

CORRIGÉ 9

Exercice 1. (a) Notons p le pourcentage recherché, et considérons $p \in (0, 1)$. Si on choisit par hasard une personne parmi les étudiant-e-s de l'EPFL/UNIL, celle-ci sera une femme avec la probabilité p et un homme avec la probabilité $1 - p$. On peut définir une variable aléatoire

$$X = \begin{cases} 1 & \text{si la personne choisie est une femme,} \\ 0 & \text{si la personne choisie est un homme.} \end{cases}$$

La loi de cette variable est $\mathcal{B}(p)$.

- (b) Le paramètre d'intérêt est p .
- (c) Puisqu'il serait difficile d'observer toutes les personnes qui étudient à l'EPFL/UNIL, on va observer un sous-ensemble. Ce sous-ensemble doit être représentatif, par exemple on peut observer un certain nombre d'étudiant-e-s qui mangent dans une grande cafétéria pendant la pause de midi.
- (d) Un choix intuitif est le pourcentage de femmes dans le sous-ensemble observé.
- (e) Même si on connaissait la valeur de p , on ne connaîtrait pas en avance la valeur de l'estimateur. Si l'on va dans la même cafétéria deux jours différents et l'on observe le même nombre d'étudiant-e-s, ce ne seront pas exactement les mêmes étudiant-e-s, donc on n'obtiendra pas le même résultat.
- (f) On suppose que $p = 0.4$ et $n = 100$. D'après la partie (a), on peut supposer que les observations x_1, \dots, x_{100} constituent une réalisation de $X_1, \dots, X_{100} \stackrel{\text{iid}}{\sim} \mathcal{B}(p)$. L'estimateur proposé dans la partie (d) s'écrit $\hat{p}_{100} = \bar{X}_{100} = (\sum_{i=1}^{100} X_i)/100$.

$$\begin{aligned} \mathbb{E}[\hat{p}_{100}] &= \mathbb{E}\left(\frac{1}{100} \sum_{i=1}^{100} X_i\right) = \mathbb{E}[X_1] = p, \\ \text{Var}[\hat{p}_{100}] &= \text{Var}\left(\frac{1}{100} \sum_{i=1}^{100} X_i\right) = \frac{1}{100} \text{Var}[X_1] = \frac{p(1-p)}{100}, \\ \text{b}(\hat{p}_{100}) &= \mathbb{E}[\hat{p}_{100}] - p = 0. \end{aligned}$$

L'estimateur \hat{p}_n est non-biaisé. Si la taille de l'échantillon augmente, la variance diminue alors que l'espérance ne change pas. Donc, avec un plus grand échantillon, on estime le pourcentage avec une plus grande précision (on s'attend à être plus proche de la vraie valeur).

- (g) Les variables X_1, \dots, X_n sont indépendantes et identiquement distribuées, d'espérance $\mu = p$ et de variance $\sigma^2 = p(1-p)$. Nous pouvons donc utiliser le théorème central limite pour approximer la loi de

$$Z_n = \sqrt{n} \frac{\bar{X}_n - \mu}{\sigma} = \sqrt{n} \frac{\hat{p}_n - p}{\sqrt{p(1-p)}}.$$

Pour trouver n tel que $\Pr(\hat{p}_n < 0.5) \geq 0.95$ on calcule (avec $p = 0.4$)

$$\begin{aligned} & \Pr(\hat{p}_n < 0.5) \geq 0.95 \\ \Rightarrow & \Pr\left(\sqrt{n} \frac{\hat{p}_n - 0.4}{\sqrt{0.4 \times 0.6}} < \sqrt{n} \frac{0.5 - 0.4}{\sqrt{0.4 \times 0.6}}\right) \geq 0.95 \\ \Rightarrow & \Phi(0.204 \sqrt{n}) \geq 0.95 \\ \Rightarrow & \sqrt{n} \geq \frac{\Phi^{-1}(0.95)}{0.204} \\ \Rightarrow & n \geq 65.42. \end{aligned}$$

Donc on a besoin d'observer au moins 66 personnes.

