

# Statistique Appliquée

I- Statistique descriptive: ≠ probabiliste

L'ensemble des instruments et de recherches mathématiques permettant de déterminer les caractéristiques d'un ensemble de données.

• But: d'extraire et de résumer des informations pertinentes d'une liste de nombres difficile à interpréter par une simple lecture.

↳ Statistiques exploratoires: on explore d'abord les données pour avoir une idée qualitative de leurs propriétés

↳ Statistiques confirmatoires: on fait des hypothèses de comportement que l'on confirme ou que l'on infirme en recourant à d'autres techniques statistiques.

• Outils:

- + La collecte des données
- + Le traitement des données collectées
  - Utiliser des représentations graphiques
  - Calculer certains paramètres et indicateurs
- + Interprétation des résultats
- Basée sur des données existantes.

⊕ Statistique probabiliste: Basée sur des données antérieures, dont je connais tout les résultats mais non pas un exact.

↳ Chercher la loi de probabilité.

⊖ Inférentielle

• Inférence: une opération logique par laquelle on admet une proposition en vertu de sa liaison avec d'autres propositions déjà faites.

• L'inférence statistique: consiste à induire des conclusions (caractéristiques) concernant un groupe auquel on ne peut pas accéder directement (grande taille - coûteux) à partir d'un sous-groupe (petite taille) auquel on a accès et que l'on considère comme un échantillon aléatoire provenant de cette population.

↳ Les caractéristiques de l'échantillon représentent une certaine marge d'erreur par rapport à celles de la population.

⇒ Consiste à partir d'un échantillon de données provenant d'une population de loi de probabilité inconnue, à déduire des propriétés fiables sur cette population.

Elle vise à:

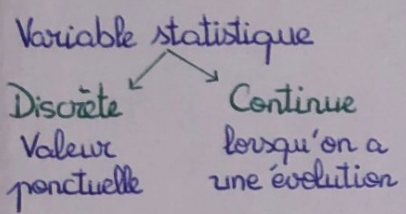
- Évaluer un paramètre ou une relation
- Prédire une valeur
- Déterminer si les différences sont dues au hasard.

Vocabulaire:

- + Population: Un ensemble de personnes, d'objets ou d'événements, base de l'étude statistique.
- + Individu: Un élément de cette population.
- + Échantillon: un sous-ensemble de la population, ayant les mêmes caractéristiques de la population mère utilisé en vue d'inférer qlq chose à cette pop.
- + Caractère: particularité ou propriété caractéristique de la population

Quantitatif:

Qualitatif:



- Ordinale (ordonnée)
- Nominale

- + Effectif d'une population: c'est le nombre total des éléments constituant cette population (N)
- + Fréquence d'un caractère: c'est le nombre d'individus possédant ce caractère divisé par N

Indicateurs de forme: Graphiques et tableaux

! Les valeurs que prendra un caractère = Modalités doivent être ordonnées.

Modalités ou classes	Eff.	Eff. cumulé ↑	Eff. cumulé ↓	freq.	freq. cumulé ↑	freq. cumulé ↓	%
$m_i$	$n_i$	→ fréquence Absolue	Ncd	$\frac{n_i}{N}$	→ fréquence relative	Fcd	100 %

- \* VSD: on les représente par un histogramme, ou des secteurs
- \* VSC: ! classes adjacentes, d'amplitude pas forcément égales.
  - On les représente par un histogramme dont les rectangles sont de largeur = amplitude.

- Dans la pratique, on prend généralement entre 5 et 20 intervalles.
- Pour déterminer l'amplitude de l'intervalle:  $\frac{\text{Grd valeur} - \text{ptt valeur}}{\text{nb d'intervalles}}$
- \* Variable qualitatif: représentation en secteurs.

Valeurs numériques:

- + Issues d'un échantillon: stat. d'échantillon
- + Issues d'une population: paramètres de la population
  - Statistique d'échantillon = Valeur numérique utilisée comme mesure d'un échantillon.
  - Paramètres de la population = Valeur numérique utilisée comme mesure de la population.
  - Estimateur ponctuel = Statistique d'échantillon utilisée pour estimer le paramètre correspondant de la population.

Indicateurs de tendance centrale:

+ Moyenne:

$$\bar{x} = \frac{\sum x_i}{N}$$

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

= échantillon

= population

+ Médiane: Valeur qui partage l'ensemble des données en 2 parties égales (! ordre !)

On calcule:  $\frac{N}{2}$

↳ décimale: on prend l'élément entier qui a la position après

↳ Entière: On somme les valeurs au position i et (i+1) et on divise par 2.

La valeur de l'observation la plus fréquente (on peut avoir 2 modes)

+ **Percentiles** : on appelle  $p$  percentile, la valeur qui partage les observations en  $p\%$  à gauche et  $(100-p)\%$  à droite.

• La médiane correspond au 50<sup>e</sup> percentile.

