

# Hoja 7 de problemas y prácticas con R

Estadística Computacional I. Grado en Estadística

Marta Venegas Pardo

## Contents

<b>1 Ejercicio 1</b>	<b>1</b>
1.1 Parte 1 (Contraste Bootstrap unilateral)	2
1.2 Paso 2 (bilateral)	3
1.3 Parte 3 (Con la librería boot)	4
<b>2 Ejercicio 2</b>	<b>6</b>
2.1 ACP	6
2.2 Histograma y p-valor	7
2.3 Con librería BOOT	7
<b>3 Ejercicio 3 IC bootstrap para un cociente</b>	<b>8</b>
<b>4 Ejercicio 4</b>	<b>12</b>
<b>5 Ejercicio 5</b>	<b>12</b>
<b>6 Ejercicio 6</b>	<b>13</b>
<b>7 Ejercicio 7 Estimar el error de clasificación</b>	<b>14</b>
<b>8 Ejercicio 8</b>	<b>19</b>
8.1 Parte 1 Leer datos	19
8.2 Parte 2 (Modelo de regresión lineal múltiple)	19
8.3 Parte 3 (Estimaciones Jackknife y bootstrap)	19
8.4 Parte 4 (generamos las muestras bootstrap)	20

## 1 Ejercicio 1

1. Realizar un contraste bootstrap unilateral de hipótesis para comparar las desviaciones típicas a partir de las siguientes muestras:

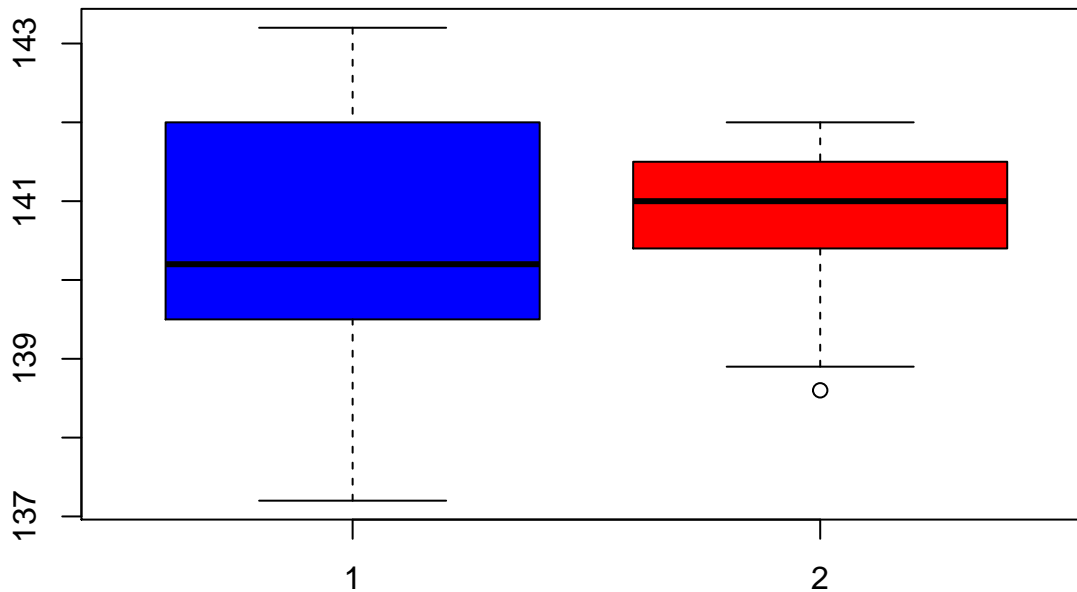
- $x=c(137.9, 143, 143.2, 140, 140.2, 139.3, 141.4, 140.1, 142, 137.2, 139.5, 142.7, 141.3)$
- $y=c(141.6, 138.9, 140, 141.9, 140.5, 138.6, 141.5, 141.5, 140.7, 141.5, 140.4, 142, 141)$

$$H_0 : \sigma_1 \leq \sigma_2$$

Escribir las instrucciones R sin y con la librería boot.

## 1.1 Parte 1 (Contraste Bootstrap unilateral)

```
x<-c(137.9, 143, 143.2, 140, 140.2, 139.3, 141.4, 140.1, 142, 137.2, 139.5, 142.7, 141.3)
y<-c(141.6, 138.9, 140, 141.9, 140.5, 138.6, 141.5, 141.5, 140.7, 141.5, 140.4, 142, 141)
nx<-length(x)
ny<-length(y)
boxplot(x,y,col=c("blue","red"))
```



```
xy<-c(x,y)
nxy<-length(xy)
```

