Scoring & Applications – Chapitre 1 – Méthodes géométriques

Contexte et objectif

On dispose de n individus décrits par d variables continues et répartis en G groupes. $\mathbf{X} = (x_{i,j}) \in \mathbb{R}^{n \times d}$ est la matrice des descripteurs : $x_{i,j}$ est la valeur de la variable j mesurée sur l'individu i. $\mathbf{Z} = (z_{i,g}) \in \{0,1\}^{n \times G}$ est la matrice des appartenances : $z_{i,g} = 1$ si l'individu i appartient au groupe g et $z_{i,g} = 0$ sinon. L'individu i est assimilé à la colonne i de \mathbf{X}' que l'on note $\mathbf{x}_i \in \mathbb{R}^d$.

On souhaite s'appuyer sur les informations contenues dans X et Z pour attribuer à l'un des groupes un nouvel individu $x_{n+1} \in \mathbb{R}^d$.

Statistiques usuelles

Statistiques par groupe

Le groupe g a pour effectif : $n_g = \sum_{i=1}^n z_{i,g}$, son centre est : $\bar{\boldsymbol{x}}_g = \sum_{i=1}^n z_{i,g} \boldsymbol{x}_i / n_g$ et sa matrice de variance : $\boldsymbol{V}_g = \sum_{i=1}^n z_{i,g} (\boldsymbol{x}_i - \bar{\boldsymbol{x}}_g) (\boldsymbol{x}_i - \bar{\boldsymbol{x}}_g)' / n_g$.

 $Statistiques\ globales$

Le centre du nuage est : $\bar{x} = \sum_{i=1}^n x_i/n$ et la matrice de variance : $V = \sum_{i=1}^n (x_i - \bar{x}) \cdot (x_i - \bar{x})'/n$.

Des statistiques par groupe aux statistiques globales

Le centre du nuage peut aussi être calculé selon : $\bar{x} = \sum_{g=1}^{G} n_g \bar{x}_g / n$.

La matrice de variance intra groupes : ${m W} = \sum_{g=1}^G n_g {m V}_g/n$ traduit la dispersion des données dans les groupes.

La matrice de variance inter groupes : $\mathbf{B} = \sum_{g=1}^{G} n_g(\bar{\mathbf{x}}_g - \bar{\mathbf{x}}).(\bar{\mathbf{x}}_g - \bar{\mathbf{x}})'/n$ traduit la dispersion des groupes autour du centre du nuage.

La matrice de variance peut être obtenue par le théorème de König-Huygens : V = B + W.

Analyse Factorielle Discriminante

Quelques prérequis mathématiques: http://alexandrelourme.free.fr/M2IREF/SCORING/CH0.pdf

Description de la méthode

Une matrice $M \in \mathbb{R}^{d \times d}$ symétrique définie positive définissant la métrique de \mathbb{R}^d , l'analyse factorielle discriminante (AFD) à un facteur se résume à quelques étapes simples.

o Conditionnement des données (centrage)

On note 1 le vecteur de \mathbb{R}^n dont toutes les composantes valent 1.

- (i) centrer le nuage de points : $X \leftarrow X 1.\bar{x}'$
- (ii) recentrer le nouvel individu : $\boldsymbol{x}_{n+1} \leftarrow \boldsymbol{x}_{n+1} \bar{\boldsymbol{x}}$

Dès lors, les statistiques de la section précédente sont calculées en prenant l'origine pour centre du nuage : $\bar{x} = \mathbf{0}_{\mathbb{R}^d}$.

- \circ Analyse spectrale de la matrice $V^{-1}B$
- (iii) déterminer la plus grande des valeurs propres de ${m V}^{-1}{m B}$ et ${m u}$ un vecteur propre associé $^{
 m a}$
- (iv) normaliser u de sorte qu'il soit M^{-1} -normé : $u \leftarrow u/\sqrt{u'M^{-1}u}$
- (v) déterminer le vecteur directeur M-unitaire de l'axe factoriel : $a = M^{-1}u$
- \circ Allocation du nouvel individu x_{n+1}
- (vi) déterminer les abscisses du nuage projeté M-orthogonalement sur l'axe factoriel : $\mathbf{s} = Xu \in \mathbb{R}^n$
- (vii) déterminer le centre de chaque groupe sur l'axe factoriel : $\bar{s}_g = \sum_{i=1}^n z_{i,g} s_i / n_g$
- (viii) attribuer \boldsymbol{x}_{n+1} au groupe dont il est le plus proche sur l'axe factoriel : \boldsymbol{x}_{n+1} est attribué au groupe g pour lequel le score $\mathfrak{s}_g(\boldsymbol{x}_{n+1}) = |\boldsymbol{x}'_{n+1}.\boldsymbol{u} \bar{s}_g|$ est minimal

