



# Cours

# Statistiques bi-variées

K. BELATTAR,  
Département Informatique - Université d'Alger 1

# Introduction aux statistiques bi-variées

## Objectif:

- Etudier sur une même population de  $N$  individus deux variables (caractères) différentes.
- Analyser la relation entre deux variables.

## Exemples:

- Taille et poids
- Taux de chômage et sexes
- Effets et doses des médicaments,
- ...etc.

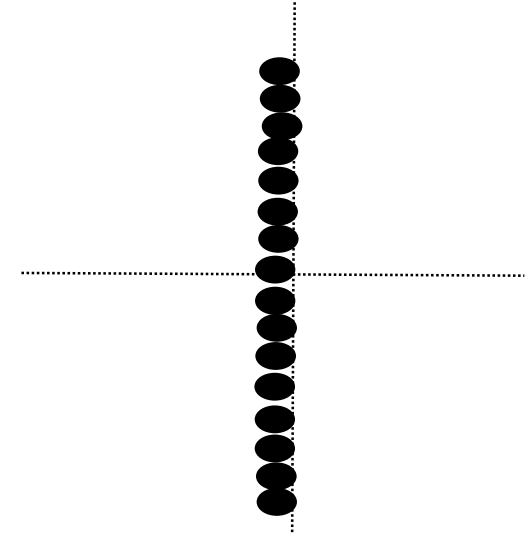
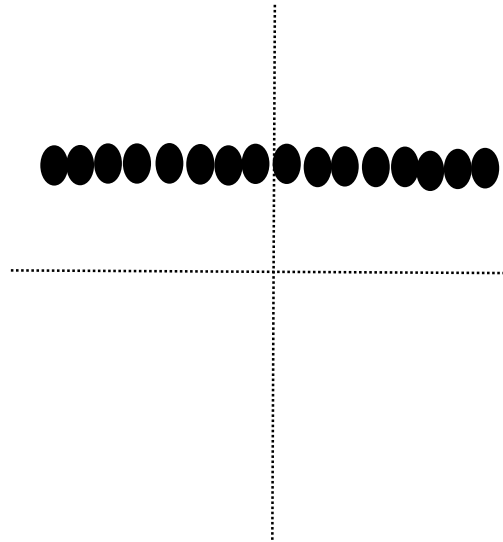
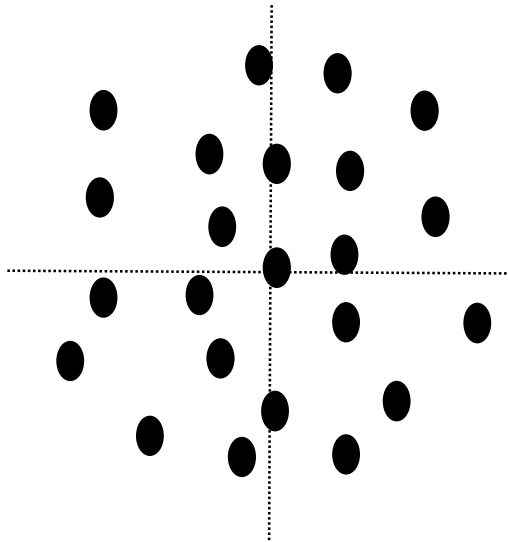
# Introduction aux statistiques bi-variées

Analyser une relation entre deux variables ***X et Y***:

- ✓ **Déterminer s'il existe une relation** entre les deux variables,
- ✓ **Caractériser la relation** entre les deux variables :
  - Forte vs faible (**selon l'intensité**)
  - Linéaire vs non linéaire (**selon la forme**)
  - Positive vs négative (**selon la direction**)
- ✓ **Tester si la liaison est statistiquement significative (tests d'hypothèse),**
- ✓ **Valider la relation identifiée (absence de biais).**

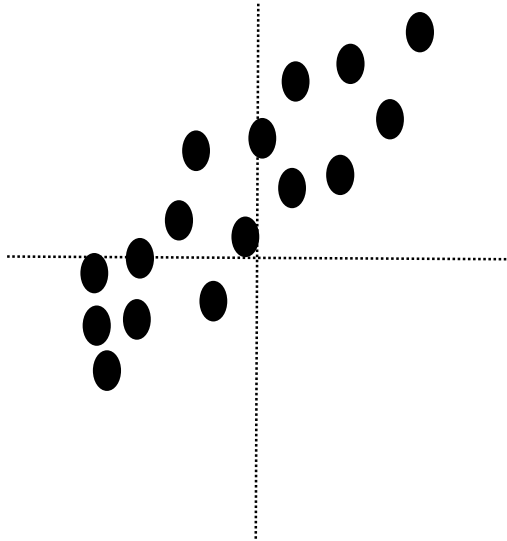
# Type de relations (selon l'intensité)