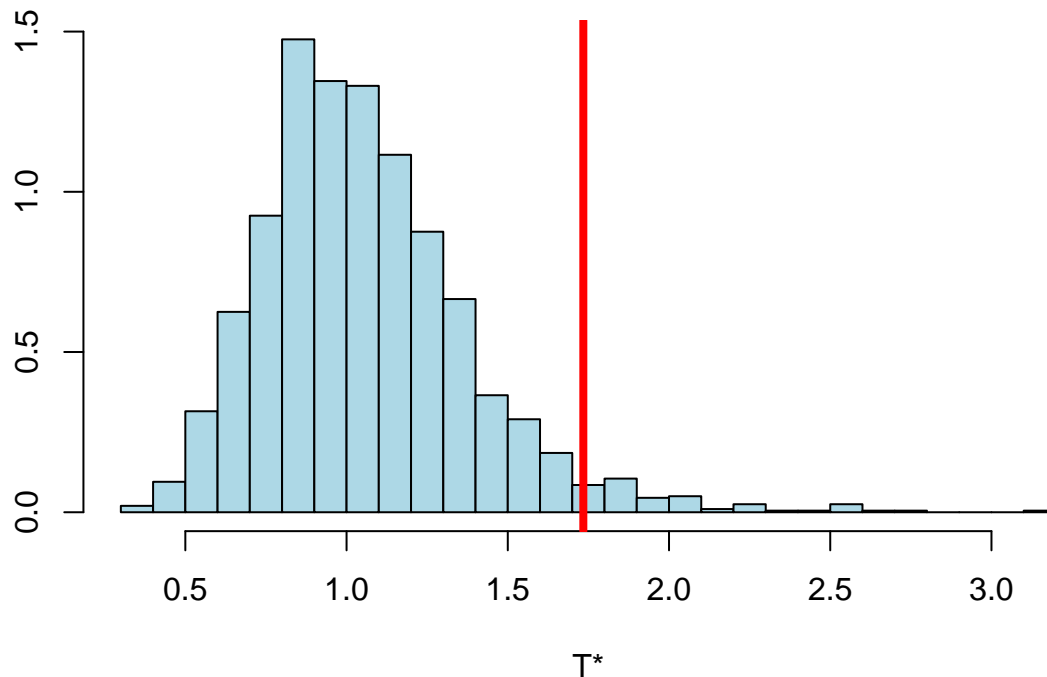
```
B<-1999
set.seed(125)
T0<-sd(x)/sd(y)
T0
```

```
## [1] 1.73431
```

```
Tast<-numeric(B)
for (b in 1:B)
  Tast[b]<-sd(sample(xy,nx,replace=TRUE))/sd(sample(xy,ny,replace=TRUE))

hist(Tast,br=30,col="lightblue",prob=TRUE,
     main=paste(B,"muestras bootstrap"),
     xlab="T*",ylab="")
abline(v=T0, lwd=4,col="red")
```

## 1999 muestras bootstrap



```
cat("p valor bootstrap unilateral (>) = ",(sum(Tast>=T0)+1)/(B+1),"\\n")
```

```
## p valor bootstrap unilateral (>) = 0.0325
```

```
# Tast>=T0 Cuántas veces el estimador T asterisco es mayor que T0
```

Es menor que 0.05, por lo que rechazo  $H_0$

### 1.2 Paso 2 (bilateral)

Como tenemos cantidades positivas mayores que uno, es una buena opción para aplicar logaritmo.

```
#Para un contraste bilateral se puede trabajar con log(sd(x)/sd(y))
```

```
hist(log(Tast),br=30,col="lightblue",
     prob=TRUE, main=paste(B,"muestras bootstrap"),
     xlab="log(T*)",ylab="")
```

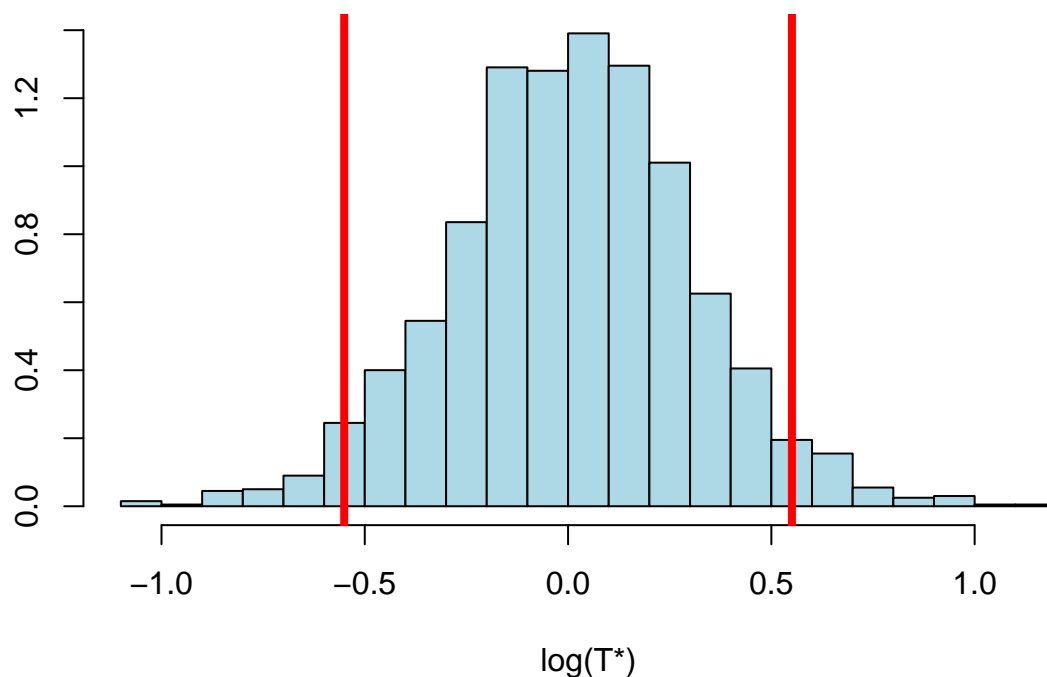
```
abline(v=log(T0), lwd=4,col="red")
```

```
cat("p valor bootstrap unilateral (>) = ",(sum(log(Tast)>=log(T0))+1)/(B+1),"\\n")
```

```
## p valor bootstrap unilateral (>) = 0.0325
```

```
abline(v=-log(T0), lwd=4,col="red")
```

## 1999 muestras bootstrap



```
cat("p valor bootstrap bilateral = ",  
    (sum(abs(log(Tast))>=abs(log(T0)))+1)/(B+1),"\n")
```

