



# Práctica 1

**Alumno:** Francisco Javier Piqueras Martínez

**Asignatura:** Modelado Estadístico de Datos

**Fecha de entrega:** 16 de diciembre de 2019

## Índice

<b>1.</b>	<b>Descripción del documento.....</b>	<b>3</b>
<b>2.</b>	<b>Ejercicios.....</b>	<b>3</b>
2.1.	Ejercicio 1.....	3
2.2.	Ejercicio 2.....	4
2.3.	Ejercicio 3.....	5
2.4.	Ejercicio 4.....	5
2.5.	Ejercicio 5.....	8

## 1. Descripción del documento

Este documento consiste en la realización de la Práctica 1 de la asignatura de MED.

## 2. Ejercicios

### 2.1. Ejercicio 1

**Enunciado:** Se ha realizado un estudio para ver si influye la metodología docente a la hora de aprobar. Para ello **75** estudiantes han recibido la metodología 1 y **25** la metodología 2. De cada estudiante se ha registrado si al final aprobaban (1) o no (2). Los datos experimentales se dan en la tabla siguiente, donde el número de individuos con perfil aprobar = 1 y metodología = 1 es 35, con perfil aprobar = 1 y metodología = 2 es 15, con perfil aprobar = 2 y metodología = 1 es 40 y con perfil aprobar = 2 y metodología = 2 es 10. ¿Hay diferencias estadísticamente significativas entre las dos metodologías?

aprobar	metodologia
1	1
...	...
1	1
1	2
...	...
1	2
2	1
...	...
2	1
2	2
...	...
2	2

### Resolución:

Para saber si hay diferencias estadísticamente significativas entre las personas que han aprobado o no para cada una de las metodologías, se va a realizar un estudio de diferencia de proporciones poblacionales entre la metodología 1 y la metodología 2.

Puesto que la variable respuesta “aprobar” solo puede tomar dos valores, así como la variable explicativa “metodología”, y como el tamaño muestral es grande ( $>5$ ), vamos a aplicar el **test z de diferencia de proporciones**, enmarcado en un esquema  $D \leftarrow D$  cuyo parámetro poblacional  $\theta$  de interés es la diferencia de proporciones poblacionales entre  $\pi_1$  y  $\pi_2$  ( $\theta = \pi_1 - \pi_2$ ). Para ello se considera el estadístico  $\hat{\theta}$  dado por la v.a.  $\hat{\Pi}_1 - \hat{\Pi}_2$  (es decir,  $\hat{\theta} = \hat{\Pi}_1 - \hat{\Pi}_2$ ).

En primer lugar, se muestran los cálculos detallados para el intervalo de confianza al 95% dado por:

$$IC_{95\%}(\pi_1 - \pi_2) = \left( \frac{35}{75} - \frac{15}{25} \mp \sqrt{\frac{35}{75} \left(1 - \frac{35}{75}\right) \frac{1}{75} + \frac{15}{25} \left(1 - \frac{15}{25}\right) \frac{1}{25}} \right) = (-0.246196, -0.020404)$$

Y para el contraste de hipótesis dado por:

$$z = \frac{\frac{35}{75} - \frac{15}{25}}{\sqrt{\frac{35+15}{75+25} \left(1 - \frac{35+15}{75+25}\right) \left(\frac{1}{75} + \frac{1}{25}\right)}} = -1.1547$$

Con lo que al ser  $|z| = 1.1547 < 1.96$  se acepta la hipótesis nula o de igualdad y se concluye en que no existe diferencias significativas entre las poblaciones. Es decir, **no hay diferencia estadísticamente significativa entre ambas metodologías.**

## 2.2. Ejercicio 2

**Enunciado:** En el modelo de regresión lineal, se define la matriz  $H$  (matriz “hat”) como aquella matriz que pone el sombrero a la  $y$ , es decir que  $\hat{y} = Hy$ , entonces se verifica que  $H$  es simétrica e idempotente.

- a) Verdadero
- b) Falso

**Resolución:**

**Verdadero.**  $H$  es simétrica e idempotente.

En primer lugar, se definen los conceptos simétrica e idempotente aplicado a matrices. Se dice que una matriz es simétrica cuando es igual a su traspuesta. Se dice que una matriz es idempotente cuando es igual a el producto por sí misma.

