

Sur l'adéquation à une loi de probabilité avec

Christophe Chesneau

<http://www.math.unicaen.fr/~chesneau/>

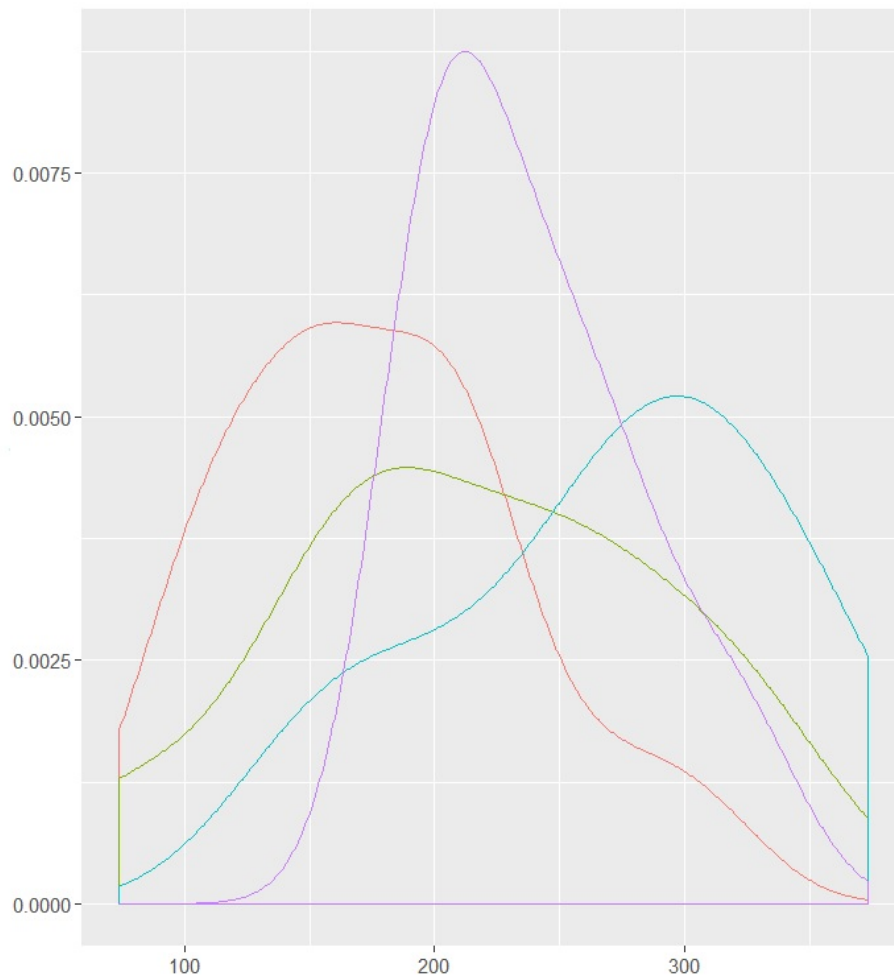


Table des matières

1	Point de départ	5
2	Analyses graphiques	7
3	Tests statistiques d'adéquation à une loi	25
3.1	Test du Chi-deux	25
3.2	Test de Kolmogorov-Smirnov	32
3.3	Test de Shapiro-Wilk	33
4	Exercices	37
5	Solutions	43

~ **Note** ~

L'objectif de ce document est de présenter les principaux outils statistiques et commandes R utilisés pour juger de l'adéquation de la distribution des valeurs d'un caractère à une loi de probabilité.

Contact : christophe.chesneau@gmail.com

Bonne lecture!

1 Point de départ

On observe la valeur d'un caractère X pour chacun des n individus d'un échantillon. Ces observations constituent les données : x_1, \dots, x_n . On modélise alors X comme une *var* (en gardant la notation X par convention). Soit \mathcal{L} une loi de probabilité étant possiblement en adéquation avec la loi inconnue de X . La problématique est la suivante :

Est-ce que ces données nous permettent d'affirmer que X ne suit pas la loi \mathcal{L} (avec un faible risque de se tromper) ?

Pour répondre à cette question, on distingue deux approches complémentaires :

- Analyses graphiques.
- Tests statistiques adaptées reposant sur les hypothèses :

$$H_0 : "X \text{ suit la loi } \mathcal{L}" \quad \text{contre} \quad H_1 : "X \text{ ne suit pas la loi } \mathcal{L}"$$

Dans ce document, nous mettons en œuvre ces tests en utilisant la p-valeur.

La p-valeur est le plus petit réel $\alpha \in]0, 1[$ calculé à partir des données tel que l'on puisse se permettre de rejeter H_0 au risque $100\alpha\%$. Autrement écrit, la p-valeur est une estimation ponctuelle de la probabilité critique de se tromper en rejetant H_0 alors que H_0 est vraie.

Les logiciels actuels travaillent principalement avec cette p-valeur.

2 Analyses graphiques

Cas de caractères qualitatifs ou quantitatifs "discrets"

Soit X un caractère non chiffré (qualitatif) ou chiffré (quantitatif) prenant un ensemble dénombrable de valeurs (possiblement infini). L'analyse graphique la plus pertinente pour juger de l'adéquation de la loi de X avec \mathcal{L} repose sur le schéma suivant :

- On trace le barplot des fréquences correspondantes aux données.
- On superpose les valeurs de la "densité" associée à la loi \mathcal{L} en estimant éventuellement les paramètres inconnus de celle-ci.

Exemple

Exemple 1. On souhaite savoir si les entrées à l'hôpital pour une certaine maladie sont réparties au hasard dans l'année ou bien si certains mois sont plus propices à la maladie. On examine le mois d'entrée d'un échantillon de 120 porteurs de la maladie étudiée. Les résultats sont :

Mois d'entrée	1	2	3	4	5	6	7	8	9	10	11	12
Nombre d'entrées	18	16	8	10	6	4	4	9	11	10	12	12

Peut-on affirmer que "les entrées ne se font pas au hasard dans l'année" (donc que "certains mois sont plus propices à la maladie") ?

Solution 1. Soit X la *var* égale au mois d'entrée à l'hôpital d'un porteur de la maladie. Par l'énoncé, on observe la valeur de X sur chacun des n individus (porteurs de la maladie) d'un échantillon avec $n = 120 : (x_1, \dots, x_n)$ (avec $x_i \in \{1, \dots, 12\}$). On forme alors un vecteur des effectifs

$$(n_1, n_2, \dots, n_{12}) = (18, 16, \dots, 12).$$

Dire que les entrées se font au hasard dans l'année signifie que X suit la loi uniforme $\mathcal{U}(\{1, \dots, 12\})$:

$$\mathbb{P}(X = i) = \frac{1}{12}, \quad i \in \{1, \dots, 12\}.$$

La problématique est la suivante : est-ce que ces données nous permettent d'affirmer que X ne suit pas la loi $\mathcal{L} = \mathcal{U}(\{1, \dots, 12\})$?

Pour tout $i \in \{1, \dots, 12\}$, une estimation (ponctuelle) de $\mathbb{P}(X = i)$ est la fréquence n_i/n .

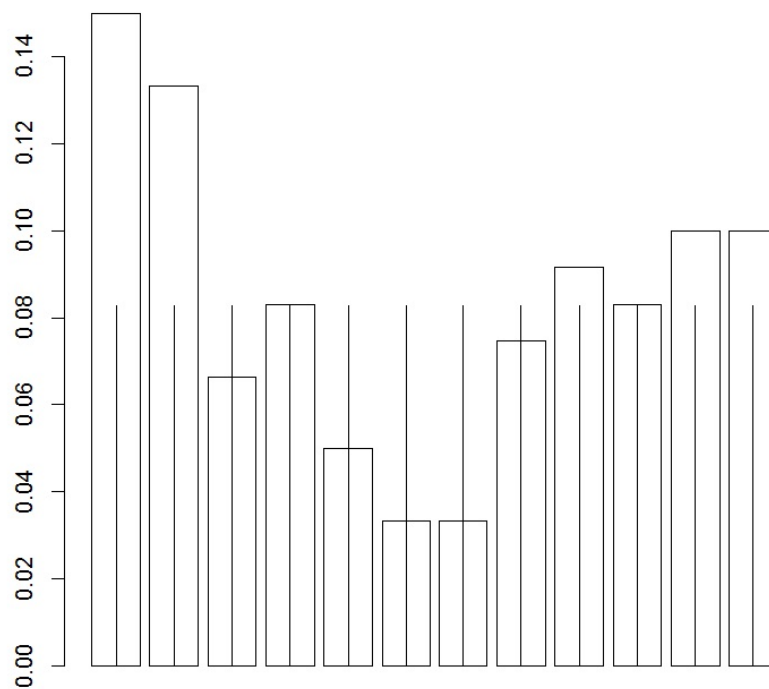
Ainsi, pour une analyse graphique, on peut

- tracer le barplot des fréquences correspondantes aux données,
- superposer les valeurs de la "densité" associée à la loi uniforme $\mathcal{U}(\{1, \dots, 12\})$.

On propose les commandes :

```
nb = c(18, 16, 8, 10, 6, 4, 4, 9, 11, 10, 12, 12)
bar = barplot(nb / 120, col = "white")
points(bar, rep(1 / 12, 12), type = "h")
```

Cela renvoie :



Les différences observées laissent penser que X ne suit pas une loi uniforme.

Exemple 2. Dans un verger, on étudie le comportement des insectes quand ceux-ci attaquent les fruits. Soit X le caractère qui dénombre le nombre d'attaques d'insectes sur un fruit pris au hasard. Une étude statistique antérieure montre que, si les attaques se font de façon indépendante les unes des autres, on peut modéliser X comme une *var* suivant une loi de Poisson $\mathcal{P}(\lambda)$ avec λ inconnu.

On considère un échantillon de 300 fruits et on compte le nombre d'attaques sur le fruit. Les résultats sont :

Nombre d'attaques	0	1	2	3	4	5	6	7
Nombre de fruits attaqués	60	105	65	47	15	4	3	1

Peut-on dire que le comportement des insectes est grégaire? (*on dit que le comportement des insectes est grégaire quand chacun d'entre eux a tendance à se comporter comme le voisin et à attaquer le même fruit; dans ce cas, X ne suit pas une loi de Poisson*).

Solution 2. Par l'énoncé, on observe la valeur de X sur chacun des n individus (fruits) d'un échantillon avec $n = 300 : (x_1, \dots, x_n)$ (avec $x_i \in \mathbb{N}$). On forme alors un vecteur des effectifs

$$(n_1, n_2, \dots, n_8) = (60, 105, \dots, 1).$$

Dire que X suit la loi de Poisson $\mathcal{P}(\lambda)$, avec λ inconnu, signifie que :

$$\mathbb{P}(X = i) = e^{-\lambda} \frac{\lambda^i}{i!}, \quad i \in \mathbb{N}.$$

Pour évaluer la loi de Poisson la plus adaptée à notre contexte, il faut estimer λ à l'aide des données.

Comme $\mathbb{E}(X) = \lambda$, la méthode des moments nous assure qu'une estimation de λ est

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{0 \times 60 + 1 \times 105 + 2 \times 65 + 3 \times 47 + 4 \times 15 + 5 \times 4 + 6 \times 3 + 7 \times 1}{300} = 1.603333.$$

Ainsi, la problématique est la suivante : est-ce que ces données nous permettent d'affirmer que X ne suit pas la loi $\mathcal{L} = \mathcal{P}(1.603333)$?

Pour tout $i \in \{0, \dots, 7\}$, une estimation de $\mathbb{P}(X = i)$ est la fréquence n_{i+1}/n et une estimation de $\mathbb{P}(X \geq 8) = 1 - \mathbb{P}(X \leq 7)$ est :

$$1 - \sum_{i=1}^8 \frac{n_i}{n}.$$

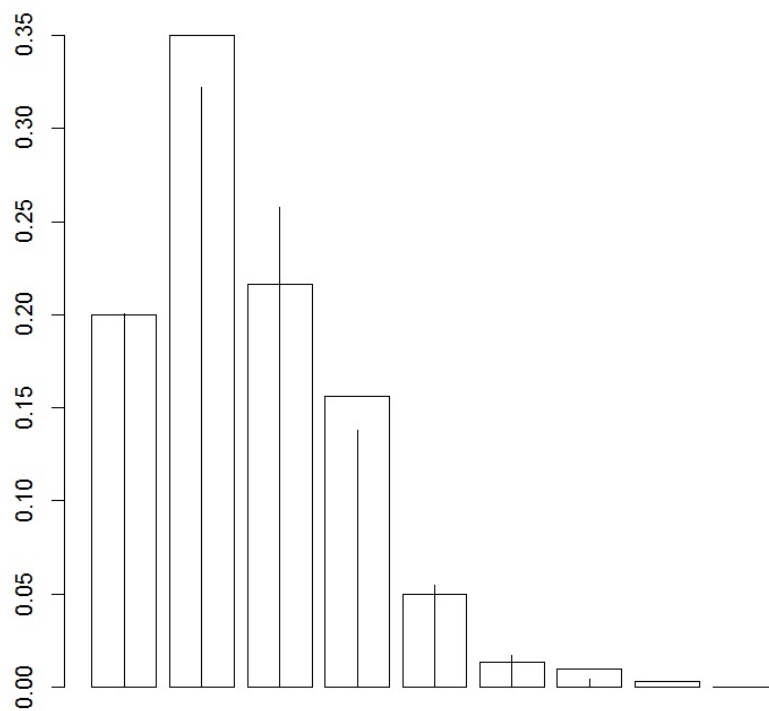
Ainsi, pour une analyse graphique, on peut

- tracer le barplot des fréquences correspondantes aux données,
- superposer les valeurs de la "densité" associée à la loi de poisson $\mathcal{P}(1.603333)$.

On propose les commandes :

```
nb = c(60, 105, 65, 47, 15, 4, 3, 1, 0)
bar = barplot(nb / 300, col = "white")
lambda = sum((0:8) * nb) / 300
prob = c(dpois(0:7, lambda), 1 - ppois(7, lambda))
points(bar, prob, type = "h")
```

Cela renvoie :



Les différences observées laissent penser que X suit une loi de Poisson, remettant ainsi en cause le comportement grégaire des insectes.

Cas de caractères quantitatifs "continus"

Soit X un caractère chiffré (quantitatif) prenant un ensemble indénombrable de valeurs. Les analyses graphiques possibles pour juger de l'adéquation de la loi de X avec \mathcal{L} sont nombreuses. Les méthodes usuelles sont :

- **Méthode de l'histogramme** : On trace l'histogramme des fréquences correspondantes aux données. Puis on superpose les valeurs de la "densité" associée à la loi \mathcal{L} en estimant éventuellement les paramètres inconnus de celle-ci.

- **Méthode de la fonction de répartition** : On trace le graphe de la fonction de répartition empirique définie par :

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{x_i \leq x\}}, \quad x \in \mathbb{R}.$$

Puis on superpose le graphe de la fonction de répartition associée à la loi \mathcal{L} en estimant éventuellement les paramètres inconnus de celle-ci.

