

Statistique

Nature des données

Contexte

Soit x une v.a. évaluée sur les objets d'un échantillon E et qui prends des valeurs dans un ensemble de référence K

$$x : E \rightarrow K : e \mapsto x(e)$$

où E est un sous-ensemble de la population P

Hypothèse: on suppose que P , E , x , K , sont bien définie en fonction des objectifs d'une étude statistique.

Nature des éléments de P, E

Les objets de P sont en lien avec le contexte de l'étude

Ça peut-être n'importe quoi: nombres, pays, chiens, personnes, etc

Ce ne sont pas des données!

Les objets de E sont de même nature que ceux de P, car $E \subset P$

Rappelons que l'ensemble P peut-être définie par

- extension: on explicite chaque élément
- compréhension: via une propriété commune aux éléments

Par exemple, considérons les pays du G7 comme population:

extension: $P = \{\text{Allemagne, Canada, E-U, France, Italie, Japon, UK}\}$

compréhension: $P = \{ a \text{ est un pays et } a \text{ est dans le G7 } \}$

Population versus temps

Le contenu d'un ensemble a souvent une dépendance sur la variable "temps".
Par exemple, le G7 pourrait ne plus exister en 2020, ou encore l'Italie pourrait être remplacée par la Suisse, etc

Si l'étude porte sur une population à un moment précis, il faut le mettre dans la définition de la population

Par exemple: $P = \{ a \text{ est un pays} \mid a \text{ est dans le G7 au 2Q14} \}$

Si les objets changent dans le temps, on doit spécifier la période

Par exemple: $Q = \{ \text{symbole coté sur le TSX en 2014} \}$

Remarque: Q contient les symboles, pas les cotes, ces dernières sont des v.a. sur Q

Échantillon versus temps

L'échantillon peut aussi dépendre du temps

Par exemple:

- population de tous les joueurs de la LNH depuis le début
- échantillon formé des joueurs de l'équipe "Canadien de Mtl"

Si l'étude porte sur un moment précis, il faut le mettre dans la définition de l'échantillon. Par exemple:

{ les joueurs de l'équipe "Canadien de Mtl" au 31 décembre 2014 }

C'est l'étude statistique qui dicte le choix de la population, l'échantillon, les v.a. et la genre d'analyse à faire sur les valeurs prise pas les variables

Ensemble des valeurs d'un v.a.

Considérons une v.a. $x : E \rightarrow K$

L'ensemble K peut-être fini ou infini, contenir des mots, des symboles, des nombres, etc

Exemples: $\{ \text{oui, non} \}$, \mathbb{N} , \mathbb{R} , $\{ \text{toutes les chaînes de caractères UTF-8} \}$, $\{ \text{Harper, Trudeau} \}$, $[0,100]$, $\{ 1, 2, 3, 5, 8, 13 \}$, etc

Variable aléatoire dynamique

Il arrive que la valeur d'une v.a. sur un élément de l'échantillon change au cours du temps. On dit qu'elle est "dynamique", par opposition à "statique"

Une suite de valeurs d'une v.a. dynamique obtenues dans le temps s'appelle une "série temporelle"

Exemple: considérons la population:

$$Q = \{ \text{symboles coté sur le TSX durant l'année 2018} \}$$

et la v.a

cote: $Q \rightarrow \mathbb{R}^+$, qui associe à un symbole de Q sa cote sur le TSX

C'est une série temporelle. En effet, si on considère le symbole IMG de Q , on peut voir les valeurs de la série pour une période récente via [cote\(IMG\)](#)

Valeurs d'un v.a. dynamique

Considérons une v.a. dynamique $\mathbf{x} : \mathbf{E} \rightarrow \mathbf{K} : \mathbf{e} \mapsto \mathbf{x}(\mathbf{e})$

Cette écriture de \mathbf{x} est incomplète car elle ne fait pas apparaître le temps \mathbf{t}

Si $\mathbf{T} = [\mathbf{t}_a , \mathbf{t}_b]$ est l'intervalle de temps qui nous intéresse, on écrit plutôt:

$$\mathbf{x} : \mathbf{E} \times \mathbf{T} \rightarrow \mathbf{K} : \mathbf{e}, \mathbf{t} \mapsto \mathbf{x}(\mathbf{e}, \mathbf{t})$$

Cette écriture fait bien voir que pour chaque élément \mathbf{e} de l'échantillon, on obtient une fonction de \mathbf{t} . C'est-à-dire une courbe sur l'intervalle \mathbf{T} .