• Les quartiles : 25<sup>e</sup> percentile, 50<sup>e</sup> percentile, 75<sup>e</sup> percentile

Calcul :  $i = \frac{P}{100} * N$    
 { décimale : val. entier ↗  
 entier :  $\frac{(i)^e + (i+1)^e}{2}$

+ Mesures de dispersion:

+ **Étendue** : La différence entre la plus grande valeur et la plus petite.

+ **Étendue interquartile (EIQ)** : La différence entre

le 3<sup>e</sup> et le 1<sup>er</sup> quartiles  $EIQ = Q_3 - Q_1$   
 = Étendue de 50% des observations qui sont au milieu.

+ **Variance** :  $\left\{ \begin{array}{l} \text{Populat} = \sigma^2 = \frac{\sum (x_i - \mu)^2}{N} \\ \text{Échantillon} = s^2 = \frac{\sum (x_i - \bar{x})^2}{n-1} \end{array} \right.$

+ **Écart type** :  $\left\{ \begin{array}{l} \sigma = \sqrt{\sigma^2} \\ s = \sqrt{s^2} \end{array} \right.$

+ **Coefficient de variation** :  $\frac{\text{Écart type}}{\text{moyenne}} * 100$

↳ Combien s'élève l'écart-type par rapport au moyen.   
 Intéret d'avoir le CV, le + petit possible,  $\downarrow CV \Rightarrow \downarrow$  dispersion   
 le CV est indépendant de l'unité de mesure.

+ Détection des valeurs singulières :

Valeur singulière : observation anormalement grande ou petite.

↳ Erreur d'enregistrement : À corriger

↳ Observation pas correctement incluse dans l'ensemble des données : À supprimer

↳ Valeur inhabituel correctement enregistrée et

qui appartient à l'ensemble des données : À conserver

→ Ces valeurs doivent être mentionner en matière d'interprétation (Biaisier l'étude)

→ **Moyenne tronquée** : calculer la moyenne de 90% de la populat<sup>on</sup> seulement (pour éviter les val. sing)

• Forme d'une distribution :

→ **Degré d'asymétrie** :  $\gamma_1 = \frac{\mu_3}{\sigma^3}$

Avec :  $\mu_3 = \frac{1}{n} \sum_{i=1}^h n_i (x_i - \bar{x})^3$

\* Des données biaisées à gauche  $\gamma_1 < 0$

\* Des données biaisées à droite  $\gamma_1 > 0$

• Variable centrée réduite z :

$z_i = \frac{x_i - \bar{x}}{s}$  → distance

↳ mesure la distance en nb d'écart type entre une observation  $x_i$  et la moyenne.

• Théorème de Chebyshev :

Au moins  $(1 - \frac{1}{z^2})$  des observations doivent se situer ds l'intervalle  $[\bar{x} - zs, \bar{x} + zs]$ ,

Avec  $z > 1$ .

Ex : Combien des étudiants ont obtenu une

note entre 60 et 80. Sachant que  $\left\{ \begin{array}{l} \bar{x} = 70 \\ s = 5 \end{array} \right.$

↳  $60 = \bar{x} - 2s = 70 - 2 \times 5$

$60 = 70 - z \times 5 \Rightarrow z = 2$

Donc  $1 - \frac{1}{2^2} = 75\%$  des étudiants ont obtenu une note entre 60 et 80.

• Règle empirique :

Pour une distribution en forme de cloche :

\* 68% des observations se situent ds  $[\bar{x} - s, \bar{x} + s]$

\* 95% " "  $[\bar{x} - 2s, \bar{x} + 2s]$

\* Presque toutes les observat<sup>ions</sup> ...  $[\bar{x} - 3s, \bar{x} + 3s]$   
 = 99,7%

⇒ Considérer que toutes les valeurs en dehors de  $[\bar{x}-3s, \bar{x}+3s]$  sont des valeurs singulières.

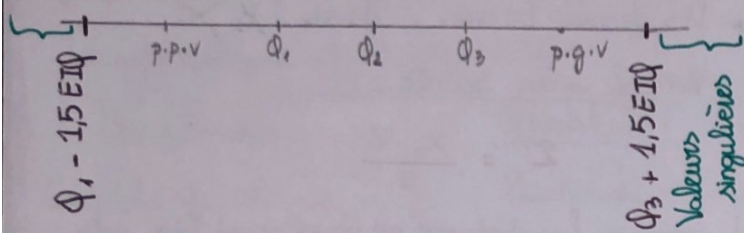
! Différence entre le th. de Chebyshev et la règle:

- $Z > 1$ 
  - $Z = 1$
  - Elle ne marche que pour les distributions normales
  - Elle donne + de précisions

• Analyse exploratoire des données:

**But:** Éliminer les valeurs singulières

**Données:** + petite valeur, + grande valeur,  $Q_1, Q_2, Q_3$



## II - Statistique bivariée:

1- Mesures de la relation entre 2 variables:

- La série statistique est représentée par un nuage de points.
- Dans la statistique bivariée, on peut calculer le point moyen  $G(x_G, y_G)$ 

$$\begin{cases} x_G = \frac{1}{n} \sum x_i = \bar{x} \\ y_G = \frac{1}{n} \sum y_i = \bar{y} \end{cases}$$
- Ajustement affine
  - + Ajustement à la règle (en passant par le pt. moyen)
  - + Méthode de Mayer = droite de Mayer (2 sous-nuages, 2 points moyens)

2- Mesures de la covariance:

↳ Mesure de la relation **linéaire** entre 2 variables.