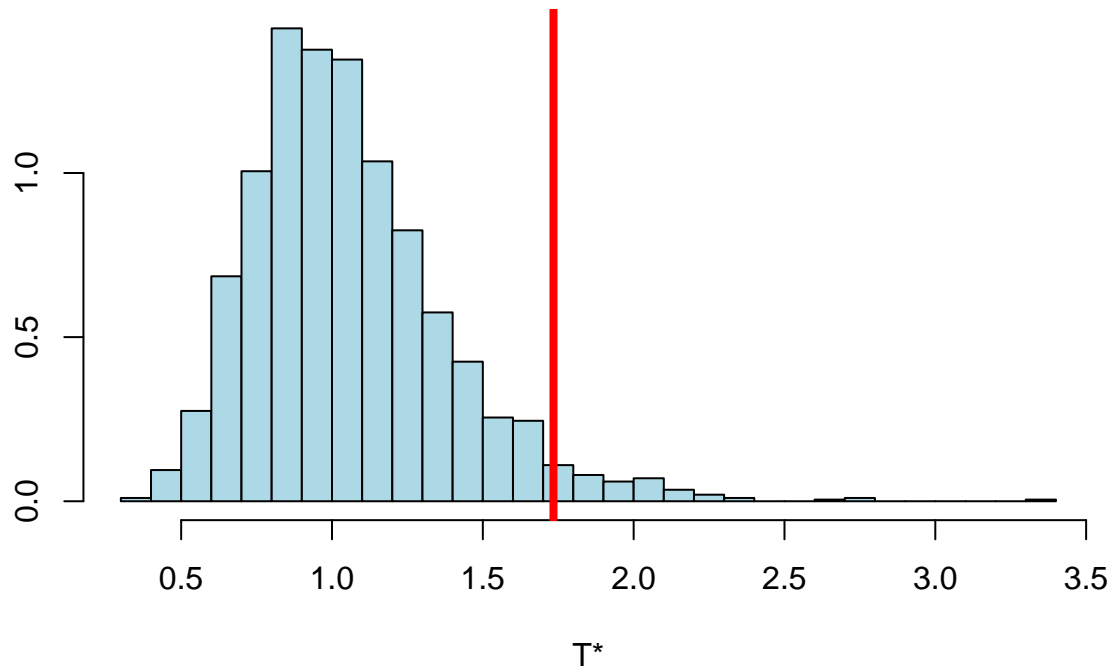
```
## p valor bootstrap bilateral = 0.063
```

Acepto  $H_0$

### 1.3 Parte 3 (Con la librería boot)

```
library(boot)  
cocboot=function(xy,indi,nx) {  
  sd(xy[indi[1:nx]])/sd(xy[indi[-c(1:nx)]])  
}  
resulboot=boot(xy,cocboot,nx=nx,R=1999)  
hist(resulboot$t,br=30,col="lightblue",prob=TRUE,  
     main=paste(B,"muestras bootstrap"),  
     xlab="T*",ylab="")  
abline(v=resulboot$t0, lwd=4,col="red")
```

## 1999 muestras bootstrap



```
cat("p valor bootstrap unilateral (>) = ",
    (sum(resulboot$t>=resulboot$t0)+1)/(B+1), "\n")
```

```
## p valor bootstrap unilateral (>) = 0.038
```

P-valor muy pequeño, rechazo  $H_0$

T0

```
## [1] 1.73431
```

```
resulboot$t0
```

```
## [1] 1.73431
```

```
summary(Tast)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.3414  0.8360  1.0135  1.0555  1.2233  3.1186
```

```
summary(resulboot$t)
```

```
##      V1
## Min.   :0.3874
## 1st Qu.:0.8283
## Median :1.0076
## Mean   :1.0558
## 3rd Qu.:1.2324
## Max.   :3.3592
```

```
cat("p valor bootstrap bilateral = ", (sum(abs(log(resulboot$t))>=abs(log(resulboot$t0)))+1)/(B+1), "\n")
```

```
## p valor bootstrap bilateral = 0.066
```

## 2 Ejercicio 2

2. Realizar mediante procedimientos bootstrap el contraste de hipótesis unilateral relativo a la posibilidad de que la primera componente principal del fichero iris de R explique más del 70% de la varianza total. Hacerlo directamente y con la ayuda de **boot**.

### 2.1 ACP

```
#Porcentaje de varianza explicada por la primera C.P.
##?Puede aceptarse que es superior al 70%?

##H0:PE1<=70; H1:PE1>70
set.seed(357)
data(iris) #?iris

ACP<- princomp(iris[, -5], cor=T)
summary(ACP)

## Importance of components:
##               Comp.1    Comp.2    Comp.3    Comp.4
## Standard deviation  1.7083611 0.9560494 0.38308860 0.143926497
## Proportion of Variance 0.7296245 0.2285076 0.03668922 0.005178709
## Cumulative Proportion 0.7296245 0.9581321 0.99482129 1.000000000

ACP$loadings

##
## Loadings:
##               Comp.1 Comp.2 Comp.3 Comp.4
## Sepal.Length  0.521  0.377  0.720  0.261
## Sepal.Width  -0.269  0.923 -0.244 -0.124
## Petal.Length  0.580          -0.142 -0.801
## Petal.Width   0.565          -0.634  0.524
##
##               Comp.1 Comp.2 Comp.3 Comp.4
## SS loadings      1.00   1.00   1.00   1.00
## Proportion Var    0.25   0.25   0.25   0.25
## Cumulative Var    0.25   0.50   0.75   1.00

varcp<- ACP$sdev^2
sum(varcp)

## [1] 4

PE1_0<- 100*varcp[1]/sum(varcp)
dife0<- PE1_0-70
n<-nrow(iris)
indices<- 1:n

B<- 4999
difeBoot<- numeric(B)
for (i in 1:B)
{ if (i%%500==0) cat("Muestra bootstrap ", i, "\n")
  ACPBoot<- princomp(iris[sample(indices, rep=T), -5], cor=T)
  varcpboot<- ACPBoot$sdev^2
  PE1Boot<- 100*varcpboot[1]/sum(varcpboot)
```

