

UNIT 6. BIVARIATE FREQUENCY DISTRIBUTIONS AND GRAPHIC PRESENTATIONS

In the previous units we have dealt with the statistical treatment of a single variable. But when we face the study of a population (for example, the socioeconomic situation of families in a city), the usual approach is to use different features of its individuals (family size, expenses and income, level of education, etc.). This not only provides a wider and more rewarding view of reality, but also makes it possible to study whether the different variables influence each other, i.e., "*Does family income affect the structure of family expenditure?*".

We begin by presenting the notation necessary for the joint analysis of two variables.

Bivariate Joint Frequency Distribution

When we jointly observe two variables X and Y from a population of size N , we have a two-dimensional statistical variable represented by (X, Y) . Let's suppose that the variable X takes k values x_i ($i=1, \dots, k$) and that the variable Y takes values y_j ($j=1, \dots, h$). We can define the following frequencies:

- **Absolute joint frequency (n_{ij}):** number of times that the pair (x_i, y_j) has been observed:

$$\sum_{i=1}^k \sum_{j=1}^h n_{ij} = N$$

- **Relative joint frequency (f_{ij}):** percentage of times that the pair (x_i, y_j) has been observed:

$$f_{ij} = \frac{n_{ij}}{N}$$

verifying that $\sum_{i=1}^k \sum_{j=1}^h f_{ij} = 1$

- **Bivariate joint absolute frequency distribution:** is defined by the three values (x_i, y_j, n_{ij}) $i = 1 \dots k; j = 1 \dots h$

The numerical representation of the data of a joint frequency distribution is done by means of a double-entry table called *cross-table*. The values of the variables X and Y are represented in the margins and the frequency of each couple of values is represented in the intersection cell. In addition, the size of the population can be calculated as:

$$\sum_{i=1}^k \sum_{j=1}^h n_{ij} = N$$

X\Y	y ₁	y ₂	...	y _j	...	y _h
x ₁	n ₁₁	n ₁₂	...	n _{1j}		n _{1h}
x ₂	n ₂₁	n ₂₂	...	n _{2j}		n _{2h}
...
x _i	n _{i1}	n _{i2}	...	n _{ij}		n _{ih}
...
x _k	n _{k1}	n _{k2}	...	n _{kj}		n _{kh}

- **Bivariate joint relative frequency distribution:** is defined by the three values $(x_i, y_j, f_{ij}) \quad i = 1 \dots k; j = 1 \dots h$

where

$$f_{ij} = \frac{n_{ij}}{N} \quad \sum_{j=1}^h f_{ij} = 1$$

The values of the variables X and Y are represented in the margins and the relative frequency of each couple of values is represented in the intersection cell.

X\Y	y ₁	y ₂	...	y _j	...	y _h
x ₁	f ₁₁	f ₁₂	...	f _{1j}		f _{1h}
x ₂	f ₂₁	f ₂₂	...	f _{2j}		f _{2h}
...
x _i	f _{i1}	f _{i2}	...	f _{ij}		f _{ih}
...
x _k	f _{k1}	f _{k2}	...	f _{kj}		f _{kh}

Example:

Given the following variables:

X = Monthly salary of the workers of a company (in euros)

Y = Education level of the workers of a company

The cross table of bivariate joint absolute frequency distribution of the example is represented as:

	Until secondary 1st stage	Secondary 2nd stage	High education
1000 - 1500	18	9	3
1500 - 2000	15	24	6
2000 - 2500	6	19	15
2500 - 3000	1	3	6

$$N = 125 \rightarrow f_{ij} = \frac{n_{ij}}{125}$$

	Until secondary 1st stage	Secondary 2nd stage	High education
1000 - 1500	14.4%	7.2%	2.4%
1500 - 2000	12.0%	19.2%	4.8%
2000 - 2500	4.8%	15.2%	12.0%
2500 - 3000	0.8%	2.4%	4.8%

By using this joint representation, we are going to be able to answer different questions related to both variables. For instance:

- How many workers are there with a higher education level?
- What percentage of workers earn more than €2,500?
- What is approximately the average salary?
- What is the most frequent level of education?

