

# Regresión múltiple y otras técnicas multivariadas

Tarea 03

*Rivera Torres Francisco de Jesús*

*Rodríguez Maya Jorge Daniel*

*Samayoa Donado Víctor Augusto*

*Trujillo Barrios Georgina*

*Febrero 27, 2019*

## Ejercicio 1

Suponer que se ajusta un modelo RLS a las observaciones  $(x_i, y_i)$  con  $i = 1, \dots, n$ . Mostrar que

$$SCE = \frac{S_{xx}S_{yy} - S_{xy}^2}{S_{xx}}$$

Donde:

- $SCE = \sum_i^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$
- $S_{yy} = \sum_i^n (y_i - \bar{y}_n)^2$

Entonces:

$$\begin{aligned} SCE &= \sum_{i=1}^n (y_i - \bar{\beta}_0 - \bar{\beta}_1 x_i)^2 = \sum_{i=1}^n \left( y_i - \left( \bar{y}_n - \frac{S_{xy}\bar{x}_n}{S_{xx}} \right) - \frac{S_{xy}}{S_{xx}} x_i \right)^2 = \sum_{i=1}^n \left( y_i - \bar{y}_n + \frac{S_{xy}\bar{x}_n - S_{xy}x_i}{S_{xx}} \right)^2 \\ &= \sum_{i=1}^n \left( (y_i - \bar{y}_n) - \frac{S_{xy}(x_i - \bar{x}_n)}{S_{xx}} \right)^2 = \sum_{i=1}^n \left( (y_i - \bar{y}_n)^2 - 2(y_i - \bar{y}_n)(x_i - \bar{x}_n) \frac{S_{xy}}{S_{xx}} + \frac{S_{xy}^2(x_i - \bar{x}_n)^2}{S_{xx}^2} \right) \\ &= \sum_{i=1}^n (y_i - \bar{y}_n)^2 - 2 \frac{S_{xy}}{S_{xx}} \sum_{i=1}^n (y_i - \bar{y}_n)(x_i - \bar{x}_n) + \frac{S_{xy}^2}{S_{xx}^2} \sum_{i=1}^n (x_i - \bar{x}_n)^2 = S_{yy} - 2 \frac{S_{xy}S_{xy}}{S_{xx}} + \frac{S_{xy}^2 S_{xx}}{S_{xx}^2} \\ &= S_{yy} - 2 \frac{S_{xy}^2}{S_{xx}} + \frac{S_{xy}^2}{S_{xx}} = \frac{S_{xx}S_{yy} - 2S_{xy}^2 + S_{xy}^2}{S_{xx}} = \frac{S_{xx}S_{yy} - S_{xy}^2}{S_{xx}} \end{aligned}$$

Por lo tanto

$$SCE = \frac{S_{xx}S_{yy} - S_{xy}^2}{S_{xx}}$$

## Ejercicio 2

Mostrar la desigualdad de Bonferroni. Si  $E_1, \dots, E_k$  son eventos en un espacio de probabilidad  $(\Omega, A, P)$  entonces:

$$P\left(\bigcap_{i=1}^k E_i\right) \geq 1 - \sum_{i=1}^k P(\Omega \setminus E_i)$$

*Demostración.* La demostración se realizará por inducción sobre el número de eventos en un espacio de probabilidad.

Base:  $k = 1$

$$P(\Omega) = 1, \quad \text{pero } \Omega = E \cup (\Omega \setminus E), \text{ entonces}$$

$$P(E \cup (\Omega \setminus E)) = P(E) + P(\Omega \setminus E) = 1, \quad \text{ya que son probabilidades mutuamente excluyentes}$$

$$P(E) = 1 - P(\Omega \setminus E), \quad \text{como se da la igualdad entonces tambien se satisface que}$$

$$P(E) \geq 1 - P(\Omega \setminus E)$$

Ahora, por hipótesis de inducción, suponemos que se vale para  $n$  eventos en el espacio de probabilidad, por lo que se satisface la desigualdad

$$P\left(\bigcap_{i=1}^n E_i\right) \geq 1 - \sum_{i=1}^n P(\Omega \setminus E_i)$$

