

LSTAT2100 - Exercices - Série 1

Solutions

Exercice 1

Nous nous intéressons à la couleur des yeux d'une certaine population. Pour $n = 300$ individus, pris au hasard, nous avons observé les chiffres suivants.

bleu	marron	noir	vert
48	122	95	35

(a) Ces couleurs sont-elles toutes réparties de manière uniforme (équiprobables) ? Répondez à cette question en utilisant le test de LR. Calculez la p -valeur de ce dernier en utilisant (i) la théorie asymptotique, et (ii) des simulations.

Solution:

Soit p_k , $k = 1, \dots, 4$, les proportions des couleurs (bleu, vert, marron, noir). L'hypothèse à tester est la suivante

$$H_0 : p_1 = p_2 = p_3 = p_4 = 1/4.$$

```
# Les données
O <- c(48, 122, 95, 35)
E <- 300*rep(1 / 4, 4)

# LR asymptotique
(2*sum(O * log(O/E))) |> print() |> pchisq(df = 3, lower.tail = FALSE)

[1] 67.4
[1] 1.51e-14

# LR via des simulations
{rmultinom(10000, size = 300, prob = rep(1 / 4, 4)) |>
  apply(2, FUN = \(O) 2*sum(O * log(O/E))) >= 67.4} |> mean()

[1] 0
```

Au seuil de 5%, nous rejetons donc l'hypothèse nulle que les proportions sont égales.

(b) Peut-on dire que les yeux foncés (marron et noir) sont deux fois plus probables que les yeux clairs (bleu et vert) ? Utilisez le test de Pearson.

Solution:

Ici, le but est de tester $H_0 : p_1 + p_4 = p_2 + p_3$.

```
chisq.test(c(O[1] + O[4], O[2] + O[3]), p = c(1 / 3, 2 / 3))
```

Chi-squared test for given probabilities

data: c(0[1] + 0[4], 0[2] + 0[3])
X-squared = 4, df = 1, p-value = 0.04

Au seuil de 5%, nous rejetons l'hypothèse H_0 .

(c) Les proportions des yeux bleus et des yeux verts sont-elles les mêmes ? Utilisez un test LR. Détaillez votre approche et vos calculs.

Solution:

Ici, le but est de tester $H_0 : p_1 = p_4$.

Pour réaliser le test, nous avons besoin des effectifs attendus sous H_0 . Pour cela, calculons $\hat{p}_{01}, \hat{p}_{02}, \hat{p}_{03}$ et \hat{p}_{04} les estimateurs de maximum de vraisemblance de p_1, p_2, p_3 et p_4 sous H_0 . La log-vraisemblance est

$$l_n = n_1 \log p_1 + n_2 \log p_2 + n_3 \log p_3 + n_4 \log p_4 + \text{const}, \text{ avec } p_1 + p_2 + p_3 + p_4 = 1.$$

Sous H_0 ,

$$l_n = (n_1 + n_4) \log p_1 + n_2 \log p_2 + n_3 \log p_3 + \text{const}, \text{ avec } p_3 = 1 - 2p_1 - p_2$$

Dès lors,

$$\frac{\partial l}{\partial p_1} = \frac{n_1 + n_4}{p_1} - 2 \frac{n_3}{p_3} \text{ et } \frac{\partial l}{\partial p_2} = \frac{n_2}{p_2} - \frac{n_3}{p_3}.$$

En annulant ces deux dérivées partielles, on obtient

$$(n_1 + n_4)\hat{p}_{03} = 2n_3\hat{p}_{01} \text{ et } n_2\hat{p}_{03} = n_3\hat{p}_{02}$$

En sommant les deux égalités, on obtient

$$(n_1 + n_2 + n_4)\hat{p}_{03} = n_3(2\hat{p}_{01} + \hat{p}_{02}) \Rightarrow (n_1 + n_2 + n_4)\hat{p}_{03} = n_3(1 - \hat{p}_{03}) \Rightarrow \hat{p}_{03} = \frac{n_3}{n}$$

Et en remplaçant dans les expressions ci-dessus, on trouve,

$$\hat{p}_{01} = \hat{p}_{04} = \frac{n_1 + n_4}{2n}; \hat{p}_{02} = \frac{n_2}{n}; \hat{p}_{03} = \frac{n_3}{n}$$

```
E <- c((0[1] + 0[4]) / 2, 0[2], 0[3], (0[1] + 0[4]) / 2)
t(cbind(0, E))
```

O	48.0	122	95	35.0
E	41.5	122	95	41.5

```
sum((E - 0)^2 / E) |> pchisq(df = 1, lower.tail = FALSE)
```

```
[1] 0.154
```

\Rightarrow Non-rejet de H_0 , à 5%.

