

GEOMETRIE EUCLIDIENNE ET STATISTIQUE

Position du problème

On considère une variable aléatoire X définie sur une population. On suppose que X est de loi normale, $N(\mu, \sigma)$.

On prélève dans cette population un échantillon aléatoire et simple de taille n . Il en résulte n variables aléatoires, X_1, X_2, \dots, X_n indépendantes et de même loi que X .

$$\text{On note : } \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \quad \text{et} \quad S^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 .$$

Nous allons démontrer que la variable aléatoire $\frac{nS^2}{\sigma^2}$ suit la loi de χ^2 à $(n-1)$ degrés de liberté.

Rappel

Par définition une variable aléatoire suit la loi de χ^2 à v degrés de liberté si elle est la somme des carrés de v variables aléatoires normales, centrées, réduites et indépendantes.

Principe

On va travailler dans l'espace euclidien \mathbf{R}^n muni de son produit scalaire usuel.

$$\text{On considère le vecteur aléatoire } V = \begin{pmatrix} \frac{X_1 - \bar{X}}{\sigma} \\ \frac{X_2 - \bar{X}}{\sigma} \\ \cdot \\ \cdot \\ \frac{X_n - \bar{X}}{\sigma} \end{pmatrix} \quad \text{et le vecteur } I = \begin{pmatrix} 1 \\ 1 \\ \cdot \\ \cdot \\ 1 \end{pmatrix}$$

$$V \cdot I = \frac{1}{\sigma} ((X_1 - \bar{X}) + (X_2 - \bar{X}) + \dots + (X_n - \bar{X})) = \frac{1}{\sigma} (X_1 + X_2 + \dots + X_n - n\bar{X}) = 0$$

V est donc orthogonal à I . Il en résulte que V appartient à l'hyperplan (sous-espace vectoriel de dimension $n-1$) de \mathbf{R}^n orthogonal à I .

$$\text{D'autre part } \|V\|^2 = \frac{1}{\sigma^2} ((X_1 - \bar{X})^2 + (X_2 - \bar{X})^2 + \dots + (X_n - \bar{X})^2) = \frac{nS^2}{\sigma^2}$$

La démonstration consiste à utiliser une base orthonormale, $(e_i)_{1 \leq i \leq n}$, de \mathbb{R}^n telle que $e_1 = \frac{1}{\sqrt{n}}I$. On sait qu'une telle base existe et, dans une telle base, V a pour coordonnées

$$(0, X'_2, X'_3, \dots, X'_n) \text{ et } \|V\|^2 = \frac{nS^2}{\sigma^2} = \sum_{i=2}^n X_i'^2.$$

Il restera à démontrer que les X'_i ($2 \leq i \leq n$) sont de loi $N(0, 1)$ et indépendantes.

Quelques rappels d'algèbre linéaire

Pour obtenir une matrice de changement de base on place en colonnes les coordonnées, dans la première base, des vecteurs de la deuxième base. Cette matrice permet de calculer les coordonnées d'un vecteur dans la première base en fonction de celles de ce vecteur dans la deuxième base.

Ce qui nous intéresse ici est la réciproque : calculer les nouvelles coordonnées en fonction des anciennes. Nous avons donc besoin de la matrice inverse de la matrice de changement de base.

Lors du passage d'une base orthonormale à une autre base orthonormale, la matrice de passage est une matrice orthogonale, son inverse est alors égale à sa transposée. Nous utiliserons donc une matrice dont les lignes seront les coordonnées des vecteurs de la nouvelle base.

Exemple de matrice convenant ici :

Considérons la matrice suivante :

$$\begin{pmatrix} 1 & 1 & 1 & 1 & \dots & \dots & 1 & 1 \\ 1 & -1 & 0 & 0 & \dots & \dots & 0 & 0 \\ 1 & 1 & -2 & 0 & \dots & \dots & 0 & 0 \\ 1 & 1 & 1 & -3 & 0 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 1 & 1 & \dots & \dots & 1 & -n+2 & 0 \\ 1 & 1 & 1 & \dots & \dots & 1 & 1 & -n+1 \end{pmatrix}$$

Ses vecteurs-lignes sont 2 à 2 orthogonaux.. pour obtenir une matrice répondant à notre problème il suffit de multiplier chacun d'eux par l'inverse de sa norme.

