

Notes de cours

Analyse Discriminante Linéaire Analyse Discriminante Quadratique

On dispose d'une population composée de k sous-populations définies par une variable catégorielle à k modalités que l'on note \mathbf{y} . Quitte à renommer les classes, on suppose que l'ensemble des modalités est $\{1, \dots, k\}$. Chaque individu de la population est aussi décrit par des variables explicatives $\mathbf{x}^1, \dots, \mathbf{x}^j, \dots, \mathbf{x}^p$.

1 Nuages et décomposition de la variance

Dans cette section, on suppose que les variables explicatives $\mathbf{x}^1, \dots, \mathbf{x}^j, \dots, \mathbf{x}^p$ sont toutes de type continu. Soit \mathbf{M} la matrice associée au nuage des individus mesurés sur chacune des p variables explicatives :

$$\mathbf{M} = \begin{bmatrix} x_1^1 & \dots & x_1^p \\ \vdots & \vdots & \vdots \\ x_n^1 & \dots & x_n^p \end{bmatrix}.$$

On appelle nuage de \mathbb{R}^p le n -uplet des observations $(\mathbf{x}_1, \dots, \mathbf{x}_n)$, où chacun des vecteurs \mathbf{x}_i de \mathbb{R}^p est un point du nuage. Par abus on notera aussi \mathbf{M} le nuage. Dans la suite on supposera toujours que le nuage est centré :

$$\left(\frac{1}{n} \sum_{i=1 \dots n} x_i^j \right)_{j=1 \dots p} = 0.$$

La matrice de variance-covariance des variables du nuage \mathbf{M} est définie par :

$$\mathbf{S} := \left[\text{cov}(\mathbf{x}^j, \mathbf{x}^{j'}) \right]_{1 \leq j, j' \leq p} = \frac{1}{n} \mathbf{M}' \mathbf{M} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i'$$

1.1 Géométrie d'un seul nuage

Avant de décrire la géométrie des k nuages, considérons le cas d'un seul nuage en laissant de côté pour le moment le problème de la classification supervisée. L'analyse en composantes principales (voir le cours dédié) nous dit que l'orientation du nuage \mathbf{M} peut être décrite à partir de la décomposition en valeurs propres de la matrice \mathbf{S} . Plus précisément, la première direction propre (pour la plus grande valeur propre) correspond à la direction de \mathbb{R}^p pour laquelle les données du nuage \mathbf{M} sont le plus dispersées. La seconde direction propre donne la direction orthogonale à la première pour laquelle les données sont le plus dispersées, etc...

1.2 Géométrie des k nuages

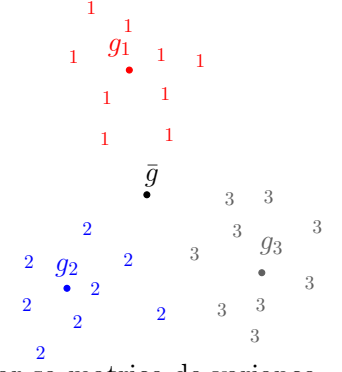
Nous supposons désormais que k nuages sont observés : la variable y présente k modalités distinctes. Pour simplifier les notations, les individus sont numérotés de telle sorte que les premiers individus sont ceux du groupe 1, puis ceux du groupe 2 :

$$\mathbf{M} = \begin{bmatrix} x_1^1 & \dots & x_1^p \\ \vdots & \vdots & \vdots \\ x_n^1 & \dots & x_n^p \end{bmatrix} = \begin{bmatrix} \mathbf{M}(1) \\ \vdots \\ \mathbf{M}(k) \end{bmatrix}$$

où $\mathbf{M}(\ell)$ est la matrice des données du groupe ℓ . Les k groupes induits par la variable catégorielle \mathbf{y} définissent k nuages dans l'espace \mathbb{R}^p des variables explicatives. On supposera de plus que $k \leq p < n$. On note :

