

TD 2 - Analyse descriptive univariante

But du TD : Utiliser des outils de base dans R pour d ecrire des donn ees univariantes quantitatives.

1 Description par des param etres num eriques

Soit $(x_i)_{i=1}^n$ un  echantillon de n donn ees. On calcule divers param etres relatifs   cet  echantillon.

1.1 Moments

La moyenne $\bar{x}_n = \frac{1}{n} \sum_{i=1}^n x_i$ s'obtient sur R gr ce   la fonction `mean`.

La variance est $s_n^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2$. Sur R, la fonction `var` calcule un estimateur sans biais de la variance : $\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}_n)^2$.

On appelle moment d'ordre p la quantit  $m_{n,p} = \frac{1}{n} \sum_{i=1}^n x_i^p$.

  Cr er une fonction des donn ees et de p qui calcule le moment d'ordre p .

On appelle moment centr  d'ordre p la quantit  $\mu_{n,p} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^p$. On retrouve la variance pour $p = 2$. Deux param etres de forme s'obtiennent   partir des moments centr s normalis s d'ordre 3 et 4. $SK_n = \frac{\mu_{n,3}}{(s_n^2)^{3/2}}$ s'appelle le coefficient de dissym trie ou *skewness*. Une grande valeur de $|SK_n|$ d note une assym trie des donn ees, le signe de SK_n montre le type d'assym trie. $K_n = \frac{\mu_{n,4}}{(s_n^2)^2}$ s'appelle le coefficient d'aplatissement ou *kurtosis*. Pour un  echantillon gaussien, le coefficient est proche de 3. Si $K_n > 3$, la distribution a des queues plus importantes que la loi normale. Si $K_n < 3$, la distribution a des queues moins importantes que la loi normale.

  Cr er une fonction des donn ees et de p qui calcule le moment centr  d'ordre p , une fonction *skewness* et une fonction *kurtosis*.

  Cr er un  echantillon N de taille 100 de loi normale de moyenne 2, variance 4, tracer un histogramme et calculer les coefficients de *skewness* et *kurtosis*. Faire la m me chose pour un  echantillon *Exp* de loi exponentielle de param tre 3. Comparer.

1.2 Quantiles

On note $x_{(1)} \leq \dots \leq x_{(n)}$ l' echantillon ordonn . La proportion de donn ees de l' echantillon qui sont inf rieures ou  gales   $x_{(1)}$ est donc $1/n$, la proportion de donn ees de l' echantillon inf rieures ou  gales   $x_{(i)}$ est i/n . La fonction de r partition empirique F_n est d finie comme suit : pour x quelconque,

$$F_n(x) = (\text{Prop. de donn ees} \leq x) = i/n \text{ pour } x_{(i)} \leq x < x_{(i+1)}.$$

La fonction quantile empirique, not e F_n^{-1} , se d finit pour $t \in (0, 1)$ par :

$$F_n^{-1}(t) = x_{(i)} \text{ lorsque } \frac{i-1}{n} < t \leq \frac{i}{n}.$$

On distingue des quantiles particuliers, appel s quartiles : $Q_1 = F_n^{-1}(1/4)$, $Q_2 = F_n^{-1}(1/2)$ (qui est aussi la **m diane**), $Q_3 = F_n^{-1}(3/4)$.

  Tracer dans une m me fen tre deux graphiques correspondant   la fonction quantile empirique pour l' echantillon N et pour l' echantillon *Exp*. Utiliser les fonctions `plot` avec le type "s", `sort`, `length`.

2 Description par des graphiques

2.1 Bo te   moustaches (box-plot)

La bo te   moustaches se calcule   partir des quartiles de l' echantillon. La bo te centrale a pour limites le premier quartile Q_1 et le troisi me quartile Q_3 . On marque la m diane (  l'int rieur de la bo te) par un trait. Les moustaches montrent la plus petite (resp. la plus grande) observation qui se situe   une distance inf rieure   $1,5(Q_3 - Q_1)$ de la bo te. Les donn ees se situant   l'ext rieur des moustaches sont consid r es comme extr mes ou *outliers*. En R, on obtient une bo te   moustaches gr ce   la commande `boxplot`.

  Tracer une bo te   moustaches des donn ees N et *Exp* de la question 1. Commenter.

2.2 Histogrammes

- Avec la fonction `hist`, former un histogramme de N et de Exp .
- En utilisant le paramètre `breaks` de `hist`, proposer un histogramme avec de nouvelles classes.

3 Application à des données disponibles dans la librairie de R

On étudiera au choix les tableaux de données suivants : `iris`, `rock`, `LifeCycleSavings`. Décrire la population étudiée, les données et leur type (qualitative, quantitative, et si quantitative, discrète ou continue). Utiliser les outils suivants pour étudier chaque variable quantitative (en choisir 3) :

- Calculer la moyenne, la médiane, le coefficient de *skewness* et de *kurtosis*. Qualifier les échantillons (symétrie, dispersion).
- Tracer une boîte à moustaches. Commenter.
- Tracer un histogramme des données, commenter. Que dire de la normalité des données, à première vue ?
- Comme dans le TD 1, tracer un graphique quantile-quantile pour étudier la gaussiannité des données. Commenter. Effectuer un test de Shapiro-Wilks et conclure.
- Choisir une variable pour laquelle il peut être intéressant de changer les classes de l’histogramme, et tracer le nouvel histogramme. Commenter.