

# Estadística Computacional

**Alumnos:**

*Córdova Proleón, Christian Therius*

**Cod:** 20173970

*Fazio Luna, Boris Manuel*

**Cod:** 20173896

*Huertas Quispe, Anthony Enrique*

**Cod:** 20173728

*Saenz Egusquiza, Miguel Angel*

**Cod:** 19921255

**Semestre:** 2017-I

**Tema:** Lista 2

PROF. CRISTIAN BAYES



Pontificia Universidad Católica del Perú  
Escuela de Posgrado  
Maestría en Estadística

## Ejercicio 1

Un estimador utilizado para reducir el efecto que puedan tener valores atípicos es la media winsorizada (winsorized mean) que es la media que se obtiene descartando el  $a\%$  de valores extremos superiores siendo reemplazados por el cuantil  $(1 - a)\%$  y el  $a\%$  de valores inferiores de los datos siendo reemplazados por el cuantil  $a\%$ .

Realice un estudio de simulación comparando la media winsorizada al 20% cada lado con la media muestral. Considere 3 escenarios:

- (a) datos simétricos sin valores atípicos: distribución normal.
- (b) datos simétricos con valores atípicos: distribución t de Student con 4 grados de libertad.
- (c) datos asimétricos: distribución gama.

En cada escenario determine los parámetros de la distribución de modo que la media y la varianza sean iguales a 1. Considere los siguientes tamaños de muestra:  $n = 30, 50$  y  $100$ . Interprete sus resultados.

**Sugerencia:** En R la función `winsor.mean` de la librería `psych` tiene el parámetro `trim` que permite calcular la media winsorizada.

Recuerde que si  $X \sim t(\mu, \sigma^2, \nu)$  su función de densidad es dada por:

$$f_X(x) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\Gamma\left(\frac{\nu}{2}\right)} \frac{1}{\sqrt{\nu\pi}\sigma} \left(1 + \frac{1}{\nu} \left(\frac{x-\mu}{\sigma}\right)^2\right)^{-\frac{\nu+1}{2}}, \quad x \in \mathbb{R}$$

donde  $\mu$  es un parámetro de localización,  $\sigma$  un parámetro de escala y  $\nu$  es un parámetro de forma que controla la curtosis de la distribución,

$$\begin{aligned} E(X) &= \mu \text{ si } \nu > 1, \\ \text{Var}(X) &= \frac{\nu}{\nu-2} \sigma^2 \text{ si } \nu > 2 \end{aligned}$$

*Solución.* El problema plantea verificar la eficiencia del estimador “media winsorizada” para reducir el efecto de los valores atípicos sobre la estimación de parámetros. En tal sentido, se medirá el desempeño de los estimadores de la “media muestral o aritmética” y “media winsorizada”, mediante el sesgo y el error cuadrático medio. Comparándolos a su vez, frente a los siguientes factores:

■ **Escenarios:**

- Datos simétricos sin valores *outlier*: Distribución Normal,  $X \sim N(\mu, \sigma^2)$ , donde

$$E[X] = \mu = 1 \text{ y } \text{Var}(X) = \sigma^2 = 1 \Rightarrow \mu = 1, \sigma = 1$$

Función de densidad,  $X \sim N(1, 1)$ :

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-1)^2}{2}}, \quad x \in \mathbb{R}$$

- Datos simétrico con valores *outlier*: Distribución t de Student con  $\nu$  grados de libertad,  $X \sim t(\mu, \sigma^2, \nu)$ , donde

$$\nu = 4, \quad E[X] = \mu = 1 \text{ y } \text{Var}(X) = \frac{4}{4-2} \sigma^2 = 1 \Rightarrow \begin{cases} \mu = 1 \\ \sigma = \frac{1}{\sqrt{2}} \\ \nu = 4 \end{cases}$$

Función de densidad,  $X \sim t(1, \frac{1}{2}, 4)$ :

$$f(x) = \frac{\Gamma(5/2)}{\Gamma(2)(2\sqrt{\pi})} \left(1 + \frac{1}{2}(x-1)^2\right)^{-5/2}, \quad x \in \mathbb{R}$$

- Datos Asimétricos: Distribución Gamma,  $X \sim \text{Gamma}(\alpha, \lambda)$ , donde

$$E[X] = \frac{\alpha}{\lambda} = 1 \text{ y } \text{Var}(X) = \frac{\alpha}{\lambda^2} = 1 \Rightarrow \alpha = 1, \lambda = 1$$

Función de densidad,  $X \sim \text{Gamma}(\alpha, \lambda)$  :

$$f(x) = e^{-x}, \quad x > 0.$$

- **Tamaños de muestra:**  $n = 30, 50$  y  $100$ .