**Nuage de points**

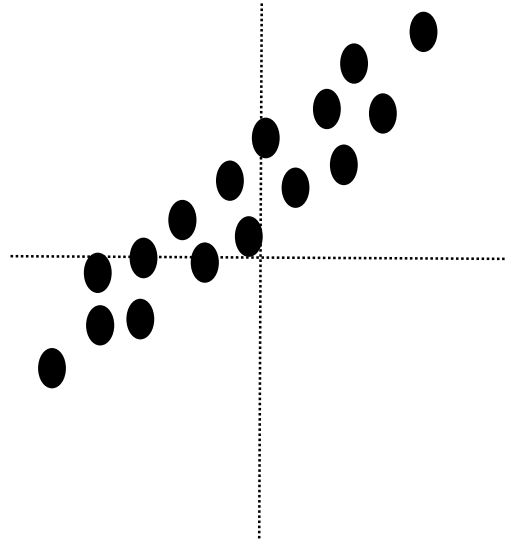


**Absence d'une relation**

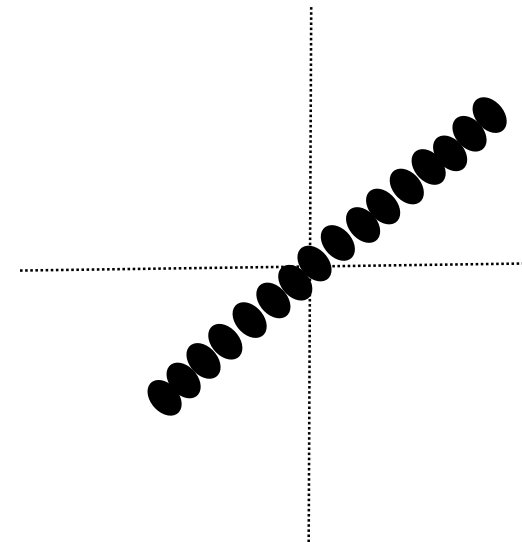
# Type de relations (selon l'intensité)



**Relation faible**

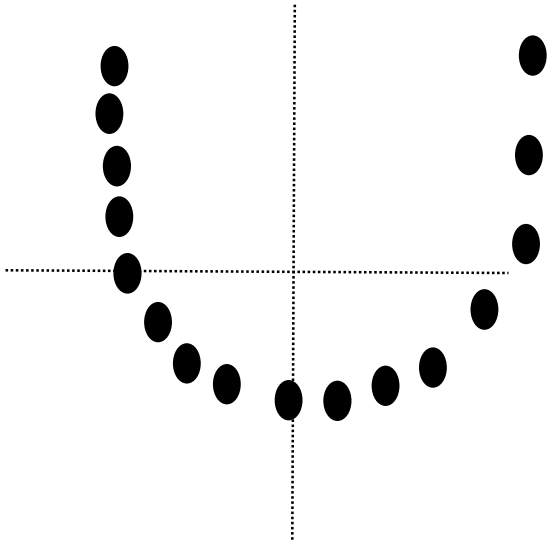


**Relation modérée**

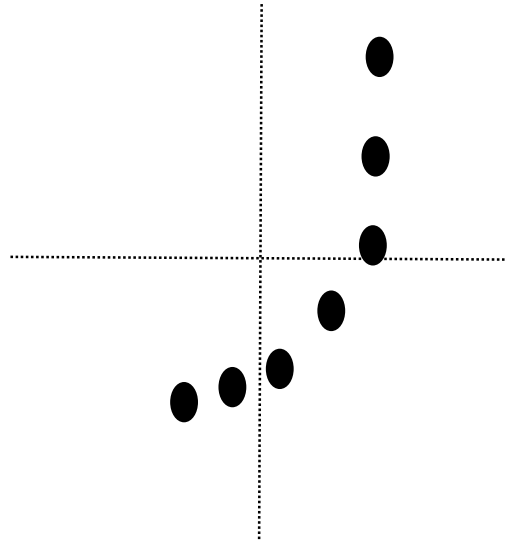


**Relation forte**

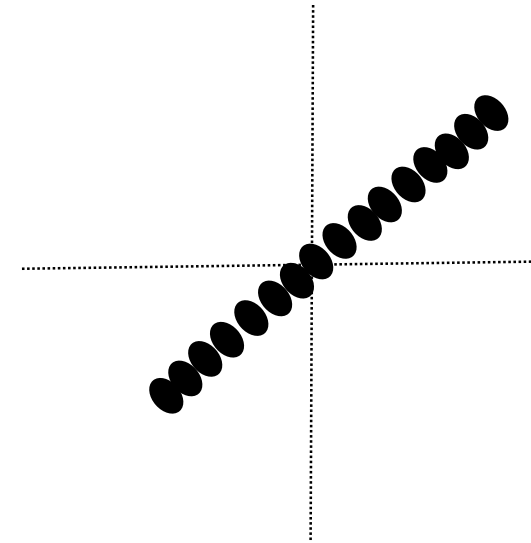
# Type de relations (selon la forme)



**Relation non linéaire  
et non monotone**

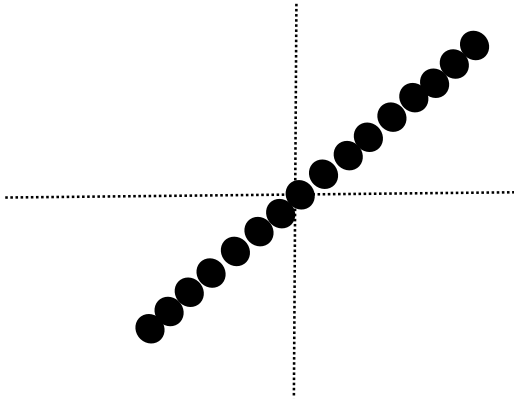


**Relation non linéaire  
et monotone**

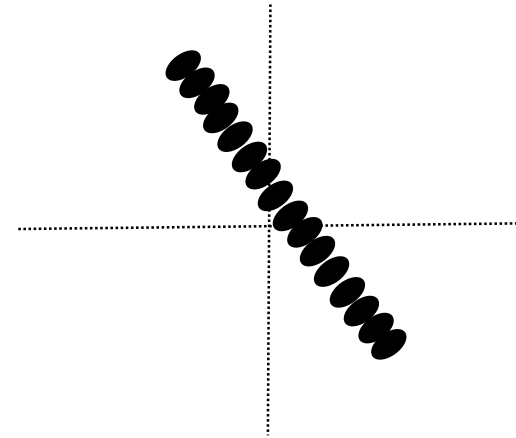


**Relation linéaire  
(toujours monotone)**

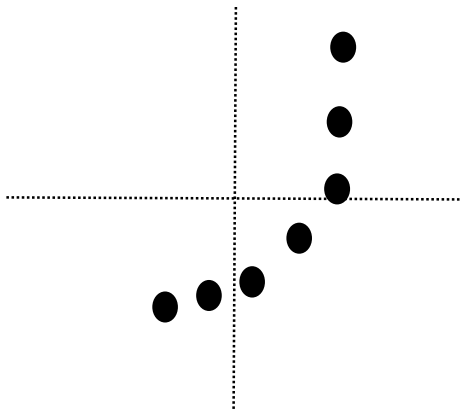
# Type de relations (selon la direction)



