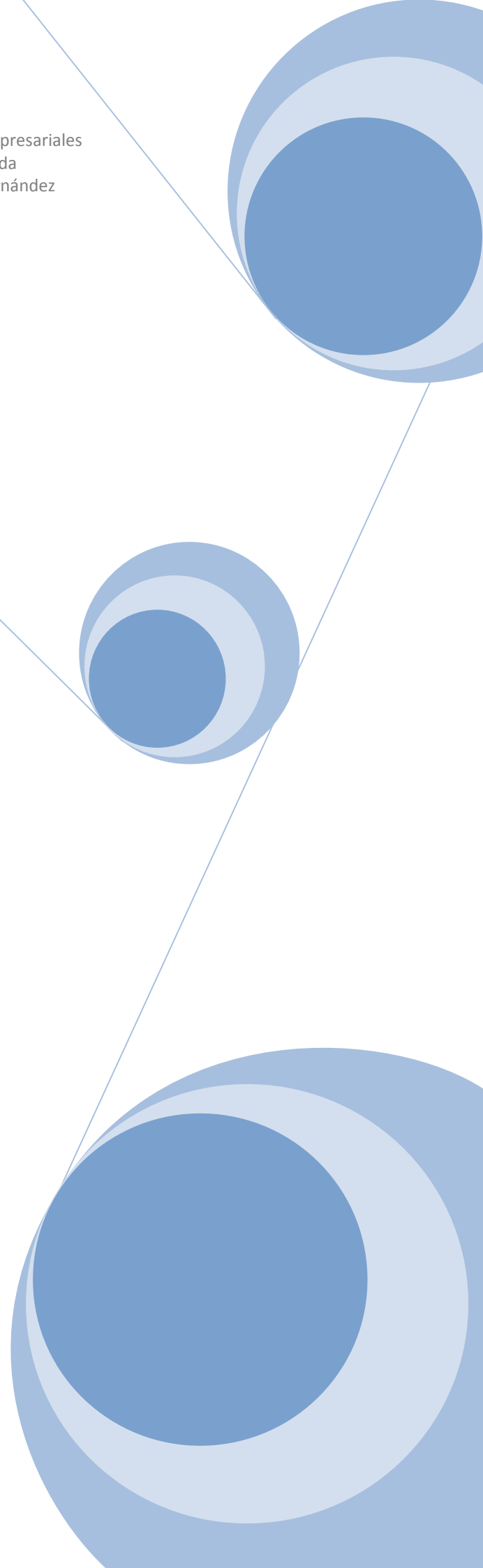




Grado Administración y Gestión
Facultad Ciencias Económicas y Empresariales
Departamento de Economía Aplicada
Profesor: Santiago de la Fuente Fernández

ESTADÍSTICA BIDIMENSIONAL



VARIABLE ESTADÍSTICA BIDIMENSIONAL

Cuando se consideran situaciones en la que el estadístico realiza la observación simultánea de dos caracteres en el individuo, se obtienen pares de resultados.

Los distintos valores de las modalidades que pueden adoptar estos caracteres forman un conjunto de pares, que representamos por (X, Y) , y llamaremos *variable estadística bidimensional*.

Los dos caracteres observados no tienen por qué ser de la misma clase, pudiendo presentarse distintas situaciones:

- **Dos caracteres cualitativos:** El sexo y color del pelo de una persona.
- **Dos caracteres cuantitativos:** El peso y la estatura de una persona.
- **Uno cuantitativo y otro cualitativo:** La profesión y los años de servicio.

Las variables (X, Y) que representan los valores de dos caracteres cuantitativos, pueden clasificarse:

- **X discreta e Y discreta:** Número de hijos y número de hermanos de una persona.
- **X continua e Y continua:** Perímetro craneal y perímetro torácico de una persona.
- **X discreta e Y continua:** Hijos de una familia y estatura del padre.
- **X continua e Y discreta:** Temperatura y pulsaciones.

ORDENACIÓN DE LOS DATOS: TABLA DE DOBLE ENTRADA

El par (X, Y) es la unidad del estudio y dos pares serán repetidos solo cuando sus respectivas componentes sean iguales. De otra parte, el número de modalidades que adopta el carácter X no tiene por qué ser el mismo que el que adopta el carácter Y:

$$X = (x_1, x_2, \dots, x_k) \quad Y = (y_1, y_2, \dots, y_m)$$

Para ordenar los datos se utiliza una tabla de doble entrada donde tengan cabida los k valores distintos de la variable X y los m valores distintos de la variable Y. En la tabla se puede expresar el número de veces que se repite cada par de valores posibles (x_i, y_j) formado en el producto cartesiano de los dos conjuntos numéricos.

TABLA DE DOBLE ENTRADA

| Y | y_1 | y_2 | ... | y_j | ... | y_m |
|-------|----------|----------|-----|------------|-----|----------|
| X | | | | | | |
| x_1 | n_{11} | n_{12} | ... | ⋮ | ... | n_{1m} |
| x_2 | n_{21} | n_{22} | ... | ⋮ | ... | n_{2m} |
| ... | ... | ... | ... | ⋮ | ... | ... |
| x_i | ⋮ | ⋮ | ⋮ | ⋮ n_{ij} | ... | n_{im} |
| ... | ... | ... | ... | | ... | ... |
| x_k | n_{k1} | n_{k2} | ... | | ... | n_{km} |

$N \equiv$ número total observaciones

$n_{ij} \equiv$ frecuencia absoluta, número de veces que aparece repetido el par (x_i, y_j) .

