

Rappel: mardi 16.12: scéance
réponse aux questions.

Rappel sur les 2-samples t-test:

$X_1, \dots, X_n, Y_1, \dots, Y_n$ deux échantillons

$$T_n(X_1, \dots, X_n, Y_1, \dots, Y_n) = \frac{\sqrt{n} (\bar{X}_n - \bar{Y}_n)}{\sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2 + \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y}_n)^2}}$$

Hypothèses: X_i, Y_i sont gaussiens.

de même variance
↳ motivation F-test

F-tests:

Déf. (lois de Fisher) Soient $h_1, h_2 \geq 1$ deux entiers.

Une variable aléatoire X suit une loi de Fisher à (h_1, h_2) -degrés de liberté, $X \sim \text{Fisher}(h_1, h_2)$, si c'est une variable continue avec densité:

$$f_X(x) = \frac{1}{\Gamma(h_1) \Gamma(h_2)} (x)^{h_1-1} \frac{1}{(1+x)^{h_1+h_2}} \frac{1}{\Gamma(h_2)} \left(\frac{h_2}{1+x} \right)^{h_2-1} \left(1 - \frac{h_2}{1+x} \right)^{h_1-1}$$

$$f_X(x) = \mathbb{1}_{[0,1)}(x) \frac{1}{B(l_1/2, l_2/2)} \frac{1}{x} \left(\frac{l_1 x}{l_1 x + l_2} \right)^{\frac{l_1}{2}} \left(1 - \frac{l_1 x}{l_1 x + l_2} \right)^{\frac{l_2}{2}}$$

où $B(a, b) = \int_0^1 x^{a-1} (1-x)^{b-1} dx$.

Intérêt: Si $X \sim \chi_{l_1}^2$, $Y \sim \chi_{l_2}^2$ sont indép., alors

$$\frac{X/l_1}{Y/l_2} \sim \text{Fisher}(l_1, l_2).$$

Vu en exo:

$$E[X] = \frac{l_2}{l_2 - 2} \quad \text{si } X \sim \text{Fisher}(l_1, l_2)$$

$$\text{Var}(X) = \frac{2 l_2^2 (l_1 + l_2 - 2)}{l_1 (l_2 - 2)^2 (l_2 - 4)} \quad \begin{matrix} l_1 = l_2 = n \\ \sim \frac{1}{n} \end{matrix}$$

2 tests dans le cadre ANOVA

test 1: F-test pour l'égalité de variance

Rappel: X_1, \dots, X_n un échantillon de $\mathcal{N}(\mu, \sigma^2)$

$$\frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \sim \chi_{n-1}^2.$$

$$\hookrightarrow \frac{1}{(n-1)\sigma^2} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \sim \frac{\chi_{n-1}^2}{n-1}$$

Cadre: 2 échantillons X_1, \dots, X_n et Y_1, \dots, Y_m
 de lois gaussiennes: $X_i \sim \mathcal{N}(\mu_x, \sigma_x^2)$
 et $Y_i \sim \mathcal{N}(\mu_y, \sigma_y^2)$.

On suppose $(X_i)_i$ et $(Y_i)_i$ sont
 indép.

But: tester H_0 : "Var(X_i) = Var(Y_i)" = " $\sigma_x^2 = \sigma_y^2$ "

$$\Theta = \mathbb{R} \times \mathbb{R} \times (0, \infty) \times (0, \infty), \quad \theta \in \Theta$$

$$\theta = (\mu_x, \mu_y, \sigma_x^2, \sigma_y^2)$$

$$\Theta_0 = \{(\mu_x, \mu_y, \sigma_x^2, \sigma_y^2) \in \Theta : \sigma_x^2 = \sigma_y^2\}$$

↳ on va utiliser la statistique
 de test:

$$T(X_1, \dots, X_n, Y_1, \dots, Y_m) = \frac{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2}{\frac{1}{m-1} \sum_{j=1}^m (Y_j - \bar{Y}_m)^2}$$