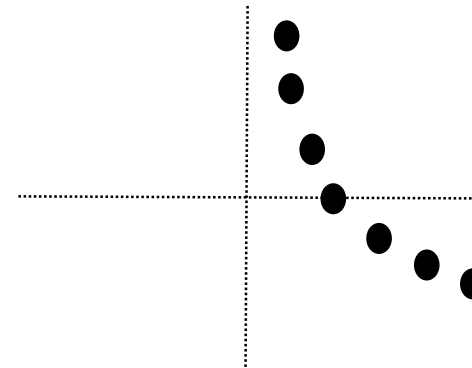
**Relation positive (et linéaire)**



**Relation négative (et linéaire)**



**Relation non linéaire positive**



**Relation non linéaire et négative**

# Séries statiques à deux variables

- Considérons deux variables (quantitatives, qualitatives, une quantitative et l'autre qualitative)  $X$  et  $Y$  définies sur la **même population** d'effectif total  $N$ .

$$N = \sum_i \sum_j n_{ij}$$

$n_{ij}$ : effectif conjoint des valeurs  $x_i, y_j$  pour lesquels  $X$  prend la valeur  $x_i$  et  $Y$  prend la valeur  $y_j$ .

- Une série statistique bi-variée est représentée par l'ensemble des couples

$$\{(x_1, y_1), (x_2, y_2), \dots, (x_k, y_l)\}$$



# Tableau de contingence

Valeurs  $y_i$

	$x_i$ / $y_j$	Valeurs $y_i$			
		$[L'_1, L'_2[$ $y_1$	.....	$[L'_{l-1}, L'_l[$ $y_l$	$n_{i\bullet}$
Valeurs $x_i$	$x_1$ ou $[L_1, L_2[$	$n_{11}$	.....	$n_{1l}$	$n_{1\bullet}$
	$x_2$ ou $[L_2, L_3[$	$n_{21}$	.....	$n_{2l}$	$n_{2\bullet}$
	.....	.....	.....	.....	.....
	$x_k$ ou $[L_{n-1}, L_n[$	$n_{n1}$	.....	$n_{nl}$	$n_{k\bullet}$
	$n_{\bullet j}$	$n_{\bullet 1}$	.....	$n_{\bullet l}$	$N$
		Effectifs marginaux de $Y$			

Effectifs marginaux de  $X$

Effectif total

Effectifs des couples  $(x_i, y_j)$

Effectif marginal de  $X$ :  $n_{i\bullet} = \sum_j n_{ij}$

Effectif marginal de  $Y$ :  $n_{\bullet j} = \sum_i n_{ij}$

Effectif des couples  $(x_i, y_j)$ :  $N = \sum_i \sum_j n_{ij} = \sum_i n_{i\bullet} = \sum_j n_{\bullet j}$

# Tableau de contingence

		Valeurs $y_i$			
Valeurs $x_i$	$x_i$ $y_j$	$[L'_1, L'_2[$ $y_1$	.....	$[L'_{l-1}, L'_l[$ $y_l$	$f_{i\bullet}$
	$x_1$ ou $[L_1, L_2[$	$f_{11}$	.....	$f_{1l}$	$f_{1\bullet}$
	$x_2$ ou $[L_2, L_3[$	$f_{21}$	.....	$f_{2l}$	$f_{2\bullet}$
	.....	.....	.....	.....	.....
	$x_k$ ou $[L_{n-1}, L_n[$	$f_{n1}$	.....	$f_{nl}$	$f_{k\bullet}$
		$f_{\bullet 1}$	.....	$f_{\bullet l}$	<b>1 (100%)</b>
		Fréquences marginales de $Y$			

Fréquences des couples  $(x_i, y_j)$

Fréquences marginales de  $X$

Fréquence marginale de  $X$ :  $f_{i\bullet} = \sum_j f_{ij}$

Fréquence marginale de  $Y$ :  $f_{\bullet j} = \sum_i f_{ij}$

Fréquence des couples  $(x_i, y_j)$ :  $\sum_i \sum_j f_{ij} = \sum_i f_{i\bullet} = \sum_j f_{\bullet j} = 1$  avec  $f_{ij} = \frac{n_{ij}}{N}$

# Tableau de contingence

## Exemple :

Soit la répartition des salaires d'une entreprise selon le nombre d'enfant (X) et le salaire mensuel(Y) en 1000DA.