La frecuencia relativa del par se

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

DISTRIBUCIONES MARGINALES

| X \ Y | y ₁ | y ₂ | ... | y _j | ... | y _m | n _{i•} |
|-----------------|-----------------|-----------------|-----|-----------------|-----|-----------------|---|
| x ₁ | n ₁₁ | n ₁₂ | ... | n _{1j} | ... | n _{1m} | n _{1•} |
| x ₂ | n ₂₁ | n ₂₂ | ... | n _{2j} | ... | n _{2m} | n _{2•} |
| ... | ... | ... | ... | ... | ... | ... | ... |
| x _i | n _{i1} | n _{i2} | ... | n _{ij} | ... | n _{im} | n _{i•} |
| ... | ... | ... | ... | ... | ... | ... | ... |
| x _k | n _{k1} | n _{k2} | ... | n _{kj} | ... | n _{km} | n _{k•} |
| n _{•j} | n _{•1} | n _{•2} | ... | n _{•j} | ... | n _{•m} | $N = \sum_{i=1}^k n_{i•} = \sum_{j=1}^m n_{•j}$ |

$$a_{11} = \frac{\sum_{i=1}^k \sum_{j=1}^m x_i \cdot y_j \cdot n_{ij}}{N} \quad s_{yx} = a_{11} - a_{10} \cdot a_{01} = a_{11} - \bar{x} \cdot \bar{y} \quad (\text{covarianza})$$

• DISTRIBUCIÓN MARGINAL DE LA VARIABLE X

| X | x ₁ | x ₂ | ... | x _i | ... | x _k |
|-----------------------------|-----------------|-----------------|-----|-----------------|-----|-----------------|
| n _{i•} | n _{1•} | n _{2•} | ... | n _{i•} | ... | n _{k•} |
| $f_{i•} = \frac{n_{i•}}{N}$ | f _{1•} | f _{2•} | ... | f _{i•} | ... | f _{k•} |

$$N = \sum_{i=1}^k n_{i•}$$

$$\sum_{i=1}^k f_{i•} = 1$$

$$a_{10} = \bar{x} = \frac{\sum_{i=1}^k x_i \cdot n_{i•}}{N} \quad a_{20} = \frac{\sum_{i=1}^k x_i^2 \cdot n_{i•}}{N} \quad m_{20} = s_x^2 = a_{20} - (a_{10})^2 = a_{20} - (\bar{x})^2$$

• DISTRIBUCIÓN MARGINAL DE LA VARIABLE Y

| Y | y ₁ | y ₂ | ... | y _j | ... | y _m |
|-----------------------------|-----------------|-----------------|-----|-----------------|-----|-----------------|
| n _{•j} | n _{•1} | n _{•2} | ... | n _{•j} | ... | n _{•m} |
| $f_{•j} = \frac{n_{•j}}{N}$ | f _{•1} | f _{•2} | ... | f _{•j} | ... | f _{•m} |

$$N = \sum_{j=1}^m n_{•j}$$

$$\sum_{j=1}^m f_{•j} = 1$$

$$a_{01} = \bar{y} = \frac{\sum_{j=1}^m y_j \cdot n_{•j}}{N} \quad a_{02} = \frac{\sum_{j=1}^m y_j^2 \cdot n_{•j}}{N} \quad m_{02} = s_y^2 = a_{02} - (a_{01})^2 = a_{02} - (\bar{y})^2$$

Las variables (X, Y) son independientes cuando: $\frac{n_{ij}}{N} = \left(\frac{n_{i•}}{N}\right) \left(\frac{n_{•j}}{N}\right) \quad \forall i, j$

Si (X,Y) independientes $\mapsto s_{yx} = 0$
 Si $s_{yx} = 0 \mapsto (X,Y)$ No independientes

• **DISTRIBUCIÓN CONDICIONADA DE LA VARIABLE X para un valor $Y = y_j$**

| | Y | | | | | | |
|-----------------|-----------------|-----------------|-----|-----------------|-----|-----------------|--|
| X | y_1 | y_2 | ... | y_j | ... | y_m | $n_{i\bullet}$ |
| x_1 | n_{11} | n_{12} | ... | n_{1j} | ... | n_{1m} | $n_{1\bullet}$ |
| x_2 | n_{21} | n_{22} | ... | n_{2j} | ... | n_{2m} | $n_{2\bullet}$ |
| ... | ... | ... | ... | ... | ... | ... | ... |
| x_i | n_{i1} | n_{i2} | ... | n_{ij} | ... | n_{im} | $n_{i\bullet}$ |
| ... | ... | ... | ... | ... | ... | ... | ... |
| x_k | n_{k1} | n_{k2} | ... | n_{kj} | ... | n_{km} | $n_{k\bullet}$ |
| $n_{\bullet j}$ | $n_{\bullet 1}$ | $n_{\bullet 2}$ | ... | $n_{\bullet j}$ | ... | $n_{\bullet m}$ | $N = \sum_{i=1}^k n_{i\bullet} = \sum_{j=1}^m n_{\bullet j}$ |

| X | x_1 | x_2 | ... | x_i | ... | x_k |
|----------------|----------|----------|-----|----------|-----|----------|
| $n(X/Y = y_j)$ | n_{1j} | n_{2j} | ... | n_{ij} | ... | n_{kj} |
| $f(X/Y = y_j)$ | f_{1j} | f_{2j} | ... | f_{ij} | ... | f_{kj} |

