

Statistique Descriptive

Mesure de tendance centrale

	Population	Echantillon
Moyenne	$\mu = \frac{\sum x_i}{N}$	$\bar{x} = \frac{\sum x_i}{n}$
Médiane	$\frac{N}{2}$ <ul style="list-style-type: none"> → Nbre décimal: Mé = nbre entier suivant → Nbre entier: Mé = $\frac{(\frac{N}{2})^e + (\frac{N}{2} + 1)^e}{2}$ 	idem
Mode	c'est l'observation la plus fréquente	
Percentile	1 ^{er} quartile = 25° percentile ($\frac{N}{4} \times 1$) 2 ^{ème} quartile = 50° percentile (Mé) 3 ^{ème} quartile = 75° percentile ($\frac{N}{4} \times 3$) 85° percentile = $\frac{N}{100} \times 85$ R st des Vls constitue x% de l'étude	

Mesure de dispersion

Étendue	La \neq entre la 1 ^{ère} et la dernière valeur Étendue interquartile (ÉIQ) \Rightarrow entre 3 ^e et 1 ^{er} quartile calcul la force de liaison entre les Vls	
Variance	$\sigma^2 = \frac{\sum (x_i - \mu)^2}{N}$	$S^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$
Ecart type	$\sigma = \sqrt{\sigma^2}$	$S = \sqrt{S^2}$
Coefficient de variation	$\frac{\text{Ecart type}}{\text{moyenne}} \times 100$ x CV $\rightarrow 0 \Rightarrow$ série homogène x CV $\rightarrow 1 \Rightarrow$ série hétérogène x bonne homogénéité = CV < 15%	

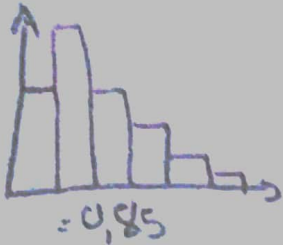
Mesure de tendance centrale et détection des valeurs singulières:

Forme d'une distribution: A partir de la forme de la degré d'asymétrie:

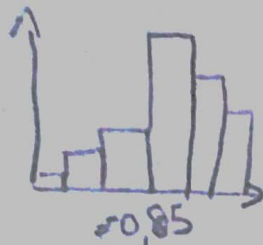
distribution on conclut que les données biaisées sont soit à droite soit à gauche.

D. Asy > 0

D. Asy < 0

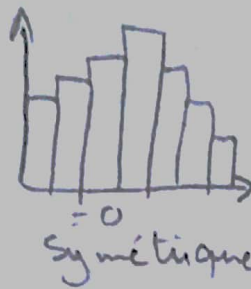


= 0,85

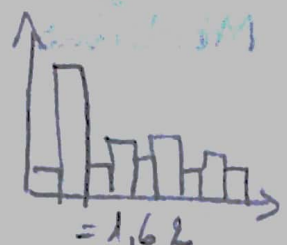


= 0,85

Moderatement asymétrique



= 0
Symétrique



= 1,62

Fortement asymétrique à droite

Degré d'asymétrie = $\gamma_1 = \frac{\mu_3}{\sigma^3}$
avec $\mu_3 = \sum_{i=1}^k \frac{1}{n} (x_i - \bar{x})^3$

Variable centrée réduite:

$$Z_i = \frac{x_i - \bar{x}}{s} \quad \text{ou} \quad Z = \frac{x_i - \mu}{\sigma}$$

x elle est calculée pour chaque valeur

x elle détermine / mesure la distance en nombre d'écart type entre x et la moyenne.