- I_ℓ la liste des individus du groupe ℓ ,
- n_ℓ l'effectif de la sous-population ℓ ,
- \mathbf{g}_ℓ le centre de gravité du groupe ℓ : $\mathbf{g}_\ell = (g_\ell^1, \dots, g_\ell^p)'$ avec

$$g_\ell^j = \frac{1}{n_\ell} \sum_{i \in I_\ell} x_i^j.$$



Chacun des k nuages $\mathbf{M}(\ell)$ a une orientation \mathbb{R}^p qui peut être décrite par sa matrice de variance-covariance \mathbf{S}_ℓ (c.f. paragraphe précédent) où

$$\begin{aligned} \mathbf{S}_\ell &:= \frac{1}{n_\ell} (\mathbf{M}(\ell) - \mathbf{e}_{n_\ell} \mathbf{g}_\ell')' (\mathbf{M}(\ell) - \mathbf{e}_{n_\ell} \mathbf{g}_\ell') \\ &= \frac{1}{n_\ell} \sum_{i \in I_\ell} (\mathbf{x}_i - \mathbf{g}_\ell)(\mathbf{x}_i - \mathbf{g}_\ell)'. \end{aligned}$$

où \mathbf{e}_k désigne le vecteur $(1, \dots, 1)'$ de \mathbb{R}^k . On définit alors la matrice de variance **intra-classe** par

$$\mathbf{W} := \sum_{\ell=1 \dots k} \frac{n_\ell}{n} \mathbf{S}_\ell \quad p \times p$$

Cette matrice est la moyenne pondérée des matrices de variance-covariance des k nuages. Elle correspond donc à “l'orientation moyenne” des k nuages. La première direction principale de \mathbf{W} correspond à la direction selon laquelle les nuages sont le plus étalés, etc ... En général, la matrice \mathbf{W} est **inversible**, ce que l'on supposera dans la suite.

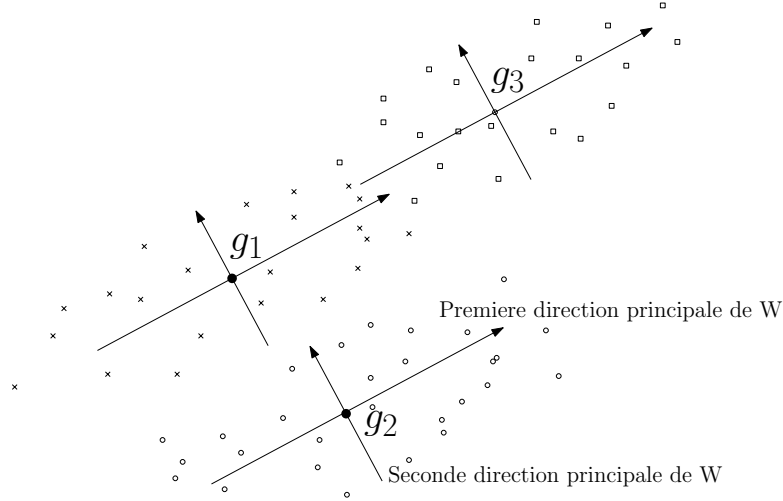


FIGURE 1 – Trois nuages avec des orientations comparables.

Soit \mathbf{G} la matrice $k \times p$ du nuage des centres des k nuages. Le centre de gravité $\bar{\mathbf{g}}$ de \mathbf{G} pondéré des n_ℓ vérifie

$$\bar{\mathbf{g}} := \frac{1}{n} \sum_{\ell=1 \dots k} n_\ell \mathbf{g}_\ell = \bar{\mathbf{x}} = 0.$$