Listing 1: Definiciones iniciales.

---

```

1 library(psych)
2
3 # Numero de simulaciones
4 M = 10000
5
6 # Tamanos de muestra
7 n = c(30,50,100)
8
9 # Parametros
10 mu = 1; sigma = c(1,1/sqrt(2))
11 nu=4
12 alpha=1; lambda=1
13
14 # Matriz de resultados
15 T = matrix(0,M,2*3*3)
16 colnames(T) = c("MedArit.30.N","MedWinsor.30.N","MedArit.30.t","MedWinsor.30.t","
    MedArit.30.G","MedWinsor.30.G","MedArit.50.N","MedWinsor.50.N","MedArit.50.t","
    MedWinsor.50.t","MedArit.50.G","MedWinsor.50.G","MedArit.100.N","MedWinsor.100.N","
    MedArit.100.t","MedWinsor.100.t","MedArit.100.G","MedWinsor.100.G")

```

---

Listing 2: Estudio de simulación.

---

```

1 set.seed(100)
2 for(h in 1:3){
3     for(j in 1:M){
4         x = rnorm(n[h],mu,sigma[1])
5         y = mu+sigma[2]*rt(n[h],4)
6         z = rgamma(n[h],alpha,lambda)
7         T[j,(6*h-5):(6*h)] = c(mean(x),winsor.mean(x,trim=0.2),mean(y),winsor.
            mean(y,trim=0.2),mean(z),winsor.mean(z,trim=0.2))
8     }
9 }

```

---

Listing 3: Análisis de resultados.

---

```

1 est = rep(c("MedArit","MedWinsor"),times=9)
2 dis = rep(c("N","N","t","t","G","G"),times=3)
3 tam = rep(n,each=4)
4
5 sesgo = colMeans(T)-mu
6 variancia = diag(var(T))
7 ecm = variancia+sesgo^2

```

---

Listing 4: Tabla Sesgo.

```

1 tab_ses = rbind(matrix(sesgo[dis=="N"], 3, 2, byrow=T), matrix(sesgo[dis=="t"], 3, 2, byrow=T)
  ), matrix(sesgo[dis=="G"], 3, 2, byrow=T) )
2 tab_ses = cbind(rep(n, 3), tab_ses)
3 colnames(tab_ses) = c("T. Muestra", "MedArit", "MedWinsor")
4 rownames(tab_ses) = rep(c("N", "t", "G"), each=3)
5 round(tab_ses, 4)

```

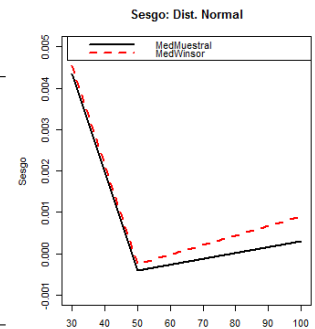
Sesgo				
Escenario	Distribución	n	Media muestral	Media winsorizada
Datos simétricos sin valores <i>outlier</i>	Normal	30	0.0043	0.0045
		50	-0.0004	-0.0003
		100	0.0003	0.0009
Datos simétricos con valores <i>outlier</i>	t-student	30	-0.0003	-0.0005
		50	0.0004	0.0009
		100	0.0002	0.0001
Datos asimétricos	Gamma	30	-0.0014	-0.1672
		50	0.0028	-0.1685
		100	-0.0003	-0.1745

Listing 5: Gráfica Sesgo (Normal).

```

1 par(mfrow=c(1, 3))
2 #Normal
3 plot(n, sesgo[(dis=="N") & (est=="MedArit")], ylim=c(-0.001, 0.005), type=
  "l", col=1, lty=1, lwd=2,
4 main="Sesgo: Dist. Normal", xlab="Tamano de Muestra", ylab="Sesgo")
5 lines(n, sesgo[(dis=="N") & (est=="MedWinsor")], col=2, lty=2, lwd=2)
6 legend("topright", c("MedMuestral", "MedWinsor"), col=c(1, 2), lty=c(1, 2)
  , lwd=2)

