

Sur la variance empirique

J'ai peut-être dit une **grosse bêtise** lors du dernier TD. Ce document est là pour, j'espère, clarifier les choses.

Définition 1. Soit X une variable aléatoire, n un entier naturel non nul et X_1, X_2, \dots, X_n, n variables aléatoires indépendantes ayant la même loi que X .

La moyenne empirique de l'échantillon (X_1, X_2, \dots, X_n) est le réel

$$\bar{X}_n := \frac{1}{n} \sum_{i=1}^n X_i.$$

La variance empirique biaisée de l'échantillon (X_1, X_2, \dots, X_n) est le réel positif

$$\hat{s}_n^2 := \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2.$$

Si $n \geq 2$, la variance empirique non biaisée de cet échantillon est le réel positif

$$\hat{\sigma}_n^2 := \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2.$$

Pour passer de l'une des variances empiriques à l'autre, il suffit bien sûr d'écrire

$$\hat{\sigma}_n^2 = \frac{n}{n-1} \hat{s}_n^2$$

et

$$\hat{s}_n^2 = \frac{n-1}{n} \hat{\sigma}_n^2.$$

Proposition 2. Avec les mêmes notations que dans la définition précédente, on a

$$\hat{s}_n^2 = \left(\frac{1}{n} \sum_{i=1}^n X_i^2 \right) - \bar{X}_n^2. \tag{1}$$

Démonstration. On a, pour tout i entre 1 et n ,

$$(X_i - \bar{X}_n)^2 = X_i^2 - 2X_i\bar{X}_n + \bar{X}_n^2.$$

En sommant pour tous les i entre 1 et n , on obtient

$$\sum_{i=1}^n (X_i - \bar{X}_n)^2 = \sum_{i=1}^n (X_i^2 - 2X_i\bar{X}_n + \bar{X}_n^2) = \left(\sum_{i=1}^n X_i^2 \right) - 2\bar{X}_n \left(\sum_{i=1}^n X_i \right) + n\bar{X}_n^2.$$

On remarque que

$$\sum_{i=1}^n X_i = n\bar{X}_n$$

donc on a

$$\sum_{i=1}^n (X_i - \bar{X}_n)^2 = \left(\sum_{i=1}^n X_i^2 \right) - 2n\bar{X}_n^2 + n\bar{X}_n^2 = \left(\sum_{i=1}^n X_i^2 \right) - n\bar{X}_n^2.$$

En en déduit le résultat, en divisant par n . □

Pour la version non biaisée, malheureusement, on a seulement

$$\hat{\sigma}_n^2 = \frac{1}{n-1} \left(\sum_{i=1}^n X_i^2 \right) - \frac{n}{n-1} \bar{X}_n^2.$$

La **grosse bêtise** (j'espère que je n'ai pas écrit ça au tableau...) est qu'on n'a surtout pas

$$\text{FAUX, FAUX, FAUX! } \hat{s}_n^2 = \frac{1}{n} \left[\left(\sum_{i=1}^n X_i^2 \right) - \bar{X}_n^2 \right] \text{ FAUX, FAUX, FAUX!}$$

Par contre, on a

$$\hat{s}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i^2 - \bar{X}_n^2),$$

et aussi

$$\hat{\sigma}_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i^2 - \bar{X}_n^2),$$

parce que

$$\sum_{i=1}^n \bar{X}_n^2 = n \bar{X}_n^2.$$

Tant qu'on y est, calculons les biais. On suppose maintenant que

$$\mathbb{E} [X^2] < \infty.$$

On va d'abord calculer $\mathbb{E} [\bar{X}_n^2]$. On commence par faire le développement suivant :

$$\left(\sum_{i=1}^n X_i \right)^2 = \sum_{i=1}^n \sum_{j=1}^n X_i X_j.$$

Il y a n^2 termes dans cette somme et parmi ces n^2 termes, il y en a n qui sont tels que $i = j$. On a la décomposition suivante

$$\left(\sum_{i=1}^n X_i \right)^2 = \sum_{i=1}^n X_i^2 + \sum_{\substack{j=1 \\ j \neq i}}^n X_i X_j. \tag{2}$$

Avant de passer à l'espérance, on souligne le fait très important que, pour tout i entre 1 et n ,

$$\mathbb{E} [X_i^2] = \mathbb{E} [X^2],$$

car les X_i ont la même loi que X . Et pour tous i et j entre 1 et n avec $i \neq j$, par indépendance

$$\mathbb{E} [X_i X_j] = \mathbb{E} [X_i] \mathbb{E} [X_j] = \mathbb{E} [X] \mathbb{E} [X] = \mathbb{E} [X]^2.$$

Maintenant, on applique l'espérance au développement (2). Par linéarité de l'espérance et le fait que les X_i ont tous la même loi que X et sont indépendants,

$$\mathbb{E} \left[\left(\sum_{i=1}^n X_i \right)^2 \right] = n \mathbb{E} [X^2] + (n^2 - n) \mathbb{E} [X]^2.$$

On en déduit, en divisant par n^2 ,

$$\mathbb{E} [\bar{X}_n^2] = \frac{1}{n} \mathbb{E} [X^2] + \mathbb{E} [X]^2 - \frac{1}{n} \mathbb{E} [X]^2 = \mathbb{E} [X]^2 + \frac{1}{n} \text{Var} [X].$$

Finalement, en utilisant la formule (1),

$$\mathbb{E}[\hat{s}_n^2] = \frac{1}{n} \times n\mathbb{E}[X^2] - \mathbb{E}[X]^2 - \frac{1}{n}\text{Var}[X] = \text{Var}[X] - \frac{1}{n}\text{Var}[X].$$

Donc le biais de \hat{s}_n^2 est

$$b(\hat{s}_n^2) = \mathbb{E}[sn^2] - \text{Var}[X] = -\frac{1}{n}\text{Var}[X].$$

Ce biais est strictement négatif (sauf si X est presque sûrement égale à une constante). Donc l'estimateur \hat{s}_n^2 de la variance est biaisé.

En revanche,

$$\mathbb{E}[\hat{\sigma}_n^2] = \mathbb{E}\left[\frac{n}{n-1}\hat{s}_n^2\right] = \frac{n}{n-1}\text{Var}[X] - \frac{1}{n-1}\text{Var}[X] = \text{Var}[X],$$

donc $\hat{\sigma}_n^2$ est non biaisé ($b(\hat{\sigma}_n^2) = 0$).