Théorème de Chebyshev:

Calculant le nombre de la variable x qui se situe entre [Y et V] de ni:

$$Y = \text{moyenne} - (Z \times \text{écart type})$$

$$V = \text{moyenne} + (Z \times \text{écart type})$$

⇒ Au moins $(1 - \frac{1}{Z^2})\%$ des observations se situent dans l'intervalle [Y, V]

Pour tous types de lois

Règle empirique:

Environ 68% des observations se situent dans $[\bar{x}-s; \bar{x}+s]$

Environ 95% " " " " " " $[\bar{x}-2s; \bar{x}+2s]$

Environ toutes les observations se situent dans $[\bar{x}-3s; \bar{x}+3s]$

Pour la loi normale

Boite à pâtes: Pour la détection des valeurs singulières

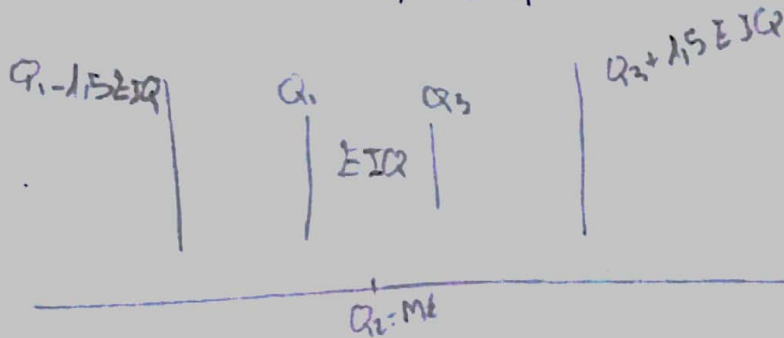
VL inhabituelle Er. d'enregistrement observation per
correctement incluse

Résumé en 5 chiffres: ⊕ petite valeur, ⊕ grande valeur, 1^{er}/2^e/3^e quartiles

Calcul Étendu interquartile: 3^eme Q - 1^{er} Q

Calcul: VL sup = $Q_3 + 1,5 \text{ EIQ}$

VL inf = $Q_1 - 1,5 \text{ EIQ}$



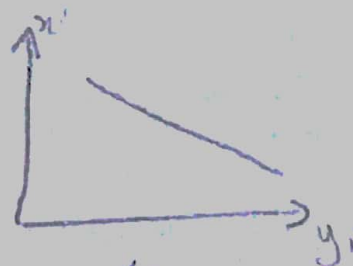
Statistique Bivarie

	Population	Echantillon
Covariance	$\sigma_{xy} = \frac{\sum (x_i - \mu_x)(y_i - \mu_y)}{N}$	$s_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1}$
Coefficient de corrélation	$\rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$	$r_{xy} = \frac{s_{xy}}{s_x s_y}$

La covariance est la mesure de la relation linéaire entre 2 variables



$cov > 0$
⇒ VL positives = Relat: linéaire ⊕



$cov < 0$
⇒ VL négatives = Relat: linéaire ⊖

Coefficient de corrélation: pas d'unité
forte rel. lin $\ominus \leftarrow -1 \leq r_{xy} \leq 1 \Rightarrow$ forte relation linéaire \oplus
 $\hookrightarrow 0 \Rightarrow$ très faible rel linéaire ou $\hat{=}$ introuvable

Méthode des moindres carrés:

Droite de régression de Y en X, $D_{Y/X} : y = ax + b$

$$\text{où } a = \frac{\sum xy}{(\sum x)^2} \quad \text{et } b = \bar{y} - a\bar{x}$$

Droite de régression de X en Y, $D_{X/Y} : x = a'y + b'$

$$\text{où } a' = \frac{\sum xy}{(\sum y)^2} \quad \text{et } b' = \bar{x} - a'\bar{y}$$

L'échantillonnage

Sondage = une étude pour un échantillon

recensement = une étude pour une population

Estimateur ponctuel = le fait de calculer un indicateur pour un échantillon afin de l'inférer à la population.

Méthodes d'échantillonnage:

1) Échantillonnages aléatoires: chaque individu a une probabilité connue et non nulle d'apparaître dans l'échantillon

* Échantillonnage aléatoire simple (EAS)

Population finie:

chaque échantillon à la n proba. d'être choisi parmi la population.

Population infinie:

x chaque élément est sélectionné de façon indépendante
x chaque élément sélectionné provient de la n population.

