

Chapter 3: Two-variable statistical series

A dataset with two variables contains what is called bivariate data. It's a statistical series where two characters are studied simultaneously, X_i and Y_j . Each one can be either quantitative or qualitative.

Example 1

- A sample of 200 families, the number of children X and the number of rooms Y are studied simultaneously.
- A sample of 20 families, we study the monthly income X in Da and the monthly expenses Y .
- Students randomly selected from a section of L2 civil engineering, we observe the hair color X and the color of the eyes Y .

A.3.1 Data tables (contingency table):

In this type of table, the numbers **n_{ij} are partial counts**. These are the numbers of individuals that present both the x_i and the y_j modality.

X \ Y	y1	y2	y3	Yq	Marginal Column ni.
x1	n11	n12	n13	n1q	n1.
x2	n21	n22	n23	n2q	n2.
x3	n31	n32	n33	n3q	n3.
.
xp	np1	np2	np3	npq	np.
Marginal line n.j	n.1	n.2	n.3	n.q	N=n..

Example 1 :

- A sample of 80 families, the number of children X and the number of rooms Y are studied simultaneously.

Number of rooms Number of children	0	1	2	3	Total
1	4	8	1	4	17
2	5	8	7	5	25
3	12	8	9	9	38
Total	21	24	17	18	80

Solution: Contingency Table

X	Y	$y_{(j=1)}=0$	$y_{(j=2)}=1$	$y_{(j=3)}=2$	$y_{(j=4)}=3$	Total
$x_{(i=1)}=1$		$n_{11}=4$	$n_{12}=8$	$n_{13}=1$	$n_{14}=4$	$n_{1.}=17$
$x_{(i=2)}=2$		$n_{21}=5$	$n_{22}=8$	$n_{23}=7$	$n_{24}=5$	$n_{2.}=25$
$x_{(i=3)}=3$		$n_{31}=12$	$n_{32}=8$	$n_{33}=9$	$n_{34}=9$	$n_{3.}=38$
Total		$n_{.1}=21$	$n_{.2}=24$	$n_{.3}=17$	$n_{.4}=18$	$N=n_{..}=80$

- n_{24} indicates that 5 families with 3 bedrooms and 2 children.

A.3.1.1. Simple Scatter plot:

The graphical representation of a two-variable statistical series is made by a scatter plot that is plotted, from the raw data of X and Y.

- If X and Y are continuous quantitative variables, we plot the graph by their class centers C_i .

The coordinates of the mean point G are the arithmetic means of X and Y (\bar{x} ; \bar{y}).

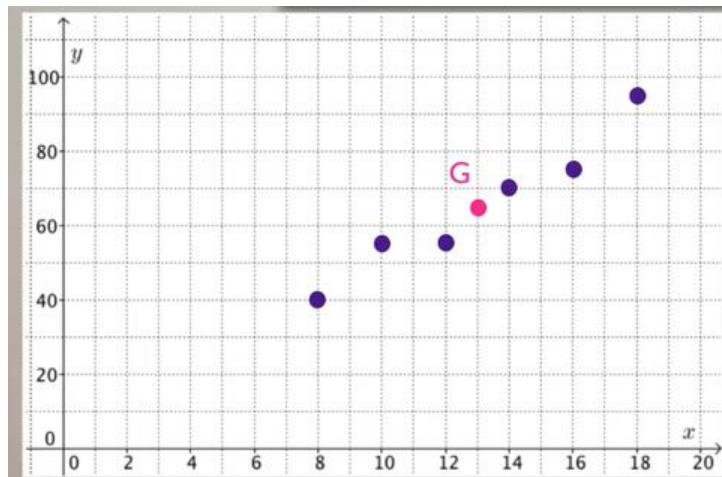
Example 2 :

The following table shows the evolution of a company's advertising budget and revenues over the last 6 years.

X Budget PUB x10 ³	8	10	12	14	16	18
Y Revenue x10 ³	40	55	55	70	75	95

1. In a coordinate system, represent the plot scatter (x,y).
2. Determine the coordinates of the mean point G of the plot scatter.

Solution:

