

Modelado estadístico de datos: Práctica 1

David Alarcón Rubio

Diciembre 2019

Ejercicio 1.

1. (CALC) (2 puntos) Se ha realizado un estudio para ver si influye la metodología docente a la hora de aprobar. Para ello 50 estudiantes han recibido la metodología 1 y 50 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 numero 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 estadadisticamente significativas entre las dos metodología?

Calcularemos el test z de diferencia de dos proporciones que esta enmarcado en el esquema $D \leftarrow D$ que indica que se esta intentado explicar la variable de respuesta Y (aprobados) dicotómica a través de la variable explicativa X (metodología) también dicotómica.

Comenzamos realizando la tabla resumen de Y=aprobar por M=metodología.

	M=1	M=2	r
Y=1	35	15	50
Y=2	40	10	50
n	75	25	100

Tabla de resultados

Ecuación

$$\begin{aligned}
 IC &= \left(\frac{35}{75} - \frac{15}{25} \right) \pm 1,96 \sqrt{\frac{35}{75} \cdot \left(1 - \frac{35}{75}\right) \cdot \frac{1}{50} + \frac{15}{25} \cdot \left(1 - \frac{35}{75}\right) \cdot \frac{1}{50}} \\
 &= (0,4666 - 0,6) \pm 1,96 \cdot 0,0988 \\
 &= (-0,3288, 0,0588)
 \end{aligned} \tag{1}$$

Ecuación

$$\begin{aligned}
 Z &= \frac{\frac{35}{75} - \frac{15}{25}}{\sqrt{\frac{35+14}{75+25} \cdot \left(1 - \frac{35+14}{75+25}\right) \cdot \left(\frac{1}{50} + \frac{1}{50}\right)}} \\
 &= \frac{-0,1334}{\sqrt{0,0099}} = \frac{-0,1334}{0,0999} = -1,33
 \end{aligned} \tag{2}$$

Podemos comprobar que $|Z| = 1,33 < 1,96$ luego comprobamos que no es estadísticamente significativa la diferencia de aprobados entre las dos metodologías.

Ejercicio 2.

2. (CALC) (1 punto) 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.

Si $\hat{y} = Hy$, y dado que:

Equation

$$H = X(X^tX)^{-1}X^t \quad (3)$$

Una matriz es simétrica si es una matriz cuadrada, la cual tiene la característica de ser igual a su traspuesta.

Dado que se comprueba que H es una matriz cuadrada $n \times n$:

Equation

$$\begin{aligned} H &= X_{nm}(X_{mn}^tX_{nm})^{-1}X_{mn}^t = H_{nn} \\ \text{Sea, } W_{mm} &= (X_{mn}^tX_{nm})^{-1} \\ \text{Entonces, } H &= X_{nm}W_{mm}X_{mn}^t = H_{nn} \end{aligned} \quad (4)$$

Aplicado las siguientes propiedades matriciales:

- $(B \cdot C)^t = C^t \cdot B^t$
- $(A \cdot B \cdot C)^t = C^t \cdot B^t \cdot A^t$
- $(B^t)^t = B$
- $(A^{-1})^t = (A^t)^{-1}$
- $A^t \cdot A = (A^tA)^t$

Obtenemos que:

Equation

$$H^t = (X(X^t X)^{-1} X^t)^t = (X^t)^t ((X^t X)^{-1})^t X^t = X((X^t X)^t)^{-1} X^t = X(X^t X)^{-1} X^t = H \quad (5)$$

Por lo que comprobamos que H es una matriz Simétrica.

Una matriz idempotente es una matriz que es igual a su cuadrado, es decir: A es idempotente si $A \times A = A^2$

Aplicado la siguiente propiedad matricial:

- $A \cdot A^{-1} = I$

Obtenemos que:

Equation

$$H^2 = (X(X^t X)^{-1} X^t)^2 = (X(X^t X)^{-1} X^t)(X(X^t X)^{-1} X^t) = X(X^t X)^{-1} (X^t X) (X^t X)^{-1} X^t = X(X^t X)^{-1} I X^t = X(X^t X)^{-1} X^t = H \quad (6)$$

Concluimos que la matriz H es una matriz Simétrica e Idempotente, por lo que la respuesta a la pregunta es Verdadero.

Ejercicio 3.

3. (CALC) (1 punto) 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.

