

Εconoméτρiε .



## Rappel:

1. Définir les lois de probabilité dans un cadre univarié
2. Espérance et variance pour un ensemble de variables aléatoires
3. Cadre multivarié
4. Présentation des lois incontournables
5. Distribution de la moyenne empirique

## 1. Définir les lois de probabilité dans un cadre univarié

### 1. probabilité, univers et variable aléatoire

issues: résultat possible d'une expérience aléatoire sachant qu'une seule issue va se réaliser et que toutes les issues n'ont pas la même chance de se réaliser

probabilité d'une issue: fréquence avec laquelle cette issue est réalisée à  $LI$

univers: ensemble des issues possibles

événement: sous-ensemble de l'univers

variable aléatoire: résumé numérique de l'issue d'une expérience aléatoire



## 2. Distribution de probabilité d'une variable discrète

Distribution de proba: énumération de toutes les valeurs possible de la variable avec leurs probabilités respective

proba des événements:  $P(H=1 \text{ ou } H=2) = p(H=1) + p(H=2)$   
 $= 0,16$

issues (nombre de panne)

0 1 2 3 4

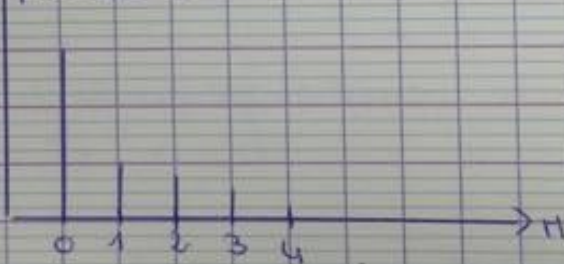
distribution 0,8 0,1 0,06 0,03 0,04  
des probas

fonction de répartition: proba que cette v.a soit  $\leq$  à une valeur

0 1 2 3 4

fct de répartition 0,8 0,9 0,96 0,99 1

~ proba



## loi de Bernoulli

G

G = 0

(femme)

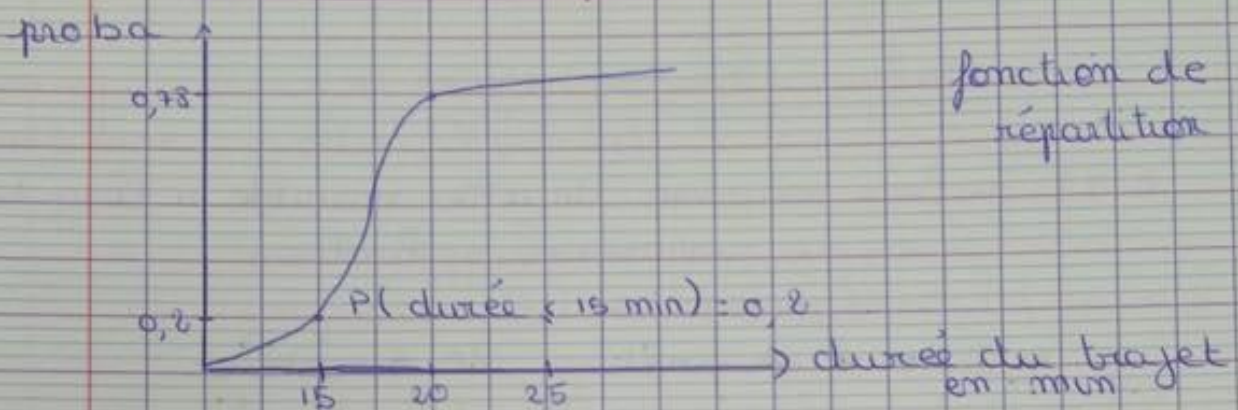
G = 1

(homme)

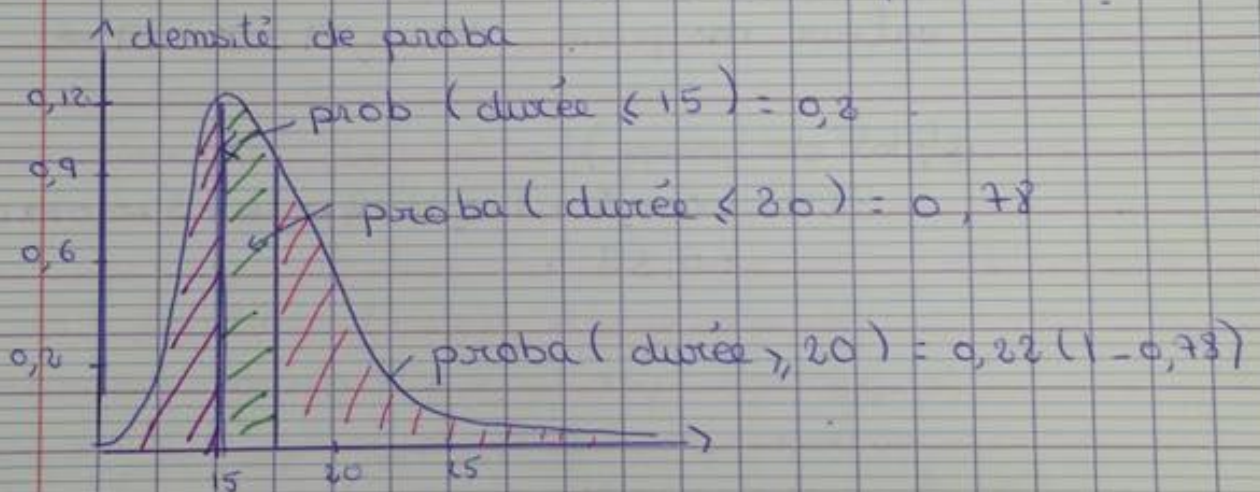
$G = \begin{cases} 1 & \text{avec une proba } p \\ 0 & = = = 1-p \end{cases}$