$$f(X/Y = y_j) = \frac{n(X_i/Y = y_j)}{n_{\bullet j}}$$

• **DISTRIBUCIÓN CONDICIONADA DE LA VARIABLE Y para un valor $X = x_i$**

| | Y | | | | | | |
|-----------------|-----------------|-----------------|-----|-----------------|-----|-----------------|--|
| X | y_1 | y_2 | ... | y_j | ... | y_m | $n_{i\bullet}$ |
| x_1 | n_{11} | n_{12} | ... | n_{1j} | ... | n_{1m} | $n_{1\bullet}$ |
| x_2 | n_{21} | n_{22} | ... | n_{2j} | ... | n_{2m} | $n_{2\bullet}$ |
| ... | ... | ... | ... | ... | ... | ... | ... |
| x_i | n_{i1} | n_{i2} | ... | n_{ij} | ... | n_{im} | $n_{i\bullet}$ |
| ... | ... | ... | ... | ... | ... | ... | ... |
| x_k | n_{k1} | n_{k2} | ... | n_{kj} | ... | n_{km} | $n_{k\bullet}$ |
| $n_{\bullet j}$ | $n_{\bullet 1}$ | $n_{\bullet 2}$ | ... | $n_{\bullet j}$ | ... | $n_{\bullet m}$ | $N = \sum_{i=1}^k n_{i\bullet} = \sum_{j=1}^m n_{\bullet j}$ |

| Y | y_1 | y_2 | ... | y_j | ... | y_m |
|----------------|----------|----------|-----|----------|-----|----------|
| $n(Y/X = x_i)$ | n_{i1} | n_{i2} | ... | n_{ij} | ... | n_{im} |
| $f(Y/X = x_i)$ | f_{i1} | f_{i2} | ... | f_{ij} | ... | f_{im} |

$$f(Y/X = x_i) = \frac{n(Y/X = x_i)}{n_{i\bullet}}$$

MOMENTOS

Se define el momento respecto al par de valores (c, v) de órdenes r y s :

$$M_{rs}(c, v) = \frac{\sum_{i=1}^k \sum_{j=1}^m (x_i - c)^r (y_j - v)^s n_{ij}}{N}$$

Tienen especial interés dos casos particulares para los valores c y v

- MOMENTOS RESPECTO AL ORIGEN $(c, v) = (0, 0)$**

$$a_{rs} = \frac{\sum_{i=1}^k \sum_{j=1}^m (x_i - 0)^r (y_j - 0)^s n_{ij}}{N} = \frac{\sum_{i=1}^k \sum_{j=1}^m x_i^r \cdot y_j^s \cdot n_{ij}}{N}$$

de interés son los particulares:

$$a_{00} = \frac{\sum_{i=1}^k \sum_{j=1}^m x_i^0 \cdot y_j^0 \cdot n_{ij}}{N} = \frac{\sum_{i=1}^k \sum_{j=1}^m n_{ij}}{N} = 1$$

$$a_{11} = \frac{\sum_{i=1}^k \sum_{j=1}^m x_i^1 \cdot y_j^1 \cdot n_{ij}}{N}$$

$$a_{10} = \frac{\sum_{i=1}^k \sum_{j=1}^m x_i^1 \cdot y_j^0 \cdot n_{ij}}{N} = \frac{\sum_{i=1}^k \sum_{j=1}^m x_i n_{ij}}{N} = \frac{\sum_{i=1}^k x_i n_{i\bullet}}{N} = \bar{x}$$

$$a_{01} = \frac{\sum_{i=1}^k \sum_{j=1}^m x_i^0 \cdot y_j^1 \cdot n_{ij}}{N} = \frac{\sum_{i=1}^k \sum_{j=1}^m y_j n_{ij}}{N} = \frac{\sum_{j=1}^m y_j n_{\bullet j}}{N} = \bar{y}$$

$$a_{20} = \frac{\sum_{i=1}^k \sum_{j=1}^m x_i^2 \cdot y_j^0 \cdot n_{ij}}{N} = \frac{\sum_{i=1}^k \sum_{j=1}^m x_i^2 n_{ij}}{N} = \frac{\sum_{i=1}^k x_i^2 n_{i\bullet}}{N}$$

$$a_{02} = \frac{\sum_{i=1}^k \sum_{j=1}^m x_i^0 \cdot y_j^2 \cdot n_{ij}}{N} = \frac{\sum_{i=1}^k \sum_{j=1}^m y_j^2 n_{ij}}{N} = \frac{\sum_{j=1}^m y_j^2 n_{\bullet j}}{N}$$

- MOMENTOS CENTRALES O RESPECTO A LAS MEDIAS $(c, v) = (\bar{x}, \bar{y})$**

$$m_{rs} = \frac{\sum_{i=1}^k \sum_{j=1}^m (x_i - \bar{x})^r (y_j - \bar{y})^s n_{ij}}{N}$$

de interés son los particulares:

$$m_{11} = s_{yx} = \frac{\sum_{i=1}^k \sum_{j=1}^m (x_i - \bar{x}) (y_j - \bar{y}) n_{ij}}{N} = s_{xy} \quad \text{covarianza}$$

$$m_{20} = \frac{\sum_{i=1}^k \sum_{j=1}^m (x_i - \bar{x})^2 (y_j - \bar{y})^0 n_{ij}}{N} = \frac{\sum_{i=1}^k \sum_{j=1}^m (x_i - \bar{x})^2 n_{ij}}{N} = \frac{\sum_{i=1}^k (x_i - \bar{x})^2 n_{i\bullet}}{N} = s_x^2 \quad \text{varianza de X}$$

$$m_{02} = \frac{\sum_{i=1}^k \sum_{j=1}^m (x_i - \bar{x})^0 (y_j - \bar{y})^2 n_{ij}}{N} = \frac{\sum_{i=1}^k \sum_{j=1}^m (y_j - \bar{y})^2 n_{ij}}{N} = \frac{\sum_{j=1}^m (y_j - \bar{y})^2 n_{\bullet j}}{N} = s_y^2 \text{ varianza de Y}$$