La matrice de variance-covariance \mathbf{B} des k centres de gravité pondérés par les n_ℓ est appelée matrice

de variance **inter-classes** :

$$\begin{aligned}
\mathbf{B} &:= \left(\frac{1}{n} \sum_{\ell=1\dots k} n_{\ell} (g_{\ell}^j - \bar{g}^j)(g_{\ell}^u - \bar{g}^u) \right)_{1 \leq j, u \leq k} \\
&= \mathbf{G}' \text{Diag}\left(\frac{n_1}{n}, \dots, \frac{n_k}{n}\right) \mathbf{G} \quad \text{car } \bar{g}^j = \bar{g}^u = 0 \\
&= \frac{1}{n} \sum_{\ell=1\dots k} n_{\ell} \mathbf{g}_{\ell} \mathbf{g}_{\ell}' \quad p \times p
\end{aligned}$$

Puisque $\sum_{\ell=1\dots k} n_{\ell} \mathbf{g}_{\ell} = 0$, le nuage des \mathbf{g}_{ℓ} est contenu dans un sous-espace vectoriel de \mathbb{R}^p de dimension $k - 1$ et la matrice \mathbf{B} n'est pas inversible. La matrice \mathbf{B} décrit la géométrie du nuage des k centres de gravité.

Les variances intra et inter classes permettent de décomposer la variance totale du nuage :

Proposition 1. *La variance du nuage \mathbf{M} se décompose en variance intra-classe et variance inter-classes :*

$$\mathbf{S} = \mathbf{W} + \mathbf{B}.$$

2 Analyses discriminantes linéaire et quadratique

On suppose dans cette section que les variables aléatoires X^j sont toutes de type continu.

2.1 Modélisation

On suppose que les variables en jeu sont des variables aléatoires notées Y et X^j , dont on observe des réalisations x_1^j, \dots, x_n^j et y_1, \dots, y_n . Supposons de plus que la distribution de Y admet $k \geq 2$ modalités. Pour chacune de ces modalités on considère la loi conditionnelle de $X = (X^1, \dots, X^p)$ sachant $Y = \ell$ et on suppose que cette loi conditionnelle admet une densité f_{ℓ} pour la mesure de Lebesgue sur \mathbb{R}^p (que l'on note ici λ^p). On considère

- $\pi_{\ell} = P(Y = \ell)$: la probabilité a priori d'appartenance au groupe ℓ .
- $P(Y = \ell | X = x)$: la probabilité a posteriori d'appartenance au groupe ℓ .

Proposition 2. *Sous les hypothèses précédentes :*

1. *La distribution du vecteur aléatoire (X, Y) admet la densité*

$$\begin{aligned}
f &: \mathbb{R}^p \times \{1, \dots, k\} \rightarrow \mathbb{R}^+ \\
(\mathbf{x}, y) &\mapsto \sum_{\ell=1}^k f_{\ell}(\mathbf{x}) \pi_{\ell} \mathbf{1}_{y=\ell}
\end{aligned}$$

par rapport à la mesure $\lambda^p \otimes \delta_k$, où δ_k désigne la mesure ponctuelle d'atomes $\{1, \dots, k\}$.

2. *La distribution du vecteur X admet la densité suivante par rapport à λ^p :*

$$f_X = \sum_{\ell=1}^k \pi_{\ell} f_{\ell}.$$

3. *Pour tout $\ell \in \{1, \dots, k\}$, la probabilité a posteriori d'appartenance au groupe ℓ vérifie :*

$$P(Y = \ell | X = x) = \frac{\pi_{\ell} f_{\ell}(x)}{\sum_{s=1}^k \pi_s f_s(x)} \quad (1)$$

Les méthodes d'analyse discriminante linéaire et quadratique sont des méthodes d'analyse discriminantes de type décisionnel. Pour prédire la variable Y à partir des variables X^1, \dots, X^p , il est naturel de s'appuyer sur les probabilités a posteriori. Plus précisément, la **règle bayésienne d'attribution** consiste à attribuer une observation au groupe le plus probable pour celle-ci, c'est-à-dire celui pour lequel la probabilité a posteriori est maximale, ce qui équivaut d'après la relation (1) à choisir

$$\hat{Y}(\mathbf{x}) = \underset{\ell=1\dots k}{\text{Argmax}} f_{\ell}(\mathbf{x}) \pi_{\ell}.$$

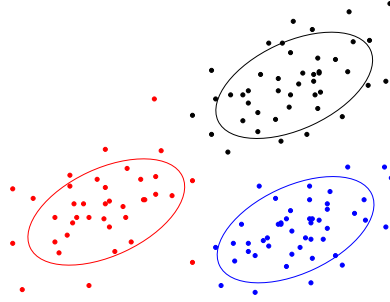
Cependant, en pratique ces quantités sont inconnues et il faut les estimer à partir des observations disponibles. Pour cela on proposera différentes hypothèses de modélisation sur la loi de X sachant Y .