Si H_0 est vraie, $\sigma_y^2 = \sigma_x^2$

$$T(X_1, \dots, X_n, Y_1, \dots, Y_m) = \frac{\frac{1}{(n-1)\sigma_x^2} \sum_{i=1}^n (X_i - \bar{X}_n)^2}{\frac{1}{(m-1)\sigma_y^2} \sum_{i=1}^m (Y_i - \bar{Y}_m)^2} \cdot \frac{\sigma_x^2}{\sigma_y^2}$$

Sous H_0 , $\frac{\sigma_x^2}{\sigma_y^2} = 1$,

\sim Fisher($n-1, m-1$)

Fischer (n-1, m-1)

Si H_1 est vraie : $\sigma_x^2 \neq \sigma_y^2$

$$T(X_1, \dots, X_n, Y_1, \dots, Y_m) \sim \frac{\sigma_x^2}{\sigma_y^2} \text{Fischer}(n-1, m-1)$$

Si $n, m \gg 1$, $\text{Fisher}(n-1, m-1)$ se concentre autour de 1

\rightsquigarrow Sous H_0 : $T \in [1-\varepsilon, 1+\varepsilon]$ avec très grande proba

Sous H_1 : $|T-1| > 0$ carif sur n, m avec grande proba.

\rightsquigarrow On va utiliser une région de rejet de la forme

$$\left\{ T(X_1, \dots, X_n, Y_1, \dots, Y_m) \notin [1 - \delta_\alpha, 1 + \delta'_\alpha] \right\}$$

avec $\delta_\alpha, \delta'_\alpha$ t.o.q.

$$- \mathbb{P}(\text{Fisher}(n-1, m-1) \leq 1 - \delta_\alpha) = \frac{\alpha}{2}$$

$$- \mathbb{P}(\text{Fisher}(n-1, m-1) \geq 1 + \delta'_\alpha) = \frac{\alpha}{2}$$

Sous H_0 :

$$\sup_{\theta \in \Theta_0} \mathbb{P}(T(X_1, \dots, X_n, Y_1, \dots, Y_m) \notin [1 - \delta_\alpha, 1 + \delta'_\alpha]) = \frac{\alpha}{2} + \frac{\alpha}{2} = \alpha$$

\rightsquigarrow test de niveau de risque α .

Commande utile : si vous utilisez R :

trouver z t.q. $\mathbb{P}(\text{Fisher}(n, m) \leq z) = \alpha$

\hookrightarrow $qf(\alpha, n, m)$

2^{em} test : "one way ANOVA"

But : tester si l'appartenance à un certain groupe modifie la valeur d'un paramètre.

Idee : on a une population Pop (ensemble fini, très grand). On la divise en groupes

$$\text{Pop} = G_1 \cup \dots \cup G_h, \quad h \geq 2$$

$$G_i \cap G_j = \emptyset \quad \text{si } i \neq j.$$

Par expl : $\text{Pop} = \{\text{étudiants EPFL}\}$

$G_i = \{\text{étudiants qui passe entre } i-1 \text{ et } i \text{ heures / jour sur Télé Toc}\}$

On s'intéresse à une certaine quantité sur Pop.

expl : - nombre heures sommeil / nuit

expl: - nombre heures sommeil / nuit
- la note moyenne.

But: déterminer si l'appartenance à différents groupes a une influence sur la valeur de la quantité d'intérêt.

Formellement: On suppose que l'on a
Pop un ens., $Pop = G_1 \cup \dots \cup G_k$,
On a X_i , $i \in Pop$ une famille de
v.a. i.i.d.