□ Se demuestra fácilmente que, $m_{11} = s_{xy} = \frac{\sum_{i=1}^k \sum_{j=1}^m (x_i - \bar{x}) (y_j - \bar{y}) n_{ij}}{N} = a_{11} - a_{10} \cdot a_{01} = a_{11} - \bar{x} \cdot \bar{y}$

$$\begin{aligned} m_{11} = s_{xy} &= \frac{\sum_{i=1}^k \sum_{j=1}^m (x_i - \bar{x}) (y_j - \bar{y}) n_{ij}}{N} = \frac{\sum_{i=1}^k \sum_{j=1}^m (x_i \cdot y_j - x_i \cdot \bar{y} - \bar{x} \cdot y_j + \bar{x} \cdot \bar{y}) n_{ij}}{N} = \\ &= \frac{\sum_{i=1}^k \sum_{j=1}^m x_i \cdot y_j \cdot n_{ij}}{N} - \bar{y} \cdot \frac{\sum_{i=1}^k \sum_{j=1}^m x_i \cdot n_{ij}}{N} - \bar{x} \cdot \frac{\sum_{i=1}^k \sum_{j=1}^m y_j \cdot n_{ij}}{N} + \bar{x} \cdot \bar{y} \cdot \frac{\sum_{i=1}^k \sum_{j=1}^m n_{ij}}{N} = \\ &= \frac{\sum_{i=1}^k \sum_{j=1}^m x_i \cdot y_j \cdot n_{ij}}{N} - \bar{y} \cdot \frac{\sum_{i=1}^k x_i \cdot n_{i\bullet}}{N} - \bar{x} \cdot \frac{\sum_{j=1}^m y_j \cdot n_{\bullet j}}{N} + \bar{x} \cdot \bar{y} \cdot \frac{\sum_{i=1}^k \sum_{j=1}^m n_{ij}}{N} = \\ &= a_{11} - \bar{y} \cdot \bar{x} - \bar{x} \cdot \bar{y} + \bar{x} \cdot \bar{y} = a_{11} - \bar{x} \cdot \bar{y} = a_{11} - a_{10} \cdot a_{01} \end{aligned}$$

DEPENDENCIA ENTRE LAS VARIABLES (X, Y)

Al observar dos caracteres en cada individuo se presenta el problema de determinar la existencia de algún tipo de dependencia entre ellos. En este sentido, conviene destacar dos tipos de dependencia:

- **Dependencia funcional:** Entre dos variables X e Y existe dependencia funcional cuando hay una expresión matemática que las relacione. Por ejemplo, los radios de una circunferencia (X) y las longitudes (Y).
- **Dependencia aleatoria:** Entre dos variables X e Y existe dependencia aleatoria cuando no existe una expresión matemática que las relacione. Por ejemplo, la edad de los niños (X) y la edad (Y).

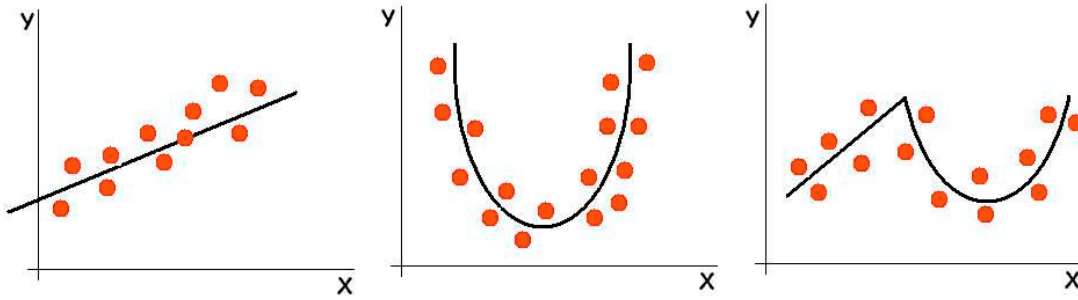
Señalar que existen variables entre las que no existe ningún tipo de dependencia, lo que conlleva a decir que los dos conceptos anteriores no son complementarios.

REGRESIÓN O AJUSTE

La observación de una variable estadística bidimensional (X, Y) comporta la representación de los puntos obtenidos en una **nube o diagrama de dispersión**. El problema general de regresión se plantea en el intento de ajustar una función de ecuación conocida (recta, parábola, exponencial, hipérbola, polinómica, etc.) a la nube de puntos con el interés de poder obtener una **predicción** aproximada de una de las variables a partir de la otra.

Naturalmente, que entre todas las funciones que se pueden elegir para ajustar a la nube de puntos, hemos de seleccionar la óptima, esto es, la que mejor encaje sobre los puntos que tenemos, para lo cual recurriremos al **método de los mínimos cuadrados**.

MÉTODO: Dependiendo de la forma que adopte la nube de puntos, en un principio sabremos si hemos de emplear una recta, una parábola, una función mixta, etc.