$x_i \backslash y_j$	[2-6[	[6-10[	[10-14[	$n_{i.}$	$f_{i.}$
1	15	8	2	25	0.42
2	13	4	1	18	0.3
3	11	3	3	17	0.28
$n_{.j}$	39	15	6	60	
$f_{.j}$	0.65	0.25	0.1		1

# Distributions marginales

- **Effectif (fréquence) marginal (e) de  $X$ :**  $n_{i.} = \sum_j n_{ij}$  ( $f_{i.} = \sum_j f_{ij}$ )
- **Effectif (fréquence) marginal (e) de  $Y$ :**  $n_{.j} = \sum_i n_{ij}$  ( $f_{.j} = \sum_i f_{ij}$ )
- **Effectifs (fréquence) des couples  $(x_i, y_j)$ :**

$$N = \sum_i \sum_j n_{ij} = \sum_i n_{i.} = \sum_j n_{.j} \left( \sum_i \sum_j f_{ij} = \sum_i f_{i.} = \sum_j f_{.j} = 1 \text{ avec } f_{ij} = \frac{n_{ij}}{N} \right)$$

- **Moyenne marginale (variables discrètes)**

$$\bar{X} = \frac{\sum_{i=1}^N n_{i.} x_i}{N} = \sum_{i=1}^N f_{i.} x_i \quad \bar{Y} = \frac{\sum_{j=1}^N n_{.j} y_j}{N} = \sum_{j=1}^N f_{.j} y_j$$

- **Variance marginale (d'un échantillon)**

$$\sigma^2(X) = \frac{\sum_{i=1}^N n_{i.} (x_i - \bar{X})^2}{N-1} \quad \sigma^2(Y) = \frac{\sum_{j=1}^N n_{.j} (y_j - \bar{Y})^2}{N-1}$$

# Distributions marginales

$x_i$	$n_{i.}$	$f_{i.}$	$\bar{X} = f_{i.} x_i$	$\sigma^2(X)$
1	25	0.42	0.42	0.14
2	18	0.3	0.6	0.6
3	17	0.28	0.84	1.34
$n_{.j}$	60		1.86	2,08 ( $\Sigma$ )

# Distributions marginales

$y_j$	$n_{.j}$	$f_{.j}$	$\bar{Y} = f_{.j} C_j$	$\sigma^2(Y)$
[2-6[	39	0.65	2.6	1.3
[6-10[	15	0.25	2	9.15
[10-14[	6	0.1	1.2	11.86

# Distributions conditionnelles

- **Fréquence conditionnelle de  $X$  sachant que  $Y = y_j$ :**

$$f_{i/j} = f_i^j = \frac{n_{ij}}{n_{.j}} \quad \sum_{i=1}^k f_i^j = 1$$

$f_i^j$ : proportion d'individus présentant la valeur  $x_i$  parmi l'ensemble des individus présentant la valeur  $y_j$  de  $Y$ .

- **Fréquence conditionnelle de  $Y$  sachant que  $X = x_i$ :**

$$f_{j/i} = f_i^j = \frac{n_{ij}}{n_{i.}} \quad \sum_{j=1}^l f_i^j = 1$$

$f_i^j$ : proportion d'individus présentant la valeur  $y_j$  parmi l'ensemble des individus présentant la valeur  $x_i$  de  $X$ .

# Distributions conditionnelles

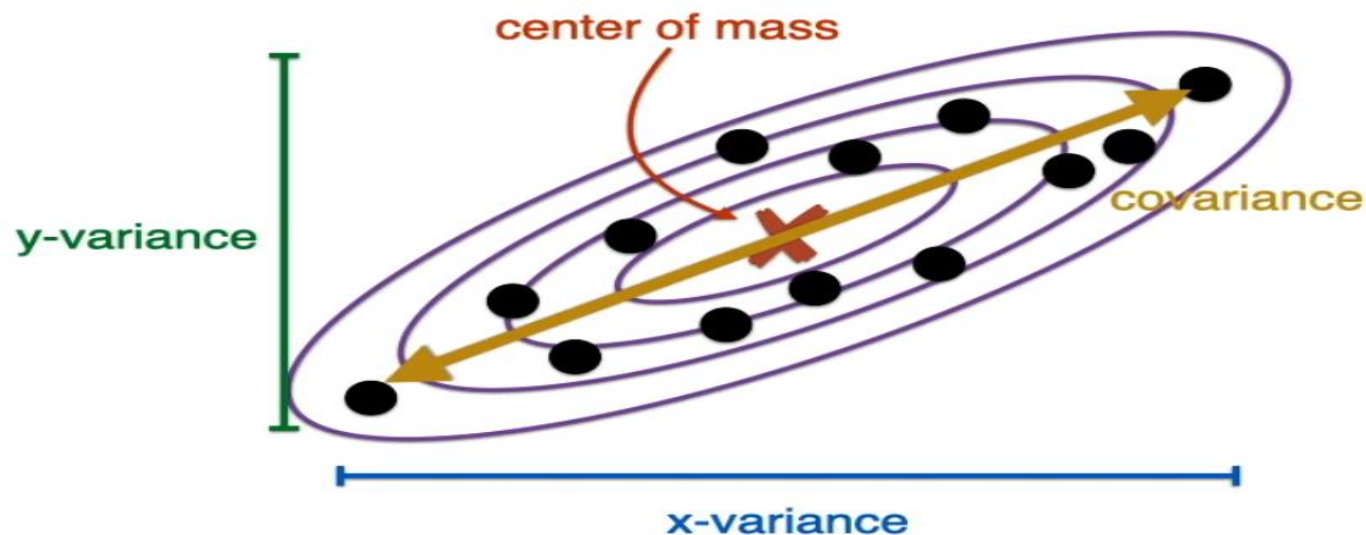
Exemple: La distribution conditionnelle de  $X$  sachant que  $Y=y_2$  ( $[6-10[$ )

$X$	$n_{ij}$	$f_{i/2}$
1	8	0.53
2	4	0.27
3	3	0.2
$\Sigma$	15	1



# Corrélation basée sur la covariance de données

- La forme du jeu de données (dataset) est déterminée par la **covariance**.
- L'élongation des points (individus) dans la direction diagonale.