- **Méthode de l'approximation de la densité** : On utilise une estimation de la densité inconnue. Puis on superpose les valeurs de la "densité" associée à la loi \mathcal{L} en estimant éventuellement les paramètres inconnus de celle-ci.

- **Méthode du QQ plot (quantile-quantile plot)** : Cette méthode consiste en la comparaison des quantiles empiriques et des quantiles théoriques. Soit $F(x) = \mathbb{P}(X \leq x)$ la fonction de répartition de X et x_p le quantile d'ordre p définie par : $x_p = \inf\{x \in \mathbb{R}; F(x) \geq p\}$. Soient $x_{(1)}, x_{(2)}, \dots, x_{(n)}$ les données rangées par ordre croissant : $x_{(1)} < x_{(2)} < \dots < x_{(n)}$. Alors on peut écrire la fonction de répartition comme

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{x_i \leq x\}} = \begin{cases} 0 & \text{si } x < x_{(1)}, \\ \frac{k}{n} & \text{si } x_{(k)} \leq x < x_{(k+1)}, \quad k \in \{1, \dots, n-1\}, \\ 1 & \text{si } x \geq x_{(n)}. \end{cases}$$

Soit (p_1, p_2, \dots, p_n) une suite strictement croissante de n réels vérifiant $p_k \in](k-1)/n, k/n[$, $k \in \{1, \dots, n\}$, de sorte que $\inf\{x \in \mathbb{R}; F_n(x) \geq p_k\} = x_{(k)}$ pour tout $k \in \{1, \dots, n\}$.

On appelle QQ plot le nuage de points \mathcal{N} dans le repère orthonormé (O, I, J) défini par :

$$\mathcal{N} = \{(x_{p_1}, x_{(1)}), (x_{p_2}, x_{(2)}), \dots, (x_{p_n}, x_{(n)})\}.$$

Si X suit la loi \mathcal{L} , les données font que F_n est une bonne estimation de F et, a fortiori, $x_{(k)}$ doit bien estimer $x_{p_k} : x_{(k)} \simeq x_{p_k}$ pour tout $k \in \{1, \dots, n\}$; les points du nuage \mathcal{N} doivent être proche de la "droite diagonale" d'équation : $y = x$.

- o **Cas d'une loi normale : méthode du QQ plot (QQ norm) avec droite de Henry** : Soit z_p le quantile d'ordre p d'une *var* Z suivant la loi normale centrée réduite $\mathcal{N}(0, 1)$.

Alors, si X suit la loi normale $\mathcal{N}(\mu, \sigma^2)$, le quantile d'ordre p de X vérifie

$$x_p = \mu + \sigma z_p.$$

Par conséquent, au lieu du QQ plot standard, on peut se contenter de construire le nuage de points \mathcal{N}_* dans le repère orthonormé (O, I, J) défini par : $\mathcal{N}_* = \{(z_{p_1}, x_{(1)}), (z_{p_2}, x_{(2)}), \dots, (z_{p_n}, x_{(n)})\}$. Si X suit la loi $\mathcal{N}(\mu, \sigma^2)$, alors les points du nuage \mathcal{N}_* doivent être proche de la droite d'équation :

$$y = \bar{x} + sx,$$

avec

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}.$$

Cette droite est appelée droite de Henry.

- o **Cas d'une loi normale : méthode de la boîte à moustaches** : On fait la boîte à moustache associée aux données. Si X suit une loi normale, il n'y a approximativement que 0,7% des points qui se trouvent en dehors des "moustaches". D'autre part, la boîte doit être à peu près symétrique par rapport à la médiane, idem pour les moustaches.

Exemples

Exemple 1. On fait passer à 50 adolescents le test psychologique de Rorschach. Les temps de passation en minutes du test sont :

43	48	65	55	51	51	44	51	59	62
45	53	55	55	49	34	52	69	45	54
59	36	36	29	52	59	41	58	54	55
72	53	52	49	57	42	70	58	42	53
57	68	40	65	54	49	32	56	50	59

On s'interroge pour savoir si la $\text{var } X$ qui à un adolescent associe son temps de passation au test suit ou non une loi normale.

Solution 1. Par l'énoncé, on observe la valeur de X sur chacun des n individus (adolescents) d'un échantillon avec $n = 50 : (x_1, \dots, x_n)$ (avec $x_i \in \mathbb{R}$).

Dire que X suit une loi normale $\mathcal{N}(\mu, \sigma^2)$, avec μ et σ inconnus, signifie qu'elle possède la densité :

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \quad x \in \mathbb{R}.$$

Pour préciser la loi normale $\mathcal{N}(\mu, \sigma^2)$ la plus adaptée à notre contexte, il faut estimer μ et σ à l'aide des données. On estime alors μ par

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = 51.94$$

et σ par

$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2} = 9.704638.$$

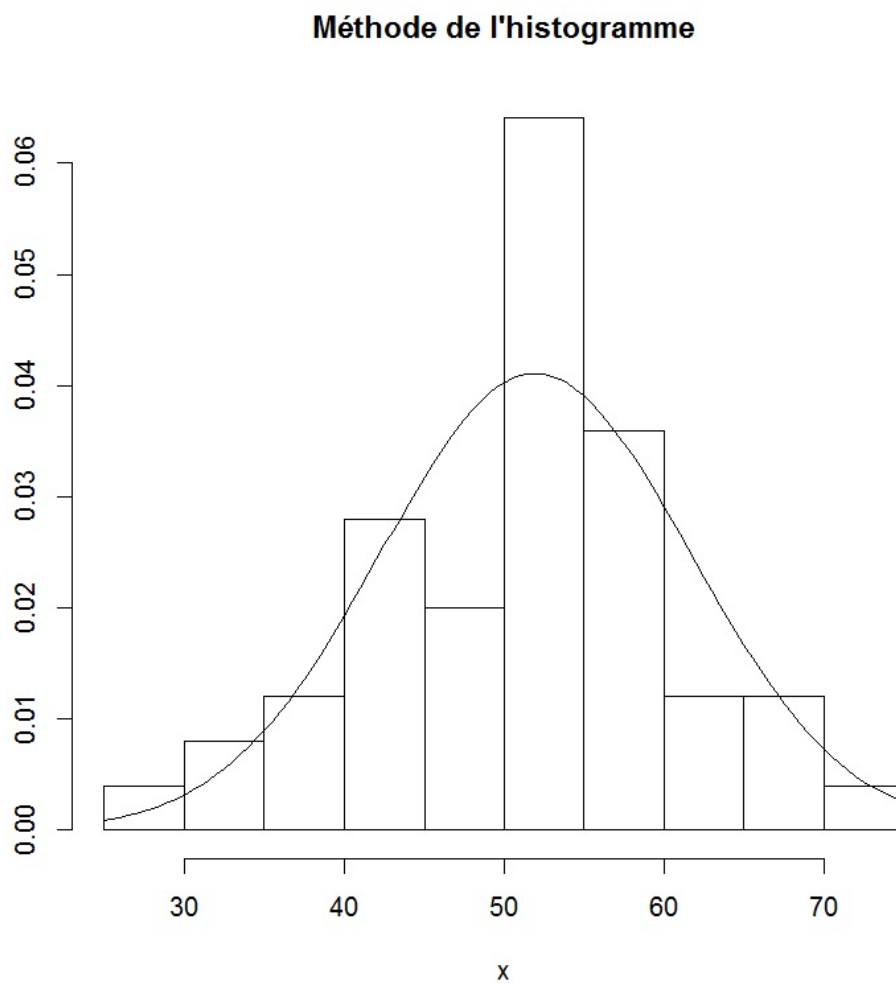
Ainsi, la problématique est la suivante : est-ce que ces données nous permettent d'affirmer que X ne suit pas la loi $\mathcal{L} = \mathcal{N}(51.94, 9.704638^2)$?

Nous allons faire une analyse graphique en utilisant les méthodes présentées précédemment.

o *Méthode de l'histogramme* : On propose les commandes :

```
x = c(43, 48, 65, 55, 51, 51, 44, 51, 59, 62, 45, 53, 55, 55, 49, 34, 52,  
69, 45, 54, 59, 36, 36, 29, 52, 59, 41, 58, 54, 55, 72, 53, 52, 49, 57, 42,  
70, 58, 42, 53, 57, 68, 40, 65, 54, 49, 32, 56, 50, 59)  
hist(x, freq = FALSE, main = "Méthode de l'histogramme", ylab = "")  
curve(dnorm(x, 51.94, 9.704638), add = TRUE)
```

Cela renvoie :

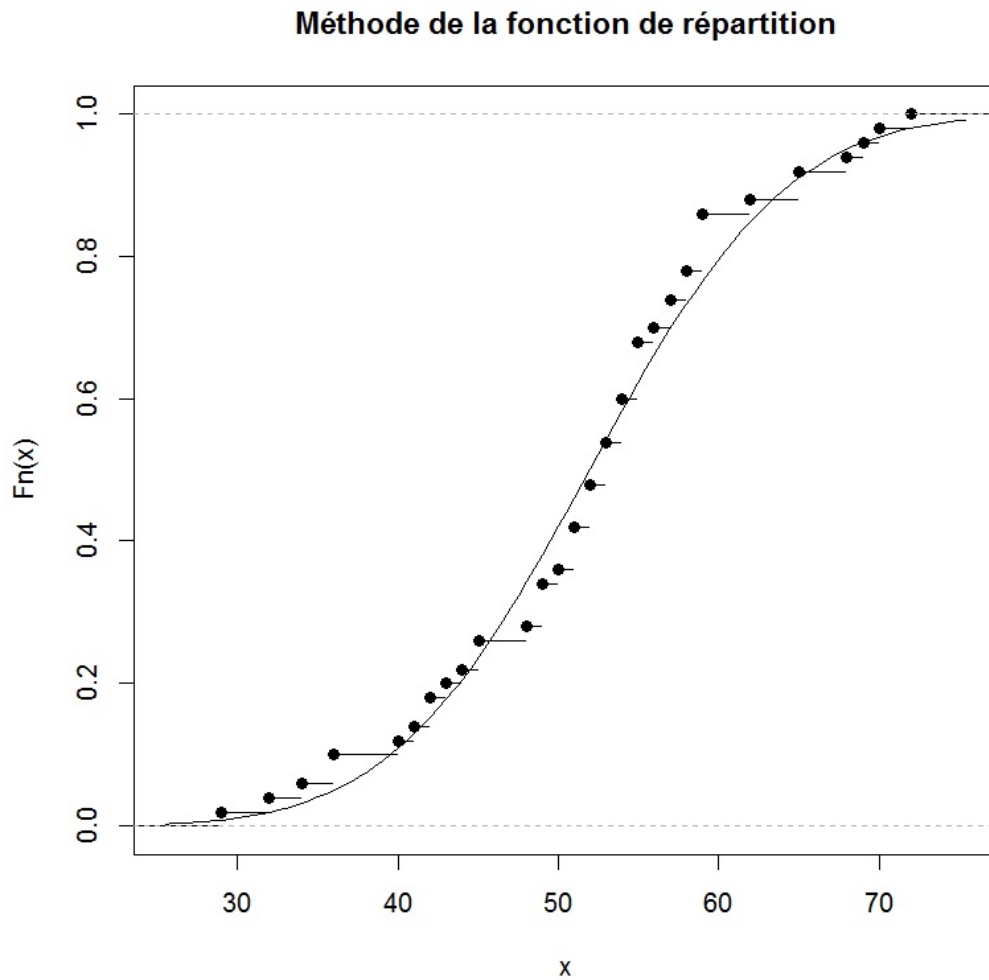


Vu les différences observées, il est difficile de conclure.

o *Méthode de la fonction de répartition* : On propose les commandes :

```
plot(ecdf(x), main = "Méthode de la fonction de répartition")
curve(pnorm(x, 51.94, 9.704638), add = TRUE)
```

Cela renvoie :

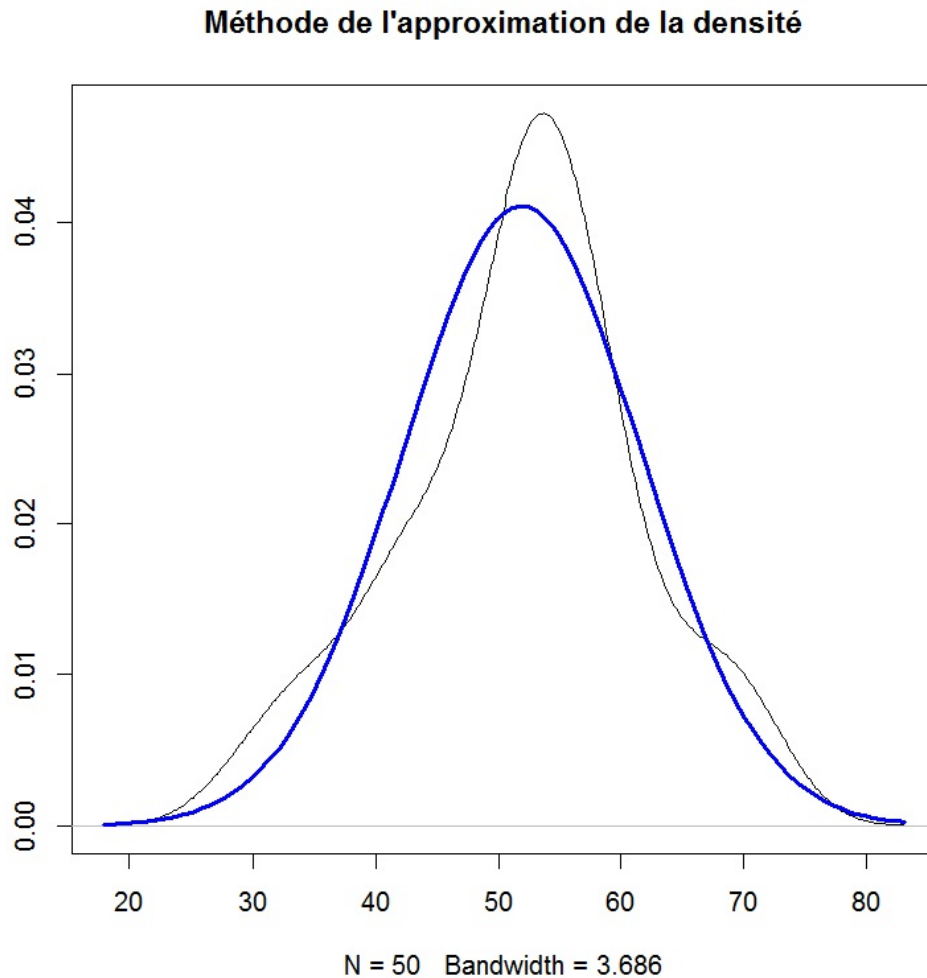


Ici, les différences observées laissent penser que X suit bien une loi normale.