### 3. Distribution de proba d'une v.a continue



fonction de densité de proba: aire de la probabilité que la v.a s'intercale entre ces 2 points









## 2. Écart-type et variance

variance :  $\text{var}(Y) = E((Y - \mu_Y)^2)$

écart type :  $\sigma(Y) = \sqrt{\text{var}(Y)}$

v.a.  $Y$  avec  $k$  valeurs  $(y_1, \dots, y_k)$

$$\sigma_Y^2 = \text{var}(Y) = E((Y - \mu_Y)^2) = \sum_{i=1}^k (y_i - \mu_Y)^2 p_i$$

(ex) de la panne :  $\text{var}(H) = (0 - 0,35)^2 \times 0,8 + (1 - 0,35)^2 \times 0,1 +$   
 $(2 - 0,35)^2 \times 0,06 + (3 - 0,35)^2 \times 0,03 +$   
 $(4 - 0,35)^2 \times 0,01$

$$= 0,6495$$

$$\sigma_H = \sqrt{\text{var}(H)} = \sqrt{0,6495} = 0,80$$

Bernoulli

$$\text{var}(G) = (0 - p)^2 \times (1 - p) + (1 - p)^2 \times p$$
$$= p(1 - p)$$

## 3. fonction linéaire d'une variable

2. v.a.  $X$  et  $Y$

taux d'imposition des salaires  $\Rightarrow 20\%$

prime annuelle non imposable 2000 €

$$Y = 2000 + 0,8 X$$

ou  $\mu_X$      $\sigma_X^2$

$X$  revenu brut incluant

$$Y = \text{net}$$

$\Rightarrow Y$  v.a.

$$E(Y) = \mu_Y = 2000 + 0,8 \mu_X$$

$$Y - \mu_Y = 0,8(X - \mu_X)$$

$$E((Y - \mu_Y)^2) = E[0,8(X - \mu_X)^2]$$

$$= 0,64 \times E[(X - \mu_X)^2]$$



$$= 0,64 \sigma_x^2$$

$$\sigma_y = 0,8 \sigma_x$$

$$y = 2000 + 0,8 X$$

$$\mu_y = 2000 + 0,8 \mu_x$$

$$y = a + b X$$

X revenus bruts individuel

v.a  $\mu_x$   $\sigma_x^2$

$$\mu_y = a + b \mu_x$$

$$\sigma_y^2 = b^2 \sigma_x^2$$

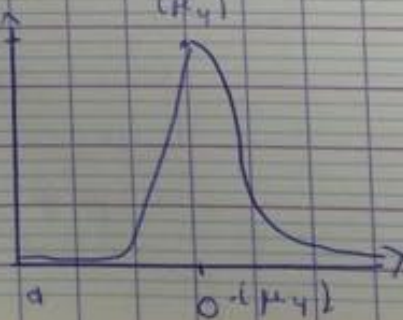
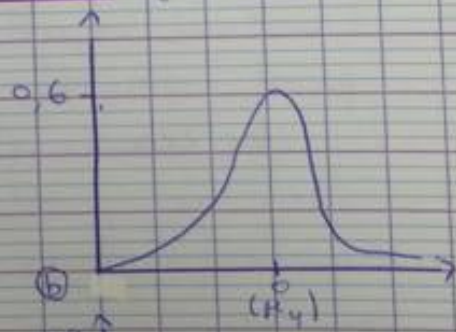
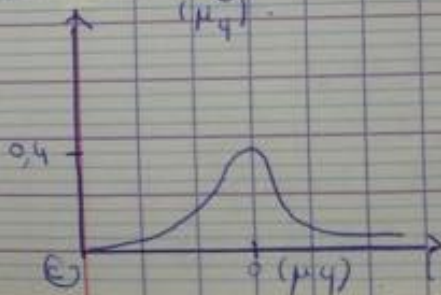
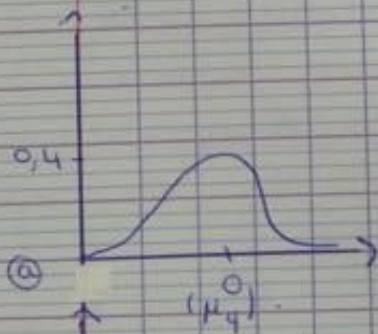
$$\sigma_y = b \sigma_x$$

#### 4. Autres mesures de la forme d'une distribution

coefficient de dissymétrie: (skewness) mesure l'asymétrie de la distribution

coefficient d'aplatissement: (kurtosis) qui mesure son aplatissement

=> basés sur les moments de distribution





coefficient de dissymétrie : 
$$\frac{E((Y - \mu_Y)^3)}{\sigma_Y^3}$$

pour une forme symétrique, le coefficient = 0

= 0  $\Rightarrow$  symétrique

$\neq 0 \Rightarrow$  asymétrique

a) skewness = 0

b) skewness = 0

c) skewness = -0,1

d) skewness = 0,6

coefficient d'aplatissement : 
$$\frac{E((Y - \mu_Y)^4)}{\sigma_Y^4}$$
  
kurtosis

$\Rightarrow$  mesure de la masse qui se trouve au niveau des queues de la distribution  
(jamais < 0)

