

# Statistiques

## Généralités :

1. Population : univers statistique (représenté par  $\Omega$ ) : ensembles étudiés par la statistique (45 clowns ont participé à une étude sur leur âge et leur taille) -> des clowns
2. Individus : unités statistiques (représentés par  $\omega$ ) : éléments des ensembles (45 clowns ont participé à une étude sur leur âge et leur taille) -> 45
3. Recensement : étude de tous les individus d'une population finie
4. Sondage : observation d'une partie de la population
5. Echantillon (représenté par  $E$ ) : le sous ensemble étudié, lors d'un sondage ( $E \subset \Omega$ )
6. Observations : nbr d'unités statistiques x nbr de variables (45 clowns ont participé à une étude sur leur âge et leur taille) ->  $2 \times 45 = 90$
7. Variables, caractères : ensemble de caractéristiques décrivant les individus d'une population. (45 clowns ont participé à une étude sur leur âge et leur taille) -> âge, taille

## Types de variables statistiques :

1. Quantitatives (nombres, chiffres, etc...) -> numériques **ATTENTION** (date de naissance)
  - Continues, assimilées : mesure, **quantité** (poids, taille, R = nombres réels)
  - Discrètes : comptage (années, N = entiers naturels, 0 1 2 3 4 5 6 ...)
2. Qualitatives (mots, expressions, etc...)
  - Nominale : aucun ordre (nom, prénom, code postal, adresse, tél, email ...)
  - Ordinale : qui peut se classer (un peu, bcp, moyen, entre 400 et 500, plus de 500, moins de 400)

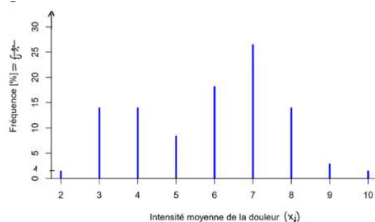
Ex : Votre genre de BD préféré : (a) humour, (b), (c), (d), (e), (f) classer de 1-6 (1=préférée)

- Préférence : (a) rang attribué aux BD d'humour, modalités : 1/2/3/4/5/6, type : numérique (b)... (c)... (d)...

## Diagramme en bâtons (variable quantitative/numérique discrète)

1. Tableau des scores bruts (tous les chiffres mélangés sous forme de tableau)
2. Tableau des scores ordonnés (mettre les chiffres dans l'ordre)
3. Tableau des données regroupées

- ⇒ J = simplement les données
- ⇒  $X_j$  = intensité de la douleur
- ⇒ n = la somme des effectifs (72)
- ⇒ f = fréquence
- ⇒  $n_j$  = effectifs (nbr de fois)



Scores bruts										Scores ordonnés									
8	4	7	7	3	7	3	7	9	6	2	3	3	3	3	3	3	3	3	3
6	3	4	4	3	3	10	7	4	7	3	4	4	4	4	4	4	4	4	4
6	3	7	4	5	7	7	6	8	6	4	5	5	5	5	5	5	5	5	5
7	6	8	7	9	8	4	6	7	8	6	6	6	6	6	6	6	6	6	6
2	6	5	5	7	3	6	4	7		7	7	7	7	7	7	7	7	7	7
8	7	5	6	5	6	5	3	4		7	7	7	7	7	7	7	7	7	7
7	8	6	7	4	7	8	3	7	8	8	8	8	8	8	8	8	8	8	8
4	8									9	10								

	A	B	C	D	E	F	G	H	I	J
1	j	1	2	3	4	5	6	7	8	9
2	$x_j$	2	3	4	5	6	7	8	9	10
3	$n_j$	1	10	10	6	13	19	10	2	1
4	$f_j$	1/72	10/72	10/72	6/72	13/72	19/72	10/72	2/72	1/72
5	$f_j$ (%)	1,4	13,9	13,9	8,3	18,1	26,4	13,9	2,8	1,4

4. Jamovi : entrer les données -> analyse -> exploration -> descriptives -> sélectionner la variable -> cocher Frequency tables -> cocher N et Missing

