

Curso : Probabilidad y Estadística
 Sigla : EYP1113
 Profesores : Ricardo Aravena C. y Ricardo Olea O.

PAUTA EXAMEN

Problema 1

Cuando los países se van desarrollando, el consumo de energía también aumenta y por esta razón es importante poder predecir el consumo de energía futuro para ir planificando con anticipación como generar dicha energía, como es la construcción de hidroeléctricas o parques eólicos por nombrar algunas formas de generación de electricidad. Para poder predecir, se proponen dos modelos de regresión simple a la demanda eléctrica (en MMWh) de la región Metropolitana 2002-01 a 2015-12, el primero utiliza como regresor el Tiempo y el segundo el PIB nacional.

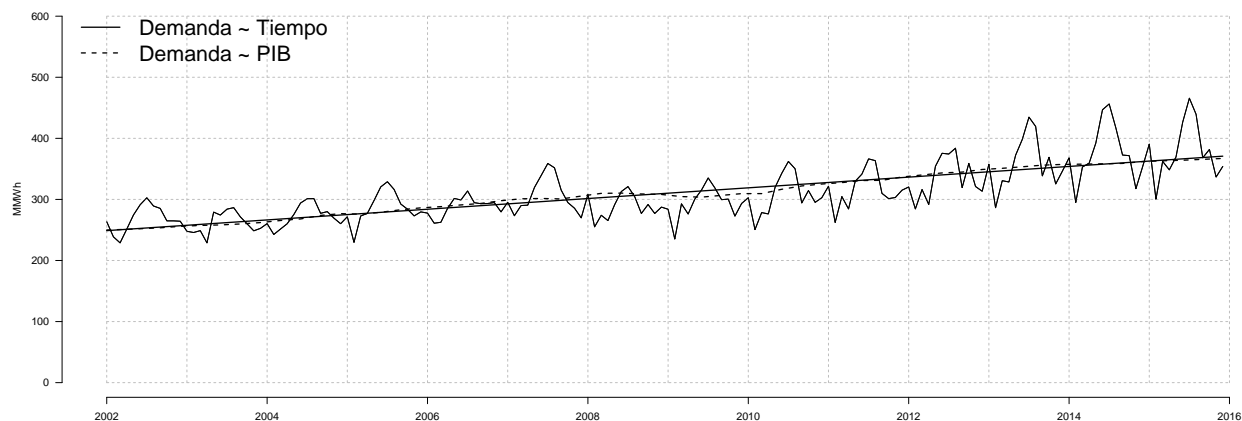


Figura 1: Demanda Electrica Región Metropolitana

```
mean(Demanda) mean(Demanda^2)
309.8397      98257.34
```

Demanda ~ Tiempo:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-1.729e+04	1.222e+03	-14.14	<2e-16 ***
Tiempo	8.759e+00	XXXXXXXXXX	XXXXX	XXXXXX

Residual standard error: 31.87 on 166 degrees of freedom
 Multiple R-squared: XXXXXX, Adjusted R-squared: XXXXXX
 F-statistic: 207.3 on 1 and 166 DF, p-value: < 2.2e-16

Demanda ~ PIB:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.766e+01	1.467e+01	6.658	3.88e-10 ***
PIB	7.445e-03	5.076e-04	14.668	< 2e-16 ***

Residual standard error: 31.54 on 166 degrees of freedom
 Multiple R-squared: 0.5645, Adjusted R-squared: 0.5618
 F-statistic: XXXXX on 1 and 166 DF, p-value: XXXXXXXXX

Complete la información faltante e indique cuál de los dos modelos ajusta mejor. Justifique su respuesta.

Respuesta

Tenemos que $n = 168$ y

$$SCT = \sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n Y_i^2 - n(\bar{Y})^2 = 168 \cdot 98257,34 - 168 \cdot (309,8397)^2 = 379125,7$$

Recordemos que

$$SCT = SCR + SCE \quad R^2 = \frac{SCR}{SCT} = 1 - \frac{SCE}{SCT} \quad r^2 = 1 - \frac{(n-1)}{(n-2)} \cdot \frac{SCE}{SCT}$$

y que en regresión simple bajo $H_0 : \beta_1 = 0$ se tiene que $T_{b_1} \sim t\text{-Student}(n-2)$ y que el estadístico $F = \frac{SCR/1}{SCE/(n-2)} = T_{b_1}^2 \sim F(1, n-2)$.