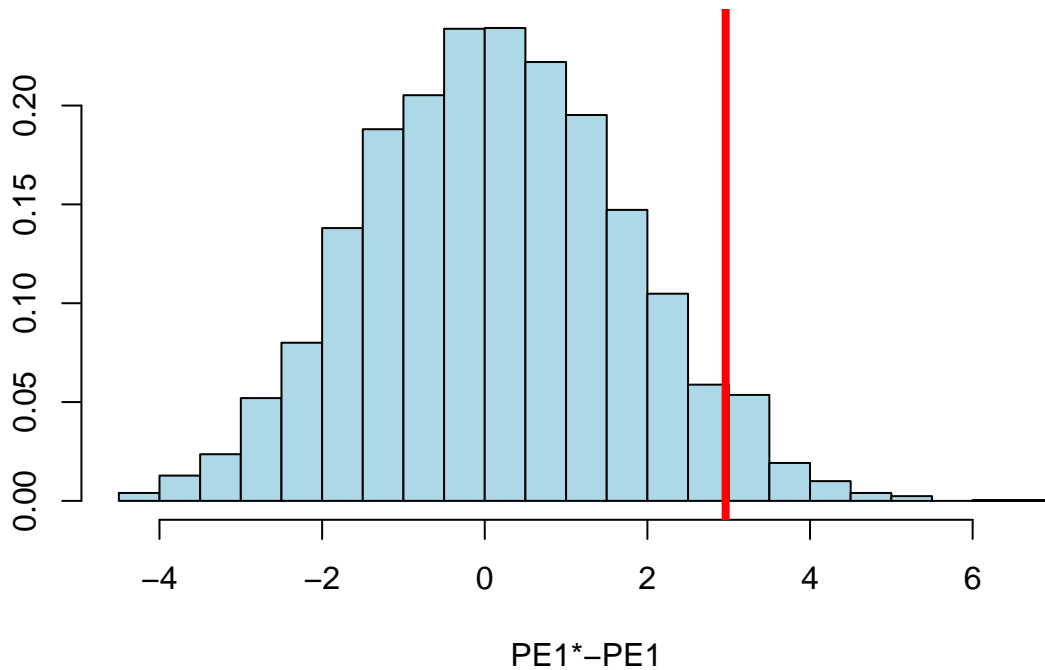
```
difeBoot[i]<- PE1Boot-PE1_0
}
```

```
## Muestra bootstrap 500
## Muestra bootstrap 1000
## Muestra bootstrap 1500
## Muestra bootstrap 2000
## Muestra bootstrap 2500
## Muestra bootstrap 3000
## Muestra bootstrap 3500
## Muestra bootstrap 4000
## Muestra bootstrap 4500
```

## 2.2 Histograma y p-valor

```
hist(difeBoot,br=30, # 30 intervalos
     col="lightblue",prob=TRUE,
     main=paste(B,"muestras bootstrap"),
     xlab="PE1*-PE1",ylab="")
abline(v=dife0, lwd=4,col="red")
```

### 4999 muestras bootstrap



```
cat("p valor bootstrap unilateral (>) = ",
    (sum(difeBoot>=dife0)+1)/(B+1),"\\n")
```

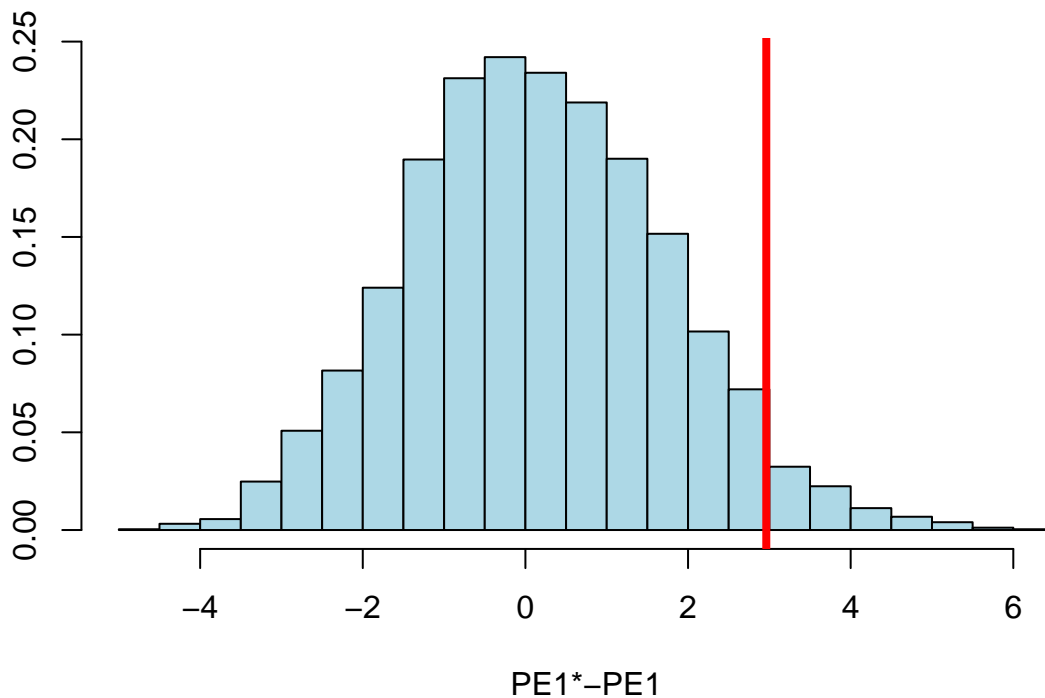
```
## p valor bootstrap unilateral (>) = 0.0472
```

