

CKM rappel

1

La statistique est la science qui étudie et analyse les Données

Statistique descriptive
décrire les données déjà existées



Statistique probabiliste
décrire les prévisions

rappel :

population $\Sigma p, \Sigma O \dots$ Base de l'étude stati	individu élément de la population	échantillon sous-ensemble de la pop	caractère propriété de la pop
-------------------------------------------------------------------	--------------------------------------	----------------------------------------	----------------------------------

- effectif
- fréquence
- effectif cumulé
- effectives cumulé

indicateurs de forme

Statistique d'échantillon

direct
ponctuel
ne change pas

qualitatif

continue

change
ou a une évaluation

quantitatif

nominale ordinal

	échantillon	pop
Moyenne	\bar{X}	μ
Variance	S^2	σ^2
Ecart type	s	σ
co-variance	S_{xy}	σ_{xy}
corrélation	r_{xy}	ρ_{xy}

Statistique de pop ou plus paramètres de la population

échantillon Mesures de tendances centrales et de dispersion **population**

$$\bar{X} = \frac{\sum x_i n_i}{N}$$

$$s = \sqrt{S^2}$$

$$\mu = \frac{\sum x_i n_i}{N}$$

$$\sigma = \sqrt{\sigma^2}$$

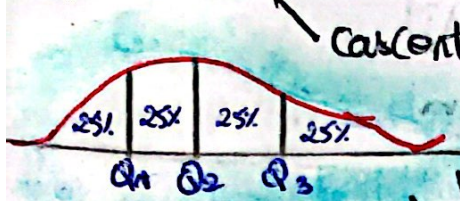
Base le nombre $\rightarrow n-1$ corrigé

$$S^2 = \frac{\sum (x_i - \bar{X})^2}{n-1}$$

$$\sigma^2 = \frac{\sum (x_i - \mu)^2}{N}$$

- Percentile : p pour cent des observations ont une valeur inférieure ou égale à p

- quartile $\left\{ \begin{array}{l} \text{cas discret } Q_1 = \frac{1}{4} N \quad Q_2 = \frac{1}{2} N \quad Q_3 = \frac{3}{4} N \\ \text{cas continue } Q_1 = \frac{N}{4} \Rightarrow \text{class } [a-b[\end{array} \right.$



$$Q_1 = a + \frac{(b-a) \left(\frac{N}{4} - n_i \right)}{(n_{i+1} - n_i)}$$



Remarques :

- La 50^{ème} percentile est la médiane
- intervalle interquartile = $Q_3 - Q_1 = IQ$
 \approx étendue interquartile

• autres mesures de dispersion

étendue = Valeur maximale - Valeur minimale

coefficient de variation = $\frac{\text{Ecart type}}{\text{moyenne}} \times 100$

↳ combien s'élève l'écart type par rapport au moyen

• Astuces

percentile = $\frac{85}{100} \times N \Rightarrow$ décimale on prend la suivante

• Mode

qui a le plus grand effectif (cas discret)

cas continue

amplitudes égales $[a, b[$
 classe modale

amplitudes inégales

$$M_o = a + \frac{(n_i - n_{i-1})}{(n_i - n_{i-1}) + (n_i - n_{i+1})} \times a$$

Classe modale →

intervalle	n_i
$[x_{i-2}, x_{i-1}[$	n_{i-1}
$[a, b[$	n_i
$[x_{i+1}, x_{i+2}[$	n_{i+1}

Valeur singulière ! What! comment on va les détecter!

Degré d'asymétrie

mesure la forme d'une distribution de données.

(-) données biaisées à gauche.

