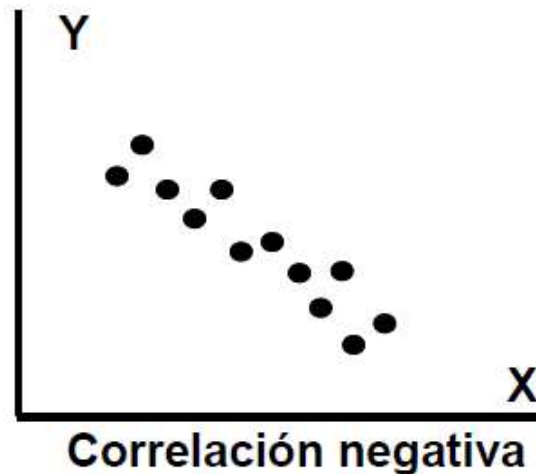
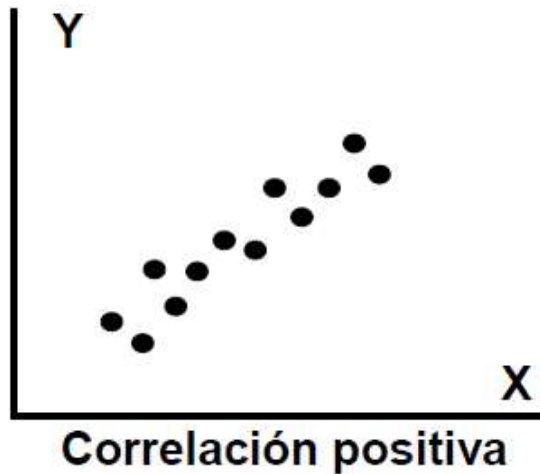


Unidad V

Análisis de Correlación



Relación entre variables aleatorias

En muchos problemas existe una relación inherente entre dos o más variables, y resulta de utilidad explorar la naturaleza de esta relación.

La **matemática** estudia variables relacionadas de manera determinística. Dos variables X e Y están relacionadas en forma **determinística** si conocido el valor de X , el valor de Y queda especificado por completo.

Ejemplo, si el alquiler de un automóvil por día tiene un costo fijo de 500\$, más 10\$ por Km. recorrido. Si X = "Km. recorridos" e Y = "costo por alquiler", entonces $Y = 500 + 10 X$. Si al cabo de un día se recorren $X = 100$ Km, el costo del alquiler queda especificado en **1500\$**.

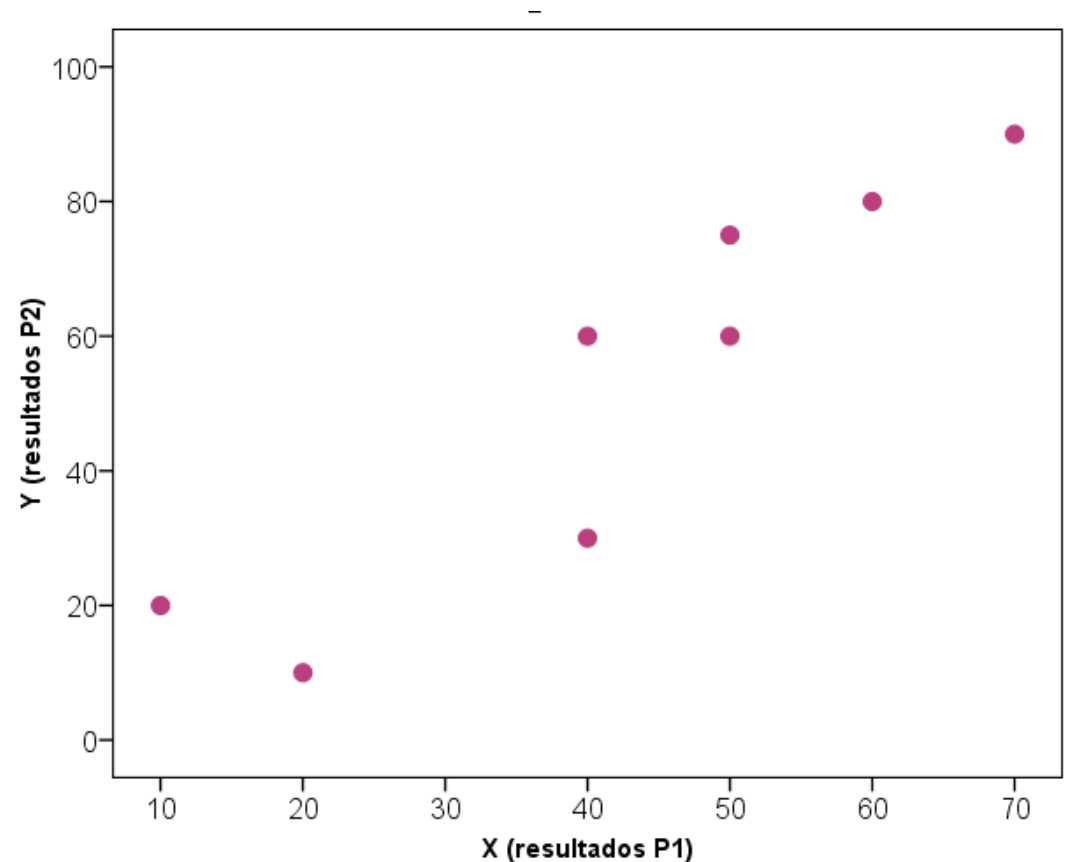
Existen muchas variables X e Y que están relacionadas entre sí, pero **no** de manera **determinística**.

Ejemplo, si X = "la altura, en m, de un estudiante de primer año de la UNS", e Y = "el peso, en Kg.", el valor de Y no queda especificado por completo si se conoce X ; dos estudiantes podrían tener el mismo valor de X , pero tener valores de peso (Y) muy diferentes.

Ejemplo: Las calificaciones obtenidas por 9 alumnos de Estadística en el primer parcial y en el segundo parcial son:

Calificaciones								
P1 (X)	10	20	40	40	50	50	60	70
P2 (Y)	20	10	30	60	60	75	80	90

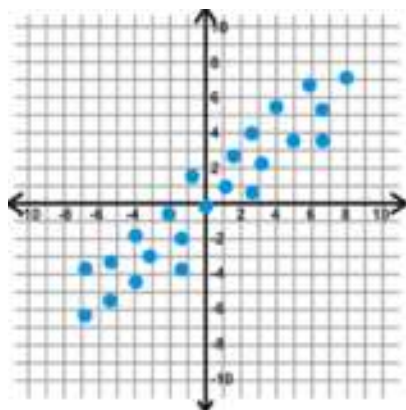
¿Existe relación entre los resultados del primer y segundo parcial?



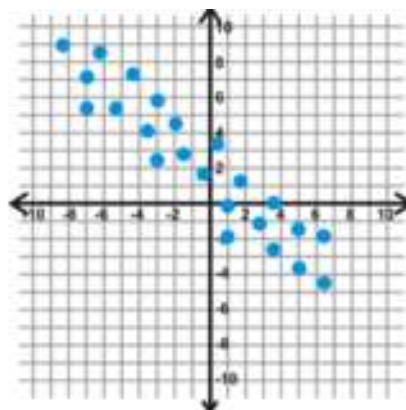
Relación lineal entre variables aleatorias

La técnica estadística que permite la investigación de la relación *no determinística lineal* entre dos o más variables es el *Análisis de Correlación*.

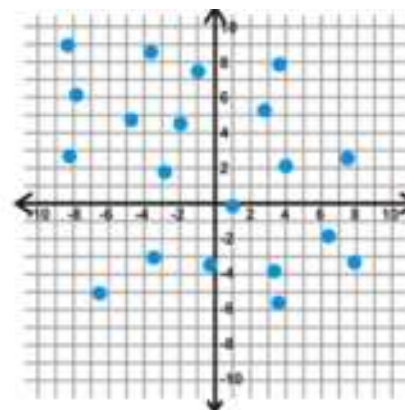
Este análisis trata de medir *el grado de asociación lineal* entre dos variables aleatorias cuantitativas **X** e **Y**. En qué medida un cambio en una de ellas afecta a la otra.