Además $\sqrt{SCE/(n-2)}$ = Residual standard error.

```
Demanda ~ Tiempo:
      Estimate Std. Error t value Pr(>|t|)
(Intercept) -1.729e+04  1.222e+03  -14.14  <2e-16 ***
Tiempo       8.759e+00  0.6083519   14.40  <2.2e-16 ***

Residual standard error: 31.87 on 166 degrees of freedom
Multiple R-squared:  0.5553, Adjusted R-squared:  0.5526
F-statistic: 207.3 on 1 and 166 DF,  p-value: < 2.2e-16
```

[3.0 Ptos.]

```
Demanda ~ PIB:
      Estimate Std. Error t value Pr(>|t|)
(Intercept)  9.766e+01  1.467e+01   6.658 3.88e-10 ***
PIB          7.445e-03  5.076e-04  14.668  < 2e-16 ***

Residual standard error: 31.54 on 166 degrees of freedom
Multiple R-squared:  0.5645, Adjusted R-squared:  0.5618
F-statistic: 215.2 on 1 and 166 DF,  p-value: < 2e-16
```

[1.2 Ptos.]

El modelo que explica la demanda de energía por PIB presenta un mayor R^2 y r^2 , un estadístico T y F más alejado de cero y por ende un menor valor-p.

[1.8 Ptos.]

+ 1 Punto Base

Problema 2

Un especialista modela, a través de una regresión múltiple, la tasa semanal de enfermedades respiratorias (ER) en función de tres contaminantes: PM10 (micro-particulado de $10 \mu g$), O₃ (Ozono) y CO₂ (dióxido de carbono), considerando además la temperatura promedio y la humedad de la semana respectiva. Selecciona al azar una muestra de 34 semanas de los dos últimos años, ajustando diversos modelos, como se muestra a continuación:

Modelo	Variables Incluidas	R ²
1	CO ₂	27.4
2	O ₃ , TEMP	31.5
3	O ₃ , CO ₂ , TEMP	34.4
4	O ₃ , CO ₂ , HUM	41.3
5	PM10, O ₃ , CO ₂ , TEMP	49.2
6	PM10, O ₃ , CO ₂ , TEMP, HUM	54.7

La desviación estándar de ER es igual a 0,45.

- (a) Determine si existe un aporte significativo del PM10 para explicar linealmente la tasa semana de ER.
- (b) Determine si existe un aporte significativo del CO₂ para explicar linealmente la tasa semana de ER.
- (c) ¿Qué aporte conjunto es más significativo: {PM10, HUM} o {PM10, TEMP}?
- (d) Para una de las dos posibles clases jerárquicas y mediante el método forward proponga un modelo de regresión múltiple para explicar linealmente la tasa semanal de ER.

Respuesta

Tenemos que

$$SCT = 0,45^2 \cdot (34 - 1) = 6,6825$$

Con este resultado y los R^2 podemos agregar a cada modelo su suma de cuadrado de errores

$$SCE = (1 - R^2) \cdot SCT$$

Modelo	Variables Incluidas	R ²	SCE
0	NULO		6.682500
1	CO ₂	27.4	4.851495
2	O ₃ , TEMP	31.5	4.577513
3	O ₃ , CO ₂ , TEMP	34.4	4.383720
4	O ₃ , CO ₂ , HUM	41.3	3.922627
5	PM10, O ₃ , CO ₂ , TEMP	49.2	3.394710
6	PM10, O ₃ , CO ₂ , TEMP, HUM	54.7	3.027172

(a) .

SOLO se puede evaluar comparando los Mod5 vs mod3:

H0: M3 vs M1: M5, lo que equivale del mod5 imponer H0: B1=0 vs H1: B1<> 0 **[0,3 ptos]**

Test F = [(4,3837 - 3,3947) / 1] / [3,3947 / (34-5)] = 8,44 **[0,5 ptos]**

Como Fc=8,44 > F(0,95; 1; 29) = 4,18 se rechaza H0, **[0,3 ptos]**

es decir PM10 aporta en presencia O3, CO2, TEMP **[0,4 ptos]**

(b) .