Interprétation de la méthode

Une fois le nuage le nuage centré, l'AFD détermine la direction (celle de a) sur laquelle le nuage doit être projeté M-orthogonalement pour que la série résultante ait une variance inter groupes maximale.

Justification de la méthode

A. Lourme, Faculté d'économie, gestion & AES, Université de Bordeaux http://alexandrelourme.free.fr

^asi le sous espace propre associé à la plus grande des valeurs propres de $V^{-1}B$ est de dimension strictement supérieure à 1, le problème d'optimisation sous jacent à l'AFD possède plusieurs solutions.

Notons \boldsymbol{a} un vecteur M-normé $(\boldsymbol{a}'M\boldsymbol{a}=1)$ de \mathbb{R}^d et $<\boldsymbol{a}>$ la droite vectorielle engendrée par \boldsymbol{a} . $(\boldsymbol{a}'M\boldsymbol{x}_i)\boldsymbol{a}$ est le projeté M-orthogonal de \boldsymbol{x}_i sur $<\boldsymbol{a}>$, $s_i=\boldsymbol{a}'M\boldsymbol{x}_i$ est son abscisse sur $<\boldsymbol{a}>$ et la série $\boldsymbol{s}=(s_1,\ldots,s_n)'=XM\boldsymbol{a}\in\mathbb{R}^n$ contient l'ensemble des abscisses des points du nuage projetés M-orthogonalement sur $<\boldsymbol{a}>$. La série \boldsymbol{s} hérite de la structure en G groupes des données de \boldsymbol{X} : si \boldsymbol{x}_i est dans le groupe \boldsymbol{g} du nuage, $s_i\boldsymbol{a}$ est dans le groupe \boldsymbol{g} des points projetés sur $<\boldsymbol{a}>$. La variance inter groupes de \boldsymbol{s} est : $\boldsymbol{a}'MBM\boldsymbol{a}$, la variance intra groupes : $\boldsymbol{a}'MWM\boldsymbol{a}$ et la variance totale : $\boldsymbol{a}'MVM\boldsymbol{a}$. On cherche \boldsymbol{a} de sorte que le rapport de corrélation de \boldsymbol{s} : $\eta^2=(\boldsymbol{a}'MBM\boldsymbol{a})/(\boldsymbol{a}'MVM\boldsymbol{a})$ soit maximal ou, de façon équivalente, $\boldsymbol{u}=M\boldsymbol{a}$ de sorte que $\eta^2(\boldsymbol{u})=(\boldsymbol{u}'B\boldsymbol{u})/(\boldsymbol{u}'V\boldsymbol{u})$ soit maximal. On peut montrerb que $\eta^2(\boldsymbol{u})$ est maximal lorsque \boldsymbol{u} est un vecteur propre associé à λ , la plus grande des valeurs propre de $\boldsymbol{V}^{-1}\boldsymbol{B}$. Or, \boldsymbol{u} est M^{-1} -normé puisque \boldsymbol{a} est M-normé. Ainsi, \boldsymbol{u} est un vecteur propre M^{-1} -normé associé à la plus grande des valeurs propres de $\boldsymbol{V}^{-1}\boldsymbol{B}$. \boldsymbol{u} est le premier facteur discriminant ; il représente la forme linéaire à appliquer aux d variables initiales pour obtenir la série \boldsymbol{s} la plus discriminante. $<\boldsymbol{a}>$ est l'axe factoriel. λ est le rapport de corrélation maximal que l'on peut obtenir en projetant le nuage M-orthogonalement sur une droite vectorielle.