correlación positiva



correlación negativa



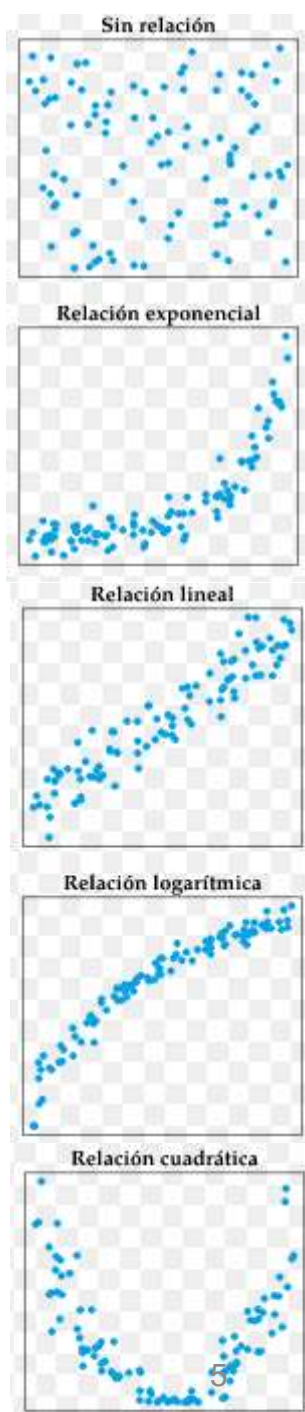
sin correlación

Herramienta gráfica: Diagrama de Dispersión

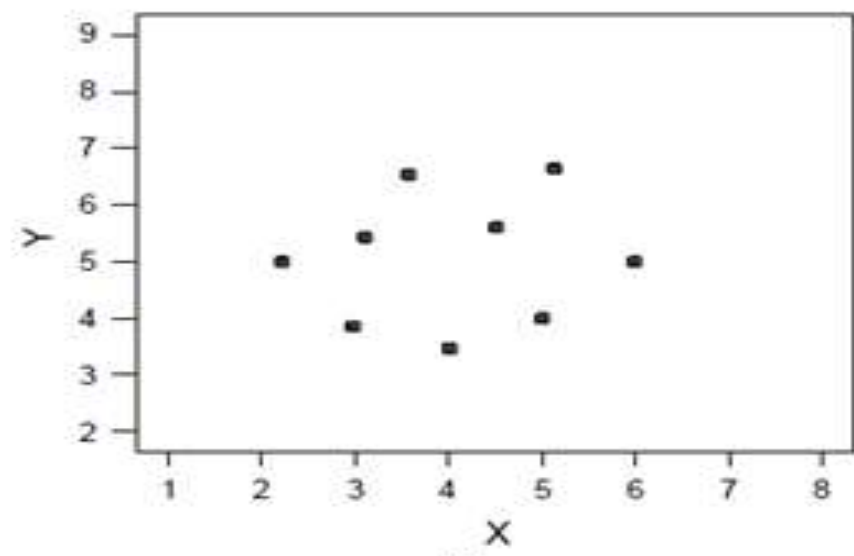
Supongamos que deseamos observar si existe **relación lineal** entre la v.a. **X** y la v.a. **Y**. Si medimos estas v.a. en **n** unidades experimentales (individuos, objetos, etc.), habremos recolectado los siguientes datos:

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n).$$

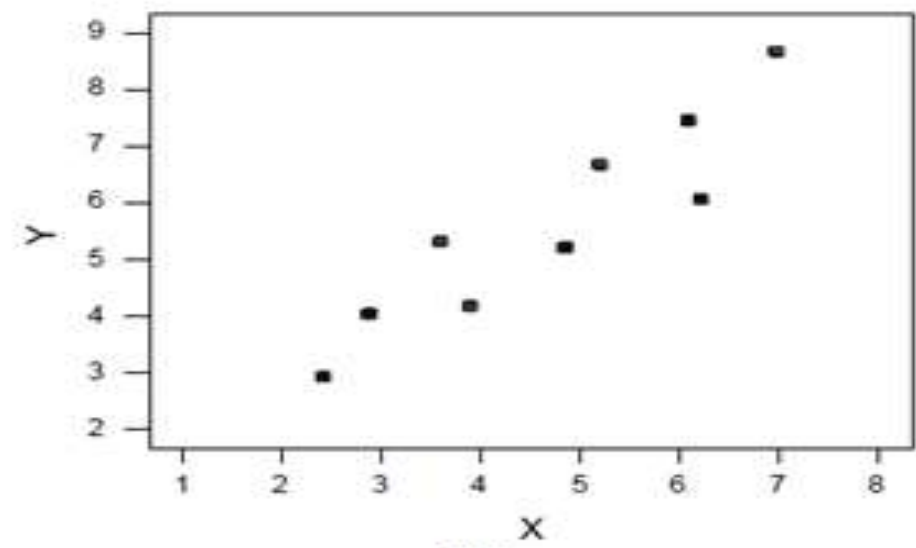
El **diagrama de dispersión** (“**scatterplot**”) permite tener una primera impresión sobre el tipo de relación existente entre dos v.a.. Este método gráfico coloca una de las v.a. (X) en el eje de las abscisas y la otra (Y), en el de las ordenadas, grafica por medio de puntos, (x_i, y_i) , los valores correspondientes de las variables bajo estudio para cada una de las unidades experimentales analizadas. Queda así determinada una **nube de puntos**. La forma de esta nube de puntos informa sobre el tipo de relación existente entre las variables. Permite establecer algún patrón de comportamiento gráfico.



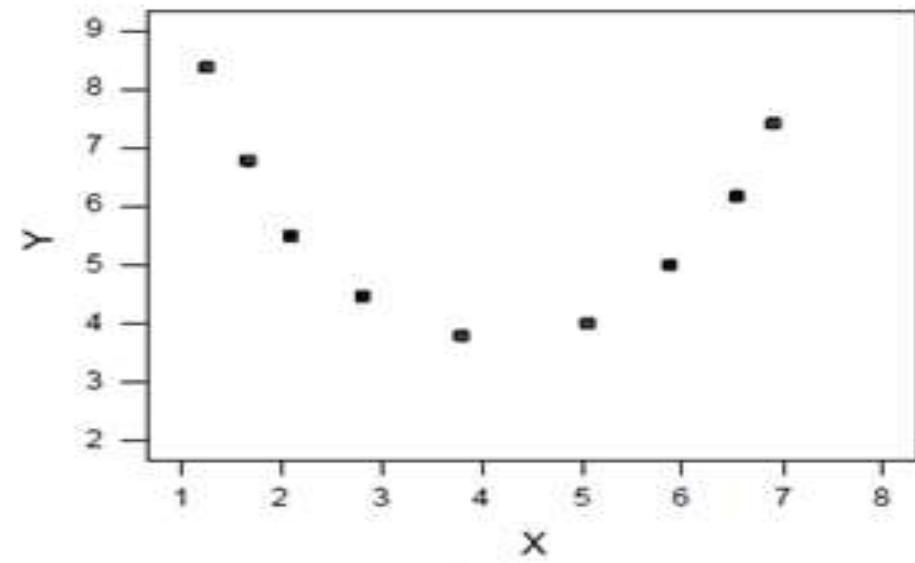
Diagramas de Dispersión que reflejan situaciones de relaciones diferentes entre las variables.



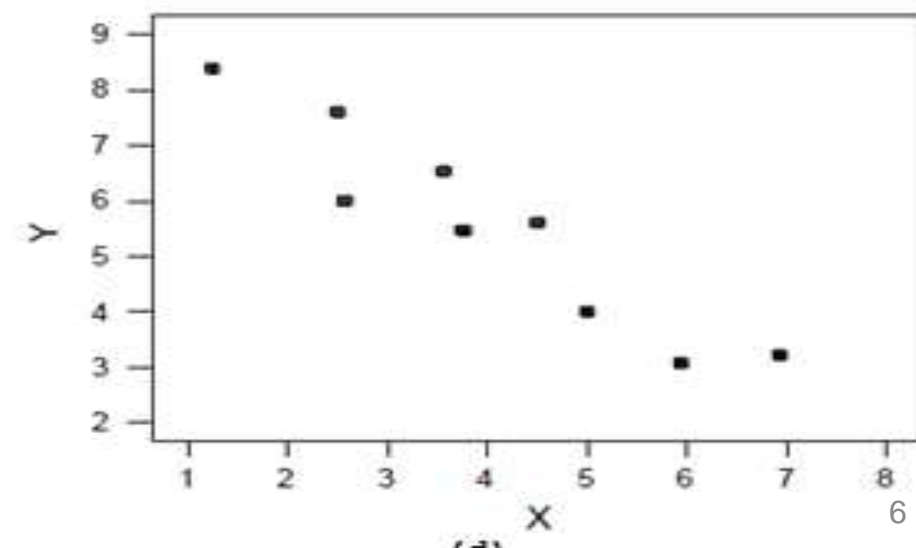
(a)



(b)



(c)



(d)

Importante!!!

El **Diagrama de Dispersión** en algunos casos puede que no resulte tan práctico para medir el **grado** de **asociación lineal**. Esto es debido a que la relación entre dos variables no siempre es **perfecta** o **nula**; habitualmente no es ni lo uno ni lo otro.

Un índice que da a conocer el **sentido de la asociación lineal** entre las v.a. X e Y, si es que ésta existe, es la **covarianza**.

$$\mathbf{Cov(X,Y)} = E[(X - E(X)).(Y - E(Y))] = E(X.Y) - E(X).E(Y)$$

**Covarianza
poblacional**

Dados los puntos $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, una muestra bidimensional obtenida al medir X e Y en **n** individuos u objetos, la **covarianza** de X e Y, denotada **Cov(X, Y)** puede ser estimada por la **Covarianza muestral**, **S_{XY}** cuya expresión es,

$$S_{XY} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{n - 1}$$

donde \bar{X} e \bar{Y} son las **medias muestrales**.

