

10 | Probabilit  s & Estimation

I – Suites de variables al  atoires discr  tes finies

1 – Ind  pendance d'une famille de variables al  atoires

D  finition 10.1 – Soient $n \geq 2$ un entier et X_1, X_2, \dots, X_n des variables al  atoires d  finies sur Ω . On dit que les variables al  atoires X_1, X_2, \dots, X_n sont **mutuellement ind  pendantes** lorsque pour tout $(x_1, x_2, \dots, x_n) \in X_1(\Omega) \times X_2(\Omega) \times \dots \times X_n(\Omega)$,

$$P([X_1 = x_1] \cap [X_2 = x_2] \cap \dots \cap [X_n = x_n]) = P(X_1 = x_1) \times P(X_2 = x_2) \times \dots \times P(X_n = x_n).$$

D  finition 10.2 – Soit $(X_n)_{n \geq 1}$ une suite de variables al  atoires d  finies sur Ω . On dit que la suite $(X_n)_{n \geq 1}$ est une suite de variables al  atoires mutuellement ind  pendantes si et seulement si pour tout $m \in \mathbb{N}^*$, les variables al  atoires X_1, X_2, \dots, X_m sont mutuellement ind  pendantes.

Exemple 10.3 – On lance une pi  ce   quilibr  e jusqu'   obtenir PILE. Pour tout $n \in \mathbb{N}^*$, on note X_n la variable al  atoire   gale    1 si on obtient PILE au n -i  me lancer et 0 sinon. Alors la suite $(X_n)_{n \geq 1}$ est une suite de variables al  atoires mutuellement ind  pendantes.

2 – Esp  rance et variance d'une famille de variables al  atoires

Th  or  me 10.4

Soient $n \geq 2$ un entier et X_1, X_2, \dots, X_n des variables al  atoires d  finies sur Ω . Alors

$$E\left(\sum_{k=1}^n X_k\right) = \sum_{k=1}^n E(X_k).$$

Autrement dit,

$$E(X_1 + X_2 + \dots + X_n) = E(X_1) + E(X_2) + \dots + E(X_n).$$

Exemple 10.5 – On consid  re une suite $(X_k)_{k \in \mathbb{N}^*}$ de variables al  atoires mutuellement ind  pendantes, dont la loi (commune) est donn  e pour tout $k \in \mathbb{N}^*$ par

$$P(X_k = 1) = \frac{1}{3} \quad \text{et} \quad P(X_k = 2) = \frac{2}{3}.$$

Pour tout $n \in \mathbb{N}^*$, on pose $S_n = \sum_{k=1}^n X_k$. Calculer $E(S_n)$.

Je commence par calculer $E(X_k)$ pour tout $k \in \mathbb{N}^*$: $E(X_k) = 1 \times \frac{1}{3} + 2 \times \frac{2}{3} = \frac{5}{3}$.

Et donc

$$E(S_n) = E\left(\sum_{k=1}^n X_k\right) = \sum_{k=1}^n E(X_k) = \sum_{k=1}^n \frac{5}{3} = \frac{5n}{3}.$$

Théorème 10.6

Soient $n \geq 2$ un entier et X_1, X_2, \dots, X_n des variables aléatoires **mutuellement indépendantes** définies sur Ω . Alors

$$V\left(\sum_{k=1}^n X_k\right) = \sum_{k=1}^n V(X_k).$$

Autrement dit,

$$V(X_1 + X_2 + \dots + X_n) = V(X_1) + V(X_2) + \dots + V(X_n).$$

Exemple 10.7 – On reprend l'exemple précédent. Calculer $V(S_n)$.

Je commence par calculer $V(X_k)$. Pour cela, il me faut calculer $E(X_k^2)$: $E(X_k^2) = 1^2 \times \frac{1}{3} + 2^2 \times \frac{2}{3} = \frac{9}{3} = 3$.
Donc d'après la formule de König-Huygens,

$$V(X_k) = E(X_k^2) - E(X_k)^2 = 3 - \left(\frac{5}{3}\right)^2 = \frac{2}{9}.$$

Ainsi, comme les variables X_1, \dots, X_n sont mutuellement indépendantes,

$$V\left(\sum_{k=1}^n X_k\right) = \sum_{k=1}^n V(X_k) = \sum_{k=1}^n \frac{2}{9} = \frac{2n}{9}.$$