* Échantillonnage aléatoire stratifié:

- c'est la division de la pop en sous groupes / strates / tranches.

- chaque élément apparaît dans une seule strate

Ex: jeunes - vieux / hommes - femmes

- On choisit un EAS dans chaque strate

- Il y a une homogénéité des sous groupes

Allocation égale

Ex: Pop \Rightarrow jeunes: 1000
vieux: 1000

Ech \Rightarrow 50% jeunes
50% vieux

Allocation

proportionnelle à la taille

Ex: Pop \Rightarrow jeunes: 2000
vieux: 1000

Ech \Rightarrow 2/3 jeunes
1/3 vieux

* Échantillonnage par grappes:

c'est une miniature de la réalité de la population entière, le meilleur résultat est d'avoir des éléments hétérogènes.

Ex: quartiers, secteurs...

On choisit un EAS dans chaque grappe.

* Échantillonnage systématique:

Au lieu de l'EAS on utilisera l'ech. systématique

Ex: Pop: 5000 Ech: 50

1- Calcul d'un pas de l'échantillon: $K = N/n = 100$

2- On choisit un nombre entre 0 et K et on rajoute le pas (100): ex: 17

17, 107, 217, 317 ... 4917

Jusqu'à l'obtention de 50 ech.

2) Échantillonnages non aléatoires: quand il y a un avis / choix raisonné / sur étude.

* Échantillonnage de commodité:

On choisit un échantillon sans étude / volontaire: ex: prof ...

* Échantillonnage subjectif:

On cible le répondant.

Ex: journaliste choisit des personnes ...

* Boule de neige: on prend quelques individus, on leur demande de remplir le questionnaire et le passer à leurs amis.

* Échantillonnage par quota:

Ech. respecte les proportions connues pour les catégories de la pop
Ex: 200 Fac 50 ENCG 40 ENSA

• Vaut mieux utiliser l'échantillonnage aléatoire

Échantillonnage aléatoire simple (EAS)

Estimation ponctuelle:

* Distribution d'échantillonnage de \bar{x} :

a - l'espérance $E(\bar{x}) = \mu$ la moyenne de l'éch moyenne de la pop

b - Écart type

$\left\{ \begin{array}{l} \times \frac{n}{N} \leq 0,05 + \text{pop finie} \\ \times \text{pop infinie} \end{array} \right. \Rightarrow \sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$

$\times \text{pop finie} + \frac{n}{N} \geq 0,05 \Rightarrow \sigma_{\bar{x}} = \sqrt{\frac{N-n}{N-1}} \times \frac{\sigma}{\sqrt{n}}$

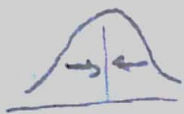
↓
Erreur type de la moyenne

c - La distribution tend vers une loi normale $N(\mu, \sigma)$

Pour \nearrow la proba, il faut:

\nearrow Échantillonnage

\searrow Écart type



Pour calculer la proba d'atteindre y écarts on reste entre ...

$$\mu - y \leq \bar{x} \leq \mu + y$$

↳ on la converge la loi normale centrée

$$\dots \leq \frac{\bar{x} - \mu}{\sigma_{\bar{x}}} \leq \dots$$

$$-v \leq \bar{z} \leq v$$

on voit la table de la loi normale

$$Ex: P(Z \leq -N) = P(Z \geq N) = 1 - P(Z < N)$$

$$P(Z \geq 2,3) = 1 - P(Z < 2,3)$$

$$P(1,2 \leq Z \leq 2,3) = P(Z \leq 2,3) - P(Z \leq 1,2)$$

$$P(-1 \leq Z \leq 1) = 2 \times P(Z \leq 1)$$

* Distribution d'échantillonnage de \bar{p} : $\bar{p} = \frac{x}{n}$

Il faut voir l'échantillon: elle suit la loi binomiale

a. Espérance: $E(\bar{p}) = p$
proportion éch. proportion pop.

