

Análisis de regresión

19 de agosto
(semana 2)

Plan de trabajo

1. Estudio de pares de variables cuantitativas
2. Uso de R para el modelo ANOVA a una vía
3. Introducción a la regresión lineal simple

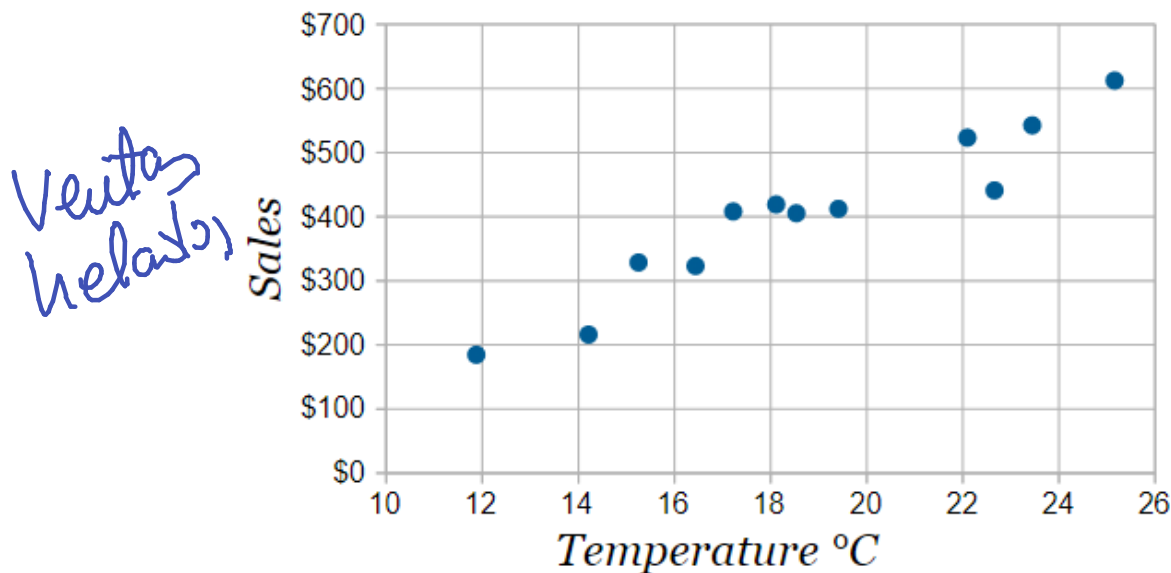
Nota:

- Ya está disponible el taller 1.
- Este fin de semana restrinjo el Drive a los inscritos y creo el Classroom y el MOODLE.
- La primera entrega del trabajo final será para el domingo 4 de septiembre (ya está en Drive).
- En la semana del 29 de agosto al 2 de septiembre no habrá clase. Pendientes de una actividad a desarrollar.

Relación entre dos variables cuantitativas

Diagrama de dispersión y coeficiente de correlación

- ❖ Vimos en estadística descriptiva que el diagrama de dispersión nos permitía visualizar si dos variables cuantitativas tenían o no una relación entre ellas.



Fuente:
<https://www.mathtsisfun.com/data/scatter-xy-plots.html>

- ❖ Pero, ¿cómo podemos cuantificar si una relación entre dos variables es fuerte o débil? Para ello, utilizamos el **coeficiente de correlación lineal** de Pearson.
- ❖ Recuerden que relación **NO** necesariamente implica causalidad entre las dos variables.

Diagrama de dispersión y coeficiente de correlación

Coeficiente de correlación lineal de Pearson

Es un parámetro que mide qué tan fuerte es la relación lineal entre dos variables. Toma valores entre -1 y 1, donde -1 indica una relación lineal inversa y 1 indica una perfecta relación lineal directa.