## Histogramme (variable quantitative continue/assimilée)

1. Estimation du nbr de classes K  
 $K = \lceil \log_2 n \rceil + 1 = \lceil \log_2 200 \rceil + 1 = 7,644 + 1 = 8 + 1 = 9$ 

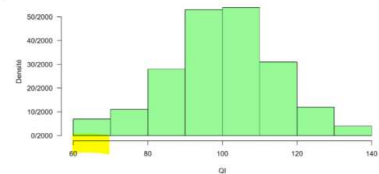
*(à arrondir)*
2. Calcul de l'empan  
 $Empan = \max(X) - \min(X) = 140 - 62 = 78$
3. Calcul de la largeur des classes  $a_j$   
 $\frac{Empan}{K} = 8,667$ 

*(à arrondir)*
4. Ajustement de la largeur des classes  $a_j$   
 Comme  $5 < 8,667 < 10$ , nous donnerons finalement à chaque classe  $j$  une largeur égale à  $a_j = 10$ .
5. Tableau des données regroupées en classe
6. Historigramme
7. Historigramme au tableau des données

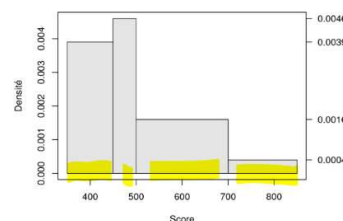
i	1	2	3	4
$[e_{i-1}; e_i[$	[350; 450[	[450; 500[	[500; 700[	[700; 850[
$c_i$	400	475	600	775
$a_i$	100	50	200	150
$d_i$	0,0039	0,0046	0,0016	0,0004
$f_i$	0,39	0,23	0,32	0,06
$n_i$	390	230	320	60

j	1	2	3	4
$[e_{j-1}; e_j[$	[60; 70[	[70; 80[	[80; 90[	[90; 100[
$a_j$	10	10	10	10
$n_j$	7	11	28	53
$f_j$	0,0350	0,0550	0,1400	0,2650
$d_j$	0,0035	0,0055	0,0140	0,0265

j	5	6	7	8
$[e_{j-1}; e_j[$	[100; 110[	[110; 120[	[120; 130[	[130; 140[
$a_j$	10	10	10	10
$n_j$	51	31	12	4
$f_j$	0,2700	0,1550	0,0600	0,0200
$d_j$	0,0270	0,0155	0,0060	0,0020



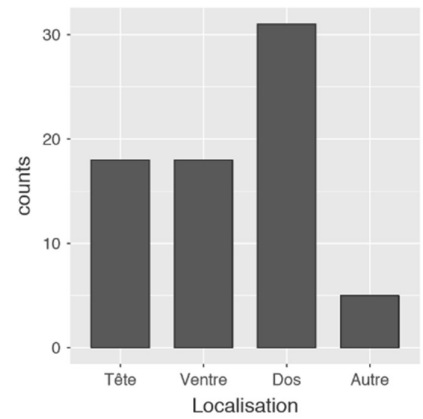
Reconstituez, à partir de l'histogramme représenté ci-dessous, le tableau des données regroupées en classes ( $n = 1000$ ).



## Diagramme en tuyau d'orgue (variable qualitative)

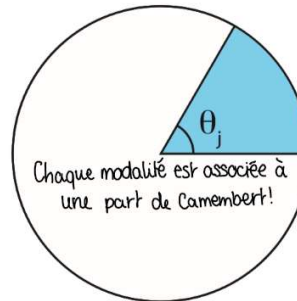
1. Tableau des effectifs
2. Jamovi : analyses -> exploration -> descriptives -> plots -> bar plot
  - La largeur des tuyau n'a pas d'importance (similaires)

localisation	tête <sup>j=1</sup>	ventre <sup>j=2</sup>	dos <sup>j=3</sup>	autre <sup>j=4</sup>
$n_j$	18	18	31	5
$f_j = \frac{n_j}{n}$	25%	25%	43,1%	6,9%



