

Statistiques Appliquées

TD 5

Théorie d'échantillonnage : simulations avec

Remarque 1 Dans ce TD on reprend le code du 4.3, « un grand nombre d'échantillons », et on l'enrichit afin d'étudier la distribution des statistiques d'un ou deux échantillons. Toutes les questions qui suivent se réfèrent alors au script du 4.3 qu'il faudra modifier astucieusement.

5.1 Distribution de la moyenne : Z

Selon la théorie d'échantillonnage, la moyenne de l'échantillon, \bar{X} , est liée à la v.a. $Z = (\bar{X} - \mu_X)/(\sigma_X/\sqrt{n})$ (normale $N(0, 1)$), selon des conditions présentées dans le cours.

5.1.1 Calcul de z et étude qualitative

- Créer un vecteur vide `z.exp.vec`, de longueur k , pour stocker les valeurs de Z obtenues expérimentalement (par simulation).
- À l'intérieur de la boucle `for`, calculer z pour chaque échantillon et stocker les résultats dans le vecteur `z.exp.vec` déjà défini.
- Après la fin des itérations, générer seulement un histogramme de *densités* (option `freq = FALSE` : pour chaque classe, afficher la fréquence relative divisée par la largeur de la classe) de `z.exp.vec`.

Remarque 2 Afin qu'un histogramme de densités soit une bonne estimation d'une ddp, il faut avoir un grand nombre d'échantillons, p.ex. $k = 100^3$ (c'est l'équivalent de laisser tourner `CenLimit` pendant très longtemps).

- Superposer à l'historogramme obtenu la ddp théorique de Z ($N(0, 1)$) : créer un vecteur `z.th.vec = seq(from=min(z.exp.vec), to=max(z.exp.vec), length=500)` et utiliser la fonction `lines` pour tracer $p_Z(z)$ sur la même fenêtre que l'historogramme.
- Est-ce que la ddp expérimentale de Z (histogramme de `z.exp.vec` avec un grand nombre d'observations) donne une bonne approximation de la ddp théorique de Z ? Quelles classes de l'histogramme suivent mieux la ddp théorique?

5.1.2 Étude quantitative

On veut maintenant *quantifier* les éventuels écarts entre l'historogramme (obtenu de façon expérimentale) et la ddp théorique. L'idée principale est de découper l'historogramme en plusieurs parties (classes) et comparer la fréquence (ou

la fréquence relative) de chaque classe à celle prévue par la théorie. Une façon naturelle de faire ce découpage est de choisir $M + 1$ classes qui, en théorie, ont les mêmes effectifs, c'est-à-dire des fréquences relatives égales à $1/(M + 1)$.

Il faut alors trouver les M valeurs de z qui partagent la ddp théorique en $M + 1$ parties, chacune contenant une probabilité égale à $1/(M + 1)$, et les utiliser dans l'argument `breaks=` de la fonction `hist`, en ajoutant, à gauche et à droite les bornes de la première et de la dernière classe. En théorie, ces limites sont $-\infty$ et $+\infty$ mais, en pratique, il suffit de prendre `min(z.exp.vec)` et `max(z.exp.vec)` pour avoir toutes les valeurs de `z.vec` représentées dans l'histogramme.