Diagrama de dispersión

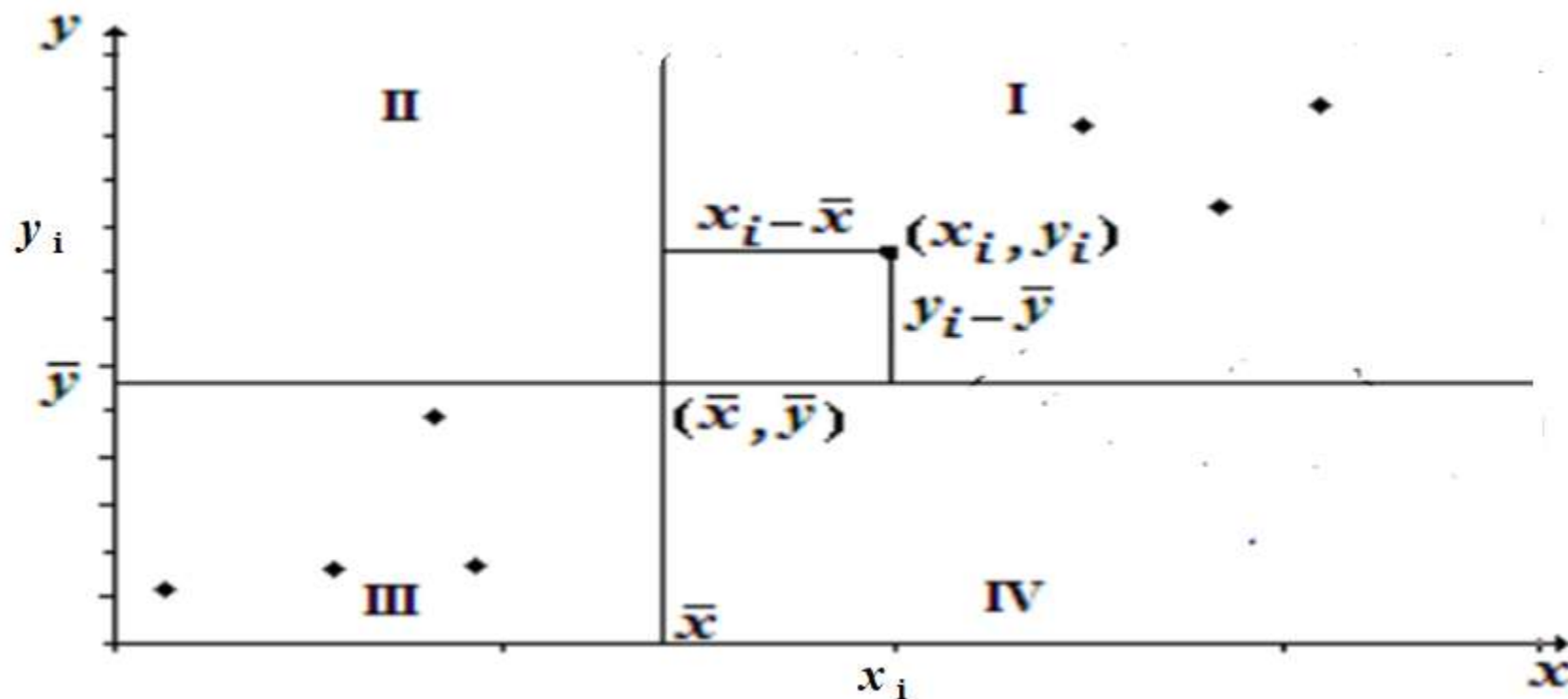


Tabla Signos de las diferencias y del producto

Diferencias y producto	I	II	III	IV
$x_i - \bar{x}$	+	+	-	-
$y_i - \bar{y}$	+	-	-	+
$(x_i - \bar{x})(y_i - \bar{y})$	+	-	+	-

Covarianza de dos variables X e Y

$$S_{xy} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{n-1} = \frac{\sum_{i=1}^n X_i Y_i - n\bar{X}\bar{Y}}{n-1}$$

La **covarianza muestral** entre dos variables nos indica si la posible **relación lineal** entre dos variables cuantitativas es **positiva (directa)** o **negativa (inversa)** ó **no existe**. El **sentido** de la relación.

- **positiva**: $S_{xy} > 0$ (cuanto mayores son las puntuaciones en la variable **X**, mayores son las puntuaciones en la variable **Y**). **Varían en igual sentido**.
- **negativa**: $S_{xy} < 0$ (cuanto mayores son las puntuaciones en la variable **X**, menores son las puntuaciones en la variable **Y**). **Varían en sentidos opuestos**.
- **Incorreladas**: $S_{xy} = 0$ (**no existe relación lineal** entre **X** e **Y**).

El **signo de la covarianza** nos dice si el aspecto de la nube de puntos es **creciente** o **no**, pero **no nos dice nada** sobre el **grado o intensidad de relación** entre las variable aleatorias. **OJO!!!**

Ejemplo: Eficiencia de programas informáticos

Un administrador de computadoras necesita saber cómo la **eficiencia de su nuevo programa de computadora** depende del **tamaño de los datos entrantes**. La eficiencia se medirá por la cantidad de solicitudes procesadas por hora. Aplicando el programa a conjuntos de datos de diferentes tamaños, obtiene los siguientes resultados

X = "Tamaño de un conjunto de datos (gigabytes)"	6	7	7	8	10	10	15
Y = "Nº de solicitudes procesadas por hora"	40	55	50	41	17	26	16



¿El Nº de solicitudes procesadas por hora depende del tamaño del conjunto de datos?

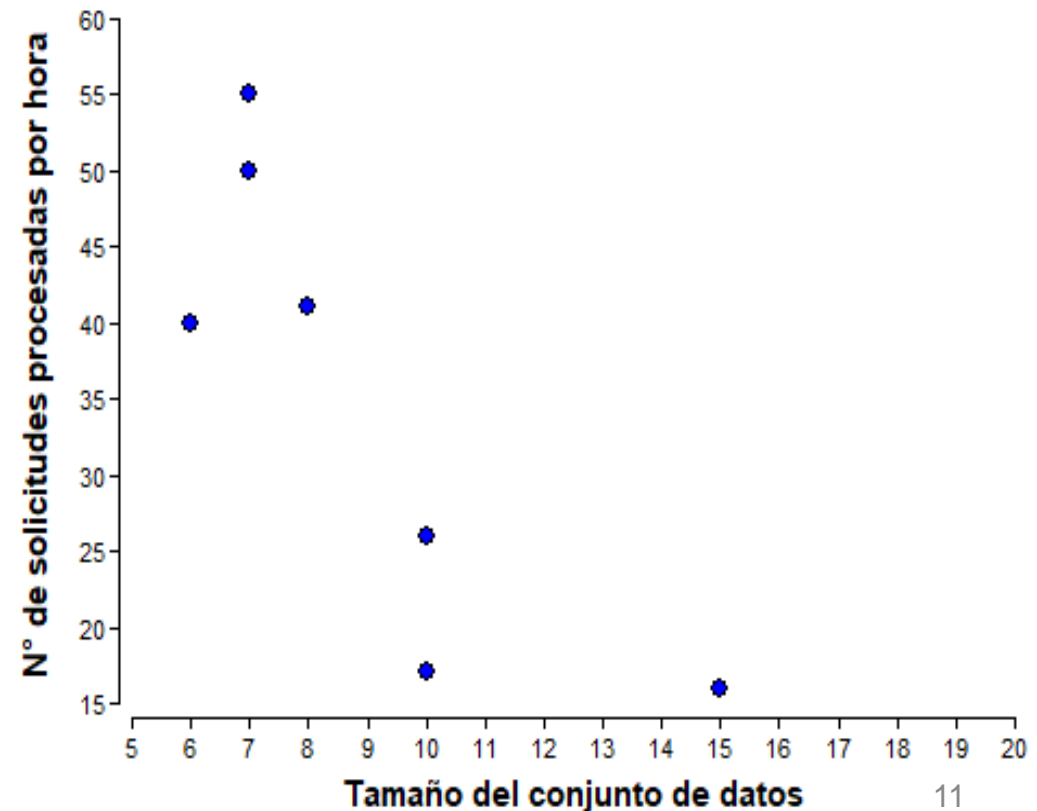
¿El N° de solicitudes procesadas por hora depende del tamaño de conjunto de datos?

$$S_{xy} = \frac{\sum_{i=1}^n X_i Y_i - n\bar{X}\bar{Y}}{n-1} = \frac{1973 - 7 \cdot 9 \cdot 35}{6} = -\frac{232}{6} = -\mathbf{38.7}$$

Relación negativa:

$S_{xy} < 0$ (cuanto mayor sea el tamaño del conjunto de datos, X , menor es la cantidad de solicitudes procesadas por hora, Y), es decir, el sentido de la asociación lineal es inverso.

¿Es una relación lineal fuerte o no?



Covarianza muestral de dos variables X e Y

$$\mathbf{S}_{\mathbf{XY}} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{n-1} = \frac{\sum_{i=1}^n X_i Y_i - n\bar{X}\bar{Y}}{n-1}$$

El **valor** de $\mathbf{S}_{\mathbf{XY}}$ depende de :

- a) las magnitudes de los valores y las unidades en que están medidas las v.a. **X** e **Y**.
- b) el número de puntos de la muestra, **n**.

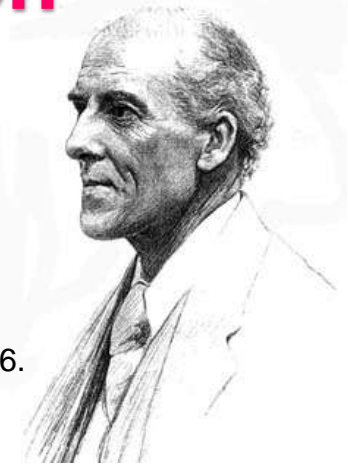
Por tal motivo, **no nos dice nada** sobre el **grado o intensidad de la relación lineal** entre las variable aleatorias.

Coefficiente de correlación de Pearson

Una medida del **grado o intensidad de relación lineal** entre dos variables cuantitativas **X** e **Y** está dada por el **coeficiente de correlación de Pearson** de las variables aleatorias **X** e **Y**, denotado: ρ_{XY} , se define:

$$\rho_{XY} = \frac{Cov(X, Y)}{\sigma_x \sigma_y}$$

Propuesto por Karl Pearson en 1896.