On suppose: $X_i \sim \mathcal{N}(\mu_i, \sigma^2)$

$\mu_i \in \mathbb{R}$, $i \in Pop$, $\sigma^2 \in (0, +\infty)$ est
la même pour tous les ind.

$$\Theta = \mathbb{R}^{Pop} \times (0, +\infty), \quad \theta \in \Theta, \quad \theta = ((\mu_i, i \in Pop), \sigma^2)$$

On va tester l'hypothèse H_0 "les moyennes de chaque groupe sont égales à la moyenne générale de la population".

$$\Theta_0 = \{(\mu_i, i \in Pop, \sigma^2) : \frac{1}{|G_1|} \sum_{i \in G_1} \mu_i = \frac{1}{|G_2|} \sum_{i \in G_2} \mu_i = \dots\}$$

$$= \frac{1}{|Pop|} \sum_{i \in Pop} \mu_i \quad \left. \vphantom{\sum} \right\}$$

On introduit :

les échantillons $X_{1,1}, \dots, X_{1,n_1}$ dans G_1
 $X_{2,1}, \dots, X_{2,n_2}$ dans G_2
 \vdots
 $X_{h,1}, \dots, X_{h,n_h}$ dans G_h

$$N = n_1 + n_2 + \dots + n_h.$$

$$\bar{X}_{l,n_l} = \frac{1}{n_l} \sum_{i=1}^{n_l} X_{l,i} \quad \text{moyenne empirique } G_l.$$

$$\bar{X} = \frac{1}{N} \sum_{l=1}^h \sum_{i=1}^{n_l} X_{l,i} \quad \text{moyenne générale.}$$

$$\text{In Grp Var} = \frac{1}{N-h} \sum_{l=1}^h \sum_{i=1}^{n_l} (X_{l,i} - \bar{X}_{l,n_l})^2$$

$$\text{Bet Grp Var} = \frac{1}{h-1} \sum_{l=1}^h n_l (\bar{X}_{l,n_l} - \bar{X})^2$$

Statistique de test : $T(-) = \frac{\text{Bet Grp Var}}{\text{In Grp Var}}$

Sous H_0 , on a que toutes les moyennes sont identiques

$$\rightsquigarrow X_{l,i} - \bar{X}_{l,n_l} \stackrel{\text{marginement}}{=} \mathcal{N}(0, \sigma^2)$$

$$\text{et } \bar{X}_{l,n_l} - \bar{X} \stackrel{\text{moyenne}}{=} \mathcal{N}\left(0, \frac{\sigma^2}{n_l}\right)$$

$$\text{et } \bar{X}_{l,n_l} - \bar{X} \stackrel{\text{normal}}{=} \mathcal{N}\left(0, \frac{\sigma^2}{n_l}\right)$$

$\Rightarrow T$ est d'ordre 1.

Il se trouve que sous H_0 ,

$$T(-) \sim \text{Fisher}(h-1, N-h)$$

Sous H_1 : on a au moins un $l^* \in \{1, \dots, h\}$

$$\text{tel que } \underbrace{\frac{1}{|G|} \sum_{i \in G} \mu_i}_{\bar{\mu}_{l^*}} \neq \underbrace{\frac{1}{|P|} \sum_{i \in P} \mu_i}_{\bar{\mu}}$$

$$\text{Bet Grp Var} = \frac{1}{h-1} \sum_{l=1}^h n_l (\bar{X}_{l,n_l} - \bar{X})^2$$

$$\geq \frac{1}{h-1} n_{l^*} \underbrace{(\bar{X}_{l^*,n_{l^*}} - \bar{X})^2}_{\longrightarrow (\bar{\mu}_{l^*} - \bar{\mu})^2 > 0}$$

$$n_1, \dots, n_h \rightarrow \infty$$

\hookrightarrow Sous H_1 , T est la mult d'un nombre > 0

$$\text{avec } n_l \approx \frac{N}{h} \xrightarrow{N \rightarrow \infty} \infty$$

Suggère la région de rejet

$$\left\{ \frac{\text{Bet Grp Var}}{\text{In Grp Var}} > S_\alpha \right\}$$

$$S_\alpha \text{ t.q. } \mathbb{P}(\text{Fisher}(h-1, N-h) > S_\alpha) = \alpha$$