Remarque: Dans ce cas particulier, on peut tester la légalité des deux proportions en ne considérant que la sous-table composée des yeux verts et bleus

bleu	vert
48	35

```
Ob <- c(48, 35)
Eb <- c((Ob[1] + Ob[2]) / 2, (Ob[1] + Ob[2]) / 2)
sum((Eb - Ob)^2 / Eb) |> pchisq(df = 1, lower.tail = FALSE)
```

```
[1] 0.154
```

Cela fonctionne car les effectifs attendus des autres modalités (marron, noir) sont identiques aux effectifs observés (voir le tableau `t(cbind(0, E))` ci-dessous).

Exercice 2

Voici les fréquences (*Freq*) du nombre de passages de bus (*nBus*) par heure à un point d'arrêt. Ces données concernent une période de 30 heures réparties sur 5 jours de la semaine.

nBus	0	1	2	3	4	5
Freq	1	5	6	10	4	4

(a) Supposons que $nBus \sim Pois(\mu)$, $\mu > 0$. Tester l'hypothèse $H_0 : \mu = 3$ vs $H_1 : \mu \neq 3$. Pour ce faire, utilisez les trois tests classiques vus dans le cours. Calculez un intervalle de confiance pour μ .

Solution:

```
nBus <- 0:5
Freq <- c(1, 5, 6, 10, 4, 4)
data <- data.frame(nBus, Freq)
t(data)
```

nBus	0	1	2	3	4	5
Freq	1	5	6	10	4	4

L'EMV de μ n'est rien d'autre que la moyenne (voir cours)

```
n <- sum(Freq)
mu.hat <- {sum(Freq * nBus) / n} |> print()
```

```
[1] 2.77
```

En utilisant les formules vues au cours, nous calculons les trois statistiques et leurs p -valeurs de la manière suivante

```
mu.0 <- 3

# Wald
{(mu.hat - mu.0)^2*(n/mu.hat)} |> print() |> pchisq(df = 1, lower.tail = FALSE)
```

```
[1] 0.59
```

```
[1] 0.442
```

```
# Score
{(mu.hat - mu.0)^2*(n/mu.0)} |> print() |> pchisq(df = 1, lower.tail = FALSE)
```

```
[1] 0.544
```

```
[1] 0.461
```

```
# LR
{2 * n * (mu.hat * log(mu.hat / mu.0) - (mu.hat - mu.0))} |> print() |>
  pchisq(df = 1, lower.tail = FALSE)
```

```
[1] 0.559
```

```
[1] 0.455
```

Pour l'intervalle de confiance, nous pouvons utiliser le fait que $\hat{\mu} \sim_a N(\mu, \mu/n)$, pour établir l'intervalle (de Wald, asymptotique) suivant

```
mu.hat + sqrt(mu.hat / n) * qnorm(0.975) * c(-1, 1)
```

```
[1] 2.17 3.36
```

(b) En supposant que la fréquence moyenne de passage est de $\mu = 3$ par heure, testez l'ajustement d'une loi de Poisson ($H_0 : nBus \sim Pois(3)$) aux données ? *Proposez un test pertinent.*

Solution:

Tout d'abord, notez que l'information dont on dispose peut être écrite sous la forme suivante.

```
cbind(t(data), c("+6", 0))
```

nBus	0	1	2	3	4	5	+6
Freq	1	5	6	10	4	4	0

Dès lors, l'hypothèse à tester est

$$H_0 : p_k = p_k^0, \text{ pour } k = 0, \dots, 5 \text{ et } p_{6+} = p_{6+}^0,$$

où $p_k = P(nBus = k)$, $p_{6+} = P(nBus \geq 6)$, $p_k^0 = P(Pois(3) = k)$, et $p_{6+}^0 = P(Pois(3) \geq 6)$.

```
# Effectifs observés/attendus
p0 <- dpois(x = 0:5, lambda = 3)
p0 <- c(p0, 1 - sum(p0))
E <- 30 * p0
O <- c(Freq, 0)
data.frame(O = O, E = E) |> t()
```