ρ_{XY} , puede estimarse por medio del **coeficiente de correlación muestral** **r**, definido por

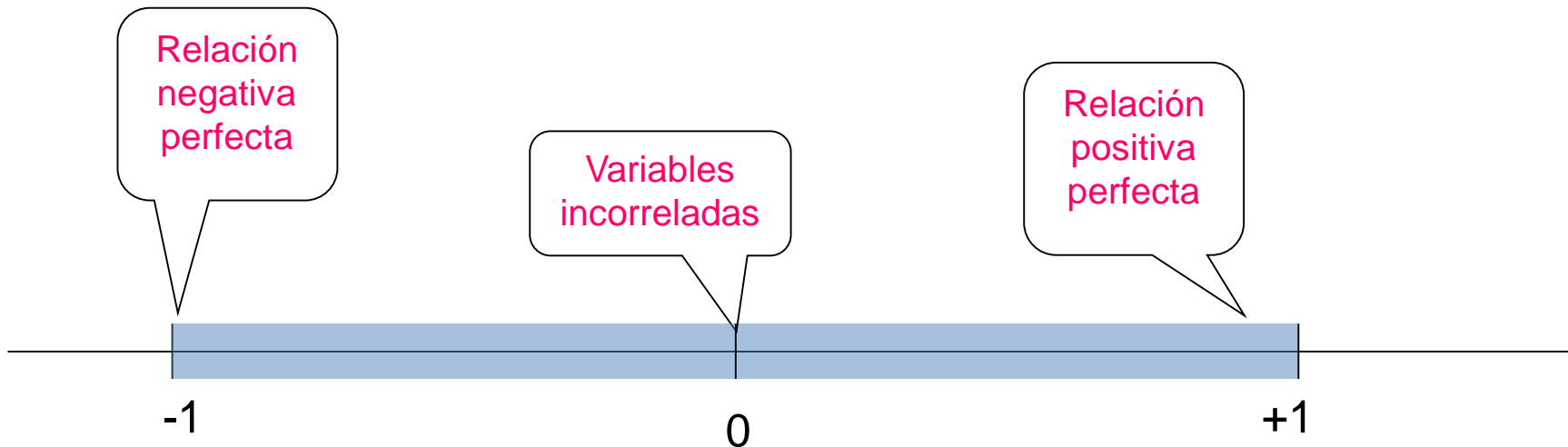
$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y - \bar{y})^2}} = \frac{S_{XY}}{\sqrt{S_X^2 S_Y^2}}$$

donde **S_{XY}** es la *covarianza muestral* y **S_X²** y **S_Y²** las *varianzas muestrales* de los valores de **X** e **Y**, respectivamente.

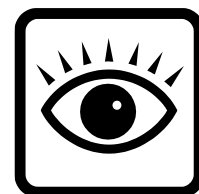
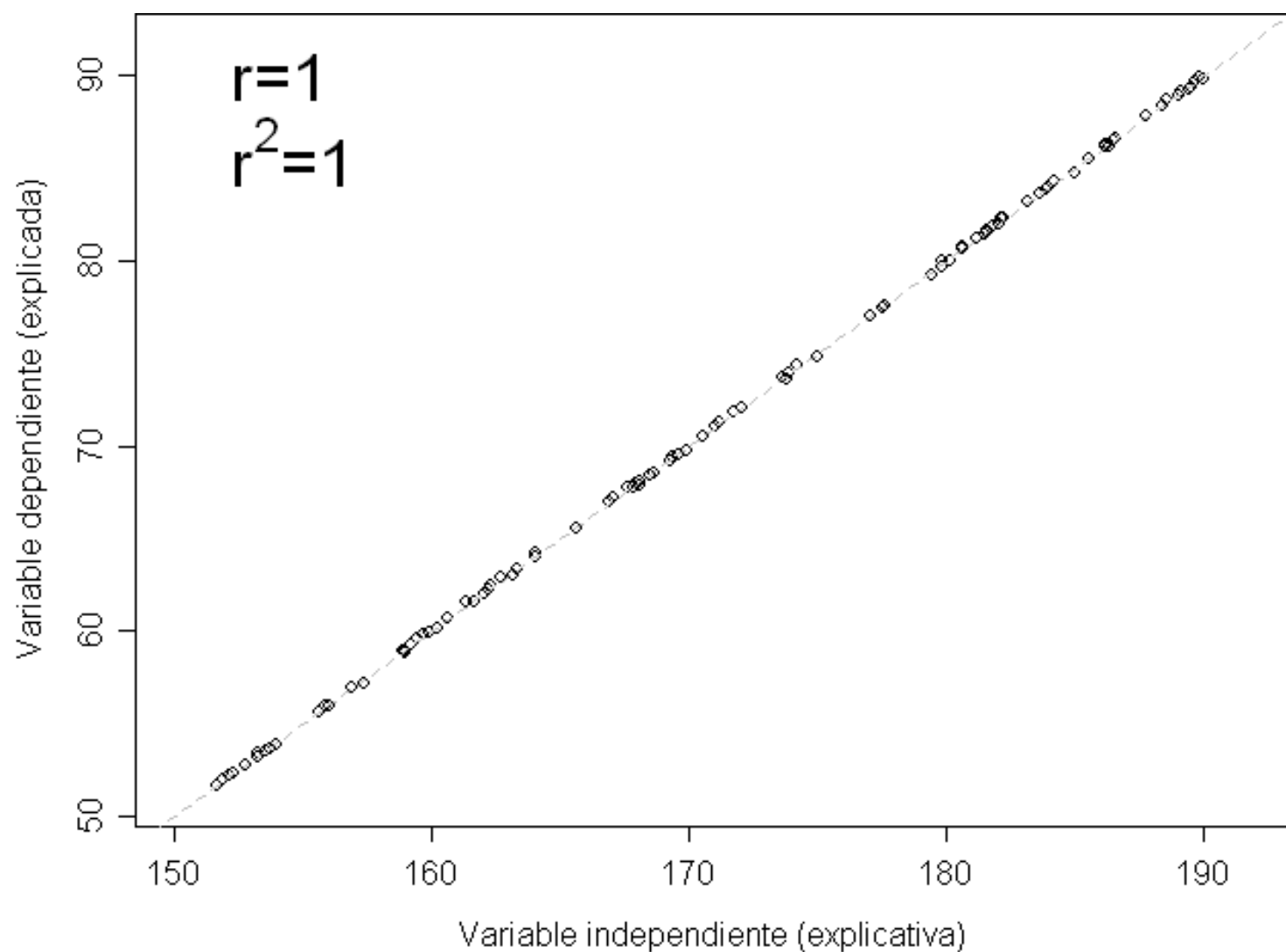
- **r** indica si los puntos tienen una tendencia a disponerse **alineadamente** (excluyendo rectas **horizontales** y **verticales**).
- Tiene el mismo signo que **S_{xy}**, por tanto, de su signo obtenemos si la posible relación es positiva o negativa.
- **r** es útil para determinar si hay **relación lineal** entre dos variables, pero no servirá para otro tipo de relaciones (cuadrática, logarítmica, etc.)

Propiedades de r

- Es **adimensional** (ya que su numerador y su denominador están medidos en las mismas escalas).
- Sólo toma valores en **$[-1,1]$**
- Las variables aleatorias cuantitativas son **incorreladas** $\Leftrightarrow r = 0$
- Relación **lineal perfecta** entre dos variables $\Leftrightarrow r = +1$ o $r = -1$
Ojo!!! Se excluyen los casos de puntos alineados **horizontalmente** o **verticalmente**.
- Cuanto más cerca esté **r** de **+1** o **-1** mayor será el **grado** o **intensidad** de relación lineal entre las variables.



Animación: Evolución de r y diagrama de dispersión



Considerando el **ejemplo** del tamaño del conjunto de datos (X), en gigabytes, y el N° de solicitudes procesadas por hora (Y),

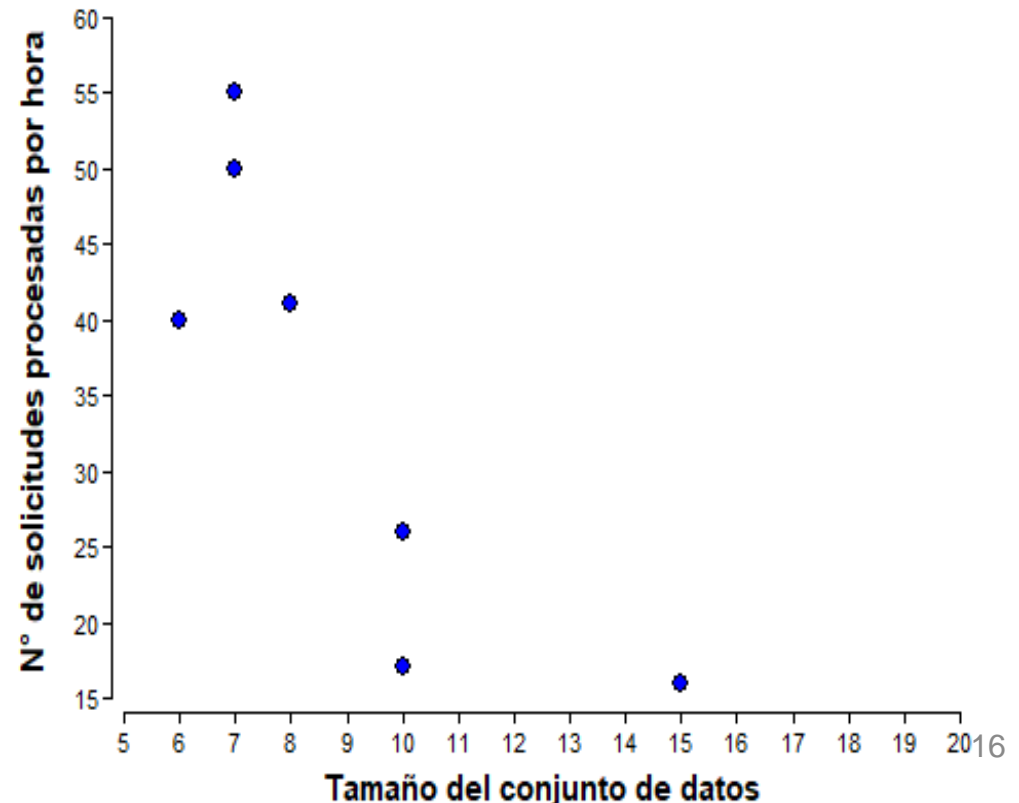
¿hay una relación lineal fuerte o no?

$$r = \frac{S_{XY}}{\sqrt{S_X^2 S_Y^2}} = \frac{-38.7}{\sqrt{9.3 * 242.1}} \cong -\mathbf{0.82}$$