$$\Sigma = \begin{pmatrix} \text{Var}(x) & \text{Cov}(x, y) \\ \text{Cov}(x, y) & \text{Var}(y) \end{pmatrix}$$

# Covariance de données

Considérons deux variables **numériques aléatoires**  $X$ ,  $Y$  et un échantillon de  $n$  individus  $\{(x_1, y_1), \dots, (x_n, y_n)\}$ .

$$E(X) = \bar{X} = \frac{\sum_{i=1}^n x_i}{n} \quad \text{et} \quad E(Y) = \bar{Y} = \frac{\sum_{i=1}^n y_i}{n}$$

$$Cov(X, Y) = E((X - \bar{X}) * (Y - \bar{Y})) = \frac{\sum_{i=1}^n (x_i - \bar{X}) * (y_i - \bar{Y})}{n}$$

$$Cov(X, Y) = E(X.Y) - \bar{X}\bar{Y}$$

**Quantifier la direction de la relation** entre les variables  $X$  et  $Y$ .

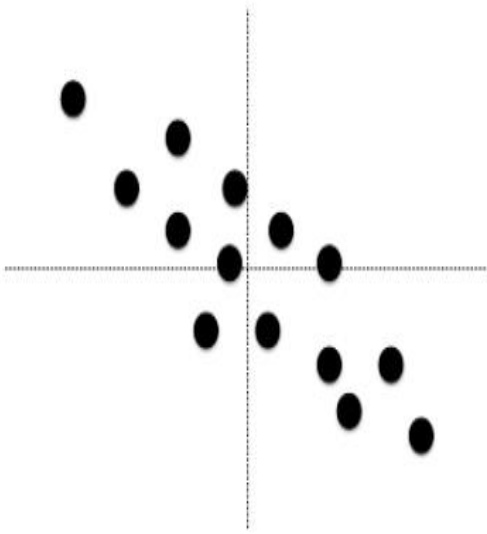
Si  $X > \bar{X}$  et  $Y > \bar{Y} \Rightarrow Cov(X, Y)$  est positive.

Si  $X > \bar{X}$  et  $Y < \bar{Y} \Rightarrow Cov(X, Y)$  est négative.

Si  $X$  et  $Y$  sont indépendants  $\rightarrow Cov(X, Y) = 0$  (ou presque null).

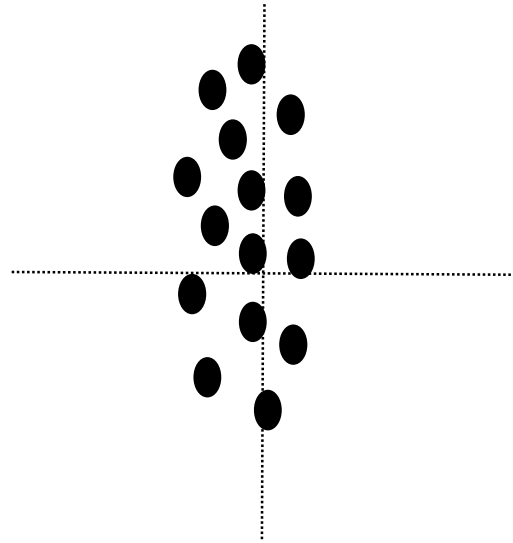
Variance est un cas spéciale de la covariance (i.e  $cov(X, X) = var(X) = \frac{\sum_{i=1}^n (x_i - \bar{X})^2}{n}$ ).

# Covariance de données



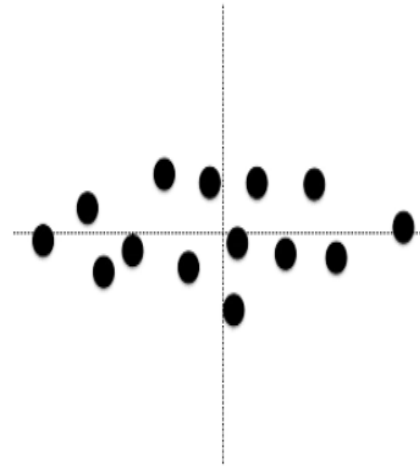
**Covariance négative**

$$Cov(X, Y) < 0$$



**Covariance nulle (ou presque nulle)**

$$Cov(X, Y) = 0$$



**Covariance positive**

$$Cov(X, Y) > 0$$

# Corrélation basée sur le coefficient de Pearson

Coefficient de corrélation linéaire de Pearson ( $r_{X,Y}$ ): quantifie **l'intensité de la relation** entre deux variables **quantitatives** (numériques: discrètes).

- Considérons deux variables numériques  $X, Y$  et un échantillon de  $n$  individus  $\{(x_1, y_1), \dots, (x_n, y_n)\}$ .

