

Analyse discriminante

I Introduction

L'analyse discriminante est une ~~branche~~ méthode statistique qui consiste à construire une variable qualitative à k modalités (ou catégories) en utilisant p variables explicatives (quantitatives).

On considère un ensemble Ω (ensemble des individus) muni d'une partition a priori $P = (P_1, \dots, P_k)$ en k classes. Chaque individu $w \in \Omega$ est caractérisé par un ensemble de p variables descriptives (ou prédicteurs).

(ou quantitatifs). Le but général, de l'analyse discriminante est de séparer les classes a priori sur Ω à partir des p prédicteurs.

Exp 1: Analyse linguistique.

Ω = ensemble des députés élus en 1981. La partition P sur Ω comprend deux classes: les députés de droite et les députés de gauche. Les descripteurs sont les fréquences d'utilisation de 53 mots figurant dans les discours des députés. Le but de l'étude est de voir si les mots utilisés permettent de retrouver la séparation droite-gauche.

Il s'agit ici d'une étude descriptive. Le but n'est pas de découvrir la tendance des députés mais de faire une étude politico-linguistique.

Exp 2: Prédiction de risque d'avalanche (des Alpes) 1982

Ω est un ensemble des stations de ski des Alpes. La partition a priori P comprend 3 classes: Avalanches à la station, coulées d'avalanches, et pas d'avalanche. Les descripteurs sont au nombre de 274 et concernent des mesures météorologiques.

Il s'agit d'analyse discriminante à but prévisionnel: à partir de l'étude du passé le but est de déterminer chaque jour la station à haut risque d'avalanche. Le but de l'analyse prévisionnelle ou décisionnelle est de construire à partir d'un échantillon d'individus connus des règles de décisions stat. Ces règles serviront à affecter dans le futur, les individus à une classe a priori et devront minimiser

les erreurs de prévision ou de diagnostic.

En résumé: affecter un ind à une classe a priori au moyen d'une fonction des vars explicatives avec une proba d'erreur minimale.

les fonctions peuvent être explicite (linéaire, quadratique ou logistique) ou implicite (estimation de densité).

Ds ce cours nous nous limitons à analyse factorielle discriminante, qui consiste à rechercher un ensemble de s ($s < p$) variables de classement (axes discriminants) qui soient des combinaisons ^{linéaires} des q vars initiales les plus discriminants possibles au sens d'un certain critère à fixer. Ce qui signifie que les valeurs des inds d'une même classe par ces vars discriminantes sont le plus concentrées possibles et celles prises par les inds de classes \neq sont les plus éloignées possibles.

Un critère naturel, utilisant les données x_i , est d'examiner les variances intraclasse (dans la classe) et variances interclasses (entre classes).

Notations et définitions.

Soit Ω l'ensemble des individus et $P(P_1, \dots, P_k)$ la partition définie a priori sur Ω . On dispose d'un échantillon de taille n (échantillon d'apprentissage) d'elts de Ω . On désigne par X la matrice (n, p) représentant le tableau de données par les p variables descriptives.

$X = (x_{ij}^d; i=1, \dots, n, j=1, \dots, p)$ où x_{ij}^d est la valeur prise par le $i^{\text{ème}}$ ind pour la $j^{\text{ème}}$ variable.

À chaque ind w_i on associe un vecteur $x_i = \begin{pmatrix} x_{i1}^d \\ \vdots \\ x_{ip}^d \end{pmatrix} \in \mathbb{R}^p$ qui correspond à la $i^{\text{ème}}$ ligne de X . \mathbb{R}^p est l'espace des inds. Et à chaque var on associe un vecteur $x_j^d = \begin{pmatrix} x_{1j}^d \\ \vdots \\ x_{nj}^d \end{pmatrix} \in \mathbb{R}^n$ qui correspond à la $j^{\text{ème}}$ colonne du tableau. \mathbb{R}^n est l'espace des vars.

x_i le $i^{\text{ème}}$ ind et x_j^d la $j^{\text{ème}}$ var.