Hypothèse gaussienne. Pour modéliser le fait que les observations de chaque groupe ℓ sont organisés en “clusters”, nous allons supposer que la loi du vecteur $X = (X^1, \dots, X^p)$ peut être modélisée par une loi normale multivariée de densité sur \mathbb{R}^p :

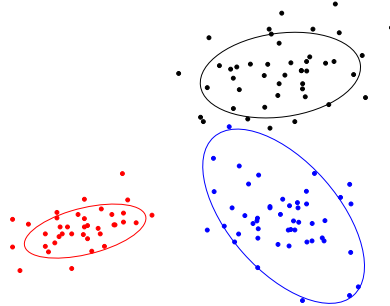
$$\mathbf{x} \in \mathbb{R}^p, \quad f_\ell(\mathbf{x}) = \frac{1}{(2\pi)^{p/2} \sqrt{\det \Sigma_\ell}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_\ell)' \Sigma_\ell^{-1} (\mathbf{x} - \boldsymbol{\mu}_\ell) \right].$$

Comme dans le cas de la régression linéaire, cette hypothèse n’est jamais rigoureusement vérifiée pour des données réelles. Cependant, cette modélisation est souvent suffisamment souple pour approcher efficacement la véritable loi des données (que l’on ne connaît évidemment pas en pratique).

ADL et ADQ. Si toutes les matrices de variance-covariance sont égales : $\Sigma_1 = \dots = \Sigma_k = \Sigma$, il s’agit de l’analyse discriminante linéaire (ADL) :



Au contraire, si les matrices de variance-covariance ne sont pas supposées égales, il s’agit de l’analyse discriminante quadratique (ADQ) :



2.2 Analyse Discriminante Linéaire (ADL)

Dans le but de déterminer les frontières entre les zones d’attribution, on considère le logarithme des rapports entre probabilités a posteriori : $\mathbf{x} \in \mathbb{R}^p$

$$\begin{aligned} \log \frac{P(Y = \ell \mid X = \mathbf{x})}{P(Y = h \mid X = \mathbf{x})} &= \log \frac{\pi_\ell}{\pi_h} + \log \frac{f_\ell(\mathbf{x})}{f_h(\mathbf{x})} \\ &= \log \frac{\pi_\ell}{\pi_h} - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_\ell)' \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_\ell) + \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_h)' \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_h) \quad (2) \\ &= \log \frac{\pi_\ell}{\pi_h} - \frac{1}{2} \boldsymbol{\mu}_\ell' \Sigma^{-1} \boldsymbol{\mu}_\ell + \frac{1}{2} \boldsymbol{\mu}_h' \Sigma^{-1} \boldsymbol{\mu}_h + \mathbf{x}' \Sigma^{-1} (\boldsymbol{\mu}_\ell - \boldsymbol{\mu}_h) \end{aligned}$$

On définit une fonction discriminante s_ℓ pour chaque groupe ℓ par :

$$\mathbf{x} \in \mathbb{R}^p, \quad s_\ell(\mathbf{x}) := \mathbf{x}' \Sigma^{-1} \boldsymbol{\mu}_\ell + \log \pi_\ell - \frac{1}{2} \boldsymbol{\mu}_\ell' \Sigma^{-1} \boldsymbol{\mu}_\ell.$$

La règle de Bayes consiste ici à affecter une observation \mathbf{x} au groupe $y(\mathbf{x})$ de score maximum :

$$y(\mathbf{x}) = \operatorname{argmax}_{\ell=1 \dots k} s_\ell(\mathbf{x}).$$

Ces fonctions discriminantes (ou scores) sont linéaires en \mathbf{x} , d’où l’appellation d’analyse discriminante linéaire.

