



Statistique descriptive

Ordonnancement des données

		caractère j				
		1	...	j	...	p
individu i	1	x_{11}	...	x_{1j}	...	x_{1p}
	\vdots	\vdots		\vdots		\vdots
	i	x_{i1}	...	x_{ij}	...	x_{ip}
	\vdots	\vdots		\vdots		\vdots
	n	x_{n1}	...	x_{nj}	...	x_{np}

Tableau individus \times caractères

i	1	...	i	...	n
x_i	x_1	...	x_i	...	x_n

Cas univarié

i	1	...	i	...	n
x_i	x_1	...	x_i	...	x_n
y_i	y_1	...	y_i	...	y_n

Cas bivarié

Analyse univariée

- Observations : x_1, \dots, x_n

individu i	1	...	i	...	n
caractère x_i	x_1	...	x_i	...	x_n

Individus \times caractères

- Observations ordonnées : $x_{(1)} \leq \dots \leq x_{(n)}$ (*avec répétition*), $x'_1 < \dots < x'_p$ ($p \leq n$) (*sans répétition*)

effectif k	n_1	...	n_k	...	n_p
caractère x'_k	x'_1	...	x'_k	...	x'_p

Série statistique

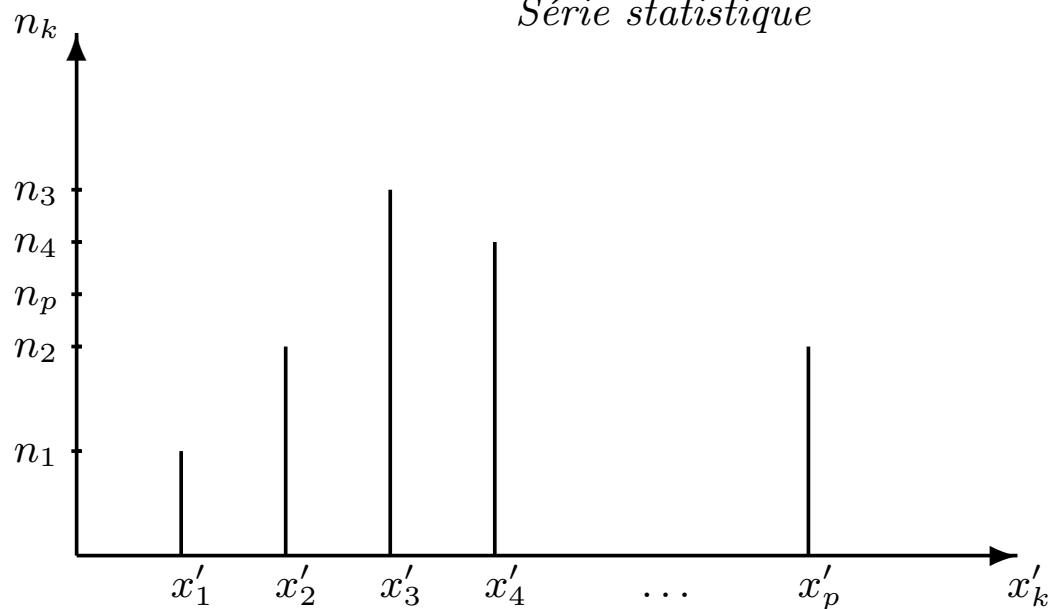


Diagramme en bâtons

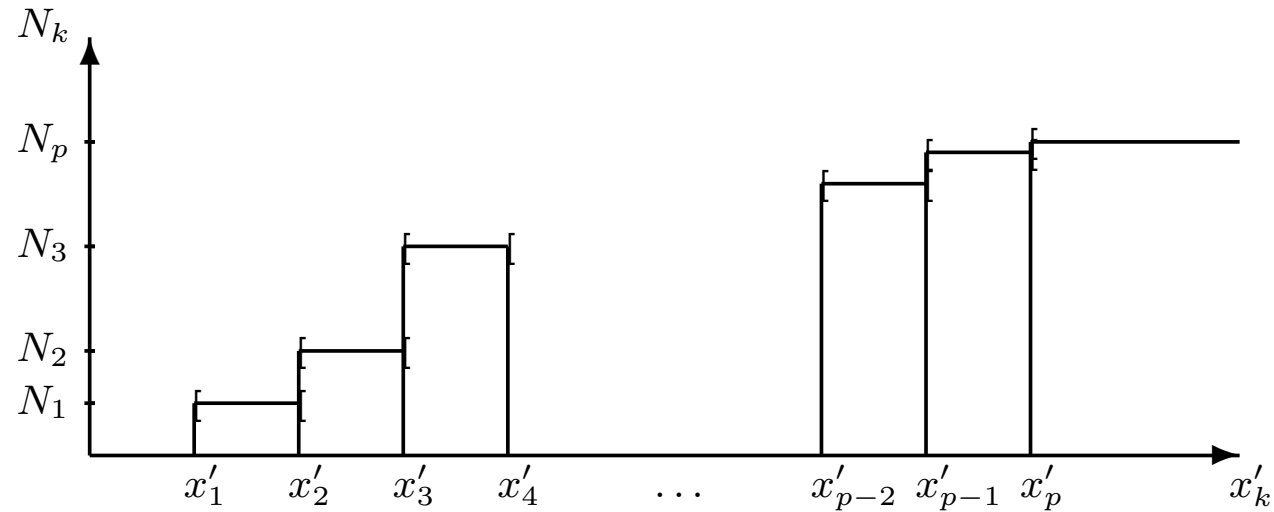
Effectifs cumulés, fréquences cumulées

- Effectifs cumulés à gauche :

$$N_k = \sum_{\ell=1}^k n_{\ell}$$

- Fréquences cumulées à gauche :

$$F_k = \sum_{\ell=1}^k \frac{n_{\ell}}{n}$$



Courbe cumulative à gauche $y = N(x)$

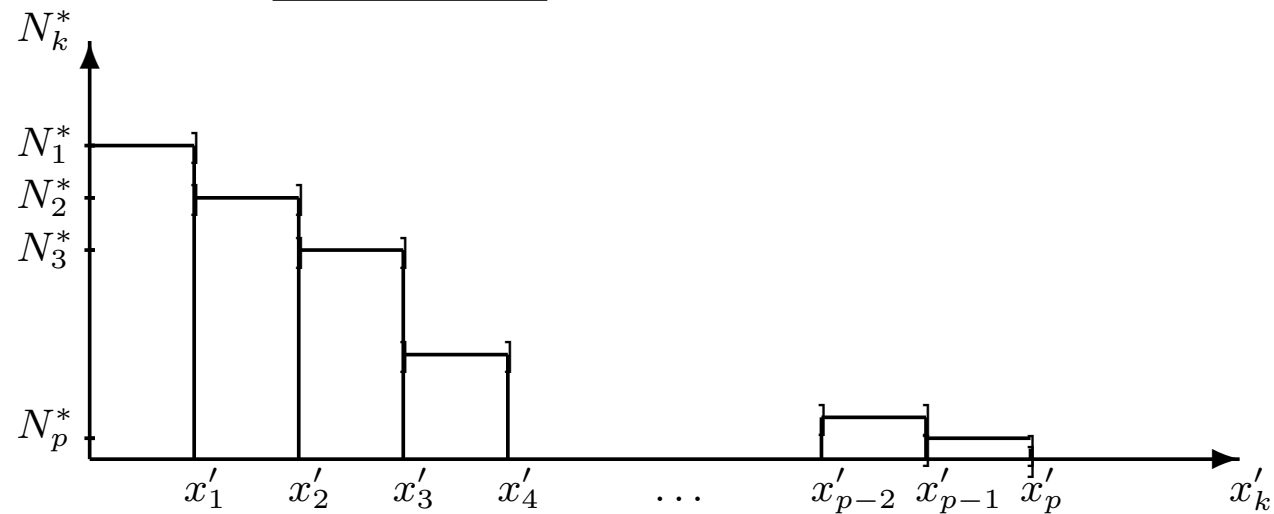
Effectifs cumulés, fréquences cumulées

- Effectifs cumulés à droite :

$$N_k^* = \sum_{\ell=k}^p n_\ell$$

- Fréquences cumulées à droite :