Aporte de CO2 → comparar Mod1 vs nulo y comparar mod3 vs mod2

Ho: mod nulo vs H1: mod 1

[0,2 ptos]

Test F = $(6,6825 - 4,8515) / (4,8515/32) = 12,077$

[0,3 ptos]

Como $F_c = 12,077 > F(0,96; 1;32) = 4,15$ se rechaza H0. CO2 aporta "solo"

[0,3 ptos]

Ho: Mod2 vs H1: Mod3

[0,2 ptos]

Test F = $(4,5775 - 4,3837) / (4,3837/30) = 1,32$

[0,3 ptos]

Como $F_c = 1,32 < F(0,95; 1;30) = 4,17$ NO se rechaza Ho.

CO2 no aporta cuando =3 y TEMP esta en el modelo.

[0,3 ptos]

(c) .

{PM10,HUM} → Mod6 vs mod3 {PM10,TEMP} → mod6 vs mod4

1er caso: Ho: Mod3 vs H1: Mod6 ↔ Ho: $B_1=B_5=0$ vs h1: uno distinto de 0

[0,2 ptos]

Test F = $[(4,3837 - 3,0272) / 2] / [3,0272 / (34 - 6)] = 6,27$

[0,2 ptos]

Como $F_c = 6,27 > F(0,95; 2; 28) = 3,34$ se rechaza Ho.

[0,2 ptos]

2do caso Ho: Mod4 vs H1: Mod6 ↔ Ho: $B_1=B_4=0$ vs h1: uno distinto de 0

[0,2 ptos]

Test F = $[(3,9226 - 3,0272) / 2] / [3,0272 / (34 - 6)] = 4,14$

[0,2 ptos]

Como $F_c = 4,14 > F(0,95; 2; 28) = 3,34$ se rechaza Ho.

[0,2 ptos]

*** PERO el aporte mas importante es de {PM10, HUM} una mayor SC

[0,3 ptos]

(d) .

Hay dos clases:

$M1 = \{\text{mod0}, \text{mod1}, \text{mod3}, \text{mod5}, \text{mod6}\}$ o bien $M2 = \{\text{Mod0}, \text{Mod2}, \text{Mod3}, \text{Mod5}, \text{Mod6}\}$

Elijo Clase $M2 = \{\text{Mod0}, \text{Mod2}, \text{Mod3}, \text{Mod5}, \text{Mod6}\} \rightarrow$ hacia adelante

[0,5 ptos]

1er paso H_0 : Modelo Nulo vs H_1 : Mod2

Test $F = [(6,6825 - 4,5775)/2] / [4,5775/31] = 7,12 > F(0,95; 2; 31) = 3,30$ se rechaza H_0 .. se itera

[0,5 ptos]

2d paso H_0 : mod2 vs H_1 : mod3 ... hecho en b) NO se rechaza H_0 .. FIN..

[0,3 ptos]

el mejor modelo es Mod2

[0,2 ptos]

ALTERNATIVAMENTE

ELIJO clase $M1 = \{\text{mod0}, \text{mod1}, \text{mod3}, \text{mod5}, \text{mod6}\}$

[0,5 ptos]

1er H_0 : Modelo nulo vs H_1 : mod1 .. hecho en b) 1era parte.. se rechaza H_0 .. se itera

[0,2 ptos]

2do H_0 : mod1 vs H_1 : mod3

[0,3 ptos]

Test $F = [(4,8515 - 4,3837) / 2] / [4,3837/(34-4)] = 1,6 < F(0,95; 2; 30) = 3,32$

[0,4 ptos]

NO se rechaza $H_0 \rightarrow$ FIN .. el mejor modelo es Mod1.

[0,1 ptos]

+ 1 Punto Base

Problema 3

Como han sabido, diversos partidos de futbol se definen en lanzamientos penales. Usted, interesado en determinar si el comportamiento observado en este tipo de campeonato (selecciones) también se da en los partidos a nivel de club, revisa las estadísticas y determina que en un total de 125 partidos a nivel de selección, en 42 se tuvo que llegar a la instancia de los penales. Mientras, que una muestra de 275 definiciones a nivel de equipo muestra que en 72 de ellos se definió por penales.

- (a) ¿Los datos anteriores permiten afirmar que es mas probable una definición a penales entre selecciones que entre clubes? Especifique hipótesis, test, valor-p y conclusión para un nivel de significancia del 5 %.

