

Reconnaissance des formes

Notions de base et notations

- Chapitre 1 -

RECONNAISSANCE DES FORMES

Objectifs

L'objectif de l'analyse de données est de synthétiser, structurer, ..., l'information véhiculée par des données multidimensionnelles :

- ▷ n : nombre d'individus
- ▷ p : nombre de variables

Les méthodes mises en œuvre relèvent essentiellement de l'*algèbre linéaire* et de la *théorie des probabilités*. En effet :

- ▷ les données sont vues comme un nuage de points dans un espace vectoriel
- ▷ La statistique inférentielle permet de fournir des résultats relatifs à une population à partir de mesures statistiques réalisées sur des échantillons.

VOCABULAIRE

Individus et variables

Population : Groupe ou ensemble d'individus que l'on analyse

Recensement : Etude de tous les individus d'une population donnée

Sondage : Etude d'une partie seulement d'une population appelée échantillon

Variables : Ensemble de caractéristiques d'une population

- *quantitatives* : nombres sur lesquels les opérations usuelles ont un sens. Elles peuvent être discrètes ou continues
- *qualitatives* : appartenance à une catégorie donnée. Elles peuvent être nominales, ou ordinales quand les catégories sont ordonnées.

VOCABULAIRE

Description de données quantitatives

Variable, individu : On appelle variable un vecteur \boldsymbol{x} de taille n . Chaque coordonnée x_i correspond à un individu.

Poids : Chaque individu a éventuellement un poids p_i , tel que $p_1 + \dots + p_n = 1$.
On choisit souvent $p_i = \frac{1}{n}$.

Analyse : On dispose d'une série d'indicateurs qui ne donne qu'une vue partielle des données : effectif, moyenne, médiane, variance, écart type, minimum, maximum, 1^{er} quartile, ...

Ces indicateurs mesurent principalement la tendance centrale et la dispersion.

On utilisera essentiellement la moyenne, la variance et l'écart type.

GRANDEURS STATISTIQUES DE BASE

Moyenne arithmétique

Définition 1. *On appelle moyenne arithmétique, que l'on note \bar{x} , la quantité suivante*

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

ou, dans le cas d'une somme pondérée

$$\bar{x} = \sum_{i=1}^n p_i x_i$$

Remarque. La moyenne arithmétique est une mesure de tendance centrale qui dépend de toutes les observations, et est sensible aux valeurs extrêmes.

GRANDEURS STATISTIQUES DE BASE

Variance et écart-type

Définition 2. *La variance de x est définie par*

$$\sigma_x^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

ou, dans le cas d'une pondération non-uniforme

$$\sigma_x^2 = \sum_{i=1}^n p_i (x_i - \bar{x})^2$$

L'écart type σ_x est la racine-carrée de la variance.

Propriété 1. *La variance satisfait la relation suivante*

$$\sigma_x^2 = \sum_{i=1}^n p_i x_i^2 - \bar{x}^2$$

L'écart-type, qui a la même unité que x , est une mesure de dispersion.

GRANDEURS STATISTIQUES DE BASE

Mesure de liaison entre deux variables

Définition 3. *La covariance observée entre deux variables x et y est définie par*

$$\sigma_{xy} = \sum_{i=1}^n p_i (x_i - \bar{x})(y_i - \bar{y})$$

Le coefficient de corrélation est donné par

$$r_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$$

GRANDEURS STATISTIQUES DE BASE

Propriétés du coefficient de corrélation

Propriété 2. *D'après l'inégalité de Cauchy-Schwartz, on a*

$$-1 \leq r_{xy} \leq 1$$

Propriété 3. *Le résultat suivant concerne des variables dites linéairement liées.*

$$|r_{xy}| = 1 \Leftrightarrow ax_i + by_i = c, 1 \leq i \leq n$$

En particulier, on a $r_{xx} = 1$.

Remarque. Si $r_{xy} = 0$, les variables sont dites décorrélées. Cela ne signifie pas qu'elles sont indépendantes.

DESCRIPTION DES DONNÉES

Notations matricielles

Matrice : De manière impropre, une matrice à p lignes et n colonnes est un tableau rectangulaire de mn nombres, rangés ligne par ligne.