Fuente:
<https://www.mathsisfun.com/data/correlation.html>

- ❖ ¡Ojo! El coeficiente de correlación solo mide relaciones lineales. Si las dos variables están relacionadas de manera no lineal, es posible que este coeficiente no detecte esa relación.
- ❖ <http://guessthecorrelation.com/>

Otras medidas de asociación

Medidas de asociación en la población

Sea $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$ una muestra aleatoria de una población bivariada continua (X, Y) con parámetros de asociación:

$$\rho_{X,Y} = \frac{\text{cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}}, \text{ (Pearson)}$$

$$\rho \in [-1, 1]$$

$$\tau_{X,Y} = p[(X_1 - X_2)(Y_1 - Y_2) > 0] - p[(X_1 - X_2)(Y_1 - Y_2) < 0],$$

(Kendall)

$$\tau \in [-1, 1]$$

$$R_{X,Y} = \frac{\text{cov}(R(X), R(Y))}{\sqrt{\text{Var}(R(X))\text{Var}(R(Y))}}, \text{ (Spearman)}$$

$$R \in [-1, 1]$$

- ❖ El primero solo mide relaciones lineales, los otros relaciones monótonas

Repaso: operador correlación

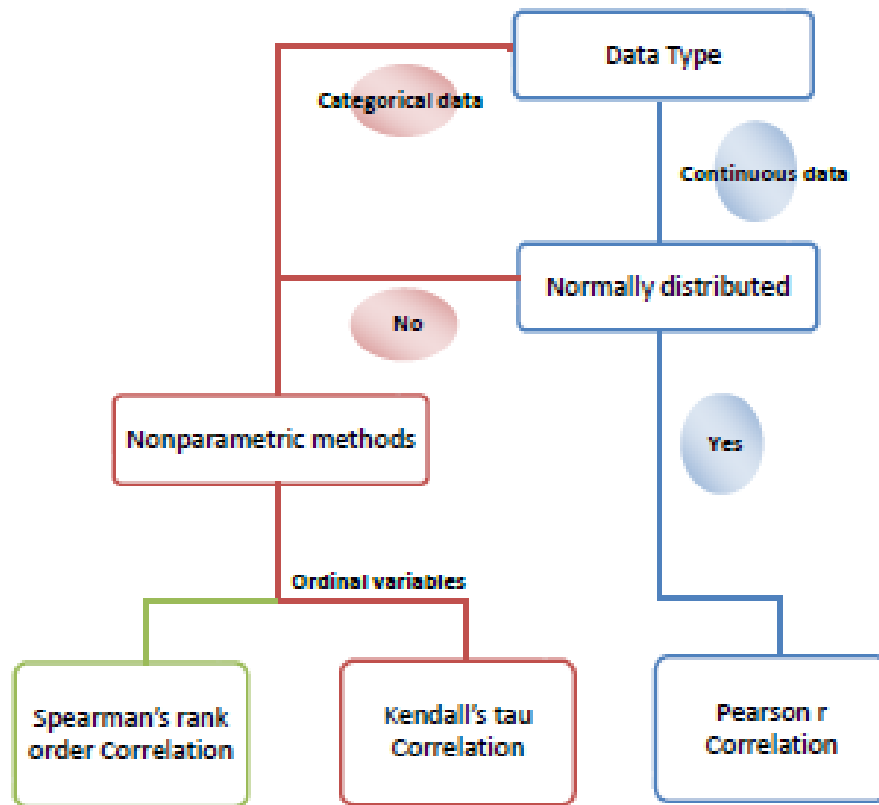
Teorema: propiedades de la correlación

- (1) $|\text{cor}(X, Y)| \leq 1$;
- (2) $\text{cor}(X, Y) = \text{cor}(Y, X)$;
- (3) $\text{cor}(X, X) = 1$;
- (3) $\text{cor}(X, -X) = -1$;
- (5) $\text{cor}(aX + b, Y) = \text{cor}(X, Y) \quad \forall a, b \in \mathbb{R}, a > 0$;
- (6) $|\text{cor}(X, Y)| = 1$ si y sólo si $\exists a, b \in \mathbb{R}$ (no simultáneamente 0) tales que $p(aX + bY = 0) = 1$.
- (7) (*Independencia* $\Rightarrow \rho = 0$) Si X, Y son independientes, $\text{cor}(X, Y) = 0$.

