a $0 \leq \rho \leq 1$ si $J(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2)$ varie dans le même sens par rapport à σ_1^2 et σ_2^2 (dérivées de même signe). Cette condition est raisonnable.

Conclusion. Ce résultat est particulièrement intéressant puisqu'il établit que l'expression de $J(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2)$ n'intervient dans \mathbf{u} que par l'intermédiaire de ρ .

CRITÈRES DU SECOND ORDRE

Choix optimum du critère

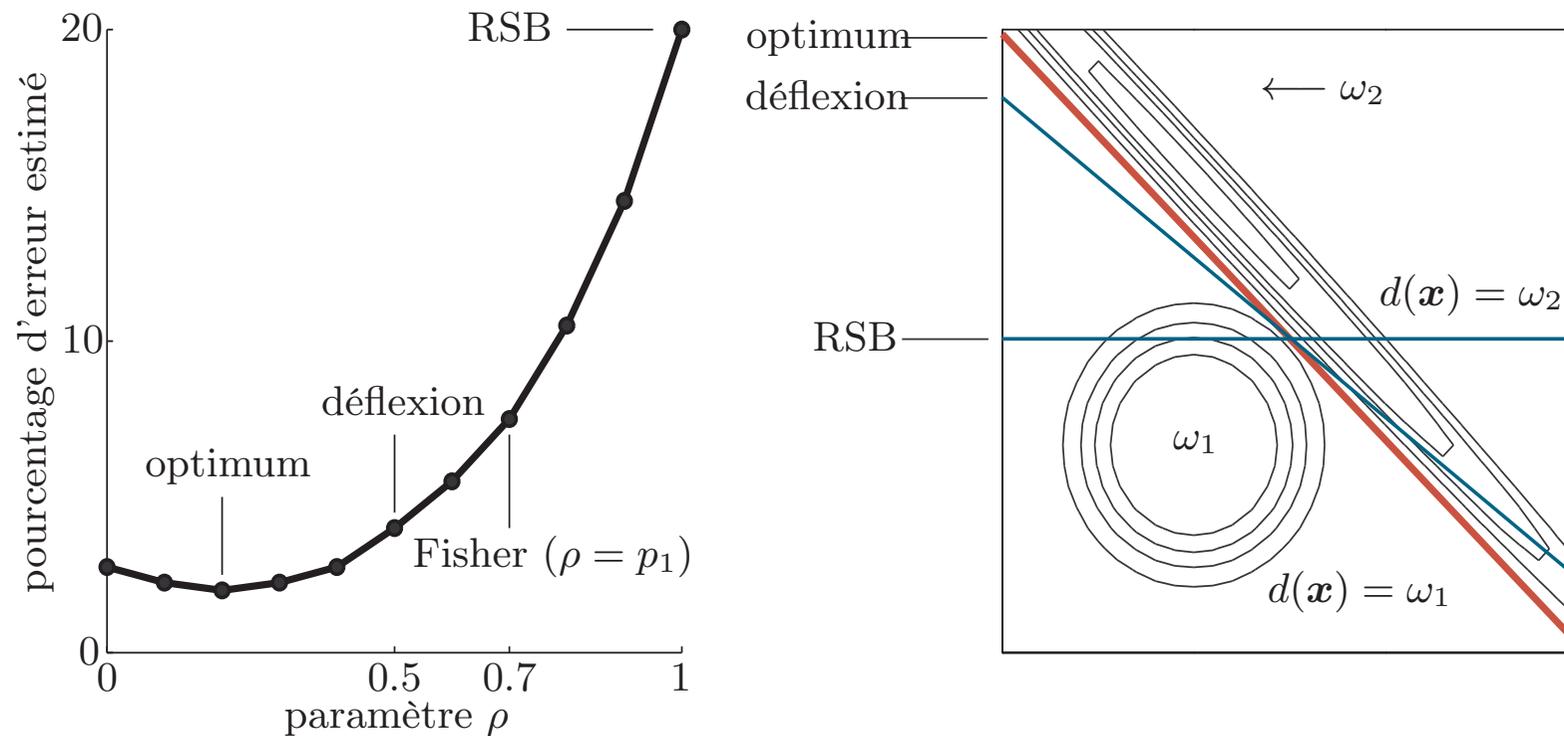
Algorithme de la méthode

1. Initialiser ρ à 0
 2. Tant que $\rho \leq 1$, répéter
 - ▷ résoudre $[\rho\Sigma_1 + (1 - \rho)\Sigma_2] \mathbf{u} = (\mathbf{m}_2 - \mathbf{m}_1)$ pour obtenir \mathbf{u}_ρ
 - ▷ déterminer le seuil $\lambda_{\rho,0}$ minimisant par exemple $P_e(d_\rho)$
 - ▷ mise à jour de ρ : $\rho \leftarrow \rho + \Delta\rho$, avec $\Delta\rho$ préalablement choisi
 3. Sélectionner le meilleur détecteur d_ρ obtenu
-

