

TD6 : ESTIMATION DANS UN PROBLÈME DE SONDAGE

**Exercice 1.** (Problème de sondages) Soit  $N$  le nombre d'habitants d'une commune. Il s'agit de faire un sondage de popularité de deux candidats (candidat  $A$  et candidat  $B$ ) qui se présentent aux élections municipales. On choisit un échantillon de  $n$  habitants auxquels on pose la question : "Pour qui voteriez-vous aux élections ?" A l'issue de ce sondage, on obtient les données  $X_1, \dots, X_n$ , où

$$X_i = \begin{cases} 1, & \text{si le } i^{\text{ème}} \text{ habitant questionné préfère le candidat } A, \\ 0, & \text{si le } i^{\text{ème}} \text{ habitant questionné préfère le candidat } B, \end{cases} \quad i = 1, \dots, n.$$

Pour les raisons évidentes, il est impossible de questionner tous les habitants. Donc  $n < N$  (dans la pratique, on a toujours  $n \ll N$ ). Notons  $\theta$  la part d'habitants de la commune qui préfèrent le candidat  $A$ . Le but du sondage est d'estimer  $\theta$  et de donner un intervalle confiance pour  $\theta$ .

Définissons les valeurs déterministes  $x_1, \dots, x_N$  par

$$x_j = \begin{cases} 1, & \text{si le } j^{\text{ème}} \text{ habitant préfère le candidat } A, \\ 0, & \text{si le } j^{\text{ème}} \text{ habitant préfère le candidat } B, \end{cases} \quad j = 1, \dots, N.$$

On a alors

$$\theta = \frac{1}{N} \sum_{j=1}^N x_j.$$

Définissons aussi

$$\sigma^2 = \frac{1}{N} \sum_{j=1}^N (x_j - \theta)^2.$$

On appelle  $\theta$  *moyenne de population* et  $\sigma^2$  *variance de population*.

1. Proposer un modèle statistique pour ce problème. Pour cela,

- (a) Calculer  $\mathbf{P}_\theta((X_1, X_2, X_3, X_4) = (0, 1, 1, 0))$  en utilisant le fait qu'il s'agit d'un tirage au hasard sans remise d'une population de taille  $N$ , car chaque habitant peut apparaître au maximum une fois dans l'échantillon.
- (b) Soit  $(a_1, \dots, a_n) \in \{0, 1\}^n$  tel que  $a_1 + \dots + a_n = a$ . Deviner la forme de

$$\mathbf{P}_\theta((X_1, \dots, X_n) = (a_1, \dots, a_n))$$

et prouver la formule obtenue par récurrence sur  $n$ .

- 2. Montrer que  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$  est un estimateur sans biais de  $\theta$ .
- 3. Montrer que

$$\mathbf{Cov}(X_i, X_j) = -\frac{\sigma^2}{N-1} \quad \text{pour } i \neq j \quad \text{et} \quad \mathbf{Var}(\bar{X}_n) = \frac{\sigma^2}{n} \left(1 - \frac{n-1}{N-1}\right).$$

4. Calculer  $\mathbf{E}[s_n^2]$ , où  $s_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2$ , et proposer un estimateur sans biais  $\hat{v}_n^2$  de la variance  $\mathbf{Var}(\bar{X}_n)$ .

5. On se place maintenant dans le cadre asymptotique où  $N \rightarrow \infty$ ,  $n = n(N) \rightarrow \infty$  et  $n/N \rightarrow 0$ .

(a) Montrer que  $\bar{X}_n$  et  $\hat{v}_n^2$  sont des estimateurs consistants de  $\theta$  et  $\sigma^2$ .

(b) On admet que  $\sqrt{n}(\bar{X} - \theta)$  converge en loi vers  $\mathcal{N}(0, \sigma^2)$ . Démontrer la normalité asymptotique

$$\sqrt{n} \frac{\bar{X}_n - \theta}{\hat{v}_n} \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{N}(0, 1).$$

(c) En utilisant le résultat ci-dessus, trouver un l'intervalle aléatoire  $[A, B]$  tel que

$$\lim_{n \rightarrow \infty} \mathbf{P}([A, B] \text{ contient } \theta) = 95\%.$$

6. *Application numérique* : donner l'intervalle de confiance de niveau asymptotique 95% pour  $\theta$  lorsque  $N = 8000$ ,  $n = 100$ ,  $n_1 = \sum_{i=1}^n \mathbb{1}\{X_i = 1\} = 65$ .

**Exercice 2. (facultatif)** Soient  $f$  et  $g$  deux densités de probabilité définies sur  $[0, 1]$ . On suppose que

$$f(x) \leq M g(x), \quad \forall x \in [0, 1]. \quad (1)$$

Soient  $\{(X_i, U_i) : i \in \mathbb{N}^*\}$  une suite de variables aléatoires iid telles que

- la loi de  $X_1$  a pour densité  $g$ ,
- $U_1$  est de loi uniforme sur  $[0, 1]$ ,
- $X_1$  et  $U_1$  sont indépendantes.

On définit les variables aléatoires

$$N = \min \{n : f(X_n) \geq M U_n g(X_n)\}, \quad Y = X_N.$$

1. Soit  $q = \mathbf{P}(f(X_n) \geq M U_n g(X_n))$ . Montrer que  $q = 1/M$ .

2. Montrer que  $N$  suit la loi géométrique de paramètre  $q$ , c'est-à-dire

$$\mathbf{P}(N = k) = (1 - q)^{k-1} q, \quad k \in \mathbb{N}^*.$$

3. Pour toute fonction mesurable bornée  $h : [0, 1] \rightarrow \mathbb{R}$ , montrer que

$$\mathbf{E}[h(Y)] = \int_0^1 h(y) f(y) dy. \quad (2)$$

En déduire que  $Y$  a pour densité  $f$ .

**Remarque** Cet exercice montre que si (1) est vrai et si on est capable de générer des variables aléatoires de densité  $g$ , alors on sera également capable de générer des variables aléatoires de densité  $f$ . Cette méthode porte le nom de méthode d'acceptation-rejet. Elle a toutefois un inconvénient : le nombre moyen de tirage de  $X_i$  nécessaire pour calculer  $Y$  est  $\mathbf{E}[N] = M$ . Ce nombre peut être excessivement grand lorsque  $M$  est trop grand.