Una vez elegida la función, se estiman los parámetros correspondientes de la misma a partir de los datos observados. Por ejemplo, si la función elegida es una parábola:

$$y = a + b x + c x^2 \quad \text{hemos de estimar } a, b, c$$

Por último, una vez realizada la estimación hay que comprobar si efectivamente el ajuste era el idóneo o no. Para ello se emplean cualquiera de los tests construidos para estudiar la bondad del ajuste. El modelo más utilizado es el de la χ^2 (chi-cuadrado).

REGRESIÓN LINEAL MÍNIMO CUADRÁTICA

En el supuesto de que sea la recta la función que mejor se comporta con arreglo a la forma de la nube de puntos, nos encontramos ante una problema de regresión lineal, distinguiendo entre:

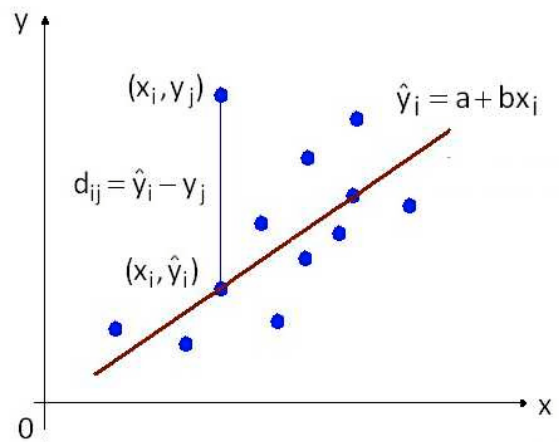
- **Recta de regresión de Y sobre X:** Obteniendo valores aproximados de la Y conocidos los de la X
- **Recta de regresión de X sobre Y:** Obteniendo valores aproximados de la X conocidos los de la Y

RECTA DE REGRESIÓN DE Y SOBRE X

En cada par (X,Y) al valor observado x_i le corresponde un valor observado y_j y otro valor teórico \hat{y}_i que sería el que le correspondería en la recta como función, es decir:

$$\hat{y}_i = a + b x_i$$

A la distancia entre estos dos valores (teórico y experimental), la denotamos por $d_{ij} = \hat{y}_i - y_j$



Para obtener los parámetros **a** y **b**, se toman las distancias (errores) al cuadrado para que no se contrarresten los signos positivos y negativos, haciendo mínima su suma: $M = \sum_{i,j} d_{i,j}^2 = \sum_{i,j} (\hat{y}_i - y_j)^2$

Por otra parte, para simplificar el mecanismo para obtener la recta de regresión de Y (*variable dependiente*) sobre X (*variable independiente*), se descartan multiplicidades y suponemos que cada par se repite una sola vez.

$$\text{Considerando que } \hat{y}_i = a + b x_i, \quad M = \sum_{i,j} d_{i,j}^2 = \sum_{i,j} (a + b x_i - y_j)^2$$

Para hallar los valores de **a** y **b** que hagan mínima esta función hemos de hallar las derivadas, igualando a cero las ecuaciones resultantes:

$$\frac{\partial M}{\partial a} = 2 \sum_{i,j} (a + bx_i - y_j) = 0 \quad \Rightarrow \quad \sum_{i,j} (a + bx_i - y_j) = 0$$

$$\frac{\partial M}{\partial b} = 2 \sum_{i,j} (a + bx_i - y_j)(x_i) = 0 \quad \Rightarrow \quad \sum_{i,j} (a + bx_i - y_j)(x_i) = 0$$

Por las propiedades del sumatorio, se obtienen las **ecuaciones normales** de la regresión:

$$\begin{cases} \sum_i a + b \sum_i x_i - \sum_j y_j = 0 \\ a \sum_i x_i + b \sum_i x_i^2 - \sum_{i,j} x_i y_j = 0 \end{cases} \Rightarrow \begin{cases} \sum_i a + b \sum_i x_i = \sum_j y_j \\ a \sum_i x_i + b \sum_i x_i^2 = \sum_{i,j} x_i y_j \end{cases}$$

Dividiendo las expresiones anteriores por N (número total de datos), habiendo supuesto que la frecuencia absoluta de cada par (X, Y) es la unidad, resulta:

$$\left. \begin{aligned} a \frac{\sum_i 1}{N} + b \frac{\sum_i x_i}{N} &= \frac{\sum_j y_j}{N} \\ a \frac{\sum_i x_i}{N} + b \frac{\sum_i x_i^2}{N} &= \frac{\sum_{i,j} x_i y_j}{N} \end{aligned} \right\} \begin{aligned} \text{Considerando los momentos, se tiene: } &\begin{cases} a + b\bar{x} = \bar{y} \\ a\bar{x} + ba_{20} = a_{11} \end{cases} \\ &a = \bar{y} - b\bar{x} \end{aligned}$$

sustituyendo en la ecuación $a\bar{x} + ba_{20} = a_{11}$, resulta:

$$(\bar{y} - b\bar{x})\bar{x} + ba_{20} = a_{11} \quad \mapsto \quad b(a_{20} - \bar{x}^2) = a_{11} - \bar{x}\bar{y} \quad \mapsto \quad b = \frac{a_{11} - \bar{x}\bar{y}}{a_{20} - \bar{x}^2} = \frac{s_{xy}}{s_x^2}$$

$$a = \bar{y} - \frac{s_{xy}}{s_x^2} \bar{x}$$

Finalmente, sustituyendo los valores obtenidos en la ecuación de la recta $y = a + bx$

$$y = \bar{y} - \frac{s_{xy}}{s_x^2} \bar{x} + \frac{s_{xy}}{s_x^2} x \quad \Rightarrow \quad y - \bar{y} = \frac{s_{xy}}{s_x^2} (x - \bar{x})$$