## Diagramme en Camembert

	Sexe		
	Fille	Garçon	Total
$n_j$	59	13	72
$f_j$	81.9%	18.1%	100%
$\theta_j$	295°	65°	360°



$$\theta_j = 360^\circ \times f_j$$

$$\theta_{\text{fille}} = 360^\circ \times 0,819 = 295^\circ$$

$$\theta_{\text{garçon}} = 360^\circ \times 0,181 = 65^\circ$$

## La Médiane (M)

- JAMOV I : Analyses -> explorations -> descriptives -> statistiques -> median -> cocher N
- Mettre les résultats dans l'ordre et prendre le résultat du centre (si impair)
- Attention si pair : prendre les deux au centre et faire la moyenne des deux  $\frac{X_k + X_{k+1}}{2}$

## La Moyenne (X) Mean

- Multiplication de tous les scores entre eux, puis diviser par le nbr total de scores multipliés
- JAMOV I : Analyses -> explorations -> descriptives -> statistiques -> mean -> cocher N

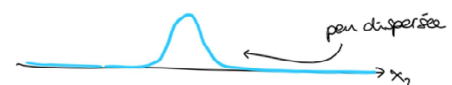
## Le Mode (Mo)

- La valeur la + fréquente de la recherche -> exemple :
- Certaines distributions peuvent avoir plusieurs modes, ou pas de mode du tout (tous pareils)

Score	1	2	3	4	5	6	7
Effectif	1	1	2	3	6	5	1

parce que c'est la qu'il y en a le +

## Caractéristiques de dispersion



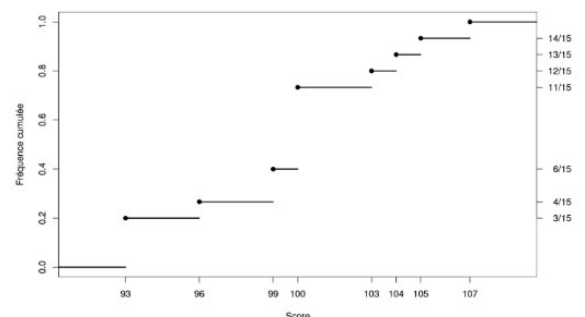
## Résumé

	mode	médiane	moyenne
nominale			
ordinaire	X	X	
numérique	X	X	X

## Fonction Quantile (à partir d'une fonction de répartition et de n)

1. Fonction de répartition F(x)
  - Chaque point =  $x_j$
  - Espace (vertical) entre chaque trait =  $n_j$

j	1	2	3	4	5	6	7	8
$x_j$	93	96	99	100	103	104	105	107
$n_j$	3	1	2	5	1	1	1	1
$f_j$	3/15	1/15	2/15	5/15	1/15	1/15	1/15	1/15
$F_j$	3/15	4/15	6/15	11/15	12/15	13/15	14/15	15/15

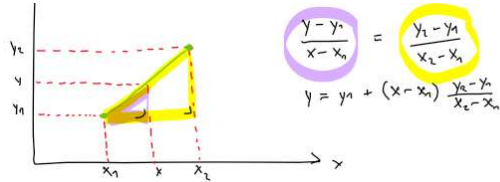


## 2. Fonction Quantile $F^{-1}(\alpha)$

- Déterminer les quartiles :  $Q1 = 0,25$  /  $Q2 = 0,5$  = Médiane /  $Q3 = 0,75$  -> **divided per 4**
- Déterminer les déciles :  $D1 = 0,1$  /  $0,2$  /  $0,3$  /  $0,4$  /  $0,5$  /  $0,6$  /  $0,7$  /  $0,8$  /  $0,9$  -> **Divided per 10**
- Déterminer les centiles :  $C1 = 0,01$  /  $C50 = 0,5$  = Médiane -> **divided per 100**
- **Jamovi** : Entrer les données, (dans une colonne) -> Mettre variable (continuous) -> Analyses -> exploration -> Descriptives -> variable (la colonne)