Inférence. En pratique, les quantités $\Sigma, \mu_1, \dots, \mu_k$ et π_1, \dots, π_k sont inconnues. On peut cependant les estimer par la méthode du maximum de vraisemblance :

Proposition 3. *En supposant que la loi conditionnelle de $(X|Y = \ell)$ est celle d'une loi normale multivariée p -dimensionnel dont la matrice de variance-covariance est inversible et ne dépend pas de ℓ (hypothèse ADL), les estimateurs du maximum de vraisemblance de Σ, μ_ℓ et π_ℓ vérifient :*

$$\hat{\Sigma} = \mathbf{W} = \sum_{\ell=1 \dots k} \frac{n_\ell}{n} \mathbf{S}_\ell, \quad \hat{\mu}_\ell = \mathbf{g}_\ell = \frac{1}{n_\ell} \sum_{i|y_i=\ell} \mathbf{x}_i \quad \text{et} \quad \hat{\pi}_\ell = \frac{n_\ell}{n}.$$

Remarque. On utilise aussi parfois l'estimateur sans biais de la matrice de variance-covariance $\frac{n}{n-k} \mathbf{W}$.

Attribution. La règle d'affectation effective pour une observation \mathbf{x} est finalement donnée par

$$\hat{y}(\mathbf{x}) = \underset{\ell=1 \dots k}{\operatorname{argmax}} \hat{s}_\ell(\mathbf{x})$$

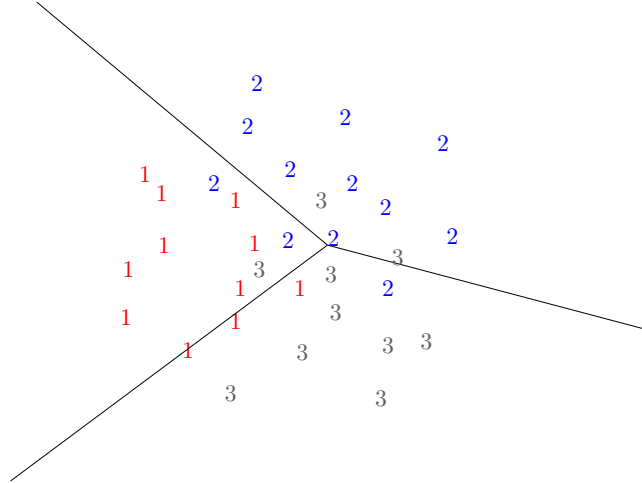
où

$$\hat{s}_\ell(\mathbf{x}) := \mathbf{x}' \mathbf{W}^{-1} \mathbf{g}_\ell + \log \frac{n_\ell}{n} - \frac{1}{2} \mathbf{g}_\ell' \mathbf{W}^{-1} \mathbf{g}_\ell.$$

Cas de deux groupes. Si $k = 2$, l'espace \mathbb{R}^2 est séparé en deux zones dont la frontière est l'hyperplan affine d'équation

$$\log \frac{n_1}{n_2} - \frac{1}{2} (\mathbf{g}_1 + \mathbf{g}_2)' \mathbf{W}^{-1} (\mathbf{g}_1 - \mathbf{g}_2) + \mathbf{x}' \mathbf{W}^{-1} (\mathbf{g}_1 - \mathbf{g}_2) = 0$$

Zones de séparation. De façon plus générale, la zone de séparation entre les régions d'attribution des classes ℓ et h est la région $\mathcal{B}_{\ell,h} \subset \mathbb{R}^p$ définie par l'équation $\hat{s}_\ell(\mathbf{x}) = \hat{s}_h(\mathbf{x})$. Cette région $\mathcal{B}_{\ell,h}$ est un hyperplan affine car les fonctions discriminantes \hat{s}_ℓ sont linéaires :