• la normale = 3

• il est leptokurtique  $> 3$

a) kurtosis = 3

b) kurtosis = 20

c) kurtosis = 5

d) kurtosis = 5

Moment : le moment d'ordre  $\pi$  de  $Y$  est donné par  $E(Y^\pi)$



### III - Cas de 2 v.a

#### 1 - Distribution jointe et marginale

x et y

c'est la proba que ces 2 variables prennent simultanément 2 valeurs données

ex: x et y

distribution des proba jointes =  $P(X=x; Y=y)$

Y v.a 1 durée course ( $\leq 20$  min)  
0 sinon

X v.a 1 il pleut  
0 sinon

	X=1 (pluie)	X=0 (pas pluie)	
long Y=0	0,15	0,07	0,22
court Y=1	0,15	0,63	0,78
	0,30	0,70	1

distribution jointe:  $P(Y=0; X=1) = 0,15$

distribution de proba marginale: distribution de Y prise de manière isolée

Si X peut prendre l' valeur  $\neq x_1, \dots, x_p$   
alors la proba que Y prenne une valeur y  
est de  $P(Y=y) = \sum_{i=1}^p P(X=x_i; Y=y)$



## 2 - Distribution conditionnelle

La probabilité conditionnelle de  $Y$  sachant  $x$   
 $\Rightarrow$  la proba que  $Y$  prenne la valeur  $y$  sachant  
que  $X$  prend la valeur  $x$   
 $P(Y=y | X=x)$

La distribution conditionnelle de  $Y$  sachant  $X=x$   
 $P(Y=y | X=x) = \frac{P(X=x, Y=y)}{P(X=x)}$

$$P(\text{Gajet long} / \text{il pleut}) = \frac{0,15}{0,30} = 0,5$$

3 - Esperance conditionnelle de  $Y$  sachant  $X$   
 $\hookrightarrow$  moyenne de la distribution conditionnelle  
de  $Y$  sachant  $x$

$$4 \quad y_1 \dots y_n$$
$$E(Y | X=x) = \sum_{i=1}^n y_i P(Y=y_i | X=x)$$

## Loi des esperances itérées

$\hookrightarrow$  moyenne  $Y$  : moyenne pondérée de  
l'esperance conditionnelle de  $Y$  sachant  $X$ ,  
ponderation  $\leftarrow$  distribution de  $X$

## taille moyenne des adultes

moyenne pondérée de la moyenne des hommes et des  
femmes pondérée par la proportion d'homme et de  
femme

$$E(Y) = \sum_{i=1}^p P \text{ valeurs de } X \quad x_1, \dots, x_p$$
$$E(Y) = \sum_{i=1}^p E(Y | X=x_i) P(X=x_i) = E(E(Y|X))$$



### Variance conditionnelle:

$$\hookrightarrow \text{var}(Y | X=x) = \sum_{i=1}^n (y_i - E(Y|X=x))^2 P(Y=y_i | X=x)$$

### 4. Indépendances

X et Y sont indépendamment distribués ou indépendantes si connaître la valeur de l'une ne donne aucune information sur la valeur de l'autre.

$$P(Y=y | X=x) = P(Y=y)$$

$$X \text{ et } Y \text{ indépendantes} \Rightarrow P(X=x, Y=y) = P(X=x) \times P(Y=y)$$

### 5. Covariance et corrélation

Covariance: indique dans quelle mesure 2 variables varient ensemble

$$\text{cov}(X, Y) \text{ ou } \sigma_{X,Y} = E[(X - \mu_X)(Y - \mu_Y)]$$

$$= \sum_{i=1}^n \sum_{j=1}^n (x_{ij} - \mu_X)(y_{ij} - \mu_Y) P(X=x_j, Y=y_i)$$

### Corrélation

$$\text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X)\text{var}(Y)}} = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}$$

X et Y sont dites corrélées si  $\text{cov}(X, Y) \neq 0$   
 $-1 \leq \text{corr}(X, Y) \leq 1$