b. Écart type:

$$\left. \begin{array}{l} \times \frac{n}{N} \leq 0,05 + \text{pop. finie} \\ \times \text{pop. infinie} \end{array} \right\} \Rightarrow \sigma_{\bar{p}} = \sqrt{\frac{p(1-p)}{n}}$$

$$\times \text{pop. finie} + \frac{n}{N} \geq 0,05 \Rightarrow \sigma_{\bar{p}} = \sqrt{\frac{N-n}{N-1}} \sqrt{\frac{p(1-p)}{n}}$$

↓
 Erreur type / Écart type proportion
Nbre d'essais

c. La distribution tend vers la loi binomiale $B(n; p)$ succès

qui peut se converger en loi normale

si $n \geq 30$; $np \geq 5$ et $nq \geq 5$ et $n(1-p) \geq 5$

$$B(n; p) \approx N(np; \sqrt{npq})$$

Pour calculer la propa

$$p-y \leq \bar{p} \leq p+y$$

↓ on la converge vers loi N.C.R

$$\dots \leq \frac{\bar{p} - p}{\sigma_{\bar{p}}} \leq \dots$$

$$u \leq Z \leq v$$

↪ on voit table de la loi normale

Estimation par intervalle:

Estimation de la moyenne d'une population:

$$\mu \in \left[\bar{x} - Z_{\alpha/2} \times \frac{s}{\sqrt{n}} ; \bar{x} + Z_{\alpha/2} \times \frac{s}{\sqrt{n}} \right] \quad \text{intervalle de confiance}$$
$$\mu \in [\bar{x} - E ; \bar{x} + E] \quad \text{↳ Marge d'erreur}$$

Si on a l'écart type de la population
⇒ loi normale centrée réduite

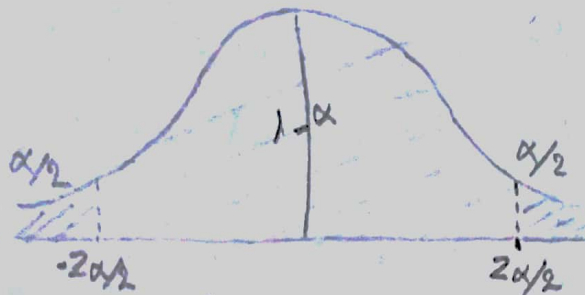
Si on n'a pas l'écart type de la population
⇒ loi de student

α → indice de l'incertitude - de risque

$1 - \alpha$ → indice de confiance

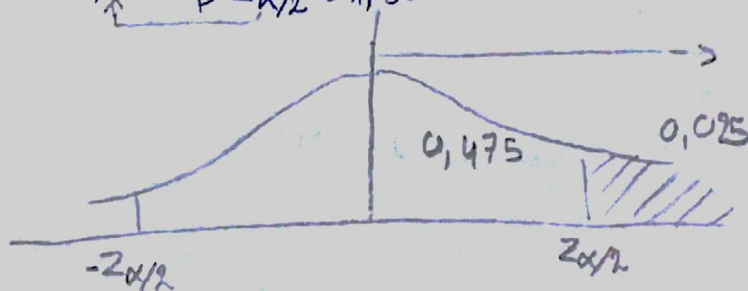
$\frac{\alpha}{2}$ → un demi-indice risque

$Z_{\alpha/2}$ → la valeur à partir de laquelle on est dans la zone d'incertitude.



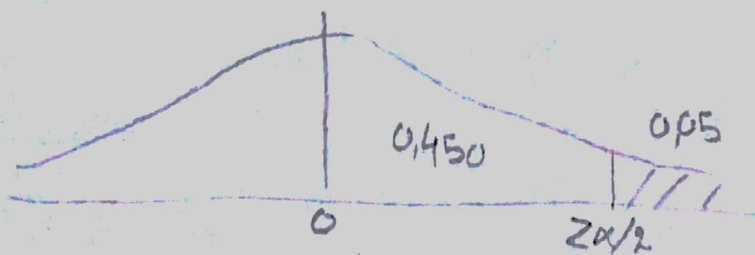
Selon L.N.C.R : on parle de $Z_{\alpha/2}$