➔ Cocher : Median, Mode, Mean, Min, Max, Box Plot, Frequency tables

## 3. (Extrapolation linéaire)



## Intervalle interquartile (IQR ou IIQ)

$$IIQ = Q3 - Q1$$

## Boîte à moustache

Adjacente inf : plus petite valeur observée supérieure ou égale à  $Q1 - 1,5 (Q3 - Q1)$

Adjacente sup : plus grande valeur observée inférieure ou égale à  $Q3 + 1,5 (Q3 - Q1)$

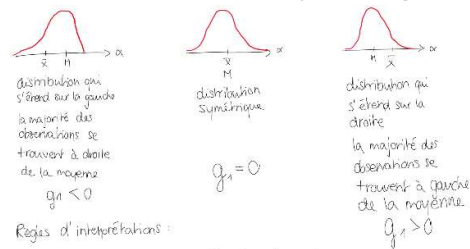
Ecart type/Standard deviation (s) : Analyse -> exploration -> descriptives -> statistics -> ...

Variance ( $s^2$ ) -> Les valeurs qui augmentent le plus la variance sont les valeurs qui sont les plus éloignées de la moyenne

Données extrêmes : point au-dessus de la moustache

## Caractéristiques de forme

### Skewness /coefficient d'asymétrie ( $g_1$ )



Règles d'interprétation :

Si  $g_1 < -1$  la distribution est asymétrique à gauche

Si  $-1 \leq g_1 \leq 1$  la distribution est symétrique

Si  $g_1 > 1$  la distribution est asymétrique à droite

### Kurtosis /coefficient d'aplatissement ( $g_2$ )



Règles d'interprétation :

Si  $g_2 < 2$  alors la distribution est platycurtique

Si  $2 \leq g_2 \leq 4$  la distribution est méso-curtique

Si  $g_2 > 4$  la distribution est leptocurtique

$$g_1 = \frac{m_3}{s^3} \quad m_3 = (x_i - \bar{x})^3 + (x_i - \bar{x})^3 \dots$$

$$g_2 = \frac{m_4}{s^4}$$

$$m_4 = (x_i - \bar{x})^4 + (x_i - \bar{x})^4 \dots$$

	Condition	Réponse	Effectif
1	Son	Juste	10
2	Son	Faux	28
3	Vidéo	Juste	17
4	Vidéo	Faux	18
5	Son + Vidéo	Juste	9
6	Son + Vidéo	Faux	23

Identification du gagnant	son	vidéo	Son + vidéo	
oui	10	17	9	36
non	28	18	23	69
	38	35	32	105

Table de contingence			
	$a_1$	$a_2$	$a_3$
$l_1$	$n_{11}$	$n_{12}$	$n_{13}$
$l_2$	$n_{21}$	$n_{22}$	$n_{23}$
	$n_{\cdot 1}$	$n_{\cdot 2}$	$n_{\cdot 3}$

$I = 2$   
n° de lignes  
 $J = 3$   
n° de colonnes

## Table de contingence = $n_{ij}$ :

1. Saisir les données
2. Chemin = Analyse – Frequencies – Contingency Tables – Independent Samples
3. Remplir -> Rows (Réponse) Columns (Condition) Counts (Effectif)

### Distribution conditionnelle de Y :

1. Faire table de contingence
  2. Cocher « Column » (sous Cells – Percentages)
- ⇒ Si les distributions conditionnelles de Y sont toutes les mêmes = Y ne dépend pas de X, sinon Y dépend de X

### Distribution conditionnelle de X :

1. Faire table de contingence
  2. Cocher « Row » (sous Cells – Percentages)
- ⇒ Si les distributions conditionnelles de X sont toutes les mêmes = X ne dépend pas de Y, sinon X dépend de Y

### Indépendance des variables X et Y :

- a. Les distributions conditionnelles de Y sont les mêmes
- b. Les distributions conditionnelles de X sont les mêmes
- c. L'effectif associé à  $(i; j) = \frac{n_{i.} \times n_{.j}}{n}$
- d. La fréquence associée à la cellule  $(i; j) = f_{ij} = f_{i.} \times f_{.j}$  (Rappel :  $f_{i.} = n_{i.}/n$  et  $f_{.j} = n_{.j}/n$ )