```

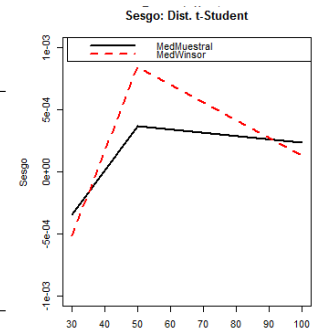


Listing 6: Gráfica Sesgo (t-Student).

```

1 #t-Student
2 plot(n, sesgo[(dis=="t") & (est=="MedArit")], ylim=c(-0.001, 0.001), type=
  "l", col=1, lty=1, lwd=2,
3 main="Sesgo: Dist. t-Student", xlab="Tamano de Muestra", ylab="Sesgo")
4 lines(n, sesgo[(dis=="t") & (est=="MedWinsor")], col=2, lty=2, lwd=2)
5 legend("topright", c("MedMuestral", "MedWinsor"), col=c(1, 2), lty=c(1, 2)
  , lwd=2)

```

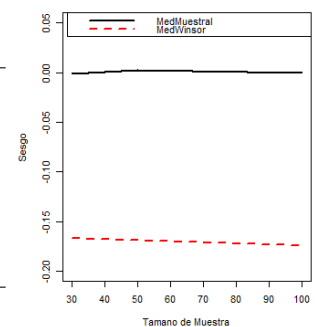


Listing 7: Gráfica Sesgo (Gamma).

```

1 #Gamma
2 plot(n, sesgo[(dis=="G") & (est=="MedArit")], ylim=c(-0.2, 0.05), type="l",
  col=1, lty=1, lwd=2,
3 main="Sesgo: Dist. Gamma", xlab="Tamano de Muestra", ylab="Sesgo")
4 lines(n, sesgo[(dis=="G") & (est=="MedWinsor")], col=2, lty=2, lwd=2)
5 legend("topright", c("MedMuestral", "MedWinsor"), col=c(1, 2), lty=c(1, 2)
  , lwd=2, inset=0)

```



Listing 8: Tabla ECM.

```

1 tab_ecm = rbind(matrix(ecm[dis=="N"],3,2,byrow=T),matrix(ecm[dis=="t"],3,2,byrow=T),
  matrix(ecm[dis=="G"],3,2,byrow=T))
2 tab_ecm = cbind(rep(n,3),tab_ecm)
3 colnames(tab_ecm) = c("T. Muestra","MedArit","MedWinsor")
4 rownames(tab_ecm) = rep(c("N","t","G"),each=3)
5 round(tab_ecm,4)

```

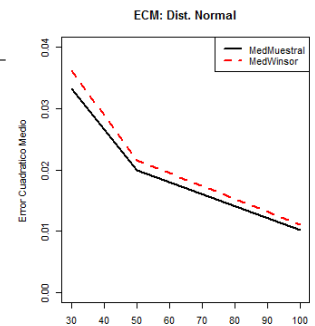
ECM				
Escenario	Distribución	n	Media muestral	Media winsorizada
Datos simétricos sin valores <i>outlier</i>	Normal	30	0.0332	0.0362
		50	0.0199	0.0216
		100	0.0102	0.0111
Datos simétricos con valores <i>outlier</i>	t-student	30	0.0331	0.0251
		50	0.0198	0.0148
		100	0.0098	0.0074
Datos asimétricos	Gamma	30	0.0331	0.0564
		50	0.0199	0.0455
		100	0.0103	0.0392

Listing 9: Gráfica Sesgo (Normal).

```

1 par(mfrow=c(1,3))
2 #Normal
3 plot(n,ecm[(dis=="N")&(est=="MedArit")],ylim=c(0,0.05),type="l",col
  =1,lty=1,lwd=2,
4 main="ECM: Dist. Normal",xlab="Tamano de Muestra",ylab="Error
  Cuadratico Medio")
5 lines(n,ecm[(dis=="N")&(est=="MedWinsor")],col=2,lty=2,lwd=2)
6 legend("topright",c("MedMuestral","MedWinsor"),col=c(1,2),lty=c(1,2)
  ,lwd=2)

```

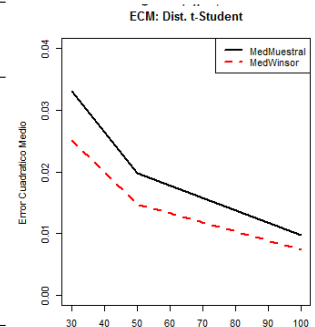


Listing 10: Gráfica Sesgo (t-Student).

```

1 #t-Student
2 plot(n,ecm[(dis=="t")&(est=="MedArit")],ylim=c(0,0.1),type="l",col
  =1,lty=1,lwd=2,
3 main="ECM: Dist. t-Student",xlab="Tamano de Muestra",ylab="Error
  Cuadratico Medio")
4 lines(n,ecm[(dis=="t")&(est=="MedWinsor")],col=2,lty=2,lwd=2)
5 legend("topright",c("MedMuestral","MedWinsor"),col=c(1,2),lty=c(1,2)
  ,lwd=2)