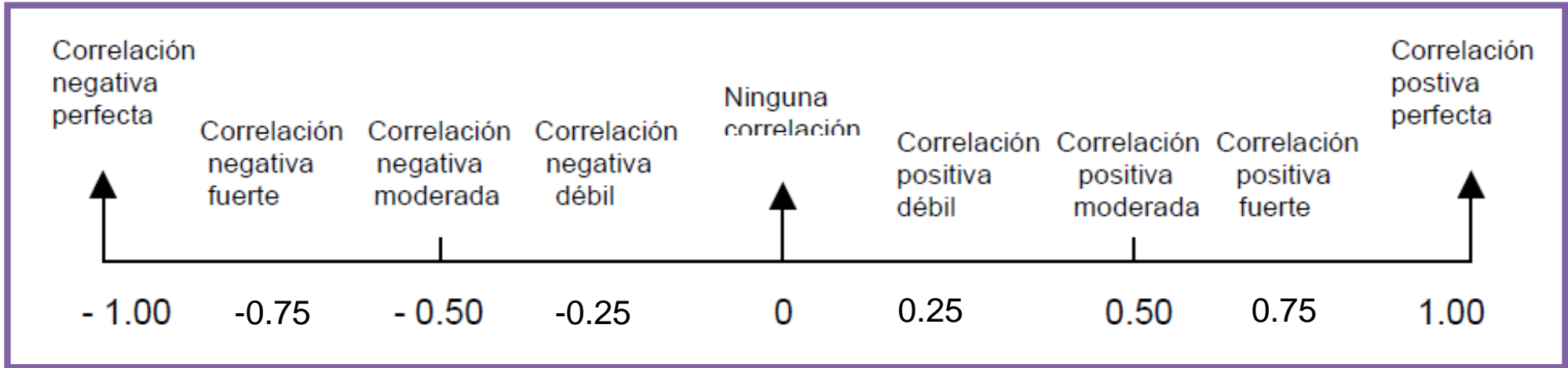


Coefficiente de correlación =
r = - 0.82

*Podemos concluir que, el tamaño del conjunto de datos (X), en gigabytes, y el N° de solicitudes procesadas por hora (Y), poseen una **relación lineal negativa fuerte**.*



Clasificación del coeficiente de correlación muestral r

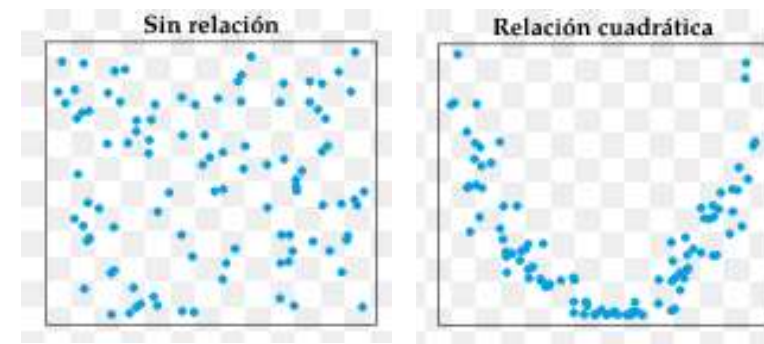


Pregunta frecuente:

¿Si $r = 0$ eso quiere decir que las variables son independientes?

Si $r = 0$, significa que **no hay relación lineal** entre las variables. Pero puede existir otro tipo de relación entre ellas.

Lo contrario, si es cierto, es decir:



Independencia



incorrelación

Cuarteto de Anscombe

El **cuarteto de Anscombe** comprende cuatro conjuntos de datos que tienen las **mismas propiedades estadísticas**, pero se **comportan distinto** al inspeccionar sus **gráficos de dispersión** respectivos.

- 1

Dato	1	2	3	4	5	6	7	8	9	10	11
$X (x_i)$	10	8	13	9	11	14	6	4	12	7	5
$Y (y_i)$	8.04	6.95	7.58	8.81	8.33	9.96	7.24	4.26	10.84	4.82	5.68

- 2

Dato	1	2	3	4	5	6	7	8	9	10	11
$X (x_i)$	10	8	13	9	11	14	6	4	12	7	5
$Y (y_i)$	9.14	8.14	8.74	8.77	9.26	8.10	6.13	3.10	9.13	7.26	4.74

- 3

Dato	1	2	3	4	5	6	7	8	9	10	11
$X (x_i)$	10	8	13	9	11	14	6	4	12	7	5
$Y (y_i)$	7.46	6.77	12.74	7.11	7.81	8.84	6.08	5.39	8.15	6.42	5.73

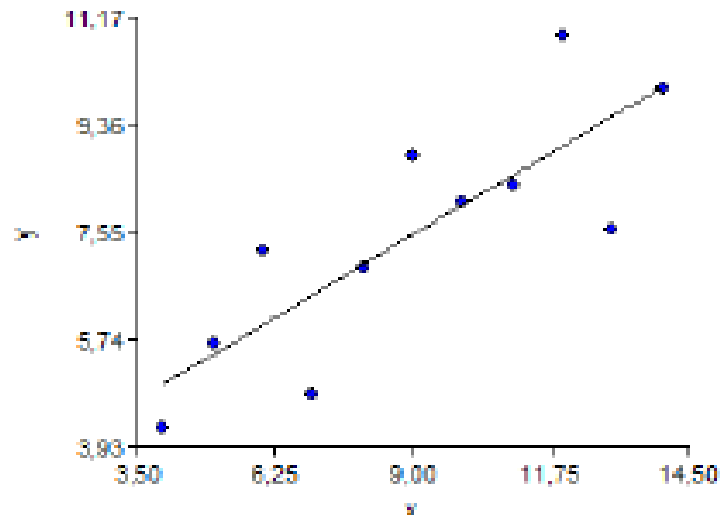
- 4

Dato	1	2	3	4	5	6	7	8	9	10	11
$X (x_i)$	8	8	8	8	8	8	8	19	8	8	8
$Y (y_i)$	6.58	5.76	7.71	8.84	8.47	7.04	5.25	12.50	5.56	7.91	6.89

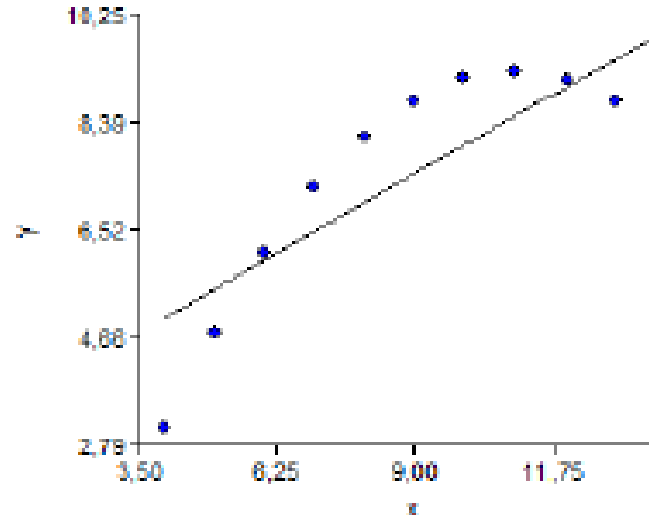
Cuarteto de Anscombe

Para los cuatro conjuntos de datos, se obtuvo:

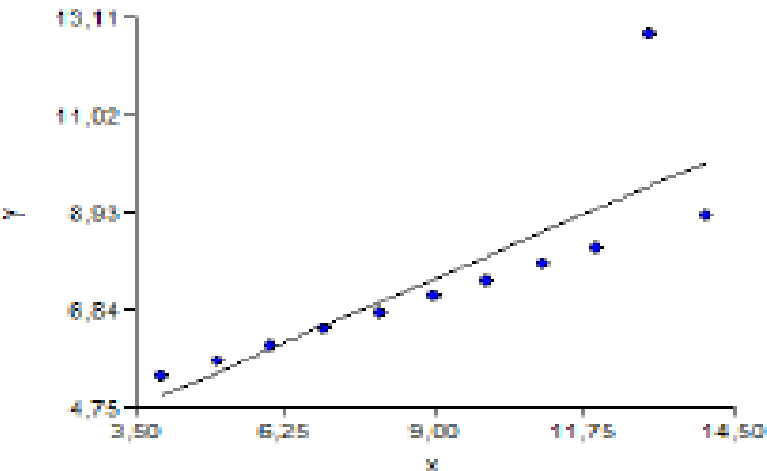
Anscombe 1



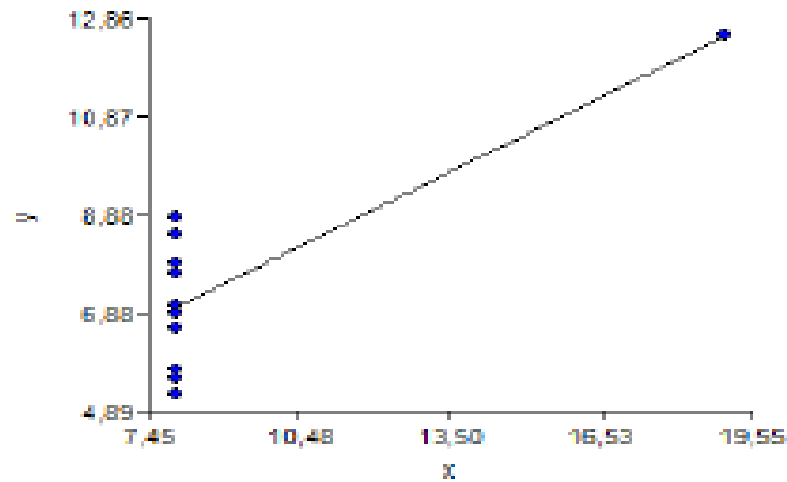
Anscombe 2



Anscombe 3



Anscombe 4




$$\bar{X} = 9 \quad \text{e} \quad \bar{Y} = 7.5$$

$$S^2_X = 11 \quad \text{y} \quad S^2_Y = 4.12$$

$$r = 0.816$$

- El primer gráfico muestra lo que parece una relación lineal simple, correspondiente a dos variables correlacionadas.
- El segundo gráfico, aunque se observa relación entre los datos, esta no es lineal y el coeficiente de correlación de Pearson, no es relevante.
- En la tercera gráfica, la distribución es lineal, pero con una línea de regresión diferente de la que se sale el dato extremo que influye lo suficiente como para alterar la recta y disminuir el coeficiente de correlación de 1 a 0.816.
- Por último, la cuarta gráfica, es un ejemplo de muestra en la que un valor atípico es suficiente para producir un coeficiente de correlación alto incluso cuando la relación entre las dos variables no es lineal.