En pratique, on dispose des données en des moments précis (pas en continu). L'écart de temps entre 2 données peut être constant ou variable (lié à des événements).

Exemple de v.a. dynamique

Reprenons l'exemple de la population:

$$Q = \{ \text{symboles coté sur le TSX durant l'année 2018} \}$$

et la v.a (en explicitant sa dépendance sur le temps):

cote: $Q \times T \rightarrow \mathbb{R}^+$, qui associe à un symbole de Q sa cote sur le TSX

C'est une série temporelle pour chacun des éléments de la population.

Le graphe suivant montre les séries temporelles pour un échantillon de 3 éléments, avec les mois de juillet et août 2018 comme intervalle de temps et une donnée par jour d'ouverture de la bourse de Toronto:

$$E = \{ \text{IMG, NA, R} \}$$

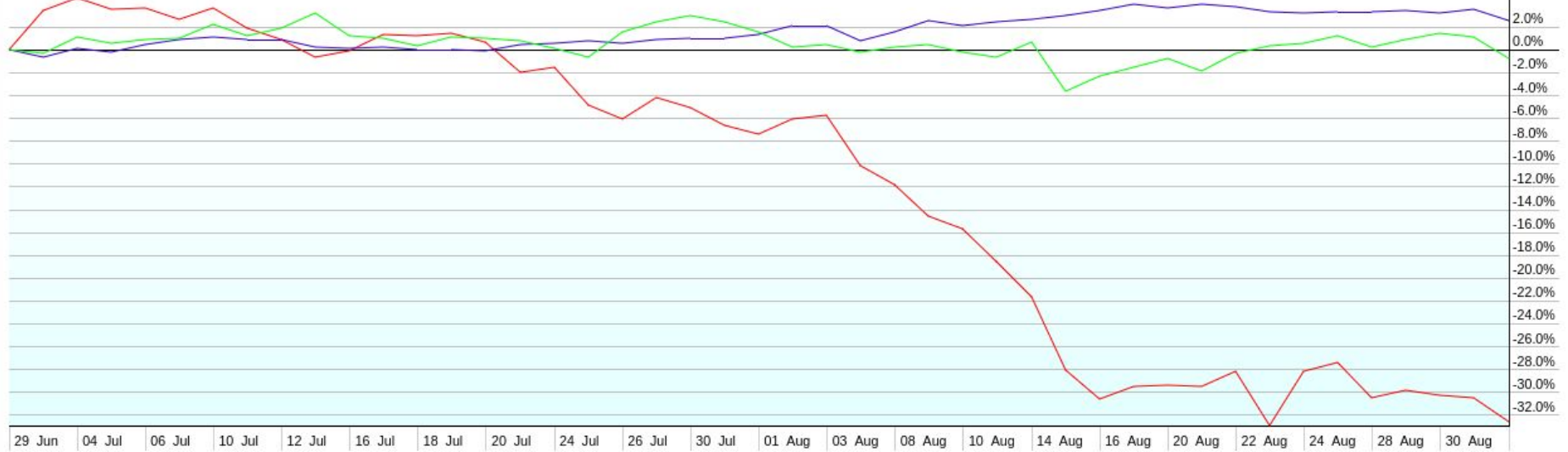
Note: Les cotes ayant des échelles différentes, les variables sont d'abord transformées via un facteur d'échelle et mis en pourcentage. Le facteur est la première valeur de la série. Par exemple, pour IMG, on affiche:

$$100 * \text{cote}(\text{IMG}, t) / \text{cote}(\text{IMG}, t_0)$$

IAMGold Natl Bank of Canada Suncor Energy

29 June 2018 - 04 September 2018

Copyright, StockCharts.com



45 days

Pratique courante

Dans un rapport d'étude, on laisse généralement au lecteur le soin de déduire la population, l'échantillon, les v.a. et les ensembles des valeurs à partir de la mise en contexte

C'est souvent une source d'erreur et de mauvaise interprétation des résultats

Il convient d'être rigoureux et fournir un contexte complet et précis