

Tests multiples

(1)

1. Introduction

On mesure le niveau d'expression d'un ensemble de gènes $i=1, \dots, m$ au sein de 2 populations (tissus sains, tissus tumoraux) de tailles n_1 et n_2 . On pose $n = n_1 + n_2$. On se place dans le cadre où m est grand devant n . On cherche à savoir quels sont les gènes différentiellement exprimés dans les deux populations.

On note $X = (X^1, \dots, X^n) = (Y^1, \dots, Y^{n_1}, Z^1, \dots, Z^{n_2}) \in \mathbb{R}^{m \times n}$ et on suppose que $\underbrace{Y^1, \dots, Y^{n_1}}_{\substack{Y \in \mathbb{R}^m \\ \text{iid de même loi que } Y}} \text{ et } \underbrace{Z^1, \dots, Z^{n_2}}_{\substack{Z \in \mathbb{R}^m}}$

On note P la loi de X (v.a. sur $\mathbb{R}^{m \times n}$) et \mathcal{P} la famille à laquelle appartient P et

$$\mu_{i,1}(P) = \mathbb{E}(Y_i) \quad \text{où } Y = \begin{pmatrix} Y_1 \\ \vdots \\ Y_m \end{pmatrix}$$

$$\mu_{i,2}(P) = \mathbb{E}(Z_i) \quad Z = \begin{pmatrix} Z_1 \\ \vdots \\ Z_m \end{pmatrix}$$

On veut décider si $P \in \mathcal{C}_{0,i}$ où $\mathcal{C}_{0,i} = \{P \in \mathcal{P}, \mu_{i,1}(P) = \mu_{i,2}(P)\}$

Pour tout $i=1, \dots, m$ on considère des hypothèses

$$H_{0,i} : \mu_{i,1}(P) = \mu_{i,2}(P) \text{ v.s. } H_{1,i} = \overline{H_{0,i}}$$

On note $\mathcal{H}_0(P) = \{i=1, \dots, m, H_{0,i} \text{ est vérifiée}\}$
= ensemble des hypothèses nulles vérifiées
(= ensemble des gènes non-différentiellement exprimés)

et $\mathcal{H}_1(P) = \{i=1, \dots, m, H_{0,i} \text{ n'est pas vérifiée}\}$

$$= \{1, \dots, m\} \setminus \mathcal{H}_0(P) \\ (= \text{ensemble des gènes différentiellement exprimés})$$

On commence par construire les tests individuels, pour chaque hypothèse $H_{0,i}$. Dans le cas contraire, dans ce paragraphe, on va construire des tests de Student.

On pose $\bar{Y}_i = \frac{1}{n_1} \sum_{j=1}^{n_1} Y_i^j$ et $\bar{Z}_i = \frac{1}{n_2} \sum_{j=1}^{n_2} Z_i^j$

EXO
En supposant que $Y_i^j \sim \mathcal{N}(\mu_{i,1}, \sigma^2)$ et $Z_i^j \sim \mathcal{N}(\mu_{i,2}, \sigma^2)$ et que $(Y^1, \dots, Y^{n_1}) \perp (Z^1, \dots, Z^{n_2})$, montrer que :

$$\sqrt{\frac{n_1 n_2}{n_1 + n_2}} (\bar{Y}_i - \bar{Z}_i) \sim \text{St}(n-2) \text{ si } H_{0,i} \text{ est v\u00e9rifi\u00e9e.}$$
$$\text{o\u00f9 } \hat{\sigma}^2 = \frac{1}{n-2} (S_i^1 + S_i^2) \text{ et } S_i^1 = \sum_{j=1}^{n_1} (Y_i^j - \bar{Y}_i)^2$$
$$S_i^2 = \sum_{j=1}^{n_2} (Z_i^j - \bar{Z}_i)^2$$

On choisit donc la statistique de Test

$$T_i(x) = \sqrt{\frac{n_1 n_2}{n_1 + n_2}} |\bar{Y}_i - \bar{Z}_i| \text{ et on rejette } H_{0,i}$$

lorsque $T_i(x) > s_\alpha$ (o\u00f9 s_α est un seuil \u00e0 fixer suivant la valeur de α)

D\u00e9finissons $T_{P,i}(s) = \mathbb{P}_{X \sim P} (T_i(X) \geq s)$. Si $H_{0,i}$ est v\u00e9rifi\u00e9e (de fa\u00e7on \u00e9quivalente si $P \in \mathcal{A}_{0,i}$) on a :

$$T_{P,i}(s) = \mathbb{P}(|Z_{n-2}| \geq s) \text{ o\u00f9 } Z_{n-2} \sim \text{St}(n-2)$$
$$= 2 \mathbb{P}(Z_{n-2} \geq s) \text{ par sym\u00e9trie de la } \text{St}(n-2)$$
$$= 2(1 - F_{\text{St}(n-2)}(s)) \text{ o\u00f9 } F_{\text{St}(n-2)} \text{ est la f.d.r de la loi } \text{St}(n-2)$$