```

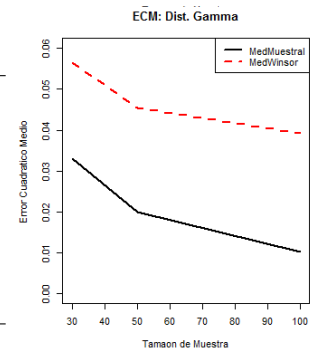


Listing 11: Gráfica Sesgo (Gamma).

```

1 #Gamma
2 plot(n,ecm[(dis=="G")&(est=="MedArit")],ylim=c(0,0.1),type="l",col
  =1,lty=1,lwd=2,
3 main="ECM: Dist. Gamma",xlab="Tamaon de Muestra",ylab="Error
  Cuadratico Medio")
4 lines(n,ecm[(dis=="G")&(est=="MedWinsor")],col=2,lty=2,lwd=2)
5 legend("topright",c("MedMuestral","MedWinsor"),col=c(1,2),lty=c(1,2)
  ,lwd=2)

```



### Interpretación:

En relación al **sesgo**, para el caso de la media muestral, notamos que este disminuye a medida que aumenta el tamaño muestral en cualquiera de los tres escenarios.

En el primer escenario para datos simétricos sin valores atípicos, se observa que los sesgos son pequeños y en la mayoría de los casos la media winzorizada tiene un menor sesgo.

En el segundo escenario para datos simétricos con valores atípicos, se observa que la media muestral tiene un menor sesgo para tamaños de muestra  $n=30$  y  $50$ , a diferencia del tamaño de muestra  $n = 100$ .

En el tercer escenario, para datos asimétricos, se observa que los sesgos para la media muestral es pequeño con diferencia significativa para el caso de la media winzorizada que parece mantenerse constante a lo lejos del sesgo de la media muestral.

En el caso del **error cuadrático medio**, notamos que a medida que el tamaño de la muestra aumenta, el ECM disminuye en los tres escenarios vistos.

En primer escenario, los datos simétricos sin valores atípicos, observamos que el ECM es menor siempre en el caso de la media muestral. Lo que confirma que la media muestral es el estimador con menos variancia en el caso de datos normales.

En el segundo escenario, datos simétricos con valores atípicos, vemos que la media winzorizada es la que presenta menor valor, y eso es debido a que elimina el efecto de los atípicos, lo que la media muestral no puede realizar.

En el tercer escenario, los datos asimétricos, vemos que el ECM es menor en la media muestral.

□

## Ejercicio 2

*Estudio de simulación sobre la influencia de valores atípicos en la estimación del modelo de regresión lineal.*

Sea un modelo de regresión lineal simple

$$\begin{aligned} y_i &\overset{\text{ind}}{\sim} N(\mu_i, \sigma^2) \\ \mu_i &= \beta_0 + \beta_1 x_i \end{aligned}$$

donde  $\beta_0, \beta_1 \in \mathbb{R}$ ,  $\sigma^2 > 0$  y  $x_i, i = 1, \dots, n$  son consideradas constantes conocidas. Para el estudio de simulación considere  $\beta_0 = 1$ ,  $\beta_1 = 2$ ,  $\sigma = 0.2$ ,  $n = 50$  y los valores de  $x_i$  deben ser simulados solamente una vez de una distribución  $Uniforme(0, 1)$ . Luego de simular un conjunto de datos se introducirán dos valores outlier mediante el siguiente procedimiento

- Reemplazar los dos primeros valores simulados por

$$y_i^* = y_i + 10 \times \sigma, \quad i = 1, 2.$$

Estudie mediante simulación el desempeño del usual estimador de mínimos cuadrados para  $\beta_0$  y  $\beta_1$  en el escenario planteado. Compare sus resultados con los que se obtendrían teóricamente en el caso que no existieran valores atípicos. Interprete sus resultados.