On suppose que chaque ind est muni d'un poids p_i ($0 \leq p_i \leq 1, \sum_{i=1}^n p_i = 1$)

$D_0 = \begin{pmatrix} p_1 & & 0 \\ & \ddots & \\ 0 & & p_n \end{pmatrix}$ métrique de poids

Elements descriptifs de l'espace des variables.

a) Moyenne $\bar{x}^j = \langle \underline{x}^j, \mathbb{1} \rangle_{D_p}$

b) Variance $\text{Var}(\underline{x}^j) = \sigma_{x^j}^2 = \sum p_i (x_i^j - \bar{x}^j)^2 = \langle \underline{x}^j - \mathbb{1}\bar{x}^j, \underline{x}^j - \mathbb{1}\bar{x}^j \rangle_{D_p}$

c) écart type $\sigma_j = \sqrt{\text{Var}(\underline{x}^j)}$

Var mesure la dispersion des autour de la moyenne.

d) covariance; $\text{cov}(\underline{x}^j, \underline{x}^{j'}) = \sum_{i=1}^n p_i (x_i^j - \bar{x}^j)(x_i^{j'} - \bar{x}^{j'}) = \langle \underline{x}^j - \mathbb{1}\bar{x}^j, \underline{x}^{j'} - \mathbb{1}\bar{x}^{j'} \rangle_{D_p}$

Rq le tableau centré $\tilde{X} = (x_i^j - \bar{x}^j; i=1, \dots, n, j=1, \dots, p)$

Elements descriptifs de l'espace des inds.

a) Le nuage des pts $N = \{ (x_i, p_i) \mid i=1, \dots, n, x_i \in \mathbb{R}^p \}$

b) Le centre de gravité ou barycentre de w_i muni de poids p_i :

$$g = \sum_{i=1}^n p_i x_i = \begin{pmatrix} x_1^p \\ \vdots \\ x_n^p \end{pmatrix} \in \mathbb{R}^p$$

Utilité de la variable qualitative.

La var quali correspond à la partition a priori $P = (P_1, \dots, P_k)$

de centres de gravité resp $G = (g_1, \dots, g_k)$ et de matrice de

var-covar (w_1, \dots, w_k) w_l désigne la dispersion à l'intérieur de la classe P_l

l'ensemble des réalisations de L est donné par

$$\Phi(\omega) = \{ y_i^l; i=1, \dots, n, l=1, \dots, k \}$$

avec $y_i^l = \begin{cases} 1 & \text{si } w_i \in P_l \\ 0 & \text{sinon.} \end{cases}$

Poids des classes

Comme chaque ind est muni d'un poids p_i , alors on attribue à

chaque classe P_l un poids $P_l = \sum_{w_i \in P_l} p_i$

Relation entre g et G .

On a $g = \sum_{l=1}^k P_l g_l$ où $g_l = \frac{1}{P_l} \sum_{w_i \in P_l} p_i x_i$

En effet:

$$\sum_{l=1}^k P_l g_l = \sum_{l=1}^k P_l \cdot \frac{1}{P_l} \sum_{w_i \in P_l} p_i x_i = \sum_{l=1}^k \sum_{w_i \in P_l} p_i x_i = \sum_{i=1}^n p_i x_i = g$$

Matrice de variance-covariance

a) Intérieur des classes: Matrice var-covar intra-classe.

On note par W (du mot anglais "within") la matrice var-covar intra-classe

$$W = \sum_{l=1}^k P_l W_l \quad \text{ou} \quad W_l = \frac{1}{P_l} \sum_{w_i \in P_l} p_i (x_i - g_l) (x_i - g_l)^t$$

W est la moyenne des matrices W_l $l=1, \dots, k$.

$$= \frac{1}{P} \sum_{l=1}^k \sum_{w_i \in P_l} p_i (x_i - g_l) (x_i - g_l)^t$$

b) Entre les classes: Matrice de var-covar inter-classes.

On note par B (du mot anglais "between") la matrice var-covar inter-classe

$$B = \sum_{l=1}^k P_l (g_l - g) (g_l - g)^t$$

Traduit l'éloignement entre les classes: B mesure la séparation des classes.