❖ La mayoría de las propiedades se tienen para el coeficiente de Kendall y de Spearman.

Otras medidas de asociación

Making a Decision of the Correlation Methods



Fuente:
<https://www.originlab.com/doc/Origin-Help/Correlation-Coefficient>

- *Independencia* $\Rightarrow \rho = 0, \tau = 0, R = 0$.
- $\rho = 0$ y normalidad bivariada \Rightarrow *Independencia*.
- $|\rho| = 1 \Rightarrow |\tau| = 1, |R| = 1$.

En distribs. Normales: ^{multiv.}
① hay relación $(\rho \neq 0)$
lineal
② Hay independ. $(\rho = 0)$

Coeficiente de correlación (Bonnett & Wright, 2000)

❖ **Estimación puntual:**

$$\hat{\rho} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \cdot \sum_{i=1}^n (y_i - \bar{y})^2}}$$

❖ **Estimación asintótica por intervalo** (Asumiendo que los datos vienen de una distribución normal):

❖ Paso 1. Calcule: $z_r = \frac{1}{2} \ln \left(\frac{1 + \hat{\rho}}{1 - \hat{\rho}} \right)$

❖ Paso 2. Identifique el límite inferior (l) y superior (u) mediante la fórmula:

$$z_r \mp z_{1-\frac{\alpha}{2}} \sqrt{\frac{1}{n-3}}$$

❖ Paso 3. Transforme l y u de vuelta a la escala original del coeficiente

$$\rho_L = \frac{\exp(2l) - 1}{\exp(2l) + 1}, \quad \rho_U = \frac{\exp(2u) - 1}{\exp(2u) + 1}$$

❖ **Pruebas de hipótesis y código R:** `cor.test` (IC e hipótesis \neq cero)

❖ Usar Bootstrap si no hay normalidad o para otros sistemas

Coeficiente Tau de Kendall (Hogg et al, 2005)

❖ **Estimación puntual:**
$$\hat{\tau} = \frac{2}{n(n-1)} \sum_{i < j} \text{signo}(x_i - x_j) \text{signo}(y_i - y_j)$$

❖ **Prueba de hipótesis de no asociación:**

`cor.test(..., method="kendall")`

Theorem 10.8.2. *Let $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$ be a random sample on the bivariate random vector (X, Y) with continuous cdf $F(x, y)$. Under the null hypothesis of independence between X and Y , i.e., $F(x, y) = F_X(x)F_Y(y)$, for all (x, y) in the support of (X, Y) , the test statistic K satisfies the following properties:*

K is distribution free with a symmetric pmf (10.8.4)

$$E_{H_0}[K] = 0 \quad (10.8.5)$$

$$\text{Var}_{H_0}(K) = \frac{2}{9} \frac{2n+5}{n(n-1)} \quad (10.8.6)$$

$\frac{K}{\sqrt{\text{Var}_{H_0}(K)}}$ has an asymptotic $N(0, 1)$ distribution. (10.8.7)

Coeficiente Rho de Spearman (Bonnett & Wright, 2000)

❖ **Estimación puntual:** $\hat{r}_S = 1 - \frac{6 \sum_{i=1}^n d_i^2}{n(n^2 - 1)}$, con $d_i = R(x_i) - R(y_i)$

❖ **Prueba de hipótesis de no asociación:**

`cor.test(..., method="spearman")`

Theorem 10.8.4. *Let $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$ be a random sample on the bivariate random vector (X, Y) with continuous cdf $F(x, y)$. Under the null hypothesis of independence between X and Y , i.e., $F(x, y) = F_X(x)F_Y(y)$, for all (x, y) in the support of (X, Y) , the test statistic r_S satisfies the following properties:*

r_S is distribution-free, symmetrically distributed about 0 (10.8.11)

$$E_{H_0}[r_S] = 0 \quad (10.8.12)$$

$$Var_{H_0}(r_S) = \frac{1}{n-1} \quad (10.8.13)$$

$$\frac{r_S}{\sqrt{Var_{H_0}(r_S)}} \text{ is asymptotically } N(0, 1). \quad (10.8.14)$$