*Solución.* Sabemos que el estimador de Monte Carlo del Error Cuadrático Medio de  $\beta_0$  y  $\beta_1$  es

$$\widehat{ECM}(T_{ij}) = \frac{1}{M} \sum_{j=1}^M (T_{ij} - \beta_i)^2 = \widehat{V}(T_{ij}) - \widehat{sesgo}(T_{ij})^2, \quad i = 1, 2$$

donde  $T_{ij}$  son estimadores de prueba de  $\beta_i$ , además

$$\begin{aligned} \widehat{V}(T_{ij}) &= \frac{1}{M} \sum_{j=1}^M (T_{ij} - \bar{T}_{ij})^2 \\ \widehat{sesgo}(T_{ij}) &= \frac{1}{M} \sum_{j=1}^M T_{ij} - \beta_i \end{aligned}$$

Primero definimos ciertas condiciones que serán implementadas en el estudio de la simulación.

Listing 12: Definiciones iniciales.

---

```

1 # Numero de simulaciones
2 M=1000
3
4 # Tamano de muestra
5 n=50
6
7 # Parametros
8 beta0=1
9 beta1=2
10 sigma=0.2
11
12 # Matriz de resultados
13 T<-matrix(0,M,2*2)
14 colnames(T)<-c("beta0.Sout","beta1.Sout","beta0.out","beta1.out")

```

---

Se procederá a realizar la simulación de datos mediante las siguientes condiciones:

■ **Escenarios:**

- (a) sin valores *outlier*:  $y_i \stackrel{ind}{\sim} N(\beta_0 + \beta_1 x_i, \sigma^2)$  donde  $x_i \sim U(0, 1)$ .
- (b) con valores *outlier*: Generamos los siguientes valores atípicos:  $y_i^* = y_i + 10 \times \sigma$ ,  $i = 1, 2$ .

■ **Tamaño de muestra:**  $n = 50$ .

Listing 13: Estudio de simulación.

---

```

1 # Generando valores x (una simulacion)
2 x=runif(n)
3
4 # Generando valores y (1000 simulaciones)
5 for(j in 1:M){
6     mu=beta0+beta1*x
7     y<-rnorm(n)*sigma+mu
8
9     # Sin valores outlier
10    regre=lm(y~x)
11
12    # Con valores outlier
13    y[1:2]=y[1:2]+10*sigma
14    regre2=lm(y~x)
15    T[j,]=c(regre$coefficients[1], regre$coefficients[2],
16            regre2$coefficients[1], regre2$coefficients[2])
17 }
```

---

Finalmente se harán las evaluaciones correspondientes del desempeño del Error Cuadrático Medio para  $\beta_0$  y  $\beta_1$ .

Listing 14: Análisis de resultados.

---

```

1 sesgo=numeric(4)
2 for(i in 1:2){
3     sesgo[2*i-1]=colMeans(T)[2*i-1]-beta0
4     sesgo[2*i]=colMeans(T)[2*i]-beta1
5 }
6 variancia=diag(var(T))
7
8 # Estimador de Monte Carlo del ECM
9 ecm=variancia+sesgo**2
10 ecm
```

---

Resultado:

Error Medio Cuadrático		
Escenario	$\beta_0$	$\beta_1$
sin valores <i>outlier</i>	0.003694595	0.010478327
con valores <i>outlier</i>	0.060819815	0.098765976

**Interpretación:**

El análisis hecho es respecto al ECM de los coeficientes de regresión de la ecuación de regresión  $y = \beta_0 + \beta_1 x$ . Notamos que tanto para el coeficiente  $\beta_0$  como para  $\beta_1$ , el ECM es menor cuando no existen valores outlier que cuando si existen comprobando que el efecto de valores outlier ocasiona un ECM mayor en ambos coeficientes.

Concluimos un mejor rendimiento en el escenario “sin valores *outlier*”, lo cual es teóricamente correcto.  $\square$



### Ejercicio 3

Una variable aleatoria  $X$ , definida en toda la recta, tiene distribución mixtura de dos normales, si su función de densidad es dada por la siguiente expresión:

$$f(x) = p\phi(x|\mu_1, \sigma_1^2) + (1-p)\phi(x|\mu_2, \sigma_2^2)$$

donde  $\mu_1, \mu_2 \in \mathbb{R}$ ,  $\sigma_1^2, \sigma_2^2 > 0$ ,  $p \in (0, 1)$  y  $\phi(\cdot|a, b^2)$  es la densidad de una distribución normal con media  $a$  y varianza  $b^2$ . Considere los datos de tiempo de espera entre erupciones del géiser Old Faithful en el Yellowstone National Park, que se encuentran en la variable `waiting` del conjunto de datos `faithful` del R, para las siguientes preguntas

- a) Estime por máxima verosimilitud los parámetros de esta distribución.