OPTIMISATION DES CRITÈRES DU SECOND ORDRE

Choix optimum du critère

Le classifieur obtenu par la méthode précédente est au moins aussi performant que ceux résultant de la maximisation du rapport signal-sur-bruit, etc. En effet, il leur correspond à chacun une valeur particulière de ρ .



OPTIMISATION DES CRITÈRES DU SECOND ORDRE

Choix optimum du critère

Le calcul du paramètre ρ correspondant au rapport signal-sur-bruit conduit au résultat connu suivant :

$$\mathbf{u}_{\text{rsb}} = \boldsymbol{\Sigma}_1^{-1}(\mathbf{m}_2 - \mathbf{m}_1),$$

car $\rho_{\text{rsb}} = 1$. De façon analogue, on établit pour la déflexion et le critère de Fisher :

$$\mathbf{u}_{\text{deflex}} = 2(\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2)^{-1}(\mathbf{m}_2 - \mathbf{m}_1)$$

$$\mathbf{u}_{\text{Fisher}} = (p_1 \boldsymbol{\Sigma}_1 + p_2 \boldsymbol{\Sigma}_2)^{-1}(\mathbf{m}_2 - \mathbf{m}_1),$$

car $\rho_{\text{deflex}} = 1/2$ et $\rho_{\text{Fisher}} = p_1$.

ANALYSE FACTORIELLE DISCRIMINANTE

Cas multi-classe

Généralisation du critère de Fisher à C classes :

On cherche $(C - 1)$ projections $\lambda_i(\mathbf{x}) = \mathbf{u}_i^\top \mathbf{x}$ à partir de $(C - 1)$ vecteurs \mathbf{u}_i

On range les vecteurs \mathbf{u}_i dans la matrice \mathbf{U} , soit $\boldsymbol{\lambda}(\mathbf{x}) = \mathbf{U}^\top \mathbf{x}$.

Notation :

On généralise la dispersion intra-classe comme ceci

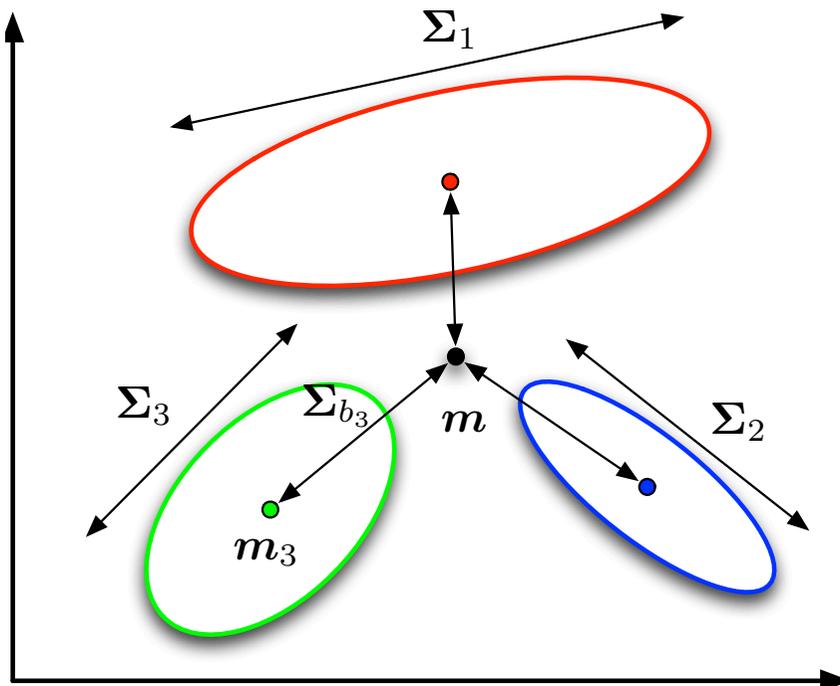
$$\boldsymbol{\Sigma}_w = \sum_{i=1}^C \boldsymbol{\Sigma}_i$$

et la dispersion inter-classe comme cela

$$\boldsymbol{\Sigma}_b = \frac{1}{n} \sum_{i=1}^C n_i (\mathbf{m}_i - \mathbf{m}) (\mathbf{m}_i - \mathbf{m})^\top$$