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

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

↳ Valeurs  $> 0$  ; Relation linéaire positive

↳ Valeurs  $< 0$  ; Relation linéaire négative

3- Mesures par le coef. de corrélation:

↳ Mesure de la relation **linéaire** entre 2 variables, dont les valeurs sont comprises entre  $-1$  et  $+1$

\* Population = 
$$r_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$$

\* Échantillon = 
$$r_{xy} = \frac{s_{xy}}{s_x s_y}$$

↳ Valeurs proches de 0 = Absence de relat. linéaire

↳ Valeurs proches de 1 = forte relat. linéaire positive

↳ Valeurs proches de -1 = forte relation linéaire négative

↳ Méthode des moindres carrés:

\* Droite de régression de Y en X

$(D_{Y/X}): y = ax + b$ 

$$\begin{cases} a = \frac{\sigma_{xy}}{(\sigma_x)^2} \\ b = \bar{y} - a\bar{x} \end{cases}$$

\* Droite de régression de X en Y } les 2 droites se coupent au point moyen

$(D_{X/Y}): x = a'y + b'$ 

$$\begin{cases} a' = \frac{\sigma_{xy}}{\sigma_y} \\ b' = \bar{x} - a'\bar{y} \end{cases}$$

↳ meilleur situat. c'est d'avoir 2 droites confondues:

Rq: • Valeur faible de  $r$  = il est possible d'avoir une forte relat. mais non linéaire

$r$  = coef. de corrélat. de la partie linéaire entre X et Y.

• 2 variables dont  $r$  proche de 0 : **décorréliées**

≠ indépendantes

• Un fort  $r$  n'implique pas une forte relation de causalité entre X et Y (Existence possible d'un 3<sup>ème</sup> var.)

• Passer d'une relat. non linéaire à une relation linéaire

Par exemple =  $y = Cx^m$  ou  $y = Ca^x$

égative  $\Rightarrow \ln y = \ln C x^m = \ln C + m \ln x$   
 $\ln x + \ln C \Rightarrow \begin{cases} m = \frac{\sigma_{xy}}{(\sigma_x)^2} \\ C = e^b = e^{\bar{y} - a\bar{x}} \end{cases}$

### III - Probabilités :

= Évaluation du caractère probable d'un événement.

#### 1- Dénombrement - Théorie des ensembles :

- 2 ensembles égaux  $A=B$  si  $A \subseteq B$   $\begin{cases} x \in A \\ x \in B \end{cases}$
- $A \subseteq B$  ( $\forall x \in A \Rightarrow x \in B$ ) : A sous ensemble de B.
- $\emptyset$  Ensemble qui ne contient aucun élément.
- $A \cap B$  (intersection)  $\{x \in A \text{ et } x \in B\}$
- $A \cup B$  (Union)  $\{x \in A \text{ ou } x \in B\}$
- $A \setminus B = \{x \in A \text{ et } x \notin B\} = A \setminus (A \cap B)$
- $\bar{A} = C_E^A = E \setminus A = \{x \notin A\}$
- $A \cap B = \emptyset$  : A et B disjoints
- $A \cap \emptyset = \emptyset$  ;  $A \cup \emptyset = A$  ;  $\overline{\bar{A}} = A$
- $A \cap \bar{A} = \emptyset$  ;  $A \cup \bar{A} = E$
- $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$
- $A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$
- $A \times B = \{a.b / a \in A, b \in B\}$
- $\overline{A \cap B} = \bar{A} \cup \bar{B}$  ;  $\overline{A \cup B} = \bar{A} \cap \bar{B}$
- $\text{Card}(\bar{A}) = \text{Card}(E) - \text{Card}(A)$
- $\text{Card}(A \cup B) = \text{Card}(A) + \text{Card}(B) - \text{Card}(A \cap B)$
- $\text{Card}(A \setminus B) = \text{Card}(A) - \text{Card}(A \cap B)$
- $\text{Card} \emptyset = 0$

#### 2- Analyse combinatoire :

- Méthodes permettant de donner le nb de tous les résultats possibles d'une expérience donnée.
- ! On doit poser 2 questions :  $\begin{cases} * \text{ L'ordonnance} \\ * \text{ La répétition} \end{cases}$