D\u00e9finition : p-value.
On d\u00e9finit la p-value $p_i(x)$ de ce test comme :

$$p_i(x) = \sup_{P \in \mathcal{A}_{0,i}} T_{P,i}(T_i(x))$$

Dans notre exemple, $\forall P \in \mathcal{A}_{0,i} T_{P,i}(T_i(x))$

$$\text{donc } p_i(x) = 2(1 - F_{\text{St}(n-2)}(T_i(x)))$$

EXO
Montrons que dans notre exemple $T_{P,i}(T_i(x))$ (le sup n'est pas n\u00e9cessaire dans cet exemple) est une loi uniforme sur $[0,1]$ quand $P \in \mathcal{A}_{0,i}$

$$\mathbb{P}(T_{P,i}(T_i(x)) \leq u) = \mathbb{P}(2(1 - F_{\text{St}(n-2)}(T_i(x))) \leq u)$$
$$= \mathbb{P}(1 - \frac{u}{2} \leq F_{\text{St}(n-2)}(T_i(x)))$$

$$\begin{aligned}
&= \mathbb{P}(F_{St(n-2)}^{-1} (1 - \frac{\alpha}{2}) \leq \mathcal{J}_i(x)) \\
&= \mathbb{Q} \left(1 - F_{St(n-2)} \left(F_{St(n-2)}^{-1} \left(1 - \frac{\alpha}{2} \right) \right) \right) \\
&= \alpha. \quad \square
\end{aligned}$$

car $P \in \mathcal{D}_0$
donc $\mathcal{J}_i(x) \sim St(n-2)$

Cette propriété est très souvent vérifiée et est fondamentale pour les tests multiples.

On passe maintenant aux tests multiples, on aura donc un vecteur $(p_i(x))_{i=1, \dots, m}$ de p-values et on va rejeter les hypothèses H_{0i} pour lesquelles les p-values associées sont "petites" (\leftarrow à définir).

2. Propriété fondamentale des p-values la propriété de l'exercice précédent

On considère dans ce paragraphe que $m=1$.

On adopte donc les notations $H_0: \theta \in \mathcal{D}_0$, la statistique de test est $\mathcal{J}(x)$, la zone de rejet $\{\mathcal{J}(x) \geq s\}$

$T_p(s) = \mathbb{P}_{X \sim P}(\mathcal{J}(x) \geq s)$

On note de plus $F_p(s) = \mathbb{P}_{X \sim P}(\mathcal{J}(x) \leq s)$

$F_p^{\circ}(s) = \mathbb{P}_{X \sim P}(\mathcal{J}(x) < s)$

ici $\mathcal{J}(x)$ est générale, par forcément la stat. d'un test de Student

(remarquez qu'on a $F_p^{\circ}(s) \leq F_p(s) \forall s \in \mathbb{R}$, avec égalité fonctionnelle si la loi admet une densité).

On peut alors écrire :

$$\begin{aligned}
T_p(s) &= \mathbb{P}_{X \sim P}(\mathcal{J}(x) \geq s) = 1 - \mathbb{P}_{X \sim P}(\mathcal{J}(x) < s) \\
&= 1 - F_p^{\circ}(s).
\end{aligned}$$

Propriété

la p-value définie par $p(x) = \sup_{P \in \mathcal{D}_0} T_p(\mathcal{J}(x))$ vérifie.

i) $p(x)$ est stochastiquement majorée par la loi $U_{[0,1]}$ ($\forall P \in \mathcal{D}_0 \forall u \in [0,1] \mathbb{P}(p(x) \leq u) \geq u$).

ii) Si F_p est continue pour tout $P \in \mathcal{D}_0$ alors $p(x) = \min \{ \alpha \in [0,1], \mathcal{J}(x) \geq \sup_{P \in \mathcal{D}_0} F_p^{-1}(1-\alpha) \}$ et si \mathcal{D}_0 est un singleton $p(x) \sim U_{[0,1]}$

Point i) Preuve
 $\forall P \in \Theta_0$ et $\alpha \in [0,1]$.