o *Méthode de l'approximation de la densité* : On propose les commandes :

```
plot(density(x), main = "Méthode de l'approximation de la densité",  
ylab = "")  
curve(dnorm(x, 51.94, 9.704638), lwd = 3, col = "blue", add = TRUE)
```

Cela renvoie :



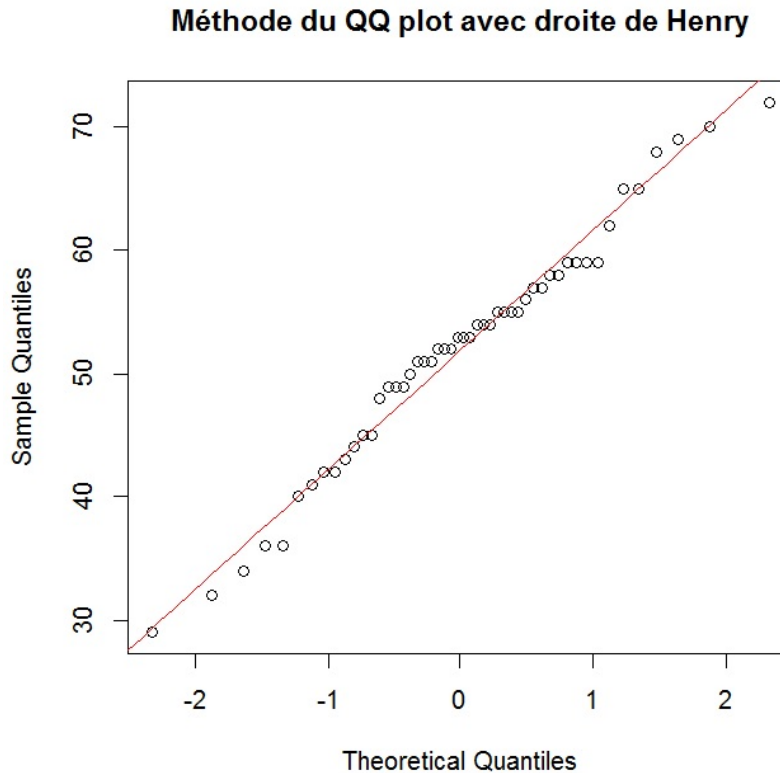
Les informations "N = 50 Bandwidth = 3.686" précisent des quantités utilisées dans l'estimateur de la densité : c'est un estimateur dit "à noyau", lequel utilise une fenêtre (Bandwidth) qui s'ajuste en fonction des données.

De nouveau, les différences observées laissent penser que X suit bien une loi normale.

o *Méthode du QQ plot (cas d'une loi normale)* : On propose les commandes :

```
qqnorm(x, main = "Méthode du QQ plot avec droite de Henry")
a = mean(x) ; b = sd(x)
curve(a + b * x, -6, 6, col = "red", add = TRUE)
```

Cela renvoie :



La droite de Henry ajuste bien le nuage de points ; on peut envisager que X suit une loi normale.

À la place de la commande `qqnorm`, on aurait pu faire le QQ plot à la main :

```
plot(qnorm(ppoints(x)), sort(x))
```

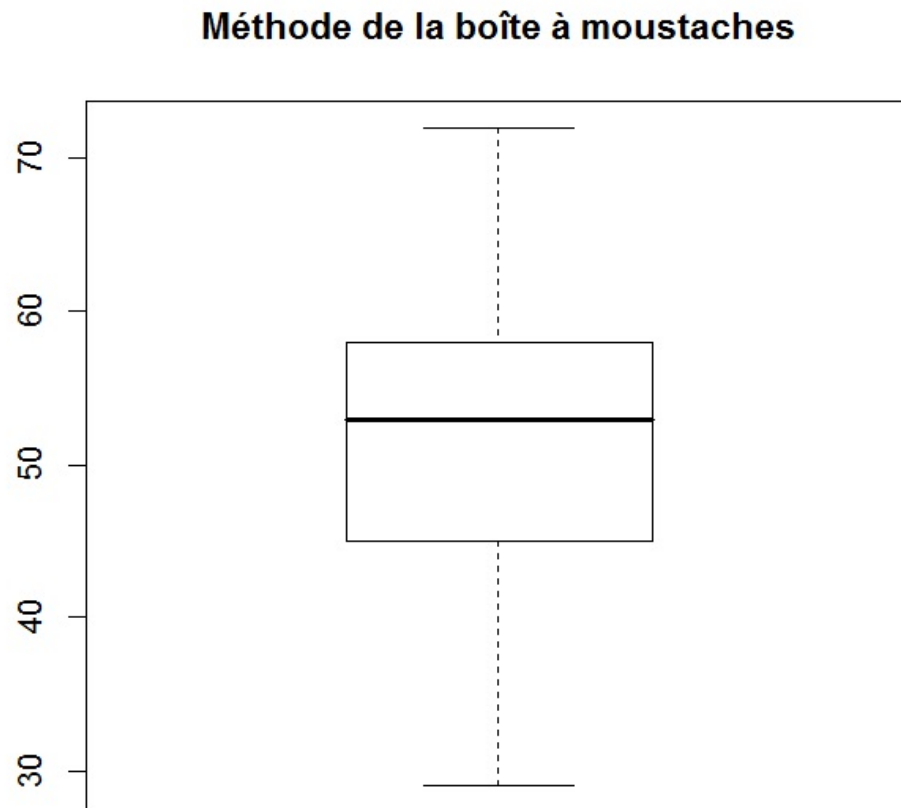
Une autre solution est de "centrer" et "réduire" les données : $x_i^* = (x_i - \bar{x})/s$, $i \in \{1, \dots, n\}$, considérer la loi normale centrée réduite $\mathcal{N}(0, 1)$, tracer le QQ plot associé et évaluer son ajustement par la droite $y = x$:

```
qqnorm(scale(x), main = "Méthode du QQ plot")
abline(0, 1, col = "red")
```

- *Cas d'une loi normale : méthode de la boîte à moustaches* : On fait :

```
boxplot(x, main = "Méthode de la boîte à moustaches")
```

Cela renvoie :



On constate qu'aucun point ne se trouve en dehors des "moustaches". D'autre part, la boîte est à peu près symétrique par rapport à la médiane, idem pour les moustaches. On peut envisager que X suit une loi normale.

Conclusion : Toutes les analyses graphiques laissent penser que X suit une loi normale.

Exemple 2. Dans une station service, un jour donné, les durée du passage en caisse de 49 clients on été mesurés. Les résultats en secondes sont :

0.96	1.45	0.42	3.69	2.58	1.95	1.74	0.01	1.02	1.12
0.17	3.19	0.85	1.27	0.68	3.60	1.23	0.34	0.31	0.16
0.07	0.79	0.02	1.20	0.05	2.09	0.24	5.46	2.57	0.89
0.74	1.67	0.88	2.27	0.22	3.39	0.12	0.06	0.78	0.32
5.79	2.09	0.39	1.82	2.96	0.20	0.08	0.37	2.58	0.30

Soit X la *var* égale à la durée de passage d'un client. On s'interroge sur le fait que X suit ou non une loi exponentielle.

Solution 1. Par l'énoncé, on observe la valeur de X sur chacun des n individus (clients) d'un échantillon avec $n = 50$: (x_1, \dots, x_n) (avec $x_i \in \mathbb{R}$).

Dire que X suit une loi exponentielle $\mathcal{E}(\lambda)$, avec $\lambda > 0$ inconnu, signifie qu'elle possède la densité :

$$f(x) = \begin{cases} \lambda e^{-\lambda x} & \text{si } x \geq 0, \\ 0 & \text{sinon.} \end{cases}$$

Pour préciser la loi exponentielle $\mathcal{E}(\lambda)$ la plus adaptée à notre contexte, il faut estimer λ à l'aide des données. Comme $\mathbb{E}(X) = 1/\lambda \Leftrightarrow \lambda = 1/\mathbb{E}(X)$, la méthode des moments nous assure qu'une estimation de λ est $1/\bar{x}$. On peut la calculer en faisant :

```
x = c(0.96, 1.45, 0.42, 3.69, 2.58, 1.95, 1.74, 0.01, 1.02, 1.12, 0.17,
3.19, 0.85, 1.27, 0.68, 3.60, 1.23, 0.34, 0.31, 0.16, 0.07, 0.79, 0.02,
1.20, 0.05, 2.09, 0.24, 5.46, 2.57, 0.89, 0.74, 1.67, 0.88, 2.27, 0.22,
3.39, 0.12, 0.06, 0.78, 0.32, 5.79, 2.09, 0.39, 1.82, 2.96, 0.20, 0.08,
0.37, 2.58, 0.30)
1 / mean(x)
```

Cela renvoie : [1] 0.7446016

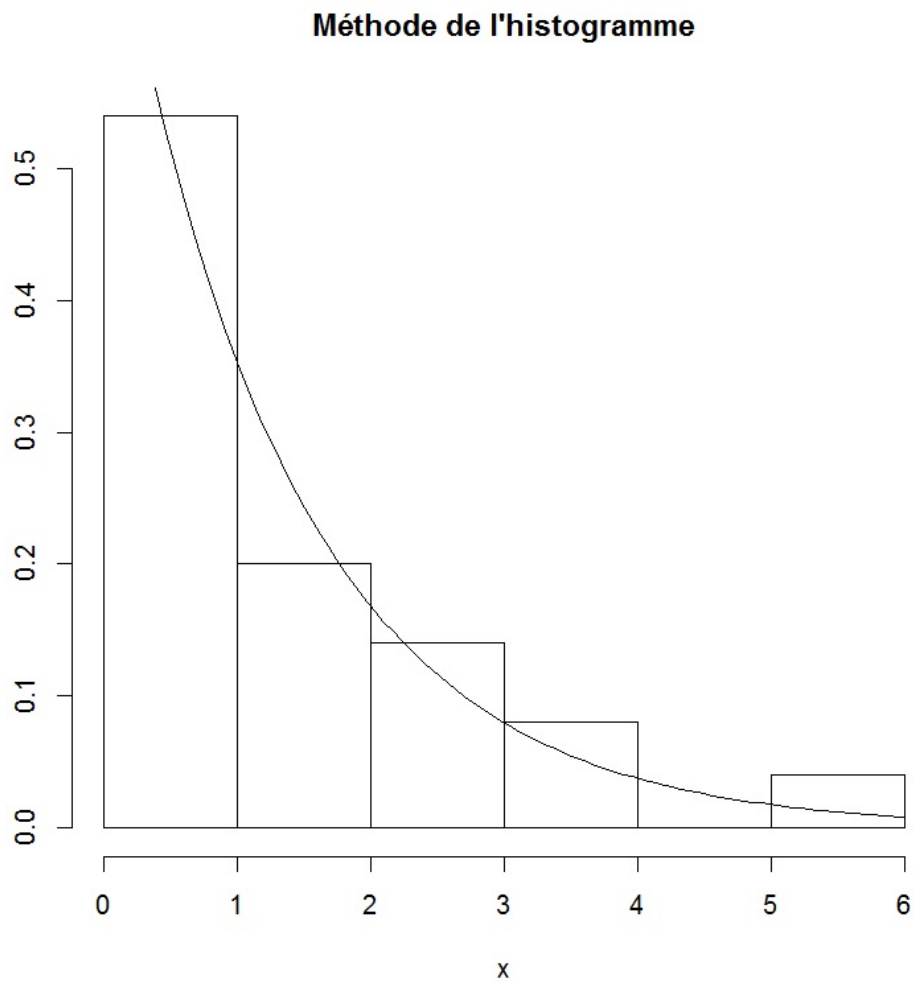
Ainsi, la problématique est la suivante : est-ce que ces données nous permettent d'affirmer que X ne suit pas la loi $\mathcal{L} = \mathcal{E}(0.7446016)$?

Nous allons faire une analyse graphique en utilisant les méthodes présentées précédemment.

- o *Méthode de l'histogramme* : On propose les commandes :

```
hist(x, freq = FALSE, main = "Méthode de l'histogramme", ylab = "")  
curve(dexp(x, 0.7446016), add = TRUE)
```

Cela renvoie :

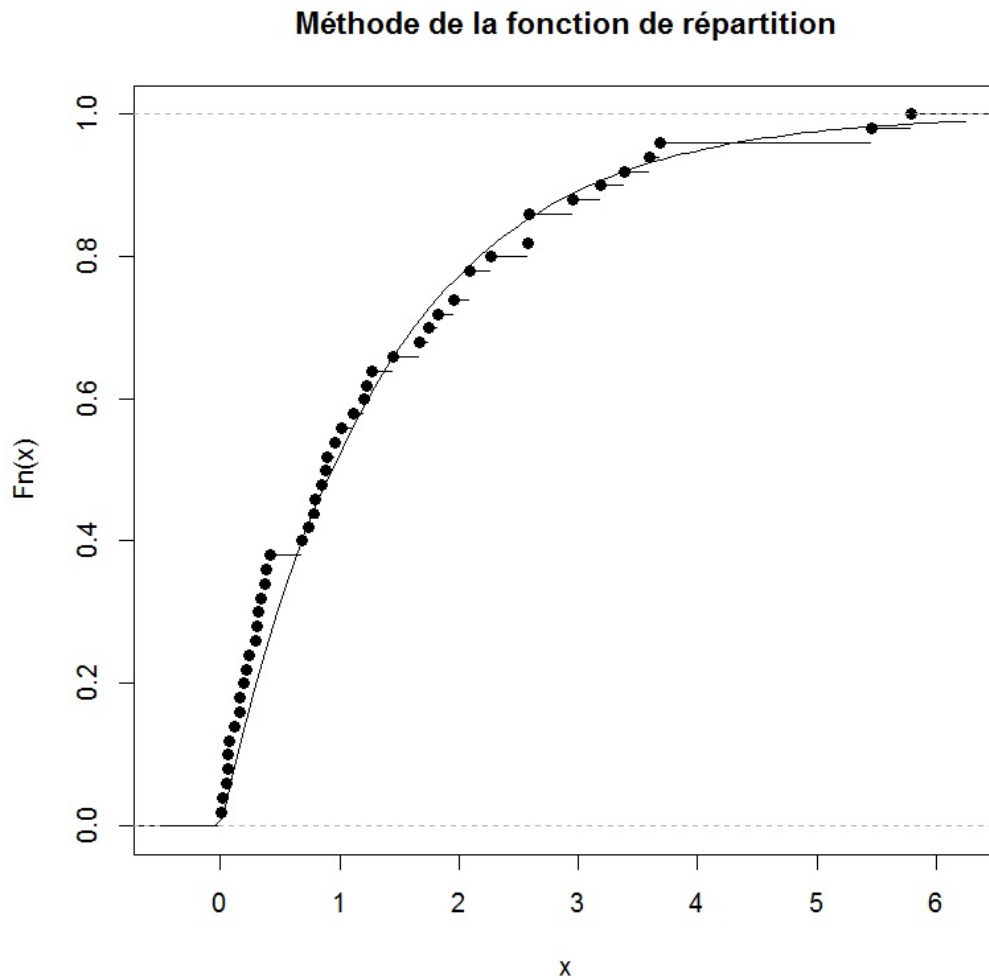


Vu les différences observées, il est vraisemblable que X suit une loi exponentielle.