	Sans répétition	Avec répétition
Ordonnées = Arrangement	$A_n^p = \frac{n!}{(n-p)!}$	$n^p$
non ordonnées = Combinaison	$C_n^p = \frac{n!}{p!(n-p)!}$	$K_n^p = \frac{(n-1+p)!}{p!(n-1)!}$

\* Si on a un seul choix : l'ordre et la répétition = non importants

#### + Vocabulaire :

- \* Expérience aléatoire : action dont le résultat ne peut être prévu avec certitude.
- ↳ Résultats possibles = Éventualités = événements élémentaires ( $\Sigma$  éventualités = univers  $\Omega$ )
- ↳ Événement = ensemble d'éventualités noté A  
 $\Omega$  = événement certain ;  $\emptyset$  événement impro

#### 3- Calcul des probabilités :

- Probabilité d'un événement : la chance que cet événement se réalise.  $! 0 \leq P(A) \leq 1$
- ! 2 éléments incompatibles = disjoints  $P(A \cap B) = 0$   
 = pas d'intersection  $\neq$  indépendants (il n'y a pas d'influence)

#### • Probabilités conditionnelles : (indépendance)

La probabilité d'un événement est influencée par le fait qu'un événement lié au premier, se soit produit.

Probabilité de B sachant A :  $P(B/A) = \frac{P(A \cap B)}{P(A)}$

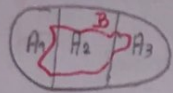
Avec :  $P(A) > 0$

- ↳ Si A n'a pas d'influence = indépendance  
 $P(A \cap B) = P(A) * P(B)$

**Théorème des probabilités totales:**

Soient  $A_1, A_2, A_3, \dots, A_n$  des événements tq:

$$\begin{cases} A_1 \cup A_2 \cup \dots \cup A_n = \Omega \\ A_i \cap A_j = \emptyset \end{cases}$$



$$P(B) = \sum_{i=1}^n P(B/A_i) \cdot P(A_i)$$

$$\begin{aligned} B \cap \Omega = B &\Rightarrow B \cap (A_1 \cup A_2 \cup A_3) = B \\ (B \cap A_1) \cup (B \cap A_2) \cup (B \cap A_3) &= B \\ P(B) &= P(B \cap A_1) + P(B \cap A_2) + P(B \cap A_3) \\ P(B) &= P(B/A_1) \cdot P(A_1) + \dots \end{aligned}$$

**Théorème de Bayes:**

$$P(A_i/B) = \frac{P(A_i) P(B/A_i)}{\sum_{i=1}^n P(A_i) P(B/A_i)}$$

**IV - Variables aléatoires et Lois de probabilité:**

Notion de variable aléatoire: une fonction qui dépend d'une expérience aléatoire, et qui associe à chaque événement une valeur  $x_i$  ( $X: \Omega \rightarrow \mathbb{R}$ )

↳ discrète = ensemble de ses valeurs est fini ou infini dénombrable.

↳ Continue = ensemble de ses valeurs est  $\mathbb{R}$  ou un intervalle de  $\mathbb{R}$ .

V.A.D et lois de probabilités discrètes:

Loi de probabilité de X: une fonction qui décrit comment sont distribuées les probabilités en fonction des valeurs de la v.a.d.

Exemple: Jet de dé;  $\Omega = \{1, 2, 3, 4, 5, 6\}$

$X =$  "Gain de jeu?"

$$\begin{cases} 1, 2, 3 \rightarrow 10 \text{ dh} \\ 4, 5 \rightarrow -5 \text{ dh} \\ 6 \rightarrow 0 \text{ dh} \end{cases}$$

$x_i$	-5	0	10
$P_i$	2/6	1/6	3/6

↳  $P(X=x_i)$

**Fonction de répartition:  $F: \mathbb{R} \rightarrow \mathbb{R}$**

$x \rightarrow F(x) = P(X \leq x)$

Ex:

$$F(x) = \begin{cases} 0 & \text{si } x < -5 \\ 0 + \frac{2}{6} & \text{si } -5 < x < 0 \\ 0 + \frac{2}{6} + \frac{1}{6} & \text{si } 0 < x < 10 \\ 0 + \frac{2}{6} + \frac{1}{6} + \frac{3}{6} & \text{si } x \geq 10 \end{cases}$$

Rq:

$$P(X > a) = 1 - P(X \leq a)$$

Espérance:  $\approx$  la moyenne

$$E(X) = \mu = \sum_{i=1}^n P_i x_i$$

Variance:

$$\begin{aligned} V(X) = \sigma^2 &= \sum_{i=1}^n P_i (x_i - \mu)^2 \\ &= \sum P_i x_i^2 - (E(x))^2 \end{aligned}$$

Loi usuelles:

+ Loi uniforme: Les valeurs prises par v.a ont la même probabilité = équiprobabilité

$$\forall i, P(X=x_i) = \frac{1}{n}$$

$n$ : le nombre des valeurs

+ loi de Bernoulli:

Variable de Bernoulli: 2 réponses  $\begin{cases} \text{Succès} = 1 \\ \text{Échec} = 0 \end{cases}$