Causalidad

X e Y están relacionadas  una es causa de la otra???

El hecho de que dos variables estén **relacionadas** **no necesariamente implica** que una sea **causa de la otra**, ya que puede ocurrir el hecho de que se esté dando una variación concomitante (simultánea), por el simple hecho de que **las dos son causa** de una **tercera**.

Por ejemplo, si se realiza un estudio en el que se analiza el número de canas (X) y la presión arterial (Y) podría encontrarse una relación lineal casi perfecta. Eso no significa que el tener canas aumente la presión arterial, lo que verdaderamente está ocurriendo es que es la edad (Z), la causante, de que se tengan más canas y una tendencia a tener más alta la presión arterial.

Correlación espuria

Se da cuando dos o más variables se **creen** estadísticamente relacionadas, porque el coeficiente de correlación presenta un valor cercano a 1 o -1, pero en realidad, **no** tienen una relación de causalidad entre ellas. Esto puede deberse a que existe un **tercer factor** no considerado, denominado “***factor de confusión***” o puede ocurrir sólo por **azar**. Este tipo de correlaciones se denominan ***correlaciones espurias***

Por ejemplo, durante el verano en cierta ciudad aumenta la venta de helados y la tasa de asesinatos. Una correlación espuria sería identificar que el aumento en la venta de helados causa un aumento en la tasa de asesinatos. En este caso, se podría considerar una tercera variable, el calor, pueda estar asociada a ambas mostrando esta falsa correlación.

Correlación espuria

Importante!!!!

Los gráficos que se muestran a continuación no pretenden dar a entender la causalidad ni están destinados a crear desconfianza en la investigación o en los datos que están bajo estudio.



Correlaciones Espurias

Un **estudiante de criminología** de la Escuela de Leyes de la Universidad de Harvard, **Tyler Vigen**, diseñó un sistema informático que busca en grandes conjuntos de datos patrones de correlación y publica en su sitio miles de **correlaciones espurias**.

Spurious correlations

tylervigen.com



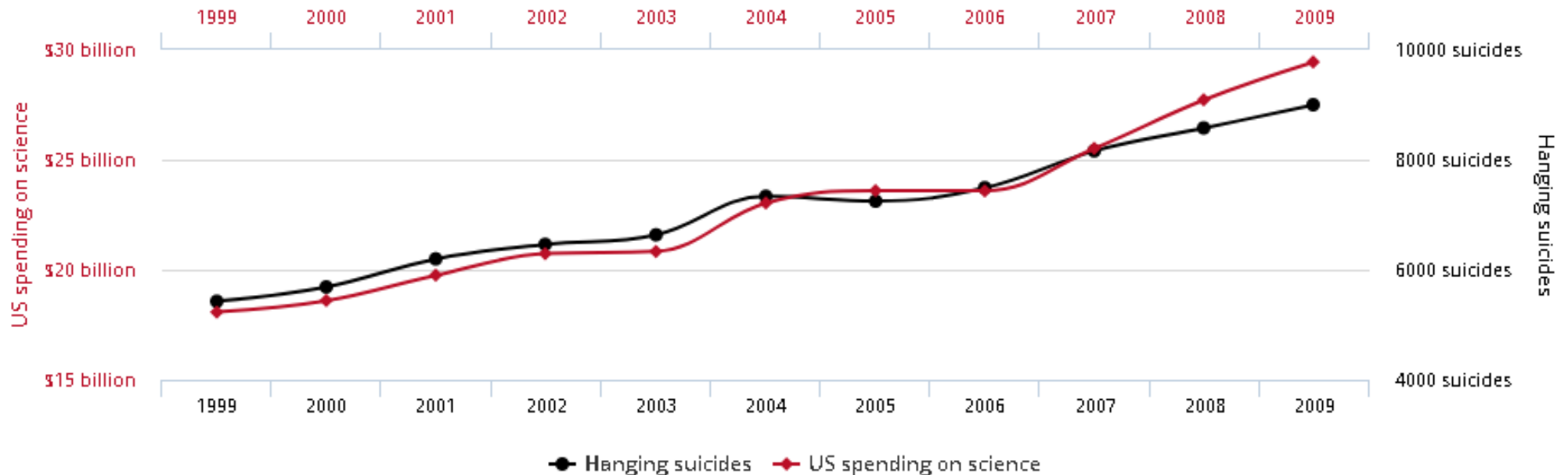
<https://www.tylervigen.com/spurious-correlations>

Correlación espuria

Gasto en ciencia, espacio, tecnología EE.UU. (rojo) vs. suicidios por asfixia, sofocación y estrangulamiento (negro). $r = 0.997$

US spending on science, space, and technology correlates with Suicides by hanging, strangulation and suffocation

Correlation: 99.79% ($r=0.99789126$)



Correlación espuria

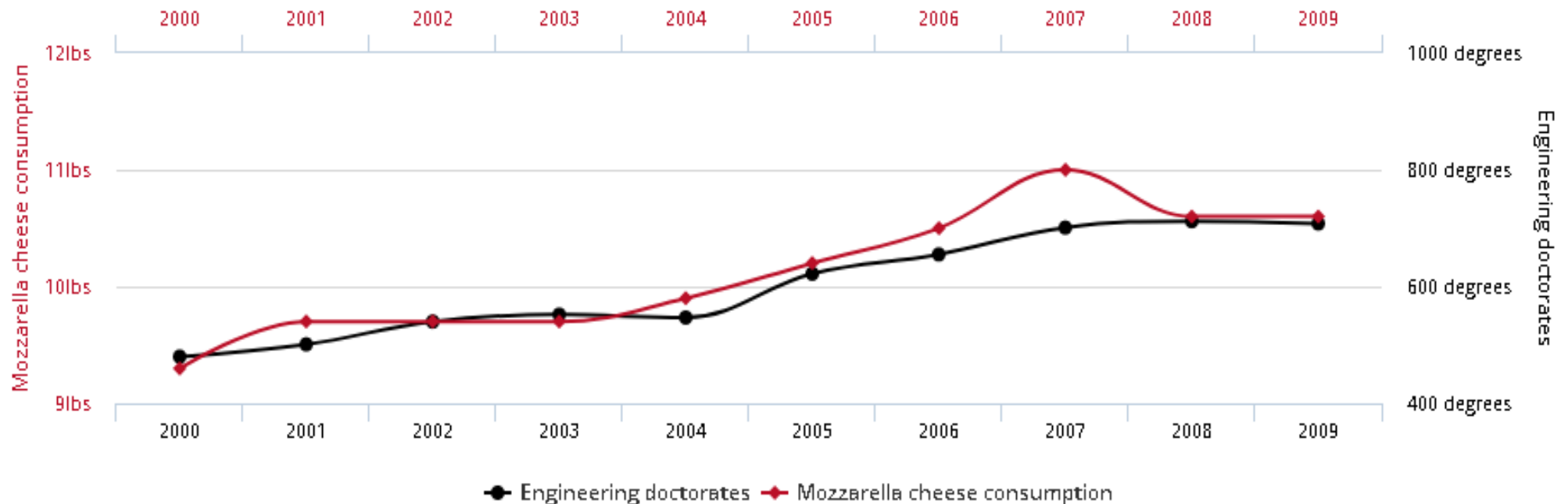
Consumo de queso mozzarella per cápita (rojo) vs. cantidad de doctores en Ingeniería civil graduados (negro), $r = 0.958$.

Per capita consumption of mozzarella cheese

correlates with

Civil engineering doctorates awarded

Correlation: 95.86% ($r=0.958648$)



IMPORTANTE!!!

- Dos variables que estén correlacionadas no significa necesariamente que exista una relación de causalidad.
- Que dos variables tengan un comportamiento similar en el tiempo, no implica que exista una relación causal entre ellas.

Uso de la Calculadora

Las calculadoras permiten calcular el coeficiente de correlación muestral r , como así también la pendiente y la ordenada al origen de la recta de ajuste a una nube de puntos, entre otras cantidades.

Por ejemplo: en una calculadora **CASIO fx - 82MS**, se debe ir:

Mode >> 3 (Reg)

>> 1 (Lin)

Luego, ingresar los datos.



CÁLCULOS ESTADÍSTICOS PARA DOS VARIABLES

- Pulsar la tecla **MODE** tantas veces hasta que en el visor aparezcan las opciones:

COMP SD REG

- Pulsar la tecla del número asociado a la opción: **REG**, es decir, la tecla: **3**. En el visor se muestran varias opciones de regresión. Pulsar la tecla del número correspondiente a la opción: **Lin** (es la opción **3**, de regresión lineal, única que tratamos en el curso).

- **Limpiar la memoria.**

Pulsar las teclas:

SHIFT+ MODE

En el visor aparece:

Scl Mode All

Pulsar la tecla correspondiente a **Scl** y en el visor se muestra el mensaje:

Stat clear

Pulsar la tecla: **=**

Después pulsar la tecla: **AC**

Ejemplo: ingresos de datos

X	Y
6	10
4	11
5	10
8	14
7	12
8	16

Pulsar

Tecclas: 6 , 1 0 M+

Tecclas : 4 , 1 1 M+

Tecclas: 5 , 1 0 M+

Tecclas: 8 , 1 4 M+

Tecclas: 7 , 1 2 M+

Tecclas: 8 , 1 6 M+

Una vez ingresados los datos, las **sumatorias** y los **índices estadísticos** se encuentran asociadas a las teclas:

- **SHIFT + 1**, para los **sumatorias**, en el visor se muestra lo siguiente:

$$\sum_{1} x^2 \quad \sum_{2} x \quad n_{3} \quad \rightarrow$$

La flecha de la derecha significa que aún hay otras sumatorias.