Si la moyenne conditionnelle de  $Y$  ne dépend pas de  $X$

$$\text{si } E(Y|X) = E(Y) = \mu_Y$$

$$\text{alors } \text{cov}(X, Y) = 0$$

$$\text{cov}(X, Y) = 0$$

$X$  et  $Y$  deux V.A de moyenne nulle

$$\text{cov}(X, Y) = E[(X - \mu_X)(Y - \mu_Y)] = E(XY)$$

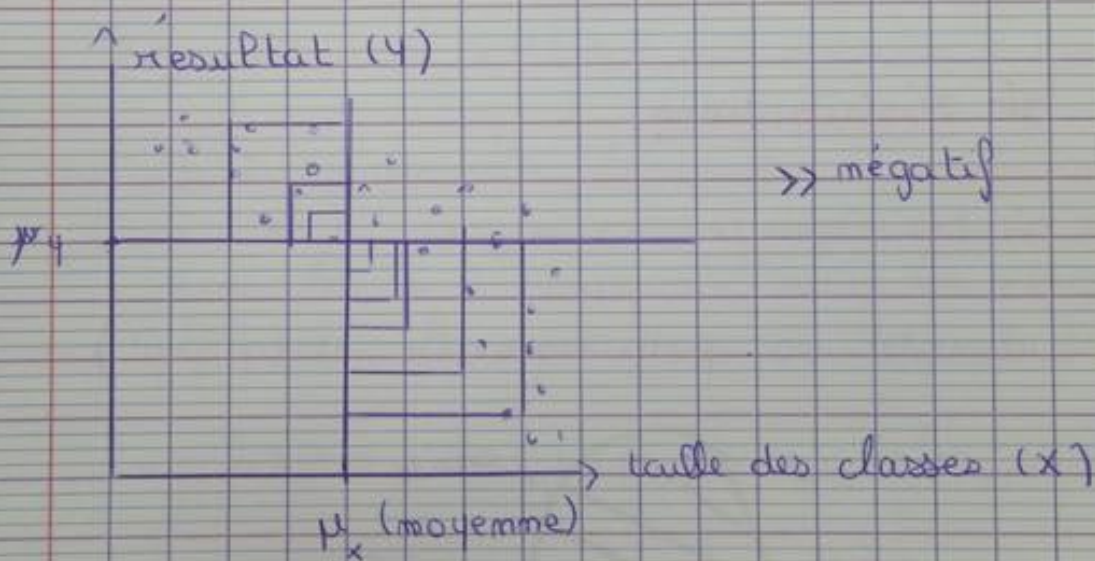
$$= E(E(XY|X))$$

$$= E(X E(Y|X))$$

$$= E(X E(Y))$$

$$= 0$$

donc si  $\text{cov}(X, Y) = 0 \Rightarrow \text{corr}(X, Y) = 0$





## 6. Moyenne, variance et somme des v.a

$$E(X+Y) = E(X) + E(Y) = \mu_x + \mu_y$$
$$\text{var}(X+Y) = \text{var}(X) + \text{var}(Y) + 2 \text{cov}(X, Y)$$

$$E(a+bX+cY) = a + b\mu_x + c\mu_y$$
$$\text{var}(a+bY) = b^2 \sigma_y^2$$

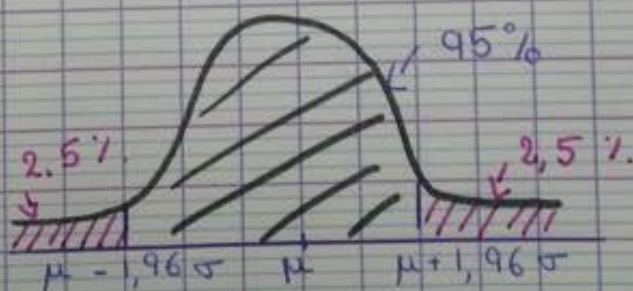
$$\text{var}(aX+bY) = a^2 \sigma_x^2 + b^2 \sigma_y^2 + 2ab \text{cov}(x, y)$$

$$E(Y^2) = \sigma_y^2 + \mu_y^2$$
$$\text{cov}(a+bX+cY, Y) = b \sigma_{x,y} + c \text{cov}$$
$$E(XY) = \sigma_{x,y} + \mu_x \mu_y$$

## IV Distributions

### 1. Loi Normale

densité de proba normale  $\mathcal{N}(\mu, \sigma^2)$



la normale centrée réduite  $\mathcal{N}(0, 1)$   $Z$

fonction de répartition  $\Phi$

$$P(Z \leq c) = \Phi(c) \quad c = \text{valeur}$$



$Y \rightsquigarrow \mathcal{D}(\mu, \sigma^2)$

$$P(Y \leq 2) = P\left(\frac{Y - \mu}{\sigma} \leq \frac{2 - \mu}{\sigma}\right) = \Phi\left(\frac{2 - \mu}{\sigma}\right)$$

$$Z = \frac{Y - \mu}{\sigma} \rightsquigarrow \mathcal{D}(0, 1)$$

$Y \rightsquigarrow \mathcal{D}(\mu, \sigma^2)$

alors  $Z = \frac{Y - \mu}{\sigma} \rightsquigarrow \mathcal{D}(0, 1)$

$$c_1 \text{ et } c_2 \quad c_1 < c_2 \quad \left\{ \begin{array}{l} d_1 = \frac{c_1 - \mu}{\sigma} \\ d_2 = \frac{c_2 - \mu}{\sigma} \end{array} \right.$$

$$P(Y \leq c_2) = P(Z \leq d_2) = \Phi(d_2)$$

$$P(Y \geq c_1) = P(Z \geq d_1) = 1 - \Phi(d_1)$$

$$P(c_1 \leq Y \leq c_2) = P(d_1 \leq Z \leq d_2) = \Phi(d_2) - \Phi(d_1)$$

distribution normale bivariée et multivariée

① si  $X$  et  $Y$  conjointement distribués ( $X$  et  $Y$  sont tous les 2) suivant une loi normale bivariée de covariance  $\sigma_{X,Y}$   $a, b$  2 constantes

$$aX + bY \rightsquigarrow \mathcal{D}(a\mu_X + b\mu_Y, a^2\sigma_X^2 + b^2\sigma_Y^2 + 2ab\sigma_{X,Y})$$

② si un ensemble de variables admettant une distribution normale multivariée la distribution marginale de chacune de ces variables est également normale

③ Si des variables avec une distribution normale multivariée ont des covariances nulles alors elles sont indépendantes



④ Si  $X$  et  $Y$  admettent une distribution normale bivariée, l'espérance conditionnelle de  $Y$  sachant  $X$  est linéaire en  $X$

$$E(Y|X) = a + bX$$

### e. loi du khi-deux

La distribution du khi-deux c'est la distribution de la somme des carrés de  $m$  v.a. normalement et indépendamment distribués de moyenne nulle et de variance unitaire