Vecteur : Un vecteur, ligne ou colonne, est une matrice ne comportant qu'une seule ligne ou qu'une seule colonne.

Transposition : Echange des lignes et des colonnes d'une matrice. On note M^T la transposée de M .

Exemples :

$$I = \begin{pmatrix} 1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 1 \end{pmatrix} \quad \mathbf{1} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$$

DESCRIPTION DES DONNÉES

Tableau de données

Dans toute la suite, pour n individus et p variables, on s'intéresse aux tableaux de données définis comme suit

$$\mathbf{X} = (\mathbf{x}^1 \quad \dots \quad \mathbf{x}^p) = \begin{pmatrix} x_1^1 & x_1^2 & \dots & x_1^p \\ x_2^1 & x_2^2 & & \\ & & \ddots & \\ \vdots & & & x_i^j & \vdots \\ & & & & \ddots & \\ x_n^1 & & \dots & & & x_n^p \end{pmatrix} \quad \mathbf{1} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$$

DESCRIPTION DES DONNÉES

Vecteurs variable et individu

Variable : Une colonne du tableau de données

$$\mathbf{x}^j = \begin{pmatrix} x_1^j \\ x_2^j \\ \dots \\ x_n^j \end{pmatrix}$$

Individu : Une ligne du tableau de données, transposées

$$\mathbf{x}_i = \left(x_i^1 \quad x_i^2 \quad \dots \quad x_i^p \right)^\top$$

DESCRIPTION DES DONNÉES

Matrice de poids

Pourquoi : Elle est nécessaire quand les individus n'ont pas la même importance.

Comment : On associe un poids p_i à chaque individu tel que :

$$p_1 + p_2 + \dots + p_n = 1$$

On regroupe ces poids dans une matrice diagonale de taille n :

$$D = \begin{pmatrix} p_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & p_n \end{pmatrix}$$

Cas uniforme : Tous les individus ont le même poids $p_i = 1/n$

PRÉ-TRAITEMENTS

Individu moyen et tableau centré

Individu moyen : L'individu moyen est obtenu à partir de la moyenne arithmétique de chaque variable

$$\mathbf{m} = \left(\bar{x}^1 \quad \bar{x}^2 \quad \dots \quad \bar{x}^p \right)^\top$$

avec $\bar{x}^j = \sum_{i=1}^n p_i x_i^j$. On peut aussi écrire

$$\mathbf{m} = \mathbf{X}^\top \mathbf{D} \mathbf{1}$$

Tableau centré : Il est obtenu en centrant l'ensemble des variables du tableau de données : $y_i^j = x_i^j - \bar{x}^j$. Sous forme matricielle, on écrit

$$\mathbf{Y} = \mathbf{X} - \mathbf{1} \mathbf{m}^\top = (\mathbf{I} - \mathbf{1} \mathbf{1}^\top \mathbf{D}) \mathbf{X}$$

PRÉ-TRAITEMENTS

Matrice de variance-covariance

Définition : Il s'agit d'une matrice de dimension p définie par

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \cdots & \sigma_{1p} \\ \sigma_{21} & \sigma_2^2 & & \\ & & \ddots & \\ \vdots & & & \sigma_i^2 & \vdots \\ & & & & \ddots & \\ \sigma_{p1} & & \cdots & & & \sigma_p^2 \end{pmatrix}$$

où σ_{ij} est la covariance des variables x^i et x^j , et σ_j^2 est la variance de x^j .

Formulation matricielle :

$$\Sigma = X^\top DX - mm^\top = Y^\top DY$$

PRÉ-TRAITEMENTS

Matrice de corrélation

Définition : Il s'agit d'une matrice de dimension p définie par

$$\mathbf{R} = \begin{pmatrix} 1 & r_{12} & \cdots & r_{1p} \\ r_{21} & 1 & & \\ & & \ddots & \\ \vdots & & & r_{ij} & \vdots \\ & & & & \ddots & \\ r_{p1} & & \cdots & & & 1 \end{pmatrix}$$

où $r_{ij} = \frac{\sigma_{ij}}{\sigma_i \sigma_j}$ est le coefficient de corrélation des variables x^i et x^j .