Pulsando el lado derecho de la tecla direccional se accede a:

$$\leftarrow \sum_{1} y^2 \quad \sum_{2} y \quad \sum_{3} xy$$

- **SHIFT + 2** para los **índices estadísticos univariados** y los relativos a la **correlación y Regresión lineal Simple**.

Cuando se pulsa **SHIFT + 2**, en el visor se muestra:

$$\begin{array}{ccc} \bar{x} & x\sigma_n & x\sigma_{n-1} \\ 1 & 2 & 3 \end{array} \rightarrow$$

Además de ésta, hay tres pantallas más:

$$\leftarrow \begin{array}{ccc} \bar{y} & y\sigma_n & y\sigma_{n-1} \\ 1 & 2 & 3 \end{array} \rightarrow$$

$$\leftarrow \begin{array}{ccc} A & B & r \\ 1 & 2 & 3 \end{array} \rightarrow$$

Los coeficientes de la recta de ajuste de la regresión lineal y el coeficiente de correlación.

$$\leftarrow \begin{array}{cc} \hat{x} & \hat{y} \\ 1 & 2 \end{array}$$

Las puntuaciones pronosticadas.

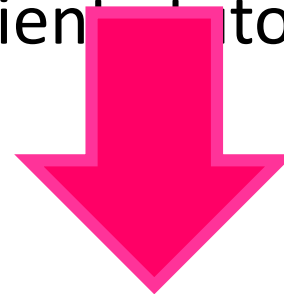
Para los datos de la tabla, algunos de los **resultados estadísticos** son:

$$\begin{aligned}\sum x^2 &= 254 \\ \sum y^2 &= 917 \\ \bar{x} &= 6.3333333 \\ \bar{y} &= 12.1666666 \\ r &= 0.799023971\end{aligned}$$

Uso de la Calculadora



Para el uso de este tipo de calculadoras en el **Análisis de Correlación** y en **Regresión Lineal Simple**, seguir paso a paso el siguiente tutorial.



https://www.google.com/search?q=como+usar+la+calculadora+en+regresion+lineal&rlz=1C1AWUC_enAR742AR742&oq=como+usar+la+calculadora+en+regresion+lineal&aqs=chrome..69i57j33.15111j0j8&sourceid=chrome&ie=UTF-8#kpvalbx=_cBDkXtSiEYuj5OUP6vqAkAw31

Prueba de Hipótesis

Significación de coeficiente de correlación de Pearson

$$\rho_{xy} = \frac{Cov(X, Y)}{\sigma_x \sigma_y}$$



Prueba de Hipótesis: Coeficiente de Correlación

Objetivo: Permite analizar la existencia de correlación entre dos v. a. y si ésta es **significativa** o **no**.

$$\rho_{XY} = \frac{Cov(X, Y)}{\sigma_x \sigma_y}$$

→ **coeficiente de correlación poblacional**

$$r = \frac{S_{XY}}{\sqrt{S_X^2 S_Y^2}}$$

→ **coeficiente de correlación muestral (estimador puntual de ρ_{XY})**

1. Planteo:

$H_0: \rho_{XY} = 0$ (No hay correlación entre las v.a.)

$H_1: \rho_{XY} \neq 0$ (Hay correlación entre las v.a.)

2. α : nivel de significación elegido.

3. Estadístico de Prueba bajo H_0 cierta:

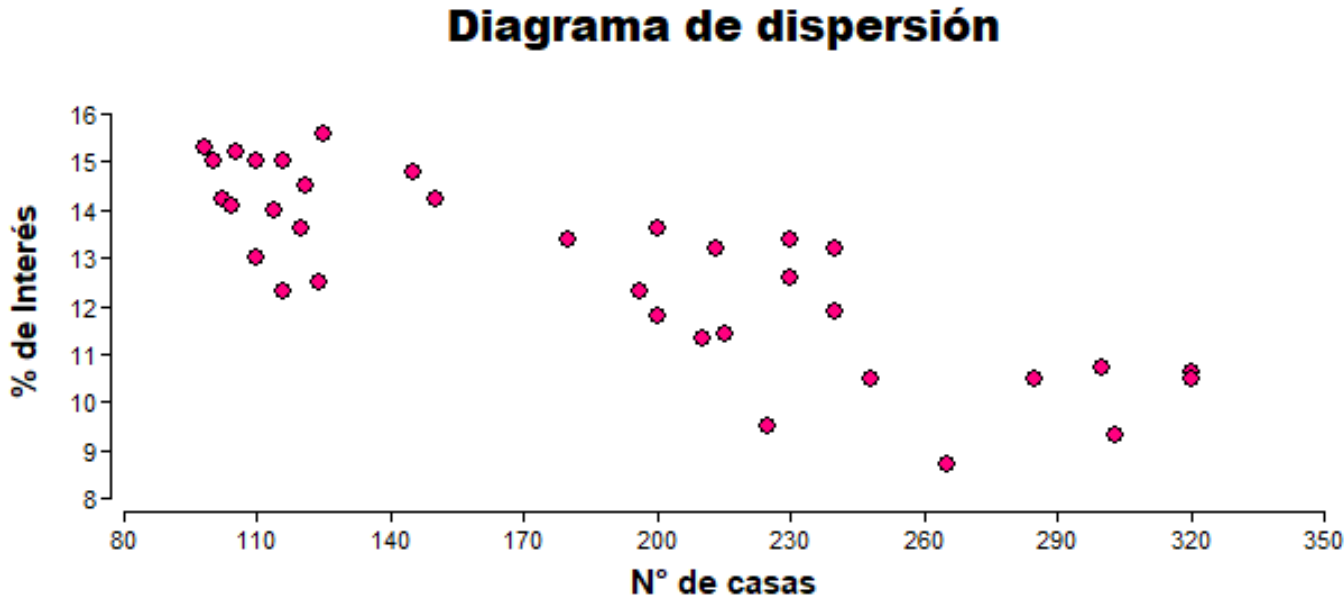
$$\frac{r - 0}{\sqrt{\frac{1 - r^2}{n - 2}}} = T \sim \text{t - student con } v = n - 2 \text{ gl}$$

Donde $\sqrt{\frac{1 - r^2}{n - 2}} = \sqrt{V(r)}$, es el **error estándar** del coeficiente de correlación muestral r .

Ejemplo:

Un banco local que se especializa en créditos para vivienda intenta analizar la relación entre *las tasas de interés mensual (%)* y el *número de casas vendidas* por mes en cierta área de la ciudad. Se compilaron los datos para un período de 35 meses elegidos al azar. La información obtenida se muestra a continuación:

Interés	Nº de casas	Interés	Nº de casas
12,3	196	10,7	300
10,5	285	13	110
15,6	125	14,1	104
9,5	225	15	110
10,5	248	13,2	240
9,3	303	13,6	200
8,7	265	13,4	230
14,2	102	13,2	213
15,2	105	12,5	124
14	114	11,9	240
15	116	10,5	320
13,6	120	13,4	180
14,5	121	15	100
12,3	116	15,3	98
11,3	210	14,8	145
11,8	200	14,2	150
11,4	215	12,6	230
		10,6	320



Ejemplo:

¿Están correladas *las tasas de interés mensual (%)*, **Y** y el *número de casas vendidas mensualmente*, **X**?

Interés	Nº de casas		Interés	Nº de casas
12,3	196		10,7	300
10,5	285		13	110
15,6	125		14,1	104
9,5	225		15	110
10,5	248		13,2	240
9,3	303		13,6	200
8,7	265		13,4	230
14,2	102		13,2	213
15,2	105		12,5	124
14	114		11,9	240
15	116		10,5	320
13,6	120		13,4	180
14,5	121		15	100
12,3	116		15,3	98
11,3	210		14,8	145
11,8	200		14,2	150
11,4	215		12,6	230
			10,6	320

$$\begin{aligned}
 r &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y - \bar{y})^2}} = \frac{S_{XY}}{\sqrt{S_X^2 S_Y^2}} \\
 &= \frac{78877.3 - 35 \cdot 12.76185 \cdot 14}{\sqrt{(5823.51 - 35 \cdot 12.763^2) (1377222 - 35 \cdot 185.1429^2)}} \\
 &= \frac{-3826.959145}{4657.518467} = \mathbf{-0.82167342}
 \end{aligned}$$

Según el valor de **r**, el número de casas vendidas mensualmente (**X**), y las tasas de interés mensual (**Y**), poseen **una relación lineal negativa o inversa, fuerte**.

Prueba de Hipótesis: Coeficiente de Correlación

A un nivel de significación de 0.01, ¿la correlación entre las tasas de interés mensual (%) y el número de casas vendidas mensualmente es significativa o no?

1. Planteo:

$H_0: \rho_{XY} = 0$ (No hay correlación entre las v.a.)

$H_1: \rho_{XY} \neq 0$ (Hay correlación entre las v.a.)