*Solución.* Determinamos primero su función de verosimilitud

$$\begin{aligned} \mathcal{L}(p, \mu_1, \sigma_1, \mu_2, \sigma_2) &= \prod_{i=1}^n f(y_i | \mu_1, \sigma_1, \mu_2, \sigma_2) \\ &= \prod_{i=1}^n \left( p \frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{1}{2} \frac{(y_i - \mu_1)^2}{\sigma_1^2}} + (1-p) \frac{1}{\sqrt{2\pi}\sigma_2} e^{-\frac{1}{2} \frac{(y_i - \mu_2)^2}{\sigma_2^2}} \right) \\ &= \prod_{i=1}^n \left( \frac{p\sigma_2 e^{-\frac{1}{2} \frac{(y_i - \mu_1)^2}{\sigma_1^2}} + (1-p)\sigma_1 e^{-\frac{1}{2} \frac{(y_i - \mu_2)^2}{\sigma_2^2}}}{\sqrt{2\pi}\sigma_1\sigma_2} \right) \end{aligned}$$

La función log-verosimilitud vendría dada por:

$$\ell(p, \mu_1, \sigma_1, \mu_2, \sigma_2) = \sum_{i=1}^n \log \left( p\sigma_2 e^{-\frac{1}{2} \frac{(y_i - \mu_1)^2}{\sigma_1^2}} + (1-p)\sigma_1 e^{-\frac{1}{2} \frac{(y_i - \mu_2)^2}{\sigma_2^2}} \right) - \log(\sqrt{2\pi}\sigma_1\sigma_2)$$

Listing 15: EMV.

---

```

1 y=faithful$waiting
2 p=runif(1)
3 u1=54
4 u2=80
5 sigma1=5.61
6 sigma2=6
7 parametros=c(p,u1,u2,sigma1,sigma2)
8 f1=function(parametros,y){
9     a1=dnorm(y,mean = parametros[2],sd=parametros[4])
10    a2=dnorm(y,mean = parametros[3],sd=parametros[5])
11    -sum(log(parametros[1]*a1+(1-parametros[1])*a2))
12 }
13 r=nlmminb(parametros,f1,y=faithful$waiting,lower = 0.1,upper =100)
14 g=r$par
15 g

```

---

Parámetros estimados

p	$\mu_1$	$\mu_2$	$\sigma_1$	$\sigma_2$
0.3608861	54.6148538	80.0910705	5.8712177	5.8677341



- b) Calcule el AIC para este modelo y compárelo con el obtenido por un modelo normal.

Listing 16: AIC (Mixtura de dos normales).

```

1 mu=mean(y)
2 sigma=sd(y)
3 parametros2=c(mu,sigma)
4 f2=function(parametros2,y){
5     a1=dnorm(y,mean = mu,sd=sigma)
6     -sum(log(a1))
7 }
8 r2=nlminb(parametros2,f2,y=faithful$waiting,lower = 43,upper =100)
9 g2=r2$par
10 AIC2 = 2*f2(g,y)+4
11 AIC2

```

Listing 17: AIC (Normal).

```

1 AIC = 2*f1(g,y)+10
2 AIC

```

Distribución	AIC
Mixtura de dos Normales	2078.003
Normal	2194.579

Dado que el AIC de la v.a con dos normales es menor que el AIC de una normal, entonces el primero es un mejor modelo.

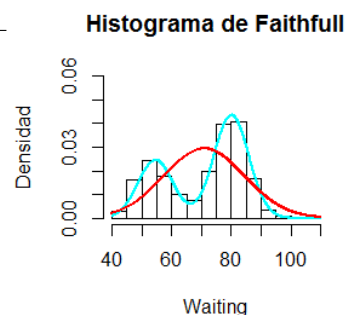
- c) Realice un histograma y compare ambos modelos. Comente sus resultados.