o *Méthode de la fonction de répartition* : On propose les commandes :

```
plot(ecdf(x), main = "Méthode de la fonction de répartition")
curve(pexp(x, 0.7446016), add = TRUE)
```

Cela renvoie :

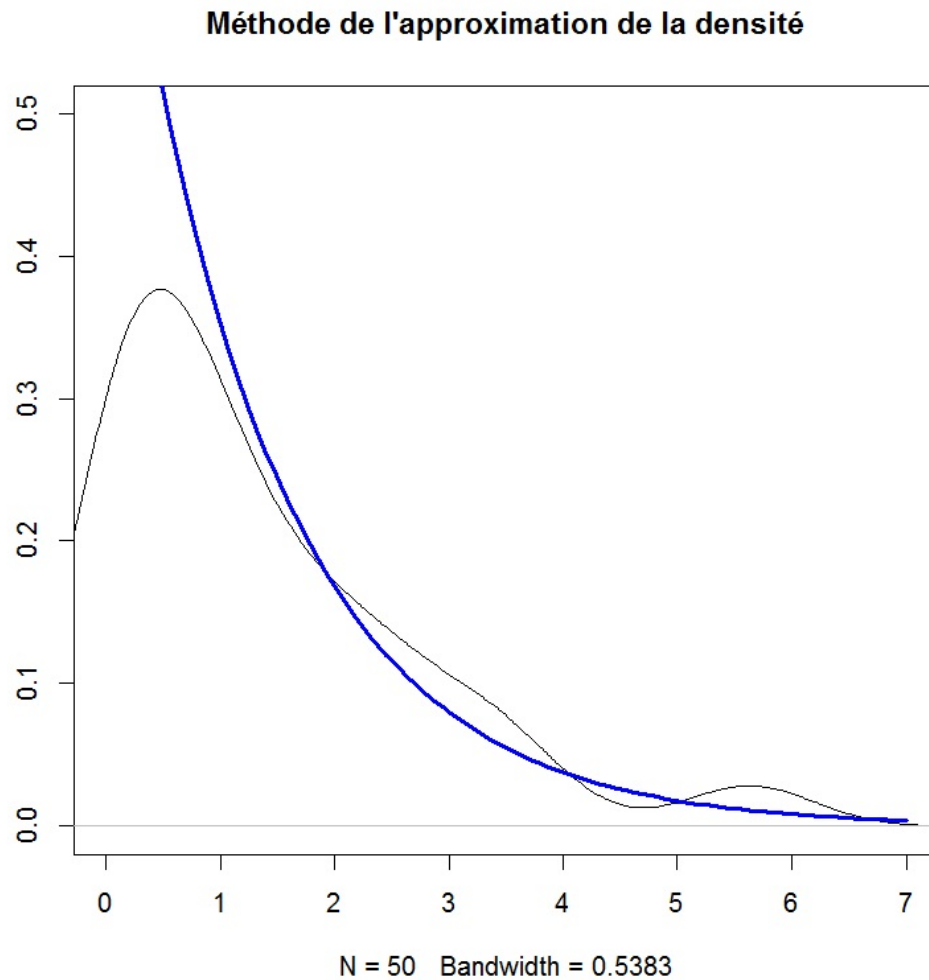


Les différences observées laissent penser que X suit bien une loi exponentielle.

o *Méthode de l'approximation de la densité* : On propose les commandes :

```
plot(density(x), main = "Méthode de l'approximation de la densité",  
xlim = c(0, 7), ylim = c(0, 0.5), ylab = "")  
curve(dexp(x, 0.7446016), lwd = 3, col = "blue", add = TRUE)
```

Cela renvoie :

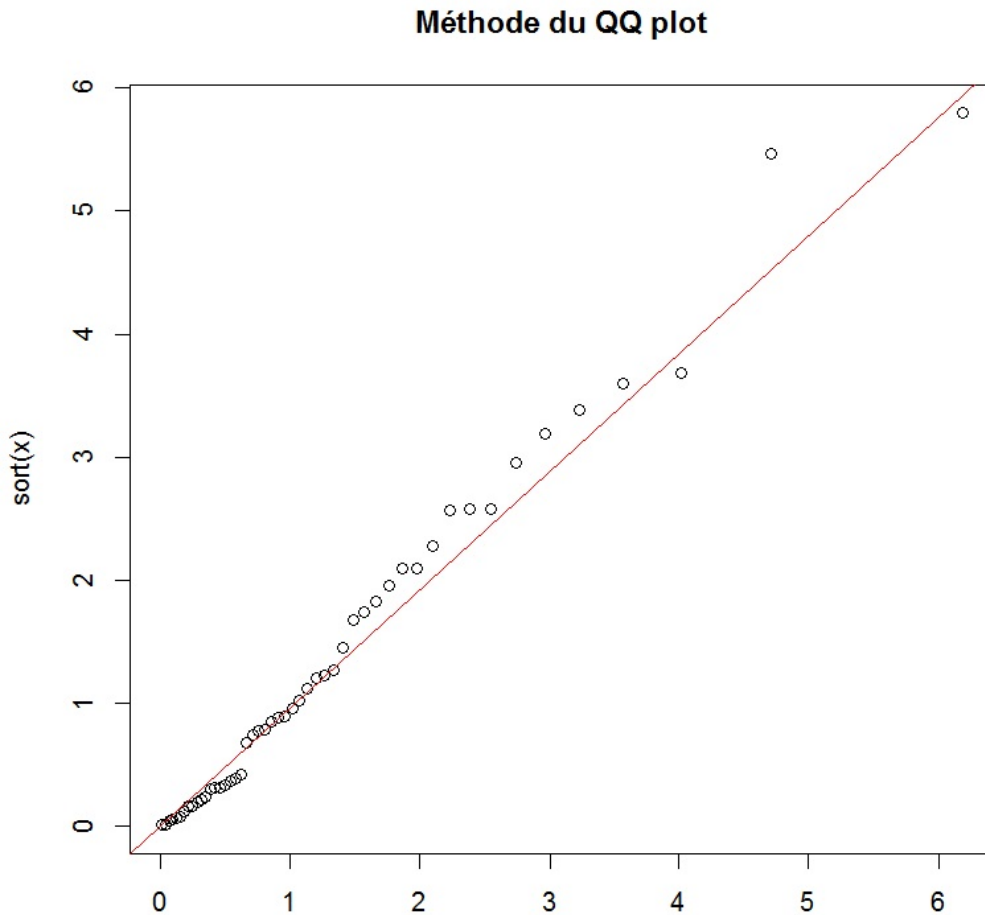


De nouveau, les différences observées laissent penser que X suit bien la loi exponentielle.

o *Méthode du QQ plot* : On propose les commandes :

```
plot(qexp(ppoints(x), 0.7446016), sort(x), main = "Méthode du QQ plot",  
xlab = "")  
abline(0,x, col = "red")
```

Cela renvoie :



La droite diagonale $y = x$ ajuste bien le nuage de points ; on peut envisager que X suit une loi exponentielle.

Conclusion : Toutes les analyses graphiques laissent penser que X suit une loi exponentielle.

3 Tests statistiques d'adéquation à une loi

3.1 Test du Chi-deux

Données

On observe la valeur d'une *var* X sur chacun des n individus d'un échantillon avec $n \geq 50$. Ces valeurs constituent les données : x_1, \dots, x_n . Elles sont généralement présentées sous la forme d'un tableau classes-effectifs :

Classes	Effectifs
C_1	n_1
\vdots	\vdots
C_k	n_k

Dans ce tableau, pour tout $i \in \{1, \dots, k\}$, C_i peut être une des valeurs x_1, \dots, x_n , ou un intervalle de valeurs et n_i est le nombre d'individus dont l'observation de X appartient la classe C_i .

Mise en oeuvre

Soit \mathcal{L} une loi de probabilité coïncidant possiblement / étant possiblement en adéquation avec la loi inconnue de X . On considère les hypothèses :

$$H_0 : "X \text{ suit la loi } \mathcal{L}" \quad \text{contre} \quad H_1 : "X \text{ ne suit pas la loi } \mathcal{L}"$$

Pour pouvoir décider du rejet de H_0 :

- Si besoin est, on estime ponctuellement les ℓ paramètres inconnus de la loi \mathcal{L} à l'aide des données, et on considère une *var* R suivant cette loi définie avec les paramètres estimés.
- Si besoin est, on ajuste la première et la dernière classes : C_1 et C_k , de sorte que $\bigcup_{i=1}^k C_i = R(\Omega)$,
- Pour tout $i \in \{1, \dots, k\}$, on calcule la probabilité :

$$p_i = \mathbb{P}(R \in C_i).$$

On vérifie que, pour tout $i \in \{1, \dots, k\}$, $np_i \geq 5$. Si tel n'est pas le cas, on crée une ou plusieurs nouvelles classes par fusion des anciennes, redéfinissant ainsi le n_i (et le k), jusqu'à obtenir de nouvelles probabilités p_i vérifiant cette hypothèse.

- On calcule

$$\chi_{obs}^2 = \sum_{i=1}^k \frac{(n_i - np_i)^2}{np_i} = \sum_{i=1}^k \frac{n_i^2}{np_i} - n.$$

- Soit $K \sim \chi^2(\nu)$, $\nu = k - 1 - \ell$. Alors la p-valeur associée au test du Chi-deux d'adéquation à une loi est

$$\text{p-valeur} = \mathbb{P}(K \geq \chi_{obs}^2).$$

Commandes

En pratique, le test du chi-deux d'adéquation à une loi est surtout utilisé pour une *var X* discrète ; il existe d'autres tests statistiques "plus puissants" dans le cas où *var X* est à densité.

Lorsqu'aucun paramètre de la loi n'est à estimer, on propose les commandes :

```
chisq.test(nb, p = proba)$p.value
```

Lorsqu'un ou plusieurs paramètres d'une loi discrète sont à estimer, en introduisant le degré de liberté ajusté *deg*, les commandes donnent la bonne p-valeur :

```
x2obs = chisq.test(nb, p = proba)$statistic
deg = nbdeclasses - 1 - nbdeparametrestimes
1 - pchisq(x2obs, deg)
```

Si un "Warning message" du type "Chi-squared approximation may be incorrect" apparaît, il vérifier que $np_i \geq 5$ pour tout $i \in \{1, \dots, k\}$ à la main et modifier les classes en fonction.

En disposant les données brutes dans un vecteur *x* , on peut essayer :

```
library(vcd)
gf = goodfit(x, type = "poisson", method = "MinChisq")
summary(gf)
plot(gf)
```

Exemples

Exemple 1. La marque Smies produit des bonbons au chocolat de six couleurs. Le responsable de la communication affirme que la proportion de chaque couleur est très précisément de 30% pour le brun (B), 20% pour le jaune (J), 20% pour le rouge (R), 10% pour l'orange (O), 10% pour le vert (V) et 10% pour le doré (D) dans tout échantillon de grande taille. Une expérience réalisée sur un échantillon de 370 bonbons donne les comptages suivants :

Couleur	B	J	R	O	V	D
Nombre de bonbons	84	79	75	49	36	47

Peut-on affirmer, au risque 5%, que le responsable de la communication à tort ?

Solution 1. Soit X la *var* égale à la couleur d'un bonbon choisi au hasard. Par l'énoncé, on observe la valeur de X sur chacun des n individus (bonbons) d'un échantillon avec $n = 370 : (x_1, \dots, x_n)$. Ces valeurs sont regroupées en $k = 6$ classes : $C_1 = \{B\}$, $C_2 = \{J\}$, $C_3 = \{R\}$, $C_4 = \{O\}$, $C_5 = \{V\}$ et $C_6 = \{D\}$, avec pour effectifs respectifs : $n_1 = 84$, $n_2 = 79$, $n_3 = 75$, $n_4 = 49$, $n_5 = 36$ et $n_6 = 47$. On considère les hypothèses :

H_0 : " X suit la loi de probabilité décrite par le responsable de la communication" contre

H_1 : " X ne suit pas la loi de probabilité décrite par le responsable de la communication".

Soit R une *var* suivant la loi de probabilité décrite par le responsable de la communication. On

a $p_1 = \mathbb{P}(R \in C_1) = 0.3$, $p_2 = \mathbb{P}(R \in C_2) = 0.2$, $p_3 = \mathbb{P}(R \in C_3) = 0.2$, $p_4 = \mathbb{P}(R \in C_4) = 0.1$, $p_5 = \mathbb{P}(R \in C_5) = 0.1$ et $p_6 = \mathbb{P}(R \in C_6) = 0.1$.

On considère les commandes :

```
nb = c(84, 79, 75, 49, 36, 47)
proba = c(0.3, 0.2, 0.2, 0.1, 0.1, 0.1)
chisq.test(nb, p = proba)$p.value
```

Cela renvoie : [1] 0.01880704

Notons qu'aucun "Warning message" n'apparaît ; les conditions d'application du test sont vérifiées.

Comme p-valeur $\in]0.01, 0.05]$, le rejet de H_0 est significatif \star .

Par conséquent, au risque 5%, on peut affirmer que X ne suit pas la loi décrite par le responsable de la communication.

Exemple 2 (retour sur un exemple déjà présenté). On souhaite savoir si les entrées à l'hôpital pour une certaine maladie sont réparties au hasard dans l'année ou bien si certains mois sont plus propices à la maladie. On examine le mois d'entrée d'un échantillon de 120 porteurs de la maladie étudiée. Les résultats sont :

Mois d'entrée	1	2	3	4	5	6	7	8	9	10	11	12
Nombre d'entrées	18	16	8	10	6	4	4	9	11	10	12	12

Peut-on affirmer, au risque 1%, que "les entrées ne se font pas au hasard dans l'année" (donc que "certains mois sont plus propices à la maladie") ?

Solution 2. Soit X la *var* égale au le mois d'entrée à l'hôpital. Par l'énoncé, on observe la valeur de X sur chacun des n individus (porteurs de la maladie) d'un échantillon avec $n = 120 : (x_1, \dots, x_n)$. Ces valeurs sont regroupées en $k = 12$ classes : $C_1 = \{1\}, C_2 = \{2\}, \dots, C_{12} = \{12\}$, avec pour effectifs respectifs : $n_1 = 18, n_2 = 16, \dots, n_{12} = 12$.

Dire que les entrées se font au hasard dans l'année signifie que X suit la loi uniforme $\mathcal{U}(\{1, \dots, 12\})$:

$$\mathbb{P}(X = i) = \frac{1}{12}, \quad i \in \{1, \dots, 12\}.$$

La problématique est la suivante : est-ce que ces données nous permettent d'affirmer, au risque 1%, que X ne suit pas la loi $\mathcal{L} = \mathcal{U}(\{1, \dots, 12\})$?