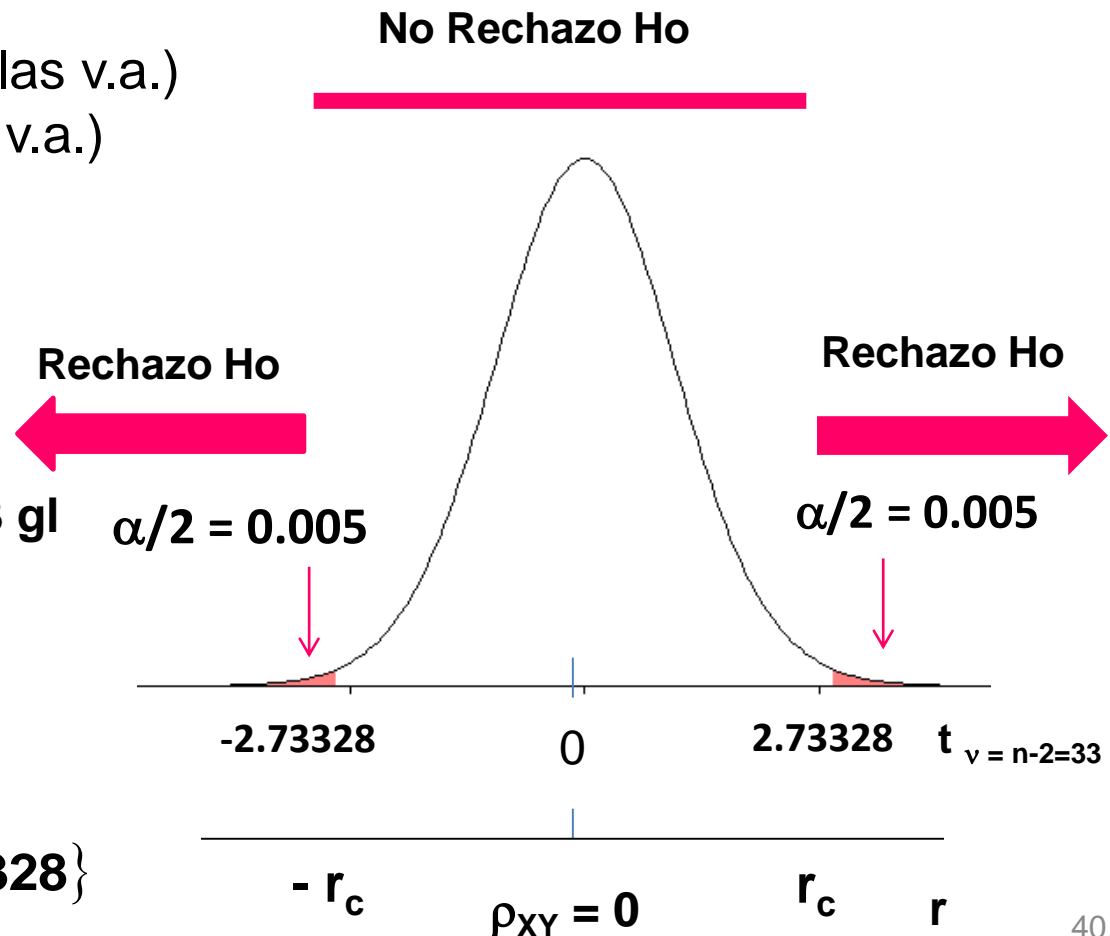
2. $\alpha = 0.01$

3. Estadístico de Prueba bajo H_0 cierta:

$$\frac{r - 0}{\sqrt{\frac{1 - r^2}{33}}} = T_{\sim} \text{ t - student con } v = n - 2 = 33 \text{ gl}$$

4. Región Crítica:

$$RC = \{ t_{33} / t_{33} < -2.73329 \text{ o } t_{33} > 2.73328 \}$$



5. El valor del estadístico es:

De los $n = 35$ datos

como $r = -0.82167342$



$$t_{\text{obs}} = \frac{-0.82167342 - 0}{\sqrt{\frac{1 - 0.82167342^2}{35 - 2}}} = -8.299$$

6. Comparación:

Como $t_{\text{obs}} = -8.299 < -2.73328$ (cae en la Región Crítica) entonces **Rechazo H_0** .

7. Conclusión:

Con una probabilidad de error de **0.01**, tengo evidencias suficientes para afirmar que existe una **relación lineal** entre la tasa de interés mensual y el N° de casas vendidas por mes en esa área de la ciudad. Las v.a. están correlacionadas. **La relación es altamente significativa (significa que es improbable que sea por azar).**

Ejemplo : Con Infostat

The screenshot displays the Infostat software interface. On the left, a data table titled 'Interes_casas_clase' is visible, containing three columns: 'Caso', 'Interés', and 'Nº de casas'. The table lists 18 rows of data. On the right, a menu titled 'Medidas resumen' is open, listing various statistical measures. The 'Análisis de correlación' option is highlighted, and a sub-menu is displayed to its right, listing correlation coefficients, partial correlations, path analysis, and Mantel's correlation.

Caso	Interés	Nº de casas
16	11,8	200
17	11,4	215
18	10,6	320
19	10,7	300
20	13,0	110
21	14,1	104
22	15,0	110
23	13,2	240
24	13,6	200
25	13,4	230
26	13,2	213
27	12,5	124
28	11,9	240
29	10,5	320
30	13,4	180
31	15,0	100
32	15,3	98
33	14,8	145

Medidas resumen

- Tablas de frecuencias
- Probabilidades y cuantiles
- Estimación de características poblacionales
- Cálculo del tamaño muestral
- Inferencia basada en una muestra
- Inferencia basada en dos muestras
- Análisis de la varianza
- Análisis de la varianza no paramétrica
- Modelos lineales generales y mixtos
- Modelos lineales generalizados mixtos (MLGM)
- Regresión lineal
- Regresión no lineal
- Modelos no lineales mixtos
- Ridge regression
- Análisis de correlación**
- Structural Model Equations
- Datos categorizados
- Análisis multivariado
- Serie de tiempo
- Suavizados y ajustes

Coeficientes de correlación Ctrl+R

- Correlaciones parciales
- Análisis de sendero (path analysis)
- Correlation entre matrices (Mantel)

Ejemplo: Con Infostat

Medidas resumen

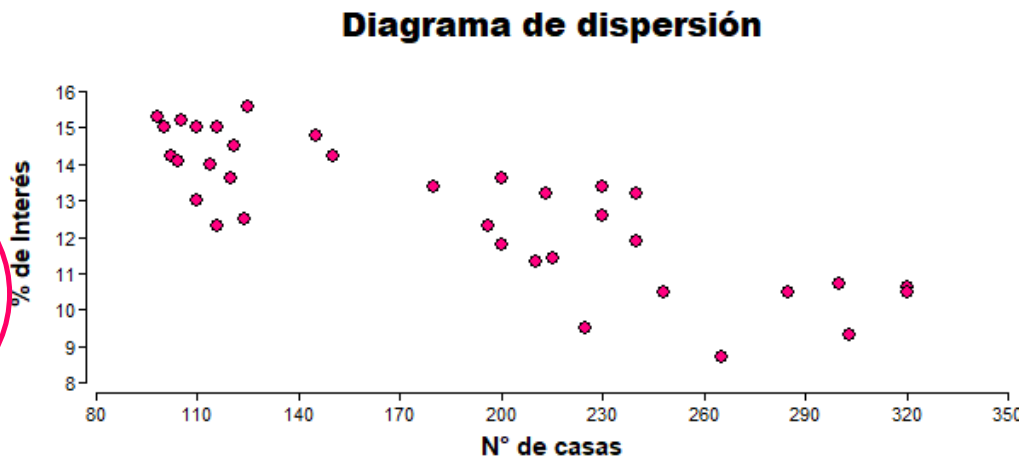
Variable	n	Media	D.E.	Var(n-1)	Mín	Máx
Interés	35	12,76	1,90	3,60	8,70	15,60
Nº de casas	35	185,14	72,25	5220,48	98,00	320,00

Salida tipo lista

Coefficientes de correlación

Correlación de Pearson

Variable(1)	Variable(2)	n	Pearson	p-valor
Interés	Nº de casas	35	-0,82	<0,0001




Conclusión: Con una probabilidad de error menor a **0.0001**, tengo evidencias suficientes para afirmar que existe una **relación lineal** entre la tasa de interés y el Nº de casas vendidas en esa área de la ciudad. Las v.a. están correlacionadas. La **relación es estadísticamente significativa** (significa que es improbable que sea por azar).

Otra forma de mostrar la salida del coeficiente de correlación muestral (Salida Matricial)

Coeficientes de correlación

Correlación de Pearson: Coeficientes\probabilidades

	Interés	Nº de casas
Interés	1,00	1,5E-09
Nº de casas	-0,82	1,00



Valor P de la prueba de
significación del análisis de
correlación.

Ejercicio: Con Infostat

Las bodegas modernas utilizan vehículos guiados computarizados y automatizados para el manejo de materiales. En consecuencia, la disposición física de la bodega debe diseñarse con cuidado a modo de evitar el congestionamiento de los vehículos y optimar el tiempo de respuesta. En *The journal of Engineering for Industry* (agosto 1993) se estudió el diseño óptimo de una bodega automatizada. La disposición empleada supone que los vehículos no se bloquean entre sí cuando viajan dentro de la bodega, es decir, no hay congestionamiento. La validez de este supuesto se verificó simulando por computadora las operaciones de la bodega. En cada simulación se varió el número de vehículos y se registró el tiempo de congestionamiento (tiempo total que un vehículo bloquea a otro). Los datos se muestran en la tabla de abajo. Los investigadores están interesados en conocer la relación entre *el tiempo de congestionamiento (Y)* y *el número de vehículos (X)*.

<i>X</i>	1	2	3	4	5	6	7	8	9	10
<i>Y</i>	0	0	0.02	0.01	0.01	0.01	0.03	0.03	0.02	0.04

- Quantificar la asociación lineal existente entre ambas variables.
- ¿Es significativa la dependencia lineal entre las variables?. Usar el **valor P** para concluir.



Curiosidad !!!

Conjuntos de datos que son idénticos en varias propiedades estadísticas, pero producen gráficos diferentes, se utilizan con frecuencia para ilustrar la importancia de las representaciones gráficas cuando se exploraran datos. Se crean mediante el algoritmo **Simulated Annealing**.

Bases de datos:

- "DataSaurus" (first created by Alberto Cairo).
- "DataSaurus Dozen" (disponible en [R package on GitHub](#)).

Página web:

<https://blog.revolutionanalytics.com/2017/05/the-datasaurus-dozen.html>

Artículo:

[Same Stats, Different Graphs: Generating Datasets with Varied Appearance and Identical Statistics through Simulated Annealing](#)