

**Modelo Lineal**  
**PRACTICA 5 Parte 1**

1. Sea  $X$  una matriz de  $\mathfrak{R}^{n \times p}$  tal que  $X = [X_1, X_2]$ , donde  $X_1 \in \mathfrak{R}^{n \times k}$  y  $X_2 \in \mathfrak{R}^{n \times (p-k)}$ . Sean  $P_1 = X_1(X_1'X_1)^{-1}X_1'$  la matriz de proyección generada por las columnas de  $X_1$  y sea  $W = (I - P_1)X_2$  la proyección de  $X_2$  sobre el complemento ortogonal de  $X_1$ . Finalmente, sea  $P_2 = W(W'W)^{-1}W'$  la matriz de proyección correspondiente a  $W$ . Probar que

$$P = P_1 + P_2 = P_1 + (I - P_1)X_2(X_2'(I - P_1)X_2)^{-1}X_2'(I - P_1).$$

2. Supongamos que  $X \in \mathfrak{R}^{n \times p}$  contiene una columna de 1's., por ejemplo, sin pérdida de generalidad, la primera, es decir  $X = [\mathbf{1}, X_2]$  y sea  $P = X(X'X)^{-1}X'$ . Pruebe que

a)  $P = P_1 + P_2$  donde  $P_1 = n^{-1}\mathbf{1}\mathbf{1}'$  ( $\mathbf{1} \in \mathfrak{R}^n$ ,  $\mathbf{1} = (1, 1, \dots, 1)'$ ) y  $P_2 = \tilde{X}(\tilde{X}'\tilde{X})^{-1}\tilde{X}'$  siendo  $\tilde{X} = (I - n^{-1}\mathbf{1}\mathbf{1}')X_2$  la matriz con las columnas centradas.

b)  $p_{ii} \geq \frac{1}{n}$ .

c)  $p_{ii} = \frac{1}{n} + \tilde{p}_{ii}$  donde  $\tilde{p}_{ii} = (P_2)_{ii}$ .

3. Pruebe que la distancia de Cook de la  $i$ -ésima observación puede calcularse como

$$D_i = \frac{1}{p} \frac{p_{ii}}{1 - p_{ii}} r_i^2$$

4. Pruebe que  $\hat{Y}_i = (1 - p_{ii}) \mathbf{x}'_i \hat{\beta}_{(i)} + p_{ii} Y_i$ . Huber (1981) interpretó de esta igualdad que  $\frac{p_{ii}}{1 - p_{ii}}$  es el cociente entre la parte de  $\hat{Y}_i$  que se explica por  $Y_i$  y la parte de  $\hat{Y}_i$  que puede predecirse a partir de  $\mathbf{x}'_i \hat{\beta}_{(i)}$ .

5. Los datos que se muestran en la Tabla 1 (y en el archivo salud.txt) corresponden a 30 miembros de un club de salud. Las variables son

$X_1$  = peso en libras

$X_2$  = pulso en reposo

$X_3$  = fuerza del brazo y pierna (número de libras que puede levantar)

$X_4$  = tiempo en segundos en que corre 1/4 de milla

$Y$  = tiempo en segundos en que corre 1 milla.

a) Estime por cuadrados mínimos los parámetros del modelo

$$E(Y) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4$$

Realice el gráfico de los residuos estandarizados versus  $\hat{Y}$  y el QQ-plot . ¿Hay algún indicio de que haya groseras violaciones de supuestos ?

- b) Realice un boxplot y un esquema de tallo y hoja para los residuos estandarizados. ¿A qué número de observación corresponde el residuo de mayor valor absoluto?
- c) Calcule el leverage de cada observación y realice un boxplot. ¿Cuál es la observación con mayor leverage? Use los criterios dados para identificar aquellas observaciones cuyo leverage puede indicar problemas.
- d) Calcule la distancia de Cook para cada observación y realice un boxplot. ¿Cuáles son los cuatro puntos que aparecen como outliers en este gráfico? ¿Cuál es la observación de mayor influencia? ¿Qué pasa con el leverage de esta observación?
- e) Según la tabla de estimación de los coeficientes, cuáles serían los coeficientes significativos? Realice los plots de residuos parciales (component plus residual plot) para verificar estos resultados. ¿A quién corresponde el punto más alejado en el grafico correspondiente a  $X_3$ ?
- f) *Plot L-R*. El L-R plot permite graficar residuos vs. leverage de manera de identificar qué tipo de patología tiene un dato. Se calculan los residuos como:

$$a_i = \frac{e_i}{\sqrt{(ete)}}$$

y se grafica  $a_i^2$  vs  $p_{ii}$ .

Realice el L-R plot correspondiente a este problema. ¿Cómo caracterizaría a los puntos correspondientes a las observaciones 23, 28 y 30?

- g) Recalcule los estimadores de mínimos cuadrados omitiendo una a la vez las observaciones 23, 28 y 30. ¿Cuándo observa el mayor cambio en los estadísticos t?