On considère alors les hypothèses :

H_0 : " X suit la loi uniforme $\mathcal{U}(\{1, \dots, 12\})$ " contre

H_1 : " X ne suit pas la loi uniforme $\mathcal{U}(\{1, \dots, 12\})$ "

Soit R une *var* suivant la loi uniforme $\mathcal{U}(\{1, \dots, 12\})$. On a $p_i = \mathbb{P}(R \in C_i) = \frac{1}{12}$ pour tout $i \in \{1, \dots, 12\}$.

On considère les commandes :

```
nb = c(18, 16, 8, 10, 6, 4, 4, 9, 11, 10, 12, 12)
proba = rep(1 / 12, 12)
chisq.test(nb, p = proba)$p.value
```

Cela renvoie : [1] 0.04267211

Notons qu'aucun "Warning message" n'apparaît ; les conditions d'application du test sont vérifiées.

Comme p-valeur $\in]0.01, 0.05]$, le rejet de H_0 est significatif \star .

Toutefois, comme p-valeur > 0.01 , au risque 1%, les données ne nous permettent pas d'affirmer que X ne suit pas la loi uniforme $\mathcal{U}(\{1, \dots, 12\})$; on ne peut donc rien conclure sur le fait que les entrées ne se font pas au hasard dans l'année.

Exemple 3 (retour sur un exemple déjà présenté). Dans un verger, on étudie le comportement des insectes quand ceux-ci attaquent les fruits. Soit X le caractère qui dénombre le nombre d'attaques d'insectes sur un fruit pris au hasard. Une étude statistique antérieure montre que, si les attaques se font de façon indépendantes les unes des autres, on peut modéliser X comme une *var* suivant une loi de Poisson $\mathcal{P}(\lambda)$ avec λ inconnu.

On considère un échantillon de 300 fruits et on compte le nombre d'attaques sur le fruit. Les résultats sont :

Nombre d'attaques	0	1	2	3	4	5	6	7
Nombre de fruits attaqués	60	105	65	47	15	4	3	1

Peut-on dire que le comportement des insectes est significativement grégaire ? (*on dit que le comportement des insectes est grégaire quand chacun d'entre eux a tendance à se comporter comme le voisin et à attaquer le même fruit ; dans ce cas, X ne suit pas une loi de Poisson*).

Solution 3. Par l'énoncé, on observe la valeur de X sur chacun des n individus (fruits) d'un échantillon avec $n = 300$: (x_1, \dots, x_n) (avec $x_i \in \mathbb{N}$). Ces valeurs sont regroupées en $k = 8$ classes : $C_1 = \{0\}$, $C_2 = \{1\}, \dots, C_8 = \{7\}$, avec pour effectifs respectifs : $n_1 = 60, n_2 = 105, \dots, n_8 = 1$.

La problématique est la suivante : est-ce que ces données nous permettent d'affirmer, a moins au risque 5%, que X ne suit pas une loi de Poisson ?

On considère les hypothèses :

$$H_0 : "X \text{ suit une loi de Poisson}" \quad \text{contre} \quad H_1 : "X \text{ ne suit pas une loi de Poisson}"$$

Soit R une *var* suivant la loi de Poisson $\mathcal{P}(\lambda)$ avec $\lambda > 0$ inconnu. On a $R(\Omega) = \mathbb{N}$.

Comme $\bigcup_{i=1}^8 C_i = \{0, \dots, 7\} \neq \mathbb{N}$, on ajuste la dernière classe comme : $C_8 = \{7, 8, \dots, \infty\}$ (de sorte à ce que $\bigcup_{i=1}^8 C_i = R(\Omega)$). Le paramètre λ étant inconnu, il faut l'estimer à l'aide des données. Comme $\mathbb{E}(X) = \lambda$, la méthode des moments nous assure qu'une estimation de λ est la moyenne \bar{x} que l'on peut calculer en faisant :

```
nb = c(60, 105, 65, 47, 15, 4, 3, 1)
lambda = sum((0:7) * nb) / 300
lambda
```

Cela renvoie : [1] 1.603333

Cette estimation sera prise en compte dans la suite. Pour tout $i \in \{1, \dots, 7\}$, on a

$$p_i = \mathbb{P}(R \in C_i) = \mathbb{P}(R = i - 1) = e^{-1.603333} \frac{1.603333^{i-1}}{(i-1)!}$$

et

$$p_8 = \mathbb{P}(R \in C_8) = \mathbb{P}(R \geq 7) = \sum_{i=7}^{\infty} e^{-1.603333} \frac{1.603333^i}{i!}.$$

On peut obtenir ces probabilités avec les commandes :

```
proba = c(dpois(0:6, lambda), 1 - ppois(6, lambda))
proba
```

Cela renvoie :

```
[1] 0.201224650 0.322630189 0.258641868 0.138229709 0.055407075 0.017767202
0.004747791 0.001351515
```

Dans un premier temps, on propose de mettre en œuvre le test du Chi-deux en faisant :

```
chisq.test(nb, p = proba)$p.value
```

Cela renvoie :

```
[1] 0.4732928
```

Warning message:

```
In chisq.test(nb, p = proba) : Chi-squared approximation may be incorrect
```

On dénombre alors deux problèmes :

- le logiciel n'a pas pris en compte le fait que λ a été estimé (il ne peut pas le savoir),
- il y a un "Warning message" nous avertissant que l'hypothèse : $np_i \geq 5$ pour tout $i \in \{1, \dots, 8\}$ n'est peut-être pas vérifiée.

Étudions ce dernier point :

```
300 * proba
```

Cela renvoie :

```
[1] 60.3673950 96.7890567 77.5925604 41.4689128 16.6221226 5.3301606  
1.4243374 0.4054545
```

Comme $np_7 < 5$ et $np_8 < 5$, nous allons fusionner les classes C_6 , C_7 et C_8 , formant ainsi une nouvelle dernière classe : $C_6 = \{5, 6, \dots\}$ avec pour effectif : $n_6 = 4 + 3 + 1 = 8$. Il y a désormais $k = 6$ classes. On vérifie alors que l'hypothèse est vérifiée avec cette nouvelle configuration :

```
nb2 = c(60, 105, 65, 47, 15, 8)  
proba2 = c(dpois(0:4, lambda), 1 - ppois(4, lambda))  
300 * proba2
```

Cela renvoie : [1] 60.367395 96.789057 77.592560 41.468913 16.622123 7.159953

Aucune valeur ne dépasse 5. Pour avoir la p-valeur associée au test du Chi-deux en prenant en compte le fait que l'on a estimé un paramètre, on considère le degré de liberté :

$\nu = k - 1 - \ell = 6 - 1 - 1 = 4$. On fait :

```
x2obs = chisq.test(nb2, p = proba2)$statistic  
deg = 4  
1 - pchisq(x2obs, deg)
```

Cela renvoie :

X-squared

0.4427608

On obtient alors la vraie p-valeur associée au test du Chi-deux (il ne faut pas faire attention au "X-squared" et notons aussi qu'aucun "Warning message" n'apparaît).

Comme p-valeur > 0.05 , les données ne nous permettent pas de rejeter H_0 . Ainsi, on ne peut pas rejeter l'hypothèse selon laquelle les insectes n'ont pas un comportement grégaire.

3.2 Test de Kolmogorov-Smirnov

Contexte. On observe la valeur d'une *var* X sur chacun des n individus d'un échantillon. Ces valeurs constituent les données : x_1, \dots, x_n . On considère les hypothèses :

$$H_0 : "X \text{ suit la loi } \mathcal{L}" \quad \text{contre} \quad H_1 : "X \text{ ne suit pas la loi } \mathcal{L}"$$

Soit $F_n(x)$ désigne la fonction de répartition empirique associée aux données et $F(x)$ la fonction de répartition associé à la loi \mathcal{L} . L'idée du test de Kolmogorov-Smirnov est que plus $F_n(x)$ diffère de $F(x)$, plus le rejet de H_0 est significatif (idée similaire à l'analyse graphique de la méthode de la fonction de répartition). La p-valeur du test de Kolmogorov-Smirnov utilise la statistique de test observée :

$$d_{obs} = \sup_{x \in \mathbb{R}} |F_n(x) - F(x)|$$

et une loi de probabilité non-usuelle que l'on arrive à évaluer.

En pratique, le test de Kolmogorov-Smirnov est surtout utilisé pour une *var* X à densité. Il est plus puissant que le test du Chi-deux quand n est petit.

Les commandes associées sont :

```
ks.test(x, "pnorm")$p.value
```

Lorsqu'un ou plusieurs paramètres de la loi sont à préciser, on fait, par exemple :

```
ks.test(x, "pexp", lambda)$p.value
```


Exemple. On mesure les durées de vie de 20 ampoules d'un même type. Les résultats, en heures, sont :

673	389	1832	570	522	2694	3683	644	1531	2916
1069	3145	2268	3574	791	1418	649	3344	1153	3922

Soit X la *var* égale à la durée de vie en heures d'une ampoule de ce type.

Est-ce que l'on peut affirmer, au risque 5%, que X ne suit pas la loi exponentielle $\mathcal{E}(1/1850)$?

Solution. Par l'énoncé, on observe la valeur de X sur chacun des n individus (ampoules) d'un échantillon avec $n = 20$: (x_1, \dots, x_n) (avec $x_i \in \mathbb{R}$).

On considère les hypothèses :

H_0 : " X suit la loi exponentielle $\mathcal{E}(1/1850)$ " contre

H_1 : " X ne suit pas la loi exponentielle $\mathcal{E}(1/1850)$ "

Nous allons utiliser le test de Kolmogorov-Smirnov. On fait :

```
x = c(673, 389, 1832, 570, 522, 2694, 3683, 644, 1531, 2916, 1069, 3145,
2268, 3574, 791, 1418, 649, 3344, 1153, 3922)
ks.test(x, "pexp", 1 / 1850)$p.value
```

Cela renvoie : [1] 0.3774748

Comme p -valeur > 0.05 , les données ne nous permettent pas de rejeter H_0 . Ainsi, on ne peut pas rejeter l'hypothèse selon laquelle X suit la loi exponentielle $\mathcal{E}(1/1850)$.

3.3 Test de Shapiro-Wilk

Contexte. On observe la valeur d'une *var* X sur chacun des n individus d'un échantillon. Ces valeurs constituent les données : x_1, \dots, x_n . On cherche à montrer que X ne suit pas une loi normale, ce qui mettrait en défaut une hypothèse cruciale pour de nombreux outils statistiques comme des intervalles de confiance (T-IntConf, ...) et des tests statistiques (T-Test, test du coefficient de corrélation, ANOVA...).

On considère donc les hypothèses :

H_0 : " X suit une loi normale \mathcal{N} " contre H_1 : " X ne suit pas une loi normale \mathcal{N} "

On peut alors utiliser le test du Chi-deux ou le test de Kolmogorov-Smirnov avec $\mathcal{L} =$ loi normale \mathcal{N} . Toutefois, dans ce cas particulier, il est fortement conseillé d'utiliser le test de Shapiro-Wilk plus puissant.

Celui-ci utilise la statistique de test observée :

$$w_{obs} = \frac{\left(\sum_{i=1}^n a_i x_{(i)} \right)^2}{(n-1)s^2},$$

où les valeurs a_1, \dots, a_n sont calculés à partir du vecteur des moyennes et de la matrice de covariance des statistiques d'ordre de n *var iid* suivant une loi normale et une loi de probabilité non-usuelle que l'on arrive à évaluer.

Les commandes associées sont :

```
shapiro.test(x)$p.value
```

Remarque. Le test de Shapiro-Wilk appartient à la grande famille des "tests de normalité". Il en existe un grands nombres, parmi lesquels : le test de Lilliefors, le test de Anderson-Darling, le test de D'Agostino et le test de Jarque-Bera.

Toutefois, s'il fallait en retenir qu'un, ca serait le test de Shapiro-Wilk.

Exemple (retour sur un exemple déjà présenté). On fait passer à 50 adolescents le test psychologique de Rorschach. Les temps de passation en minutes du test sont :

43	48	65	55	51	51	44	51	59	62
45	53	55	55	49	34	52	69	45	54
59	36	36	29	52	59	41	58	54	55
72	53	52	49	57	42	70	58	42	53
57	68	40	65	54	49	32	56	50	59

Soit X la *var* qui à un adolescent associe son temps de passation au test. Peut-on affirmer, au risque 5%, que X ne suit pas une loi normale ?

Solution. Par l'énoncé, on observe la valeur de X sur chacun des n individus (adolescents) d'un échantillon avec $n = 50 : (x_1, \dots, x_n)$ (avec $x_i \in \mathbb{R}$). On considère les hypothèses :

$$H_0 : "X \text{ suit une loi normale}" \quad \text{contre} \quad H_1 : "X \text{ ne suit pas une loi normale}"$$

On utilise le test de Shapiro-Wilk :

```
x = c(43, 48, 65, 55, 51, 51, 44, 51, 59, 62, 45, 53, 55, 55, 49, 34, 52,  
69, 45, 54, 59, 36, 36, 29, 52, 59, 41, 58, 54, 55, 72, 53, 52, 49, 57, 42,  
70, 58, 42, 53, 57, 68, 40, 65, 54, 49, 32, 56, 50, 59)  
shapiro.test(x)$p.value
```

Cela renvoie : [1] 0.4801461

Comme p-valeur > 0.05 , on ne rejette pas H_0 ; l'hypothèse selon laquelle X suit une loi normale n'est pas rejetée.

4 Exercices

Exercice 1. On lance 100 fois un dé à 6 faces numérotées de 1 à 6. On obtient les résultats suivants :

Numéro	1	2	3	4	5	6
Nombre de fois	18	23	19	12	11	15

Peut-on affirmer que le dé est truqué ? (*on fera une analyse graphique convenable, puis un test statistique adapté au risque 5%*).

Exercice 2. Lindström, spécialiste de la génétique et de l'hybridation du maïs, a croisé deux types récessifs de maïs : le type vert-zébré et le type doré. Si les lois de la génétique sont respectées obtient :

- "vert" avec la probabilité $9/16$,
- "doré" avec la probabilité $3/16$,
- "vert-zébré" avec la probabilité $3/16$,
- "doré-vert-zébré" avec la probabilité $1/16$.

On effectue 1301 croisements. On obtient les résultats suivants :

Type	"vert"	"doré"	"vert-zébré"	"doré-vert-zébré"
Nombre de fois	773	231	238	59

Peut-on dire que les lois de la génétique ne sont pas respectées ? (*on fera une analyse graphique convenable, puis un test statistique adapté au risque 5%*).

Exercice 3. L'agence immobilière Centurion a étudié le nombre de biens vendus par agent le mois de Juin 2016. Les résultats obtenus sont :

Nombre de biens vendus	0	1	2	≥ 3
Nombre d'agents	15	18	11	8

Soit X la *var* égale au nombre de biens vendu par un agent en Juin 2016. Peut-on affirmer que X ne suit pas la loi de Poisson $\mathcal{P}(2)$? (*on fera une analyse graphique convenable, puis un test statistique adapté au risque 5%*).