### Tableau des effectifs théoriques = $e_{ij}$ (correspond à la situation d'indépendance) :

1. Reprendre les données de la table de contingence (même chemin, remplissage)
2. Cocher « Expected » (sous Cells – Counts)

### Tableau des résidus :

1. Permet de comparer le tableau des effectifs théoriques et la table de contingence
2. Contingence  $(n_{ij}) - \text{Effectif}(e_{ij}) \rightarrow (10 - 13 = -3)$

### Tableau des résidus standardisés :

1. -----  $\rightarrow \frac{n_{ij} - e_{ij}}{\sqrt{e_{ij}}}$
- ⇒ Intéressant si la valeur est supérieure ou égale à 2

*taille d'effet*

df*	small	medium	large
1	.10	.30	.50
2	.07	.21	.35
3	.06	.17	.29
4	.05	.15	.25
5	.04	.13	.22

### Règle d'interprétation de Cohen : ----- $\rightarrow$

### Valeur du Khi carré (distance entre $n_{ij}$ et $e_{ij}$ ):

1. Après avoir fait une table de contingence et un tableau des effectifs théoriques
  2. Cocher «  $X^2$  » (sous Statistics)
  3. Chiffre recherché = croisement  $X^2$  et Value
- ⇒  $X^2_{\max} = n [\min(I; J) - 1] \Rightarrow 105 [\min(2; 3) - 1] \Rightarrow 105 [2 - 1] \Rightarrow 105 [1] = 105$  (dépendance<sub>max</sub> fonctionnelle)
- ⇒  $X^2_{\min} = 0$  (situation d'indépendance)

### Coefficient de contingence et V de Cramer :

- ⇒ Indice ne dépendant plus ni de « n » ni de la dimension de la table de contingence
- Cocher « Contingency coefficient et Phi and Cramer's V » sous (Statistics – Nominal)
1. C (coefficient de contingence) =  $\sqrt{\frac{X^2}{X^2 + n}}$  →  $(C < 1)$  et en situation d'indépendance  $(C = 0)$
  2.  $\Phi$  (phi) =  $\sqrt{\frac{X^2}{n}}$  → situation d'indépendance  $(\Phi = 0)$ , Dépendance fonctionnelle  $(\Phi \geq 1)$
  3. V (V de Cramer) =  $\sqrt{\frac{X^2}{X^2_{\max}}}$  →  $X^2 = 0, V = 0 / X^2_{\max} = X^2, V = 1$

### Liaison entre deux variables ordinales (Mesure d'un degré d'accord) :

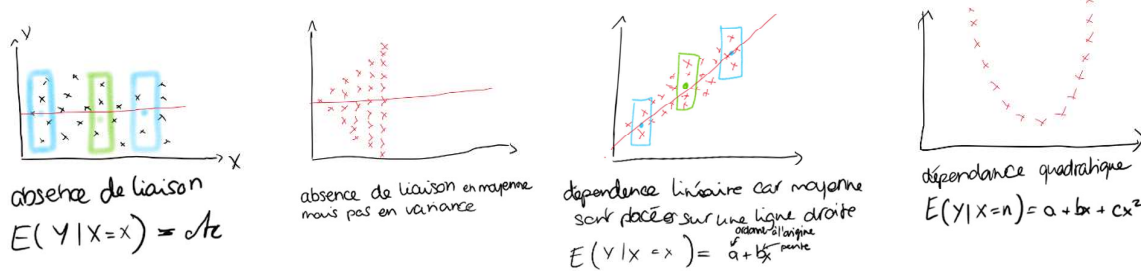
1. Calculer la concordance (en dessous à droite) → C
2. Calculer la discordance (en dessous à gauche) → D
3. Calculer les paires ex-aequo 1' (lignes) →  $N_x$
4. Calculer les paires ex-aequo 2' (colonnes) →  $N_y$
5.  $\text{Gamma}_{\text{Goodman}}$  -----  $\rightarrow$
6.  $\text{Tau}_b$  : -----  $\rightarrow$

$$\hat{\gamma} = \frac{C - D}{C + D}$$
$$\hat{\tau}_b = \frac{C - D}{\sqrt{(C + D + N_x)(C + D + N_y)}}$$