$$F_k^* = \sum_{\ell=k}^p \frac{n_\ell}{n}$$

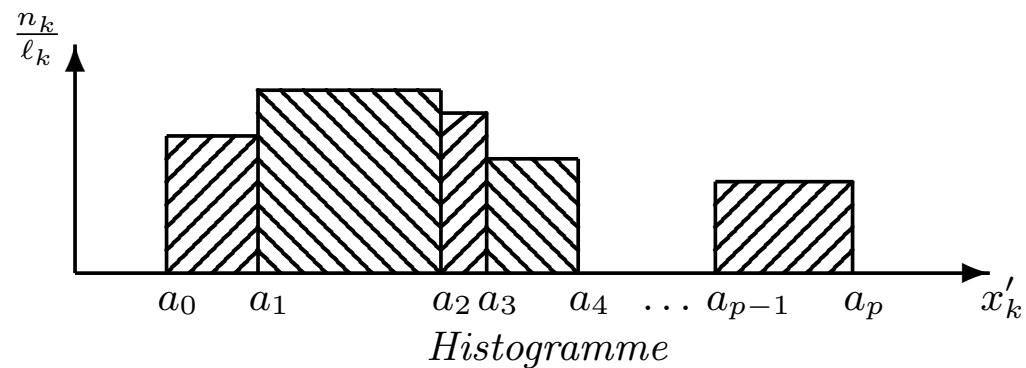


Courbe cumulative à droite $y = N^(x)$*

$$N_k(x) + N_k^*(x) = \begin{cases} n & \text{si } x \notin \{x'_1, \dots, x'_p\} \\ n + n_k & \text{si } x = x'_k \end{cases}$$

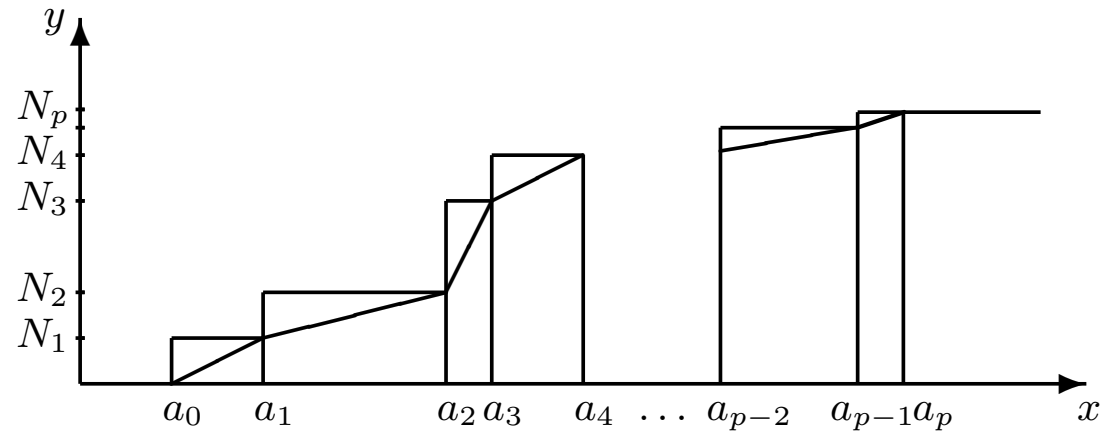
Observations groupées

- Données réparties en classes $[a_0, a_1[$, $[a_1, a_2[$, \dots , $[a_{p-2}, a_{p-1}[$, $[a_{p-1}, a_p]$
- Effectif et longueur de la classe $[a_{k-1}, a_k[$, $1 \leq k \leq p$: n_k et ℓ_k
- Densité d'effectif : $\frac{n_k}{\ell_k}$, $1 \leq k \leq p$

