

Statistique

I. Statistique Descriptive :

Mesure de Tendance centrale et De Dispersion.

- ↳ Collecte Des Données
- ↳ Traitement Des données collectées
- ↳ Interprétation des données.

Effectif : c'est le Nombre total des élt's constituant cette population

Fréquence : c'est le nbre d'individus possédant ce caractère divisé par l'effectif total de la population : N

	Echantillon	Population
Moyenne	\bar{x}	μ
Variance	s^2	σ^2
Ecart type	s	σ
Covariance	s_{xy}	σ_{xy}
corrélat	r_{xy}	ρ_{xy}

Echantillon

Population

Moyenne :

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

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

Médiane : Il s'agit de la valeur centrale de l'ensemble des données, classés en ordre croissant.

Mode : Défini comme la valeur de l'observat la plus fréquente.

Percentile : au moins pour cent des Observat ont cette valeur

quartile : 25^e, 50^e, 75^e ⇒ Mesure de Tendance centrale

Mesure De Dispersion :

Étendue : égale à la Différence entre la plus grande et la plus petite valeurs.

Étendue Interquartile : (EIQ) : égale à la Différence entre le 3^eme et la 1^{ère} quartile :

Variance : basé sur les écarts au carré des Observat par rapport à la moyenne :

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

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

Écart type :

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

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

Coefficient de variation: Mesure de Dispersion relative, égale au rapport de l'écart type à la moyenne, multiplié par 100,

$$\frac{\text{Ecart type}}{\text{Moyenne}} \times 100$$

Valeur singulière: Observatⁿ anormalement grande ou petite.

Degré d'asymétrie: + Des données biaisées à gauche sont

caractérisées par un degré d'asy négatif

+ Des données bi à dte sont " " " positive.

Variable centrée réduite z : $z_i = \frac{x_i - \bar{x}}{s}$

$$\begin{aligned} \text{Variance} &= 4 + 100 + 4 + 4 + 144 \\ &= \frac{256}{4} = 64 \end{aligned}$$

$$\text{Ecart type} = \sqrt{s} = 8$$

	s	Ecart par rap à la Moyenne	v. de la vari centré réduite
46		0,25	0,25
54		10/8	1,25
42		-2/8	-0,25
46		2/8	0,25
32		12/8	-1,50

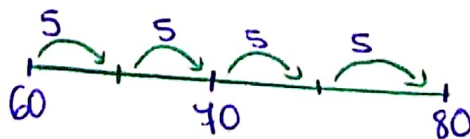
Théorème de Chebyshev

Théo utilisé pour réduire le pourcentage d'Observatⁿ qui se situe dans un intervalle de z écart type de part et d'autre de la Moyenne.

\Rightarrow Au moins $(1 - \frac{1}{z^2})$ des Observatⁿ doivent se situer au plus à $|z|$ écart type de part et d'autres de la Moyenne (càd dans l'intervalle $[\bar{x} - zs, \bar{x} + zs]$) avec $z > 1 \gg$.

Exple e supposons que la Moyenne des notes de 100 Étudiants de l'ENGT soit égale à 70 et que l'écart type = 5,

1/ combien d'étudiants ont obtenu une note entre 60 et 80?



$$[70 - 2s; 70 + 2s] \Rightarrow 1 - \frac{1}{2^2} = 75\%$$

75% des étudiants ont obtenus une note entre 60 et 80.

$\Rightarrow \Delta$ On peut pas l'utiliser $1s$;

Règle Empirique Règle qui donne le % d'Observat situées dans des Intervalles 1, 2, 3 écarts type autour de la Moyenne, pour une Distribut en forme de cloche (Distribut normale).



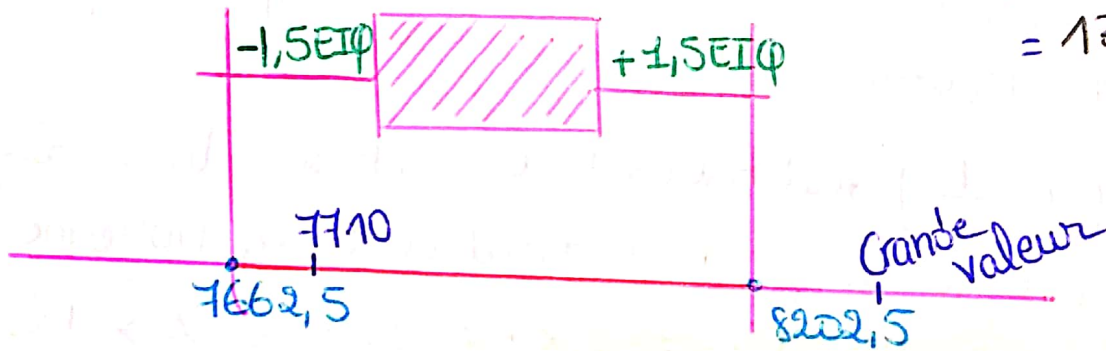
- Environ 68% des Observat se situe dans $[\bar{x} - s; \bar{x} + s]$
- Environ 95% des Observat se situe dans $[\bar{x} - 2s; \bar{x} + 2s]$
- presque toutes les Observat se situent dans $[\bar{x} - 3s; \bar{x} + 3s]$