II – Inégalités classiques en théorie des probabilités

1 – Inégalité de Markov

Proposition 10.8 – Inégalité de Markov

Soit X une variable aléatoire **positive** (discrète ou à densité) admettant une espérance.
Alors pour tout réel a strictement positif,

$$P(X \geq a) \leq \frac{E(X)}{a}.$$

Remarque 10.9 – L'inégalité stricte est aussi vérifiée :

$$P(X > a) \leq \frac{E(X)}{a}.$$

Corollaire 10.10

Soit X une variable aléatoire (discrète ou à densité). On suppose que X^2 admet une espérance.
Alors pour tout réel a strictement positif,

$$P(|X| \geq a) \leq \frac{E(X^2)}{a^2}.$$

2 – Inégalité de Bienaymé-Tchebychev

Proposition 10.11

Soit X une variable aléatoire (discrète ou à densité). On suppose que X^2 admet une espérance. Alors pour tout réel ε strictement positif,

$$P(|X - E(X)| \geq \varepsilon) \leq \frac{V(X)}{\varepsilon^2}.$$

Remarque 10.12 – Souvent, on reconnaît qu'il faut se servir de l'inégalité de Bienaymé-Tchebychev grâce aux valeurs absolues présentes dans la probabilité.

3 – Loi faible des grands nombres

Théorème 10.13 – Loi faible des grands nombres

Soit $(X_n)_{n \in \mathbb{N}^*}$ une suite de variables aléatoires mutuellement indépendantes, ayant chacune la même espérance m et la même variance σ^2 . On pose $\bar{X}_n = \frac{1}{n} \sum_{k=1}^n X_k$. Alors pour tout réel ε strictement positif,

$$\lim_{n \rightarrow +\infty} P(|\bar{X}_n - m| \geq \varepsilon) = 0.$$

III – Estimation

Les statisticiens connaissent, en général, le type de loi qui décrit tel ou tel phénomène, par l'observation. Mais souvent ils ne connaissent pas tous les paramètres de la dite loi. Ils doivent donc les estimer : c'est l'objectif de ce que l'on appelle la statistique inférentielle.

On considère une variable aléatoire X , dont le type de loi est connu et dépend d'un paramètre réel θ inconnu (ce peut être le paramètre λ d'une variable exponentielle, l'étendue $b - a$ d'une variable uniforme sur $[a, b]$, le paramètre p d'une variable de Bernoulli, l'espérance m d'une loi normale, etc).

L'objectif est de donner une *estimation* de la vraie valeur du paramètre θ .

Il y a deux types d'estimation : l'estimation ponctuelle et l'estimation par intervalle de confiance.

1 – Échantillons et estimateurs

Exemple 10.14 – On suppose que la durée de vie (en heures) des ampoules produites par l'entreprise Lumilux suit une loi exponentielle X dont le paramètre λ est inconnu.

1. Rappeler la formule donnant l'espérance de X en fonction de λ .

Comme X suit une loi exponentielle, je sais que $E(X) = \frac{1}{\lambda}$.

2. Des tests ont été effectués sur 10 ampoules. On obtient les durées de vie suivantes :

62.1	75.4	73.1	81	68.7	73.6	64.2	78.5	74.4	63
------	------	------	----	------	------	------	------	------	----

À l'aide de la série statistique ci-dessus, donner une estimation du paramètre λ .

Je sais que l'espérance correspond à la *moyenne* des valeurs prises par X .

Or la moyenne de la série statistique ci-dessus vaut

$$\frac{62.1 + 75.4 + \cdots + 63}{10} \approx 71.4.$$

Ainsi je peux estimer que $\frac{1}{\lambda} \approx 71.4$. Et donc $\lambda \approx \frac{1}{71.4} \approx 0.014$.

3. Peut-on affirmer que λ est égal à la valeur précédente? Si non, comment pourrait-on tenter d'améliorer la qualité de l'estimation?

Je ne peux pas affirmer que λ est égal à la valeur précédente car il s'agit ici uniquement d'un échantillon de valeurs prises par X . Une autre série de 10 mesures donnerait une autre estimation, *a priori* proche, du paramètre λ . Pour tenter d'améliorer la qualité de l'estimation, il faudrait obtenir un échantillon de taille plus importante.

Définition 10.15 – Soient X une variable aléatoire (discrète ou à densité) et $n \in \mathbb{N}^*$ un entier. On appelle **n -échantillon** de X tout n -uplet (X_1, \dots, X_n) de variables aléatoires indépendantes et de même loi que X .

Définition 10.16 – Soient X une variable aléatoire (discrète ou à densité) et $n \in \mathbb{N}^*$ un entier. On appelle **réalisation de l'échantillon** (x_1, \dots, x_n) (ou **échantillon observé**) tout n -uplet (x_1, \dots, x_n) de valeurs prises par (X_1, \dots, X_n) (x_1 est la valeur prise par X_1 , x_2 est la valeur prise par X_2 , ..., x_n est la valeur prise par X_n).