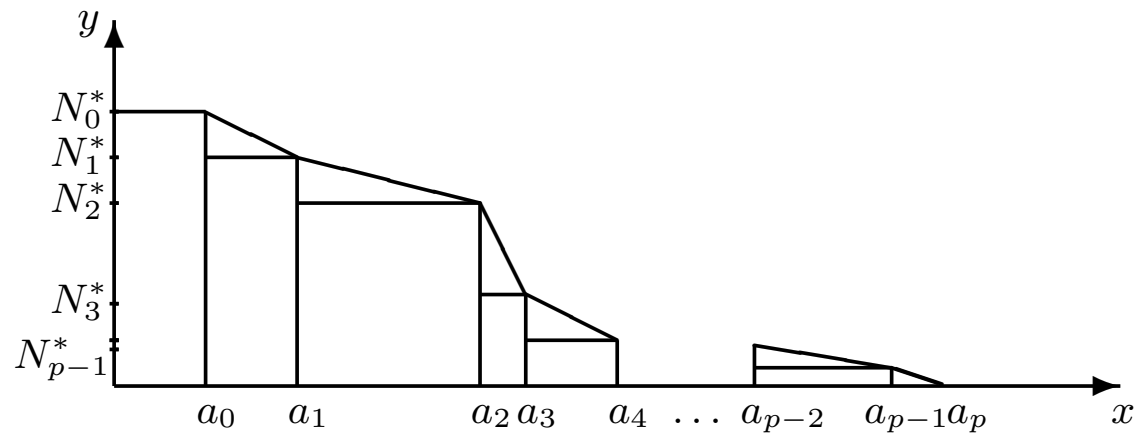


$$N(x) = \begin{cases} 0 & \text{si } x < a_0 \\ \frac{n_1}{\ell_1}(x - a_0) & \text{si } a_0 \leq x < a_1 \\ \vdots & \vdots \\ N_{k-1} + \frac{n_k}{\ell_k}(x - a_{k-1}) & \text{si } a_{k-1} \leq x < a_k \\ \vdots & \vdots \\ n & \text{si } x \geq a_p \end{cases}$$

Observations groupées



Courbe cumulative à gauche $y = N(x)$



Courbe cumulative à droite $y = N^(x)$*

$$\forall x, N(x) + N^*(x) = n.$$

Moyenne

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{n} \sum_{k=1}^p n_k x'_k$$

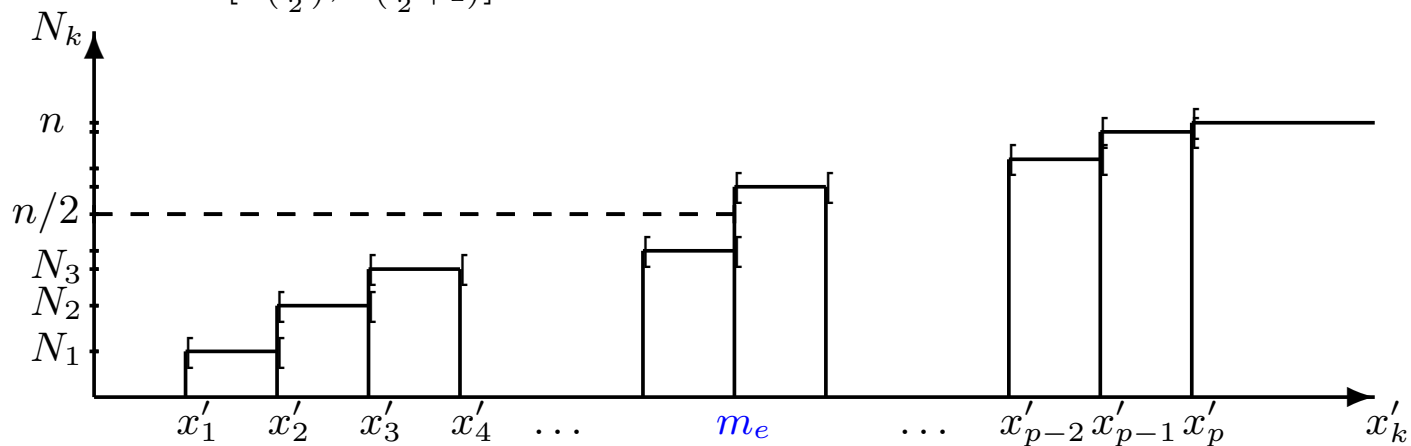
- Pour les observations groupées : x'_k est le centre de la classe $[a_{k-1}, a_k[$ soit $x'_k = \frac{1}{2}(a_{k-1} + a_k)$
- Cas de deux séries statistiques $\{x_i, 1 \leq i \leq n_x\}$ et $\{y_j, 1 \leq j \leq n_y\}$:
la série regroupée $\{z_k, 1 \leq k \leq n_z\} = \{x_i, 1 \leq i \leq n_x\} \cup \{y_j, 1 \leq j \leq n_y\}$ a pour moyenne