Y procedemos a demostrar que siempre que se cumpla para  $n$  eventos, se debe de cumplir para  $n + 1$  eventos.

$$\begin{aligned} P\left(\bigcap_{i=1}^{n+1} E_i\right) &= P\left(\left(\bigcap_{i=1}^n E_i\right) \cap E_{n+1}\right) \\ &= P\left(\bigcap_{i=1}^n E_i\right) + P(E_{n+1}) - P\left(\left(\bigcap_{i=1}^n E_i\right) \cup E_{n+1}\right) \end{aligned}$$

Pero notemos que  $P((\bigcap_{i=1}^n E_i) \cup E_{n+1}) \leq 1$ , por lo que  $-P((\bigcap_{i=1}^n E_i) \cup E_{n+1}) \geq -1$ , entonces

$$P\left(\bigcap_{i=1}^{n+1} E_i\right) \geq P\left(\bigcap_{i=1}^n E_i\right) + P(E_{n+1}) - 1$$

aplicando la hipótesis de inducción, se tiene

$$\begin{aligned} &\geq 1 - \sum_{i=1}^n P(\Omega \setminus E_i) + P(E_{n+1}) - 1 \\ &\geq 1 - \sum_{i=1}^n P(\Omega \setminus E_i) + (1 - P(\Omega \setminus E_{n+1})) - 1 \\ &\geq 1 - \sum_{i=1}^n P(\Omega \setminus E_i) + P(\Omega \setminus E_{n+1}) \\ &\geq 1 - \sum_{i=1}^{n+1} P(\Omega \setminus E_i) \end{aligned}$$



## Ejercicio 3

Considerar los datos de ingreso y escolaridad utilizados en los ejemplos de intervalos de confianza de las notas. Reportar intervalos simultáneos de confianza 95% para las medias del ingreso por hora para 9, 15 y 19 años de escolaridad a) con el método de Bonferroni y b) con el método de Hotelling–Scheffé

En el caso del modelo RLS para datos de ingresos por hora tenemos lo siguiente:

```
rls <- function(x, y) {
  x.bar <- mean(x, na.rm = TRUE)
  y.bar <- mean(y, na.rm = TRUE)
  S.xx <- sum((x - x.bar)^2, na.rm = TRUE)
  S.xy <- sum((x - x.bar) * (y - y.bar), na.rm = TRUE)

  b0.hat <- y.bar - x.bar * S.xy / S.xx
  b1.hat <- S.xy / S.xx

  y.adj <- b0.hat + b1.hat * x
  residuos <- y - y.adj
  sigma2.hat <- sum(residuos^2) / (length(residuos) - 2)

  return(list(b0.hat = b0.hat, b1.hat = b1.hat, sigma2.hat = sigma2.hat, sxx = S.xx))
}
```

```
beta0 <- rls(ingresos$ANIOS_ESC, ingresos$ING_X_HRS)$b0.hat
beta1 <- rls(ingresos$ANIOS_ESC, ingresos$ING_X_HRS)$b1.hat
sigma2 <- rls(ingresos$ANIOS_ESC, ingresos$ING_X_HRS)$sigma2.hat
sxx <- rls(ingresos$ANIOS_ESC, ingresos$ING_X_HRS)$sxx
n <- nrow(ingresos) # pues no hay valores perdidos
sum_x2 <- sum((ingresos$ANIOS_ESC)^2)
varBeta0 = (sigma2*sum_x2)/(n*sxx)
varBeta1 = sigma2/sxx
```

$$\begin{aligned}\hat{\beta}_0 &= -10.0022517, \\ \hat{\beta}_1 &= 4.8425593, \\ \hat{\sigma}^2 &= 1823.1875432, \\ S_{xx} &= 4363.502008, \\ n &= 249, \\ \sum_{i=1}^n x_i^2 &= 3.3489 \times 10^4, \\ V(\beta_0) &= 56.195171, \\ V(\beta_1) &= 0.4178267\end{aligned}$$