O	1.00	5.00	6.00	10.00	4.00	4.00	0.00
E	1.49	4.48	6.72	6.72	5.04	3.02	2.52

Nous pouvons alors calculer la statistique de Pearson et sa p -valeur comme suit

```
sum((O - E)^2 / E) |> print() |> pchisq(df = 6, lower.tail = FALSE)
```

```
[1] 4.95
```

```
[1] 0.551
```

ou, directement via la fonction `chisq.test()`:

```
chisq.test(x = 0, p = p0)
```

Warning in `chisq.test(x = 0, p = p0)`: Chi-squared approximation may be incorrect

Chi-squared test for given probabilities

```
data: 0
X-squared = 4.95, df = 6, p-value = 0.55
```

Ce résultat doit être considéré avec prudence, car les effectifs attendus sont faibles (< 5). Au lieu de nous référer à la distribution asymptotique de χ_6^2 , nous pouvons refaire le test en calculant la p -valeur par simulation

```
chisq.test(x = 0, p = p0, simulate.p.value = TRUE, B = 10000)
```

Chi-squared test for given probabilities with simulated p-value (based on 10000 replicates)

```
data: 0
X-squared = 4.95, df = NA, p-value = 0.56
```

(c) Refaites le même test que celui de la question précédente mais, cette fois, sans supposer que $\mu = 3$; c'est-à-dire testez $H_0 : nBus \sim Pois(\mu)$, pour un μ fixe mais inconnue.

Solution:

À la différence du cas précédent il faudra ici estimer μ par la méthode de maximum de vraisemblance. Nous savons que l'EMV de μ est

```
mu.hat
```

```
[1] 2.7667
```

Par la suite on applique la même démarche que pour la question précédente, mais, attention, le degré de liberté change de 6 à 5 (à cause de l'estimation de μ).

```
# Effectifs observés/attendus
p0 <- dpois(x = 0:5, lambda = mu.hat)
p0 <- c(p0, 1 - sum(p0))
E <- 30 * p0
data.frame(O = 0, E = E) |> t()
```

O	1.0000	5.0000	6.0000	10.0000	4.0000	4.0000	0.0000
E	1.8861	5.2183	7.2187	6.6572	4.6046	2.5479	1.8672

```
# Test de Pearson - asymptotique
sum((O - E)^2 / E) |> print() |> pchisq(df = 5, lower.tail = FALSE)
```

```
[1] 5.084
```

```
[1] 0.40572
```

```
# Test de Pearson - simulation
{rmultinom(10000, size = 30, prob = p0) |>
  apply(2, FUN = \(O) sum((O - E)^2 / E)) >= 5.084} |> mean()
```

[1] 0.5214

Exercice 3

Dans cet exercice, nous allons utiliser le jeu de données `mdata.csv`. Ce jeu de données comprend des observations liées à des clients qui ont contracté un crédit.

`mdata.csv` contient de nombreuses variables, mais nous n'en utiliserons que quelques-unes dans cette série et dans la série qui suit. Parmi ces variables, il y a

- **repay**: variable binaire prenant les valeurs "Default" ou "NotInDefault" pour un crédit remboursé ou non à temps.
- **account**: flux mensuels moyens sur le compte courant du client; "<0", "[0-200)", ">=200", ou "No acc", si le client n'a pas de compte courant auprès de la banque.
- **tel**: le client dispose-t-il d'un numéro de téléphone fixe ? "Yes" ou "No".
- **employ**: la situation professionnelle du client en termes de durée de l'emploi (actuel) en années; "No or <1", pour un client sans-emploi ou qui travaille mais depuis moins d'un an, "[1-4)", "[4-7)", ou ">=7", pour un client qui exerce son activité depuis au moins 7 ans.