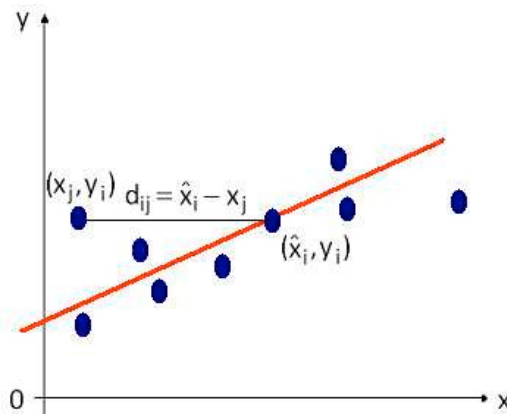
NOTA.- En el supuesto de que no hubiéramos partido de las hipótesis iniciales para el desarrollo, es decir, si hay multiplicidades de (x_i, y_j) y si cada par se repite n_{ij} veces, la ecuación a minimizar sería

$$M = \sum_{i,j} d_{i,j}^2 n_{ij} = \sum_{i,j} (a + bx_i - y_j)^2 n_{ij}$$

RECTA DE REGRESIÓN DE X SOBRE Y

Si en lugar de tomar las distancias d_{ij} sobre las verticales (esto es, sobre la Y) se toman sobre las horizontales (sobre la X) y se utiliza el mismo método de los mínimos cuadrados, por un proceso idénticamente igual se llega a la **ecuación de regresión de X sobre Y**:

$$x - \bar{x} = \frac{s_{xy}}{s_y^2} (y - \bar{y})$$



COEFICIENTES DE REGRESIÓN LINEAL

- La **recta de regresión de Y sobre X**: $y - \bar{y} = \frac{s_{xy}}{s_x^2} (x - \bar{x})$, donde el coeficiente de regresión lineal

$$b_{yx} = \frac{s_{xy}}{s_x^2} \text{ es la pendiente de la recta.}$$

Recta de regresión de Y sobre X, según el coeficiente de regresión $b_{yx} \equiv \begin{cases} > 0 & \text{creciente} \\ = 0 & \text{horizontal} \\ < 0 & \text{decreciente} \end{cases}$

- La **recta de regresión de X sobre Y**: $x - \bar{x} = \frac{s_{xy}}{s_y^2} (y - \bar{y})$, donde el coeficiente de regresión lineal

$$b_{xy} = \frac{s_{yx}}{s_y^2} \text{ es la pendiente de la recta.}$$

Recta de regresión de X sobre Y, según el coeficiente de regresión $b_{xy} \equiv \begin{cases} > 0 & \text{creciente} \\ = 0 & \text{vertical} \\ < 0 & \text{decreciente} \end{cases}$

CORRELACIÓN

Así como la regresión estudia la posible predicción de los valores de una variable a partir de la otra, la correlación estudia el tipo de dependencia que existe entre ambas variables, intentando cuantificarla mediante el cálculo de los coeficientes de correlación.

A continuación se estudian los coeficientes de determinación y correlación lineal.

COEFICIENTE DE CORRELACIÓN LINEAL

- El **coeficiente de correlación lineal** R es un número abstracto que determina el grado de ajuste entre una nube de puntos y una recta de regresión. Se define como la media geométrica de los coeficientes de correlación lineal

$$r = \sqrt{b_{yx} \cdot b_{xy}} = \sqrt{\frac{s_{xy}}{s_x^2} \frac{s_{xy}}{s_y^2}} = \frac{s_{xy}}{s_x s_y}$$

RELACIÓN ENTRE LOS COEFICIENTES DE REGRESIÓN Y DE CORRELACIÓN

- ◆ **Recta de regresión de Y sobre X:** $y - \bar{y} = \frac{s_{xy}}{s_x^2} (x - \bar{x})$, coeficiente de regresión lineal $b_{yx} = \frac{s_{xy}}{s_x^2}$

(pendiente de la recta), coeficiente de correlación $r = \frac{s_{xy}}{s_x s_y}$ (grado de ajuste)

$$\left. \begin{array}{l} b_{yx} = \frac{s_{xy}}{s_x^2} \\ r = \frac{s_{xy}}{s_x s_y} \end{array} \right\} \mapsto \left\{ \begin{array}{l} s_{xy} = b_{yx} s_x^2 \\ s_{xy} = r s_x s_y \end{array} \right\} \mapsto b_{yx} s_x^2 = r s_x s_y \mapsto \boxed{b_{yx} = r \frac{s_y}{s_x}}$$

se observa que los dos coeficientes (regresión lineal y correlación lineal) tienen el mismo signo.

- ◆ **Recta de regresión de X sobre Y:** $x - \bar{x} = \frac{s_{xy}}{s_y^2} (y - \bar{y})$, coeficiente de regresión lineal $b_{xy} = \frac{s_{xy}}{s_y^2}$

(pendiente de la recta), coeficiente de correlación $r = \frac{s_{xy}}{s_x s_y}$ (grado de ajuste)

$$\left. \begin{array}{l} b_{xy} = \frac{s_{xy}}{s_y^2} \\ r = \frac{s_{xy}}{s_x s_y} \end{array} \right\} \mapsto \left\{ \begin{array}{l} s_{xy} = b_{xy} s_y^2 \\ s_{xy} = r s_x s_y \end{array} \right\} \mapsto b_{xy} s_y^2 = r s_x s_y \mapsto \boxed{b_{xy} = r \frac{s_x}{s_y}}$$