### Inciso 3.a Método de Bonferroni

Los intervalos simultáneos de confianza en el método Bonferroni se definen como sigue (con el  $100 \times (1 - \frac{\alpha}{2})$  de confianza para cada intervalo):

$$\hat{\beta}_0 \pm t_{n-2}^{1-\frac{\alpha}{4}} \sqrt{\hat{V}(\beta_0)}$$

```
LI_beta0_bonf = beta0 - qt(1-(.5/4), n-2) * sqrt(varBeta0)
LS_beta0_bonf = beta0 + qt(1-(.5/4), n-2) * sqrt(varBeta0)

for(i in c(9, 15, 19)){
  LI_beta1_bonf = i*beta1 - qt(1-(.5/4), n-2) * sqrt(varBeta1)
  LS_beta1_bonf = i*beta1 + qt(1-(.5/4), n-2) * sqrt(varBeta1)

  cat(
    paste0("\n\r", " Los intervalos simultáneos de confianza para beta0 y beta_1,",
      "calculados por\nel método de Hotelling-Scheffé, para la media de",
      "ingreso por hora para\n", i,
      " años de escolaridad son:\n(",
      LI_beta0_bonf, ", ", LS_beta0_bonf, ") y\n(",
      LI_beta1_bonf, ", ", LS_beta1_bonf, "), respectivamente."))
}

##
##
  Los intervalos simultáneos de confianza para beta0 y beta_1,calculados por
## el método de Hotelling-Scheffé, para la media de ingreso por hora para
## 9 años de escolaridad son:
## (-18.6459928475289, -1.35851045480005) y
## (42.8377006734473, 44.3283666345806), respectivamente.
##
  Los intervalos simultáneos de confianza para beta0 y beta_1,calculados por
```

```
## el método de Hotelling-Scheffé, para la media de ingreso por hora para
## 15 años de escolaridad son:
## (-18.6459928475289, -1.35851045480005) y
## (71.8930564427899, 73.3837224039233), respectivamente.
##
Los intervalos simultáneos de confianza para beta0 y beta_1, calculados por
## el método de Hotelling-Scheffé, para la media de ingreso por hora para
## 19 años de escolaridad son:
## (-18.6459928475289, -1.35851045480005) y
## (91.2632936223517, 92.753959583485), respectivamente.
```

### Inciso 3.b Método de Hotelling-Scheffé

En este caso, tenemos que los intervalos simultáneos de confianza están dados por la siguiente expresión:

$$a_0\hat{\beta}_0 + a_1\hat{\beta}_1 \pm \sqrt{2F_{(2,n-2)}^{1-\alpha} + \frac{a_0^2}{n} + \frac{a_1 - a_0\bar{x}_n}{S_{xx}} \hat{\sigma}_{mco}}$$

```
#options(Encoding="UTF-8")
LI_beta0_hs = beta0 - sqrt(2*qf(1-0.25, 2, n-2))* sqrt(varBeta0)
LS_beta0_hs =beta0 + sqrt(2*qf(1-0.25, 2, n-2))* sqrt(varBeta0)

for(i in c(9, 15, 19)){
  LI_beta1_hs = i*beta1 - sqrt(2*qf(1-0.25, 2, n-2))* sqrt(varBeta1)
  LS_beta1_hs = i*beta1 + sqrt(2*qf(1-0.25, 2, n-2))* sqrt(varBeta1)

  cat(
    paste0("\n\r", " Los intervalos simultáneos de confianza para beta0 y beta_1",
      "calculados por\nel método de Hotelling-Scheffé, para la media de",
      "ingreso por hora para\n", i,
      "años de escolaridad son:\n(",
      LI_beta0_hs, ", ", LS_beta0_hs, ") y\n(",
      LI_beta1_hs, ", ", LS_beta1_hs, "), respectivamente.))
}
```

```
##
##
Los intervalos simultáneos de confianza para beta0 y beta_1 calculados por
## el método de Hotelling-Scheffé, para la media de ingreso por hora para
## 9 años de escolaridad son:
## (-22.5195934053441, 2.51509010301522) y
## (42.5036876363037, 44.6623796717242), respectivamente.
##
Los intervalos simultáneos de confianza para beta0 y beta_1 calculados por
## el método de Hotelling-Scheffé, para la media de ingreso por hora para
## 15 años de escolaridad son:
## (-22.5195934053441, 2.51509010301522) y
```

```
## (71.5590434056463, 73.7177354410669), respectivamente.
##
  Los intervalos simultáneos de confianza para beta0 y beta_1 calculados por
## el método de Hotelling-Scheffé, para la media de ingreso por hora para
## 19 años de escolaridad son:
## (-22.5195934053441, 2.51509010301522) y
## (90.929280585208, 93.0879726206286), respectivamente.
```

## Ejercicio 4

El conjunto de datos `airquality`, de paquete `datasets` de R contiene información sobre la calidad del aire en Nueva York registrada de Mayo a Septiembre de 1973 (se puede consultar más información con el comando `help("airquality")`). Para responder este ejercicio, descartar las observaciones con valores perdidos.