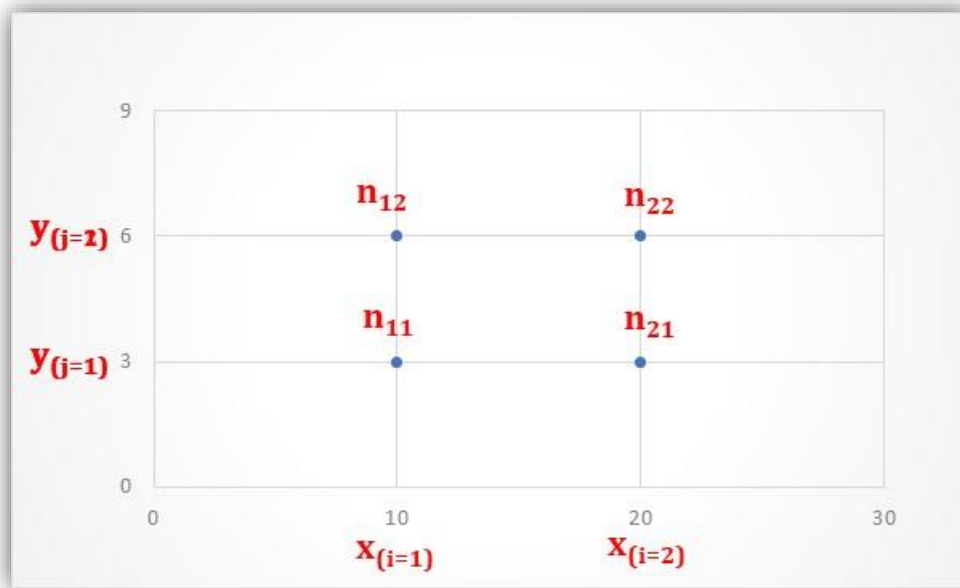


$$\bar{x} = \frac{8 + 10 + 12 + 14 + 16 + 18}{6} = 13$$

$$\bar{y} = \frac{40 + 55 + 55 + 70 + 75 + 95}{6} = 65$$

A.3.1.2. Weighted Scatter plot:

X	Y	y _(j=1) =3	y _(j=2) =6	Total
x _(i=1) =10		n11	n12	n1.
x _(i=2) =20		n21	n22	n2.
Total		n.1	n.2	N=n..



A.3.2 Marginal and conditional distributions:

By examining the marginal distributions of the contingency table, we can obtain data solely for variable X or variable Y

A.3.2.1 Marginal distributions of X :

We take the first column and the last column. The counts ($n_{i.}$) represent individuals with modality x_i independently of the modalities of the second character studied Y. They are reported to define the marginal distribution of X. (This statistical series is a single-character statistical series).

Counts and marginal frequencies with respect to X : we have, for $i = 1 \dots p$,

$$f_{i.} = \frac{n_{i.}}{N = n_{..}}$$

X	$n_{i.}$	$f_{i.}$
x1	n1.	f1.
x2	n2.	f2.
x3	n3.	f3.
.	.	.
xp	np.	fp.
Total	$N=n_{..}$	$F=f_{..}=1$

A.3.2.2 Marginal distributions of Y:

We take the first line and the last line. Counts and marginal frequencies with respect to Y : we have, for $i = 1 \dots q$,

$$f \cdot j = \frac{n \cdot j}{N = n \dots}$$

Y	n.j	f.j
y1	n.1	f.1
y2	n.2	f.2
y3	n.3	f.3
.	.	.
yq	n.q	f.q
Total	N=n..	F=f..=1

A.3.2.3 Conditional distributions of X:

We take the first column and the 2nd column will be taken according to the value of j . Notation $X/(Y=y_j)$. It is said to define the "conditional distribution of X given that $Y=y_j$ ".

X/(Y=y _j)	n _{i/j}	f _{i/j}
x1	n1j	f1j
x2	n2j	f2j
x3	n3j	f3j
.	.	.
x _p	npj	fpj
Total	N=n.j	F=f.j

A.3.2.4 Conditional distributions of Y:

We take the first line and the 2nd line will be taken according to the value of i. Notation $Y/(X=x_i)$. It is said to define the "Conditional distribution of Y given that $X=x_i$ ".

Example 3:

- Déterminer la distribution conditionnelle de X pour $Y_{(j=2)}$.

$\begin{matrix} Y \\ X \end{matrix}$	3	6	Total
10	n11	n12	n1.
20	n21	n22	n2.
Total	n.1	n.2	N=n..

Solution :

X	$\begin{matrix} n_{i/j=2} \\ y_{(j=2)=6} \end{matrix}$	$f_{i/j=2}$
$x_{(i=1)}=10$	n12	f12
$x_{(i=2)}=20$	n22	f22
Total	n.2	f..

A.3.2.5 Statistical independence:

X and Y are independent if the counts observed n_{ij} are identical to the counts of Independence.

$$n_{ij} = \frac{(n_{i.} * n_{.j})}{n_{..}}$$

or

$$f_{ij} = f_{i.} * f_{.j}$$

Example 3 :

A survey was carried out on 100 families by observing "monthly expenditure" X and "monthly income" Y (in thousands of DA), the results are given in the table below.

X \ Y	[25-40[[40-55[[55-70]	Total
[20-25[$n_{11}=30$	$n_{12}=20$	$n_{13}=10$	$n_{1\bullet}=60$
[25-30[$n_{21}=10$	$n_{22}=30$	$n_{23}=10$	$n_{2\bullet}=50$
[30-35]	$n_{31}=10$	$n_{32}=10$	$n_{33}=20$	$n_{3\bullet}=40$
Total	$n_{\bullet 1}=50$	$n_{\bullet 2}=60$	$n_{\bullet 3}=40$	$n_{\bullet\bullet}=150$

- Check the dependence of the variables X and Y for $i=2$ and $j=1$.

Solution:

For the independence of the variables X and Y, for $i=2$ and $j=1$, we obtain:

$$n_{21} = \frac{(n_{2\bullet} * n_{\bullet 1})}{n_{\bullet\bullet}}$$

$$n_{21} * n_{\bullet\bullet} = n_{2\bullet} * n_{\bullet 1}$$

$$n_{21} * n_{\bullet\bullet} = 10 * 150 = 1500$$

$$n_{2\bullet} * n_{\bullet 1} = 50 * 50 = 2500$$

Since $n_{21} * n_{\bullet\bullet} \neq n_{2\bullet} * n_{\bullet 1}$, Therefore, X and Y are not independent