$$\bar{z} = \frac{n_x \bar{x} + n_y \bar{y}}{n_x + n_y}$$

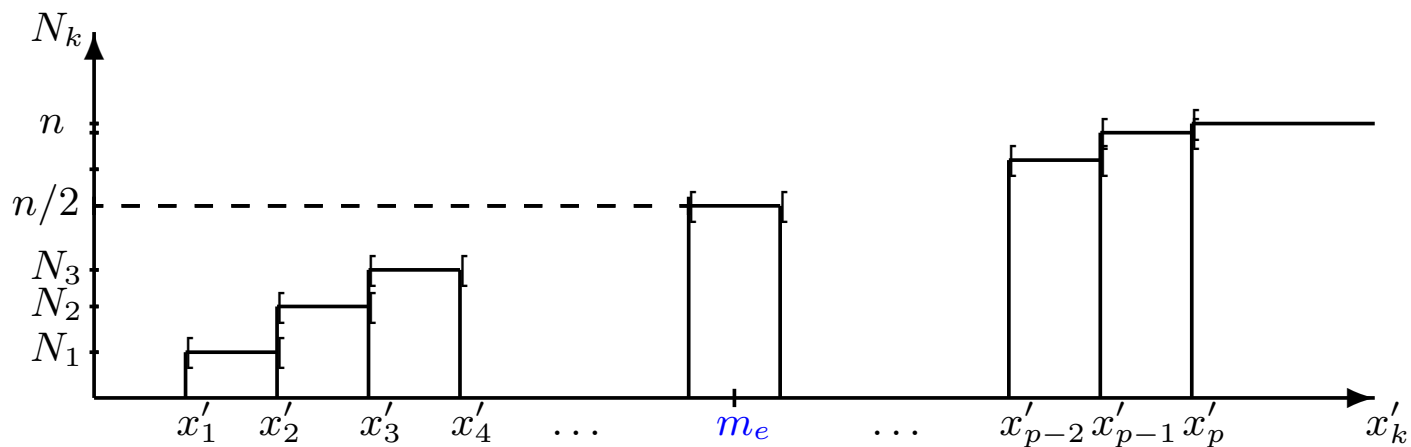
Médiane

$$m_e = \begin{cases} x_{(\frac{n+1}{2})} & \text{si } n \text{ est impair} \\ \frac{1}{2} [x_{(\frac{n}{2})} + x_{(\frac{n}{2}+1)}] & \text{si } n \text{ est pair} \end{cases}$$

- Cas n pair : intervalle médian $[x_{(\frac{n}{2})}, x_{(\frac{n}{2}+1)}]$



Courbe cumulative à gauche $y = N(x)$ (n impair)



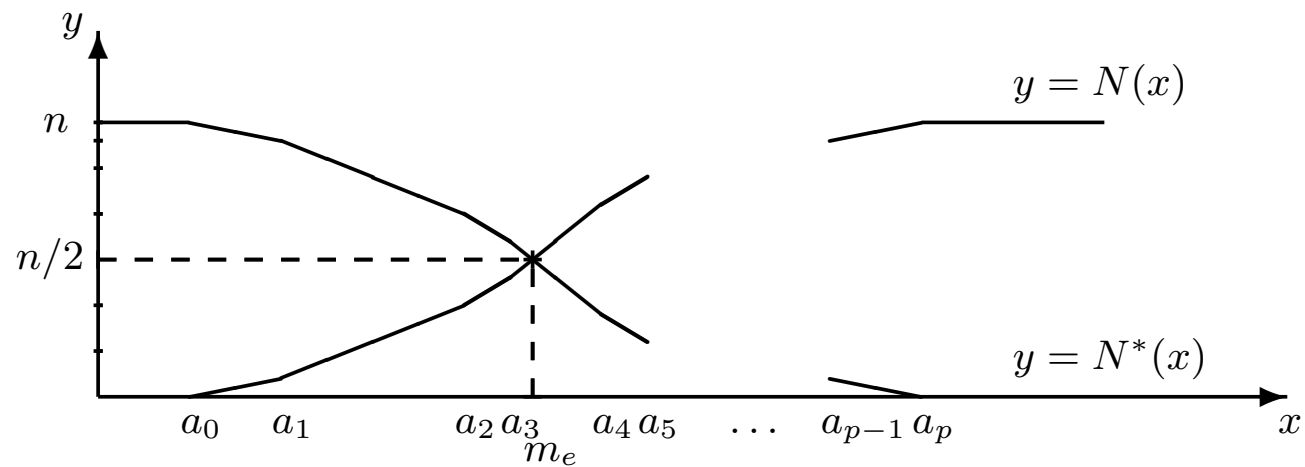
Courbe cumulative à gauche $y = N(x)$ (n pair)

Médiane

- Cas d'observations groupées :

$$m_e = a_{i-1} + \frac{\ell_k}{n_k} \left(\frac{n}{2} - N_{k-1} \right)$$

k étant l'indice de la classe contenant m_e



Courbes cumulatives $y = N(x)$ et $y = N^*(x)$

Quantiles

Soit $p \in]0, 1[$

$$Q_p = \begin{cases} x_{([np]+1)} & \text{si } np \text{ n'est pas entier} \\ \frac{1}{2}[x_{([np])} + x_{([np]+1)}] & \text{si } np \text{ est entier} \end{cases}$$