Charger et examiner les données. Pour ce faire, vous pouvez utiliser la commande suivante (on présume que votre répertoire de travail comprend `Data/mdata.csv`).

```
mdata <- read.csv(file = "Data/mdata.csv")
str(mdata)
```

(a) Donnez le tableau des fréquences pour **repay**. Représentez ce tableau à l'aide d'un graphique approprié.

Solution:

Pour commencer, il convient de transformer les variables catégorielles en facteurs.

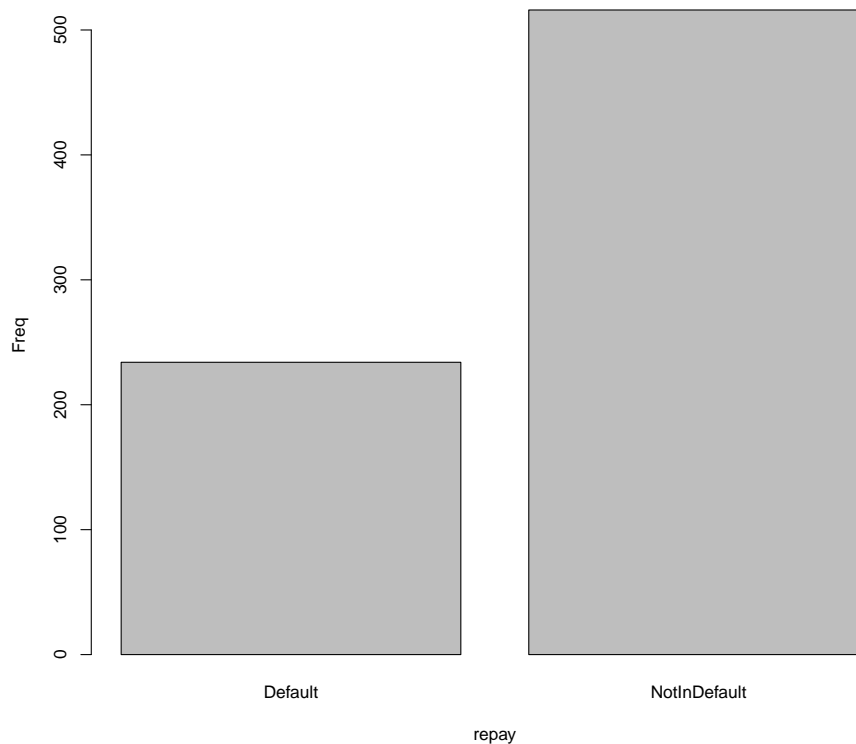
```
mdata <- transform(mdata, repay = factor(repay),
                  account = factor(account, levels = c("<0", "[0-200)", ">=200", "No acc")),
                  tel = factor(tel),
                  employ = factor(employ, levels = c("No or <1", "[1-4)", "[4-7)", ">=7")))
str(mdata)
```

Voici le tableau de contingence demandé

```
tbl <- xtabs(~ repay, data = mdata) |> print()      #ou table(mdata$repay)
```

```
repay
      Default NotInDefault
      234         516
```

```
barplot(tbl, xlab = "repay", ylab = "Freq")
```



(b) Construisez la table de contingence croisant les variables `repay` et `account`, et estimez les proportions $P(\text{repay}|\text{account})$, pour les différentes valeurs de `(repay, account)`.

Solution:

```
tbl <- xtabs(~ account + repay, data = mdata) |> print()
```

```
      repay
account Default NotInDefault
<0          105          103
[0-200)      81          115
>=200        10           33
No acc       38          265
```