Relation entre B et G .

$$G = \begin{pmatrix} g_1 \\ \vdots \\ g_k \end{pmatrix} \quad D_p = \begin{pmatrix} P_1 & & 0 \\ & \ddots & \\ 0 & & P_k \end{pmatrix}$$

On a $B = G D_p G^t$

Matrice de variance-covariance totale

Soit $V = (cov(x^{(i)}, x^{(j)}))_{\substack{1 \leq i, j \leq p \\ 1 \leq i, j \leq p}}$ la matrice de var-covar.

$$= \sum_{i=1}^n p_i (x_i - g) (x_i - g)^t$$

La matrice de var-covar totale V mesure la dispersion totale de toutes les ind.

Relation entre V et X

Si les var sont centrées $V = \sum_{i=1}^n p_i x_i \cdot x_i^t \quad g = 0_{R^p}$.

$$V = X D_p X^t \quad (\text{sinon on pose } \bar{X} \text{ la matrice centrée}) \quad D_p \begin{pmatrix} P_1 \\ \vdots \\ P_k \end{pmatrix}$$

Théorème de Huygens

On a $V = W + B$.

Preuve. On suppose $g = 0_{R^p}$

$$V = \sum_{i=1}^n p_i (x_i - g) (x_i - g)^t = \sum_{i=1}^n p_i (x_i - g) (x_i - g)^t = \sum_{l=1}^k \sum_{w_i \in P_l} p_i (x_i - g_l + g_l) (x_i - g_l + g_l)^t$$
$$= \sum_{l=1}^k \sum_{w_i \in P_l} p_i [(x_i - g_l) + (g_l - g)] [(x_i - g_l) + (g_l - g)]^t$$

$$\begin{aligned}
 V &= \sum_{l=1}^k \sum_{w_i \in P_l} p_i (x_i - \mu_l) (x_i - \mu_l)^T + \sum_{l=1}^k \sum_{w_i \in P_l} p_i (\mu_l - \mu) (\mu_l - \mu)^T + \\
 &\quad \sum_{l=1}^k \sum_{w_i \in P_l} p_i (\mu_l - \mu) (x_i - \mu_l) (x_i - \mu_l)^T + \sum_{l=1}^k \sum_{w_i \in P_l} p_i (\mu_l - \mu) (\mu_l - \mu)^T \\
 &= W + \sum_{l=1}^k (\mu_l - \mu) \sum_{w_i \in P_l} p_i (x_i - \mu_l) + \sum_{l=1}^k \left(\sum_{w_i \in P_l} p_i (x_i - \mu_l) \right) (\mu_l - \mu)^T + B \\
 &= W + B + \sum_{l=1}^k (\mu_l - \mu) \left[\sum_{w_i \in P_l} p_i \underset{0}{x_i} - P_l \mu_l \right] + \sum_{l=1}^k \left[\sum_{w_i \in P_l} p_i \underset{0}{x_i} - P_l \mu_l \right] (\mu_l - \mu)^T
 \end{aligned}$$

avec $\mu_l = \frac{1}{P_l} \sum_{w_i \in P_l} p_i x_i$

D'où $V = W + B$

Commentaire : La décomposition d'hygens produit un nouvel espace de représentation qui permet de distinguer au mieux les k classes en construisant une suite de vars discriminantes z^h (donc de nouveaux axes discriminants) telle que les vars z^h soient non corrélés entre et la projection de l'indiv d'une m^e classe sur les axes discriminant soient le plus proches possible et ceux d'une classe \neq soient les plus éloignés possible.

Les nouvelles vars z^h sont des combinaisons linéaires de vars initiales.

Pour des raisons de simplicité de notation nous supposons ds ce qui suit que les vars sont centrées $\mu = 0_{\mathbb{R}^p}$.

Variances et covariances de vars discriminantes.