Démonstration

- Soit une base orthonormale, $(e_i)_{1 \leq i \leq n}$, de \mathbb{R}^n telle que $e_1 = \frac{1}{\sqrt{n}}I$

Pour $1 \leq i \leq n$ et $1 \leq j \leq n$ on notera a_{ij} la $j^{\text{ème}}$ coordonnée de e_i .

Pour tout i , $1 \leq i \leq n$, on a : $\|e_i\| = 1$ d'où $\sum_{j=1}^n a_{ij}^2 = 1$. (1)

D'autre part, pour $i \neq k$, $e_i \cdot e_k = 0$ d'où : $\sum_{j=1}^n a_{ij} \cdot a_{kj} = 0$. (2)

En particulier, pour tout i , $2 \leq i \leq n$, $e_i \cdot e_1 = 0$ d'où : $\sum_{j=1}^n a_{ij} = 0$. (3)

• Avec les notations définies précédemment on obtient :

$$\begin{pmatrix} X'_1 \\ X'_2 \\ \vdots \\ X'_n \end{pmatrix} = \begin{pmatrix} \frac{1}{\sqrt{n}} & \frac{1}{\sqrt{n}} & \dots & \dots & \frac{1}{\sqrt{n}} \\ a_{21} & a_{22} & \dots & \dots & a_{2n} \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & \dots & a_{nn} \end{pmatrix} \begin{pmatrix} \frac{X_1 - \bar{X}}{\sigma} \\ \frac{X_2 - \bar{X}}{\sigma} \\ \vdots \\ \frac{X_n - \bar{X}}{\sigma} \end{pmatrix}$$

$X'_1 = \frac{1}{\sqrt{n}} V \cdot I = 0$, d'où, puisque la nouvelle base est aussi orthonormale, l'égalité :

$$\|V\|^2 = \sum_{i=2}^n X_i'^2 \text{ c'est-à-dire } \frac{nS^2}{\sigma^2} = \sum_{i=2}^n X_i'^2.$$

• D'autre part pour $2 \leq i \leq n$ on a :

$$X'_i = \frac{1}{\sigma} \sum_{j=1}^n a_{ij} (X_j - \bar{X}) = \frac{1}{\sigma} \sum_{j=1}^n a_{ij} X_j - \frac{\bar{X}}{\sigma} \sum_{j=1}^n a_{ij} = \frac{1}{\sigma} \sum_{j=1}^n a_{ij} X_j \text{ d'après (3).}$$

X'_i est donc de loi normale comme combinaison linéaire de variables aléatoires normales indépendantes. De plus :

$$E(X'_i) = \frac{1}{\sigma} \sum_{j=1}^n a_{ij} E(X_j) = \frac{\mu}{\sigma} \sum_{j=1}^n a_{ij} = 0 \text{ d'après (3).}$$

et, les X_j étant indépendantes,

$$V(X'_i) = \frac{1}{\sigma^2} \sum_{j=1}^n a_{ij}^2 V(X_j) = \frac{\sigma^2}{\sigma^2} \sum_{j=1}^n a_{ij}^2 = 1 \text{ d'après (1).}$$

Il en résulte que, pour tout i , $2 \leq i \leq n$, X'_i est de loi normale centrée et réduite.

- Il reste à montrer que les X'_i , pour $2 \leq i \leq n$, sont indépendantes.

