<u>Def statistique : la stat (science des données) est la science de la collecte, de l'analyse, de la prez et de l'interprétation des données.</u>

#### Objectifs:

- maitrise outils de l'analyse stat pr synthétiser l'info des données écos
- maitrise formalisation math pr dev capacités de trad d'un pb en formalisation math.

#### Livres:

- statistics for business and economics (amphi sieges verts)
- statistiques pour l'économie et la gestion (blanc/violet/bleu)

## CHAPITRE 1: INTRODUCTION, GRAPHES ET TABLEAUX

## Section 1. Définitions fondamentales

- 1. Une population est l'E des éléments sur lesquels une étude se porte.
- 2. <u>L'unité statistique</u>, ou élément, est l'unité pour laquelle des données sont recueillies.

  Diff natures : prsn, pays, objet, bactérie...
- 3. Un échantillon est un sous-ensemble de la pop (noté N ou T)
- 4. <u>L'échantillonage aléatoire (simple)</u> «est une procédure utilisée pr select un échantillon de N individus ds une pop de telle way que chaque indiv de la pop est choisi au hasard, aucune influence, chacun peut être choisi et chaque échantillon possible d'une taille donnée N a la mm chance d'être selectionné.
- 5. <u>Une variable stat</u> est une caractéristique des éléments à laquelle on s'intéresse.

  Variables notées y ; x ou z
- 6. Une modalité est une valeur prise par une valeur stat
- 7. <u>Une observation</u> est un E de mesures obtenues pour chaque élément d'un E de données

**Exercice 1**. Pr connaître la fréquentat° des theatres de l'UB, supposons que 200 étudiants sont selected pr rep à un qst sur le nb de fois ou ils sont allés au theatre Population : les étudiants de l'UB

Echantillon : 200 etudiants Unité stat : 1 étudiant

Variable : nb de fois au theatre

Modalités: 2,3..... fois

Exercice 2 : Effet du lvl d'éducation sur les salaires des f mariées. 428 femmes

Pop : femmes mariées Ech : 428 femmes

Unité stat : 1 femme mariée

Variables: niveau d'études // salaire

Modalités : diplôme, nb années d'études // dollars/euros...

- 8. <u>Un paramètre</u> est une mesure numérique qui décrit une caractéristique spécifique d'une pop. Une statistique/estimation est une mesure numérique qui décrit une caract d'un échantillon.
- 9. (L'erreur d'échantillonage est la diff entre le paramètre d'intérêt et son estimation.)
- 10. <u>La statistique descriptive</u> focus sur les procédures graphiques et numériques used pr résumer et traiter les données. <u>L'inférence statistique</u> focus sur l'utilisation des données pr faire des prévisions/ estimations/ tester les théories pr prendre de meilleures décisions.

## **ETAPES D'UN PROJET DE RECHERCHE:**

Question de recherche



Déterminer la pop + paramètres d'intérêt



Echantillonage + sélection variables + collecte des données



Analyse données(stats descriptives vs stats inférentielles)



Interprétation et discussion des résultats

# **SECTION 2. Variables catégoriques**

11. <u>Une variable catégorique/qualitative/facteur est une variable stat dont les</u> modalités sont des catégories (classes)

+classifient les unités stat tq chaque indiv est mis ds une seule catégorie

**REMARQUE 1 :** Une variable indicatrice est une variable catégorique qui prend seulement 2 diff valeurs : 0 et 1. Une variable cat avk L cats peut être ré-exprimée avk L variables indicatrices.

#### Exemple 1:

#### Exemple 2:



12. <u>Une variable numérique/quantitative</u> est une variable stat dont les modalités sont mesurables. Elle est numérique discrète si elle peut prendre un nb fini/infini de valeurs dénombrables, tq 0,1,2...

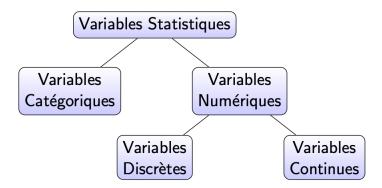
<u>Une variable qui peut prendre ses valeurs nu mds un intervalle/ suite d'intervalles est appelée variable continue.</u>

#### **Exemple 3 : variables discrètes :**

- Nb appels reçus en 5min ds un centre d'appel
- Nb voitures arrivant au péage

#### **Exemple 4 : variables continues**

- Temps écoulé entre les appels reçus
- Tx inflation en France



## Section 3 : les types de données

13. <u>Les données en coupe transversale</u> sont les données liées aux diff indivs collectées au mm moment (ou approximativemnt au mm moment). Les <mark>séries temporelles sont des données collectées sur plusieurs périodes de temps diff.</mark>

## Section 4. représentation des variables catégoriques

<u>Déf 14. La fréquence d'une catégorie</u> (ou classe) est le nb d'unités stats présentant cette dernière.

