

# Simplifions la statistique

Par Louis Blais, stat., ASSQ

## MESURES DE DISPERSION

Les mesures de tendance centrale sont très importantes en statistique, employées seules, elles ne permettent pas de bien décrire les données. Les mesures qui les complètent bien sont les mesures de dispersion. Une mesure de dispersion est une mesure qui quantifie l'étalement des observations entre leur minimum à la valeur maximum.

En effet, deux échantillons différents peuvent avoir la même moyenne arithmétique des dispersions très différentes. Ou encore, l'étalement des observations est tellement vaste que la moyenne seule ne veut pratiquement rien dire. Il existe plusieurs mesures de dispersion. Nous allons nous attarder aux principales, soit l'étendue, la variance, l'écart type et le coefficient de variation. Nous verrons aussi les centiles. Les centiles ne constituent pas une mesure proprement dite, mais un ensemble de valeurs qui décrit la dispersion des données observées.

Pour illustrer ces concepts, nous allons utiliser les statistiques des hauteurs de plants d'épines noirs produits en pépinière, prêts à être mis en terre et dont la qualité doit être vérifiée. Pour l'exercice, on a mesuré la hauteur de 360 plants en centimètres. Le tableau 1 présente les statistiques de l'inventaire.

**TABLEAU 1** Statistiques des hauteurs des plants produits en pépinière

Hauteur moyenne :	14,3 cm
Variance :	11,33 cm <sup>2</sup>
Écart type :	3,37 cm
Minimum :	8 cm
Maximum :	24 cm

## L'ÉTENDUE

L'étendue, ou *range* en anglais, est la différence entre la valeur maximale et la valeur minimale des observations. Cette mesure était très utilisée avant l'avènement des ordinateurs, car elle est très facile à calculer. Cependant, c'est la moins recommandée de garantir les valeurs extrêmes, puisque le calcul ne prend pas compte des autres observations. Elle peut donc donner une idée erronée de la dispersion des observations. Dans l'exemple présenté en introduction, l'étendue est égale à 16 cm (24-8).

## LA VARIANCE

La variance est la moyenne arithmétique des écarts, mis au carré, entre les observations et leur moyenne arithmétique. Elle se calcule ainsi :

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}$$

La moyenne de la population étant pas connue, on divise par « n-1 » plutôt que par « n » pour éviter un biais dans le calcul de la variance de l'échantillon.

Plus les observations sont proches de la moyenne, moins elles sont dispersées, plus la variance sera petite. Plus les observations sont dispersées, plus la variance sera grande. À noter que les unités de la variance sont toujours des « unités » au carré. Par exemple, si les unités des observations sont en millimètres, on obtient une variance en millimètres carrés. La variance est très utilisée dans les tests statistiques, par exemple dans les tests statistiques à intervalles de confiance.

déterminer la taille d'un échantillon. Dans l'exemple précédent, la variance est de 11,33 cm<sup>2</sup>. Les unités étant au carré, il est important de savoir si cette variance est élevée ou non.

La variance sert au calcul de l'écart type, de l'erreur type et du coefficient de variation. Elle est également utilisée pour calculer l'intervalle de confiance et à déterminer le nombre d'unités d'échantillon nécessaires selon le niveau de précision et de confiance dans les tests statistiques, notamment lorsque l'on compare deux moyennes de deux échantillons.

Dans Excel 2013, quatre fonctions calculent la variance. Ce sont « VAR.P.N », « VAR.P », « VAR.S » et « VAR ». Les fonctions « VAR.P.N » et « VAR.P » sont identiques et calculent la variance de la population en supposant que la moyenne de la population est connue. Dans les faits, le diviseur de la fonction est « n » plutôt que « n-1 » et on ne connaît jamais, ou très rarement, la moyenne de la population. Les fonctions « VAR.S » ou « VAR », qui sont équivalentes.