### Inciso 4.a

Ajustar un modelo RLS para explicar el nivel de ozono como función del  $\log_2$  de la velocidad del viento. Reportar las estimaciones de los parámetros.

```
library(tidyverse)

modelo_air <- airquality %>%
  as_tibble() %>%
  filter(!is.na(Ozone), !is.na(Wind)) %>%
  mutate(log2_Wind = log2(Wind)) %>%
  lm(formula = Ozone ~ log2_Wind)

coefficients(modelo_air)
```

```
## (Intercept)    log2_Wind
##    166.64640    -38.92431
```

### Inciso 4.b

Mostrar una gráfica de dispersión de los datos utilizados para ajustar el modelo del inciso anterior, la recta de regresión ajustada y bandas de confianza 95%. Anexar el código relacionado con el cómputo de las bandas de confianza.

```
# Código para el cálculo de los intervalos de confianza para el modelo RLS
bandas_confianza_rls <- function(datos, formula_rls, alpha = 0.95) {
  # Función que calcula las bandas de confianza para el modelo RLS
  # usando el método de Hotelling-Scheffé

  datos <- datos %>%
    as_tibble()
```

```

# Se extrae la variable independiente (X) de la formula de RLS
variable_x <- formula_rls[[3]]
variable_x <- enquos(variable_x)

# Se calcula el modelo RLS
modelo <- datos %>%
  lm(formula = formula_rls)

# estimador de beta0
b0.hat <- coefficients(modelo)[1]

# estimador de beta1
b1.hat <- coefficients(modelo)[2]

# estimador sigma2
s2.hat <- (summary(modelo)$sigma)^2

# Se calculan las bandas de confianza
resultado <- datos %>%
  mutate(# Se obtienen el número de observaciones
    n = n(),
    # Se obtiene el quantil de la distribución F(2, n - 2),
    # con un nivel de confianza alpha
    fa = qf(alpha, 2, n - 2),
    # Se calcula la estimación de y
    y.hat = b0.hat + b1.hat * !! variable_x,
    # Se calcula la media x
    x.bar = mean(!! variable_x),
    # Se calcula la varianza de x
    S.xx = sum((!! variable_x - mean(!! variable_x))^2),
    # Se calculan los límites de la banda de confianza
    banda_superior = y.hat + sqrt(2*fa*s2.hat)*
      sqrt(1/n + (!! variable_x-x.bar)^2/S.xx),
    banda_inferior = y.hat - sqrt(2*fa*s2.hat)*
      sqrt(1/n + (!! variable_x-x.bar)^2/S.xx)) %>%
  select(-c(n, fa, x.bar, S.xx))

return(resultado)
}

```

A continuación se muestra la tabla con los resultados de  $\hat{y}$  y los límites inferior y superior de la banda de confianza para el modelo RLS.

```

airquality %>%
  as_tibble() %>%
  select(Ozone, Wind) %>%
  filter(!is.na(Ozone), !is.na(Wind)) %>%

```

```
mutate(log2_Wind = log2(Wind)) %>%
bandas_confianza_rls(formula_rls = Ozone ~ log2_Wind)
```