(+) données biaisées à droite



$$y_1 = \frac{\mu_3}{\sigma^3}$$

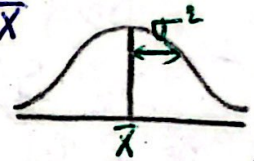
distribution symétrique = 0

$$\mu_3 = \frac{1}{N} \sum n_i (x_i - \bar{x})^3$$

Variation centrée réduite Z_i

$$Z_i = \frac{X_i - \bar{x}}{S}$$

mesure la distance en nombre d'écart type entre x_i et \bar{x}



théorème de Chebyshev

$$[\bar{x} - zS, \bar{x} + zS]$$

pour connaître le pourcentage de données dans un intervalle

$$1 - \frac{1}{z^2} \Rightarrow \text{pourcentage}$$

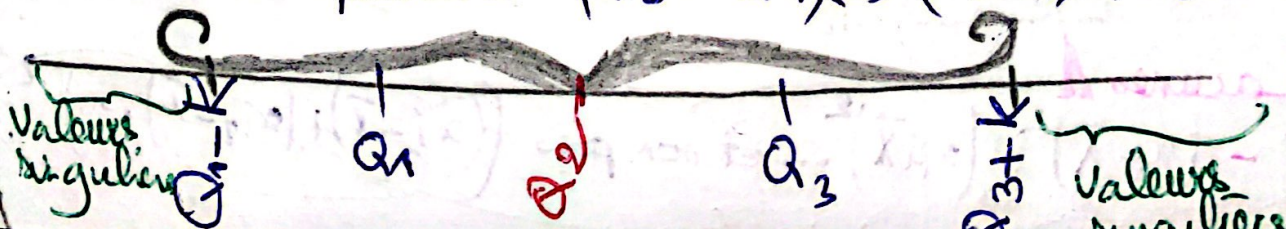
Règle empirique

1. 68% des observations se situent dans $[\bar{x} - S, \bar{x} + S]$
2. 95% des observations se situent dans $[\bar{x} - 2S, \bar{x} + 2S]$
3. presque toutes les observations se situent dans $[\bar{x} - 3S, \bar{x} + 3S]$

Analyse exploratoire des données

- ① Résumé en 5 chiffres
 - valeur la plus petite
 - Q_1 • Q_2 • Q_3
 - valeur la plus grande

- ② constitution de la Boîte à Moustaches
- étendu inéquartils : $(Q_3 - Q_1) \Rightarrow (EIQ) \times 1,5$



Statistiques bivariées mesure la relation entre 2 variables
échantillon covariance population

$$S_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

$$\sigma_{xy} = \frac{\sum (x_i - \mu)(y_i - \mu_y)}{N}$$

on se retrouve dans 2 cas :

- valeurs positives → relation linéaire positive
- valeurs négatives → relation linéaire négative

mesure relation linéaire

coefficient de corrélation

$$r = \frac{S_{xy}}{S_x S_y}$$

$$\rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \cdot \sigma_y}$$

mesure relation linéaire

on se retrouve dans 3 cas :

- valeurs proches de +1 → forte relation linéaire positive
- valeurs proches de -1 → forte relation linéaire négative
- valeurs proches de 0 → absence de relation linéaire

$$-1 < r < 1$$

méthode des moindres carrés partie du CH3

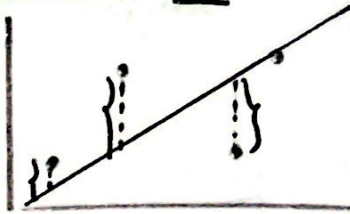
$$\hat{y} = b_0 + b_1 x$$

$$b_0 = \bar{y} - b_1 \bar{x}$$

$$b_1 = \frac{\text{Covariance}(x, y)}{\text{Variance}(x)}$$

$$b_1 = \frac{\frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1}}{\frac{\sum (x_i - \bar{x})^2}{n-1}}$$

pourquoi on utilise la méthode des moindres carrés



J'ai ces points là !
 et je veux

$$b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

une droite qui va mieux passer avec ces points → pour minimiser les écarts !

Lacunes

$(- (x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 \dots$ et non pas $((x_1 - \bar{x}) + (x_2 - \bar{x}) \dots)^2$

Statistique appliquée # Révision 55

Ch 2 Échantillonnage

9

but : avoir un échantillon suffisamment représentatif de la pop

I. Échantillonnage aléatoire

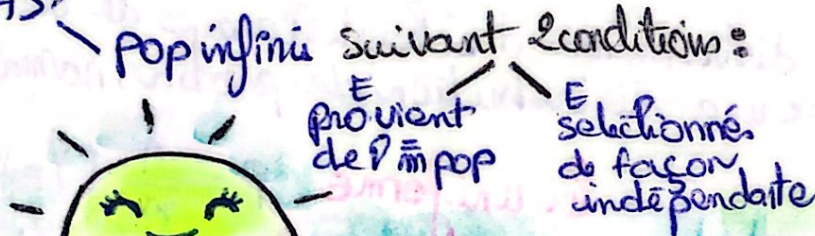
- Échantillon aléatoire simple (EAS) $f = p$
- Échantillon aléatoire stratifié
- Échantillonnage par grappes
- Échantillonnage systématique

