

TD 4 - Statistiques - corrections

Exercice 1 - Elections

(a) Définissons X_i une variable aléatoire suivant une loi de Bernoulli, valant 1 si la personne i vote Erika, 0 sinon. $X_i \sim \text{Bern}(0.5)$; $E(X_i) = p = 0,5$; $V(X_i) = p(1-p) = 0,25$
Ensuite, $N_{\text{Erika}} = X_1 + \dots + X_{400}$, donc $E(N_{\text{Erika}}) = 400 * E(X_i) = 200$ et $V(N_{\text{Erika}}) = 400 * V(X_i) = 100$

Le théorème central limite :

- s'applique pour une somme de n variables indépendantes **et de même loi X_i** (c'est le cas)
- alors cette somme suit une loi normale $N(n * E(X_i), n * V(X_i))$

Donc $N_{\text{Erika}} \sim N(200, 100)$. On centre et on réduit pour se ramener à une loi $Z \sim N(0,1)$:
 $p(N_{\text{Erika}} > 210) = p((N_{\text{Erika}} - 200)/10 > (210 - 200)/10) = p(Z > 1) \approx 0,16$

(b) Définissons Y_i une variable aléatoire suivant une loi de Bernoulli, valant 1 si la personne i vote pour Pierre, Jeanne, ou Jerry, et 0 sinon.

$$Y_i \sim \text{Bern}(0.3) ; E(Y_i) = p = 0,3 ; V(Y_i) = p(1-p) = 0,3 * 0,7 = 0,21$$

Les Y_i sont indépendantes et de même loi, donc par le théorème central limite, la somme $S = Y_1 + \dots + Y_{400}$ suit une loi normale $N(400 * 0.3, 400 * 0.21) = N(120, 84)$.

$$\text{Donc } p(S < 100) = p\left(\frac{S - 120}{\sqrt{84}} < \frac{100 - 120}{\sqrt{84}}\right) = p\left(Z < \frac{-20}{\sqrt{21 \times 4}}\right) = p\left(Z < \frac{-10 \times 2}{\sqrt{21 \times 2}}\right) = p\left(Z < \frac{-10}{\sqrt{21}}\right) \approx 0,0145$$

Exercice 2 - Ajustement de courbe de données

(a) Vraisemblance = $p(y_i | a, b, x_i, \sigma)$, et $y_i - ax_i - b \sim N(0, \sigma^2)$ donc $y_i \sim N(ax_i + b, \sigma^2)$, on écrit sa densité de probabilité:

$$p(y_i | a, b, x_i, \sigma) = f(y_i) = \frac{1}{\sigma \sqrt{2\pi}} \times e^{-\frac{(y_i - ax_i - b)^2}{2\sigma^2}}$$

(b) On a les données $x_i = [1, 3, 5]$ et $y_i = [8, 2, 1]$. Les données étant indépendantes,

$$p(y_1, y_2, y_3 | a, b, x_1, x_2, x_3, \sigma) = \prod_{i=1}^3 p(y_i | a, b, x_i, \sigma)$$

$$\text{En développant, on obtient } \left(\frac{1}{\sigma \sqrt{2\pi}}\right)^3 \times e^{-\frac{[(8-a-b)^2 + (2-3a-b)^2 + (1-5a-b)^2]}{2\sigma^2}}.$$

On a donc $\ln(p(8, 2, 1 | a, b, 1, 3, 5, \sigma))$

$$-3 \ln(\sigma) - \frac{3}{2} \ln(2\pi) - \frac{[(8-a-b)^2 + (2-3a-b)^2 + (1-5a-b)^2]}{2\sigma^2}$$

$$-3 \ln(\sigma) - \frac{3}{2} \ln(2\pi) - \frac{[35a^2 + 3b^2 - 38a - 22b + 18ab + 69]}{2\sigma^2}$$

Dans le cas général, on écrirait:

$$p(y_1, \dots, y_n | a, b, x_1, \dots, x_n, \sigma) = \prod_{i=1}^n p(y_i | a, b, x_i, \sigma)$$

$$\left(\frac{1}{\sigma \sqrt{2\pi}} \right)^n \times e^{-\sum_{i=1}^n (y_i - ax_i - b)^2 / (2\sigma^2)} = -n \ln(\sigma) - \frac{n}{2} \ln(2\pi) - \frac{\sum_{i=1}^n (y_i - ax_i - b)^2}{2\sigma^2}$$

(c) Chercher les dérivées partielles de la vraisemblance par rapport à a et b. Chercher ensuite quand elles s'annulent, et en déduire la meilleure estimation de a et b.

$$\frac{\partial \ln(p(8, 2, 1 | a, b, 1, 3, 5, \sigma))}{\partial a} = \frac{70a + 18b - 38}{-2\sigma^2} = 0 \Leftrightarrow 70a + 18b = 38$$

$$\frac{\partial \ln(p(8, 2, 1 | a, b, 1, 3, 5, \sigma))}{\partial b} = \frac{6b + 18a - 22}{-2\sigma^2} = 0 \Leftrightarrow 18a + 6b = 22$$

$$\begin{cases} 70a + 18b = 38 \\ 18a + 6b = 22 \end{cases} \Leftrightarrow \begin{cases} 70a + 18b = 38 \\ 54a + 18b = 66 \end{cases} \Leftrightarrow \begin{cases} 70a + 18b = 38 \\ a = -7/4 \end{cases}$$

$$18 \times \frac{-7}{4} + 6b = 22 \Rightarrow b = \frac{1}{6} \times \left(\frac{88 + 7 \times 18}{4} \right) = \frac{88 + 70 + 56}{24} = \frac{214}{24} = \frac{107}{12}$$

Ainsi, la droite qui passe le mieux par les données a pour équation $y = \frac{-7}{4}x + \frac{107}{12}$.

Exercice 3 - Maximum de vraisemblance

(a) La densité de probabilité d'une loi uniforme sur [a,b] est

$$f(x_i | a, b) = \frac{1}{b-a} \text{ si } x \in [a, b], 0 \text{ sinon.}$$

Les données sont indépendantes, donc la vraisemblance L vaut

$$L = \prod_{i=1}^5 p(x_i | a, b) = \frac{1}{(b-a)^5} \text{ si } (\forall i \text{ alors } x_i \in [a, b]), \text{ sinon } 0$$

On veut trouver a et b tels que L soit maximum. Donc on veut déjà que toutes les données soient dans [a,b] : on a une première condition $a \leq \min(x_i) \leq \max(x_i) \leq b$.

Ensuite, il faut que (b-a) soit le plus petit possible.

On prendra donc $a = \min(x_i) = 1.2$ et $b = \max(x_i) = 10.5$

Ainsi, la distribution uniforme qui maximise la vraisemblance de ces données est la distribution uniforme sur l'intervalle [1.2, 10.5].

Note: il n'y a aucune raison que ce soit bien celle de laquelle on a tiré ces données. C'est juste celle pour laquelle notre observation de données est la plus probable.

(b) Dans le cas général, on prendra toujours $a = \min(\text{observations})$ et $b = \max(\text{observations})$.