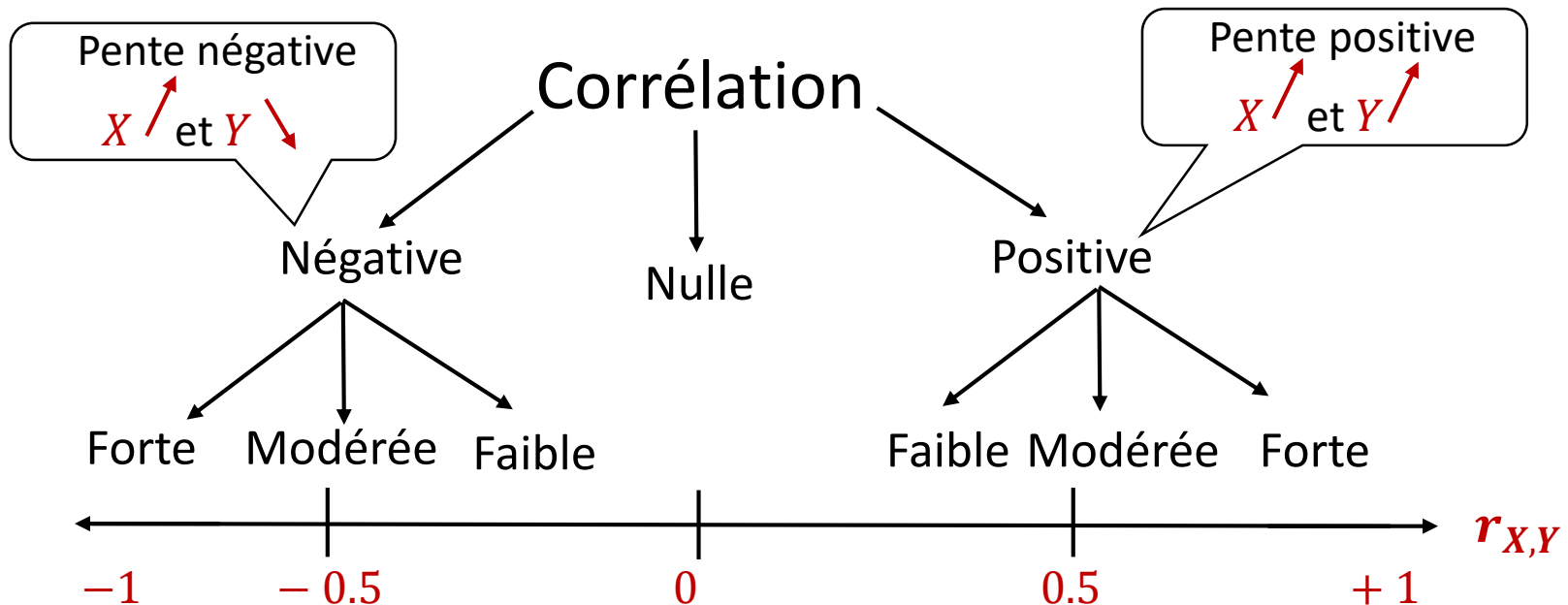
$$r_{X,Y} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} \quad \sigma_X = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})^2} \quad \sigma_Y = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{Y})^2}$$

$$r_{X,Y} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

$$-1 \leq r_{X,Y} \leq 1$$

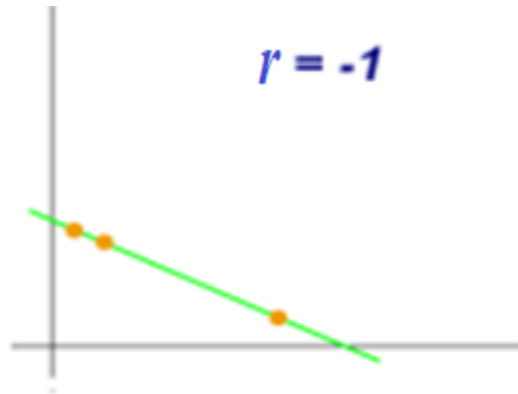
# Coefficient de corrélation linéaire

- Si  $r_{X,Y} \gg 0$  (valeur large proche de 1)  $\Rightarrow$  Corrélation linéaire forte positive (variables redondantes) entre  $X$  et  $Y$
- Si  $r_{X,Y} = 0$  (ou proche de 0)  $\Rightarrow$  Les variables  $X$  et  $Y$  ne sont pas corrélés linéairement (variables indépendantes)
- Si  $r_{X,Y} < 0$  (proche de  $-1$ )  $\Rightarrow$  Corrélation linéaire forte négative (variables redondantes) entre les variables  $X$  et  $Y$ .

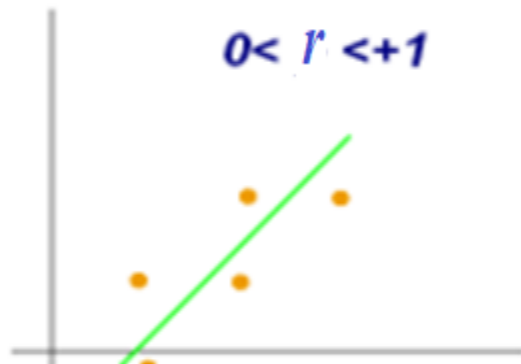
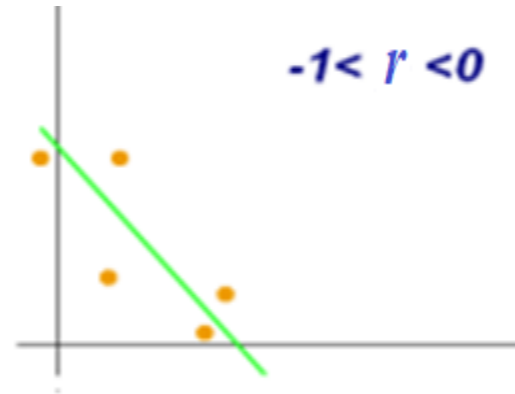