vc risque 0,05 ⇒ $Z_{\alpha/2} = 1,96$



Donc on cherche les coordonnées de $P = 0,475 \Rightarrow 1,96$

vc risque 0,1 ⇒ $Z_{\alpha/2}$



$Z_{\alpha/2}$ est entre 1,64 et 1,65

on fait la moyenne $Z_{\alpha/2} = 1,645$

Selon L.S. on parle de $t_{\alpha/2}$

on cherche degré de liberté

$$\text{Nech. } -1 = n$$

On cherche la proba de α degré de liberté

$$\text{et } \underbrace{(0,10/0,05/0,025/0,01/0,005)}_{t_{\alpha/2}} ?$$

$$\text{Erreur} = t_{\alpha/2} \times \frac{s}{\sqrt{n}} = z_{\alpha/2} \times \frac{\sigma}{\sqrt{n}}$$

* Pour $z_{\alpha/2}$ l'écart type est celui de la pop.

* Pour $t_{\alpha/2}$ l'écart type est celui de l'échantillon

Estimation de la proportion d'une pop:

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

on n'utilise pas la loi de student

$$\text{Erreur} = z_{\alpha/2} \times \frac{\sqrt{\bar{p}(1-\bar{p})}}{\sqrt{n}}$$

Estimation de la taille de l'échantillon:

Moyenne:

$$E = z_{\alpha/2} \times \frac{\sigma}{\sqrt{n}}$$

$$\sqrt{n} = \frac{z_{\alpha/2} \times \sigma}{E}$$

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

Proportion:

$$E = z_{\alpha/2} \times \frac{\sqrt{\bar{p}(1-\bar{p})}}{\sqrt{n}}$$

$$\sqrt{n} = \frac{z_{\alpha/2} \times \sqrt{\bar{p}(1-\bar{p})}}{E}$$

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

$$\bar{p}(1-\bar{p}) \Rightarrow p^2 = 0,5$$

ou bien la valeur de p est basée sur des études anciennes

Test d'hypothèse pour Regression

$$H_0 \rightarrow \beta_1 = 0$$

$$H_1 \rightarrow \beta_1 \neq 0$$

on veut que notre modèle soit significatif
 \rightarrow Il y a une relat: entre x et y

Donc on cherche à rejeter H_0

Test student: Statistique du test: $t = \frac{b_1}{S_{b_1}}$

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

ou bien $E(b_1) = \beta_1$ Espérance

$$S_{b_1} \rightarrow S_{b_1} = \frac{S}{\sqrt{\sum (x_i - \bar{x})^2}}$$

Écart type estimé de b_1

α avec $S =$ Écart type de l'estimation

$$S = \sqrt{MC_{res}} = \sqrt{\frac{SC_{res}}{n-2}}$$

Erreur type de l'estimation

Rejet H_0 :

App. par vl critique: $t \leq -t_{\alpha/2}$ ou $t \geq t_{\alpha/2}$

\leftarrow vc ddl (n-2)

App. par vl p: $p \leq \alpha$

Intervalle de confiance pour β_1 :

$$\beta_1 \in [b_1 \pm t_{\alpha/2} \times S_{b_1}]$$

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

avec: $MC_{reg} = \frac{SC_{reg}}{1}$

$$MC_{res} = \frac{SC_{res}}{N-2}$$

\leftarrow 1 le nbre ddl qui est: nbre de vl indép.