Por el método mínimos cuadrados, sabemos:

$$\beta = (X^T X)^{-1} X^T y \quad \hat{y} = X\beta = X(X^T X)^{-1} X^T y$$

Por lo tanto:

$$\hat{y} = Hy \rightarrow H = \frac{\hat{y}}{y} = \frac{X(X^T X)^{-1} X^T y}{y} = X(X^T X)^{-1} X^T$$

Demostración de que  $H$  es simétrica:

$$\begin{aligned} H = H^T &\rightarrow H^T = (X(X^T X)^{-1} X^T)^T = (X^T)^T [(X^T X)^{-1}]^T X^T = X[(X^T X)^T]^{-1} X^T \\ &= X(X^T X)^{-1} X^T = H \end{aligned}$$

Demostración de que H es idempotente:

$$H = HH \rightarrow HH = X(X^T X)^{-1} X^T X (X^T X)^{-1} X^T = X(X^T X)^{-1} (X^T X) (X^T X)^{-1} X^T \\ = X(X^T X)^{-1} I X^T = X(X^T X)^{-1} X^T = H$$

### 2.3. Ejercicio 3

**Enunciado:** En el modelo de regresión lineal, se define la matriz H (matriz “hat”) como aquella matriz que pone el sombrero a la y, es decir que  $\hat{y} = Hy$ , entonces se verifica que los elementos  $h_{ii}$  de la diagonal de H, vienen dados por  $h_{ii} = x_i^T (X^T X)^{-1} x_i$ , siendo  $x_i^T = (1 \ x_{i1} \dots \ x_{ip})$ .

- a) Verdadero
- b) Falso

**Resolución:**

Verdadero.

Asumiendo, tal y como hemos calculado en el ejercicio anterior:

$$H = X(X^T X)^{-1} X^T$$

Vamos a asumir que  $D = (X^T X)^{-1}$ .

$$H = \begin{pmatrix} 1 & x_{11} & \dots & x_{1p} \\ 1 & x_{21} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{pmatrix} * \begin{pmatrix} d_{11} & d_{12} & \dots & d_{1(p+1)} \\ d_{21} & d_{22} & \dots & d_{2(p+1)} \\ \vdots & \vdots & \ddots & \vdots \\ d_{(p+1)1} & d_{(p+1)2} & \dots & d_{(p+1)(p+1)} \end{pmatrix} * \begin{pmatrix} 1 & 1 & \dots & 1 \\ x_{11} & x_{21} & \dots & x_{n1} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1p} & x_{2p} & \dots & x_{np} \end{pmatrix} \\ = \begin{pmatrix} x_1^t d_1 & x_1^t d_2 & \dots & x_1^t d_{p+1} \\ x_2^t d_1 & x_2^t d_2 & \dots & x_2^t d_{p+1} \\ \vdots & \vdots & \ddots & \vdots \\ x_n^t d_1 & x_n^t d_2 & \dots & x_n^t d_{p+1} \end{pmatrix} * \begin{pmatrix} 1 & 1 & \dots & 1 \\ x_{11} & x_{21} & \dots & x_{n1} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1p} & x_{2p} & \dots & x_{np} \end{pmatrix}$$

Por lo tanto, los elementos de la matriz H que pertenecen a la diagonal vienen dados por la siguiente fórmula:

$$h_{ii} = (x_i^t d_1 \ x_i^t d_2 \ \dots \ x_i^t d_{p+1}) * x_i = x_i^t * (d_1 \ d_2 \ \dots \ d_{p+1}) * x_i = x_i^t * (X^T X)^{-1} * x_i$$

### 2.4. Ejercicio 4

**Enunciado:** Dado el código en R, se pide rellenar el mayó número de ‘xxx’ posible:

```
rm(list=ls())

datos=read.table('c_d_1.txt',header=T)
attach(datos)

ind1=which(exp==1)

ind2=which(exp==2)
n1=length(rta[ind1]); n1
n2=length(rta[ind2]); n2
tapply(rta,exp,mean)
tapply(rta,exp,sd)

t.test(rta[ind1],rta[ind2],var.equal=TRUE)
```

proporciona el siguiente resultado

```
> n1=length(rta[ind1]); n1
[1] 7
> n2=length(rta[ind2]); n2
[1] 10
> tapply(rta,exp,mean)
      1      2
25.85714 26.20000
> tapply(rta,exp,sd)
      1      2
9.856108 8.866917
```