$Z_1 \rightsquigarrow \mathcal{D}(0, 1)$

$Z_2 \rightsquigarrow \mathcal{D}(0, 1)$  alors  $Z_1^2 + Z_2^2 + Z_3^2 \rightsquigarrow \chi^2_3$

$Z_3 \rightsquigarrow \mathcal{D}(0, 1)$

### 3. loi de student

La distribution de Student à  $m$  d.l. est définie par la distribution du ratio de 2 v.a. indépendantes

une v.a.  $\rightsquigarrow \mathcal{D}(0, 1)$

une v.a.  $\rightsquigarrow \chi^2_m$

$Z \rightsquigarrow \mathcal{D}(0, 1)$

$W \rightsquigarrow \chi^2_m$

alors  $\frac{Z}{\sqrt{\frac{W}{m}}}$   $\rightsquigarrow$  Student  $t$

Pour  $m > 30$ , cela se rapproche de la loi normale



#### 4 - loi de Fisher (distribution F)

La distribution F à m et n ddl notée  $F_{m,n}$  est définie par le ratio suivant :

$$\frac{\chi^2_{m|m}}{\chi^2_{n|n}}$$

si W et V 2 v.a. indépendamment distribuées  
 $\hookrightarrow \chi^2$  à m et n ddl

$$\text{alors } \frac{W|m}{V|n} \hookrightarrow F_{m,n}$$

#### Cas particulier

Si n est très grand alors  $F_{m,n}$  peut être approché par  $F_{m,+\infty}$  alors  $V_m$  correspond à la moyenne d'une infinité de  $\chi^2$  ; cette moyenne vaut 1. donc  $F_{m,+\infty}$  est la distribution d'une v.a. khi-deux à m ddl, divisé par m.



## Échantillonnage aléatoire et distribution d'échantillonnage des moyennes

### 1. Échantillonnage aléatoire

$n$  objets aléatoirement sélectionnés à partir d'une pop et dont chaque membre a la même proba d'être sélectionné

les  $n$  observations de l'échantillon sont notés  $y_1, \dots, y_n$ .

la sélection de manière aléatoire de membre de la pop conduit à des r.a.  $y_1 \rightarrow y_n$   
la r.a.  $y_i$  correspondant au  $i$ -ème objet aléatoirement tiré noté  $y_i$ .

- distribution commune à tous les  $y_i$
- équi probabilité de tirage
- $y_1, \dots, y_n$  sont indépendamment et identiquement distribués (i.i.d).

### 2. Moyenne et distribution d'échantillonnage

moyenne d'échantillonnage (empirique)  $\bar{y}$  d'un ensemble d'observation de  $y_1, \dots, y_n$ .

$$\bar{y} = \frac{1}{n} (y_1 + \dots + y_n) = \sum_{i=1}^n \frac{y_i}{n}$$

Distribution d'échantillonnage de  $\bar{y}$ .



$Y_1, Y_2, \dots, Y_m$  sont iid  $Y_i \sim \mu_Y \quad \sigma_Y^2$   
 $m=2$

$$E(Y_1 + Y_2) = \mu_Y \times 2$$
$$E\left(\frac{Y_1 + Y_2}{2}\right) = \mu_Y$$

$$\text{var}(Y_1 + Y_2) = 2 \sigma_Y^2$$

$$\text{var}\left(\frac{Y_1 + Y_2}{2}\right) = \frac{1}{2} \sigma_Y^2$$

$$\text{var}(\bar{y}) = \text{var}\left(\frac{1}{m} \sum_{i=1}^m Y_i\right)$$

$$= \frac{1}{m^2} \sum \text{var}(Y_i) + \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m \text{cov}(Y_i, Y_j)$$

$$\text{var}(\bar{y}) = \frac{1}{m} \sigma_Y^2$$

$\Rightarrow$  écart-type de  $\bar{y}$   $\frac{1}{\sqrt{m}} \sigma_Y$

Quand  $y$  est normalement distribué

$Y_1, \dots, Y_m$  sont iid de loi  $\mathcal{N}(\mu_Y, \sigma_Y^2)$   
 $\bar{y} \rightsquigarrow \mathcal{N}\left(\mu_Y, \frac{\sigma_Y^2}{m}\right)$



## VI approximation asymptotique de distribution d'échange

approche exacte = distribution exacte  
à distance finie

$$\frac{\bar{y} - \mu_y}{\sigma_{\bar{y}}} \xrightarrow{P} \mathcal{N}(0, 1)$$

La des grands nombres si  $Y_i$ ,  $i=1, \dots, n$   
des v.a. i.e.d.  
avec  $E(Y_i) = \mu_y$  et qu'il n'y a pas de valeurs  
extrêmes  $\text{var}(Y_i) = \sigma_y^2 < +\infty$   
alors  $\bar{y} \xrightarrow{P} \mu_y$

### TCL

$Y_1, \dots, Y_n$  sont i.i.d. de moyenne  $E(Y_i) = \mu_y$  et  
de variance  $\text{var}(Y_i) = \sigma_y^2$ ,  $0 < \sigma_y^2 < +\infty$   
alors quand  $n \rightarrow +\infty$ , la distribution de  
 $\frac{\bar{y} - \mu_y}{\sigma_{\bar{y}}}$  est arbitrairement bien approchée par  
une loi normale  $\mathcal{N}(0, 1)$ .