Rejet H_0 :

App par vl p: $p < \alpha$

App par vl critique: $F > F_{\alpha}$ \rightarrow voir table fisher

	Source de la Variation	Somme des Carrés: SC	ddl	Moyenne des Carrés	F
Fisher	Régression	SC _{reg}	1	MC _{reg}	$F = \frac{MC_{reg}}{MC_{res}}$
Student	Résidu	SC _{res}	n-2	MC _{res}	
	Total	SCT	n-1	-	

Précision et estimation: ← utilise Student avec n-1 ddl

de la Vb moyenne de y $E(y)$

Ecart type estimé de \hat{y}_p :

$$S_{\hat{y}_p} = S \sqrt{\frac{1}{n} + \frac{(x_p - \bar{x})^2}{\sum (x_i - \bar{x})^2}}$$

Intervalle de confiance de $E(y_p)$

$$E(y) \in [\hat{y}_p \pm t_{\alpha/2} \times S_{\hat{y}_p}]$$

de la Vb individuelle de y:

Ecart type estimé de y_{ind}

$$S_{ind} = S \sqrt{1 + \frac{1}{n} + \frac{(x_p - \bar{x})^2}{\sum (x_i - \bar{x})^2}}$$

Intervalle de prévision de y:

$$y \in [\hat{y}_p \pm t_{\alpha/2} \times S_{ind}]$$

Régression linéaire simple

$$Y = \beta_1 x + \beta_0 + \varepsilon \leftarrow \text{erreur}$$

Droite linéaire

équation estimée de la régression

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

Méthode des moindres carrés:

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

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

Somme des carrés résidus:

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

Somme des carrés total:

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

$$SCT = SC_{res} + SC_{reg}$$

Somme des carrés de la régression:

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

Coefficient de détermination:

$$0 \leq r^2 = \frac{SC_{reg}}{SCT} \leq 1$$

Plus on s'approche de 1 \rightarrow \ominus erreurs

0 \rightarrow \oplus erreurs

\rightarrow X% est expliqué par la droite $\hat{y} = b_1 x + b_0$ (Rlt linéaire)

\rightarrow s'utilise ds le cas des rlt linéaires et non linéaires.

\rightarrow teste et évalue l'adéquation de l'équation estimée de la régression.

Coefficient de corrélation:

$$-1 \leq r = (\text{signe } b_1) \sqrt{r^2} \leq 1$$

\rightarrow teste et évalue l'intensité de la relati^o linéaire

\rightarrow utilisé juste ds le cas des relations linéaires.

$$\begin{array}{l} x_i \\ \bar{x} \\ y_i \\ \bar{y} \\ (x_i - \bar{x})(y_i - \bar{y}) \\ (x_i - \bar{x})^2 \\ (y_i - \bar{y})^2 \end{array}$$

$$\begin{array}{l} \hat{y} = b_1 x + b_0 \\ y_i - \hat{y} \\ (y_i - \hat{y})^2 \\ \hat{y}_i - \bar{y} \\ (\hat{y}_i - \bar{y})^2 \end{array}$$

$$\begin{array}{l} \bar{x} \\ \bar{y} \\ b_1 \\ b_0 \end{array}$$

⚠ Même si on obtient un coeff de détermination assez élevé il faut faire un test d'hypothèse de signification pour analyser le modèle supposé (y).

On suppose que:

$$E(\varepsilon) = 0 \text{ donc } E(y) = \beta_1 x + \beta_0$$

$\sigma^2(\varepsilon)$ est le même pour tous x

$\sum \varepsilon$ sont indépendants

ε suit la loi normale

on voit si $\beta_1 = 0$
y et x sont indépendants
 H_0

ou $\beta_1 \neq 0$
y et x sont dépendants
 H_1

Test de signification:

$$E(y) = \beta_0 + \beta_1 x$$

$$H_0: \beta_1 = 0$$

$$H_1: \beta_1 \neq 0$$

Si $\beta_1 = 0 \rightarrow E(y) = \beta_0 \rightarrow x$ et y ne sont pas liés

sinon: les 2 variables le sont (si $\beta_1 \neq 0$)

Estimation de σ^2 :

MCres fournit une estimation de σ^2 :

σ inconnu

$$s^2 = MCres = \frac{SCres}{n-2}$$

$$s = \sqrt{MCres}$$

Test d'hypothèse de la moyenne

Si σ est connu:

on utilise la loi normale CR

statistique du test:
$$Z = \frac{\bar{x} - \mu_0}{\frac{\sigma}{\sqrt{n}}}$$

$H_0: \mu \geq \mu_0$

$H_1: \mu < \mu_0$

Test unilatéral inf.