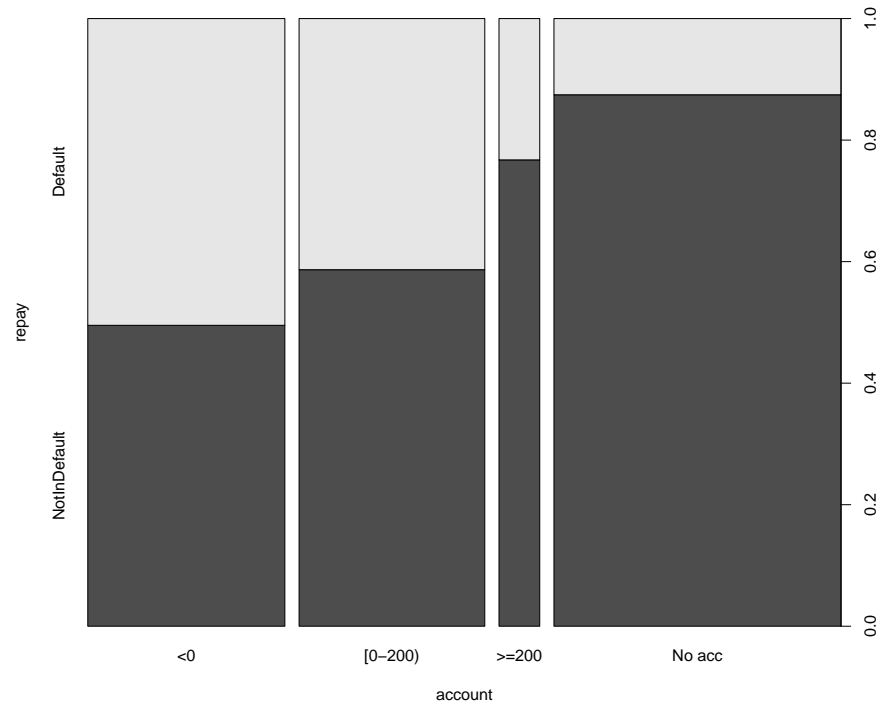
```
proportions(tbl, margin = "account")
```

account/repay	Default	NotInDefault
<0	0.50481	0.49519
[0-200)	0.41327	0.58673
>=200	0.23256	0.76744
No acc	0.12541	0.87459

(c) Faites un graphique pour représenter la distribution (marginale) de `account` et la distribution (conditionnelle) de `repay|account`. Que suggère ce graphique quant à l'association entre ces deux variables ?

Solution:

```
spineplot(tbl)
```



Ce graphique montre une forte association entre les deux variables étudiées.

(d) Testez l'indépendance entre `account` et `repay`.

Solution:

Nous avons plusieurs choix/fonctions pour effectuer le test:

```
summary(tbl)
# ou
chisq.test(tbl)
# ou
vcd::assocstats(tbl)$chisq_tests
```

	X^2	df	P(> X^2)
Likelihood Ratio	101.471	3	0
Pearson	95.793	3	0

⇒ il existe bien une association très significative (à 5%) entre les deux variables étudiées.

(e) Calculez le rapport de cotes entre `tel` et `repay`. Quelle information peut-on en tirer sur l'association entre ces deux variables ? Compléter par un test statistique (d'indépendance) fondé sur le ratio calculé.

Solution:

Nous pouvons utiliser la fonction `loddsratio` du package `vcd`

```
tbl <- xtabs(~ tel + repay, data = mdata) |> print()
```

```
      repay
tel    Default NotInDefault
No      142      305
Yes      92      211
```

```
or <- loddsratio(tbl, log = FALSE) |> print()
```

odds ratios for tel and repay

```
[1] 1.0678
```

Le rapport de cotes est proche de 1, ce qui indique que la proportion de personnes qui remboursent leur crédit à temps (ou non) est similaire chez les détenteurs et les non-détenteurs d'un fixe. En effet,

```
tbl |> proportions("tel")
```

tel/repay	Default	NotInDefault
No	0.31767	0.68233
Yes	0.30363	0.69637

Nous pouvons construire un intervalle de confiance pour le rapport de cotes. Ensuite, il suffira de vérifier si cet intervalle contient la valeur 1 ou non pour confirmer ou infirmer l'indépendance.

```
confint(or)
```

	2.5 %	97.5 %
No:Yes/Default:NotInDefault	0.77878	1.464

L'intervalle du rapport des cotes contient 1, ce qui signifie que le fait de rembourser son crédit à temps est indépendant de la possession ou non d'un fixe.

Exercice 4

Soit le tableau de de contingence suivant

X/Y	1	2
1	5	48
2	34	251

La signification des variables X et Y n'est pas importante pour la suite.

(a) Tester l'indépendance entre X et Y .