Por otra parte, Bravo se lució atajando penales, Usted constata que en los 42 partidos de mundiales que se han definidos por penales, solo han sido atajados 21 penales (de más de 150 que efectuaron), mientras que en los 72 partidos de otras competencias han sido atajados 58 penales (de más de 250 que se efectuaron).

- (b) ¿Para qué valor de significancia existe evidencia que permita afirmar que la tasa de penales atajados difiere según tipo de competencia? Asuma que el número de penales atajados en las definiciones se comporta como una variable aleatorio con distribución Poisson.

Respuesta

- (a) .

Hipótesis: $H_0: P_s = P_c$ vs $H_1: P_s > P_c$ [0.5 ptos]

Selección: 42 de 125 $\rightarrow P_s = 0,336$ Clubes: 72 de 275 $\rightarrow P_c = 0,262$ [0.5 ptos]

$P_{comun} = (42+72) \text{ de } (125+275) \rightarrow P = 0,285$ [0.5 ptos]

Test $z = 1,5196$ [0.5 ptos]

$\rightarrow \text{valor-p} = P(Z > 1,52) = 0,064$ [0.5 ptos]

Como valor-p no es menor que alfa, no hay evidencia para rechazar H_0 . Es decir, no se puede afirmar que sea más probable una definición a penales entre selecciones que entre clubes.

[0.5 ptos]

(b) .

DOS PROCESOS POISSON INDEPENDIENTES. – COMPARACION DE TASAS

Ho: $L1 = L2$ vs H1: $L1 \neq L2$

[0.5 ptos]

Estimaciones $L1 = 21/42 = 0,5$ $L2 = 58/72 = 0,806$,

[0.5 ptos]

Con L-Comun = $(21+58)/(42+72) = 0,693$

[0.5 ptos]

Test $z = (0,5 - 0,806) / (\text{raíz}(0,693 \times (1/42+1/72))) = -1,893$

[0.5 ptos]

Valor-p = $P(Z < -1,893) = 0,029$

[0.5 ptos]

Por tanto, para todo $\alpha > 0,029$ existe evidencia que permite afirmar que la tasa de penales atajados difiere según tipo de competencia.

[0.5 ptos]

+ 1 Punto Base

Problema 4

El mejor ajuste de regresión lineal simple para la demanda eléctrica en la Región Metropolitana visto en el problema 1, entregó los siguientes residuos

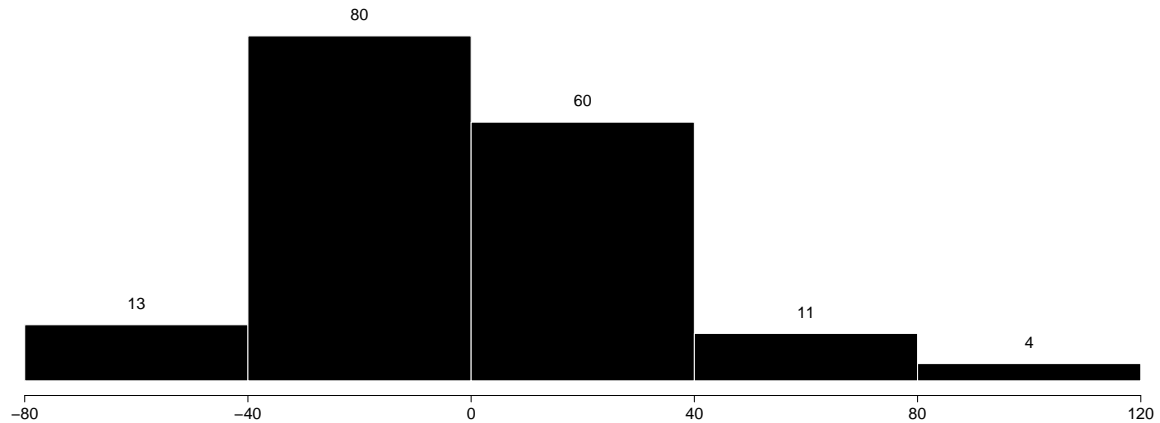


Figura 2: Histograma de frecuencia para residuos mejor ajuste demanda eléctrica región metropolitana

Min. 1st Qu. Median
-70.700 -21.420 -3.266

¿Entre una distribución Normal o Log-Normal, cuál ajusta mejor?