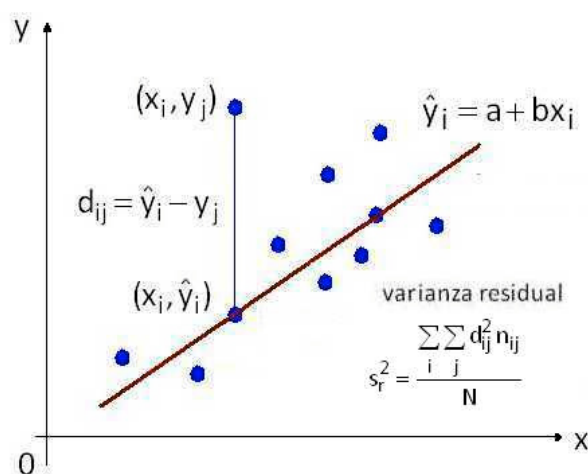
los dos coeficientes (regresión lineal y correlación lineal) tienen el mismo signo.

VARIANZA RESIDUAL

Es la dispersión de los errores cometidos entre los residuos, dispersión de la suma de las distancias de los valores observados (o experimentales) y los valores teóricos (en la recta de regresión).

* Las diferencias se toman al cuadrado para que no se puedan contrarrestar los signos positivos y negativos.

$$s_r^2 = \frac{\sum_{i,j} d_{ij}^2 n_{ij}}{N}$$



Para simplificar el mecanismo suponemos que el centro de gravedad se encuentra en el origen ($\bar{x} = 0, \bar{y} = 0$), con lo que la ecuación de la recta $y = a + bx$ se reduce a $y = bx$, partiendo que cada par (x_i, y_j) se repite una sola vez (descartando multiplicidades).

Con las hipótesis planteadas:

$$(\bar{x}=0, \bar{y}=0) \quad b = \frac{s_{xy}}{s_x^2} \quad \begin{cases} s_{xy} = a_{11} - \bar{x}\bar{y} = a_{11} \\ s_x^2 = a_{20} - \bar{x}^2 = a_{20} \\ s_y^2 = a_{02} - \bar{y}^2 = a_{02} \end{cases}$$

Con lo cual,

$$\sum_{i,j} d_{i,j}^2 = \sum_{i,j} (\hat{y}_i - y_j)^2 = \sum_{i,j} (bx_i - y_j)^2 = \sum_{i,j} (b^2 x_i^2 + y_j^2 - 2bx_i y_j) = b^2 \sum_i x_i^2 + \sum_j y_j^2 - 2b \sum_{i,j} x_i y_j$$

$$s_r^2 = \frac{\sum_{i,j} d_{i,j}^2}{N} = b^2 \frac{\sum_i x_i^2}{N} + \frac{\sum_j y_j^2}{N} - 2b \frac{\sum_{i,j} x_i y_j}{N} = b^2 a_{20} + a_{02} - 2ba_{11} = b^2 s_x^2 + s_y^2 - 2bs_{yx}$$

$$s_r^2 = b^2 s_x^2 + s_y^2 - 2bs_{yx} = b(bs_x^2 - 2s_{xy}) + s_y^2 = \frac{s_{xy}}{s_x^2} \left[\frac{s_{xy}}{s_x^2} s_x^2 - 2s_{xy} \right] + s_y^2 =$$

$$= \frac{s_{xy}}{s_x^2} [s_{xy} - 2s_{xy}] + s_y^2 = -\frac{s_{xy}^2}{s_x^2} + s_y^2 = s_y^2 \left[1 - \frac{s_{xy}^2}{s_x^2 s_y^2} \right] = s_y^2 [1 - r^2]$$

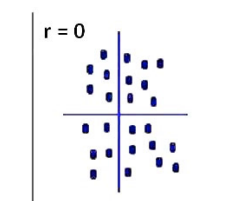
- ✓ La cota máxima de la varianza residual s_r^2 es la varianza que tratamos de explicar mediante el modelo de regresión, es decir, la varianza de la variable dependiente. En este caso, $s_r^2 = s_y^2$, hecho que sucede cuando $r=0$, esto es cuando las **variables son incorreladas**.
- ✓ La cota mínima de la varianza residual s_r^2 se obtendrá cuando las variables tienen una **dependencia funcional** $r^2=1$
- ✓ % variaciones no explicado = $100 \frac{s_r^2}{s_y^2}$
- ✓ Una forma de definir el **coeficiente de determinación**: $r^2 = 1 - \frac{s_r^2}{s_y^2} \quad 0 \leq r^2 \leq 1$

INTERPRETACIÓN COEFICIENTE DE CORRELACIÓN LINEAL

Se hace una interpretación a partir de la relación con la varianza residual $s_r^2 = s_y^2(1-r^2)$:

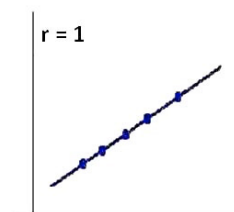
- Si $r=0 \Rightarrow s_r^2 = s_y^2$ y $b_{yx}=0$ y $b_{xy}=0$.

Las dos rectas son perpendiculares y las variables son **INCORRELADAS**



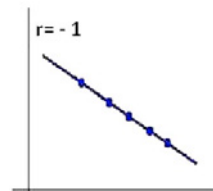
- Si $r=1 \Rightarrow s_r^2 = 0$.

Todos los puntos se encuentran situados sobre la recta de regresión, existiendo entre las dos variables una **DEPENDENCIA FUNCIONAL** (recta de regresión creciente).