### Faire diagramme de dispersion avec Jamovi :

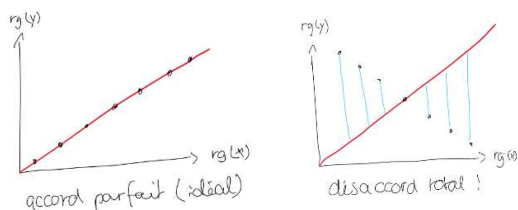
- Tableau des données (attentions variables numériques) → exploration → Scatterplot

## Etude de diagramme de dispersion (corrélation graphique) :



## Coefficient de Spearman ( $r_s$ ) :

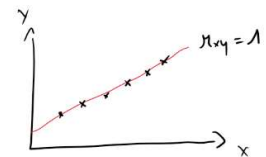
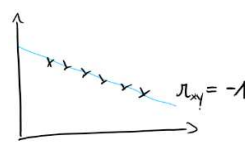
1. Calculer les rangs ( $rg_x$  et  $rg_y$ )  
 $\Rightarrow$  La donnée la plus basse = 1, puis la suivante = 2, etc. (si plusieurs fois la même valeur, on fait la moyenne)
2. Calculer  $d_i$  et  $d_i^2$   
 $\Rightarrow d_i = (rg_y - rg_x)$  et  $d_i^2 = d_i \times d_i$
3. Analyse  $\rightarrow$  regression  $\rightarrow$  correlation Matrix  
 $\Rightarrow$  Force de corrélation  $\rightarrow$  0-0.19 très faible / 0.2-0.39 faible / 0.4-0.59 modéré / 0.6-0.79 forte / 0.8-1 très forte



## Liaison entre deux variables numériques, Bravais-Pearson ( $r_{xy}$ )

- Jamovi  $\rightarrow$  Analyse  $\rightarrow$  regression  $\rightarrow$  correlation Matrix
- $\Rightarrow$  Interprétation : valeur absolue entre 0-1 (-1)  $\rightarrow$  0-0.5 = faible / 0.5-1 = forte / 0 et indé = réciproque fausse !!!

$|r_{xy}| = 1$  si et seulement si tous les points sont placés sur une même droite



## Estimation d'une proportion (ponctuelle)

1. Individus statistiques :  $n$  = nbr de quartiers (97) / nbr total d'individu (160)
2.  $n_a$  = nbr de pépins (26) / nbr d'individu avec tel caractéristique (30)
3.  $f_A = 26/97 = 0,268$  //  $30/160 = 0.188$
4.  $\hat{\pi} = f_A = 26,8\%$
5. Estimation de  $\hat{\pi}$  par intervalle de confiance  
 $\Rightarrow$  Méthode exacte/binomiale  $\rightarrow$  Jamovi  $\rightarrow$  Analyses  $\rightarrow$  Frequencies  $\rightarrow$  2 outcomes (binomial test)
6. Précision  $\rightarrow$  inversement proportionnelle à la variance  $F \rightarrow$  accroître la précision = augmenter  $n$  (taille d'échantillon)
7. SE : (erreur standard  $\rightarrow$  écart type de la distribution)

# Statistiques 6

	Condition	Réponse	Effectif
1	Son	Juste	10
2	Son	Faux	28
3	Vidéo	Juste	17
4	Vidéo	Faux	18
5	Son + Vidéo	Juste	9
6	Son + Vidéo	Faux	23

Identification du gagnant	Son	vidéo	Son + vidéo	
oui	10	17	9	36
non	28	18	23	69
	38	35	32	105