## 2.3 Con librería BOOT

```
library(boot)
calculoPE1boot=function(datos,indi) {
  ACPBoot<- princomp(iris[indi,-5],cor=T)
  varcpboot<- ACPBoot$sdev^2
  PE1<- 100*varcpboot[1]/sum(varcpboot)
  PE1
}
resulboot=boot(iris,calculoPE1boot,R=4999)

hist(resulboot$t-resulboot$t0, # hacemos la diferencia a la hora de representar
     br=30,
     col="lightblue",
     prob=TRUE,
     main=paste(B,"muestras bootstrap (boot)"),
     xlab="PE1*-PE1",ylab="")
abline(v=resulboot$t0-70, lwd=4,col="red")
```

### 4999 muestras bootstrap (boot)



## 3 Ejercicio 3 IC bootstrap para un cociente

- Las siguientes instrucciones permiten definir dos variables en R que contienen las medidas de la corrosión (y) en 13 aleaciones de níquel-cobre, cada una de ellas con un contenido de hierro x. Es de interés el cambio en la corrosión cuando aumenta el nivel de hierro, comparado con la pérdida de corrosión cuando no hay hierro:  $\beta_1/\beta_0$  en el modelo de regresión lineal.

Calcular intervalos de confianza bootstrap para dicho cociente.

- $x \leftarrow c(0.01, 0.48, 0.71, 0.95, 1.19, 0.01, 0.48, 1.44, 0.71, 1.96, 0.01, 1.44, 1.96)$
- $y \leftarrow c(127.6, 124, 110.8, 103.9, 101.5, 130.01, 122, 92.3, 113.1, 83.7, 128, 91.4, 86.2)$

Se sabe que un estimador de la varianza de  $\beta_1/\beta_0$  mediante el método delta es:



$$\left(\frac{\hat{\beta}_1}{\hat{\beta}_0}\right)^2 \left(\frac{\hat{v}(\hat{\beta}_1)}{\hat{\beta}_1^2} + \frac{\hat{v}(\hat{\beta}_0)}{\hat{\beta}_0^2} - \frac{2cov(\hat{\beta}_0, \hat{\beta}_1)}{\hat{\beta}_0 \hat{\beta}_1}\right)$$

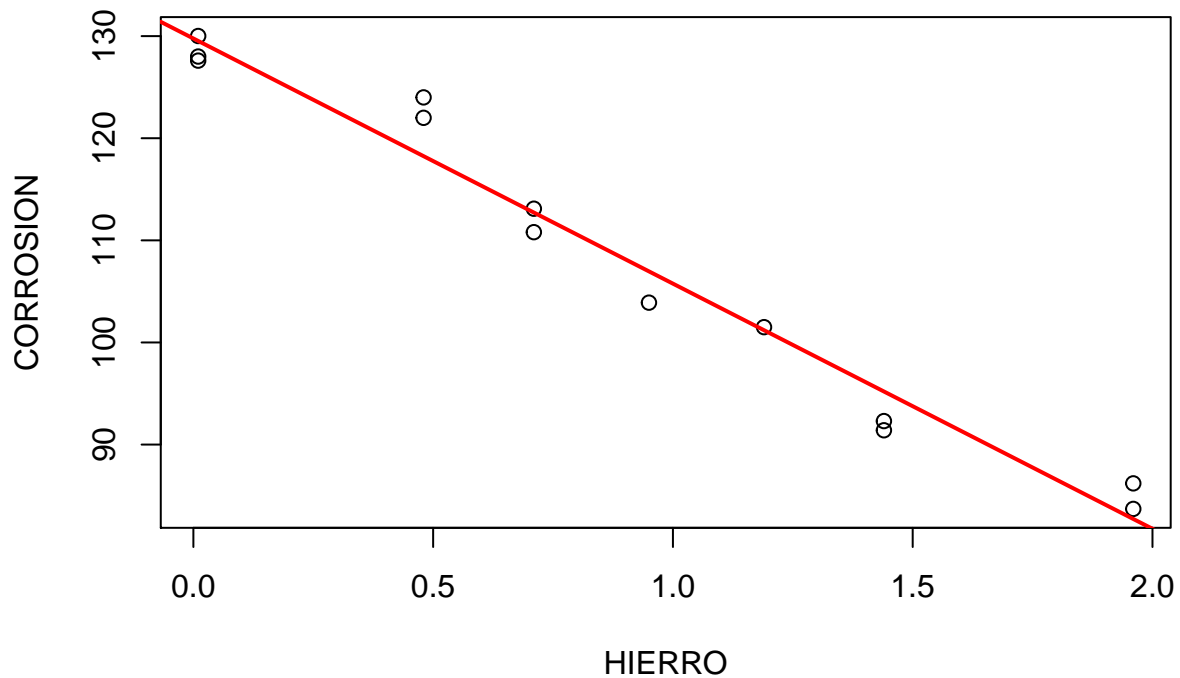
La función utilizada para calcular el estadístico de interés debe incluir también esta estimación de su varianza, para usarla en el método bootstrap-t.

```
x<- c(0.01,0.48,0.71,0.95,1.19,0.01,0.48,1.44,0.71,1.96,0.01,1.44,1.96)
y<- c(127.6,124,110.8,103.9,101.5,130.01,122,92.3,113.1,83.7,128,91.4,86.2)
cor(x,y)
```

```
## [1] -0.9847401
```

```
plot(x,y,xlab="HIERRO",ylab="CORROSION",main="")
regrelin<- lm(y~x)
summary(regrelin)
```

```
##
## Call:
## lm(formula = y ~ x)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.7987 -1.9277  0.2997  0.9845  5.7552
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  129.768      1.402    92.55 < 2e-16 ***
## x           -24.006      1.279   -18.77 1.06e-09 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.056 on 11 degrees of freedom
## Multiple R-squared:  0.9697, Adjusted R-squared:  0.967
## F-statistic: 352.2 on 1 and 11 DF, p-value: 1.057e-09
abline(regrelin,col="red",lwd=2)
```



```
coefici<- coef(regrelin)
print(T0<- coefici[2]/coefici[1])
```

```
##          x
## -0.1849942
```