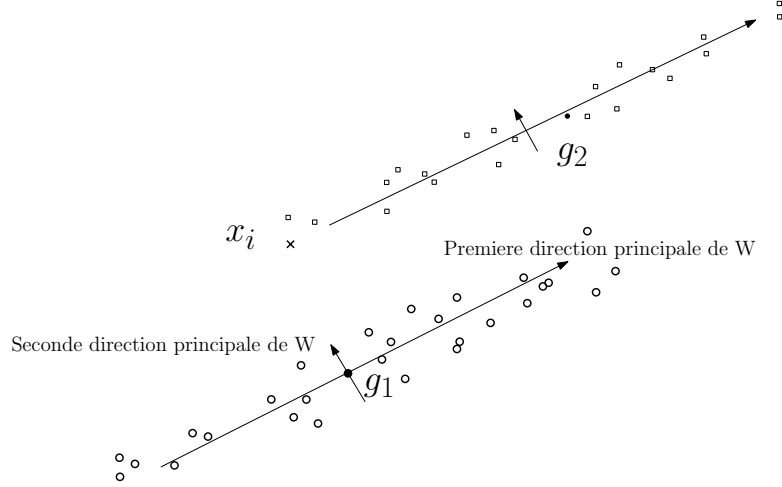
2.3 Métrique de Mahalanobis et ADL

Si les probabilités a priori sont égales, l'égalité (2) montre que la règle d'attribution revient à affecter une observation \mathbf{x}_i au groupe pour lequel la quantité $(\mathbf{x} - \mathbf{g}_\ell)' \Sigma^{-1} (\mathbf{x} - \mathbf{g}_\ell)$ est minimale. Cette quantité peut être interprétée comme une distance pour une métrique particulière : la métrique de Mahalanobis. Celle-ci est définie dans \mathbb{R}^p par le produit scalaire

$$\langle \mathbf{u}, \mathbf{v} \rangle_w := \mathbf{u}' \mathbf{W}^{-1} \mathbf{v}.$$

Pour cette métrique, les points situés sur un ellipsoïde d'équation $(\mathbf{x} - \mathbf{g}_l)' \mathbf{W}^{-1} (\mathbf{x} - \mathbf{g}_l) = c$ sont tous équidistants du point \mathbf{g}_l . Cette normalisation par \mathbf{W}^{-1} permet d'éviter que les directions relatives aux

grandes valeurs propres de \mathbf{W} soient trop prépondérantes dans le calcul des distances. Par exemple, dans l'exemple ci-dessous, les deux nuages ont des orientations comparables et un étalement important selon la première direction de \mathbf{W} .



Le point \mathbf{x}_i est plus proche de \mathbf{g}_1 que de \mathbf{g}_2 pour la métrique euclidienne mais en réalité il est plus naturel d'affecter ce point au groupe 2. En effet l'étalement dans la première direction de \mathbf{W} est tel que certains points du groupe 2 sont à proximité de \mathbf{x}_i , alors que ce n'est pas le cas pour les points du groupe 1. Pour la métrique de Mahalanobis \mathbf{x}_i est plus proche de \mathbf{g}_2 que de \mathbf{g}_1 .

Pour mieux comprendre l'effet de cette métrique, considérons le cas où il n'y a qu'une seule classe, on a alors $\mathbf{W} = \mathbf{S}$. Soit $\mathbf{Z} = \mathbf{M}\mathbf{W}^{-1/2}$ le nuage renormalisé : la matrice de variance-covariance vaut l'identité et le nuage a une forme sphérique. On peut vérifier que

$$d_{\mathbf{W}}^2(\mathbf{x}_i; \mathbf{x}_s) = d^2(\mathbf{z}_i; \mathbf{z}_s)$$

où d désigne la distance euclidienne. La métrique de Mahalanobis revient à considérer la distance euclidienne pour un nuage associé pour lequel les variables sont non corrélées et de même variance 1.