$$\{T_P(\mathcal{Y}(x)) \leq \alpha\} = \{1 - \overset{\circ}{F}_P(\mathcal{Y}(x)) \leq \alpha\} = \{1 - \alpha \leq \overset{\circ}{F}_P(\mathcal{Y}(x))\} \text{ par définition}$$

on veut montrer que $\{T_P(\mathcal{Y}(x)) \leq \alpha\} \subset \{\mathcal{Y}(x) \geq F_P^{-1}(1-\alpha)\}$. En effet, si c'est vérifié on a alors :

$$\begin{aligned} \mathbb{P}(p(x) \leq \alpha) &= \mathbb{P}\left(\sup_{P \in \Theta_0} T_P(\mathcal{Y}(x)) \leq \alpha\right) \quad \left. \begin{array}{l} \text{car sup.} \\ \text{grâce à} \\ \text{l'inclusion} \end{array} \right\} \\ &\leq \mathbb{P}_{x \sim P \in \Theta_0} (T_P(\mathcal{Y}(x)) \leq \alpha) \\ &\leq \mathbb{P}_{x \sim P \in \Theta_0} (\mathcal{Y}(x) \geq F_P^{-1}(1-\alpha)) \\ &\leq \mathbb{P}_{x \sim P \in \Theta_0} (\mathcal{Y}(x) \geq F_P^{-1}(1-\alpha)) = 1 - F_P(F_P^{-1}(1-\alpha)) = \alpha. \end{aligned}$$

où $F_P^{-1}(u) = \inf\{x \in \mathbb{R}, F_P(x) \geq u\}$. (c'est une inverse généralisée).

Montrons donc l'inclusion

$$\{T_P(\mathcal{Y}(x)) \leq \alpha\} = \{1 - T_P(\mathcal{Y}(x)) \geq 1 - \alpha\} = \{\overset{\circ}{F}_P(\mathcal{Y}(x)) \geq 1 - \alpha\}.$$

$$\subset \{F_P(\mathcal{Y}(x)) \geq 1 - \alpha\} = \{\mathcal{Y}(x) \geq F_P^{-1}(1 - \alpha)\} \quad \begin{array}{l} \text{car } \overset{\circ}{F}_P \leq F_P \\ \text{par définition de } F_P^{-1} \end{array}$$

Point ii).

Si tous les $P \in \Theta_0$ sont absolument continues $\overset{\circ}{F}_P = F_P$ l'inclusion précédente est donc une égalité.

Comme $p(x) = \sup_{P \in \Theta_0} T_P(\mathcal{Y}(x))$ on peut également la définir

$$\begin{aligned} \text{comme } p(x) &= \min\{\alpha \in [0,1], \forall P \in \Theta_0, T_P(\mathcal{Y}(x)) \leq \alpha\} \\ &= \min\{\alpha \in [0,1], \forall P \in \Theta_0, \mathcal{Y}(x) \geq F_P^{-1}(1-\alpha)\} \quad \left. \begin{array}{l} \text{grâce} \\ \text{au} \\ \text{fait que} \\ \text{l'inclusion} \\ \text{est une} \\ \text{égalité dans} \\ \text{ce cas.} \end{array} \right\} \\ &= \min\{\alpha \in [0,1], \mathcal{Y}(x) \geq \sup_{P \in \Theta_0} F_P^{-1}(1-\alpha)\}. \end{aligned}$$

Si Θ_0 est un singleton, on a

$$\begin{aligned} \mathbb{P}(p(x) \leq u) &= \mathbb{P}(\min\{\alpha \in [0,1], \mathcal{Y}(x) \geq F_P^{-1}(1-\alpha)\} \leq u) \\ &= \mathbb{P}(\min\{\alpha \in [0,1], F_P(\mathcal{Y}(x)) \geq 1-\alpha\} \leq u) \\ &= \mathbb{P}(F_P^{-1}(F_P(1-u)) \leq u) \\ &= \mathbb{P}(\sup_{P \in \Theta_0} T_P(\mathcal{Y}(x)) \leq u) = \mathbb{P}_{x \sim P} (T_P(\mathcal{Y}(x)) \leq u) \\ &= \dots = u. \end{aligned}$$

donc $p(x) \sim U_{[0,1]}$ D.

③ Cadre général des tests multiples. ($m > 1$)

On note $\#(\mathcal{H}_0(P)) = m_0(P)$ nombre d'hypothèse nulle vérifiées par P. A partir de l'observation de X on veut dé couvrir $\mathcal{H}_1(P)$.

Note $\#(A) = \text{cardinal}(A)$

On va donc utiliser le vecteur $p = (p_i(x))_{i=1, \dots, m} \in [0, 1]^m$. (5)

avec $\forall P \in \mathcal{P} \forall i \in \mathcal{H}_0(P) \forall u \in [0, 1] \mathbb{P}(p_i(x) \leq u) \leq u$.

