

Analyse des données

10 janvier 2022

Table des matières

1	Introduction	2
1.1	Les différents types de variables	2
1.2	Description d'une variable quantitative	3
1.3	Description d'une variable qualitative	3
1.4	Relation entre deux variables quantitatives	4
1.5	Relation entre une variable quantitative expliquée et un ensemble de variables quantitatives explicatives	4

1 Introduction

Les méthodes statistiques sont utilisées dans presque tous les secteurs de l'activité humaine et font parti des connaissances de base de l'ingénieur, du gestionnaire, de l'économiste, du chercheur . . .

A l'issue de la phase de recueil de données, la démarche statistique consiste à traiter et interpréter les informations recueillies. Cette démarche comporte suivant deux grands aspects, d'abord l'aspect descriptif ou exploratoire et dans un second temps l'aspect inférentiel ou décisionnel.

La statistique exploratoire a pour objet de synthétiser, résumer et structurer l'information contenue dans les données et utilise pour cela des représentations des données sous forme de tableau, graphiques et ou indicateurs.

La phase décisionnelle arrive après la phase exploratoire, mais nous n'allons pas la traiter

Depuis plus de 50 ans les méthodes d'analyse de données ont largement démontré leur efficacité dans l'étude de gros volumes de données grâce à l'informatique. Les méthodes dites **multi-dimensionnelles** telles que l'analyse en composante principale (ACP) ou l'analyse factorielle des correspondances (AFC) permettent la mise en relation de nombreuses variables entre elles.

Les méthodes multi-dimensionnelles permettent d'obtenir des représentations graphiques qui constituent le meilleur résumé possible de l'information dans un gros tableau de données.

Pour cela le statisticien consent à une perte d'information afin de gagner en lisibilité, il est alors en mesure de faire apparaître les principaux phénomènes qu'il cherche à analyser.

Il est possible de diviser les principales méthodes d'analyse de données en deux grands groupes :

Les méthodes de classification qui ont pour objet de réduire la taille de l'ensemble des individus en formant des groupes homogènes

Les méthodes factorielles cherchent à réduire le nombre de variables en les résumant par un petit nombre de composants synthétiques. Selon que les variables soient quantitatives ou qualitatives, on utilisera soit l'analyse en composante principale, soit l'analyse factorielle des correspondances.

Le choix d'une méthode statistique dépend de l'objectif que l'on se fixe, mais également de la nature des variables, ainsi que le nombre de variables.

1.1 Les différents types de variables

Une variable statistique décrit une caractéristique pour les différents individus. On distingue deux grands types de variables, les variables qualitatives et quantitatives.

- Les variables quantitatives qualifient des quantités, elles peuvent se classer en variables continues ou discrètes.

- Les variables qualitatives définissent une caractéristique (des qualités), elles peuvent être de deux natures : nominales ou ordinales.

1.2 Description d'une variable quantitative

Une variable quantitative est décrite par l'ensemble des valeurs qu'elle prend pour n individus. Afin de synthétiser l'information, on peut utiliser la moyenne ou la variance.

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

Par définition une variable est centrée si sa moyenne est nulle et est réduite et, si sa variance est égale à 1. L'intérêt de centrer et réduire permet de donner une interprétation géométrique au coefficient de corrélation linéaire.

1.3 Description d'une variable qualitative

Pour décrire les variables qualitatives, on utilise des fréquences relatives et les tableaux disjonctifs complets.

Considérons la variable qualitative «couleur des yeux» avec trois modalités : bleu, vert et marron.

Individu	1	2	3	4	5	6
Couleur des Yeux	Bleu	Vert	Marron	Vert	Bleu	Bleu

L'effectif est le nombre d'individus ayant une modalité dans l'échantillon. On peut donc à partir de l'effectif, calculer la fréquence relative qui est l'effectif divisé par le nombre total d'individus.

La présentation d'une variable qualitative sous sa forme disjonctive complète est celle qui se prête le mieux à des calculs statistiques. Cette dernière s'obtient en définissant une variable indicatrice pour chacune des modalités de la variable («B» pour bleu ...).

Tableau disjonctif complet des données ci-dessus :

$$X = \begin{matrix} & \begin{matrix} B & V & M \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{matrix} & \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{bmatrix} \end{matrix}$$

Propriétés du tableau disjonctif :

- La somme des colonnes du tableau complet est égale à un vecteur colonne de dimension n dont tous les éléments sont égaux à 1. Chaque individu possède une seule modalité, cela signifie que sur une ligne de données ne figure que des 0, à l'exception d'un élément unique égal à 1.

- Le produit matriciel de la transposée de la matrice X fois X est une matrice diagonale dont les éléments sont les effectifs de chacune des modalités. Exemple avec les données précédentes :

$$X'X = \begin{bmatrix} 1 & 0 & 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 3 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

1.4 Relation entre deux variables quantitatives

Une variable x prenant n valeurs peut être représentée par un vecteur dans un ensemble \mathbb{R}^n . Dans cet espace, le produit scalaire entre deux vecteurs \vec{x} et \vec{y} , est égal à la somme pour l'ensembles des individus de x et y

$$\vec{x} \cdot \vec{y} = \sum_{i=1}^n x_i y_i$$

En statistiques, le produit scalaire utilisé est :

$$\vec{x} \cdot \vec{y} = \frac{1}{n} \sum_{i=1}^n x_i y_i$$

Ce produit scalaire permet de donner une interprétation géométrique du coefficient de corrélation linéaire. Le cosinus de l'angle formé entre les deux variables est égal au coefficient de corrélation entre ces deux variables tel que : $\cos(X,Y) = r_{X,Y}$.

Si le coefficient de corrélation est égal a 1, les deux vecteurs sont colinéaires (les valeurs prises par x_i et y_i sont proportionnelles). L'absence de corrélation se traduit par un coefficient de corrélation nul et par un angle droit entre x et y .

1.5 Relation entre une variable quantitative expliquée et un ensemble de variables quantitatives explicatives