Dans le cas de k classes, \mathbf{W} est la matrice de variance-covariance moyenne des k nuages. La correspondance précédente entre $d_{\mathbf{W}}$ et d n'est plus rigoureusement exacte car en général les matrices \mathbf{S}_ℓ ne coïncident pas exactement. Cependant si ces dernières ne diffèrent pas trop, à l'intérieur d'une même nuage la métrique de Mahalanobis revient en première approximation à considérer la distance euclidienne pour un nuage renormalisé (dont les matrices de variance-covariance sont des matrices identités).

2.4 Analyse discriminante quadratique (ADQ)

Contrairement à l'ADL, l'analyse discriminante quadratique autorise les matrices de variance-covariance des nuages à être différentes. Le logarithme des rapports entre probabilités a posteriori de deux classes ℓ et h vérifie

$$\begin{aligned} \mathbf{x} \in \mathbb{R}^p, \quad \log \frac{P(Y = \ell \mid X = \mathbf{x})}{P(Y = h \mid X = \mathbf{x})} &= \log \frac{\pi_\ell}{\pi_h} + \log \frac{f_\ell(\mathbf{x})}{f_h(\mathbf{x})} \\ &= s_\ell(\mathbf{x}) - s_h(\mathbf{x}) \end{aligned}$$

où s_ℓ est la fonction discriminante du groupe ℓ définie par

$$s_\ell = \log \pi_\ell - \frac{1}{2} \log \det \Sigma_\ell - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_\ell)' \Sigma_\ell^{-1} (\mathbf{x} - \boldsymbol{\mu}_\ell).$$

Inférence. Comme dans le cas de l'ADL, les quantités Σ , $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k$ et π_1, \dots, π_k peuvent être estimées par la méthode du maximum de vraisemblance :

Proposition 4. En supposant que la loi conditionnelle de $(X|Y = \ell)$ est celle d'une loi normale multivariée p -dimensionnel (sans supposer ici que les k matrices de variances-covariances sont égales : hypothèse ADQ), les estimateurs du maximum de vraisemblance de $\boldsymbol{\mu}_\ell$, $\boldsymbol{\Sigma}_\ell$ et π_ℓ vérifient :

$$\hat{\pi}_\ell = \frac{n_\ell}{n}, \quad \hat{\boldsymbol{\mu}}_\ell = \mathbf{g}_\ell \quad \text{et} \quad \hat{\boldsymbol{\Sigma}}_\ell = \mathbf{S}_\ell.$$

Remarque. On utilise aussi parfois les estimateurs sans biais des matrices de variance-covariance $\frac{n_\ell}{n_\ell - 1} \mathbf{S}_\ell$.

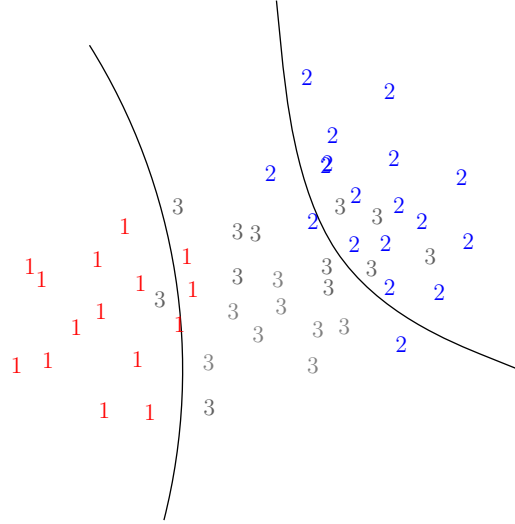
Attribution. La règle d'affectation effective pour une observation \mathbf{x} est comme précédemment donnée par

$$\hat{y}(\mathbf{x}) = \operatorname{argmax}_{\ell=1\dots k} \hat{s}_\ell(\mathbf{x})$$

où

$$\hat{s}_\ell = \log \frac{n_\ell}{n} - \frac{1}{2} \log \det \mathbf{S}_\ell - \frac{1}{2} (\mathbf{x} - \mathbf{g}_\ell)' \mathbf{S}_\ell^{-1} (\mathbf{x} - \mathbf{g}_\ell).$$

Zones de séparation. La zone de séparation entre les régions d'attribution des classes ℓ et h est l'hypersurface de \mathbb{R}^p définie par l'équation $\hat{s}_\ell(\mathbf{x}) = \hat{s}_h(\mathbf{x})$.