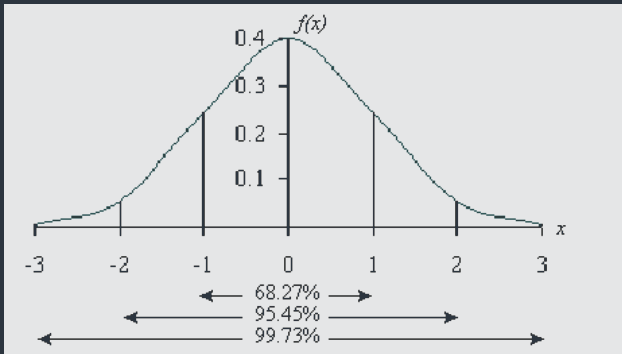
## L'ÉCART TYPE

Mathématiquement, l'écart type n'est rien d'autre que la racine carrée de la variance. Il se calcule ainsi :

$$s = \sqrt{s^2} = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{(n - 1)}}$$

C'est la moyenne quadratique des écarts par rapport à la moyenne arithmétique des observations par « n-1 » plutôt que par « n ». L'écart type sert à exprimer l'écart moyen des observations par rapport à la moyenne arithmétique. L'avantage ici est que ses unités (p.ex. : cm) sont les mêmes que celles des observations.

Pour désigner l'écart type, certains utilisent à tort le terme « déviation standard » qui est une traduction littérale de l'anglais « *standard deviation* ». Il ne faut pas non plus le confondre avec « l'erreur type ».



Si les observations ont une loi normale, 68,27 % des observations sont à  $\pm 1$  écart type, 95,45 %, à  $\pm 2$  écarts types et 99,73 %, à  $\pm 3$  écarts types. Cette connaissance permet de bien évaluer la dispersion des observations (voir ci-dessus).

Dans l'exemple des plants produits en pépinière, l'écart type est de 3,37 cm. Cela signifie qu'en moyenne chaque plant s'écarte de 3,37 cm de la moyenne, mais cet énoncé est incomplet. En fait, 68,27 % des observations ont un écart de  $\pm 3,37$  cm, 95,45 % de  $\pm 6,74$  cm ( $2 \times 3,37$  cm) et 99,73 % de  $\pm 10,1$  cm ( $3 \times 3,37$  cm).

Dans Excel 2013, plusieurs fonctions calculent l'écart type. Ce sont « ECARTYPE », « ECARTYPE.STANDARD », « ECARTYPE.PEARSON », « ECARTYPE.PEARSON » et « ECART.MOYEN ». Les fonctions « ECARTYPE » et « ECARTYPE.STANDARD » sont identiques et calculent l'écart type présenté dans ce chapitre. « ECARTYPE.PEARSON » et « ECARTYPE.PEARSON » sont identiques et calculent l'écart type d'une population que l'on connaît la vraie moyenne. Comme pour la variance, elle divise par « n » au lieu de « n-1 », mais puisque l'on connaît rarement la vraie moyenne, il ne faut pas l'utiliser.

« ECART.MOYEN » calcule l'écart moyen par rapport à la moyenne en valeur absolue. Elle peut être intéressante, bien qu'elle est généralement l'écart type qui peut servir à faire d'autres calculs.

## L'ERREUR TYPE

**L'erreur type est l'écart type divisé par la racine**

**carrée du nombre d'unités d'échantillonnage**, c'est-à-dire :  $\frac{s}{\sqrt{n}}$

Cette mesure a un lien direct avec la précision des sondages. Plus la taille de l'échantillon est élevée, plus l'erreur type sera petite, et plus le sondage sera précis que l'écart type, c'est-à-dire les mêmes unités que les observations. L'écart type et l'erreur type sont deux mesures. Prenons l'exemple de la pépinière qui produit des plants pour le reboisement. Les plants doivent-ils se ressembler le plus possible ou doivent-ils avoir une hauteur moyenne le plus près possible de l'objectif ? Si l'objectif est de produire des plants qui se ressemblent le plus possible, l'écart type est une excellente mesure. Par contre, si l'objectif est de produire des plants avec une moyenne de « x » cm, indépendamment de l'écart entre les plants, c'est l'erreur type qui est importante. Il faut cependant noter que, dans le deuxième cas, si l'écart type est élevé et que l'échantillon doit avoir une certaine précision, le nombre de plants à inventorier sera plus grand et les coûts de l'échantillonnage seront plus élevés.