$[np]$ étant la partie entière de np

- Médiane : $p = \frac{1}{2} \longrightarrow m_e = Q_{\frac{1}{2}}$
- Quartiles : $p \in \{\frac{1}{4}, \frac{1}{2}, \frac{3}{4}\} \longrightarrow$ intervalle interquartile : $Q_{\frac{3}{4}} - Q_{\frac{1}{4}}$
- Déciles : $p \in \{\frac{1}{10}, \frac{2}{10}, \dots, \frac{9}{10}\}$
- Percentiles : $p \in \{\frac{1}{100}, \frac{2}{100}, \dots, \frac{99}{100}\}$

Variance, écart-type

$$v_x = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \sum_{k=1}^p n_k (x'_k - \bar{x})^2 \quad \sigma_x = \sqrt{v_x}$$

- Propriétés :

$$v_x = \overline{x^2} - \bar{x}^2 \quad \text{avec} \quad \overline{x^2} = \frac{1}{n} \sum_{i=1}^n x_i^2$$

$$v_x = \frac{1}{2n^2} \sum_{i=1}^n \sum_{j=1}^n (x_i - x_j)^2$$

$$\forall m \in \mathbb{R}, \quad \frac{1}{n} \sum_{i=1}^n (x_i - m)^2 = v_x + (\bar{x} - m)^2$$

- Cas de deux séries statistiques $\{x_i, 1 \leq i \leq n_x\}$ et $\{y_j, 1 \leq j \leq n_y\}$:

la série regroupée $\{z_k, 1 \leq k \leq n_z\} = \{x_i, 1 \leq i \leq n_x\} \cup \{y_j, 1 \leq j \leq n_y\}$ a pour variance

$$v_z = v_{\text{intra}} + v_{\text{inter}}$$

avec

$$v_{\text{intra}} = \frac{n_x v_x + n_y v_y}{n_x + n_y} \quad v_{\text{inter}} = \frac{n_x (\bar{z} - \bar{x})^2 + n_y (\bar{z} - \bar{y})^2}{n_x + n_y} = \frac{n_x n_y}{(n_x + n_y)^2} (\bar{x} - \bar{y})^2$$

v_{intra} est la variance *dans* les groupes, v_{inter} est la variance *entre* les groupes

Analyse bivariée

- Observations : $(x_1, y_1), \dots, (x_n, y_n)$

individu i	1	...	i	...	n
caractère x_i	x_1	...	x_i	...	x_n
caractère y_i	y_1	...	y_i	...	y_n

Individus \times caractères bivariés

- Observations ordonnées : $x'_1 < \dots < x'_p$ et $y'_1 < \dots < y'_q$ ($p, q \leq n$)

- Effectifs marginaux : $n_{i.} = \sum_{j=1}^q n_{ij}$ et $n_{.j} = \sum_{i=1}^p n_{ij} \rightarrow \sum_{i=1}^p n_{i.} = \sum_{j=1}^q n_{.j} = n$

caractère i	caractère j					effectifs marginaux
	y'_1	...	y'_j	...	y'_q	
x'_1	n_{11}	...	n_{1j}	...	n_{1q}	$n_{1.}$
\vdots	\vdots		\vdots		\vdots	\vdots
x'_i	n_{i1}	...	n_{ij}	...	n_{iq}	$n_{i.}$
\vdots	\vdots		\vdots		\vdots	\vdots
x'_p	n_{p1}	...	n_{pj}	...	n_{pq}	$n_{p.}$
effectifs marginaux	$n_{.1}$...	$n_{.j}$...	$n_{.q}$	n

Tableau de contingence

Covariance, corrélation

$$\sigma_{xy} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \frac{1}{n} \sum_{k=1}^p \sum_{\ell=1}^q n_{k\ell} (x'_k - \bar{x})(y'_\ell - \bar{y}) \quad \rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$$

- Propriétés :

$$\sigma_{xx} = \sigma_x^2$$

$$\sigma_{xy} = \overline{xy} - \bar{x}\bar{y} \text{ avec } \overline{xy} = \frac{1}{n} \sum_{i=1}^n x_i y_i$$

$$|\rho_{xy}| \leq 1$$

- Les données x_1, \dots, x_n et y_1, \dots, y_n sont *non-corrélées* lorsque $\sigma_{xy} = 0$ (ou encore $\rho_{xy} = 0$)

- Matrice de covariance :
$$\begin{pmatrix} \sigma_x^2 & \sigma_{xy} \\ \sigma_{xy} & \sigma_y^2 \end{pmatrix}$$