Soit xi, i=1,2...N, une variable cat avec les cats CI, I=1,2...L Nous notons la fréquence de la cat I, nI La fréquence de la cat CI est le nb d'unités stats i (tq xi appartient à CI) Nous avons l'identité :  $\sum_{I=1}^{L} nI = N$ 

$$\sum_{i=1}^{N} xi = x1 + x2 + \dots + xN$$

<u>Déf 15 : Une</u> distribution de fréquence est un résumé des données sous forme de tableau décrivant la fréquence des ds diff catégories déf par une variable.

Les fréquences ne sont pas direct présentables pr une variable continue/discrète présentant bcp de modalités. On doit parfois regrouper les données.

# Déf 16. La fréquence relative d'une catégorie correspond à la proportion des observations appartenant à cette catégorie :

Fréquence relative d'une catégorie =  $\frac{fréquence d'une \ catégorie}{N}$ 

Nous notons la fréquence relative de la classe CI, fI

Nous avons alors :  $fI = \frac{nI}{N}$ 

Proposition 1 : la somme des fréquences relatives est égale à 1 :  $\sum_{l=1}^{L} fl = 1$ 

### Section 5 : rpz variables numériques

Tableaux de fréquence :

Table 5: Distribution des fréquences d'une variable numérique

Classes	Fréquence	Fréquence relative	Fréquence relative cumulée
$[a_1; b_1)$	$n_1$	$f_1$	$F_1$
$[a_2;b_2)$	<b>n</b> <sub>2</sub>	$f_2$	$F_2$
:	:	:	:
$[a_l;b_l)$	$n_I$	$f_l$	$F_{l}$
:	:	:	:
$\left[ a_{L-1};b_{L-1} ight)$	$n_{L-1}$	$f_{L-1}$	$\mathcal{L}_{\mathcal{K}-1}$
$[a_L;b_L)$	$n_L$	$f_L$	1.00
Total	N	1.00	

Les trois étapes nécessaires à la définition des classes d'une distribution de fréquence pr des données qtt :

- 1- Déter le nb de classes (L) en fonction du nb d'observations
- 2- Déter la largeur approximative de chaque classe : d =  $\sim \frac{max min}{L}$
- 3- Déter les bornes al, bl de chaque classe. Une méthode est de déf la borne inférieure de la 1<sup>e</sup> classe comme min et calculer itérativement les autres bornes en using la largeur approximative trouvée ds l'étape précédente.
- 4- A1 = min, bI = aI + d, et aI= bI 1,  $\sqrt{I=2,...,L}$

<u>Déf 17. La fréquence relative</u> cumulée d'une classe est la somme des fréquences de ttes les classes dont les observations sont inférieures ou égales à la borne supérieure de cette classe.

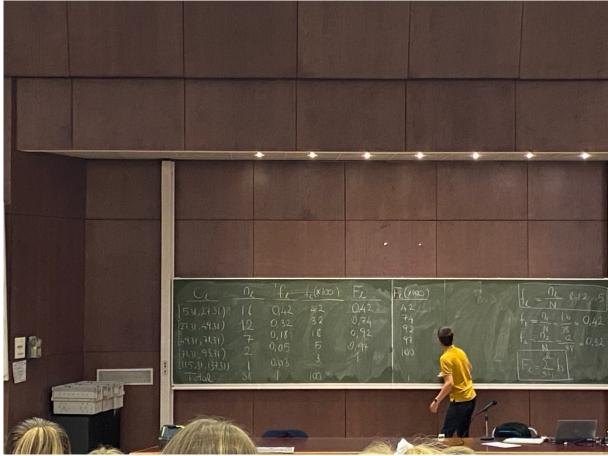
$$\begin{aligned} & \text{FI} = \sum_{S=1}^{l} fs \ l = 1 \\ & \text{F1} = \sum_{S=1}^{1} fs = f1 \\ & \text{F2} = \sum_{s=1}^{s} f3 = f1 + f2 \\ & \text{FL} = \sum_{S=1}^{L} f3 = \text{f1} + \text{f2} + \dots + \text{fL} = 1 \end{aligned}$$