For this, we need to introduce a new concept: marginal distributions

Marginal distributions

From the joint frequency distribution of two variables, univariate distributions can be obtained for each of the variables separately. These distributions are called **marginal distributions** and are obtained by considering the values taken by one of the variables with their respective frequencies, independently of the values of the other variable. Therefore, marginal distributions can be characterised by all the measures of position, dispersion and shape studied in the previous lessons.

Marginal frequencies are obtained by **adding** the **joint frequencies** (absolute or relative) by rows or columns.

X \ Y	y ₁	y ₂	...	y _j	...	y _h	n _{i.}
x ₁	n ₁₁	n ₁₂	...	n _{1j}	...	n _{1h}	n _{1.}
...
x _i	n _{i1}	n _{i2}	...	n _{ij}	...	n _{ih}	n _{i.}
...
x _k	n _{k1}	n _{k2}	...	n _{kj}	...	n _{kh}	n _{k.}
n _{.j}	n _{.1}	n _{.2}	...	n _{.j}	...	n _{.h}	N

- **Absolute marginal frequency of X (n_{i.}):** number of times the value x_i is observed, regardless of the values of Y

$$n_{i.} = \sum_{j=1}^h n_{ij} \quad \sum_{i=1}^k n_{i.} = N$$

In the cross-table, the marginal distribution of X is obtained by adding the values of each row:

X\Y	y ₁	y ₂	...	y _j	...	y _h	n _{i.}
x ₁	n ₁₁	n ₁₂	...	n _{1j}	...	n _{1h}	n _{1.}
x ₂	n ₂₁	n ₂₂	...	n _{2j}	...	n _{2h}	n _{2.}
...
x _i	n _{i1}	n _{i2}	...	n _{ij}	...	n _{ih}	n _{i.}
...
x _k	n _{k1}	n _{k2}	...	n _{kj}	...	n _{kh}	n _{k.}
n _{.j}	n _{.1}	n _{.2}	...	n _{.j}	...	n _{.h}	N

- **Absolute marginal frequency of Y (n_{.j}):** number of times the value y_j is observed, regardless of the values of X

$$n_{.j} = \sum_{i=1}^h n_{ij} \quad \sum_{j=1}^k n_{.j} = N$$

In the cross-table, the marginal distribution of Y is obtained by adding the values of each column:

X\Y	y ₁	y ₂	...	y _j	...	y _h	n _{i.}
x ₁	n ₁₁	n ₁₂	...	n _{1j}	...	n _{1h}	n _{1.}
x ₂	n ₂₁	n ₂₂	...	n _{2j}	...	n _{2h}	n _{2.}
...
x _i	n _{i1}	n _{i2}	...	n _{ij}	...	n _{ih}	n _{i.}
...
x _k	n _{k1}	n _{k2}	...	n _{kj}	...	n _{kh}	n _{k.}
n _{.j}	n _{.1}	n _{.2}	...	n _{.j}	...	n _{.h}	N

From the previous *example*:

	Until secondary 1st stage	Secondary 2nd stage	High education	
1000 - 1500	14.4%	7.2%	2.4%	24.0%
1500 - 2000	12.0%	19.2%	4.8%	36.0%
2000 - 2500	4.8%	15.2%	12.0%	32.0%
2500 - 3000	0.8%	2.4%	4.8%	8.0%
	32.0%	44.0%	24.0%	100.0%

We can also express the marginal distributions in relative terms:

$X \backslash Y$	Y_1	Y_2	...	Y_j	...	Y_h	$f_{i.}$
x_1	f_{11}	f_{12}	...	f_{1j}	...	f_{1h}	$f_{1.}$
x_2	f_{21}	f_{22}	...	f_{2j}	...	f_{2h}	$f_{2.}$
...
x_i	f_{i1}	f_{i2}	...	f_{ij}	...	f_{ih}	$f_{i.}$
...
x_k	f_{k1}	f_{k2}	...	f_{kj}	...	f_{kh}	$f_{k.}$
$f_{.j}$	$f_{.1}$	$f_{.2}$...	$f_{.j}$...	$f_{.h}$	1