Coeficientes Tau y Rho (Hogg et al, 2005)

❖ Intervalos de confianza

- ❖ Es necesario usar técnicas de Bootstrap. Por ejemplo, para intervalos bilaterales:

```
library(boot); library(npsm)  
cor.boot.ci(...,method="spearman")  
cor.boot.ci(...,method="kendall")
```

- ❖ Para intervalos unilaterales o pruebas de hipótesis a una cola, deberán escribir su propio código.

Coeficiente Xi (Chatterjee, 2021)

- ❖ **Parámetro:** Coeficiente Xi de dependencia funcional de Y en función de X , en la población objetivo.

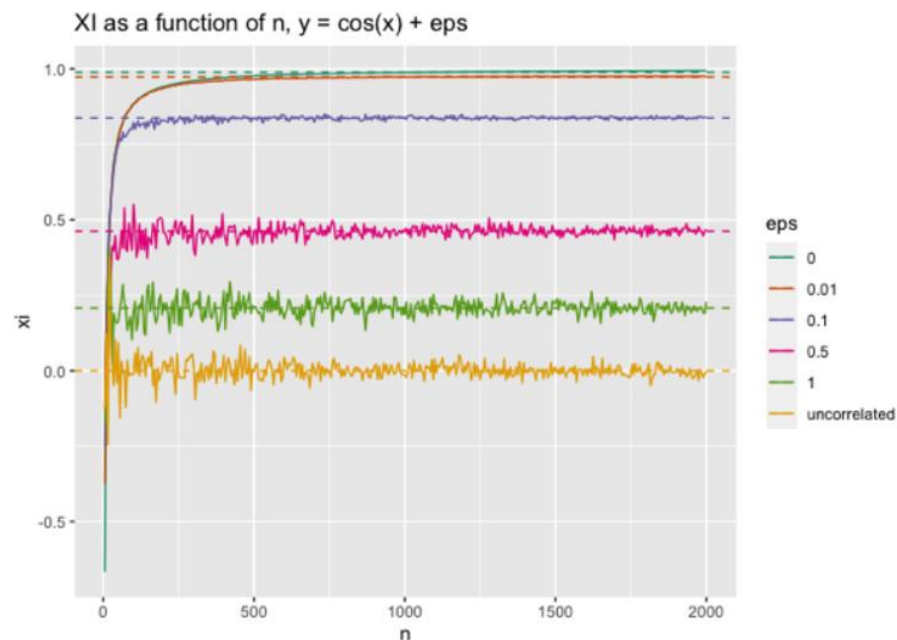
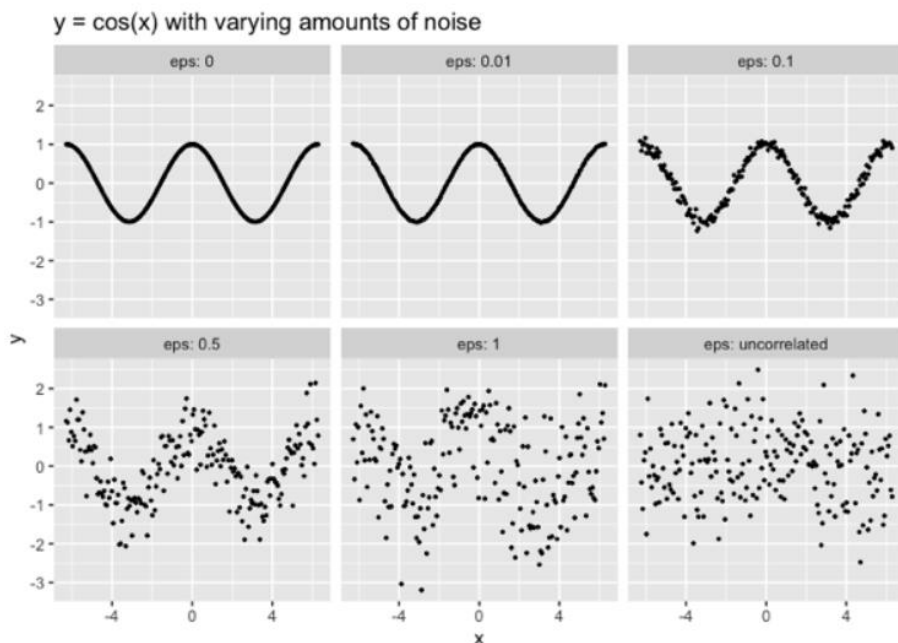
Let (X, Y) be a pair of random variables, where Y is not a constant. Let $(X_1, Y_1), \dots, (X_n, Y_n)$ be i.i.d. pairs with the same law as (X, Y) , where $n \geq 2$. The new coefficient has a simpler formula if the X_i 's and the Y_i 's have no ties. This simpler formula is presented first, and then the general case is given. Suppose that the X_i 's and the Y_i 's have no ties. Rearrange the data as $(X_{(1)}, Y_{(1)}), \dots, (X_{(n)}, Y_{(n)})$ such that $X_{(1)} \leq \dots \leq X_{(n)}$. Since the X_i 's have no ties, there is a unique way of doing this. Let r_i be the rank of $Y_{(i)}$, that is, the number of j such that $Y_{(j)} \leq Y_{(i)}$. The new correlation coefficient is defined as