$x_i$	0	1
$P_i$	$1-p$	$p$

$$\begin{aligned} E(X) &= \sum P_i x_i = p \\ \sigma(X)^2 &= p(1-p) \end{aligned}$$

+ loi binomiale:

Variable binomiale: le nb de succès obtenus lors de la répétition de  $n$  épreuves identiques et indépendantes (loi de Bernoulli répété plusieurs fois).

$$S_n = X_1 + X_2 + \dots + X_n$$

$X_i$  = variable de Bernoulli

$X) = P(X)$  Résultats de  $x$  succès pendant  $n$  tirages

$$C_n^x = \frac{n!}{x!(n-x)!}$$

+ Fonction de proba.

$$f(x) = C_n^x p^x (1-p)^{n-x}$$

+ Espérance  $E(X) = np$

+ Variance  $V(X) = np(1-p)$

\* Loi de poisson: (utiliser ds les tx d'arrivée ds une file d'attente)

Variable de poisson: décrire le nb d'occurrences d'un événement au cours d'un intervalle de temps ou d'espace bien défini.

+ La probabilité d'occurrence est la même dans 2 intervalles de même longueur

+ Fonction de proba.:

$$f(x) = \frac{\mu^x e^{-\mu}}{x!}$$

nb moyen d'arrivée ds une file pendant un intervalle:  $\mu$

+ Variance:  $V(X) = \mu$

+ V.A.C et lois de probabilités continues:

V.A.C est définie par sa fonction de répartition.

Fonction densité de proba.  $f$ :  $f = F'$

ou encore:  $F(x) = \int_{-\infty}^x f(y) dy$

Donc: l'aire située sous le graphique de  $f(x)$  dans un intervalle particulier, donne la proba.

que la v.a.c  $X$  prenne une valeur ds cet intervalle.

• La proba. que la v.a.c prenne une valeur particulière est nulle.

Espérance:  $E(X) = \int_{-\infty}^{+\infty} x f(x) dx$

Variance:  $V(X) = \int_{-\infty}^{+\infty} x^2 f(x) dx - (E(X))^2$

Ecart type:  $\sigma(X) = \sqrt{V(X)}$

! L'aire totale de la courbe est égale à  $1 \left( \int_{-\infty}^{+\infty} f(x) dx = 1 \right)$

+ Loi normale centrée réduite:

Si  $\forall x \in \mathbb{R}$ ,  $f(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$

;  $\mu = 0$ ,  $\sigma = 1 \Rightarrow N(0,1)$

+ Approximation normale des proba. binomiales:

Loi binomiale, dans le cas où:

$$np \geq 5 \text{ et } n(1-p) \geq 5$$

→ La loi normale permet d'estimer facilement des proba. binomiales, en posant:

$$\mu = np \text{ et } \sigma = \sqrt{np(1-p)}$$

+ Loi exponentielle:

Si  $f(x) = \frac{1}{\mu} e^{-x/\mu}$ ;  $x > 0$ ;  $\mu > 0$

## Ch 2: Introduction à la théorie de l'échantillonnage

⊖ Théorie de sondage: l'optimisation de la collecte de données.

- Sondage = l'observation partielle d'une partie d'une population statistique.
- Méthode d'échantillonnage: procédure par laquelle on choisit dans une population un sous-ensemble représentatif.

↳ **Obj:** Avoir un échantillon suffisamment représentatif pour que les données puissent être extrapolées à la population.

### 1- Échantillonnages aléatoires:

Lorsque chaque individu a une proba. connue et non nulle d'être inclus dans l'échantillon.

\* Échantillon aléatoire simple: de taille  $n$  issu d'une population finie de taille  $N$ .

Est un échantillon sélectionné de manière à ce que chaque échantillon possible de taille  $n$  ait la même probabilité d'être sélectionné. ( $p = \frac{n}{N}$ )

+ **Avantage:** Conception simple

• Base de la théorie statistique.

- **Incon:** Difficile sur des grands échantillons.

• Prob de la disponibilité de la base de sondage.

Ex: tirage au sort

\* **EAS** issu d'une population infinie.

Conditions

- Chaque élément sélectionné provient de la même population
- Chaque élément est sélectionné de façon indépendante

Ex: contrôle de qualité

• Base de sondage: liste de toutes les enquêtes.

\* Échantillonnage aléatoire stratifié:

Il consiste à diviser la population en sous-groupes appelés **Strates**, tq chaque élément apparaît dans une et une seule strate.

↳ Sélection d'un échantillon aléatoire simple dans chaque strate (+ homogénéité des strates)

Comment faire la répartition de sujets de chaque strate?

\* **Allocat égale** = m<sup>ême</sup> nb de sujets dans chaque strate.

\* **Allocat proportionnelle** = à la taille de strate

Ex: Étudiants 1<sup>ère</sup> année / 2<sup>ème</sup> 3<sup>ème</sup> ...