Exercice 4. Dans le garage AutoLover, sur une période de 320 jours, on a relevé le nombre journalier d'accidents du travail. Les résultats obtenus sont :

Nombre d'accidents journalier	0	1	2	3	4	5	6	7
Nombre de jours	65	110	70	48	16	5	4	2

On peut modéliser le nombre journalier d'accidents du travail dans ce garage par une *var* X . Peut-on affirmer, au risque 5%, que X ne suit pas une loi de Poisson? (*on fera une analyse graphique convenable, puis un test statistique adapté au risque 5%*).

Exercice 5. On mesure les durées de vie en heures de 7 appareils. Les résultats obtenus sont :

Appareil	1	2	3	4	5	6	7
Durée de vie	145	110	170	48	116	95	74

On peut modéliser la durée de vie en heures d'un appareil par une *var* X . Peut-on affirmer que X , au risque 5%, que ne suit pas une loi exponentielle $\mathcal{E}(0.01)$?

Exercice 6. Dans un magasin, un jour donnée, on mesure les temps d'attente en minutes entre 2 clients à une caisse. Les résultats obtenus sont :

25.12	12.36	24.35	12.19	5.27	18.35	19.11	27.08	21.09	17.19	8.45	13.27	15.17
-------	-------	-------	-------	------	-------	-------	-------	-------	-------	------	-------	-------

On peut modéliser le temps d'attente entre 2 clients par une *var* X . Peut-on affirmer que X , au risque 5%, que ne suit pas une loi exponentielle $\mathcal{E}(1/\bar{x})$, \bar{x} où désigne la moyenne des valeurs obtenues?

Exercice 7. On étudie l'accroissement de poids en kilogrammes chez des pourceaux pendant une période de 15 jours. On prélève au hasard 50 pourceaux. Les résultats, sous la forme d'un tableau classes-effectifs, sont :

Classes	Effectifs
]0,4]	2
]4,8]	5
]8,12]	12
]12,16]	14
]16,20]	11
]20,24]	5
]24,28]	1

Soit X la *var* égale à l'accroissement de poids en kilogrammes d'un pourceau pendant 15 jours. À l'aide d'une analyse graphique, montrer qu'il est vraisemblable que X suit une loi normale.

Exercice 8. On pèse 20 plaquettes de beurre pris au hasard dans une production normande. Les résultats, en grammes, sont :

247.0	247.8	250.2	251.3	251.9	249.4	248.8	247.1	255.0	247.0
254.8	244.8	250.7	250.7	252.6	251.1	254.1	249.2	252.0	254.0

On suppose que le poids en grammes d'une plaquette de beurre de cette production peut être modélisé par une *var* X . Peut-on affirmer que X suit une loi normale ? (*on fera une analyse graphique convenable, puis un test statistique adapté au risque 5%*).

Exercice 9. Une ferme de Bay of Plenty en Nouvelle-Zélande produit des kiwis. On pèse 16 kiwis choisis au hasard dans cette ferme. Les résultats, en grammes, sont :

65.06	71.44	67.93	69.02	67.28	62.34	66.23	64.16
68.56	70.45	64.91	69.90	65.52	66.75	68.54	67.90

Soit X la *var* égale au poids en grammes d'un kiwi. Peut-on affirmer que X suit une loi normale? (on fera une analyse graphique convenable, puis un test statistique adapté au risque 5%).

Exercice 10. Soient X et Y deux *var* indépendantes telles que X suit la loi de Poisson $\mathcal{P}(5)$ et Y suit la loi de Poisson $\mathcal{P}(3)$. Alors on sait que $Z = X + Y$ suit la loi de Poisson $\mathcal{P}(5 + 3)$.

Illustrer ce résultat en simulant des *var* et en utilisant la commande `qqplot`.

Exercice 11. Soient X_1 et X_2 deux *var* indépendantes. Illustrer les résultats ci-dessous avec la commande `qqplot` :

$X_i \sim$	$\mathcal{E}(\lambda)$	$\Gamma(m_i, \lambda)$	$\mathcal{N}(\mu_i, \sigma_i^2)$	$\chi^2(\nu_i)$
$X_1 + X_2 \sim$	$\Gamma(2, \lambda)$	$\Gamma(m_1 + m_2, \lambda)$	$\mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$	$\chi^2(\nu_1 + \nu_2)$
Numérique	$\lambda = 3.8$	$m_1 = m_2 = 4.2, \lambda = 2.1$	$\mu_1 = \mu_2 = 1.6, \sigma_1 = \sigma_2 = 1.5$	$\nu_1 = \nu_2 = 3.2$

Exercice 12. Soient X et Y deux *var* indépendantes. Illustrer les résultats ci-dessous avec la commande `qqplot` :

- Caractérisation de la loi du chi-deux $\chi^2(2)$: Si $X \sim \mathcal{N}(0, 1)$ et $Y \sim \mathcal{N}(0, 1)$, alors

$$X^2 + Y^2 \sim \chi^2(2).$$

- Caractérisation de la loi de Student $\mathcal{T}(\nu)$: Si $X \sim \mathcal{N}(0, 1)$ et $Y \sim \chi^2(\nu)$, alors

$$\frac{X}{\sqrt{\frac{Y}{\nu}}} \sim \mathcal{T}(\nu).$$

Prendre $\nu = 3.9$.

- Caractérisation de la loi de Fisher $\mathcal{F}(\nu_1, \nu_2)$: Si $X \sim \chi^2(\nu_1)$ et $Y \sim \chi^2(\nu_2)$, alors

$$\frac{\frac{X}{\nu_1}}{\frac{Y}{\nu_2}} = \frac{\nu_2 X}{\nu_1 Y} \sim \mathcal{F}(\nu_1, \nu_2).$$

Prendre $(\nu_1, \nu_2) = (2.1, 8.3)$.

Exercice 13. On considère les commandes :

```
x = rnorm(100)
a = numeric()
for (i in 1:6) {
  a[i] = ks.test(x, "pnorm", 0, 1 + (i - 1) / 10)$p.value
}
a
```

Cela renvoie : [1] 0.32488689 0.17631961 0.09775220 0.05602855 0.03332187 0.02056516

Commenter ces résultats numériques.

5 Solutions

Solution 1. Soit X la *var* égale au numéro affiché par le dé après un lancer. Par l'énoncé, on observe la valeur de X sur chacun des n individus (dés) d'un échantillon avec $n = 100$: (x_1, \dots, x_n) (avec $x_i \in \{1, \dots, 6\}$). On forme alors un vecteur des effectifs $(n_1, n_2, \dots, n_6) = (18, 23, \dots, 15)$.

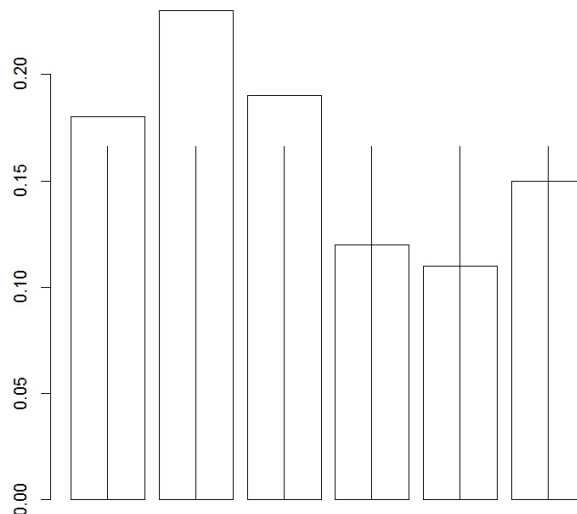
Dire que le dé n'est pas truqué signifie que X suit la loi uniforme $\mathcal{U}(\{1, \dots, 6\})$:

$$\mathbb{P}(X = i) = \frac{1}{6}, \quad i \in \{1, \dots, 6\}.$$

La problématique est la suivante : est-ce que ces données nous permettent d'affirmer que X ne suit pas la loi $\mathcal{U}(\{1, \dots, 6\})$? Ainsi, pour une analyse graphique, on propose les commandes :

```
nb = c(18, 23, 19, 12, 11, 15)
bar = barplot(nb / 100, col = "white")
points(bar, rep(1 / 6, 6), type = "h")
```

Cela renvoie :



Il est difficile de conclure au vu des différences observées. On a besoin d'un test statistique pour trancher. On considère alors les hypothèses :

H_0 : "X suit la loi uniforme $\mathcal{U}(\{1, \dots, 6\})$ " contre

H_1 : "X ne suit pas la loi uniforme $\mathcal{U}(\{1, \dots, 6\})$ "

Les valeurs sont regroupées en $k = 6$ classes : $C_1 = \{1\}$, $C_2 = \{2\}$, ..., $C_6 = \{6\}$. On considère les commandes :

```
proba = rep(1 / 6, 6)
chisq.test(nb, p = proba)$p.value
```

Cela renvoie : [1] 0.2757299

Notons qu'aucun "Warning message" n'apparaît ; les conditions d'application du test sont vérifiées.

Comme p-valeur > 0.05 , les données ne nous permettent pas de rejeter H_0 . Ainsi, on ne peut pas affirmer que le dé est truqué.

Solution 2.

Soit X la *var* égale au type de maïs obtenu avec un croisement. On utilise le codage :

X	"vert"	"doré"	"vert-zébré"	"doré-vert-zébré"
Codage	0	1	2	3

Par l'énoncé, on observe la valeur de X sur chacun des n individus (croisements) d'un échantillon avec $n = 1301$: (x_1, \dots, x_n) (avec $x_i \in \{0, 1, 2, 3\}$). On forme alors un vecteur des effectifs $(n_1, n_2, n_3, n_4) = (773, 231, 238, 59)$.

Dire que les lois de la génétique sont respectées signifie que X suit la loi décrite dans l'énoncé :

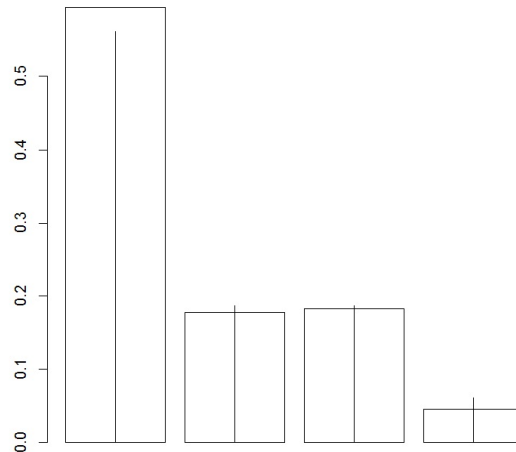
i	0	1	2	3
$\mathbb{P}(X = i)$	$\frac{9}{16}$	$\frac{3}{16}$	$\frac{3}{16}$	$\frac{1}{16}$

La problématique est la suivante : est-ce que ces données nous permettent d'affirmer que X ne suit pas la loi décrite dans l'énoncé ?

Ainsi, pour une analyse graphique, on propose les commandes :

```
nb = c(773, 231, 238, 59)
bar = barplot(nb / 1301, col = "white")
points(bar, c(9 / 16, 3 / 16, 3 / 16, 1 / 16), type = "h")
```

Cela renvoie :



Il est difficile de conclure au vu des différences observées. On a besoin d'un test statistique pour trancher. On considère alors les hypothèses :

H_0 : " X suit la loi décrite dans l'énoncé" contre

H_1 : " X ne suit pas la loi décrite dans l'énoncé"

Les valeurs sont regroupées en $k = 4$ classes : $C_1 = \{0\}$, $C_2 = \{1\}$, $C_3 = \{2\}$ et $C_4 = \{3\}$. On considère les commandes :

```
proba = c(9 / 16, 3 / 16, 3 / 16, 1 / 16)
chisq.test(nb, p = proba)$p.value
```

Cela renvoie : [1] 0.02589168

Notons qu'aucun "Warning message" n'apparaît ; les conditions d'application du test sont vérifiées.

Comme p-valeur $\in]0.01, 0.05]$, le rejet de H_0 est significatif \star . Par conséquent, on peut affirmer, au risque 5%, que X ne suit pas la loi décrite dans l'énoncé ; les lois de la génétique ne sont pas respectées.

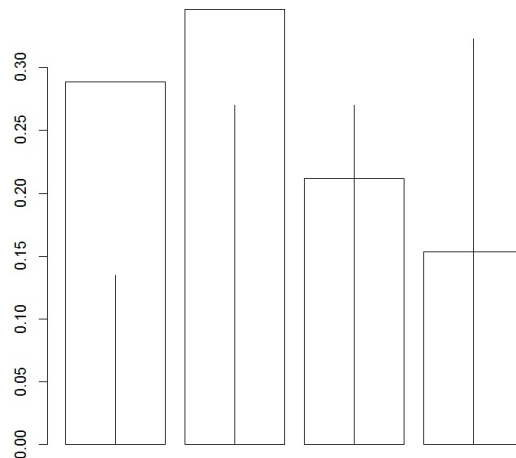
Solution 3. Par l'énoncé, on observe la valeur de X sur chacun des n individus (agents) d'un échantillon avec $n = 52$: (x_1, \dots, x_n) (avec $x_i \in \mathbb{N}$). On forme alors un vecteur des effectifs $(n_1, n_2, n_3, n_4) = (15, 18, 11, 8)$. Dire que X suit la loi de Poisson $\mathcal{P}(2)$ signifie que

$$\mathbb{P}(X = i) = e^{-2} \frac{2^i}{i!}, \quad i \in \mathbb{N}.$$

La problématique est la suivante : est-ce que ces données nous permettent d'affirmer que X ne suit pas la loi $\mathcal{P}(2)$? Ainsi, pour une analyse graphique, on propose les commandes :

```
nb = c(15, 18, 11, 8)
bar = barplot(nb / 52, col = "white")
lambda = sum((0:4) * nb) / 52
prob = c(dpois(0:3, lambda), 1 - ppois(3, lambda))
points(bar, prob, type = "h")
```

Cela renvoie :



Les différences observées laissent penser que X ne suit pas la loi de Poisson $\mathcal{P}(2)$. Confirmons cela avec un test statistique. On considère alors les hypothèses :

$$H_0 : "X \text{ suit la loi de Poisson } \mathcal{P}(2)" \quad \text{contre} \quad H_1 : "X \text{ ne suit pas la loi de Poisson } \mathcal{P}(2)"$$

Les données sont regroupées en $k = 4$ classes : $C_1 = \{0\}$, $C_2 = \{1\}$, $C_3 = \{2\}$ et $C_4 = \{3, \dots, \infty\}$. Soit R une *var* suivant la loi de Poisson $\mathcal{P}(2)$. On a $R(\Omega) = \mathbb{N}$ et $\bigcup_{i=1}^4 C_i = \mathbb{N} = R(\Omega)$ (il n'y a pas d'ajustement à faire). On met en œuvre le test du Chi-deux en faisant :

```
chisq.test(nb, p = proba)$p.value
```

Cela renvoie : [1] 0.001508304

Notons qu'aucun "Warning message" n'apparaît ; les conditions d'application du test sont vérifiées.