- Relative marginal frequency of X ($f_{i.}$):

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

$$f_{i.} = \sum_{j=1}^h f_{ij}$$

$$\sum_{i=1}^k f_{i.} = 1$$

- Relative marginal frequency of Y ($f_{.j}$):

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

$$f_{.j} = \sum_{i=1}^k f_{ij}$$

$$\sum_{j=1}^h f_{.j} = 1$$

Conditional distributions

Conditional distributions can be defined as the distribution of one of the two variables when a particular value/modality has been observed in the other variable. Like the marginal distributions, they are univariate distributions to which the statistical techniques studied in the previous topics can be applied. We can differentiate between:

- Distributions of X given Y
 - When a **value of Y is fixed** and the **distribution of the variable X** is analysed.
 - Example:* analyse the salary (X) for workers with a certain level of education (a given value of Y).
- Distributions of Y given X

- When a **value of X is fixed** and the **distribution of the variable Y** is analysed.
- *Example:* analyse the level of education (Y) for workers with a given amount of salary (a given value of X).

Distributions of X given that Y=y_j

For a bivariate distribution $(x_i, y_j, n_{ij}) \quad i = 1 \dots k; j = 1 \dots h$, the distribution of X given y_j is given by:

$$(x_i, n_{ij}) \quad i = 1 \dots k \text{ or } (x_i, f_{X=x_i|Y=y_j}) \quad i = 1 \dots k$$

X\Y	y ₁	y ₂	...	y _j	...	y _h	n _{i.}
x ₁	n ₁₁	n ₁₂	...	n _{1j}	...	n _{1h}	
x ₂	n ₂₁	n ₂₂	...	n _{2j}	...	n _{2h}	
...	
x _i	n _{i1}	n _{i2}	...	n _{ij}	...	n _{ih}	
...	
x _k	n _{k1}	n _{k2}	...	n _{kj}	...	n _{kh}	
n _{.j}				n _{.j}			

➔

X\Y	y _j
x ₁	n _{1j}
x ₂	n _{2j}
...	...
x _i	n _{ij}
...	...
x _k	n _{kj}
n _{.j}	n _{.j}

$$f_{X=x_i|Y=y_j} = \frac{n_{ij}}{n_{.j}} \quad \forall i = 1, \dots, k \quad \sum_{i=1}^k f_{i/j} = 1$$

X\Y	y ₁	y ₂	...	y _j	...	y _h	Marginal X
x ₁	f _{X=x1 Y=y1}	f _{X=x1 Y=y2}	...	f _{X=x1 Y=yj}	...	f _{X=x1 Y=yh}	f _{1.}
x ₂	f _{X=x2 Y=y1}	f _{X=x2 Y=y2}	...	f _{X=x2 Y=yj}	...	f _{X=x2 Y=yh}	f _{2.}
...
x _i	f _{X=xi Y=y1}	f _{X=xi Y=y2}	...	f _{X=xi Y=yj}	...	f _{X=xi Y=yh}	f _{i.}
...
x _k	f _{X=xk Y=y1}	f _{X=xk Y=y2}	...	f _{X=xk Y=yj}	...	f _{X=xk Y=yh}	f _{k.}
Total	1	1	...	1	...	1	1

The set of distributions of X, using relative frequencies, for each value/modality of the variable Y (Y=y_j j=1,...,h), is known as the **column profiles** (the total sum of each column equals 1).

Example of column profiles:

Distributions of salaries given the level of education.