• Population 3000 { 2000 jeunes → Éch: 200 { 200 × 66,6% = 133  
1000 vieux → 200 × 33,3% = 67

\* Échantillonnage par groupes:

Partager par rapport à un critère hétérogène.

Il consiste à diviser la population en sous-groupes = **groupes**, tq chaque élément de la pop. appartient à une seule groupe. Puis on sélectionne de manière aléatoire **des groupes**.

Ex: division de Masc en 12 régions / hétérogénéité de la pop. Choisir au hasard de 2 cartons parmi 10 pour vérifier la qualité des prod.

Cas idéal: chaque groupe est une version à petite échelle de la population entière.

+ Les éléments de chaque groupe sont hétérogènes.

\* Pas de l'échantillon =  $\frac{\text{Population}}{\text{taille d'échantillon}} = K$

↳ Avancer par un pas ↓



## Echantillonnage systématique:

que la populat<sup>i</sup> est très importante.

## 2- Echantillonnages non aléatoires:

③ par choix raisonné.

Utilisés dans les études qualitatives.

### \* Echantillonnage de commodité:

les éléments sont inclus dans l'échantillon sans proba. connue ou préspecifiée d'être choisis.

Ex: les volontaires qui participent à une expérience à l'école.

↳ coût très faible ou n<sup>o</sup> nul.

### \* Echantillonnage subjectif:

Sélectionner des éléments de la populat<sup>i</sup> qui semblent représenter la populat<sup>i</sup>

Ex: Un journaliste qui peut choisir 3 pers<sup>o</sup> (Opinion générale).

**Avantage:** Échantillon facilement construit

• Données facilement collectées

**Inc:** Mauvaise représentativité de l'échantillon

• Aucune procédure stat. ne permet de faire une analyse proba. ou de l'inférence stat.

\* E. à participat<sup>i</sup> volontaire:

Faire appel à des volontaires pour constituer l'É.

\* Enquête boule de neige:

On choisit des individus qui sont pertinents pour l'étude, ensuite on leur demande de proposer d'autres individus.

\* Enquête par Quota:

basée sur la répartit<sup>i</sup> connue de la pop. pour un certain nombre de caractères (sexes, âge, ...)

## \* Echantillonnage aléatoire simple:

• Estimation ponctuelle:

la proport<sup>i</sup> de l'échantillon

$\bar{x}$ ,  $s$ ,  $\bar{p}$  sont des estimateurs ponctuelles de  $\mu$ ,  $\sigma$ ,  $p$ .

• Distribut<sup>i</sup> d'échantillonnage: nb d'échantillon possible =  $C_N^n$

Soit:

• **Expérience aléatoire:** Sélection d'un E.A.S

• **V.A:** Valeur de la moyenne d'échantillon  $\bar{x}$ .

$\bar{x}$  a une espérance, une variance, un loi de proba = **distribut<sup>i</sup> d'échantillonnage de  $\bar{x}$**

La connaissance de cette distribut<sup>i</sup> nous permet de tirer des conclusions en termes de proba.,

Écart-type de l'échantillon, moyenne de la pop.

Ésérance mathématique de  $\bar{x}$ :  $E(\bar{x}) = \mu$

⇒ La moyenne de tous les échantillons = moy. de tous les moy.

Moyenne de la populat<sup>i</sup>

Écart-type:

Erreur type de la moyenne

\* Populat<sup>i</sup> finie =  $\sigma_{\bar{x}} = \sqrt{\frac{N-n}{N-1}} \left( \frac{\sigma}{\sqrt{n}} \right)$

\* Populat<sup>i</sup> infinie =  $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$

! Lorsque la taille de l'échantillon est inf. à 5% de la populat<sup>i</sup>:  $\left( \frac{n}{N} \leq 0,05 \right)$

\* **Théorème central limite:**

Lorsque la taille de l'échantillon devient importante (+30) la distribution d'échantillonnage<sup>de  $\bar{x}$</sup>  peut être approchée par une distribut<sup>i</sup> de proba. normale.

• **Distribution d'échantillonnage de  $\bar{p}$**

$$\bar{p} = \frac{x}{n}$$

le nombre d'éléments dans l'échantillon possédant la caractéristique à laquelle on s'intéresse.

Espérance:  $E(\bar{p}) = p$

Écart type:

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

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

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

$np > 5$  et  $n(1-p) > 5$

$\beta(n, p) \rightarrow \text{LNCR}$

### Conclusion:

•  $\nearrow$  la taille de l'échantillon  $\rightarrow$

$\downarrow$  l'écart type  $\downarrow \rightarrow$  Probabilité  $\nearrow$   
[ moyenne  
[ proportionnalité

### \* Estimation par intervalle:

\* De la moyenne d'une population:

• Lorsque  $\sigma$  est connu, l'estimation par intervalle est donnée par

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

$\alpha$ : coeff. de risque

$1 - \alpha$ : indice de confiance

$z_{\alpha/2} \frac{\sigma}{\sqrt{n}} = \text{marge d'erreur}$

$\Rightarrow \downarrow$  indice de confiance  $\rightarrow \downarrow$  intervalle  
 $\rightarrow$  moyenne bonne

• Lorsque  $\sigma$  est inconnu

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

T de student (doi) =  $n - 1$  degré de liberté.