Tabla 1: Estimaciones y bandas de confianza al 95%

Ozone	Wind	$\log_2(\text{Wind})$	$\widehat{\text{Ozone}}$	Banda de confianza	
				Inferior	Superior
41	7.4	2.89	54.25	60.67	47.83
36	8.0	3.00	49.87	55.84	43.91
12	12.6	3.66	24.36	31.58	17.15
18	11.5	3.52	29.49	35.98	23.01
28	14.9	3.90	14.95	23.86	6.03
23	8.6	3.10	45.81	51.52	40.10
19	13.8	3.79	19.26	27.35	11.16
8	20.1	4.33	-1.86	10.66	-14.38
7	6.9	2.79	58.18	65.14	51.23
16	9.7	3.28	39.05	44.74	33.37
11	9.2	3.20	42.03	47.66	36.39
14	10.9	3.45	32.50	38.64	26.36
18	13.2	3.72	21.75	29.40	14.10
14	11.5	3.52	29.49	35.98	23.01
34	12.0	3.58	27.10	33.91	20.30
6	18.4	4.20	3.10	14.51	-8.31
30	11.5	3.52	29.49	35.98	23.01
11	9.7	3.28	39.05	44.74	33.37
1	9.7	3.28	39.05	44.74	33.37
11	16.6	4.05	8.88	19.04	-1.27
4	9.7	3.28	39.05	44.74	33.37
32	12.0	3.58	27.10	33.91	20.30
23	12.0	3.58	27.10	33.91	20.30
45	14.9	3.90	14.95	23.86	6.03
115	5.7	2.51	68.91	77.74	60.07
37	7.4	2.89	54.25	60.67	47.83
29	9.7	3.28	39.05	44.74	33.37
71	13.8	3.79	19.26	27.35	11.16
39	11.5	3.52	29.49	35.98	23.01
23	8.0	3.00	49.87	55.84	43.91
21	14.9	3.90	14.95	23.86	6.03
37	20.7	4.37	-3.51	9.38	-16.41
20	9.2	3.20	42.03	47.66	36.39
12	11.5	3.52	29.49	35.98	23.01
13	10.3	3.36	35.68	41.55	29.82
135	4.1	2.04	87.41	100.23	74.60
49	9.2	3.20	42.03	47.66	36.39
32	9.2	3.20	42.03	47.66	36.39
64	4.6	2.20	80.95	92.31	69.59