Résumé en 5 chiffres Technique d'analyse Exploratoire des données qui utilise 5 chiffres pour résumer les données, la plus Grande valeur, le 1^{er} quartile, la médiane, le 3^e quartile et La plus Grande valeur, Exple

7710 - 7755 7850 7880 7880 7890 7920 7940 7950 8050
 - 8130 8325

1/ 7710 2/ $Q_1 = 7865$ 3/ $Q_3 = 7905$
 4/ 8000 5/ 8325

$$EIQ = Q_3 - Q_1 = 135$$



$$Q_1 - 1,5EIQ = 7662,25$$

$$Q_3 + 1,5EIQ = 8202,5$$

Donc 8325 est une Valeur Singulière

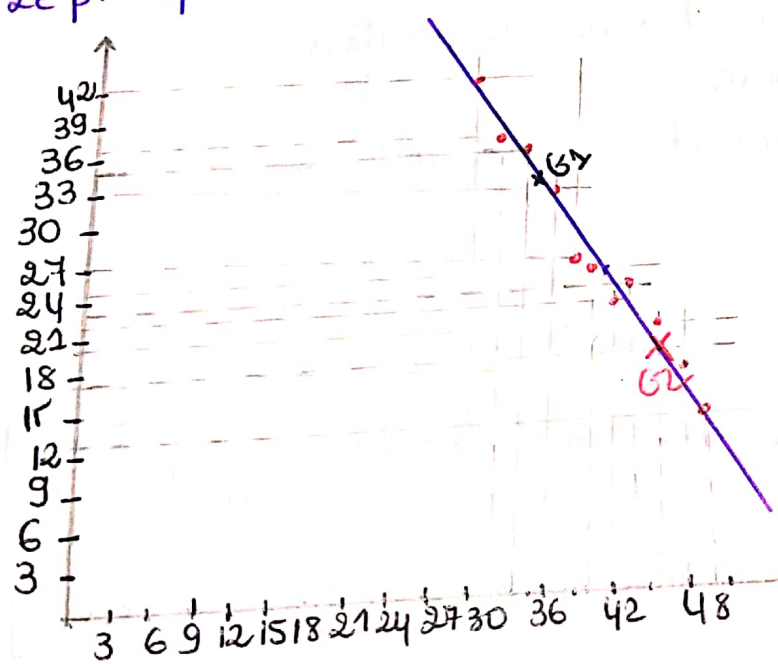
II - Statistique Bivariée :

Nuage des pts :

Exple : Une chaîne franchisé Désire Déterminer le prix idéal, puis il a proposé les prix Différents ... à relever le % des clls interressé par cette prestation.

Prix proposé	30	32	34	36	38	40	42	44	46	48	50
% des clls interressé	42	37	35	33	27	25	22	23	20	17	13

- 1/ Présenter Graphiquement puis tracer la dté d'ajustement.
- 2/ Le prix pour avoir 30% des clls interressés ?



Mthode De Mayer

$$G_1 = (35, 33,16)$$

$$G_2 = (46, 19)$$

$$y = ax + b$$

$$a = \frac{y_{G_2} - y_{G_1}}{x_{G_2} - x_{G_1}} = \frac{19 - 33,16}{46 - 35} = -1,28$$

$$b = y - ax \Rightarrow 33,16 - (-1,28 \times 35) = 78,28$$

$$y \Rightarrow -1,28x + 78,28$$

d'où $x = \frac{30 - 78,28}{-1,28} = 37,69 \Rightarrow$ Le prix pour avoir 30% des clls interressé est : 37,69 dh.

Mesure par la Covariance :

Mesure de la relat linéaire entre deux variables.

$$\text{Covariance populat: } \sigma_{ny} = \frac{\sum (x_i - \mu_x)(y_i - \mu_y)}{N}$$

$$\text{Covariance Échantillon: } S_{ny} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

Coefficient de corrélat : Mesure de la relat linéaire entre deux variables, dont les valeurs sont comprises entre -1 et +1

$$r_{xy} = \frac{S_{ny}}{S_n S_y} \quad \text{ou} \quad \rho_{ny} = \frac{\sigma_{ny}}{\sigma_n \sigma_y}$$

- proche de +1 \Rightarrow forte relat linéaire positive
- " -1 \Rightarrow forte relat " négative
- proche de 0 \Rightarrow absence de relat linéaire

Estimation Ponctuelle :

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

Estimation de la Moyenne \nearrow σ connu $\mu \in [\bar{x} - z_{\alpha/2} \times \frac{\sigma}{\sqrt{n}}, \bar{x} + \dots]$
 \searrow σ inconnu $\mu \in [\bar{x} - t_{\alpha/2} \times \frac{\sigma}{\sqrt{n}}, \bar{x} + \dots]$

Distribution D'échantillonnage par \bar{p} :

$$\text{population finie} \rightarrow \sqrt{\frac{N-n}{N-1}} \times \sqrt{\frac{P(1-P)}{n}}$$

$$\text{population infinie} \rightarrow \sqrt{\frac{P(1-P)}{n}}$$

$$\rightarrow \text{Si } \frac{N}{n} < 0,05$$

$$\text{La Marge D'erreur : } m = \frac{(z_{\alpha/2})^2 \times \sigma^2}{E^2}$$

Chapitre III = Régression linéaire simple :

$$y = \beta_0 + \beta_1 x + \varepsilon$$

y : variable à expliquer

x : variable explicative (indépendante).