2.5 Choisir entre ADL et ADQ

Il est possible de construire un test sur l'égalité des matrices de variance-covariance à l'aide de la statistique :

$$Z := \left(1 - \frac{2p^2 + 3p - 1}{6(p+1)(k-1)}\right) \left[\left(\sum_{\ell=1\dots k} \frac{1}{n_\ell - 1} - \frac{1}{n - k} \right) (n - k) \log \left| \frac{n}{n - k} \mathbf{W} \right| - \sum_{\ell=1\dots k} (n_\ell - 1) \log \left| \frac{n_\ell}{n_\ell - 1} \mathbf{S}_\ell \right| \right].$$

On peut en effet montrer (admis) que sous l'hypothèse $H_0 : \mathbf{S}_1 = \dots \mathbf{S}_k$ et sous de bonnes conditions, la statistique Z converge vers une loi du χ^2 à $\frac{p(p+1)(k-1)}{2}$ degrés de liberté. Cette propriété permet ainsi de construire le test de Box (voir par exemple [Anderson, 2002], chap 10).

Attention cependant : même si le test rejette H_0 , l'ADQ ne donne pas nécessairement une meilleure classification que l'ADL car l'ADQ nécessite d'estimer beaucoup plus de coefficients que l'ADL.

Une stratégie parfois intéressante consiste à utiliser l'ADL en enrichissant la famille des variables explicatives de variables quadratiques $(\mathbf{x}^j)^2$ et de variables d'interaction $\mathbf{x}^j \times \mathbf{x}^{j'}$. Cette méthode est en effet moins « consommatrice » en paramètres. Dans tous les cas, on évaluera les erreurs de classement pour comparer les méthodes (voir plus loin).

2.6 Une version non paramétrique

Les estimateurs à noyau sont des estimateurs non paramétriques couramment utilisés en statistique. Ils permettent notamment d'estimer une densité sans hypothèse d'appartenance à une famille paramétrique de loi. On les définit par :

$$\mathbf{x} \in \mathbb{R}^p, \quad \hat{f}(\mathbf{u}) = \frac{1}{nh} \sum_{i=1 \dots n} K\left(\frac{\mathbf{x} - \mathbf{x}_i}{h}\right)$$

où $h > 0$ est la fenêtre d'estimation et $K : \mathbb{R}^p \mapsto \mathbb{R}^+$ est un noyau i.e. une fonction symétrique, à valeurs positives ou nulles et d'intégrale 1 (ex : noyau gaussien, $K = \frac{1}{2} \mathbf{1}_{[-1,1]}$). On parle alors d'estimation non paramétrique.

Dans le contexte de l'analyse discriminante, on estime donc pour chaque groupe ℓ la densité jointe des variables explicatives par

$$\mathbf{x} \in \mathbb{R}^p, \quad \hat{f}_\ell(\mathbf{x}) = \frac{1}{n_\ell h} \sum_{i \in I_\ell} K\left(\frac{\mathbf{x} - \mathbf{x}_i}{h}\right).$$

Comme pour les modèles gaussiens précédents, on utilise ensuite la formule de Bayes pour estimer $\hat{P}(Y = \ell \mid X = \mathbf{x}) = \frac{\hat{f}_\ell(\mathbf{x}) \hat{\pi}_\ell}{\sum_{j=1}^k \hat{\pi}_j \hat{f}_j(\mathbf{x})}$ et une observation est attribuée au groupe le plus probable selon la règle de Bayes.

Références

[Anderson, 2002] Anderson, T. W. (2002). *An introduction to multivariate statistical analysis, Third edition*. Wiley, New Jersey.