ANALYSE FACTORIELLE DISCRIMINANTE

Cas multi-classe



ANALYSE FACTORIELLE DISCRIMINANTE

Cas multi-classe

Définition : Dans le cas multi-classe, le critère de Fisher s'exprime ainsi

$$J(\mathbf{U}) = \frac{|\mathbf{U}^\top \boldsymbol{\Sigma}_b \mathbf{U}|}{|\mathbf{U}^\top \boldsymbol{\Sigma}_w \mathbf{U}|}$$

Maximisation : La solution $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_{C-1}]$ du problème est obtenue en résolvant le problème aux valeurs propres généralisé suivant

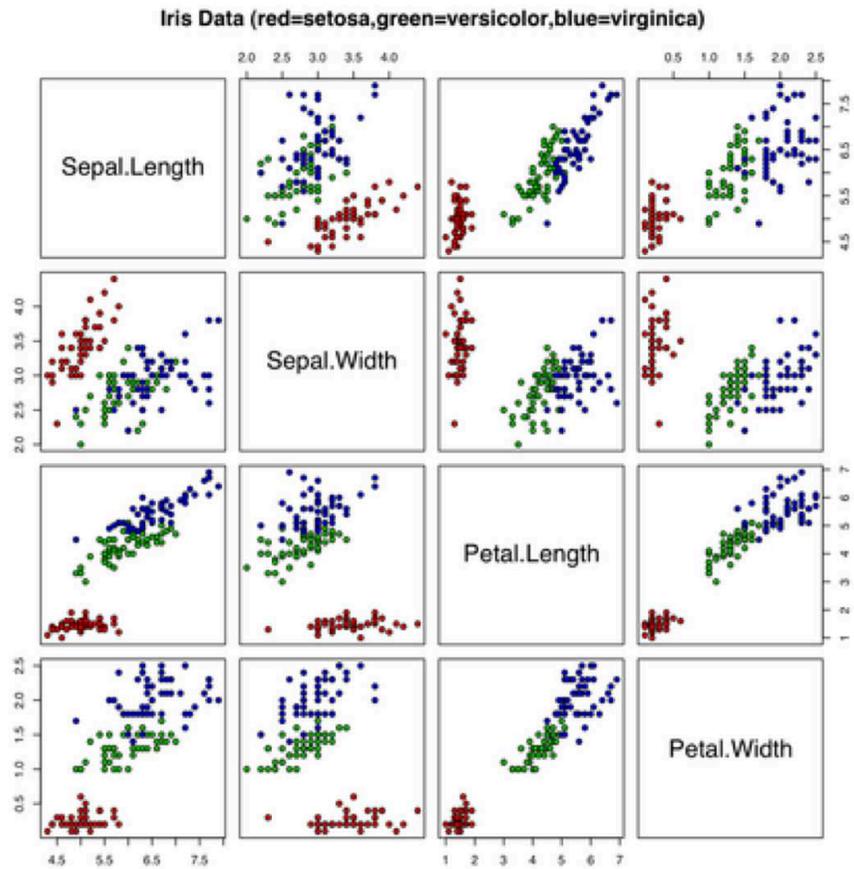
$$\boxed{(\boldsymbol{\Sigma}_b - \lambda_i \boldsymbol{\Sigma}_w) \mathbf{u}_i = 0}$$

pour $i = 1, \dots, C - 1$.

Notes : Les directions de projection \mathbf{u}_i correspondent aux vecteurs propres associés aux plus grandes valeurs propres de $\boldsymbol{\Sigma}_w^{-1} \boldsymbol{\Sigma}_b$. Au plus $C - 1$ valeurs propres sont non-nulles car $\text{rang}(\boldsymbol{\Sigma}_b) = C - 1$.

ANALYSE FACTORIELLE DISCRIMINANTE

Cas multi-classe



ANALYSE FACTORIELLE DISCRIMINANTE

Cas multi-classe

LDA: iris projection onto the first 2 linear discriminants