# Rappel Stat.

## 4) Estimation

### 1) Estimateurs et propriétés

Quelques de ceux-ci :

#### a) absence de biais

plusieurs tirages aléatoires et répétitifs d'échantillons

si à l'issue des répétitions, la moyenne de la distribution d'échantillonnage de l'estimateur est  $\mu_y$  alors c'est sans biais.

$\hat{\mu}_y$  est dit sans biais si  $E(\hat{\mu}_y) = \mu_y$ .

#### b) Convergence

quand  $n \rightarrow +\infty$ , la proba que  $\hat{\mu}_y$  soit très proche de  $\mu_y$  tend vers 1.

$$\hat{\mu}_y \xrightarrow{P} \mu_y$$

#### c) efficacité

Si 2 estimateurs  $\hat{\mu}_y$  et  $\tilde{\mu}_y$  sans biais  $\hat{\mu}_y$  sera + efficace que  $\tilde{\mu}_y$  si  $\text{var}(\hat{\mu}_y) < \text{var}(\tilde{\mu}_y)$

sans biais car  $E(\bar{Y}) = \mu_y$  convergent  
à des grands nombres  $\bar{Y} \xrightarrow{P} \mu_y$

efficace?  $\text{var}(\bar{Y}) = \frac{\sigma_y^2}{n}$  pour  $n \gg 1$



La moyenne empirique est le meilleur estimateur linéaire sans biais.

$\bar{y}$  : moyenne empirique

\* estimateur des moindres carrés

critère à minimiser  $\sum_{i=1}^n (y_i - m)^2$   
écart de prédiction

$$\frac{\partial}{\partial m} \sum (y_i - m)^2 = 0$$

$$\Leftrightarrow -2 \sum (y_i - m) = 0$$

$$\Leftrightarrow \sum y_i - n \cdot m = 0$$

$$\Leftrightarrow \sum y_i - n \times m = 0$$

$$\Leftrightarrow m = \frac{1}{n} \sum y_i = \bar{y}$$

II tests d'hypothèse

1) hyp nulle ou alternative  
 $H_0$  ou  $H_1$

$H_0$  " la moyenne théorique de  $Y$ ,  $E(Y)$ , prend une valeur spécifique  $\mu_{Y,0}$  "  
 $E(Y) = \mu_{Y,0}$

$H_1$   $E(Y) \neq \mu_{Y,0}$



## 2) La proba marginale de rejet

si  $\bar{y} \neq \mu_{y,0}$  2 raisons

→ la vraie moyenne n'est pas  $\mu_{y,0}$   
 $H_0$  est fautive

→ la vraie moyenne est bien  $\mu_{y,0}$  ( $H_0$  est vraie)  
mais aspect aléatoire dans  
l'échantillonnage.

La probabilité marginale de rejet

proba de tirer une valeur au moins aussi éloignée  
de la vraie valeur que la valeur estimée à  
partir de l'échantillon quand on suppose que  
 $H_0$  est vraie.

ex: Diplômés

salairé moyen: 22,64 € observé

$H_0$ : salairé la moyenné du salairé est de  
20 €

$$\mu_w = \underbrace{\mu_{w,0}}_{20 \text{ €}} \quad E(w) = 20 \text{ €}$$

Avec quelle proba, observera-t-on une  
valeur du salairé moyen  $> 22,64$   
ou  $\bar{w} < 17,36$  €

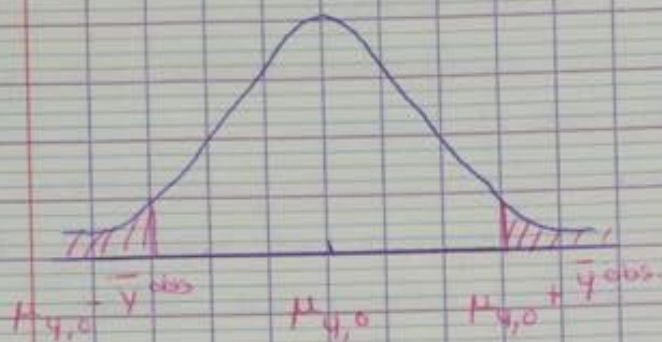
sachant que  $\mu_w = 20$  €

$\bar{y}_{obs}$  = valeur de la moyenné empirique  
effectivement calculé

$P_{H_0}$  = proba calculée sans  $H_0$  (ie  $E(Y_i) = \mu_{y,0}$ )



valeur  $p = P_{H_0} (|\bar{Y} - \mu_{y,0}| > |\bar{y}^{obs} - \mu_{y,0}|)$



$\bar{y}^{obs}$  pas compatible avec  $H_0$  si  $p$  est très faible

→ nécessaire de connaître la distribution d'échantillonnage de  $\bar{Y}$  sous  $H_0$ .