Formulation matricielle :

$$\mathbf{R} = \mathbf{D}_{1/\sigma} \mathbf{\Sigma} \mathbf{D}_{1/\sigma}$$

où $\mathbf{D}_{1/\sigma}$ est la matrice diagonale de termes diagonaux $\frac{1}{\sigma_i}$.

ANALYSE DES DONNÉES

Métrie

Motivation : Il est nécessaire d'introduire une métrique afin de caractériser la topologie du nuage de points.

Définition : On appelle distance sur E une application $d : E \times E \rightarrow \mathbb{R}^+$ vérifiant les propriétés suivantes

- symétrie : $\forall x, y \in E, \quad d(x, y) = d(y, x)$
- séparation : $\forall x, y \in E, \quad d(x, y) = 0 \Leftrightarrow x = y$
- inégalité triangulaire : $\forall x, y, z \in E, \quad d(x, z) \leq d(x, y) + d(y, z)$

Exemple : La distance euclidienne entre 2 points \mathbf{u} et \mathbf{v} de \mathbb{R}^p est définie par

$$d^2(\mathbf{u}, \mathbf{v}) = \sum_{j=1}^p (u_j - v_j)^2 = \|\mathbf{u} - \mathbf{v}\|^2$$

ANALYSE DES DONNÉES

Métrie

Matrice définie positive : Il s'agit d'une matrice symétrique \mathbf{M} telle que, pour tout \mathbf{u} non nul, on a $\mathbf{u}^\top \mathbf{M} \mathbf{u} > 0$.

Définition : Soit \mathbf{M} une matrice définie positive de dimension p . La fonction suivante $d_M : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}^+$ définit une métrie

$$d_M^2(\mathbf{u}, \mathbf{v}) = \|\mathbf{u} - \mathbf{v}\|_M^2 \quad \text{avec} \quad \|\mathbf{u}\|_M^2 = \sum_{i,j=1}^p m_{ij} u_i u_j$$

Cette distance est appelée distance de Mahalanobis lorsque $\mathbf{M} = \Sigma^{-1}$, où Σ est la matrice de variance-covariance des données.

Produit scalaire : La métrie définie ci-dessus dérive du produit scalaire

$$\langle \mathbf{u}, \mathbf{v} \rangle_M = \sum_{i,j=1}^p m_{ij} u_i v_j$$

On dit que \mathbf{u} et \mathbf{v} sont orthogonaux si $\langle \mathbf{u}, \mathbf{v} \rangle_M = 0$.

ANALYSE DES DONNÉES

Métriques particulières

Métrique euclidienne : Elle est obtenue pour $M = I$.

L'une des difficultés rencontrées avec la métrique euclidienne est qu'elle privilégie les variables les plus dispersées et dépend donc de leur unité de mesure.

Métrique réduite : Elle consiste à prendre $M = D_{1/\sigma^2}$, où D_{1/σ^2} est la matrice diagonale de termes diagonaux les inverses $\frac{1}{\sigma_i}$ des variances des variables.

$$D_{1/\sigma^2} = \begin{pmatrix} \frac{1}{\sigma_1^2} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{1}{\sigma_p^2} \end{pmatrix}$$

ANALYSE DES DONNÉES

Inertie

Définition : L'inertie du nuage de points $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ en un point quelconque \mathbf{a} est donnée par

$$I_a = \sum_{i=1}^n p_i \|\mathbf{x}_i - \mathbf{a}\|_M^2$$

Définition : L'inertie totale du nuage de points $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ est donnée par

$$I_m = \frac{1}{2} \sum_{i,j=1}^n p_i p_j \|\mathbf{x}_i - \mathbf{x}_j\|_M^2$$

ANALYSE DES DONNÉES

Inertie

Propriété :

$$I_m = \text{Trace}(\Sigma M)$$

Démonstration : On introduit le vecteur moyen \mathbf{m} , et on déroule le calcul ainsi

$$\begin{aligned} I_m &= \frac{1}{2} \sum_{i,j=1}^n p_i p_j \|\mathbf{x}_i - \mathbf{x}_j\|_M^2 \\ &= \frac{1}{2} \sum_{i,j=1}^n p_i p_j \|(\mathbf{x}_i - \mathbf{m}) - (\mathbf{x}_j - \mathbf{m})\|_M^2 \\ &= \underbrace{\sum_{i=1}^n p_i \|\mathbf{x}_i - \mathbf{m}\|_M^2}_{\text{Trace}(\Sigma M)} - \underbrace{\sum_{i,j=1}^n p_i p_j \langle \mathbf{x}_i - \mathbf{m}, \mathbf{x}_j - \mathbf{m} \rangle_M}_0 \\ &= \text{Trace}(\Sigma M) \end{aligned}$$