Ex : présenter la distrib du pib/h (1000\$) des pays de l'OCDE en 2020 sous la forme d'un tableau de fréquences

Nb observations dépend du nb de valeurs ds l'intervalle voulu de l'échantillon

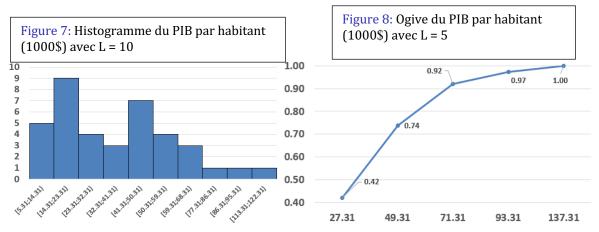
CI (5,31; 27,31) 
$$\rightarrow$$
 nI= 16  $\rightarrow$  fI= $\frac{nl}{N}$ =0,42  $\rightarrow$  fI= x100= 42%  $\rightarrow$  FI= $\sum_{S=1}^{l} fl = 42$ 

42% des pays ont un pib/h entre 5,31 et 27,31



Un histogramme est construit en plac, ant les classes la variable considérée sur l'axe horizontal et la fréquence ou la fréquence (relative) sur l'axe vertical. La fréquence (relative) de chaque classe est représentée par un rectangle dont la base est d'eterminée par les limites de classes et dont la hauteur correspond `a la fréquence (relative).

Un ogive est construit en plaçant les bornes sup'erieurs des classes de la variable consid'er'ee sur l'axe horizontal et les fr'equences relative cumul'ees sur l'axe vertical.



# SECTION 6 : REPRÉSENTATION DES DONNÉES RELATIVES À DEUX VARIABLES

Definition 18 = Un tableau de contingence, ou une tabulation crois ee, est un resume sous forme de tableau des donnees relatives a deux variables.

Les deux variables peuvent etre categoriques ou numeriques.

Supposons que nous voulons étudier la distribution des frequences dans un echantillon selon deux variables  $x_i$  et  $y_i$ , i=1,2,...,N.

Notons les categories (ou` classes) de  $x_i$  et  $y_i$ ,  $C_l$ , l = 1,2,...,L et  $G_k$ , k = 1,2,...,K, respectivement.

Table 6: La forme générale d'un tableau de contingence

$x_i \setminus y_i$	$G_1$	$G_2$		$G_K$	Total
$C_1$	<i>n</i> <sub>11</sub>	<i>n</i> <sub>12</sub>	• • •	$n_{1K}$	$n_{1.}$
$C_2$	<i>n</i> <sub>21</sub>	<i>n</i> <sub>22</sub>	• • •	$n_{2K}$	$n_{2.}$
:	÷	÷	٠	:	:
$C_L$	$n_{L1}$	$n_{L2}$	• • •	$n_{LK}$	$n_{L.}$
Total	n <sub>.1</sub>	n <sub>.2</sub>		n <sub>.K</sub>	N

**Definition 19 =** Une frequence partielle, notee  $n_{lk}$ , l = 1,2,...,L, k = 1,2,...,K, est le nombre d'unites statistiques i telle que  $x_i \in C_l$  et  $y_i \in G_k$ .

Autrement dit,  $n_{lk}$  est le nombre d'unites statistiques presentant a la fois la classe l de x et la classe k de y.

Definition 20 = Une frequence marginale de la variable x, notee n, l = 1,2,...,L, est le nombre d'unites statistiques i telle que x;  $\in C$ ]. Cette definition est identique a la Definition

14.  $\frac{f_{LK}}{N} = \frac{nlk}{N}$  avk nlk l'effectif d'une case et N l'effectif total de l'échantillon (en bas à droite)

## Exemple 7

Table 7: Fréquences de 217 pays classés selon la région et le niveau de revenu

	AME	AEP	EAC	MOA	AS	Total
FR	0	1	0	26	1	28
RE	20	14	38	9	0	81
RITI	5	14	4	25	6	54
RITS	19	9	16	9	1	54
Total	44	38	58	69	8	217

**Fréquence conditionnelle :** eff de chaque case eff total ligne