Listing 18: AIC (Normal).

```

1 S = seq(from = 40, to = 110, by = 0.1)
2 d1 = dnorm(S, mean = g[2], sd = g[4])
3 d2 = dnorm(S, mean = g[3], sd = g[5])
4 f = g[1] * d1 + (1 - g[1]) * d2
5 hist(y, probability = TRUE, xlim = range(S), ylim = c
6     (0, 0.06), xlab = "Waiting", ylab = "Densidad",
7     main = "Histograma de Faithfull")
8 lines(S,f, lwd = 2,col=5)
9 lines(S, dnorm(S, mean = mean(y), sd = sd(y)),lwd = 2,
10     col=2)

```



Según el gráfico se ratifica lo que indicaba el AIC, que la función que es una mixtura de dos normales ajusta mejor los datos y es un mejor modelo que la función normal para explicar el comportamiento de este tipo de variables que tienen dicho comportamiento, como por ejemplo la variable waiting que mide el tiempo de espera entre erupciones del geiser Old Faithful en el Yellowstone National Park.

La función normal podría quizás explicar en cierta medida el comportamiento de las observaciones que se encuentra en ambas colas (derecha e izquierda), pero no en la parte central, donde el error de estimación sería grande.

## Ejercicio 4

*Modelo de regresión lineal t de Student.*

El modelo de regresión lineal t de Student es dado por

$$\begin{aligned} y_i &\stackrel{ind}{\sim} t(\mu_i, \sigma^2, \nu) \\ \mu_i &= \beta_0 + \beta_1 x_i \end{aligned}$$

donde  $\beta_0, \beta_1 \in \mathbb{R}$ ,  $\sigma^2 > 0$  y  $x_i, i = 1, \dots, n$  son consideradas constantes conocidas.

- a) Estime por máxima verosimilitud el modelo de regresión lineal t de Student, considerando que el parámetro de grados de libertad es fijo e igual a 3 ( $\nu = 3$ ), para las variables  $x_3$  y  $y_3$  del conjunto de datos de `anscombe` de R. Compare sus resultados con el usual modelo de regresión lineal en un gráfico de dispersión. Interprete sus resultados.

*Solución.* Para una variable  $Y \sim t(\mu, \sigma^2, \nu)$ , su función de densidad está definida como:

$$f(Y = y \mid \mu, \sigma^2, \nu) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\Gamma\left(\frac{\nu}{2}\right) \sqrt{\pi\nu\sigma^2}} \left(1 + \frac{(y - \mu)^2}{\nu\sigma^2}\right)^{-\frac{\nu+1}{2}}$$

Entonces la función de verosimilitud, siendo  $\nu$  fijo, para este caso viene dada por:

$$\begin{aligned} \mathcal{L}(\beta_0, \beta_1, \sigma^2) &= f(x_1, \dots, x_n, y_1, \dots, y_n \mid \beta_0, \beta_1, \sigma^2) \\ &= \prod_{i=1}^n \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\Gamma\left(\frac{\nu}{2}\right) \sqrt{\pi\nu\sigma^2}} \left(1 + \frac{(y_i - \beta_0 - \beta_1 x_i)^2}{\nu\sigma^2}\right)^{-\frac{\nu+1}{2}} \end{aligned}$$

Para obtener la estimación de máxima verosimilitud de cada parámetro basta maximizar la función log-verosimilitud:

$$\ell(\beta_0, \beta_1, \sigma) = -\left(\frac{\nu+1}{2}\right)^n \sum_{i=1}^n \left(1 + \frac{(y_i - \beta_0 - \beta_1 x_i)^2}{\nu\sigma^2}\right) - n \log(\sigma) + cte$$

Procedemos a maximizar  $\ell(\beta_0, \beta_1, \sigma)$  usando *optimización directa*:

Listing 19: EMV (Método de Optimización directa).

---

```

1 # Tomando datos x3, y3 de anscombe
2 x=anscombe$x3; y=anscombe$y3
3 nu=3;n=length(x)
4
5 # Funcion logVerosimilitud
6 llike =function(params){
7     b0=params[1]; b1=params[2]; sigma=params[3]
8     ll=n*log(sigma)+0.5*(nu+1)*sum(log(1+((y-b0-b1*x)**2)/(nu*sigma**2)))
9 }
10
11 # Tomamos estimadores iniciales de un modelo regresion lineal
12 mln=lm(y~x)
13 inicial=c(mln$coefficients[1],mln$coefficients[2],sd(mln$residuals))
14
15 # Optimizacion directa
16 mlt=nlminb(inicial,llike)
17
18 # Comparacion
19 inicial; mlt$par

```

---

Resultado:

Estimadores			
Modelo	$\beta_0$	$\beta_1$	$\sigma$
Regresión Lineal	3.0024545	0.4997273	1.1728679
t-student con 3 grados de libertad	4.006308050	0.345313169	0.003071092