- a. Créer le vecteur `q.m.th` contenant les M valeurs de z qui définissent les frontières entre les classes. On prendra d'abord $M = 1$ ou $M = 3$ et on gardera par la suite la valeur $M = 99$.
- b. Créer le vecteur `bornes=c(min(z.exp.vec), q.m.th, max(z.exp.vec))` contenant les bornes des $M + 1$ classes.
- c. Créer l'histogramme de densités de `z.exp.vec` et récupérer les paramètres calculées dans la variable `hist.data` avec la commande
`hist.data = hist(z.exp.vec, freq=FALSE, breaks=bornes)`
 (Comme les classes n'ont pas la même largeur, R crée par défaut un histogramme de densités, même sans l'option `freq=FALSE`.)
- d. La variable `hist.data$counts` contient la fréquence de chaque classe ; il suffit de la diviser par le nombre total `sum(hist.data$counts)` (ici égal à k) pour obtenir les fréquences relatives dans un vecteur `freq.rel.exp`.
- e. Les fréquences relatives théoriques sont égales à
`freq.rel.th = 1 / (M+1)`
 alors que les fréquences théoriques sont égales à
`freq.th = k / (M+1)`
 puisqu'on possède au total k valeurs de `z.exp`.
- f. Sur un histogramme de *fréquences*, les classes ainsi définies devraient donner toutes la même fréquence. Comme les classes n'ont pas la même largeur, il faut demander explicitement à R d'afficher les fréquences avec l'option `freq=TRUE`, sinon il affiche par défaut les densités (ce qui est le bon choix normalement ; en plus il génère un *warning* pour indiquer à l'utilisateur que ce n'est pas une bonne idée de représenter des fréquences alors que les largeurs des classes ne sont pas toutes égales !).
 Afficher l'histogramme de fréquences de `z.exp.vec` avec les bornes définies précédemment et ajouter une ligne horizontale au niveau de la fréquence théorique, égale pour toutes les classes.
- g. On peut finalement créer le vecteur de l'écart relatif entre les fréquences relatives expérimentales et les fréquences relatives théoriques :
`ecart.rel.freq.rel = (freq.rel.exp - freq.rel.th) / freq.rel.th`
- h. Afficher les valeurs absolues de l'écart relatif
`plot(abs(ecart.rel.freq.rel))`
- i. Calculer la norme du vecteur ; elle représente une mesure de la distance entre les fréquences relatives expérimentales et théoriques
`norm.ecart.rel.freq.rel = sqrt(sum(ecart.rel.freq.rel^2))`
 le maximum de sa valeur absolue

```
max.abs.ecart.rel.freq.rel = max( abs( ecart.rel.freq.rel ) )  
et la position de ce maximum  
which.max.abs.ecart.rel.freq.rel = which.max( abs( ecart.rel.freq.rel ) )
```

- j. Quelles classes de l'histogramme présentent les écarts les plus importants entre fréquences relatives expérimentales et théoriques ? Comparer avec la réponse de la question e. de 5.1.1. Expliquer la différence.

Remarque 3 *On verra plus tard, dans le cadre des « tests du χ^2 », une autre façon de traiter ces écarts entre les fréquences expérimentales et théoriques.*

5.1.3 Étude qualitative (deuxième approche) : quantiles vs. quantiles

L'approche suivie dans 5.1.2 n'est pas toujours applicable. Elle suppose qu'on sait non seulement quelle est la loi théorique utilisée pour la comparaison (ici la loi normale) mais, en plus, quels sont ses paramètres de position (ici $\mu = 0$) et de dispersion (ici $\sigma = 1$). Le plus souvent, on n'a pas d'information sur ce deuxième point et on veut juste savoir si une ddp expérimentale suit une loi, en général (p.ex. la loi normale).

De façon très simple, on veut examiner si une variable aléatoire Y (dont on obtient par mesure un échantillon de k valeurs) est liée à une variable aléatoire X (la loi théorique, avec μ et σ « standard », choisis par nous) par l'intermédiaire d'une relation linéaire $Y = aX + b$. Car si une telle relation existe, alors Y suit la même loi que X mais avec une espérance et un écart-type potentiellement différents de ceux choisis pour X (cf. cours, transformations linéaires et v.a. centrée réduite).

Donc le problème d'origine (comparer une ddp expérimentale à une autre théorique) revient à examiner l'existence d'une relation linéaire entre deux variables aléatoires. Or, on sait que si une telle relation existe, alors les valeurs des v.a. la suivent aussi, $y = ax + b$, et, par conséquent, les valeurs « spéciales » $y_\alpha = ax_\alpha + b$, c'est-à-dire leurs *quantiles*, se trouvent sur une ligne.