Si les vars discriminantes z^h sont données par $z^h = \sum_{j=1}^p u_{hj} x^j \in \mathbb{R}^h$ (combinaisons linéaires des vars initiales), alors

$\text{Var}(z^h) = U_h^T V U_h$ avec $U_h = \begin{pmatrix} u_{h1} \\ \vdots \\ u_{hp} \end{pmatrix}$

et $\text{cov}(z^h, z^l) = U_h^T V U_l = U_l^T V U_h$

Démonstration

~~$$\begin{aligned}
 U_h^T V U_h &= \begin{pmatrix} u_{h1} & \dots & u_{hp} \end{pmatrix} \left(\sum_{i=1}^n p_i (x_i - \mu) (x_i - \mu)^T \right) \begin{pmatrix} u_{h1} \\ \vdots \\ u_{hp} \end{pmatrix} \\
 &= \begin{pmatrix} u_{h1} & \dots & u_{hp} \end{pmatrix} \begin{pmatrix} \sum p_i x_i^1 x_i^1 & & \\ & \sum p_i x_i^2 x_i^2 & \\ & & \ddots & \\ & & & \sum p_i x_i^p x_i^p \end{pmatrix} \begin{pmatrix} u_{h1} \\ \vdots \\ u_{hp} \end{pmatrix} \\
 &= \sum_{j=1}^p \left(\sum_{i=1}^n p_i x_i^j x_i^j \right) u_{hj}^2
 \end{aligned}$$~~

$$\begin{aligned} \text{Var}(z^h) &= \sum_{i=1}^n p_i (z_i^h)^2 = \sum_{i=1}^n p_i \left(\sum_{j=1}^p u_{jh} x_i^j \right)^2 \\ &= \sum_{i=1}^n p_i \begin{pmatrix} u_{1h} & \dots & u_{ph} \end{pmatrix} \begin{pmatrix} x_i^1 \\ \vdots \\ x_i^p \end{pmatrix} \begin{pmatrix} x_i^1 & \dots & x_i^p \end{pmatrix} \begin{pmatrix} u_{1h} \\ \vdots \\ u_{ph} \end{pmatrix} \\ &= \begin{pmatrix} u_{1h} & \dots & u_{ph} \end{pmatrix} \sum_{i=1}^n p_i \begin{pmatrix} x_i^1 \\ \vdots \\ x_i^p \end{pmatrix} \begin{pmatrix} x_i^1 & \dots & x_i^p \end{pmatrix} \begin{pmatrix} u_{1h} \\ \vdots \\ u_{ph} \end{pmatrix} \\ &= \begin{pmatrix} u_{1h} & \dots & u_{ph} \end{pmatrix} V \begin{pmatrix} u_{1h} \\ \vdots \\ u_{ph} \end{pmatrix} \end{aligned}$$

Pour la cov

$$\begin{aligned} \text{cov}(z^h, z^l) &= \sum_{i=1}^n p_i z_i^h z_i^l = \sum_{i=1}^n p_i \left(\sum_{j=1}^p u_{jh} x_i^j \right) \left(\sum_{k=1}^p u_{kl} x_i^k \right) \\ &= \sum_{i=1}^n p_i \begin{pmatrix} u_{1h} & \dots & u_{ph} \end{pmatrix} \begin{pmatrix} x_i^1 \\ \vdots \\ x_i^p \end{pmatrix} \begin{pmatrix} x_i^1 & \dots & x_i^p \end{pmatrix} \begin{pmatrix} u_{1l} \\ \vdots \\ u_{pl} \end{pmatrix} \\ &= \begin{pmatrix} u_{1h} & \dots & u_{ph} \end{pmatrix} \sum_{i=1}^n p_i \begin{pmatrix} x_i^1 \\ \vdots \\ x_i^p \end{pmatrix} \begin{pmatrix} x_i^1 & \dots & x_i^p \end{pmatrix} \begin{pmatrix} u_{1l} \\ \vdots \\ u_{pl} \end{pmatrix} = \begin{pmatrix} u_{1h} & \dots & u_{ph} \end{pmatrix} V \begin{pmatrix} u_{1l} \\ \vdots \\ u_{pl} \end{pmatrix} \end{aligned}$$

Rp: comme on a $V = W + B$ (var totale = variance intra classe + var inter class)

$$\text{alors } \begin{pmatrix} u_{1h} & \dots & u_{ph} \end{pmatrix} V \begin{pmatrix} u_{1h} \\ \vdots \\ u_{ph} \end{pmatrix} = \begin{pmatrix} u_{1h} & \dots & u_{ph} \end{pmatrix} W \begin{pmatrix} u_{1h} \\ \vdots \\ u_{ph} \end{pmatrix} + \begin{pmatrix} u_{1h} & \dots & u_{ph} \end{pmatrix} B \begin{pmatrix} u_{1h} \\ \vdots \\ u_{ph} \end{pmatrix}$$

ce qui signifie que la variance totale de z^h est égale à la variance de z^h intra classe + variance de z^h inter classe.