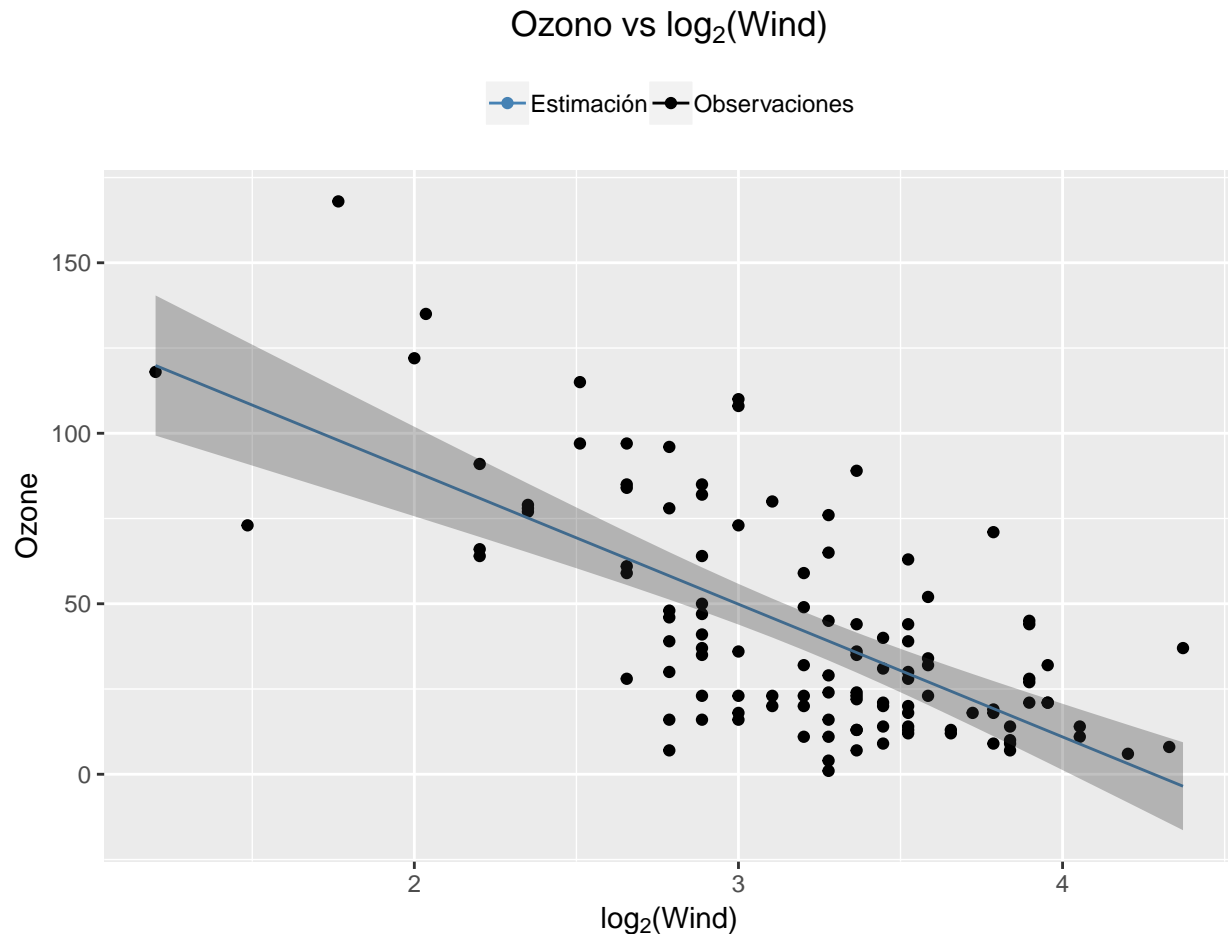


Tabla 1: Estimaciones y bandas de confianza al 95% (*continued*)

Ozone	Wind	$\log_2(\text{Wind})$	$\widehat{\text{Ozone}}$	Banda de confianza	
				Inferior	Superior
40	10.9	3.45	32.50	38.64	26.36
77	5.1	2.35	75.15	85.26	65.05
97	6.3	2.66	63.29	71.08	55.50
97	5.7	2.51	68.91	77.74	60.07
85	7.4	2.89	54.25	60.67	47.83
10	14.3	3.84	17.26	25.72	8.79
27	14.9	3.90	14.95	23.86	6.03
7	14.3	3.84	17.26	25.72	8.79
48	6.9	2.79	58.18	65.14	51.23
35	10.3	3.36	35.68	41.55	29.82
61	6.3	2.66	63.29	71.08	55.50
79	5.1	2.35	75.15	85.26	65.05
63	11.5	3.52	29.49	35.98	23.01
16	6.9	2.79	58.18	65.14	51.23
80	8.6	3.10	45.81	51.52	40.10
108	8.0	3.00	49.87	55.84	43.91
20	8.6	3.10	45.81	51.52	40.10
52	12.0	3.58	27.10	33.91	20.30
82	7.4	2.89	54.25	60.67	47.83
50	7.4	2.89	54.25	60.67	47.83
64	7.4	2.89	54.25	60.67	47.83
59	9.2	3.20	42.03	47.66	36.39
39	6.9	2.79	58.18	65.14	51.23
9	13.8	3.79	19.26	27.35	11.16
16	7.4	2.89	54.25	60.67	47.83
78	6.9	2.79	58.18	65.14	51.23
35	7.4	2.89	54.25	60.67	47.83
66	4.6	2.20	80.95	92.31	69.59
122	4.0	2.00	88.80	101.93	75.67
89	10.3	3.36	35.68	41.55	29.82
110	8.0	3.00	49.87	55.84	43.91
44	11.5	3.52	29.49	35.98	23.01
28	11.5	3.52	29.49	35.98	23.01
65	9.7	3.28	39.05	44.74	33.37
22	10.3	3.36	35.68	41.55	29.82
59	6.3	2.66	63.29	71.08	55.50
23	7.4	2.89	54.25	60.67	47.83
31	10.9	3.45	32.50	38.64	26.36
44	10.3	3.36	35.68	41.55	29.82
21	15.5	3.95	12.73	22.09	3.37
9	14.3	3.84	17.26	25.72	8.79

Tabla 1: Estimaciones y bandas de confianza al 95% (*continued*)