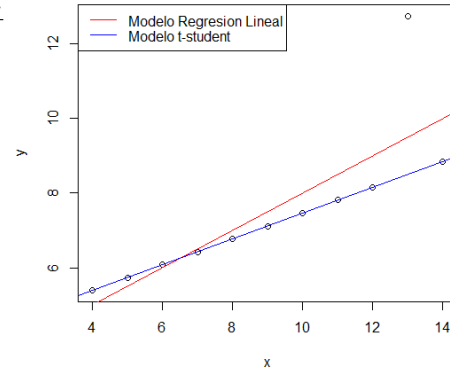
Mostramos el gráfico de dispersión de los datos y comparamos los modelos resultantes:

Listing 20: Gráfico de dispersión y ajuste por modelos.

```

1 # Grafica de dispersion
2 plot(x,y)
3
4 # Modelo de Regresion Lineal
5 abline(inicial[1],inicial[2],col="red")
6
7 # Modelo t-student con 3 grados de libertad
8 abline(mlt$par[1],mlt$par[2],col="blue")
9
10 legend("topleft",legend=c("Modelo Regresion
    Lineal","Modelo t-student"),lty=1,col=c("
    red","blue"))

```



La pendiente estimada bajo el supuesto de normalidad, mediante el modelo de regresión lineal, es influenciada por el valor atípico. Por su parte, el modelo t-student logra capturar perfectamente la tendencia lineal mostrada por la mayoría de observaciones.  $\square$

- b) Repite el estudio de simulación de la pregunta 2, estimando ahora los coeficientes de regresión utilizando el modelo de regresión *t*-Student y compare sus resultados con los obtenidos en la pregunta 2.

*Solución.* Nos situaremos en el Escenario con valores *outlier*, en donde se generarán valores  $y_i \stackrel{\text{ind}}{\sim} N(\beta_0 + \beta_1 x_i, \sigma^2)$  donde  $x_i \sim U(0, 1)$  y los valores atípicos  $y_i^* = y_i + 10 \times \sigma$ ,  $i = 1, 2$ . Simularemos  $\beta_0, \beta_1$ , con el procedimiento siguiente.

- Usando el método de regresión lineal. (Desarrollada en el problema 2).
- Usando el método de regresión *t*-student con 3 grados de libertad.

Listing 21: Definiciones Adicionales.

```

1 # Numero de simulaciones
2 M=1000
3
4 # Tamano de muestra
5 n=50
6
7 # Parametros
8 beta0=1; beta1=2; sigma=0.2
9
10 # Matriz de resultados
11 T<-matrix(0,M,2)
12 colnames(T)<-c("beta0.t","beta1.t")

```

Listing 22: Estudio de simulación.

---

```

1 # Generando valores x (una simulacion)
2 x=runif(n)
3
4 # Generando valores y (1000 simulaciones)
5 for(j in 1:M){
6     mu=beta0+beta1*x
7     y<-rnorm(n)*sigma+mu
8     y[1:2]=y[1:2]+10*sigma
9
10    # EMV de beta0 y beta1 (Optimizacion directa para t-student)
11    mln=lm(y~x)
12    inicial=c(mln$coefficients[1],mln$coefficients[2])
13    mlt=nlminb(inicial,llike)
14
15    T[j,]=c(mlt$par[1],mlt$par[2])
16 }

```

---

Listing 23: Análisis de resultados.

---

```

1 sesgo=numeric(2)
2 sesgo[1]=colMeans(T)[1]-beta0
3 sesgo[2]=colMeans(T)[2]-beta1
4 variancia=diag(var(T))
5
6 # Estimador de Monte Carlo del ECM
7 ecm=variancia+sesgo**2
8 ecm

```

---

Resultado: Tabla comparativa de ECM, en Escenario con valores *outlier* respecto a los métodos de regresión lineal y regresión *t*-student.

Error Medio Cuadrático		
Modelos	$\beta_0$	$\beta_1$
Modelo de regresión lineal	0.060819815	0.098765976
Modelo <i>t</i> -student	0.01668002	0.01541430

### Interpretación:

Nos situamos en un escenario con valores outlier, y observamos que el ECM obtenido tanto para  $\beta_0$  como para  $\beta_1$  mediante el modelo de regresión *t*-student es menor en ambos casos con respecto a lo obtenido usando el modelo de regresión lineal vista en el Ejercicio 2. Esto nos indica un mejor rendimiento por parte del Modelo de regresión *t*-student.

□