Choix des variables discriminantes

c'est la détermination des axes discriminants.

On cherche les variables discriminantes z^h de variance totale $V(z^h)$ donnée dont la variance interclass B est maximale. (de manière à séparer au mieux les classes).

Par la relation $\begin{pmatrix} u_{1h} & \dots & u_{ph} \end{pmatrix} V \begin{pmatrix} u_{1h} \\ \vdots \\ u_{ph} \end{pmatrix} = \begin{pmatrix} u_{1h} & \dots & u_{ph} \end{pmatrix} W \begin{pmatrix} u_{1h} \\ \vdots \\ u_{ph} \end{pmatrix} + \begin{pmatrix} u_{1h} & \dots & u_{ph} \end{pmatrix} B \begin{pmatrix} u_{1h} \\ \vdots \\ u_{ph} \end{pmatrix}$

$$\max_{u_h} \begin{pmatrix} u_{1h} & \dots & u_{ph} \end{pmatrix} B \begin{pmatrix} u_{1h} \\ \vdots \\ u_{ph} \end{pmatrix} \Leftrightarrow \min_{u_h} \begin{pmatrix} u_{1h} & \dots & u_{ph} \end{pmatrix} W \begin{pmatrix} u_{1h} \\ \vdots \\ u_{ph} \end{pmatrix} \quad (\text{classes les plus homogènes minimum})$$

On peut formuler le pb d'optimisation de plusieurs façons

Pb 1: Trouver $u_h \in \mathbb{R}^p$ tp
$$\begin{cases} \begin{pmatrix} u_{1h} & \dots & u_{ph} \end{pmatrix} B \begin{pmatrix} u_{1h} \\ \vdots \\ u_{ph} \end{pmatrix} \text{ max} \\ \begin{pmatrix} u_{1h} & \dots & u_{ph} \end{pmatrix} V \begin{pmatrix} u_{1h} \\ \vdots \\ u_{ph} \end{pmatrix} = C \end{cases} \text{ sous la contrainte}$$

La discrimination est alors impossible.

La statistique de décision est Λ de Wilks

$$\Lambda = \frac{|W|}{|V|} \quad \text{1.1 déterminant.}$$

sous H_0 , Λ suit la loi de Wilks de paramètres $(n-1)$, p , $k-1$ notée $\Lambda(n-1, p, k-1)$.

La région critique du test $RC = [\Lambda < \Lambda_{\alpha, tab}^{\alpha}]$ est tabulée.
au niveau α fixé.

Si H_0 est rejeté, il est possible de projeter les inds ds un nouvel espace fact permettant de distinguer au mieux les classes.

Variance expliquée

Le pouvoir discriminant d'un axe factoriel est évalué par le rapport $\frac{\lambda_j}{\sum_{j=1}^k \lambda_j}$ inertie ou proportion de variance expliquée par l'axe.

L'inertie ou proportion de variance expliquée par le plan est obtenue par les deux premiers axes factoriels est : $\frac{\lambda_1 + \lambda_2}{\sum_{j=1}^k \lambda_j}$.

de manière générale l'inertie expliquée par le s. espace factoriel de dim s

$$\text{est } \frac{\lambda_1 + \lambda_2 + \dots + \lambda_s}{\sum_{j=1}^{k-1} \lambda_j} \quad s \leq p-1$$

Rq!

La somme est jusqu'à $k-1$ car il y a au plus $(k-1)$ vps non nulles si $k \leq p$

En effet, la matrice B de variance-covariance intra-classe est la matrice de associé aux centres de gravité g_1, \dots, g_k

Or on a $\sum_{k=1}^k \Pi_k g_k = g$. d'où B est de rang $\leq k-1$

$$\Rightarrow \text{rg}(V^{-1}B) \leq \min(p, k-1)$$

$$\text{si } k < p \Rightarrow k-1 < p \Rightarrow \min(p, k-1) = k-1$$

$$\Rightarrow \text{rg}(V^{-1}B) \leq (k-1)$$