	Until secondary 1st stage	Secondary 2nd stage	High education	f_i
1000-1500	45.0%	16.4%	10.0%	24.0%
1500-2000	37.5%	43.6%	20.0%	36.0%
2000-2500	15.0%	34.5%	50.0%	32.0%
2500-3000	2.5%	5.5%	20.0%	8.0%
	100.0%	100.0%	100.0%	100.0%

Distributions of Y given that $X=x_i$

For a bivariate distribution $(x_i, y_j, n_{ij}) \quad i = 1 \dots k; j = 1 \dots h$, the distribution of Y given x_i is given by:

$$(y_j, n_{ij}) \quad j = 1, \dots, h \quad \text{or} \quad (y_j, f_{Y=y_j|X=x_i}) \quad j = 1 \dots h$$

X \ Y	y_1	y_2	...	y_j	...	y_h	n_i
x_1	n_{11}	n_{12}	...	n_{1j}	...	n_{1h}	
x_2	n_{21}	n_{22}	...	n_{2j}	...	n_{2h}	
...	
x_i	n_{i1}	n_{i2}	...	n_{ij}	...	n_{ih}	n_i
...	
x_k	n_{k1}	n_{k2}	...	n_{kj}	...	n_{kh}	
n_j							



X \ Y	y_1	y_2	...	y_j	...	y_h	n_i
x_i	n_{i1}	n_{i2}	...	n_{ij}	...	n_{ih}	n_i

$$f_{j|i} = f_{Y=y_j|X=x_i} = \frac{n_{ij}}{n_i} \quad \forall j = 1, \dots, h \quad \sum_{j=1}^h f_{j|i} = 1$$

$X \backslash Y$	y_1	y_2	...	y_j	...	y_h	Total
x_1	$f_{Y=y_1 X=x_1}$	$f_{Y=y_2 X=x_1}$...	$f_{Y=y_j X=x_1}$...	$f_{Y=y_h X=x_1}$	1
x_2	$f_{Y=y_1 X=x_2}$	$f_{Y=y_2 X=x_2}$...	$f_{Y=y_j X=x_2}$...	$f_{Y=y_h X=x_2}$	1
...
x_i	$f_{Y=y_1 X=x_i}$	$f_{Y=y_2 X=x_i}$...	$f_{Y=y_j X=x_i}$...	$f_{Y=y_h X=x_i}$	1
...
x_k	$f_{Y=y_1 X=x_k}$	$f_{Y=y_2 X=x_k}$...	$f_{Y=y_j X=x_k}$...	$f_{Y=y_h X=x_k}$	1
Marginal Y	$f_{.1}$	$f_{.2}$...	$f_{.j}$...	$f_{.h}$	1

The set of distributions of X, using relative frequencies, for each value/modality of the variable X ($X=x_i$ $i=1, \dots, k$), is known as the **row profiles** (the total sum of each row equals 1).

Example of row profiles:

Distributions of the level of education for each range of salaries.

	Until secondary 1st stage	Secondary 2nd stage	High education	
1000-1500	60,0%	30.0%	10.0%	100.0%
1500-2000	33.3%	53.3%	13.3%	100.0%
2000-2500	15.0%	47.5%	37.5%	100.0%
2500-3000	10.0%	30.0%	60.0%	100.0%
f_j	32.0%	44.0%	24.0%	100.0%

Statistical Dependence and Independence

Two variables X and Y are said to be **independent** if the values of one of them are not affected by the values of the other

X and Y are statistically independent if one of the following conditions is fulfilled:

- | | | |
|--|---|---|
| <ol style="list-style-type: none"> 1. $f_{i j} = f_i \quad \forall j = 1 \dots h$ 2. $f_{j i} = f_j \quad \forall i = 1 \dots k$ 3. $f_{ij} = f_i \times f_j \quad \forall i = 1 \dots k, \quad \forall j = 1 \dots h$ | } | <p>Row and column profiles are equal to the corresponding marginal distributions</p> |
|--|---|---|

All the three conditions are equivalent. If one of them is verified, the other two are verified. These are their meanings:

1. All the conditional distributions of X are the same as the marginal distribution of X
2. All the conditional distributions of Y are the same as the marginal distribution of Y
3. The relative joint frequencies coincide with the product of the corresponding relative marginal frequencies

Types of dependence

When variables are not independent, there are various possible degrees of dependence between two variables X and Y :

- **Functional dependence:** one of the variables can be expressed as a function of the other: $Y = f(X)$
- **Statistical dependence:** there exists a certain relationship between the variables (intermediate situations)