#### Two Sample t-test

```
data: rta[ind1] and rta[ind2]
t = -0.075009, df = 15, p-value = 0.9412
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -10.085497  9.399783
sample estimates:
mean of x mean of y
 25.85714  26.20000
```

A continuación se escribe el siguiente código:

```
exp2=1*(exp==2)
summary(lm(data = datos,formula = rta ~ exp2))
```

que proporciona el siguiente resultado.

	Estimate	Std. Error	t value	$Pr(>  t )$
(Intercept)	xxx	xxx	xxx	xxx
exp2	xxx	xxx	xxx	xxx

Tabla 2: Coeficientes de RL con  $p = 1$  sin información rellena

Residual standard error: xxx on xxx degrees of freedom

Multiple R-squared: xxx, Adjusted R-squared: xxx

F-statistic: xxx on xxx and xxx, p-value: xxx

**Resolución:**

	Estimate	Std. Error	T value	Pr(>  t )
(Intercept)	<b>25.85714</b>	<b>3.5057</b>	<b>7.376</b>	<b>0.000000321</b>
Exp2	<b>0.3429</b>	<b>4.572</b>	<b>0.075</b>	<b>0.9412</b>

Residual standard error: xxx on **15** degrees of freedom

Multiple R-squared: xxx, Adjusted R-squared: xxx

F-statistic: **0.005626** on **1** and **15**, p-value: **0.9412**

Intercept  $\rightarrow$  (X=0), por lo tanto, exp1.

Estimate

La media de exp1 es 25.85714

$$\text{exp2: } |\widehat{y}_1 - \widehat{y}_2| = |25.85714 - 26.2| = |-0.34286| = 0.34286$$

Asumiendo homocedasticidad ( $\sigma_1^2 = \sigma_2^2$ ):

$$\hat{s}_c^2 = \frac{(n_1 - 1)\hat{s}_1^2 - (n_2 - 1)\hat{s}_2^2}{(n_1 - 1) + (n_2 - 1)} = 86.0305$$

Std. Error:

$$\text{Intercept} \rightarrow \sqrt{\hat{s}_c^2 \left(\frac{1}{n_1}\right)} = 3.5057$$

$$\text{exp2} \rightarrow \sqrt{\hat{s}_c^2 \left(\frac{1}{n_1} + \frac{1}{n_2}\right)} = 4.572$$

t-Student:

$$\text{exp2} \rightarrow t = \frac{\widehat{y}_1 - \widehat{y}_2}{4.572} = 0.075$$

$$\text{Intercept} \rightarrow t = \frac{\widehat{y}_1}{3.5057} = 7.376$$

p-value:

$$\text{exp2} \rightarrow pval = 2(1 - pt(|0.075|, 15)) = 0.9412$$

$$\text{Intercept} \rightarrow pval = 2(1 - pt(|7.376|, 15)) = 0.000000321$$

Grados de libertad:

$$df = (n_1 - 1) + (n_2 - 1) = 15$$

f-statistic:

$$f_{statistic} = t^2 = 0.005626$$

## 2.5. Ejercicio 5

**Enunciado:** Dado el código en R, se pide rellenar el mayor número de 'xxx' posible:

```
rm(list=ls())

datos=read.table('c_n_1.txt',header=T)
attach(datos)

ind1=which(exp==1);

ind2=which(exp==2);
ind3=which(exp==3);
n1=length(rta[ind1]); n1
n2=length(rta[ind2]); n2
n3=length(rta[ind3]); n3
tapply(rta,exp,mean); tapply(rta,exp,sd)
summary(aov(rta~factor(exp)))
```

proporciona el siguiente resultado

```
> n1=length(rta[ind1]); n1
[1] 7
> n2=length(rta[ind2]); n2
[1] 10
> n3=length(rta[ind3]); n3
[1] 5
> tapply(rta,exp,mean); tapply(rta,exp,sd)
      1      2      3
25.85714 26.20000 22.60000
      1      2      3
9.856108 8.866917 8.876936
      Df Sum Sq Mean Sq F value Pr(>F)
factor(exp) 2  46.7   23.35   0.276  0.762
Residuals 19 1605.7   84.51
```