**6.** La Tabla 2 (y el archivo salario.txt) corresponde a un estudio para relacionar el salario mensual de una muestra aleatoria de 31 empleados y un conjunto de factores que se piensa pueden determinar las diferencias en los salarios. Las variables observadas son:

$X_1$  = evaluación de trabajo

$X_2$  = sexo (1=hombre, 0=mujer)

$X_3$  = número de años en la compañía

$X_4$  = número de años en el mismo cargo

$X_5$  = ranking de rendimiento (1=no satisfactorio, 5=muy bueno)

$Y$  = salario mensual.

- a) Estime por cuadrados mínimos los parámetros del modelo

$$E(Y) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \beta_5 X_5$$

Realice el gráfico de los residuos estandarizados versus  $\hat{y}$  y el QQ-plot . ¿Hay algún residuo que le llame la atención ?

- b) Realice un boxplot y un esquema de tallo y hoja para los residuos estandarizados. ¿A qué observaciones corresponden los residuos mayores?
- c) Para cada variable independiente realice un gráfico de los residuos estandarizados versus  $X_i$ . Observe qué ocurre con la observación número 6 en cada uno de estos gráficos.
- d) Calcule el leverage de cada observación y realice un boxplot. ¿Qué observación se destaca?
- e) Realice el L-R plot para este conjunto de datos. ¿A qué conclusiones llega?

**7.** La Tabla 3 (gener2.txt) contiene un conjunto de 30 datos generados correspondientes a dos variables:  $X$  e  $Y$ .

- a) Ajuste un modelo lineal utilizando la variable  $Y$  como respuesta y la variable  $X$  como explicativa. Realice el QQ-plot de los residuos. En base a este gráfico, ¿cree que es válida la hipótesis de normalidad de los residuos? Aplique el test de Kolmogorov-Smirnov a los residuos. ¿Cuál es su conclusión?
- b) Seleccione una transformación adecuada en la familia de transformaciones de Box-Cox y repita el análisis realizado en a) con la variable transformada como respuesta.

**Tabla 1.** Club de Salud

obs	$X_1$	$X_2$	$X_3$	$X_4$	$Y$
1	217	67	260	91	481
2	141	52	190	66	292
3	152	58	203	68	338
4	153	56	183	70	357
5	180	66	170	77	396
6	193	71	178	82	429
7	162	65	160	74	345
8	180	80	170	84	469
9	205	77	188	83	425
10	168	74	170	79	358
11	232	65	220	72	393
12	146	68	158	68	346
13	173	51	243	56	279
14	155	64	198	59	311
15	212	66	220	77	401
16	138	70	180	62	267
17	147	54	150	75	404
18	197	76	228	88	442
19	165	59	188	70	368
20	125	58	160	66	295
21	161	52	190	69	391
22	132	62	163	59	264
23	257	64	313	96	487
24	236	72	225	84	481
25	149	57	173	68	374
26	161	57	173	65	309
27	198	59	220	62	367
28	245	70	218	69	469
29	141	63	193	60	252
30	177	53	183	75	338

**Tabla 2.** Salarios

obs	$X_1$	$X_2$	$X_3$	$X_4$	$X_5$	$Y$
1	350	1	2	2	5	1000
2	350	1	5	5	5	1400
3	350	0	4	4	4	1200
4	350	1	20	20	1	1800
5	425	0	10	2	3	2800
6	425	1	15	10	3	4000
7	425	0	1	1	4	2500
8	425	1	5	5	4	3000
9	600	1	10	5	2	3500
10	600	0	8	8	3	2800
11	600	0	4	3	4	2900
12	600	1	20	10	2	3800
13	600	1	7	7	5	4200
14	700	1	8	8	1	4600
15	700	0	25	15	5	5000
16	700	1	19	16	4	4600
17	700	0	20	14	5	4700
18	400	0	6	4	3	1800
19	400	1	20	8	3	3400
20	400	0	5	3	5	2000
21	500	1	22	12	3	3200
22	500	1	25	10	3	3200
23	500	0	8	3	4	2800
24	500	0	2	1	5	2400
25	800	1	10	10	3	5200
26	475	1	10	4	3	2400
27	475	0	3	3	4	2400
28	475	1	8	8	2	3000
29	475	1	6	6	4	2800
30	475	0	12	4	3	2500
31	475	0	4	2	5	2100

**Tabla 3.** Datos generados

obs	$X$	$Y$	obs	$X$	$Y$
1	0.2775889	696.20174	16	-0.886568	6.8234601
2	-0.517632	72.295359	17	0.4935309	72.466261
3	-0.303792	113.55840	18	0.7128258	2621.0541
4	-2.078444	1.8890014	19	1.0458933	293.29129
5	1.4366185	650.52223	20	-1.559868	5.5906742
6	-0.235768	96.996475	21	0.9805923	339.38955
7	-0.868704	35.627923	22	-0.597517	39.637404
8	1.8890238	31566.051	23	-1.448257	36.087594
9	-1.743802	3.9250497	24	0.443873	446.68841
10	-0.305831	445.65564	25	0.1877068	520.12349
11	0.2756379	290.04645	26	-0.295656	62.706673
12	1.0474655	402.32246	27	-0.284534	116.83186
13	1.3234957	930.59459	28	1.140876	3917.2343
14	-0.237969	56.925312	29	0.1878306	179.25413
15	0.7173053	800.51180	30	-0.153660	551.13103