Exercice 2. (a) On sait que $\int_{-\infty}^{\infty} f(x) dx = 1$. Donc

$$1 = \int_0^1 c x^{\theta-1} dx = c \left[\frac{x^\theta}{\theta} \right]_0^1 = \frac{c}{\theta},$$

et on voit bien que $c = \theta$. On a donc la densité

$$f(x) = \begin{cases} \theta x^{\theta-1} & \text{si } x \in (0, 1) \\ 0 & \text{sinon.} \end{cases}$$

(b) On a

$$\mathbb{E}[X_1] = \int_{-\infty}^{\infty} x f(x) dx = \int_0^1 x \theta x^{\theta-1} dx = \theta \int_0^1 x^\theta dx = \theta \left[\frac{x^{\theta+1}}{\theta+1} \right]_0^1 = \frac{\theta}{\theta+1}.$$

(c) Les variables X_i sont continues, donc la fonction de vraisemblance est

$$L(\theta) = f_1(x_1; \theta) \times f_2(x_2; \theta) \times \dots \times f_n(x_n; \theta),$$

où $f_i(x_i; \theta) = f(x_i; \theta)$ est la densité pour chaque X_i . On trouve

$$L(\theta) = \theta x_1^{\theta-1} \theta x_2^{\theta-1} \dots \theta x_n^{\theta-1} = \theta^n \left(\prod_{i=1}^n x_i \right)^{\theta-1}.$$

Donc

$$\ell(\theta) = \log(L(\theta)) = n \log(\theta) + (\theta - 1) \sum_{i=1}^n \log(x_i).$$

Pour trouver la valeur de θ qui maximise $\ell(\theta)$ on résout

$$\begin{aligned} & \ell'(\theta) = 0 \\ \Leftrightarrow & \frac{n}{\theta} + \sum_{i=1}^n \log(x_i) = 0 \\ \Leftrightarrow & \frac{1}{\theta} = -\frac{1}{n} \sum_{i=1}^n \log(x_i) \\ \Leftrightarrow & \theta = -\frac{n}{\sum_{i=1}^n \log(x_i)}. \end{aligned}$$

Il s'agit bien d'un maximum puisque

$$\ell''(\theta) = -\frac{n}{\theta^2} < 0,$$

pour tout $\theta > 0$. Donc la valeur $\theta = -n/(\sum_{i=1}^n \log(x_i))$ maximise la fonction $L(\theta)$ et $\hat{\theta}_{ML} = -n/(\sum_{i=1}^n \log(X_i))$ est l'estimateur du maximum de vraisemblance. Remarquons que puisque $x_i \in (0, 1)$, on a $\log(x_i) < 0$ et par conséquent $-n/(\sum_{i=1}^n \log(x_i)) > 0$.

(d) Pour la méthode des moments on doit résoudre l'équation

$$\bar{X}_n = \frac{\hat{\theta}}{\hat{\theta} + 1}.$$

Ceci donne $\hat{\theta}_{MOM} = \bar{X}_n/(1 - \bar{X}_n)$.

Exercice 3. (a) Les variables X_i sont discrètes, donc la fonction de vraisemblance est

$$L(p) = f_1(x_1; p) \times f_2(x_2; p) \times \dots \times f_n(x_n; p),$$

où $f_i(x_i; p) = P(X_i = x_i) = p^{x_i}(1-p)^{1-x_i}$ est la fonction de fréquences pour chaque X_i . On trouve

$$L(p) = p^{x_1}(1-p)^{1-x_1} p^{x_2}(1-p)^{1-x_2} \dots p^{x_n}(1-p)^{1-x_n} = p^{\sum_{i=1}^n x_i} (1-p)^{n - \sum_{i=1}^n x_i}.$$

(b) Comme $\mathbb{E}[X_i] = p$, l'estimateur des moments se trouve en résolvant l'équation $p = \bar{X}_n$. Cette équation est déjà résolue, donc $\hat{p}_{MOM} = \bar{X}_n$.

(c) L'estimateur du maximum de vraisemblance est la valeur de p qui maximise $L(p)$, ou, de manière équivalente, la valeur qui maximise la fonction $\ell(p) = \log(L(p))$.

On a

$$\ell(p) = \sum_{i=1}^n x_i \log(p) + \left(n - \sum_{i=1}^n x_i \right) \log(1-p).$$

Pour trouver le maximum on résout

$$\begin{aligned} \ell'(p) &= 0 \\ \Leftrightarrow \frac{\sum_{i=1}^n x_i}{p} - \frac{n - \sum_{i=1}^n x_i}{1-p} &= 0 \\ \Leftrightarrow (1-p) \sum_{i=1}^n x_i - p \left(n - \sum_{i=1}^n x_i \right) &= 0 \\ \Leftrightarrow \sum_{i=1}^n x_i &= pn \\ \Leftrightarrow p &= \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}_n. \end{aligned}$$

Il s'agit bien d'un maximum, étant donné que

$$\ell''(p) = -\frac{\sum_{i=1}^n x_i}{p^2} - \frac{n - \sum_{i=1}^n x_i}{(1-p)^2} < 0,$$

quel que soit $p \in (0, 1)$. Donc la valeur $p = \bar{x}_n$ maximise la fonction $L(p)$ et \bar{X}_n est l'estimateur du maximum de vraisemblance, $\hat{p}_{ML} = \bar{X}_n$.