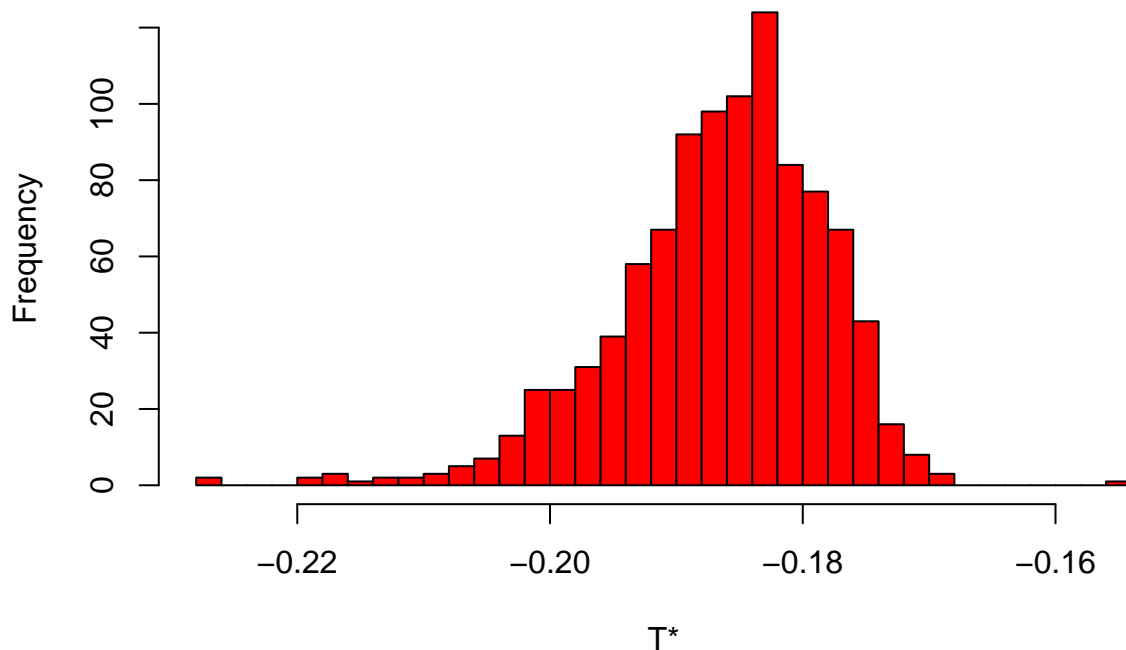
```
alfa<-0.05
xydat<- data.frame(x,y)
library(boot)
#Para bootstrap de pares
```

```
calcuT<-function(xydat,indi)
{
  regre<-lm(y~x,data=xydat[indi,])
  coefi<- regre$coefficients
  T<-coefi[2]/coefi[1]
  V<-vcov(regre) #cov(beta0,beta1)
  auxi<- (V[1,1]/(coefi[1]^2)) + (V[2,2]/(coefi[2]^2))-
    2*V[1,2]/(coefi[1]*coefi[2])
  varT<- (T^2) * auxi #var(T), método delta para bootstrap stud (t)
  c(T,varT)
}
```

```
Tbootpares<-boot(xydat,calcuT,1000) #Datos, estadístico, B
```

```
hist(Tbootpares$t[,1],br=30,main="Bootstrap de pares", col="red",xlab="T*")
```

## Bootstrap de pares



```
boot.ci(Tbootpares, conf=1-alfa, type=c("norm","perc","stud","bca"),
        var.t = Tbootpares$t[,2]) # a boot.ci le paso lo que me ha devuelto
```