Il suffit donc, pour chacune des valeurs de Y (on en a obtenu un ensemble de k) de trouver la valeur de α correspondant (c'est-à-dire la proportion des valeurs de Y qui sont inférieures à celle examinée), calculer ensuite x_α (la valeur de X qui laisse la même probabilité α à sa gauche) et tracer y_α en fonction de x_α pour les différentes valeurs de α . Si le résultat est une ligne droite, alors Y suit la même loi que X .¹

Avec R, tout ça se résume en *trois commandes* ; si y est un vecteur contenant des valeurs de Y ,

- `qqnorm(y)` affiche le plot des quantiles expérimentales de Y en fonction des quantiles théoriques de X , où X est la loi normale centrée réduite.
- `qqline(y)` ajoute par la suite une ligne qui passe par les points $(x_{0.25}, y_{0.25})$ et $(x_{0.75}, y_{0.75})$, les premiers et les troisièmes quartiles.
- Si on active la librairie “car” (Companion to Applied Regression) avec la commande `library(car)` on peut ensuite utiliser la commande `qq.plot(y, distribution="norm")` où, à la place de `norm` (par défaut) on peut demander une autre loi connue par R. La commande `qq.plot` (appelée aussi `qqp`) ajoute une ligne qui, par défaut, passe par les premiers

¹En plus, comme $\mu_Y = a\mu_X + b$, $\sigma_Y = |a|\sigma_X$ et μ_X , σ_X sont connus, on peut utiliser l'intercepte b et la pente a de la ligne pour obtenir une estimation des paramètres μ_Y et σ_Y !

et troisièmes quartiles. Il s'agit d'une commande beaucoup plus puissante et flexible que les deux premières.²

- Ajouter les commandes qui affichent un Q-Q plot des quantiles du vecteur `z.exp.vec` par rapport aux quantiles de la distribution théorique.
- Est-ce que les quantiles expérimentaux et théoriques s'alignent bien sur une droite ?
- Dans le cas traité ici, on connaît l'espérance et la variance de Z (qui joue le rôle de Y dans la discussion précédente). On peut donc calculer les valeurs de a et b qu'il faudrait obtenir.
 - Avec `abline(h=0)` et `abline(v=0)` ajouter des axes. Estimer (à l'œil!) la valeur de l'intercepte b .
 - Avec `abline(a=..., b=...)` tracer la ligne droite sur laquelle devraient se trouver les points. (Attention! la commande `abline` inverse les a et b , en partant de la relation $y = a + bx$).
- Est-ce que Z suit bien la ddp théorique ? Quelles quantiles expérimentales de Z correspondent mieux aux quantiles théoriques ? Comparer avec les dernières questions de 5.1.1 et 5.1.2.

Remarque 4 *Pour transformer cette étude qualitative en étude quantitative très puissante, il suffit de calculer le coefficient de corrélation linéaire, ρ , entre les y_α et x_α . Des valeurs de ρ proches de l'unité indiquent qu'une relation linéaire existe effectivement entre les quantiles de Y et de X . Cela concernera la partie regression linéaire du cours.*

5.1.4 Étude paramétrique

On veut étudier l'influence de la taille n de l'échantillon sur les écarts entre les fréquences relatives expérimentales et théoriques. Pour cela, on va envelopper la quasi-totalité du code déjà produit dans une boucle qui va balayer plusieurs valeurs du paramètre étudié (taille de l'échantillon).

- Modifier la définition de la variable `n` : au lieu de lui donner une seule valeur, la définir comme un vecteur, p.ex.
`n = seq(from=3, to=33, by=5)`
- Initialiser les variables
`norm.ecart.rel.freq.rel`
`max.abs.ecart.rel.freq.rel`
`which.max.abs.ecart.rel.freq.rel`
comme des vecteurs vides (initialisation à 0) de longueur égale à `length(n)`. Ils serviront à stocker les valeurs respectives pour chaque valeur de la taille de l'échantillon.
- Commencer une boucle `for (index.n in 1:length(n))` just'avant la boucle sur `k`; elle finira après les calculs de la question h. de 5.1.1.
- À l'intérieur de cette boucle remplacer `n` par `n[index.n]`.
- De la même façon, les calculs de la question h. de 5.1.1 seront modifiés.
P.ex. :
`norm.ecart.rel.freq.rel[index.n] = sqrt(sum(ecart.rel.freq.rel^2))`

²Donc tout se resume en une seule commande.

- f. Après la fin de la boucle sur `index.n`, générer des graphiques montrant l'évolution de
- ```
norm.ecart.rel.freq.rel
max.abs.ecart.rel.freq.rel
which.max.abs.ecart.rel.freq.rel
```
- en fonction des valeurs de `n`. P.ex. :
- ```
plot( n, norm.ecart.rel.freq.rel, type="b" )
```
- g. Conclusions, commentaires, etc.