Ozone	Wind	$\log_2(\text{Wind})$	$\widehat{\text{Ozone}}$	Banda de confianza	
				Inferior	Superior
45	9.7	3.28	39.05	44.74	33.37
168	3.4	1.77	97.92	113.18	82.66
73	8.0	3.00	49.87	55.84	43.91
76	9.7	3.28	39.05	44.74	33.37
118	2.3	1.20	119.87	140.42	99.32
84	6.3	2.66	63.29	71.08	55.50
85	6.3	2.66	63.29	71.08	55.50
96	6.9	2.79	58.18	65.14	51.23
78	5.1	2.35	75.15	85.26	65.05
73	2.8	1.49	108.83	126.69	90.96
91	4.6	2.20	80.95	92.31	69.59
47	7.4	2.89	54.25	60.67	47.83
32	15.5	3.95	12.73	22.09	3.37
20	10.9	3.45	32.50	38.64	26.36
23	10.3	3.36	35.68	41.55	29.82
21	10.9	3.45	32.50	38.64	26.36
24	9.7	3.28	39.05	44.74	33.37
44	14.9	3.90	14.95	23.86	6.03
21	15.5	3.95	12.73	22.09	3.37
28	6.3	2.66	63.29	71.08	55.50
9	10.9	3.45	32.50	38.64	26.36
13	11.5	3.52	29.49	35.98	23.01
46	6.9	2.79	58.18	65.14	51.23
18	13.8	3.79	19.26	27.35	11.16
13	10.3	3.36	35.68	41.55	29.82
24	10.3	3.36	35.68	41.55	29.82
16	8.0	3.00	49.87	55.84	43.91
13	12.6	3.66	24.36	31.58	17.15
23	9.2	3.20	42.03	47.66	36.39
36	10.3	3.36	35.68	41.55	29.82
7	10.3	3.36	35.68	41.55	29.82
14	16.6	4.05	8.88	19.04	-1.27
30	6.9	2.79	58.18	65.14	51.23
14	14.3	3.84	17.26	25.72	8.79
18	8.0	3.00	49.87	55.84	43.91
20	11.5	3.52	29.49	35.98	23.01



## Ejercicio 5

(Sheater) Un estadístico colaboró en un proyecto de investigación con dos entomólogos. El análisis involucró el ajuste de modelos de regresión a grandes conjuntos de datos. Entre los tres escribieron y sometieron un manuscrito a una revista de entomología. El escrito contenía varias gráficas de dispersión mostrando la recta de regresión ajustada y las bandas de confianza 95% para la verdadera recta de regresión calculadas con los IC individuales, así como los datos observados. Uno de los revisores del manuscrito hizo la siguiente observación:

*No puedo entender cómo el 95% de las observaciones cae fuera de las bandas de confianza 95% que se muestran en las figuras*

Dar una respuesta breve del porqué es posible que el 95% de las observaciones caigan fuera de las bandas de confianza 95% que se muestran en las figuras.

Esto sucede por el carácter de las bandas, estas nos indican la posibilidad de que la verdadera línea de regresión esté dentro con un 95% de confianza, es decir, que la verdadera recta de regresión debe caer un 95% de las veces dentro de esta banda. Por lo que interpretar las observaciones con respecto a esta banda carece de sentido y no necesariamente estos valores corresponderán a estas bandas.

## Ejercicio 6

(Ross) Suponer que se tiene el siguiente conjunto de datos donde  $x$  representa la humedad de una mezcla fresca de un determinado producto y  $y$  la densidad del producto terminado.

$x$	5	6	7	10	12	15	18	20
$y$	7.4	9.3	10.6	15.4	18.1	22.2	24.1	24.8

Ajustar un modelo RLS a los datos anteriores y responder lo siguiente.

### Inciso 6.a

Reportar la estimación puntual de  $\sigma^2$  e interpretar el resultado en cuanto a la utilidad del modelo RLS ajustado.

### Inciso 6.b

Reportar el IC 90% para  $\sigma^2$  con los cuantiles simétricos y su longitud.

### Inciso 6.c

Indicar cuáles son los cuantiles que proporcionan el IC 90% para  $\sigma^2$  de menor longitud.

### Inciso 6.d

Reportar el IC 90% para  $\sigma^2$  de menor longitud y compararlo con el intervalo del inciso a).