

Simplifions la statistique

Par Louis Blais, stat., ASSQ

MESURES DE TENDANCE CENTRALE

Les valeurs des unités d'échantillon observées ont souvent tendance à se concentrer autour d'une même valeur. Ce sont des mesures de tendance centrale. Les mesures les plus fréquemment utilisées sont la moyenne arithmétique, la médiane et le mode d'un échantillon. La moyenne géométrique et la moyenne quadratique sont aussi parfois utilisées. Il existe également une moyenne dite harmonique, disponible dans Excel 2013. La moyenne arithmétique est la plus fréquemment utilisée. Elle se trouve aussi bien dans les rapports statistiques que dans les logiciels de statistiques. Plusieurs tests statistiques seraient impossibles à réaliser avec les autres mesures, ce qui ne veut pas dire qu'elles sont dénuées d'intérêt, au contraire. Toutes ces mesures — y compris les mesures de dispersion¹ — servent à caractériser l'ensemble des valeurs d'une variable.

1. Voir la fiche *Simplifions la statistique – Mesures de dispersion* (Blais, 2016)

LA MOYENNE ARITHMÉTIQUE

C'est la mesure la plus fréquemment utilisée. Elle se trouve aussi bien dans les rapports statistiques que dans les logiciels de statistiques.

Plusieurs décisions sont prises en fonction de cette mesure qui se calcule comme suit :

$$\bar{x} = \sum x_i / n$$

Pour des données qui sont mesurées, mais notées par classe (par exemple, le recouvrement), on suggère d'utiliser la moyenne pour calculer la moyenne.

Facile à calculer, cette mesure a comme principal inconvénient d'être sensible aux données extrêmes, surtout si le nombre d'unités d'échantillon est petit. Elle est peu utilisée si les données extrêmes sont très nombreuses.

LA MÉDIANE

La médiane est la valeur du rang qui se situe au-dessus d'elle, l'autre correspond au rang précédent. Si la distribution est symétrique par rapport à la moyenne, la médiane et la moyenne, la médiane et la

moyenne auront des valeurs pratiquement égales. Dans le cas contraire, ces valeurs seront différentes, parfois même très différentes. L'avantage de la médiane par rapport à la moyenne est que sa valeur est complètement indépendante des valeurs extrêmes qui auront été observées. La médiane peut être différente de la moyenne lorsque le nombre d'observations est faible. En effet, dans ce cas, une ou deux observations extrêmes peuvent grandement influencer la moyenne.

Un excellent exemple de la différence entre la moyenne et la médiane est le revenu annuel des Québécois. En juin 2013, le revenu annuel moyen avant impôt était de 51 000 \$ alors que la médiane était de 42 400 \$ (Statistique Canada). Cette différence vient du fait que les salaires très élevés font augmenter la moyenne. Si on veut connaître le revenu total de tous les Québécois, la moyenne est la meilleure mesure. Par contre, pour obtenir le portrait réel du revenu des Québécois, la médiane est plus appropriée.

Dans Excel 2013, il existe deux fonctions nommées « Censile » et « Centile ». La fonction d'inclusion ou d'exclusion concerne l'inclusion de la valeur du 50^e centile ou son exclusion (exclure). Si le nombre d'observations est très grand, l'inclure ou l'exclure n'aura pas d'influence s'il y a peu d'observations (généralement moins de 30).

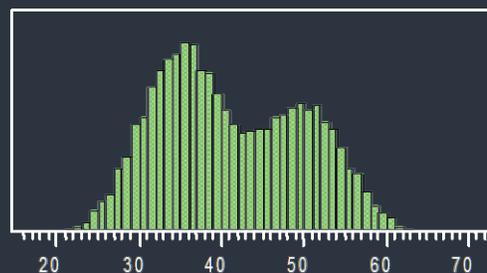
LE MODE D'UN ÉCHANTILLON

Pour des variables discrètes, dont les valeurs sont un code (par exemple, le drainage) ou une classe (par exemple, le recouvrement), le mode d'un échantillon est la valeur la plus fréquemment observée. Il peut ne pas être unique. Par exemple, lorsque le drainage a été noté et que les drainages « rapide » et « modéré » sont les plus souvent observés et reviennent un même nombre de fois, on aura alors

deux modes pour le drainage. Ce n'est pas une situation privilégiée fréquemment, mais, pour une variable discrète, elle permet d'avoir rapidement un aperçu de la valeur « centrale », c'est-à-dire celle qu'on trouve le plus fréquemment dans un échantillon.

Pour les variables continues, le mode d'un échantillon n'est pas très intéressant. Avec ce type de variable, il est difficile de déterminer la région du ou des pics de la distribution de la variable. Un histogramme de la distribution de la variable permet de dire si on a affaire à une distribution unimodale (p. ex. : en forme de cloche). On dit alors que la distribution est unimodale. L'histogramme peut également nous dire aussi s'il y a plusieurs pics dans la distribution (p. ex. : en forme de cloche). C'est une indication que la variable peut avoir été mesurée sur plusieurs populations ou selon plusieurs stratégies différentes. Une attention particulière doit être apportée si aucune autre variable notée ne permet de décrire les groupes que l'histogramme ne peut pas distinguer. Ici, au contraire, il y aurait deux groupes, par exemple, deux strates). Certains logiciels donnent toujours le mode, mais celui-ci ne constitue pas une indication valable.

Le problème, c'est qu'une même mesure soit un peu plus fréquente pour qu'elle devienne le mode, peu importe si elle est près du centre de la distribution. On ne peut rien dire dans ce cas. Pour des variables continues, le mode n'a pas de valeur.