$H_0: \mu \leq \mu_0$

$H_1: \mu > \mu_0$

Test unilatéral sup

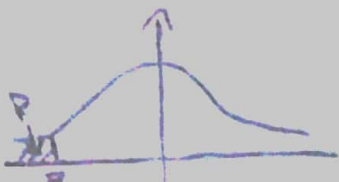
$H_0: \mu = \mu_0$

$H_1: \mu \neq \mu_0$

Test bilatéral

Pour rejeter H_0 :

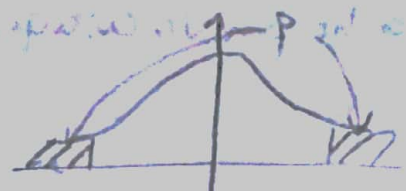
Approche par Vls p: Rejet de H_0 si $p \leq \alpha$ (p = seuil de signification: observé)



Test unilatéral inf $Z < 0$



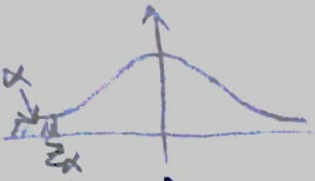
Test unilatéral sup $Z > 0$



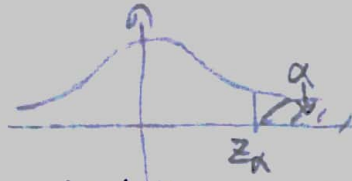
Test bilatéral

$P = 0,5 \cdot (Praba Z) \Rightarrow \alpha = \text{Risque}$

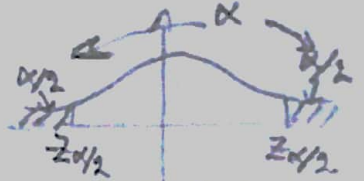
Approche par Vls critique z: Rejet de H_0 si $Z \leq Z_\alpha$



Test unilatéral inf $Z \leq Z_\alpha$



Test unilatéral sup $Z \geq Z_\alpha$



Test bilatéral $Z \leq -Z_{\alpha/2}$ ou $Z \geq Z_{\alpha/2}$

Z = statistique du test

Z_α = Z du risque complet

$Z_{\alpha/2}$ = Z de 1/2 du risque

Si σ est inconnu:

on utilise loi de student degré de liberté $(n-1)$

$H_0: \mu \geq \mu_0$

$H_1: \mu < \mu_0$

Test unilatéral sup

Pour rejeter H_0 :

$H_0: \mu \leq \mu_0$

$H_1: \mu > \mu_0$

Test unilatéral inf

$H_0: \mu = \mu_0$

$H_1: \mu \neq \mu_0$

Test bilatéral

Approche par Vls p: $p \leq \alpha$

Approche par Vls critique:

$t \leq -t_\alpha$ $t \geq t_\alpha$ $t \leq -t_{\alpha/2}$
ou $t \geq t_{\alpha/2}$

t : statistique du test

t_α : statistique t du risque complet

$t_{\alpha/2}$: statistique t de 1/2 du risque

Test d'hypothèse de la proportion

Statistique du test
$$Z = \frac{\bar{P} - P_0}{\sqrt{\frac{P_0(1-P_0)}{n}}}$$

$H_0: P \geq P_0$

$H_1: P < P_0$

Test unilatéral inf.

$H_0: P \leq P_0$

$H_1: P > P_0$

Test unilatéral sup

$H_0: P = P_0$

$H_1: P \neq P_0$

Test bilatéral

Pour rejeter H_0 :

Approche par vl p:

$P < \alpha$

Approche par vl critique:

$Z \leq -Z_\alpha$

$Z \geq Z_\alpha$

$Z \leq -Z_{\alpha/2}$

ou

$Z \geq Z_{\alpha/2}$