Exemple 10.17 – Dans l'exemple précédent, l'échantillon observé est

(62.1, 75.4, 73.1, 81, 68.7, 73.6, 64.2, 78.5, 74.4, 63).

Bien évidemment, l'échantillon observé est *aléatoire*. On aurait par exemple pu, en effectuant les tests avec un autre jeu de 10 ampoules, obtenir

74.1	82.5	68.5	70.3	84	77.2	69.6	73.8	76.3	68.7
------	------	------	------	----	------	------	------	------	------

Définition 10.18 – Soient θ un réel, X une variable aléatoire dont la loi dépend d'un paramètre θ , $n \in \mathbb{N}^*$ un entier et (X_1, \dots, X_n) un échantillon de X .

On appelle **estimateur** de θ toute variable aléatoire T_n de la forme $\varphi(X_1, \dots, X_n)$.

Exemple 10.19 – Dans l'exemple précédent, le paramètre que l'on cherche à estimer est le paramètre λ de la loi exponentielle suivie par X . L'estimateur considéré pour l'espérance est

$$T_n = \frac{1}{n} \sum_{k=1}^n X_k = \frac{1}{n} (X_1 + X_2 + \dots + X_n).$$

Cet estimateur m'a permis d'estimer **ponctuellement**.

2 – Biais et risque quadratique d'un estimateur

Définition 10.20 – Soient θ un réel, X une variable aléatoire dont la loi dépend d'un paramètre θ et T_n un estimateur de θ .

- Si T_n admet une espérance, on appelle **biais** de T_n le réel $b(T_n)$ défini par

$$b(T_n) = E(T_n) - \theta.$$

- On dit que T_n est un **estimateur sans biais** de θ si et seulement si le biais de T_n est nul, *i.e.* si et seulement si $E(T_n) = \theta$. Dans le cas contraire, on dit que T_n est un estimateur **biaisé** de θ .

Exemple 10.21 – Soient X une variable aléatoire de loi de Bernoulli de paramètre p et $(X_n)_{n \in \mathbb{N}^*}$ une suite de variables aléatoires indépendantes qui suivent toutes la même loi que X .

On pose pour tout $n \in \mathbb{N}^*$, $\bar{X}_n = \frac{1}{n} \sum_{k=1}^n X_k$. Montrer que \bar{X}_n est un estimateur sans biais de p .

Je calcule l'espérance de \bar{X}_n :

$$E(\bar{X}_n) = E\left(\frac{1}{n} \sum_{k=1}^n X_k\right) = \frac{1}{n} \sum_{k=1}^n E(X_k) = \frac{1}{n} \sum_{k=1}^n p = \frac{1}{n} \times np = p.$$

Comme $E(\bar{X}_n) = p$, j'en déduis que \bar{X}_n est bien un estimateur sans biais de p .

Définition 10.22 – Soient θ un réel, X une variable aléatoire dont la loi dépend d'un paramètre θ et T_n un estimateur de θ .

Si T_n admet une variance, on appelle **risque quadratique** de T_n le réel $r(T_n)$ défini par

$$r(T_n) = E((T_n - \theta)^2).$$

Proposition 10.23

Soient θ un réel, X une variable aléatoire dont la loi dépend d'un paramètre θ et T_n un estimateur de θ qui admet une variance. Alors

$$r(T_n) = b(T_n)^2 + V(T_n).$$

En particulier, si T_n est un estimateur sans biais de θ , alors

$$r(T_n) = V(T_n).$$

3 – Estimation par intervalle de confiance

Les estimations ponctuelles ne fournissent pas d'information sur la précision des estimations, c'est-à-dire qu'elles ne tiennent pas compte de l'erreur possible attribuable aux fluctuations d'échantillonnage. Or deux échantillons distincts donnent presque certainement des valeurs distinctes pour l'estimation. Ici, il s'agit toujours d'estimer un paramètre inconnu, mais au lieu de lui attribuer une valeur unique en faisant appel à un estimateur ponctuel, on construit un intervalle aléatoire qui permet de "recouvrir" avec une certaine fiabilité, la vraie valeur du paramètre estimé.

Définition 10.24 – Soient θ un réel, X une variable aléatoire dont la loi dépend d'un paramètre θ , U_n et V_n deux estimateurs de θ et $\alpha \in [0, 1]$ un réel.

On dit que $[U_n, V_n]$ est un **intervalle de confiance** de θ au niveau de confiance $1 - \alpha$ si

$$P(U_n \leq \theta \leq V_n) \geq 1 - \alpha.$$