Siendo: $x_i^t = (1x_{i1}...x_{ip})$ Entonces:

$$X_{np+1}^t = \begin{pmatrix} 1 & x_{11} & \dots & x_{1p} \\ 1 & x_{21} & \dots & x_{2p} \\ \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots \\ 1 & x_{n1} & \dots & x_{np} \end{pmatrix} \quad (7)$$

Si se denota por $(q_{ij}) = (X^t X)^{-1}$, que tiene dimensión $(p + 1) \times (p + 1)$ y se realiza el producto matricial, se tiene que:

$$\begin{aligned} H &= \begin{pmatrix} 1 & x_{11} & \dots & x_{1p} \\ 1 & x_{21} & \dots & x_{2p} \\ \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots \\ 1 & x_{n1} & \dots & x_{np} \end{pmatrix} \begin{pmatrix} q_{11} & q_{12} & \dots & q_{1p} \\ q_{21} & q_{22} & \dots & q_{2p} \\ \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots \\ q_{p+11} & q_{p+12} & \dots & q_{p+1p+1} \end{pmatrix} \begin{pmatrix} 1 & 1 & \dots & 1 \\ x_{11} & x_{21} & \dots & x_{n1} \\ \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots \\ x_{1p} & x_{2p} & \dots & x_{np} \end{pmatrix} \\ &= \begin{pmatrix} x_1^t q_{11} & x_1^t q_{12} & \dots & x_1^t q_{1p} \\ x_2^t q_{21} & x_2^t q_{22} & \dots & x_2^t q_{2p} \\ \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots \\ x_n^t q_{p+11} & x_n^t q_{p+12} & \dots & x_n^t q_{p+1p+1} \end{pmatrix} \begin{pmatrix} 1 & 1 & \dots & 1 \\ x_{11} & x_{21} & \dots & x_{n1} \\ \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots \\ x_{1p} & x_{2p} & \dots & x_{np} \end{pmatrix} \quad (8) \end{aligned}$$

Por lo que:

$$h_{ii} = (x_i^t q_{11} \quad x_i^t q_{12} \quad \dots \quad x_i^t q_{p+1}) x_i = x_i^t (q_{11} \quad q_{12} \quad \dots \quad q_{p+1}) x_i = x_i^t (X^t X)^{-1} x_i \quad (9)$$

Ejercicio 4.

4.(CALC) (2 puntos) El siguiente código en R

```
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

Cuadro 1: Tabla 2: Coeficientes de RL con p = 1 sin informacion 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
Se pide rellena el mayor numero posible de valores marcados con xxx.

Ejercicio 5.

4.(CALC) (2 puntos) El siguiente código en R

```
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
```

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

Cuadro 2: Tabla 2: Coeficientes de RL con $p = 2$ sin informacion rellena

```
[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 numero posible de valores marcados con xxx.

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

Ejercicio 6.

6. (2 puntos) Se ha realizado un estudio para ver si el peso en kg (rta) de unos deportistas depende de su cintura en cm (exp1), del

numero de km de entrenamiento (exp2) y del tipo de entrenamiento (exp3=1: Body building, exp3=2: Fitness). Han participado en el estudio 26 individuos. Los datos experimentales estan en el

chero c ccd.txt alojado en el curso virtual y se muestran en la tabla 6. Se pide: Interpretar los resultados del modelo de regresion lineal con todas las variables. Repetir el analisis quitando las variables no signi

cativas. >Que sucede? Crear una variable interaccion entre exp1 y exp3 e incorporarla al modelo anterior. >Que ocurre? Elegir de los tres modelos anteriores el mejor. >Se cumplen las condiciones de aplicabilidad de la regresion lineal? Elaborar otro enunciado para estos datos. En el documento que se entregue habra que incluir el codigo utilizado.

rta	exp1	exp2	exp3
69.3	83	8	1
69.6	84	7	1
71.5	86.5	4	1
71.5	84.5	32	1
70.6	86.4	15	1
69.2	82.5	6	1
65	82	10	2
65.4	81.8	17	2
63.7	80	6	2
69	82.5	18	1
65.8	84	0	2
68.7	87.2	3	2
64.8	84	10	2
70	86	11	1
65.9	84.2	18	2
63.9	84	4	2
62.1	79	12	2
73.1	97.2	18	2
75.4	91	0	1
72.6	89.5	9	1
69.6	89.5	11	2
72.3	87.5	7	1
67.3	87.5	15	2
68	87.5	5	2
68.1	86.5	14	2
71.3	87	9	1