Example:

Distributions of X given Y (column profiles)

	Until secondary 1st stage	Secondary 2nd stage	High education	$f_{.i}$
1000-1500	45.0%	16.4%	10.0%	24.0%
1500-2000	37.5%	43.6%	20.0%	36.0%
2000-2500	15.0%	34.5%	50.0%	32.0%
2500-3000	2.5%	5.5%	20.0%	8.0%
	100.0%	100.0%	100.0%	100.0%

Column profiles \neq Marginal distribution of X \rightarrow Variables X and Y are NOT statistically independent

Distributions of Y given X (row profiles)

	Until secondary 1st stage	Secondary 2nd stage	High education	
1000-1500	60.0%	30.0%	10.0%	100.0%
1500-2000	33.3%	53.3%	13.3%	100.0%
2000-2500	15.0%	47.5%	37.5%	100.0%
2500-3000	10.0%	30.0%	60.0%	100.0%
$f_{.i}$	32.0%	44.0%	24.0%	100.0%

Row profiles \neq Marginal distribution of Y \rightarrow Variables X and Y are NOT statistically independent

Joint distribution of salary and level of education

	Until secondary 1st stage	Secondary 2nd stage	High education	
1000 - 1500	14.4%	7.2%	2.4%	24.0%
1500 - 2000	12.0%	19.2%	4.8%	36.0%
2000 - 2500	4.8%	15.2%	12.0%	32.0%
2500 - 3000	0.8%	2.4%	4.8%	8.0%
	32.0%	44.0%	24.0%	100.0%

It is necessary to check if $f_{ij} = f_{i.} \times f_{.j} \quad \forall i = 1, \dots, k \quad j = 1, \dots, h$

$$i = 1, j = 1 \quad f_{11} = 14.4 \% \quad f_{1.} \times f_{.1} = 32.0 \% \times 24.0 \% = 7.68 \%$$

$$i = 1, j = 2 \quad f_{12} = 7.2 \% \quad f_{1.} \times f_{.2} = 44.0 \% \times 24.0 \% = 10.56 \%$$

$$i = 1, j = 3 \quad f_{13} = 2.4 \% \quad f_{1.} \times f_{.3} = 24.0 \% \times 24.0 \% = 5.76 \%$$

$$i = 2, j = 1 \quad f_{21} = 12.0 \% \quad f_{2.} \times f_{.1} = 32.0 \% \times 36.0 \% = 11.52 \%$$



$$\begin{aligned}i = 2, j = 2 & \quad f_{22} = 19.2 \% \quad f_{2.} \times f_{.2} = 44.0 \% \times 36.0 \% = 15.84 \% \\i = 2, j = 3 & \quad f_{12} = 4.8 \% \quad f_{1.} \times f_{.2} = 24.0 \% \times 36.0 \% = 8.64 \% \\i = 3, j = 1 & \quad f_{21} = 4.8 \% \quad f_{2.} \times f_{.1} = 32.0 \% \times 32.0 \% = 10.24 \% \\i = 3, j = 2 & \quad f_{22} = 15.2 \% \quad f_{2.} \times f_{.2} = 44.0 \% \times 32.0 \% = 14.08 \% \\i = 3, j = 3 & \quad f_{12} = 12.0 \% \quad f_{1.} \times f_{.2} = 24.0 \% \times 32.0 \% = 7.68 \% \\i = 4, j = 1 & \quad f_{21} = 0.8 \% \quad f_{2.} \times f_{.1} = 32.0 \% \times 8.0 \% = 2.56 \% \\i = 4, j = 2 & \quad f_{22} = 2.4 \% \quad f_{2.} \times f_{.2} = 44.0 \% \times 8.0 \% = 3.52 \% \\i = 4, j = 3 & \quad f_{12} = 4.8 \% \quad f_{1.} \times f_{.2} = 24.0 \% \times 8.0 \% = 1.92 \%\end{aligned}$$

As can be seen, it is proven that the variables X and Y are NOT statistically independent.