Toute combinaison linéaire des X'_i est une combinaison linéaire des X_i et est donc de loi normale. La définition donnée en annexe permet donc de dire que le vecteur aléatoire $(X'_2, X'_3, \dots, X'_n)$ est gaussien. Pour que les X'_i ($2 \leq i \leq n$) soient indépendantes il suffit, d'après le théorème donné en annexe, qu'elles soient 2 à 2 non corrélées.

Calculons donc, pour h différent de k (h et k compris entre 2 et n), la covariance de X'_h et X'_k :

$$\text{Cov}(X'_h, X'_k) = \frac{1}{\sigma^2} \sum_{i=1}^n \sum_{j=1}^n a_{hi} \cdot a_{kj} \text{Cov}(X_i, X_j)$$

Or, pour $i \neq j$, $\text{Cov}(X_i, X_j) = 0$ puisque les X_i sont indépendantes et, pour tout i compris entre 1 et n, $\text{Cov}(X_i, X_i) = V(X_i) = \sigma^2$. Il en résulte :

$$\text{Cov}(X'_h, X'_k) = \frac{1}{\sigma^2} \sum_{j=1}^n a_{hj} \cdot a_{kj} \cdot \sigma^2 = \sum_{j=1}^n a_{hj} a_{kj} = 0 \text{ d'après (2).}$$

Les X'_i ($2 \leq i \leq n$) sont donc indépendantes.

Conclusion : La variable aléatoire $\frac{nS^2}{\sigma^2}$ est la somme des carrés de n - 1 variables aléatoires normales centrées réduites et indépendantes, elle est donc de loi de χ^2 à n - 1 degrés de liberté.

Prolongements possibles :

- On peut également montrer que $\frac{nS^2}{\sigma^2}$ est indépendante de \bar{X} :

$$\text{Cov}(X_i', \bar{X}) = \text{Cov}\left(\frac{1}{\sigma} \sum_{j=1}^n a_{ij} X_j, \frac{1}{n} \sum_{k=1}^n X_k\right) = \frac{1}{\sigma n} \sum_{j=1}^n \sum_{k=1}^n a_{ij} \text{Cov}(X_j, X_k)$$

$$\text{Cov}(X_i', \bar{X}) = \frac{1}{\sigma n} \sum_{j=1}^n a_{ij} \text{Cov}(X_j, X_j) = \frac{\sigma^2}{\sigma n} \sum_{j=1}^n a_{ij} = 0 \text{ d'après (3).}$$

Il en résulte que, pour tout i compris entre 2 et n , X_i' et \bar{X} sont non corrélées et donc indépendantes d'après le théorème cité précédemment. On en déduit que $\frac{nS^2}{\sigma^2}$ qui est la somme des carrés de ces variables est indépendante de \bar{X} .

- En analyse de variance à un facteur :

On démontre de la même manière que, lorsque les moyennes des p populations sont égales, la variable aléatoire $\frac{SCE_F}{\sigma^2}$ suit la loi de χ^2 à $p - 1$ ddl (le même type de matrice est utilisable

avec $n = p$ lorsque les échantillons sont de même taille sinon on doit avoir $a_{1j} = \sqrt{\frac{n_j}{n}}$, les n_j étant les effectifs des échantillons et n l'effectif total).

- La matrice donnée en exemple (ou toute autre construite de la même manière) est utilisée pour les comparaisons multiples de moyennes par une méthode appelée méthode des contrastes lorsque les effectifs des échantillons sont égaux.

Mais tout cela sera pour une autre fois si certains en manifestent le souhait.

Annexe

Définition : Une variable aléatoire définie sur un espace de probabilité et à valeurs dans \mathbb{R}^n est dite normale ou gaussienne si toute combinaison linéaire de ses coordonnées est de loi normale. (remarque : une variable aléatoire constante est considérée comme étant de loi normale d'écart type 0 ou dégénérée)

Théorème : Pour que la suite des coordonnées d'un vecteur normal soit indépendante, il faut et il suffit que la covariance de ce vecteur (matrice des covariances de ses coordonnées) soit une matrice diagonale.

- - - - -