Un cas particulier

Lorsque les observations sont réparties en deux groupes (G=2), la métrique de Mahalanobis $(\boldsymbol{M}=\boldsymbol{V}^{-1})$ permet d'obtenir le facteur discriminant et l'axe discriminant directement^c ; en effet, \boldsymbol{a} est colinéaire à $\bar{\boldsymbol{x}}_1 - \bar{\boldsymbol{x}}_2$ et \boldsymbol{u} est colinéaire à $\boldsymbol{V}^{-1}(\bar{\boldsymbol{x}}_1 - \bar{\boldsymbol{x}}_2)$.

Extension à plusieurs facteurs

L'analyse factorielle discriminante s'étend à p facteurs $(1 \le p \le \min\{G-1,d\})$ en adoptant la métrique de Mahalanobis : $\mathbf{M} = \mathbf{V}^{-1}$.

 $V^{-1}B$ est V-symétrique, donc diagonalisable dans une base V-orthonormée : (u_1, \ldots, u_d) . En notant λ_i la valeur propre associée au vecteur propre u_i on peut supposer sans perte de généralité : $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d$.

Rq. Les valeurs propres λ_i sont positives ou nulles : $V^{-1/2}BV^{-1/2}$ est une matrice réelle, symétrique, positive ; elle est diagonalisable dans une base orthonormée et ses valeurs propres sont positives ou nulles ; elle possède les mêmes valeurs propres que $V^{-1}B$.

On note U la matrice $d \times p$ dont la colonne j est composée des coordonnées du $j^{\rm e}$ facteur discriminant : u_j .

La matrice $S = XU \in \mathbb{R}^{n \times p}$ est alors constituée des données factorielles : la ligne i de S représente l'individu i dans l'espace factoriel et la colonne j de S la j^e variable discriminante.

La matrice de variance des données factorielles : $U'VU \in \mathbb{R}^{p \times p}$ est l'identité. La matrice de variance inter groupes : $U'BU \in \mathbb{R}^{p \times p}$ est diagonale ; ses coefficients diagonaux sont : $\lambda_1, \ldots, \lambda_p$. Donc : $\operatorname{trace}(U'BU)/\operatorname{trace}(U'VU) = \sum_{j=1}^p \lambda_j/p$.

L'analyse factorielle discriminante à p facteurs en quelques étapes

Les étapes (i) et (ii) de conditionnement des données étant réalisées :

- (a) déterminer la matrice $U \in \mathbb{R}^{d \times p}$ formée en colonne des facteurs discriminants u_1, \dots, u_n
- (b) calculer $S = XU \in \mathbb{R}^{n \times p}$; la colonne j de S est la variable discriminante j; elle contient les coordonnées des points du nuage projetés sur l'axe factoriel j
 - (c) déterminer le centre de chaque groupe dans l'espace factoriel : $U'\bar{x}_q$
 - (d) déterminer les coordonnées du nouvel individu dans l'espace factoriel : $U'x_{n+1}$
- (e) affecter le nouvel individu au groupe dont le centre est le plus proche dans l'espace factoriel : x_{n+1} est attribué au groupe g pour lequel le score $\mathfrak{s}_g(x_{n+1}) = \|U'x_{n+1} U'\bar{x}_g\|$ est minimal

Application

Le fichier client donne pour cent clients bancaires : la liquidité (cash), le flux (flow), l'épargne (saving), le niveau de consommation (consume) et la classe de risque (risk).

A quelle classe de risque l'Analyse Factorielle Discriminante affecte-t-elle un nouveau client avec les caractéristiques suivantes : liquidité = 6, 2; flux = 2, 9; épargne = 4, 9; consommation = 1, 7?

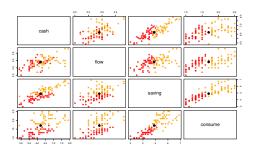
References

[1] Gilbert Saporta. Probabilités, analyse des données et statistique. Editions Technip, 2006.

^bvoir [1], p.484

 $^{^{\}text{c}}$ sans avoir à déterminer le spectre de $V^{-1}B$

 $^{^{\}rm d} disponible \ sous: \ {\tt http://alexandrelourme.free.fr/M2IREF/SCORING/client.csv}$



 $Fig. \ 1: \ \textit{Cent clients bancaires décrits par quatre variables continues (cash, flow, saving, consume) et répartis en deux classes de risque (0 en rouge, 1 en orange) ; à quelle classe de risque doit-on affecter le nouveau client (point noir) ? \\$