Dans notre exemple, l'erreur type est de 0,18 cm. Cela veut dire que s'il était possible de réaliser une indépendance des observations de 360 plants de la même population, on obtiendrait alors (si on les mesurait tous) s'écarterait, en moyenne, d'environ  $\pm 0,18$  cm de la vraie valeur.

## LE COEFFICIENT DE VARIATION

Le coefficient de variation est calculé en faisant le **rapport de l'écart type sur la moyenne**, c'est-à-dire :  $CV = 100s/\bar{x}$

C'est une mesure sans unité, exprimée en pourcentage (%), mais ce n'est pas une mesure de précision. Le coefficient de variation est supérieur à 100 % si l'écart type est supérieur à la moyenne. Il est très important de l'importance de l'écart type en tenant compte de sa moyenne. Un faible coefficient de variation (exemple, inférieur à 10 %) indique une faible dispersion des données.

autour de la valeur centrale. Comme on l'a mentionné précédemment, deux échantillons peuvent avoir la même moyenne, tout en ayant des variances différentes. Aussi, deux échantillons peuvent avoir la même variance, mais une moyenne différente. Dans ce dernier cas, le coefficient de variabilité a plus grand pour l'échantillon dont la moyenne est plus petite, ce qui indique que la variabilité est plus importante. C'est ce qu'on fait lorsqu'on fait un échantillon et qu'on recherche une certaine précision (par exemple, pour les échantillons d'une population). Plus le coefficient de variabilité est élevé, plus on aura besoin d'unités d'échantillon pour obtenir cette précision.

Dans l'exemple sur la population, le coefficient de variabilité est de 23,6%. Si l'objectif est de diminuer la variabilité individuelle de la hauteur de chaque plant par rapport à la moyenne, il faudrait alors viser un coefficient de variabilité plus faible. Les façons de faire pour le diminuer : augmenter la hauteur des plants sans augmenter l'écart type ou diminuer l'écart type tout en conservant la même moyenne.

## LES CENTILES

Un centile est chacune des 99 valeurs qui divisent les données triées, de la plus petite valeur à la plus grande, en 100 parts égales, de sorte que chaque centile représente 1/100 de l'échantillon.

Un centile est le 50<sup>e</sup> centile, c'est-à-dire la médiane. Mais d'autres centiles sont importants, à savoir le 25<sup>e</sup>, le 75<sup>e</sup> et le 95<sup>e</sup>. L'ensemble de ces centiles est appelé l'écart type de la dispersion des données. Ainsi, les données situées entre le 25<sup>e</sup> et le

75<sup>e</sup> centile comprennent environ 50% des observations autour de la médiane. Entre le 5<sup>e</sup> et le 95<sup>e</sup> centile sont 90% des observations comprises entre ces valeurs.

**TABLEAU 2** Centiles et échantillons de l'épave noire

Centile	5 <sup>e</sup>	25 <sup>e</sup>	50 <sup>e</sup>	75 <sup>e</sup>	95 <sup>e</sup>
Hauteur (cm)	9	12	15	17	21

Ces données montrent que la médiane (tableau 2) et la moyenne (tableau 1) sont relativement semblables. Elles traduisent bien la symétrie des données autour de la moyenne. Notons que 50% des plants ont entre 12 et 17 cm de hauteur et que 90% d'entre eux ont entre 9 et 21 cm de hauteur.

Les centiles sont indépendants des valeurs extrêmes qui peuvent être présentes dans un échantillon. Ils risquent d'être influencés par la moyenne, même si dans certains cas ils sont robustes. Lorsque l'on prend une décision. En effet, pour déterminer des seuils, les centiles sont utiles. Les décisions seront prises, les centiles permettent de connaître l'ensemble des cas touchés par ces seuils.

D'autre part, comme il n'est pas certain que les observations suivent une loi normale, les centiles sont utiles si la distribution est asymétrique, les centiles peuvent alors servir de valeurs de référence dans la distribution.

## CONCLUSION

En conclusion, à la fin de chaque mesure de dispersion, il est important de savoir ce que cela signifie, qui peut aller d'une simple description jusqu'à la détermination d'un nombre d'unités d'échantillon.

### 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