- Si $r = -1 \Rightarrow s_r^2 = 0$.

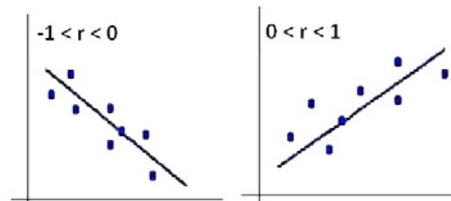
Todos los puntos se encuentran situados sobre la recta de regresión, existiendo entre las dos variables una **DEPENDENCIA FUNCIONAL** (recta de regresión decreciente).



- Si $-1 < r < 0$ ó $0 < r < 1$

Las variables están tanto más correladas en cuanto el coeficiente se aproxima más a -1 ó 1, respectivamente.

En ambos casos, existe una **DEPENDENCIA ALEATORIA** entre las variables.



Resaltar que no puede darse el caso de que una recta de regresión sea creciente y la otra decreciente

$$\begin{cases} b_{yx} > 0 \Leftrightarrow r > 0 \Leftrightarrow b_{xy} > 0 \\ b_{yx} < 0 \Leftrightarrow r < 0 \Leftrightarrow b_{xy} < 0 \end{cases}$$

DESCOMPOSICIÓN DE LA VARIABILIDAD: COEFICIENTE DE CORRELACIÓN - VARIANZA RESIDUAL

Sea \hat{y}_i el valor teórico que correspondería en la recta de regresión de Y sobre X: $\hat{y}_i = a + bx_i$. Elevando al cuadrado la descomposición $(y_i - \bar{y}) = (y_i - \hat{y}_i) + (\hat{y}_i - \bar{y})$, se obtiene:

$$SCT = \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n [(y_i - \hat{y}_i) + (\hat{y}_i - \bar{y})]^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + 2 \sum_{i=1}^n (y_i - \hat{y}_i) \cdot (\hat{y}_i - \bar{y})$$

Observemos que, $\sum_{i=1}^n (y_i - \hat{y}_i) \cdot (\hat{y}_i - \bar{y}) = \sum_{i=1}^n (y_i - a - bx_i) \cdot (a + bx_i - \bar{y}) =$

$$= a \underbrace{\sum_{i=1}^n (y_i - a - bx_i)}_{=0} + b \underbrace{\sum_{i=1}^n x_i (y_i - a - bx_i)}_{=0} + \bar{y} \underbrace{\sum_{i=1}^n (y_i - a - bx_i)}_{=0}$$

con lo cual,

| |
|--|
| $\underbrace{\sum_{i=1}^n (y_i - \bar{y})^2}_{\text{SCT}} = \underbrace{\sum_{i=1}^n (y_i - \hat{y}_i)^2}_{\text{SCR}} + \underbrace{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}_{\text{SCE}}$ |
| <small>suma cuadrados total suma cuadrados residual suma cuadrados explicada</small> |

Por otro lado, $\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \Rightarrow 1 = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} + \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$

$\underbrace{\hspace{10em}}_{\text{SCR/SCT}} \quad \underbrace{\hspace{10em}}_{r^2 = \text{SCE/SCT}}$

Una vez estimado el modelo es conveniente obtener una medida acerca de la bondad del ajuste realizado. Un estadístico que facilita esta medida es el Coeficiente de Determinación (r^2), que se

define:
$$r^2 = \frac{\text{SCE}}{\text{SCT}} = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

El Coeficiente de Determinación permite, además, seleccionar entre modelos clásicos que tengan el mismo número de regresores, ya que la capacidad explicativa de un modelo es mayor cuanto más elevado sea el valor que tome este coeficiente.

De otra parte,
$$r^2 = 1 - \frac{\text{SCR}}{\text{SCT}} = 1 - \frac{\overbrace{\sum (y_i - \hat{y}_i)^2 / N}^{s_r^2}}{\underbrace{\sum (y_i - \bar{y})^2 / N}_{s_y^2}} = 1 - \frac{s_r^2}{s_y^2} \mapsto \overbrace{s_r^2 = s_y^2 (1 - r^2)}^{\text{varianza residual}}$$

Considerando la recta de regresión de Y sobre X, el coeficiente de determinación r^2 puede expresarse:

$$r^2 = \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} = \frac{\sum \left[\frac{s_{xy}}{s_x^2} (x_i - \bar{x}) \right]^2}{\sum (y_i - \bar{y})^2} = \frac{\left[\frac{s_{xy}}{s_x^2} \right]^2 \frac{\sum (x_i - \bar{x})^2}{N}}{\frac{\sum (y_i - \bar{y})^2}{N}} = \frac{s_{xy}^2}{s_x^2 s_y^2} \mapsto \overbrace{r = \frac{s_{xy}}{s_x s_y}}^{\text{coeficiente correlación}}$$

El coeficiente de correlación lineal r es un número abstracto que determinará el grado de ajuste entre una nube de puntos y una recta de regresión. Se define como la media geométrica de los coeficientes de regresión lineal:

$$r = \sqrt{b_{yx} b_{xy}} = \sqrt{\frac{s_{yx}}{s_x^2} \frac{s_{yx}}{s_y^2}} = \frac{s_{yx}}{s_x s_y}$$

Adviértase que si la varianza residual es cero, $s_r^2 = 0$, se tiene, $s_r^2 = s_y^2 (1 - r^2) = 0 \mapsto 1 - r^2 = 0$

con lo cual, $r^2 = 1 \Rightarrow s_{xy} = \pm s_x s_y$