Comme p-valeur $\in]0.001, 0.01]$, le rejet de H_0 est très significatif $\star\star$.

Par conséquent, au risque 5%, on peut affirmer que X ne suit pas la loi de Poisson $\mathcal{P}(2)$.

Solution 4. Par l'énoncé, on observe la valeur de X sur chacun des n individus (jours) d'un échantillon avec $n = 320$: (x_1, \dots, x_n) (avec $x_i \in \mathbb{N}$). Ces valeurs sont regroupées en $k = 8$ classes : $C_1 = \{0\}$, $C_2 = \{1\}, \dots, C_8 = \{7\}$, avec pour effectifs respectifs : $n_1 = 65, n_2 = 110, \dots, n_8 = 2$.

On considère les hypothèses :

H_0 : " X suit une loi de Poisson" contre H_1 : " X ne suit pas une loi de Poisson"

Soit R une *var* suivant la loi de Poisson $\mathcal{P}(\lambda)$ avec $\lambda > 0$ inconnu. On a $R(\Omega) = \mathbb{N}$.

Comme $\bigcup_{i=1}^8 C_i = \{0, \dots, 7\} \neq \mathbb{N}$, on ajuste la dernière classe comme : $C_8 = \{7, 8, \dots, \infty\}$ (de sorte à ce que $\bigcup_{i=1}^8 C_i = R(\Omega)$). Le paramètre λ étant inconnu, il faut l'estimer à l'aide des données. Comme $\mathbb{E}(X) = \lambda$, la méthode des moments nous assure qu'une estimation de λ est la moyenne \bar{x} . On peut la calculer en faisant :

```
nb = c(65, 110, 70, 48, 16, 5, 4, 2)
lambda = sum((0:7) * nb) / 320
lambda
```

Cela renvoie : [1] 1.628125

Cette estimation sera prise en compte dans la suite. Pour tout $i \in \{1, \dots, 7\}$, on a

$$p_i = \mathbb{P}(R \in C_i) = \mathbb{P}(R = i - 1) = e^{-1.628125} \frac{1.628125^{i-1}}{(i-1)!}.$$

De plus, on a

$$p_8 = \mathbb{P}(R \in C_8) = \mathbb{P}(R \geq 7) = \sum_{i=7}^{\infty} e^{-1.628125} \frac{1.628125^i}{i!}.$$

On peut obtenir ces probabilités avec les commandes :

```
proba = c(dpois(0:6, lambda), 1 - ppois(6, lambda))
proba
```

Dans un premier temps, on propose de mettre en œuvre le test du Chi-deux en faisant :

```
chisq.test(nb, p = proba)$p.value
```

Cela renvoie :

```
[1] 0.1057037
```

Warning message:

```
In chisq.test(nb, p = proba) : Chi-squared approximation may be incorrect
```

On dénombre alors deux problèmes :

- le logiciel n'a pas pris en compte le fait que λ a été estimé,
- il y a un "Warning message" nous avertissant que l'hypothèse : $np_i \geq 5$ pour tout $i \in \{1, \dots, 8\}$ n'est peut-être pas vérifiée.

Étudions ce dernier point :

```
320 * proba
```

Cela renvoie :

```
[1] 62.8151317 102.2708864 83.2548934 45.1831245 18.3909436 5.9885510
```

```
[7] 1.6250183 0.4714511
```

Comme $np_7 < 5$ et $np_8 < 5$, nous allons fusionner les classes C_6 , C_7 et C_8 , formant ainsi une nouvelle dernière classe : $C_6 = \{5, 6, \dots\}$ avec pour effectif : $n_6 = 5 + 4 + 2 = 11$. Il y a désormais $k = 6$ classes.

On vérifie alors que l'hypothèse est vérifiée avec cette nouvelle configuration :

```
nb2 = c(65, 110, 70, 48, 16, 11)
proba2 = c(dpois(0:4, lambda), 1 - ppois(4, lambda))
320 * proba2
```

Cela renvoie : [1] 62.81513 102.27089 83.25489 45.18312 18.39094 8.08502

Aucune valeur ne dépasse 5. Pour avoir la p-valeur associée au test du Chi-deux en prenant en compte le fait que l'on a estimé un paramètre, on considère le degré de liberté :

$\nu = k - 1 - \ell = 6 - 1 - 1 = 4$. On fait :

```
x2obs = chisq.test(nb2, p = proba2)$statistic
deg = 4
1 - pchisq(x2obs, deg)
```

Cela renvoie :

X-squared

0.3659448

On obtient alors la vraie p-valeur associée au test du Chi-deux (notons aussi qu'aucun "Warning message" n'apparaît).

Comme p-valeur > 0.05 , les données ne nous permettent pas de rejeter H_0 . Ainsi, on ne peut pas rejeter l'hypothèse selon laquelle X ne suit pas une loi de Poisson.

Solution 5. Par l'énoncé, on observe la valeur de X sur chacun des n individus (appareils) d'un échantillon avec $n = 7$: (x_1, \dots, x_n) (avec $x_i \in \mathbb{R}$).

On considère les hypothèses :

H_0 : " X suit la loi exponentielle $\mathcal{E}(0.01)$ " contre H_1 : " X ne suit pas la loi exponentielle $\mathcal{E}(0.01)$ "

Nous allons utiliser le test de Kolmogorov-Smirnov. On fait :

```
x = c(145, 110, 170, 48, 116, 95, 74 )
ks.test(x, "pexp", 0.01)$p.value
```

Cela renvoie : [1] 0.2005605

Comme p -valeur > 0.05 , les données ne nous permettent pas de rejeter H_0 . Ainsi, on ne peut pas rejeter l'hypothèse selon laquelle X suit la loi exponentielle $\mathcal{E}(0.01)$.

Solution 6. Par l'énoncé, on observe la valeur de X sur chacun des n individus (attentes) d'un échantillon avec $n = 7 : (x_1, \dots, x_n)$ (avec $x_i \in \mathbb{R}$). On détermine $1/\bar{x}$ en faisant :

```
x = c(25.12, 12.36, 24.35, 12.19, 5.27, 18.35, 19.11, 27.08, 21.09, 17.19,
      8.45, 13.27, 15.17)
1 / mean(x)
```

Cela renvoie : [1] 0.05936073

On considère les hypothèses :

H_0 : " X suit la loi exponentielle $\mathcal{E}(0.05936073)$ " contre

H_1 : " X ne suit pas la loi exponentielle $\mathcal{E}(0.05936073)$ "

Nous allons utiliser le test de Kolmogorov-Smirnov. On fait :

```
ks.test(x, "pexp", 0.05936073)$p.value
```

Cela renvoie : [1] 0.05028391

Comme p -valeur > 0.05 (de justesse), les données ne nous permettent pas de rejeter H_0 . Ainsi, on ne peut pas rejeter l'hypothèse selon laquelle X suit la loi exponentielle $\mathcal{E}(0.05936073)$.

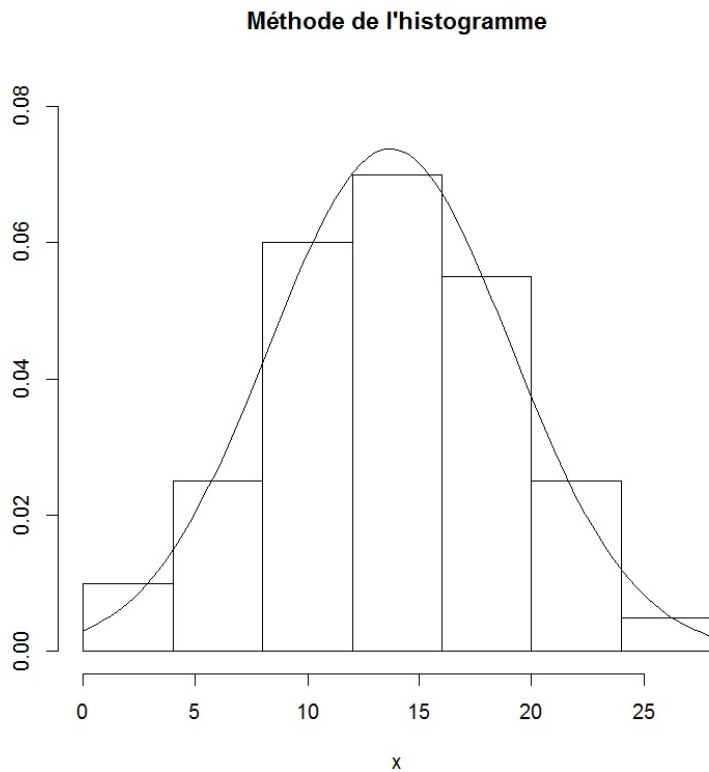
Solution 7. Par l'énoncé, on observe la valeur de X sur chacun des n individus (pourceaux) d'un échantillon avec $n = 50 : (x_1, \dots, x_n)$ (avec $x_i \in \mathbb{R}$). Celles-ci sont présentées sous la forme d'un tableau classes-effectifs. On propose d'extrapoler les données en utilisant les centres des classes et les effectifs respectifs :

```
xinf = c(0, 4, 8, 12, 16, 20, 24)
xsup = c(4, 8, 12, 16, 20, 24, 28)
centre = (xinf + xsup) / 2
n = c(2, 5, 12, 14, 11, 5, 1)
x = rep(centre, n)
```

Pour une analyse graphique, on peut utiliser la méthode de l'histogramme. On propose les commandes :

```
hist(x, freq = FALSE, breaks = c(0, 4, 8, 12, 16, 20, 24, 28),
     main = "Méthode de l'histogramme", ylim = c(0, 0.082), ylab = "")
a = mean(x) ; b = sd(x)
curve(dnorm(x, a, b), add = TRUE)
```

Cela renvoie :



Les différences observées laissent penser que X suit une loi normale.

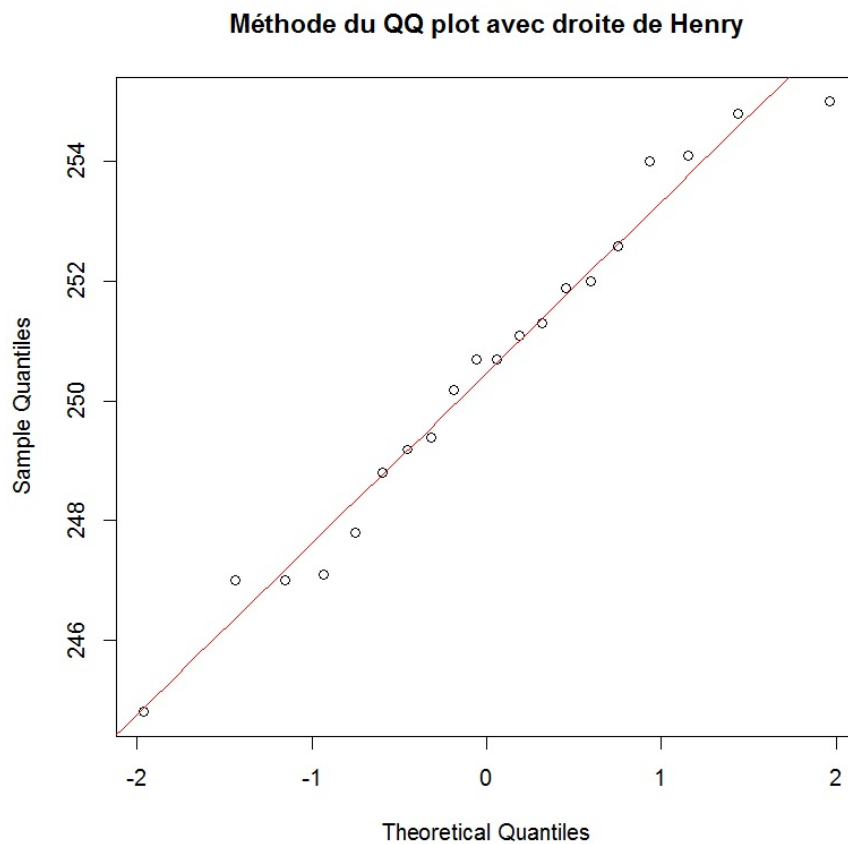
Solution 8. Par l'énoncé, on observe la valeur de X sur chacun des n individus (plaquettes de beurre) d'un échantillon avec $n = 20 : (x_1, \dots, x_n)$ (avec $x_i \in \mathbb{R}$).

Pour une analyse graphique, on peut utiliser la méthode du QQ plot avec la droite de Henry.

On propose les commandes :

```
x = c(247.0, 247.8, 250.2, 251.3, 251.9, 249.4, 248.8, 247.1, 255.0, 247.0,  
254.8, 244.8, 250.7, 250.7, 252.6, 251.1, 254.1, 249.2, 252.0, 254.0)  
qqnorm(x, main = "Méthode du QQ plot avec droite de Henry")  
a = mean(x) ; b = sd(x)  
curve(a + b * x, -6, 6, col = "red", add = TRUE)
```

Cela renvoie :



On constate que la droite de Henry ajuste bien le nuage de points, ce qui traduit le fait que X suit une loi normale. Confirmons cela avec un test statistique. On considère les hypothèses :

$$H_0 : "X \text{ suit une loi normale}" \quad \text{contre} \quad H_1 : "X \text{ ne suit pas une loi normale}"$$

On utilise le test de Shapiro-Wilk :

```
shapiro.test(x)$p.value
```

Cela renvoie : [1] 0.751598

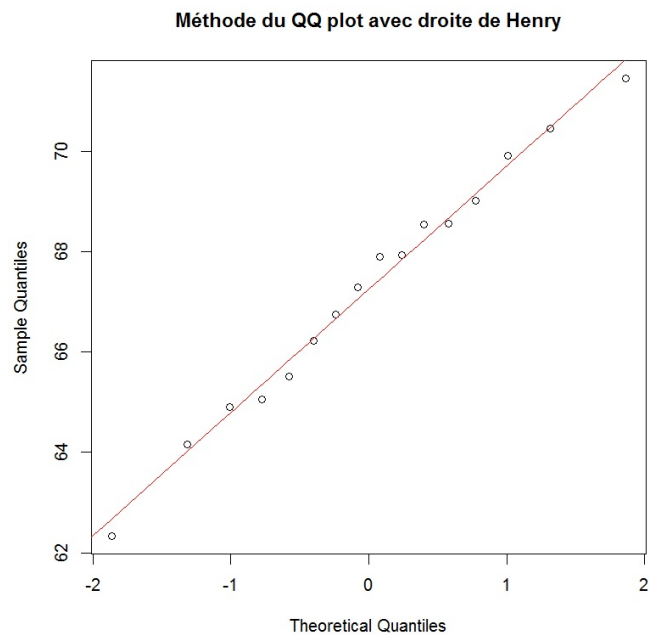
Comme $p\text{-valeur} > 0.05$, on ne rejette pas H_0 ; l'hypothèse selon laquelle X suit une loi normale n'est pas rejetée.