On note le vecteur $(p_i(x))_{i=1, \dots, m}$ quand la dépendance en X n'est pas importante.

[Définition: une procédure de tests multiples est une fonction $R : q = (q_i)_{i=1, \dots, m} \in [0, 1]^m \rightarrow R(q) \in \{1, \dots, m\}$

l'ensemble $R(p)$ est donc l'ensemble des hypothèses nulles $\mathcal{H}_{0,i}$ rejetées par la procédure.

On considère dans la suite des procédures particulières, les procédures de seuillage qui prennent la forme.

$$R(q) = \{i = 1, \dots, m, q_i \leq \underline{t}(q)\}$$

Exemple procédure de Bonferroni (de niveau α).

$$R(q) = \{i = 1, \dots, m, q_i \leq \frac{\alpha}{m}\}$$

(4) Erreurs de type I

On consacre l'approche de Neyman-Pearson qui consiste à contrôler l'erreur de type I en essayant de rendre celle de type II la plus petite possible.

Dans un cadre multi-tests, il y a plusieurs notions d'erreurs de type I.

a) k -FWER: k -family error rate.

$$k\text{-FWER}(R(p), P) = \mathbb{P}(\#(R(p) \cap \mathcal{H}_0(P)) \geq k)$$

= proba que R fasse plus de k faux rejets quand P_0 doit être P .

b). FDP = false discovery proportion.

$$FDR(R(p), P) = \frac{\#(R(p) \cap \mathcal{H}_0(P))}{\#(R(p)) \vee 1}$$

\leftarrow car $\#(R(p))$ peut être nul!

= proportion d'erreurs dans l'ensemble des hypothèses rejetées.

⚠ la FDR est une quantité aléatoire: elle dépend de $p = p(x)$

c) FDR = false discovery rate

$$FDR(R(p), R) = \mathbb{E}(FDR(R(p), P))$$

Exemple procédure de Bonferroni $R^{\text{Bonf}}(p) = \{i=1, \dots, m, p_i = p_i(x) \leq \frac{\alpha}{m}\}$

Dans ce cas $\#(R(p) \cap \mathcal{H}_0(p)) = \#\{i \in \mathcal{H}_0(p), p_i(x) \leq \frac{\alpha}{m}\}$

On a donc

$$\text{FDR}(R^{\text{Bonf}}(p), P) = \mathbb{E}\left(\frac{\#(R(p) \cap \mathcal{H}_0(p))}{\#(R(p))} \vee 1\right)$$

$$\leq \mathbb{E}(\#(R(p) \cap \mathcal{H}_0(p)))$$

$$= \sum_{i \in \mathcal{H}_0(p)} \mathbb{P}(p_i(x) \leq \frac{\alpha}{m})$$

$$\leq \alpha \frac{m_0(p)}{m} \leq \alpha$$

par minoration
par 1 du
dénominateur

Exemple: procédure de Holm.

A partir du vecteur des p-values $(p_i)_{i=1, \dots, m}$, on crée celui des p-values ordonnées $p^{(1)} \leq \dots \leq p^{(m)}$ auquel on associe les hypothèses correspondantes $\mathcal{H}_{0,(1)} \dots \mathcal{H}_{0,(m)}$

On définit l'indice minimale k pour lequel $p^{(k)} \geq \frac{\alpha}{m+1-k}$.
On rejette les hypothèses $\mathcal{H}_{0,(1)}, \dots, \mathcal{H}_{0,(k-1)}$ (et on accepte $\mathcal{H}_{0,(k)}, \dots, \mathcal{H}_{0,(m)}$).

On note k^0 la première hypothèse de $\mathcal{H}_0(p)$ rejetée. Avec ces définitions, les hypothèses $\mathcal{H}_{0,(1)}, \dots, \mathcal{H}_{0,(k^0-1)}$ sont toutes rejetées à raison. On a donc $k^0 - 1 \leq m - m_0(p)$

$$\Leftrightarrow \frac{1}{m - k^0 + 1} \leq \frac{1}{m_0(p)} \quad \text{et} \quad p^{(k^0)} \leq \frac{\alpha}{m+1+k^0} \leq \frac{\alpha}{m_0(p)}$$

car (k^0) est rejetée.

On peut alors calculer:

$$\mathbb{P}(\#(R(p) \cap \mathcal{H}_0(p)) \geq 1) = \mathbb{P}(\exists i \in \mathcal{H}_0(p) p_i \leq \frac{\alpha}{m_0(p)})$$

$$\leq \alpha$$

Donc la procédure de Holm contrôle le 1-FWER.

un peu rapide, récrire les détails