### \* De la proportion d'une population

$$\left[ \bar{p} \pm z_{\alpha/2} \sqrt{\frac{\bar{p}(1-\bar{p})}{n}} \right]$$

### \* Déterminant de la taille d'échantillon:

$$n = \frac{(z_{\alpha/2})^2 \sigma^2}{E^2}$$

$$n = (z_{\alpha/2})^2 \frac{p^*(1-p^*)}{E^2}$$

$p^*$ : proportion d'une ancienne étude  
 $p^* = 0,5$

$z_{\alpha/2}$ : la variable centrée réduite qui correspond à la valeur qui permet de négliger  $\alpha/2$  des informations dans la queue droite sup.

# Régression linéaire

## Simple

Une droite qui peut s'approcher de plus à l'ensemble des points (x et y); afin de savoir s'il existe une relation linéaire entre les 2 variables.

### + Modèle de régression linéaire simple:

$$y = \beta_0 + \beta_1 x + \epsilon \rightarrow \text{représente la partie non linéaire de la relat.}$$

y: variable à expliquer (dépendante)

x: variable explicative (indépendante)

$\epsilon$ : terme d'erreur; la différence entre les valeurs réelles et les valeurs obtenues si la relation était exacte.

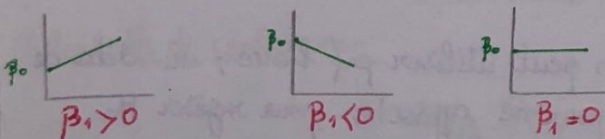
\* **Erreur de spécification:** la seule variable explicative n'est pas suffisante pour expliquer la totalité du phénomène.

\* **Erreur de mesure:** les données ne représentent pas exactement le phénomène.

\* **Erreur de fluctuation d'échantillonnage:** d'un échantillon à l'autre les observations et donc les estimations sont différentes.

### + Équation de régression linéaire simple:

moyenne de y  $\leftarrow E(y) = \beta_0 + \beta_1 x$

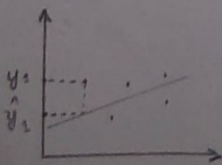


$\Rightarrow$  Relat linéaire  $\oplus$  R. linéaire  $\ominus$  Pas de relat

$\hookrightarrow$  En pratique,  $\beta_0$  et  $\beta_1$  ne sont pas connus donc on doit les estimer en utilisant les données d'un échantillon; d'où la droite de régression estimée:

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

droite de l'échantillon



### \* La méthode des moindres carrés:

Permet d'utiliser les données de l'échantillon pour estimer l'équation de la régression ( $b_0, b_1$ )

Elle consiste à minimiser l'écart au carré entre  $y_i$  et  $\hat{y}_i$ .

$$\text{Min } \sum (y_i - \hat{y}_i)^2$$

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

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

### \* Coefficient de détermination:

#### + Somme des carrés des résidus:

$$SC_{res} = \sum (y_i - \hat{y}_i)^2$$

$\Rightarrow$  La valeur de  $SC_{res}$  est une mesure de l'erreur commise en utilisant l'équation estimée de la régression pour estimer les valeurs de la variable dépendante dans l'échantillon.

La valeur de  $SC_{res}$  est une mesure de l'erreur commise en utilisant l'équation estimée de la régression pour estimer les valeurs de la variable dépendante dans l'échantillon.

#### + Somme des carrés totale:

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

$\hookrightarrow$  Pour estimer  $y_i$  sans avoir les  $x_i$ , on utilise  $\bar{y}$

$\Rightarrow$   $SC_T$  est une mesure de l'ajustement autour des observations de la droite  $\bar{y}$ .

#### + Somme des carrés de la régression:

$$SC_{reg} = \sum (\hat{y}_i - \bar{y})^2$$

$$\Rightarrow SC_T = SC_{res} + SC_{reg}$$

$$r^2 = \frac{SC_{reg}}{SC_T} \Rightarrow \text{Coef. de détermination}$$

$$\frac{SC_T}{SC_T} = \frac{SC_{res} + SC_{reg}}{SC_T} \xrightarrow{\text{cas idéal: } y_i - \hat{y}_i = 0} 1 = \frac{SC_{reg}}{SC_T} = r^2$$

$$r_{xy} = \text{signe de } b_1 \sqrt{r^2} \Rightarrow \text{coef. de corrélation}$$

\* Différence entre coef. de corrélation et coef. de détermination:

1- Dans le cas d'une relation linéaire entre deux variables le coef. de détermination et le coef. de corrélation fournissent une mesure de robustesse de la relation.