β_0 et β_1 sont les paramètres du modèle

ε est une aléatoire appelée : terme d'erreur, ce terme prend en compte la variabilité de y qui n'est pas expliquée par la relation linéaire entre x et y , ce terme regroupe 3 erreurs :

→ E. de spécification : le fait que la seule variable explicative n'est pas suffisante pour rendre compte de la totalité du phénomène expliqué.

→ E. de mesure : les données ne représentent pas exactement le phénomène.

→ E. fluctuation d'échantillonnage : d'un échantillon à l'autre les observations, et donc les estimations sont légèrement différentes.

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

En pratique la valeur du paramètre n'est pas connue et doit être estimée en utilisant les données d'un échantillon :

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

Méthode Des Moindres Carrées :

est une procédure qui permet d'utiliser les données de l'échantillon pour estimer l'équation de la régression (b_0 et b_1). Elle consiste à minimiser la somme des écarts au carré :

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

$$\text{Or: } b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \quad \text{et} \quad b_0 = \bar{y} - b_1 \bar{x}$$

Exple : Restaurant :

$$r = \frac{\text{cov}}{v}$$

x_i	y_i	$x_i - \bar{x}$	$y_i - \bar{y}$	$\hat{y}_i = 60 + 5x_i$	$(y_i - \hat{y}_i)$	$(y_i - \hat{y}_i)^2$	$(x_i - \bar{x})^2$
2	58	-12	-72	70	-12	144	144
6	105	-8	-25	90	15	225	64
8	88	-6	-42	100	-12	144	36
8	148	-6	-12	100	18	324	36
12	147	-2	-13	120	-3	9	4
16	137	2	7	140	-3	9	4
20	157	6	27	160	-3	9	36
20	169	6	39	160	9	81	36
22	149	8	19	170	-21	441	64
26	202	12	72	190	12	144	144

$$\bar{x} = \frac{140}{40} = 3.5 \quad \text{et} \quad \bar{y} = 130 \quad ; \quad b_1 = 5 \quad ; \quad b_0 = 60$$

$$\text{donc } \hat{y}_i = 60 + 5x_i$$

$$+ \text{Scres} = \sum (y_i - \hat{y}_i)^2 \rightarrow \text{erreur}$$

$$+ \text{Screg} = \sum (\hat{y}_i - \bar{y})^2 \rightarrow \text{la partie expliquée}$$

$$\left\{ \text{SCT} = \sum (y_i - \bar{y})^2 \right.$$

$$\text{SCT} = \text{Scres} + \text{Screg}$$

$$\text{Coeff de détermination : } \frac{\text{Screg}}{\text{SCT}} = 1 \Rightarrow \text{Meilleure situation}$$

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

$$\text{On a } r = \frac{\text{Screg}}{\text{SCT}} = \frac{14200}{15700} = 0,9027$$

en d'autre terme 90,27% de la variation des ventes mensuelle peut s'expliquer par la relation linéaire de la pop et ventes, Une telle adéquation de l'équation estimée de la reg est satisfaisante

V. hypothèse du Modèle

H_0 = Hypothèse nulle \rightarrow il faut l'a rejeter.

H_a = hyp alternative

\rightarrow test uni infé

$$H_0 = \mu \geq 2$$

$$H_a = \mu < 2$$

test uni supé

$$H_0 = \mu \leq 2$$

$$H_a = \mu > 2$$

Test Bilaterale

$$H_0 = \mu = 2$$

$$H_a = \mu \neq 2$$

Seuil de signification α est la probabilité de faire une erreur de 1^{ère} espèce lorsque l'hypothèse nulle est vraie et satisfaite avec égalité.

$$Z = \frac{\bar{x} - \mu_0}{\sigma / \sqrt{n}}$$

Rejet de H_0 si

1 $Z \leq -z_{\alpha}$

2 $Z \geq z_{\alpha}$

3 $Z \leq -z_{\alpha/2}$ ou $Z \geq z_{\alpha/2}$

I. Estimation du σ^2

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

Espérance : $E(b_1) = \beta_1$ et $\sigma_{b_1} = \frac{\sigma}{\sqrt{\sum (x_i - \bar{x})^2}}$

$$MC_{reg} = \frac{SC_{reg}}{\text{Nbr de V. Indép}}$$

$F = \frac{MC_{reg}}{MC_{res}}$; suit une loi de Fisher avec 1 ddl en numé et $n-2$ ddl en déno