TCL de la normale est une bonne approximation de la distribution de  $\bar{Y}$

Sous  $H_0$ ,  $\bar{Y} \sim N(\mu_{y,0}, \sigma_{\bar{Y}}^2)$  avec  $\sigma_{\bar{Y}}^2 = \frac{\sigma_Y^2}{n}$

3. Calcul de la  $p$  valeur quand  $\sigma_{\bar{Y}}^2$  est connue

$\frac{\bar{Y} - \mu_{y,0}}{\sigma_{\bar{Y}}} \Rightarrow$  la version centrée réduite de  $\bar{Y}$



elle admet une distribution normale standard sous  $H_0$

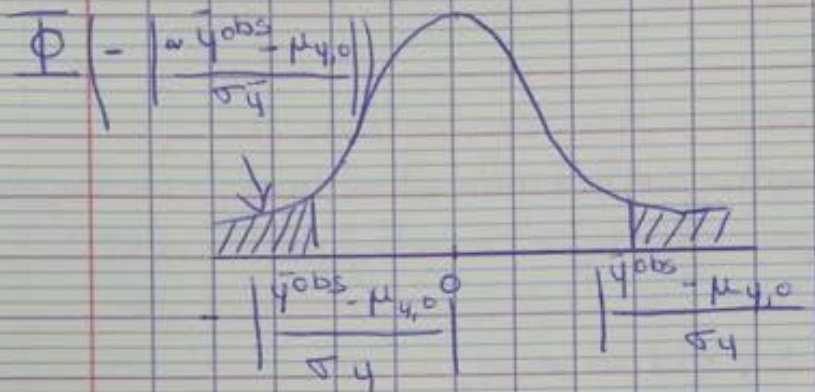
$p$  valeur = proba d'obtenir un  $\bar{Y}$  tq sous  $H_0$

$$P\left(\frac{|\bar{Y} - \mu_{y,0}|}{\sigma_{\bar{Y}}} > \frac{|\bar{y}^{obs} - \mu_{y,0}|}{\sigma_{\bar{Y}}}\right)$$



$$= P\left(|Z| > \left| \frac{\bar{y}^{obs} - \mu_{y,0}}{\sigma_{\bar{y}}} \right| \right) \quad \text{avec } Z \text{ c.v. } \mathcal{N}(0, 1)$$

soit  $\Phi$  la fonction de répartition de  $\mathcal{N}(0, 1)$   
 $P(d_1 < Z < d_2)$



valeur  $p = 2 \times \Phi\left(-\left|\frac{\bar{y}^{obs} - \mu_{y,0}}{\sigma_{\bar{y}}}\right|\right)$

$$\begin{aligned} P\left(|Z| > \underbrace{\left|\frac{\bar{y}^{obs} - \mu_{y,0}}{\sigma_{\bar{y}}}\right|}_c\right) &= 1 - P(|Z| < c) \\ &= 1 - P(-c < Z < c) \\ &= 1 - (\Phi(c) - \Phi(-c)) \\ &= \Phi(-c) + \Phi(-c) \\ &= 2 \times \Phi(-c) \end{aligned}$$

surface des queues de la distribution normale centrée réduite au delà de  $\pm \left|\frac{\bar{y}^{obs} - \mu_{y,0}}{\sigma_{\bar{y}}}\right|$



#### 4 - Variance empirique / écart type empirique $\left\{ \begin{array}{l} \text{erreur type} \\ \text{écart type estimé} \end{array} \right.$

##### Variance empirique

$$S_y^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$$

##### écart-type $s_{\bar{y}}$

La variance empirique est un estimateur convergent de la variance théorique.  $S_y^2 \xrightarrow{p} \sigma_y^2$

$y_1, \dots, y_n$  iid  $E(y^2) < \infty$

$$\begin{aligned} S_y^2 &= \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \mu_y + \mu_y - \bar{y})^2 \\ &= \frac{1}{n-1} \sum_{i=1}^n (y_i - \mu_y)^2 - \frac{2}{n-1} \sum_{i=1}^n (y_i - \mu_y)(\bar{y} - \mu_y) + \\ &\quad \frac{1}{n-1} \sum_{i=1}^n (\bar{y} - \mu_y)^2 \end{aligned}$$

$$\begin{aligned} &= \left( \frac{n}{n-1} \right) \left( \frac{1}{n} \sum_{i=1}^n (y_i - \mu_y)^2 \right) - \left( \frac{n}{n-1} \right) \left( \frac{2}{n} \sum_{i=1}^n (y_i - \mu_y)(\bar{y} - \mu_y) \right) + \\ &\quad \left( \frac{n}{n-1} \right) \left( \frac{1}{n} \sum_{i=1}^n (\mu_y - \bar{y})^2 \right) \end{aligned}$$

$$= \sigma_y^2$$

##### erreur type de $\bar{y}$

on a  $\sigma_{\bar{y}} = \frac{\sigma_y}{\sqrt{n}}$

( ) est un estimateur de  $\sigma_{\bar{y}}$

On l'appelle  $\bullet$  erreur type de  $\bar{y}$   
 $\bullet$  écart type estimé  $\cdot \hat{\sigma}_{\bar{y}}$   
 $\bullet$  "standard error"  $SE(\bar{y})$



$Y_1, \dots, Y_m$  i.i.d.  $\text{Bernoulli}(p)$

$$\text{var}(\bar{Y}) = \frac{p(1-p)}{m}$$

$$\hat{\sigma}_{\bar{Y}} = \text{SE}(\bar{Y}) = \sqrt{\frac{p(1-p)}{m}}$$

### 5. Calcul de la valeur p quand $\sigma_Y$ inconnue.

$S_Y^2$  est un estimateur convergent de  $\sigma_Y^2$   
pour calculer la valeur p, on remplace  
simplement  $\sigma_{\bar{Y}}$  par  $\hat{\sigma}_{\bar{Y}}$

$\sigma_Y$  inconnue,  $Y_1, \dots, Y_m$  i.i.d.:

$$\text{valeur } p = 2 \Phi\left(-\left|\frac{\bar{Y}^{\text{obs}} - \mu_{Y,0}}{\text{SE}(\bar{Y})}\right|\right) = 2\bar{\Phi}\left(-\left|\frac{\bar{Y}^{\text{obs}} - \mu_{Y,0}}{\hat{\sigma}_{\bar{Y}}}\right|\right)$$