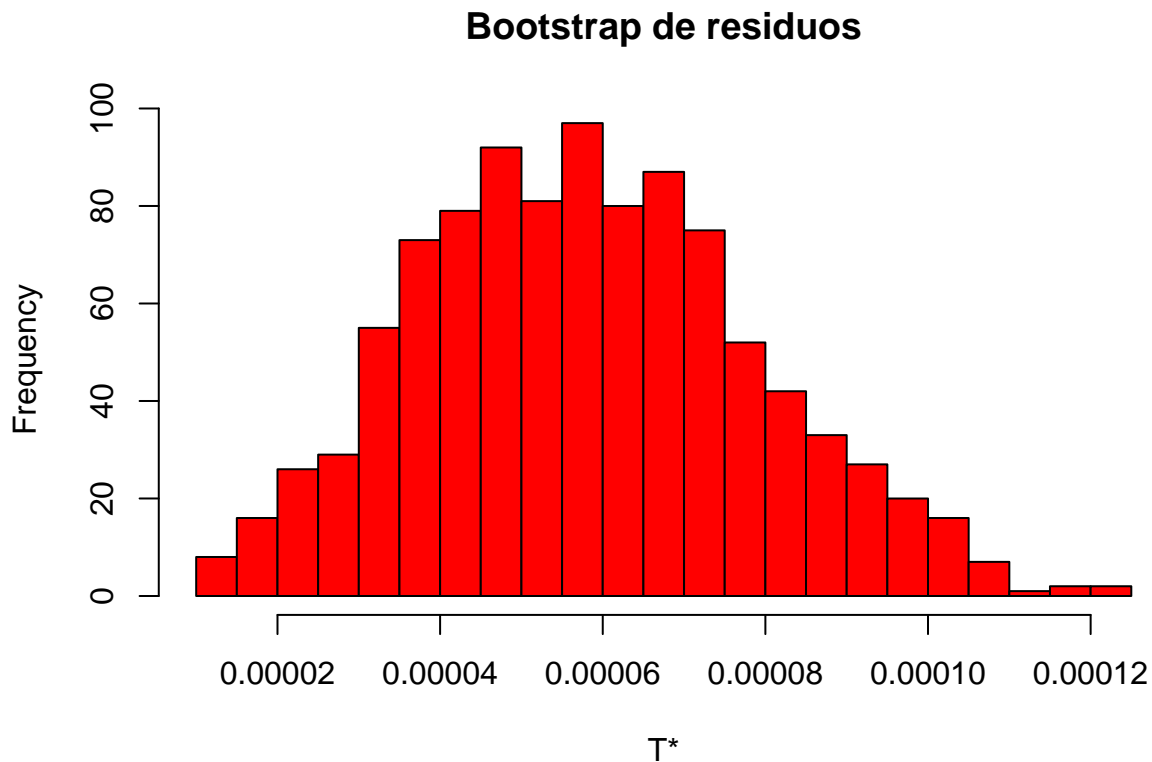
```
## BOOTSTRAP CONFIDENCE INTERVAL CALCULATIONS
## Based on 1000 bootstrap replicates
##
## CALL :
## boot.ci(boot.out = Tbootpares, conf = 1 - alfa, type = c("norm",
##      "perc", "stud", "bca"), var.t = Tbootpares$t[, 2])
##
## Intervals :
## Level      Normal          Studentized
## 95%  (-0.1996, -0.1673 )  (-0.2020, -0.1658 )
##
## Level      Percentile      BCa
## 95%  (-0.2048, -0.1737 )  (-0.2012, -0.1718 )
## Calculations and Intervals on Original Scale
```

*# la función boot*

```
#Bootstrap de residuos
calcuTres<-function(xydat,indi)
{ xydat$y<- predict(regrelin)+residuals(regrelin)[indi]
  regre<-lm(y~x,data=xydat)
  coefi<- regre$coefficients
  T<-coefi[2]/coefi[1]
  V<-vcov(regre)
  auxi<- (V[1,1]/(coefi[1]^2)) + (V[2,2]/(coefi[2]^2))-
    2*V[1,2]/(coefi[1]*coefi[2])
  varT<- (T^2) * auxi #var(T), m?todo delta para bootstrap stud (t) c(T,varT)
}
```

```
Tbootresi<-boot(xydat,calcuTres,1000)
```

```
hist(Tbootresi$t[,1],br=30,main="Bootstrap de residuos", col="red",xlab="T*")
```



```
#boot.ci(Tbootresi,  
#       conf=1-alfa,  
#       type=c("norm", "perc", "stud", "bca"),  
#       var.t = Tbootresi$t[,2])
```

## 4 Ejercicio 4

4. El fichero “salmon.dat” contiene datos anuales sobre una población de salmones. Las variables que aparecen son: R = “recruits” número de salmones que entran, y S = “spawners” número de salmones que están poniendo huevos, que mueren en cuanto lo hacen. Para que la población se estabilice se requiere: R = S (en otro caso, o hay demasiados salmones para los mismos recursos, o bien no se repone la población). Se desea calcular un I.C. bootstrap para el punto donde R = S, trabajando con el modelo de Beverton- Holt:

$$R = \frac{1}{\beta_0 + \beta_1/s}$$

## 5 Ejercicio 5

5. Diseñar un estudio empírico para analizar la efectividad de la transformación Z de Fisher en el cálculo de intervalos de confianza para el coeficiente de correlación lineal poblacional. Por ejemplo, generar 100 muestras de tamaño 10 de una Normal bivalente de media 0, desviaciones típicas 1 y correlación 0.6 y calcular los I.C. al 95% con los métodos Percentil, Normal y BCa. Comparar los cubrimientos y las longitudes medias utilizando procedimientos numéricos y gráficos. 1

## 6 Ejercicio 6

6. Estimar el sesgo de la razón (cociente de las medias de las variables  $x$  y  $u$ ) en el fichero “city” de R mediante el bootstrap balanceado (fichero disponible en la librería “boot”):
  - a. Escribiendo directamente las instrucciones.
  - b. Empleando la librería boot.

## 7 Ejercicio 7 Estimar el error de clasificación

7. Estimar el error de clasificación para el modelo de análisis discriminante lineal sobre los datos “cats” de la librería MASS.