A continuación se escribe el siguiente código:

```
exp2=1*(exp==2)
exp3=1*(exp==3)
summary(lm(data=datos, formula=rta ~ exp2+exp3))
```

que proporciona el siguiente resultado donde se pide rellenar el mayor número posible de valores marcados con xxx.

	Estimate	Std. Error	t value	Pr(>  t )
(Intercept)	xxx	xxx	xxx	xxx
exp2	xxx	xxx	xxx	xxx
exp3	xxx	xxx	xxx	xxx

Tabla 3: Coeficientes de RL con  $p = 2$  sin información rellenada

Residual standard error: xxx on xxx degrees of freedom

Multiple R-squared: xxx, Adjusted R-squared: xxx

F-statistic: xxx on xxx and xxx, p-value: xxx



**Resolución:**

	Estimate	Std. Error	T value	Pr(>  t )
(Intercept)	<b>25.85714</b>	<b>3.4746</b>	<b>7.442</b>	<b>0.0000000482</b>
Exp2	<b>0.3429</b>	<b>4.5303</b>	<b>0.076</b>	<b>0.940</b>
Exp3	<b>3.2571</b>	<b>5.3828</b>	<b>-0.605</b>	<b>0.552</b>

Residual standard error: xxx on **19** degrees of freedom

Multiple R-squared: xxx, Adjusted R-squared: xxx

F-statistic: xxx on **2** and **19**, p-value: xxx

Intercept  $\rightarrow$  (X=0), por lo tanto, exp1.

Estimate

La media de exp1 es 25.85714

$$\text{exp2: } |\widehat{y}_1 - \widehat{y}_2| = |25.85714 - 26.2| = |-0.34286| = 0.34286$$

$$\text{exp3: } |\widehat{y}_1 - \widehat{y}_3| = |25.85714 - 22.6| = |-3.25714| = 3.25714$$

Asumiendo homocedasticidad ( $\sigma_1^2 = \sigma_2^2 = \sigma_3^2$ ):

$$\hat{s}_c^2 = \frac{(n_1 - 1)\hat{s}_1^2 + (n_2 - 1)\hat{s}_2^2 + (n_3 - 1)\hat{s}_3^2}{(n_1 - 1) + (n_2 - 1) + (n_3 - 1)} = 84.4807$$

Std. Error:

$$\text{Intercept} \rightarrow \sqrt{\hat{s}_c^2 \left(\frac{1}{n_1}\right)} = 3.4746$$

$$\text{exp2} \rightarrow \sqrt{\hat{s}_c^2 \left(\frac{1}{n_1} + \frac{1}{n_2}\right)} = 4.5303$$

$$\text{exp3} \rightarrow \sqrt{\hat{s}_c^2 \left(\frac{1}{n_1} + \frac{1}{n_2}\right)} = 5.3828$$

t-Student:

$$\text{exp3} \rightarrow t = \frac{\widehat{y}_1 - \widehat{y}_3}{5.3828} = -0.605$$

$$\text{exp2} \rightarrow t = \frac{\widehat{y}_1 - \widehat{y}_2}{4.5303} = 0.076$$

$$\text{Intercept} \rightarrow t = \frac{\widehat{y}_1}{3.4740} = 7.442$$

p-value:

$$\text{exp3} \rightarrow pval = 2(1 - pt(|-0.605|, 19)) = 0.552$$

$$\text{exp2} \rightarrow pval = 2(1 - pt(|0.076|, 19)) = 0.940$$

$$\text{Intercept} \rightarrow pval = 2(1 - pt(|7.442|, 19)) = 0.0000000482$$

Grados de libertad:

$$df = (n_1 - 1) + (n_2 - 1) + (n_3 - 1) = 19$$