### 6. la statistique t.

$$t = \frac{\bar{Y} - \mu_{Y,0}}{\text{SE}(\bar{Y})} \text{ suit une loi de Student}$$

distribution asymptotique suit une loi normale  $\mathcal{N}(0,1)$   
pour  $n$  grand (sous  $H_0$ )

Si on note  $t^{\text{obs}}$  la valeur obs de  $t$

$$t^{\text{obs}} = \frac{\bar{Y}^{\text{obs}} - \mu_{Y,0}}{\hat{\sigma}_{\bar{Y}}}$$

alors si  $n$  grand, la valeur  $p = 2\bar{\Phi}(-|t^{\text{obs}}|)$



### Exemple

$n = 200$  observés

$w$  balance hebdomadaire moyen

$$\bar{w}^{obs} = 22,64 \text{ €}$$

$$s_w = 18,14 \text{ €} = \sqrt{\frac{1}{200-1} \sum_{i=1}^{200} (w_i - \bar{w})^2}$$

écart-type estimé de  $\bar{w}$

$$\sigma_{\bar{w}} = \frac{s_w}{\sqrt{n}} = \frac{18,14}{\sqrt{200}} = 1,28$$

$$t^{obs} = \frac{22,64 - 20}{1,28} = 2,06$$

$$\begin{aligned} \text{valeur } p: & 2 \Phi(-|2,06|) \\ & = 2 \times 0,0193 \\ & = 0,039 = 3,9\% \end{aligned}$$

7. tester une hypothèse avec un seuil de signification préspécifié

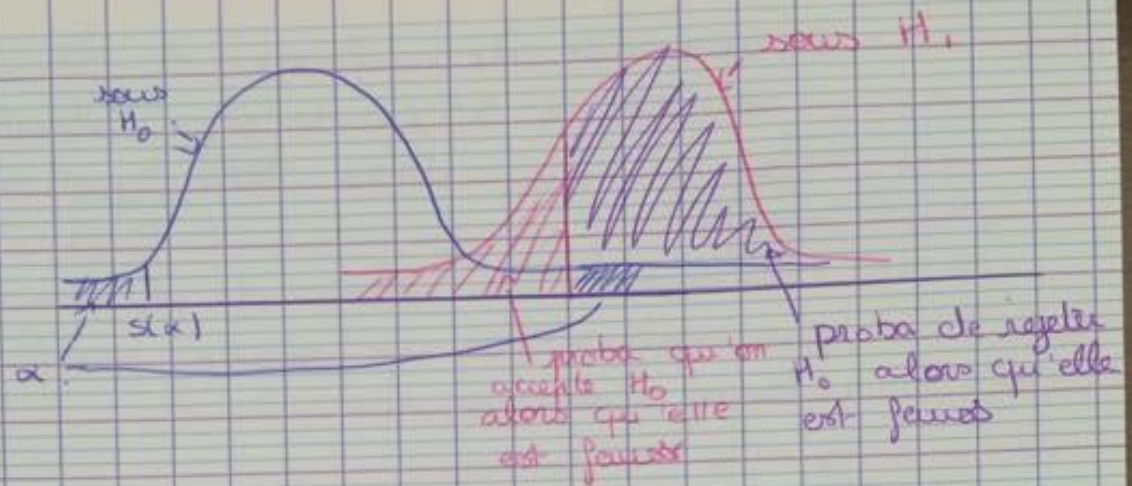
2 types d'erreur

- rejeter  $H_0$  alors qu'elle est vraie
  - erreur de première espèce  $\alpha(a)$
- accepter  $H_0$  alors qu'elle est fautive
  - erreur de deuxième espèce  $\beta(b)$

puissance d'un test

aptitude à rejeter  $H_0$  alors qu'elle est fautive





test d'hypothèse, en utilisant un seuil de signification fixé

règle { on rejette  $H_0$  si la proba marginale de rejet est inférieure à un certain seuil par ex 5%.

rejet de  $H_0$  si  $|t^{obs}| > 1,96$  (pour  $n$  grand).

Ainsi le risque d'erreur de première espèce est de 5%

• niveau de signification 5%

valeur critique du test bilatéral est 1,96

Région de rejet correspond à la valeur de la statistique  $t$  en dehors de  $\pm 1,96$



## 8. intervalle de confiance

Un IC à 95% de  $\mu_Y$  contient la vraie valeur de  $\mu_Y$  dans 95% des échantillons répétés.

L'ensemble des valeurs de  $\mu_Y$  qui ne sont pas rejetées par un test au seuil de signification de 5%

$$\left\{ \mu_Y = \left| \frac{Y - \mu_Y}{s_Y / \sqrt{n}} \right| \leq 1,96 \right\}$$

$$\Leftrightarrow \left\{ \mu_Y - 1,96 \leq \frac{\bar{Y} - \mu_Y}{s_Y / \sqrt{n}} \leq \bar{Y} - \mu_Y \leq 1,96 \frac{s_Y}{\sqrt{n}} \right\}$$

$$IC_{95\%} = \left[ \bar{Y} \pm 1,96 \frac{s_Y}{\sqrt{n}} \right] \Leftrightarrow \bar{Y} - 1,96 \frac{s_Y}{\sqrt{n}} \leq \mu_Y \leq \bar{Y} + 1,96 \frac{s_Y}{\sqrt{n}}$$