II. Échantillonnage N. aléa

- n de commodo
- n subjectif
- n à participation volontaire
- enquête Boule de neige
- n par Quota



EAS / pop finis (+)



avantages

- échantillon facilement constitués
- données facilement collectées

inconvénients

- pas d'extrapolation / la pop
- pas d'analyse probabiliste ou de l'inférence
- y'a pas un critère standard de choix

utilisés dans les études qualitatives pas d'extrapolation POP

I. Échantillonnage aléatoire : EAS Estimation pour

1) Distributions d'échantillonnage de \bar{x} (Nombre d'échantillon par

$E(\bar{x}) = \mu$

Ecart type / pop fini $\sqrt{\bar{x}} = \sqrt{\frac{N-n}{N-1}} \left(\frac{\sigma}{\sqrt{n}} \right)$ * sinon

pop infini $\sqrt{\bar{x}} = \frac{\sigma}{\sqrt{n}}$

facteur de correction $\sqrt{\bar{x}} = \frac{\sigma}{\sqrt{n}}$

$\frac{n}{N} < 0,05$

- La moyenne de toutes les moyennes = la moyenne de tous les échantillons
- chaque loi de proba se caractérise par $Z(E(\bar{x}), \sigma(\bar{x}))$
- une pop infinie \Rightarrow on connaît pas le fin

Echantillon aléatoire simple

taille: n (Récapitulatif de P/E)

pop finie

distribution de \bar{x}

pop infinie

$$\sigma_{\bar{x}} = \sqrt{\frac{N-n}{N-1}} \left(\frac{\sigma}{\sqrt{n}} \right)$$

erreur type de la moyenne

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

→ si la taille de EAS $\leq 0,05$ de la taille de pop

$$n \leq 5\% N$$

Astuce:

$\bar{x}, s, p \rightarrow$ sont des estimateurs ponctuelles de μ, σ, p

théorème centrale limite

1 - sélectionner EAS

2 - when $n > 30$, la distribution d'échantillonnage de \bar{x} peut être approchée par une distribution de proba, normale.

• $n < 30$ individus 

loi uniforme $P = \frac{1}{n} (X=X_i)$

• $n > 30$ individus 

loi normale on va le voir

e) Distribution d'échantillonnage de \bar{p} → proportion de pop possédant un caractère donné

• $\bar{p} = \frac{x}{n}$ (nombre d'élément)

• $E(\bar{p}) = p$

Khtaha obent 3an

Ecart type:
 → pop finie
 → pop infinie

$$\sigma_{\bar{p}} = \sqrt{\frac{N-n}{N-1}} \cdot \sqrt{\frac{p(1-p)}{n}}$$

$$\sigma_{\bar{p}} = \sqrt{\frac{p(1-p)}{n}}$$

Exception → si la $n \leq 5\% N \rightarrow$ pop infinie

→ Aussi dans ce cas la loi binomiale peut être approchée par une loi normale si :

• $E(X) = np > 5$
 • $\sigma = \sqrt{np(1-p)} > 5$

question d'examen 2017

$$P(X=k) = C_n^k p^k (1-p)^{n-k}$$

Astuce:

$(n \rightarrow$ l'écart type \rightarrow probabilité \rightarrow moyenne conventionnelle)

Estimation par intervalle de confiance

de μ (1)

de p (2)

1. La moyenne de la pop

Pors que le seuil de confiance augmente, ME \downarrow et les intervalles s'élargissent

$1 - \alpha$: indice de confiance

lorsque σ est connu, l'estimation par intervalle est donnée par:

$$\left[\bar{x} - Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} ; \bar{x} + Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right]$$

\Rightarrow intervalle de confiance

où :