(d) On a $\hat{p}_{MOM} = \hat{p}_{ML} = \bar{X}_n$. Donc

$$\mathbb{E}[\hat{p}_{MOM}] = \mathbb{E}[\hat{p}_{ML}] = \mathbb{E}[\bar{X}_n] = \mathbb{E}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} \mathbb{E}\left(\sum_{i=1}^n X_i\right) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_i] = p,$$

parce que les variables X_i sont toutes $\mathcal{B}(p)$. Donc les estimateurs sont non-biaisés. Pour la variance on a

$$\begin{aligned} \text{Var}[\hat{p}_{MOM}] &= \text{Var}[\hat{p}_{ML}] = \text{Var}[\bar{X}_n] = \\ &= \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \text{Var}\left(\sum_{i=1}^n X_i\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}[X_i] = \frac{p(1-p)}{n}, \end{aligned}$$

parce que les variables X_i sont indépendantes et toutes $\mathcal{B}(p)$.

Exercice 4. (a) $\mathbb{E}(Y) = \int_{-\infty}^{\infty} y f_Y(y) dy = 1/\lambda$, et $\text{Var}(Y) = \int_{-\infty}^{\infty} y^2 f_Y(y) dy - (\mathbb{E}(Y))^2 = 1/\lambda^2$.

(b) On sait que $\mathbb{E}(\bar{Y}_n) = \mathbb{E}(Y) = 1/\lambda$ et $\text{Var}(\bar{Y}_n) = \text{Var}(Y)/n$. Par le théorème central limite,

$$\Pr\left(\frac{\bar{Y}_n - \mathbb{E}(\bar{Y}_n)}{\sqrt{\text{Var}(\bar{Y}_n)}} \leq z\right) = \Pr\left(\frac{\sqrt{n}\bar{Y}_n - \mathbb{E}(Y)}{\sqrt{\text{Var}(Y)}} \leq z\right) = \Pr\left(\sqrt{n}\frac{\bar{Y}_n - 1/\lambda}{1/\lambda} \leq z\right) \rightarrow \Pr(Z \leq z),$$

où $Z \sim \mathcal{N}(0, 1)$.

(c) La vraisemblance

$$L(\lambda) = \prod_{i=1}^n f_Y(y_i; \lambda) = \prod_{i=1}^n \lambda \exp(-\lambda y_i) = \lambda^n \exp\left(-\sum_{i=1}^n \lambda y_i\right).$$

La log-vraisemblance est

$$l(\lambda) = n \log(\lambda) - \sum_{i=1}^n y_i \lambda$$

On a

$$\frac{dl}{d\lambda}(\lambda) = \frac{n}{\lambda} - \sum_{i=1}^n y_i$$

et

$$\frac{d^2l}{d\lambda^2}(\lambda) = -\frac{n}{\lambda^2} < 0$$

pour tout $\lambda > 0$. L'estimateur du maximum de vraisemblance satisfait donc

$$\frac{dl}{d\lambda}(\hat{\lambda}) = 0 \Leftrightarrow \hat{\lambda} = 1/\bar{Y}_n.$$

C'est bien un maximum car la deuxième dérivée est négative.

(d) On applique la méthode Delta avec $g(x) = 1/x$. On a

$$\bar{Y}_n \stackrel{\text{app.}}{\approx} \mathcal{N}\left(\frac{1}{\lambda}, \frac{1}{\lambda^2} \frac{1}{n}\right) = \mathcal{N}(\mu, \sigma^2/n),$$

avec $\mu = 1/\lambda$ et $\sigma^2 = 1/\lambda^2$. Ainsi $g(\mu) = 1/\mu = \lambda$ et $g'(\mu) = -1/\mu^2 = -\lambda^2$ de sorte que

$$\hat{\lambda}_{ML} = g(\bar{Y}_n) \stackrel{\text{app.}}{\approx} \mathcal{N}(g(\mu), g'(\mu)^2 \sigma^2/n) = \mathcal{N}(\lambda, \lambda^2/n).$$