Solution:

```
0 <- c(5, 34, 48, 251)
dt <- data.frame(X = c(1, 2, 1, 2), Y = c(1, 1, 2, 2), Freq = 0)
tab <- xtabs(Freq ~ X + Y, data = dt)
tab
```

X/Y	1	2
1	5	48
2	34	251

```
summary(tab)$p.value
```

```
[1] 0.6015
```

(c) Soit $p_{ij} = P(X = i, Y = j)$, $i, j = 1, 2$. Considérons l'hypothèse suivante

$$H_0 : p_{11} = \theta^2, p_{12} = p_{21} = \theta(1 - \theta), \text{ et } p_{22} = (1 - \theta)^2$$

Montrer que sous H_0 , X et Y sont indépendants et identiquement distribués.

Solution:

X et Y sont identiquement distribuées puisque

$$p_{1.} = p_{.1} = \theta \text{ et } p_{2.} = p_{.2} = 1 - \theta.$$

Ces variables sont indépendantes car

$$p_{ij} = p_{i.} p_{.j} \quad i, j = 1, 2.$$

(d) En supposant un échantillonnage multinomial simple, donner l'estimateur du maximum de vraisemblance de θ et calculer le.

Solution:

Nous avons que le log-vraisemblance est donné par

$$l_n(\theta) = 2n_{11} \ln \theta + 2n_{22} \ln(1 - \theta) + (n_{12} + n_{21}) \ln \theta(1 - \theta) + \text{const}$$

Dès lors,

$$\begin{aligned} \frac{dl_n}{d\theta}(\theta) &= 0 \\ \iff 2n_{11}(1 - \theta) - 2n_{22}\theta + (n_{12} + n_{21})(1 - 2\theta) &= 0 \\ \iff 2n\theta &= 2n_{11} + n_{12} + n_{21} \\ \iff \theta &= \frac{n_{1.} + n_{.1}}{2n} = \frac{p_{1.} + p_{.1}}{2} \end{aligned}$$

Donc

$$\hat{\theta} = \frac{\hat{p}_{1.} + \hat{p}_{.1}}{2}.$$

On peut calculer cet estimateur “manuellement” ou en utilisant R à l'aide du code suivant

```
mptab <- tab |> proportions() |> addmargins() |> print()
```

	Y		
X	1	2	Sum
1	0.014793	0.142012	0.156805
2	0.100592	0.742604	0.843195
Sum	0.115385	0.884615	1.000000

```
theta <- {(mptab[1, 3] + mptab[3, 1]) / 2} |> print()
```

```
[1] 0.13609
```

(e) Proposer une statistique de test pour tester H_0 . Effectuez le test et concluez.

Solution:

On peut utiliser le test du rapport de vraisemblance. Sa statistique est donnée par (voir cours)

$$G^2 = 2 \sum O \log \left(\frac{O}{E} \right),$$

avec $O = (N_{11}, N_{12}, N_{21}, N_{22})$ et $E = (n\hat{\theta}^2, n\hat{\theta}(1 - \hat{\theta}), n\hat{\theta}(1 - \hat{\theta}), n(1 - \hat{\theta})^2)$. çàd

$$G^2 = 2 \left(N_{11} \ln \frac{N_{11}}{n\hat{\theta}^2} + N_{12} \ln \frac{N_{12}}{n\hat{\theta}(1 - \hat{\theta})} + N_{21} \ln \frac{N_{21}}{n\hat{\theta}(1 - \hat{\theta})} + N_{22} \ln \frac{N_{22}}{n(1 - \hat{\theta})^2} \right).$$

Sous H_0 , Cette variable suit asymptotiquement une distribution chi-deux de $3 - 1$ degrés de liberté. En effet, pour la vraisemblance non contrainte il y a trois paramètres à estimer (à savoir p_{11}, p_{12} , et p_{21}) alors qu'il n'y a qu'un seul paramètre à estimer sous H_0 (à savoir θ).

```
n <- sum(tab) # Nombre d'observations
E <- n * c(theta^2, theta * (1 - theta), theta * (1 - theta), (1 - theta)^2)
# la statistique de rapport de vraisemblance
g2 <- 2 * sum(O * log(O / E))
g2
```

```
[1] 2.7601
```

```
# pvalueur
pchisq(g2, df = 2, lower.tail = FALSE)
```

```
[1] 0.25156
```