Solution 9. Par l'énoncé, on observe la valeur de X sur chacun des n individus (kiwis) d'un échantillon avec $n = 16$: (x_1, \dots, x_n) (avec $x_i \in \mathbb{R}$).

Pour une analyse graphique, on peut utiliser la méthode du QQ plot avec la droite de Henry. On propose les commandes :

```
x = c(65.06, 71.44, 67.93, 69.02, 67.28, 62.34, 66.23, 64.16, 68.56, 70.45,  
64.91, 69.90, 65.52, 66.75, 68.54, 67.90)  
qqnorm(x, main = "Méthode du QQ plot avec droite de Henry")  
a = mean(x) ; b = sd(x)  
curve(a + b * x, -6, 6, col = "red", add = TRUE)
```

Cela renvoie :



On constate que la droite de Henry ajuste bien le nuage de points, ce qui traduit le fait que X suit une loi normale. Confirmons cela avec un test statistique. On considère les hypothèses :

$$H_0 : "X \text{ suit une loi normale}" \quad \text{contre} \quad H_1 : "X \text{ ne suit pas une loi normale}"$$

On utilise le test de Shapiro-Wilk :

```
shapiro.test(x)$p.value
```

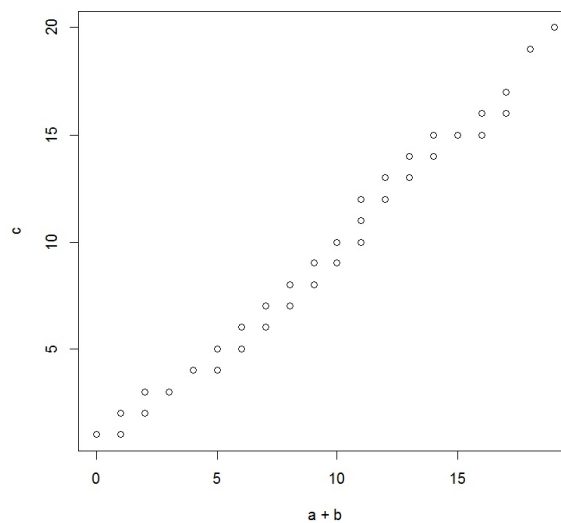
Cela renvoie : [1] 0.9971045

Comme $p\text{-valeur} > 0.05$, on ne rejette pas H_0 ; l'hypothèse selon laquelle X suit une loi normale n'est pas rejetée.

Solution 10. On fait :

```
n = 1000
a = rpois(n, 5)
b = rpois(n, 3)
c = rpois(n, 8)
qqplot(a + b, c)
```

Cela renvoie :



On constate que le nuage de points peut être ajusté par la droite $y = x$, ce qui illustre le résultat.

Solution 14.

– $X_1 \sim \mathcal{E}(\lambda)$, $X_2 \sim \mathcal{E}(\lambda)$, X_1 et X_2 indépendantes entraînent $X_1 + X_2 \sim \Gamma(\lambda, 2)$, avec $\lambda = 3.8$:

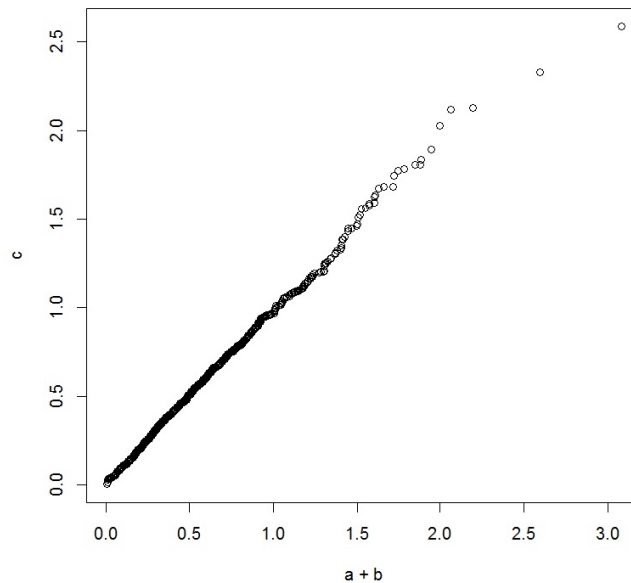
```
a = rexp(1000, 3.8)
```

```
b = rexp(1000, 3.8)
```

```
c = rgamma(1000, 2, 3.8)
```

```
qqplot(a + b, c)
```

Cela renvoie :



On constate que le nuage de points peut être ajusté par la droite d'équation $y = x$, ce qui illustre le résultat.

– $X_1 \sim \Gamma(m_1, \lambda)$, $X_2 \sim \Gamma(m_2, \lambda)$, X_1 et X_2 indépendantes entraînent $X_1 + X_2 \sim \Gamma(m_1 + m_2, \lambda)$, avec $m_1 = m_2 = 4.2$ et $\lambda = 2.1$:

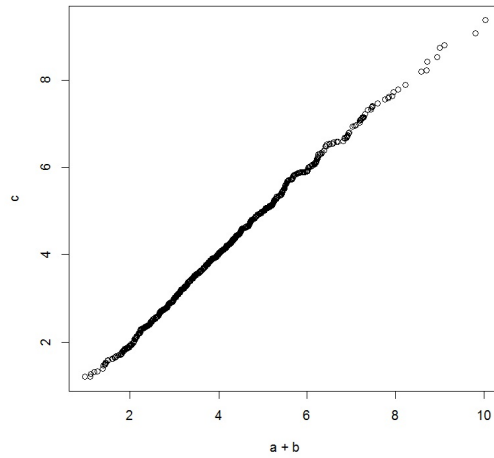
```
a = rgamma(1000, 4.2, 2.1)
```

```
b = rgamma(1000, 4.2, 2.1)
```

```
c = rgamma(1000, 8.4, 2.1)
```

```
qqplot(a + b, c)
```

Cela renvoie :



De nouveau, on constate que le nuage de points peut être ajusté par la droite d'équation $y = x$, ce qui illustre le résultat.

– $X_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$, $X_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$, X_1 et X_2 indépendantes entraînent

$X_1 + X_2 \sim \mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$, avec $\mu_1 = \mu_2 = 1.6$ et $\sigma_1 = \sigma_2 = 1.5$:

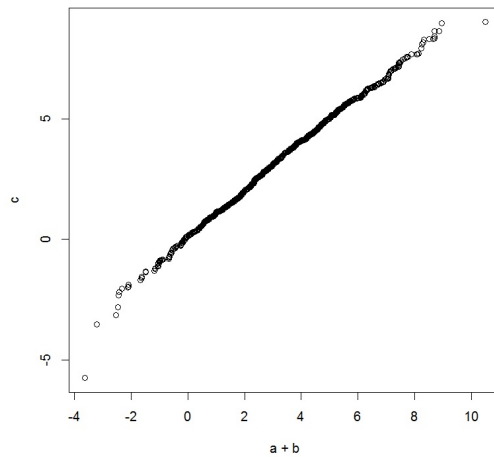
```
a = rnorm(1000, 1.6, 1.5)
```

```
b = rnorm(1000, 1.6, 1.5)
```

```
c = rnorm(1000, 3.2, sqrt(1.5^2 + 1.5^2))
```

```
qqplot(a + b, c)
```

Cela renvoie :



De nouveau, on constate que le nuage de points peut être ajusté par la droite d'équation $y = x$, ce qui illustre le résultat.

- $X_1 \sim \chi^2(\nu_1)$, $X_2 \sim \chi^2(\nu_2)$, X_1 et X_2 indépendantes entraînent $X_1 + X_2 \sim \chi^2(\nu_1 + \nu_2)$, avec $\nu_1 = \nu_2 = 3.2$:

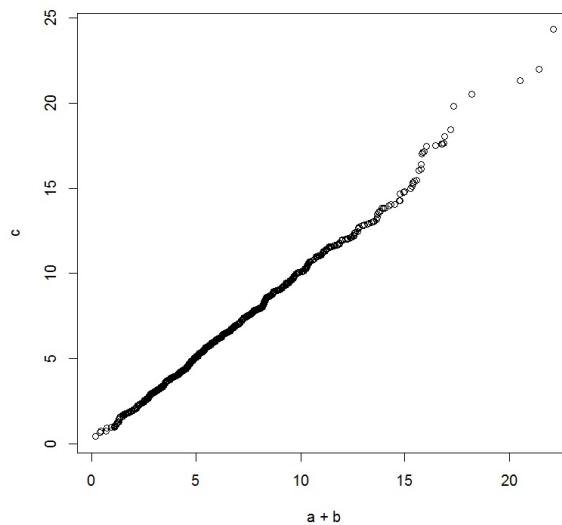
```
a = rchisq(1000, 3.2)
```

```
b = rchisq(1000, 3.2)
```

```
c = rchisq(1000, 6.4)
```

```
qqplot(a + b, c)
```

Cela renvoie :



De nouveau, on constate que le nuage de points peut être ajusté par la droite d'équation $y = x$, ce qui illustre le résultat.

Solution 15.

- Caractérisation de la loi du chi-deux $\chi^2(2)$: Si $X \sim \mathcal{N}(0, 1)$ et $Y \sim \mathcal{N}(0, 1)$, alors

$X^2 + Y^2 \sim \chi^2(2)$:

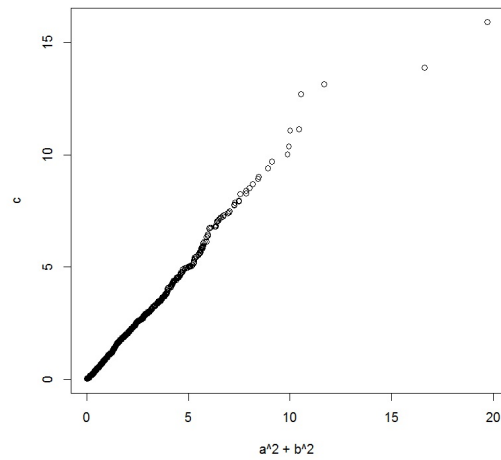
```
a = rnorm(1000)
```

```
b = rnorm(1000)
```

```
c = rchisq(1000, 2)
```

```
qqplot(a^2 + b^2, c)
```

Cela renvoie :



On constate que le nuage de points peut être ajusté par la droite d'équation $y = x$, ce qui illustre le résultat.

- Caractérisation de la loi de Student $\mathcal{T}(\nu)$: Si $X \sim \mathcal{N}(0, 1)$ et $Y \sim \chi^2(\nu)$, alors $\frac{X}{\sqrt{\frac{Y}{\nu}}} \sim \mathcal{T}(\nu)$, avec $\nu = 3.9$:

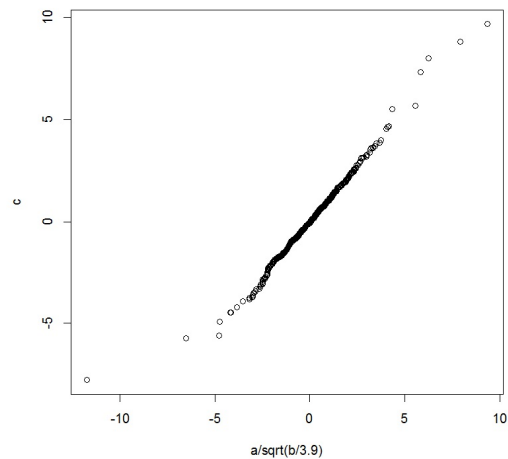
```
a = rnorm(1000)
```

```
b = rchisq(1000, 3.9)
```

```
c = rt(1000, 3.9)
```

```
qqplot(a / sqrt(b / 3.9), c)
```

Cela renvoie :



De nouveau, on constate que le nuage de points peut être ajusté par la droite d'équation $y = x$, ce qui illustre le résultat.

– Caractérisation de la loi de Fisher $\mathcal{F}(\nu_1, \nu_2)$: Si $X \sim \chi^2(\nu_1)$ et $Y \sim \chi^2(\nu_2)$, alors

$$\frac{\frac{X}{\nu_1}}{\frac{Y}{\nu_2}} = \frac{\nu_2 X}{\nu_1 Y} \sim \mathcal{F}(\nu_1, \nu_2), \text{ avec } (\nu_1, \nu_2) = (2.1, 8.3) :$$

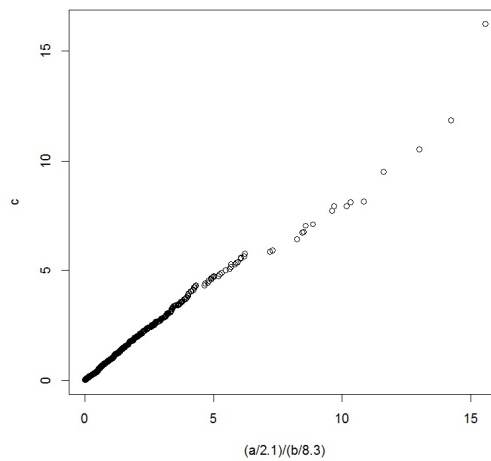
```
a = rchisq(1000, 2.1)
```

```
b = rchisq(1000, 8.3)
```

```
c = rf(1000, 2.1, 8.3)
```

```
qqplot((a / 2.1) / (b / 8.3), c)
```

Cela renvoie :



De nouveau, on constate que le nuage de points peut être ajusté par la droite d'équation $y = x$, ce qui illustre le résultat.

Solution 16. Dans les commandes, on génère 100 observations d'une *var* X suivant la loi normale centrée réduite $\mathcal{N}(0, 1)$. On les place dans un vecteur \mathbf{x} . Ensuite, on crée un vecteur numérique \mathbf{a} . Dans la boucle `for`, pour tout $i \in \{1, \dots, 6\}$, on considère les hypothèses suivantes :

$$H_0(i) : "X \text{ suit la loi normale } \mathcal{N}\left(0, \left(1 + \frac{i-1}{10}\right)^2\right)" \text{ contre}$$

$$H_1(i) : "X \text{ ne suit pas la loi normale } \mathcal{N}\left(0, \left(1 + \frac{i-1}{10}\right)^2\right)",$$

on utilise le test de Kolmogorov-Smirnov, on calcule la p-valeur associée à $H_0(i)$ et on met cette valeur au i -ème élément de \mathbf{a} . Ensuite, on affiche \mathbf{a} .

On remarque alors que, logiquement, plus i est grand, plus $1 + (i - 1)/10$ s'éloigne de 1, plus la p-valeur est petite, moins on a de certitude en affirmant que la loi de X est en adéquation avec la loi normale $\mathcal{N}\left(0, \left(1 + \frac{i-1}{10}\right)^2\right)$.