# Coefficient de corrélation linéaire

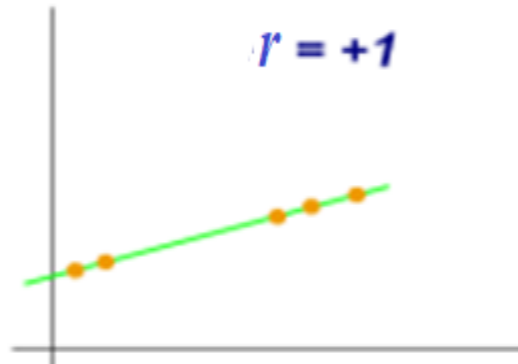
Corrélation linéaire forte négative



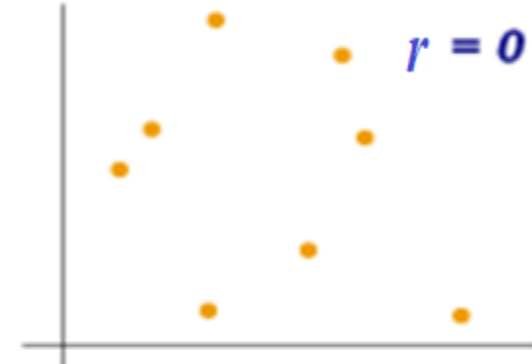
Corrélation linéaire modérée négative



Corrélation linéaire modérée positive



Corrélation linéaire forte positive



Absence de corrélation

# Corrélation basée sur le test de carré de Pearson

## Test de carré de Pearson

- ✓ Test statistique où la statistique de test (variable aléatoire) suit une loi du  $\chi^2$  sous l'hypothèse nulle.
- ✓ Il permet d'évaluer :
  - (1) La qualité de l'ajustement ,
  - (2) L'homogénéité et
  - (3) L'indépendance des variables (les individus de deux variables, exprimées dans un tableau de contingence)  
→ Accepter ou rejet d'une hypothèse

# Corrélation basée sur le test de carré de Pearson

## Test d'indépendance de carré de Pearson

- Supposant deux variables **qualitatives aléatoires**  $X = \{x_i, \dots, x_c\}$  et  $Y = \{y_j, \dots, y_r\}$  et un échantillon de  $n$  individus.
- Vérifier l'indépendance de ces deux variables ? (**hypothèse à tester**)
- Créer la table de contingence avec les distributions jointes  $(X_i, Y_j)$

$$\chi^2 = \sum_{i=1}^c \sum_{j=1}^r \frac{(o_{ij} - e_{ij})^2}{e_{ij}} \quad e_{ij} = \frac{\text{fréquence}_{marginale}(X = x_i) * \text{fréquence}_{marginale}(Y = y_j)}{n}$$

$o_{ij}$ : valeur Observée (empirique) de la distribution jointe  $(X_i, Y_j)$

$e_{ij}$ : valeur Espérée (théorique) de la distribution jointe  $(X_i, Y_j)$

$\text{fréquence}_{marginale}(X = x_i)$ : nombre des individus ayant la valeur  $x_i$

$\text{fréquence}_{marginale}(Y = y_j)$ : nombre des individus ayant la valeur  $y_j$



# Corrélation basée sur le test de carré de Pearson

## Test d'indépendance de carré de Pearson

### Principe:

- Définir **les hypothèses nulle et alternative**.
- Sélectionner un **niveau de confiance** souhaité (niveau de signification , valeur p ou niveau alpha correspondant ) pour le résultat du test.
- Calculer la statistique du test du chi carré  $\chi^2$  .
- Déterminer les **degrés de liberté**  $(rows-1)*(cols-1)$  , de cette statistique.
- Comparer avec **la valeur critique** de la distribution du chi carré avec les degrés de liberté et le niveau de confiance sélectionné.
- **Accepter ou rejeter l'hypothèse** nulle.

# Corrélation basée sur le test de carré de Pearson

## Exemple :

Table des valeurs observées:

Exemple :

Table des valeurs observées:

Diplôme

	Collège	Lycée	Bac	Master	Doctorat	Total	
Etat civil	Non marié	18	36	21	9	6	90
	Marié	12	36	45	36	21	150
	Divorcé	6	9	9	3	3	30
	Veuf	3	9	9	6	3	30
	Total	39	90	84	54	33	300

**Objectif:** vérifier la relation de corrélation entre les deux variables aléatoires « diplôme » et « état civil ».

# Corrélation basée sur le test de carré de Pearson

## Test d'indépendance de carré de Pearson

**Hypothèse nulle:** Pas de relation entre les variable diplôme et état civil.

**Hypothèse alternative:** relation significative entre les variable diplôme et état civil.  
niveau de signification=0.05

Table des valeurs espérées:

	Collège	Lycée	Bac	Master	Doctorat
Non marié	$\frac{90*39}{300} = 11.7$	$\frac{90*90}{300} = 27$	25.2	16.2	9.9
Marié	19.5	45	42	27	16.5
Divorcé	3.9	9	8.4	5.4	3.3
Veuf	3.9	9	8.4	5.4	3.3

# Corrélation basée sur le test de carré de Pearson

Calcul  $\chi^2$

Valeurs observées (o)	Valeurs espérées (e)	$(o_{ij} - e_{ij})$	$(o_{ij} - e_{ij})^2$	$\frac{(o_{ij} - e_{ij})^2}{e_{ij}}$
18	11.7	6.3	39.69	3.39
36	27	9	81	3
....	.....	.....	.....	...
3	3.3	-0.3	0.09	0.02

$$\chi^2 = \sum_{i=1}^c \sum_{j=1}^r \frac{(o_{ij} - e_{ij})^2}{e_{ij}} = 3.39 + 3 + \dots + 0.02 = 23.57$$