- α : coeff de risque

- $Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$: marge d'erreur

- longueur : 2x marge d'erreur

Loi normale

Astuce 3

$(1 - \alpha) \downarrow \rightarrow$ intervalle \Rightarrow bonne moyenne

lorsque σ est inconnu, l'estimation par intervalle est donnée par:

$$\left[\bar{x} - t_{\alpha/2} \frac{s}{\sqrt{n}} ; \bar{x} + t_{\alpha/2} \frac{s}{\sqrt{n}} \right]$$

TEST STUDENT

à Retenir

Moyenne μ

ponctuelle
 $\hat{\mu} = \bar{X}$

Vous trouverez plus de détails dans la page précédente

normal
 σ connu

$$\left[\bar{X} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right]$$

$X \sim N(\mu; \frac{\sigma}{\sqrt{n}})$

intervalle de conf
student

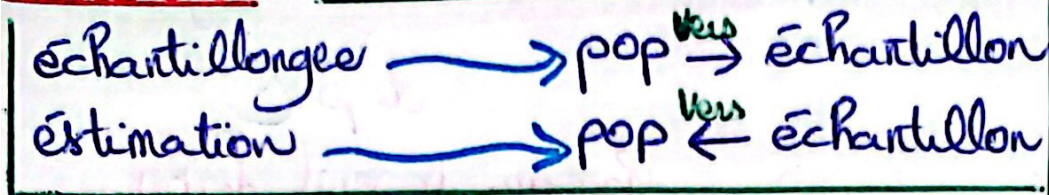
σ non connu

$$\left[\bar{X} \pm t_{\alpha/2} \frac{s}{\sqrt{n}} \right]$$

$n \geq 30$
loi normale

$n < 30$
loi de Student

à Retenir



Regression linéaire simple

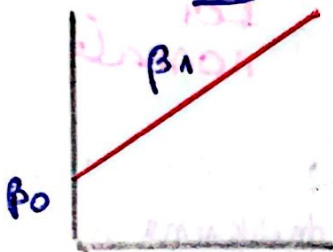
CH3

on est dans la sphère d'étudier une pop par rapport à 2 caractères!
 \rightarrow on essaye de trouver une droite qui s'approche le plus possible

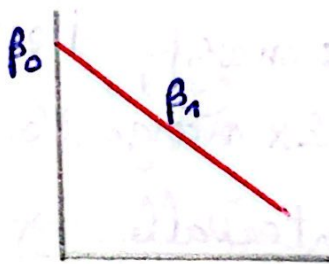
Modèle de regression linéaire simple

l'équation de $y = \beta_0 + \beta_1 x + E$

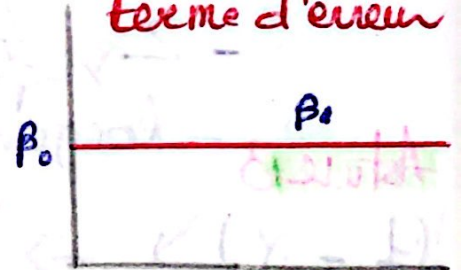
représente la partie non linéaire dans la relation
 terme d'erreur



⇒ R. linéaire positive



⇒ R. linéaire négative



⇒ R. linéaire de pente nulle

droit de regression estimé : $\hat{y} = b_0 + b_1 x$ droit d'échantillon

comment calculer le coefficient de détermination ?

$$r^2 = \frac{SC \text{ rég}}{SCT}$$

Somme des carrés de régression

$$SC \text{ rég} = \sum (\hat{y}_i - \bar{y})^2$$

d'après $\hat{y} = b_0 + b_1 x$

qui permet de juger la qualité du fit

Somme des carrés Totales

$$SCT = \sum (y_i - \bar{y})^2$$

$$SCT = SC \text{ rég} + SC \text{ rési}$$

Somme des carrés résiduels

$$SC \text{ rés} = \sum (y_i - \hat{y}_i)^2$$

Aussi

$$r^2 = 1 - \frac{SC \text{ Rési}}{SCT}$$

comment calculer le coefficient de corrélation ?

$$r_{xy} = \frac{S_{xy}}{S_{yx} \cdot S_x}$$

$$r_{xy} = \frac{\sigma_{xy}}{\sigma_x \cdot \sigma_y}$$

$$r_{xy} = (\text{signe de } b_1) \cdot \sqrt{r^2}$$