LA MOYENNE GÉOMÉTRIQUE

La moyenne géométrique est la racine n^e du produit de n valeurs prises par une variable.

Elle se calcule comme suit :

$$\bar{x}_g = \sqrt[n]{\prod_{i=1}^n x_i} = \sqrt[n]{x_1 * x_2 * \dots * x_n}$$

S'il y a deux données, on calcule la racine carrée de leur produit. S'il y en a dix, on calcule la racine 10^e du produit des dix mesures. La moyenne géométrique est moins sensible aux données extrêmes et peut représenter une bonne es

Elle s'emploie occasionnellement lorsqu'on a une distribution symétrique, par exemple dans le cas des dénombrements où il y a une très grande variabilité ou dans celui des mesures environnementales sur les polluants atmosphériques. Pour pouvoir utiliser cette moyenne, les valeurs doivent être positives et non nulles. Très peu de logiciels calculent la moyenne géométrique. En effet, la racine n^e d'un produit est en général fastidieux à calculer, même pour un ordinateur puissant. Excel ne le fait pas mais on peut utiliser la fonction « Moyenne.geometrique »).

LA MOYENNE QUADRATIQUE

La moyenne quadratique est définie comme la racine carrée de la moyenne arithmétique des carrés des valeurs des observations. On peut aussi dire :

$$\bar{x}_Q = \sqrt{\left(\frac{\sum x_i^2}{n}\right)}$$

TABEAU 1 DHP et surface terrière de 10 arbres dans une parcelle à rayonnement égal

Arbres	DHP (mm)	ST (m ²)	DHP (mm ²)
1	110	0,009503	12 100
2	250	0,049087	62 500
3	350	0,096211	122 500
4	140	0,015394	19 600
5	150	0,017671	22 500
6	200	0,031416	40 000
7	180	0,025447	32 400
8	400	0,125664	160 000
9	300	0,070686	90 000
10	120	0,011310	14 400
Sommes	2 200	0,452389	576 000
Moyenne	220		240

Ce type de moyenne est utilisée pour les variables reliées à des unités de surface. En foresterie, puisque la surface terrière a un lien direct avec le DHP, la moyenne quadratique est la plus appropriée. En effet, si vous calculez la moyenne quadratique des diamètres et que vous vous en servez pour calculer la surface terrière, vous allez obtenir exactement le même résultat que si vous calculez la surface terrière par hectare prise individuellement. L'exemple suivant l'illustre très bien. La moyenne arithmétique des diamètres est de 220 mm. La moyenne quadratique est de 240 mm. La surface terrière est la somme de toutes les surfaces individuelles. La surface terrière totale des dix arbres donne 0,452389 m².

On calcule la surface terrière de la manière suivante :

$$ST = \sum \left(\pi (DHP_{mm}/2) \left[\frac{1^m}{1000_{mm}} \right]^2 \right)$$

$$= \sum \left(\pi (DHP_{mm}/2000) \right)^2$$

(On doit diviser le diamètre par deux, puisqu'une surface se calcule avec un rayon et non avec un diamètre, et convertir les mètres en mètres en divisant par 1 000, d'après le tableau 1.)

On obtient :

$$10\pi(240/2000)^2 = 0,452389m^2$$

On a dix arbres, ce qui donne exactement le même résultat que dans le tableau 1.)

Si on utilise la moyenne arithmétique pour calculer de la surface terrière, on obtient :

$$10\pi(220/2000)^2 = 0,38013m^2$$

Ce dernier calcul ne donne pas la bonne surface terrière.

Comme pour la moyenne géométrique, beaucoup de logiciels statistiques ne calculent pas cette moyenne. On doit la programmer.

LA MOYENNE HARMONIQUE

La moyenne harmonique est rarement utilisée même s'il est possible de la calculer dans Excel 2013. On l'obtient de la manière suivante :

$$\bar{x}_h = \frac{n}{\frac{1}{x_1} + \frac{1}{x_2} + \dots + \frac{1}{x_n}}$$

Voici un exemple d'application de la moyenne harmonique : si vous roulez sur une distance de 10 km à une vitesse de 40 km/h, puis que vous parcourez les 10 km suivants à 60 km/h, votre vitesse moyenne pour les 20 km sera de 48 km/h. La moyenne harmonique permet de calculer précisément la vitesse moyenne sur cette distance. Il ne faut pas confondre avec la vitesse moyenne basée sur le temps. Si les kilomètres parcourus étaient remplacés par un temps parcouru (10 minutes plutôt que 10 km), votre vitesse moyenne serait alors de 50 km/h pendant 20 minutes, soit la moyenne arithmétique des vitesses.

CONCLUSION

En conclusion, chaque type de moyenne correspond à un besoin précis. Si la moyenne géométrique est utilisée pour calculer une quantité totale, par exemple, un volume total, la moyenne arithmétique est le meilleur choix. Elle est également plus appropriée lors d'un test de comparaison entre deux moyennes. Lorsqu'il y a peu d'unités d'échantillon, ou que la distribution n'est pas symétrique autour de la moyenne, la médiane permet alors de fournir un meilleur portrait des données que la moyenne arithmétique. Pour déterminer une surface terrière, le diamètre quadratique du DHP est la place est à réserver. Pour les autres types de moyennes, l'intérêt est de savoir qu'elles existent même si elles sont rarement utilisées. Les logiciels ne les fournissent pas par défaut.

POUR EN SAVOIR PLUS...

Ministère des Forêts, de la Faune et des Parcs
Direction de l'aménagement et de l'environnement forestiers
5700, 4e Avenue Ouest
Québec (Québec) G1H 6R1
daef@mffp.gouv.qc.ca