2- Le coef. de détermination est compris entre 0 et 1, tandis que le coef. de corrélation est compris entre -1 et 1.

3- Alors que le coef. de corrélation est restreint à des relations linéaires entre 2 variables, le coef. de détermination peut être utilisé dans le cas de relation non linéaire et de relations comprenant plus de 2 variables indépendantes.

⇒ Même si le coefficient de détermination est assez élevé, une analyse approfondie de la robustesse du modèle supposé doit être faite.

D'où: Le test d'hypothèses

1- Établir le test d'hypothèse adéquat:

\* test d'hypothèse unilatéral supérieur  $\begin{cases} H_0: \mu \leq \mu_0 \\ H_a: \mu > \mu_0 \end{cases}$

\* Test unilatéral inférieur  $\begin{cases} H_0: \mu \geq \mu_0 \\ H_a: \mu < \mu_0 \end{cases}$

\* Test bilatéral  $\begin{cases} H_0: \mu = \mu_0 \\ H_a: \mu \neq \mu_0 \end{cases}$

$H_0 =$  Hypothèse nulle  $\neq$

$H_a =$  Hypothèse alternative = ce que je dois prouver

	$H_0$ vraie	$H_a$ vraie
Accepter $H_0$	Conclusion correcte	Erreur de 2 <sup>de</sup> espèce
Rejeter $H_0$	Erreur de 1 <sup>ère</sup> espèce	Conclusion correcte

2\* Seuil de signification: La probabilité de faire une erreur de 1<sup>ère</sup> espèce  $\ominus$  Seuil d'acceptation

3\* test de signification:

+ Si  $\sigma$  connu: 
$$Z = \frac{\bar{x} - \mu_0}{\sigma / \sqrt{n}}$$

	test unilatéral inf.	test unilatéral Sup.	test bilatéral
Règle de rejet $H_0$	$Z \leq -Z_\alpha$	$Z \geq Z_\alpha$	$Z \leq -Z_{\alpha/2}$ $Z \geq Z_{\alpha/2}$

+ Si  $\sigma$  est inconnu: 
$$Z = \frac{\bar{x} - \mu_0}{s / \sqrt{n}}$$

	test unilatéral inf.	test unilatéral Sup.	test bilatéral
Rejet de $H_0$	$t \leq -t_\alpha$	$t \geq t_\alpha$	$t \leq -t_{\alpha/2}$ $t \geq t_{\alpha/2}$

⇒ Marge de liberté (n-1)

$Z_\alpha, t_\alpha =$  valeur critique

N.B: On peut utiliser p (l'aire; au delà de z) comme approche pour rejeter  $H_0$

\* Règle de rejet:  $p \leq \alpha$

$$Z = \frac{\bar{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$$
 : test d'hypothèse relatif à la proportion.

de signification (sur la régression linéaire simple):

+ Test de Student :

\* Le test d'hypothèse adéquat :

$$\begin{cases} H_0: \beta_1 = 0 \\ H_1: \beta_1 \neq 0 \end{cases}$$

$E(y) = \beta_0 + \beta_1 x$  (ca concerne la populat<sup>n</sup>)

→ Les propriétés de  $b_1$  : pop → échantillon

+ Espérance =  $E(b_1) = \beta_1$

+ Ecart-type =  $\sigma_{b_1} = \frac{\sigma}{\sqrt{\sum(x_i - \bar{x})^2}}$

→ Distribut<sup>n</sup> normale:

+ Ecart-type estimé de  $b_1$  =

Erreur de  $b_1$  ←  $S_{b_1} = \frac{s}{\sqrt{\sum(x_i - \bar{x})^2}}$

tel que :

$$s^2 = MC_{res} = \frac{SC_{res}}{n-2}$$

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

test :  $t = \frac{b_1}{S_{b_1}}$

\* Règle de rejet :  $t < -t_{\alpha/2}$  ou  $t > t_{\alpha/2}$   
! ddl = n-2

Intervalle de confiance:

$$b_1 \pm t_{\alpha/2} S_{b_1}$$

+ Test F de Fisher :

$$s^2 = MC_{reg} = \frac{SC_{reg}}{1}$$

Statistique de test :  $F = \frac{MC_{reg}}{MC_{res}}$

\* Règle de rejet : (de  $H_0$ )

$$F \geq F_{\alpha}$$

$F_{\alpha}$  est basé sur la distribut<sup>n</sup> de Fisher à 1 ddl au numérateur et (n-2) ddl au dénominateur.

\* Tableau ANOVA:

Source de la variat <sup>n</sup>	Somme des variat <sup>n</sup>	ddl	Moyenne des variat <sup>n</sup>	F
Régression	SC <sub>reg</sub>	1	SC <sub>reg</sub>	$\frac{MC_{reg}}{MC_{res}}$
Résidu	SC <sub>res</sub>	n-2	$\frac{SC_{res}}{n-2}$	
Totale	SC <sub>T</sub>	n-1		