ANALYSE DES DONNÉES

Inertie

Métrie euclidienne :

$$I_m = \text{Trace}(\mathbf{\Sigma}) = \sum_{i=1}^p \sigma_i^2$$

Métrie réduite :

$$\begin{aligned} I_m &= \text{Trace}(\mathbf{\Sigma} \mathbf{D}_{1/\sigma^2}) \\ &= \text{Trace}(\mathbf{D}_{1/\sigma} \mathbf{\Sigma} \mathbf{D}_{1/\sigma}) \\ &= \text{Trace}(\mathbf{R}) = p \end{aligned}$$

ANALYSE DES DONNÉES

Métrie et tableau de données

Utiliser la métrie $\mathbf{M} = \mathbf{T}^\top \mathbf{T}$ sur le tableau de données \mathbf{X} est équivalent à travailler avec la métrie euclidienne sur le tableau transformé $\mathbf{X}\mathbf{T}^\top$.

Tableau transformé : Lorsqu'on travaille sur le tableau transformé comme ci-dessus, il convient d'utiliser la norme euclidienne. En effet,

$$\langle \mathbf{T}\mathbf{x}_i, \mathbf{T}\mathbf{x}_j \rangle = (\mathbf{T}\mathbf{x}_i)^\top (\mathbf{T}\mathbf{x}_j) = \mathbf{x}_i^\top (\mathbf{T}^\top \mathbf{T}) \mathbf{x}_j = \langle \mathbf{x}_i, \mathbf{x}_j \rangle_M$$

Réciproque : Pour toute matrice définie positive \mathbf{M} , il existe une matrice définie positive \mathbf{T} telle que $\mathbf{M} = \mathbf{T}^\top \mathbf{T}$. On notera improprement $\mathbf{T} = \mathbf{M}^{\frac{1}{2}}$.

RAPPELS ÉLÉMENTAIRES D'ALGÈBRE LINÉAIRE

Valeurs et vecteurs propres

Définition : Une matrice A à coefficients dans un corps \mathbb{K} est diagonalisable sur ce corps \mathbb{K} s'il existe une matrice inversible P et une matrice diagonale D à coefficients dans \mathbb{K} telles que

$$A = P D P^{-1}$$

Chaque colonne p de P est un vecteur propre de M , c'est à dire qu'il existe λ sur la diagonale de D tel que

$$A p = \lambda p$$

RAPPELS ÉLÉMENTAIRES D'ALGÈBRE LINÉAIRE

Valeurs et vecteurs propres

Propriété : Toute matrice symétrique réelle est diagonalisable sur \mathbb{R} par une matrice orthogonale P , c'est à dire telle que

$$P^{\top} P = I$$

Propriété : Toute matrice M -symétrique réelle ($A^{\top} M = M A$) est diagonalisable sur \mathbb{R} par une matrice M -orthogonale P , c'est à dire telle que

$$P^{\top} M P = I$$

RAPPELS ÉLÉMENTAIRES D'ALGÈBRE LINÉAIRE

Valeurs et vecteurs propres : cas de la matrice ΣM

Valeurs propres : La matrice ΣM est M -symétrique. Elle est donc diagonalisable sur \mathbb{R} . Ses valeurs propres sont positives, et l'on note

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$$

Vecteurs propres : Les vecteurs propres de ΣM sont M -orthogonaux.

Lien avec l'inertie : On sait que

$$\text{Trace}(\Sigma M) = \lambda_1 + \lambda_2 + \dots + \lambda_k + \dots + \lambda_p$$

En conservant l'information relative au sous-espace propre $\{\lambda_1, \dots, \lambda_k\}$, on conserve l'inertie $\lambda_1 + \lambda_2 + \dots + \lambda_k$.