**Degré de liberté** = (rows-1)\*(cols-1) = (4-1)\*(5-1)=12

# Corrélation basée sur le test de carré de Pearson

## Test d'indépendance de carré de Pearson

Niveau de signification=0.05

Percentage Points of the Chi-Square Distribution									
Degrees of Freedom	Probability of a larger value of $\chi^2$								
	0.99	0.95	0.90	0.75	0.50	0.25	0.10	0.05	0.01
5	0.554	1.145	1.610	2.675	4.351	6.63	9.24	11.07	15.09
6	0.872	1.635	2.204	3.455	5.348	7.84	10.64	12.59	16.81
7	1.239	2.167	2.833	4.255	6.346	9.04	12.02	14.07	18.48
8	1.647	2.733	3.490	5.071	7.344	10.22	13.36	15.51	20.09
9	2.088	3.325	4.168	5.899	8.343	11.39	14.68	16.92	21.67
10	2.558	3.940	4.865	6.737	9.342	12.55	15.99	18.31	23.21
11	3.053	4.575	5.578	7.584	10.341	13.70	17.28	19.68	24.72
12	3.571	5.226	6.304	8.438	11.340	14.85	18.55	21.03	26.22
13	4.107	5.892	7.042	9.299	12.340	15.98	19.81	22.36	27.69
14	4.660	6.571	7.790	10.165	13.339	17.12	21.06	23.68	29.14

$$\chi^2_{\text{Critique}} = 21.03$$

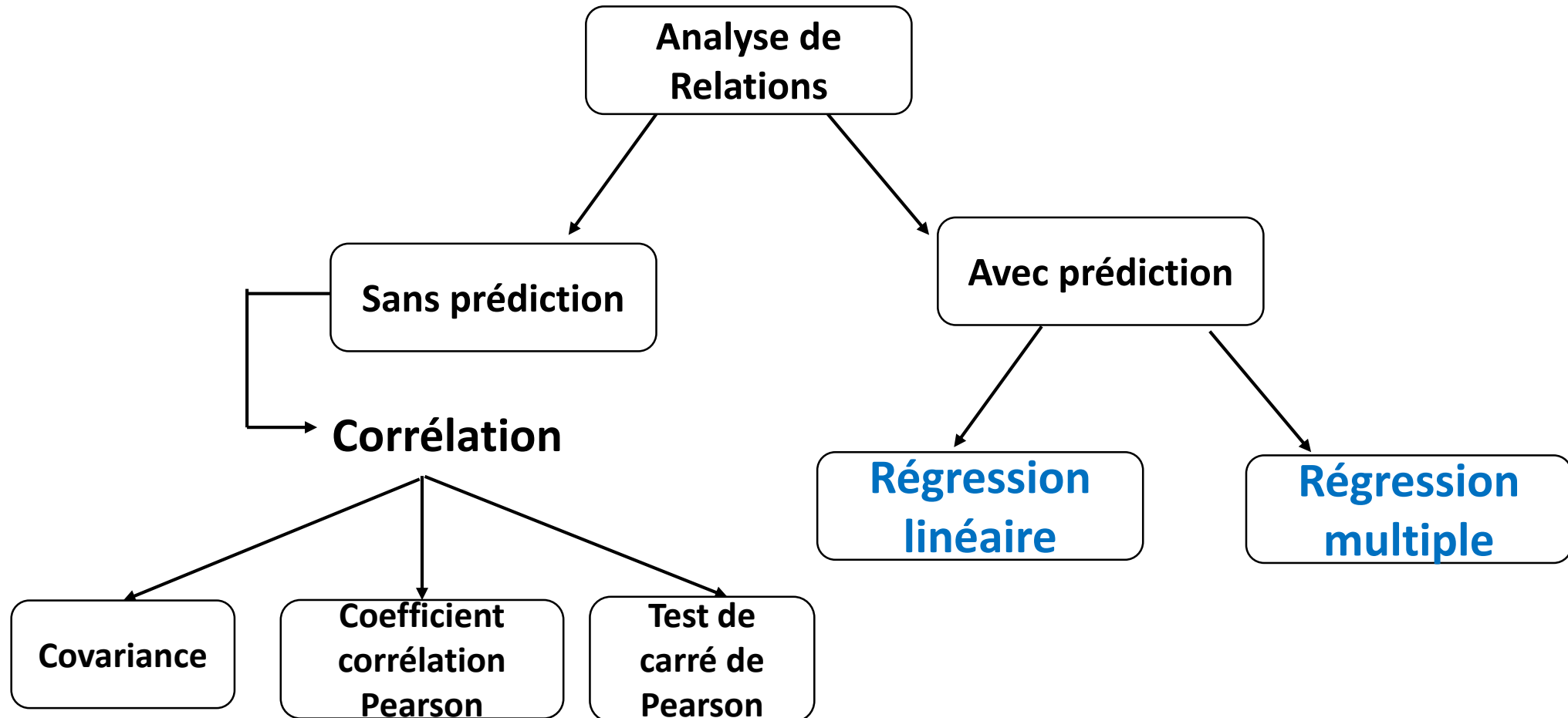
$$\chi^2 = 23.57$$

$$\chi^2 > \chi^2_{\text{Critique}}$$

<https://www.youtube.com/watch?v=f53nXHoMXx4&list=RDCMU3gEcM1157Q&index=1>

→ Rejeter l'hypothèse nulle et  
→ Accepter l'hypothèse alternative

# Introduction aux statistiques bi-variées



# Conclusion

- Etude statistique d'une série bi-variée.
- **Analyse des relations** entre **deux variables aléatoires** :
  - **Corrélation**
    - 1) **Variables quantitatives**:
      - Covariance (**quantifier la direction de la relation**)
      - Coefficient de corrélation linéaire (**quantifier l'intensité de la relation**)
    - 2) **Variables qualitatives**:
      - Test de carré de Pearson (**tester si la liaison est statistiquement significative**).
  - **Régression**