<u>Table de contingence</u>				
	$a_1$	$a_2$	$a_3$	
$l_1$	$n_{11}$	$n_{12}$	$n_{13}$	$n_{1.}$
$l_2$	$n_{21}$	$n_{22}$	$n_{23}$	$n_{2.}$
	$n_{.1}$	$n_{.2}$	$n_{.3}$	$n_{..}$

$I = 2$   
n° de lignes  
 $J = 3$   
n° de colonnes

## Table de contingence = $n_{ij}$ :

- Saisir les données
- Chemin = Analyse – Frequencies – Contingency Tables – Independent Samples
- Remplir -> Rows (Réponse) Columns (Condition) Counts (Effectif)

## Distribution conditionnelle de Y :

- Faire table de contingence
  - Cocher « Column » (sous Cells – Percentages)
- ⇒ Si les distributions conditionnelles de Y sont toutes les mêmes = Y ne dépend pas de X, sinon Y dépend de X

## Distribution conditionnelle de X :

- Faire table de contingence
  - Cocher « Row » (sous Cells – Percentages)
- ⇒ Si les distributions conditionnelles de X sont toutes les mêmes = X ne dépend pas de Y, sinon X dépend de Y

## Indépendance des variables X et Y :

- Les distributions conditionnelles de Y sont les mêmes
- Les distributions conditionnelles de X sont les mêmes
- L'effectif associé à  $(i; j) = \frac{n_{i.} \times n_{.j}}{n}$
- La fréquence associée à la cellule  $(i; j) = f_{ij} = f_i \times f_j$  (Rappel :  $f_i = n_{i.}/n$  et  $f_j = n_{.j}/n$ )

## Tableau des effectifs théoriques = $e_{ij}$ (correspond à la situation d'indépendance) :

- Reprendre les données de la table de contingence (même chemin, remplissage)
- Cocher « Expected » (sous Cells – Counts)

## Tableau des résidus :

- Permet de comparer le tableau des effectifs théoriques et la table de contingence
- Contingence ( $n_{ij}$ ) – Effectif ( $e_{ij}$ ) ->  $(10 - 13 = -3)$

## Tableau des résidus standardisés :

- >  $\frac{n_{ij} - e_{ij}}{\sqrt{e_{ij}}}$

⇒ Intéressant si la valeur est supérieure ou égale à 2

taille d'effet

df*	small	medium	large
1	.10	.30	.50
2	.07	.21	.35
3	.06	.17	.29
4	.05	.15	.25
5	.04	.13	.22

## Règle d'interprétation de Cohen : ----->

## Valeur du Khi carré (distance entre $n_{ij}$ et $e_{ij}$ ):

- Après avoir fait une table de contingence et un tableau des effectifs théoriques
  - Cocher «  $X^2$  » (sous Statistics)
  - Chiffre recherché = croisement  $X^2$  et Value
- ⇒  $X^2_{\max} = n [\min(I; J) - 1] \Rightarrow 105 [\min(2; 3) - 1] \Rightarrow 105 [2 - 1] \Rightarrow 105 [1] = 105$  (dépendance<sub>max</sub> fonctionnelle)
- ⇒  $X^2_{\min} = 0$  (situation d'indépendance)

## Coefficient de contingence et V de Cramer :

- ⇒ Indice ne dépendant plus ni de « n » ni de la dimension de la table de contingence
- Cocher « Contingency coefficient et Phi and Cramer's V » sous (Statistics – Nominal)
- C (coefficient de contingence) =  $\sqrt{\frac{X^2}{X^2 + n}}$  -> ( $C < 1$ ) et en situation d'indépendance ( $C = 0$ )
  - $\Phi$  (phi) =  $\sqrt{\frac{X^2}{n}}$  -> situation d'indépendance ( $\Phi = 0$ ), Dépendance fonctionnelle ( $\Phi \geq 1$ )
  - V (V de Cramer) =  $\sqrt{\frac{X^2}{X^2_{\max}}}$  ->  $X^2 = 0, V = 0 / X^2_{\max} = X^2, V = 1$