Respuesta

Definamos como Z a los residuos del mejor modelo del problema 1.

$H_0 : Z \sim \text{Normal}(\mu, \sigma)$ vs $H_a : Z \not\sim \text{Normal}(\mu, \sigma)$

Si z_p es el percentil $p \times 100\%$, entonces

$$P(Z \leq z_p) = \Phi\left[\frac{z_p - \mu}{\sigma}\right] = p \rightarrow z_p = \mu + \sigma \Phi^{-1}(p) \quad [0.2 \text{ Ptos.}]$$

es decir,

$$[0.2 \text{ Ptos.}] \quad z_{0,25} = \mu - \sigma 0,6745 \quad \text{y} \quad z_{0,50} = \mu \quad [0.2 \text{ Ptos.}]$$

Despejando se tiene que

$$[0.2 \text{ Ptos.}] \quad \mu = -3,266 \quad \text{y} \quad \sigma = 26,91475 \quad [0.2 \text{ Ptos.}]$$

Al calcular los valores esperados para cada intervalo, tenemos que el 5to tiene una esperanza menor a 5.

	Observado	Probabilidad	Esperado
[-80, -40)	13	0.086153597	14.4738043
[-40, 0)	80	0.462137947	77.6391751
[0, +40)	60	0.397738276	66.8200304
[+40, +80)	11	0.052981772	8.9009378
[+80, +120]	4	0.000988407	0.1660524
Total	168	1.000000000	168.0000000

[1.0 Ptos.]

Dado lo anterior, se procede a colapsar los últimos dos intervalos y rehacer la tabla

Observado	Probabilidad	Esperado	(O-E)^2/E
[-80, -40)	13	0.08615360	14.47380
[-40, 0)	80	0.46213795	77.63918
[0, +40)	60	0.39773828	66.82003
[+40,+120]	15	0.05397018	9.06699
Total	168	1.00000000	168.00000

[0.5 Ptos.]

El estadístico de prueba $X^2 = 4,8$ y como distribuye $\chi^2(4 - 1 - 2)$ se tiene que

$$2,5 \% < \text{valor-p} < 5 \% \quad [0.5 \text{ Ptos.}]$$

$H_0 : Z \sim \text{Log-Normal}(\lambda, \zeta)$ trasladada en α vs $H_a : Z \not\sim \text{Log-Normal}(\lambda, \zeta)$ trasladada en α

Si z_p es el percentil $p \times 100 \%$, entonces

$$P(Z \leq z_p) = P(Z - \alpha \leq z_p - \alpha) = \Phi \left[\frac{\ln(z_p - \alpha) - \lambda}{\zeta} \right] = p \rightarrow z_p = e^{\lambda + \zeta \Phi^{-1}(p)} + \alpha \quad [0.2 \text{ Ptos.}]$$

es decir,

$$z_{0,00} = \alpha \text{ (mínimo valor posible)} \quad [0.2 \text{ Ptos.}]$$

$$z_{0,25} = e^{\lambda - \zeta 0,6745} + \alpha \quad [0.2 \text{ Ptos.}]$$

$$z_{0,50} = e^{\lambda} + \alpha \quad [0.2 \text{ Ptos.}]$$

Despejando se tiene que

$$\alpha = -70,7; \quad [0.2 \text{ Ptos.}] \quad \lambda = 4,211149; \quad [0.2 \text{ Ptos.}] \quad \zeta = 0,464983 \quad [0.2 \text{ Ptos.}]$$

Al calcular los valores esperados para cada intervalo, tenemos que

Observado	Probabilidad	Esperado	(O-E)^2/E
[-80, -40)	13	0.04529522	7.609596
[-40, 0)	80	0.49521375	83.195910
[0, +40)	60	0.31628049	53.135123
[+40, +80)	11	0.10134181	17.025424
[+80,+120]	4	0.04186873	7.033946
Total	168	1.00000000	168.00000

[1.0 Ptos.]

El estadístico de prueba $X^2 = 8,27$ y como distribuye $\chi^2(5 - 1 - 3)$ se tiene que

$$\text{valor-p} < 0,5 \% \quad [0.2 \text{ Ptos.}]$$

Aunque ambos modelos son rechazados al 5 % de significancia, un mejor ajusta se logra con el modelo Normal, por tener un mayor valor-p. [0.4 Ptos.]

+ 1 Punto Base