$$\xi_n(X, Y) := 1 - \frac{3 \sum_{i=1}^{n-1} |r_{i+1} - r_i|}{n^2 - 1}. \quad (1.1)$$

No relac
func
[0,1]
Relac
Func
perf

- ❖ No es simétrico. Si su valor es cercano a 1 indica relación funcional entre las variables. Si es 0, da evidencia de independencia.
- ❖ Requiere tamaños de muestra grandes para ser concluyente.
- ❖ Es más potente cuando la relación funcional es suave y no monótona

Coeficiente Xi (Chatterjee, 2021) (II)



Theorem 2.1. Suppose that X and Y are independent and Y is continuous. Then $\sqrt{n}\xi_n(X, Y) \rightarrow N(0, 2/5)$ in distribution as $n \rightarrow \infty$.

❖ **Código en R:** `calculateXI` o `xicor` (librería XICOR)

Implementación de un modelo

ANOVA

en R

Relac. var. cualit. y cuantit.

1. Identificación

Gráficos por cada categoría {
Boxplot
Histogramas

- Relación que afecte únicamente la tendencia central

- Se mantienen iguales:

- Variabilidad

- Simetría

- Curtosis

ANOVA

iid $N(0, \sigma^2)$

$$Y_{ij} = \mu + \tau_i + \epsilon_{ij}$$

2. Estimación e inferencia

Cambio medio en el grupo i

Prueba F

$$\begin{cases} H_0: \tau_1 = \tau_2 = \dots = \tau_k = 0 \\ H_1: \exists \tau_i \neq 0 \end{cases}$$

$$\textcircled{1} \sum_i \sum_j (Y_{ij} - \bar{Y}_{..})^2 = \sum_i \sum_j (Y_{ij} - \bar{Y}_{i.})^2 + \sum_i n_i (\bar{Y}_{i.} - \bar{Y}_{..})^2$$

SC Tot SC error SC trat

Si H_0 es cierta:

Grande

Pequeño

Si H_0 es falsa:

Pequeño

grande

$$F = \frac{SC_{\text{trat}} / (k-1)}{SC_{\text{error}} / (N-k)} \stackrel{H_0}{\sim} F(k-1, N-k)$$

τ : "Rechazar H_0 si $f_c > f_{1-\alpha}(k-1, N-k)$ "