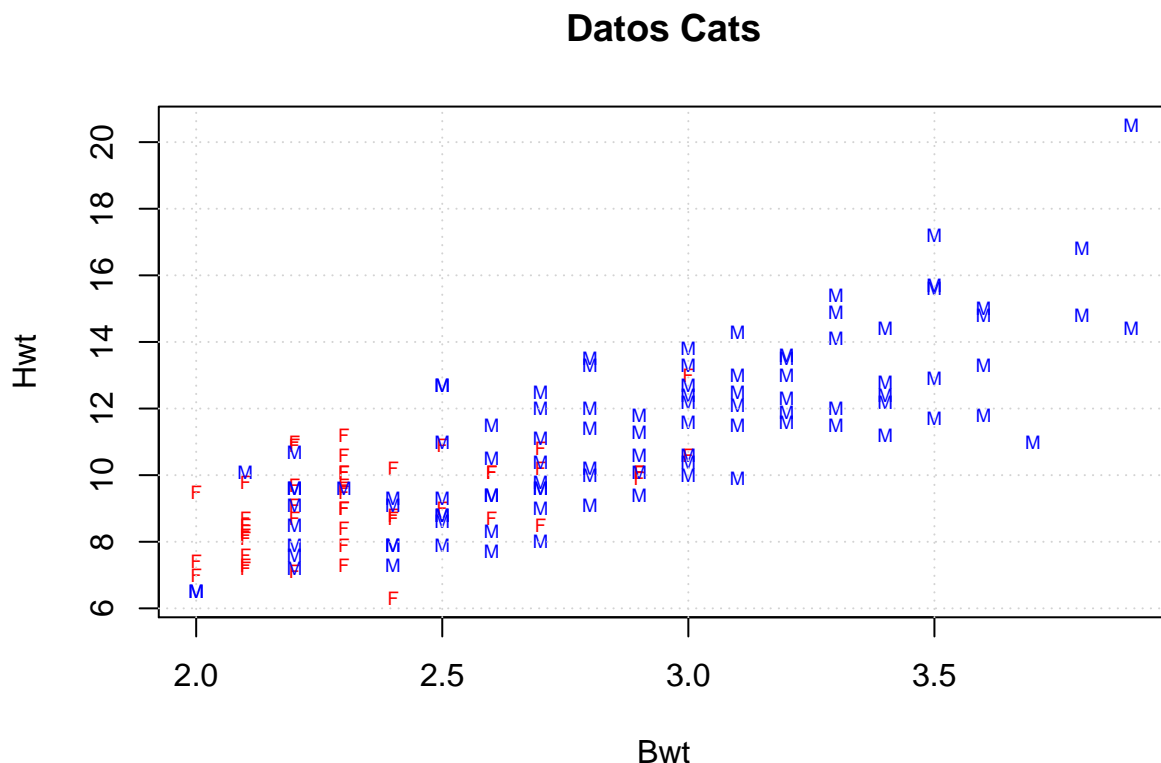
```
library(MASS)
data(cats)
summary(cats)
```

```
## Sex      Bwt      Hwt
## F:47  Min.   :2.000  Min.   : 6.30
## M:97  1st Qu.:2.300  1st Qu.: 8.95
##      Median :2.700  Median :10.10
##      Mean   :2.724  Mean   :10.63
##      3rd Qu.:3.025  3rd Qu.:12.12
##      Max.   :3.900  Max.   :20.50
```

Vamos a aplicar análisis discriminante lineal.

```
##cats
colores<- c("red","blue")
plot(cats[,2:3],type="n",main="Datos Cats")
text(cats[,2:3],as.character(cats$Sex),
     col=colores[cats$Sex],cex=0.6)

grid()
```



Hemos podido discriminar y ver a que categoría pertenecen, no va a ser un modelo perfecto y vamos a cometer errores probablemente.

Para obtener los datos del modelo discriminante (función lda)

```
cats.lda<-lda(Sex~.,cats)
cats.lda
```

```
## Call:
## lda(Sex ~ ., data = cats)
##
## Prior probabilities of groups:
##      F      M
## 0.3263889 0.6736111
##
## Group means:
##      Bwt      Hwt
## F 2.359574  9.202128
## M 2.900000 11.322680
##
## Coefficients of linear discriminants:
##      LD1
## Bwt  2.53019769
## Hwt -0.02986042
```

En LD1 está la línea que separa a las categorías.

```
table(cats$Sex,predict(cats.lda)$class)
```

```
##
##      F  M
## F 31 16
## M 12 85
```

```
#  $g_{\{\lambda\}}(li)=predict(cats.lda)$class$ 
```

Se tiene que:

- Aciertos: Diagonal
- Errores: diagonal opuesta

```
erroremp<- mean(cats$Sex!=predict(cats.lda)$class)
# Número de unos / número total de observaciones
cat("Error empírico=",100*erroremp ,"% \n")
```

```
## Error empírico= 19.44444 %
```

```
# Error empírico ó error de entrenamiento
```

```
#lda admite CV=TRUE que implementa n-VC (calcula el error esperado del modelo)
##o sea, Jackknife
cats.ldaJ<-lda(Sex~.,cats,CV=TRUE)
str(cats.ldaJ)
```

```
## List of 5
## $ class      : Factor w/ 2 levels "F","M": 1 1 1 1 1 1 1 1 1 1 ...
## $ posterior: num [1:144, 1:2] 0.769 0.772 0.783 0.707 0.708 ...
## .. attr(*, "dimnames")=List of 2
## .. ..$ : chr [1:144] "1" "2" "3" "4" ...
## .. ..$ : chr [1:2] "F" "M"
## $ terms      :Classes 'terms', 'formula' language Sex ~ Bwt + Hwt
## .. .. attr(*, "variables")= language list(Sex, Bwt, Hwt)
## .. .. attr(*, "factors")= int [1:3, 1:2] 0 1 0 0 0 1
```

```
## .. ..- attr(*, "dimnames")=List of 2
## .. ..$ : chr [1:3] "Sex" "Bwt" "Hwt"
## .. ..$ : chr [1:2] "Bwt" "Hwt"
## .. ..- attr(*, "term.labels")= chr [1:2] "Bwt" "Hwt"
## .. ..- attr(*, "order")= int [1:2] 1 1
## .. ..- attr(*, "intercept")= int 1
## .. ..- attr(*, "response")= int 1
## .. ..- attr(*, ".Environment")=<environment: R_GlobalEnv>
## .. ..- attr(*, "predvars")= language list(Sex, Bwt, Hwt)
## .. ..- attr(*, "dataClasses")= Named chr [1:3] "factor" "numeric" "numeric"
## .. ..- attr(*, "names")= chr [1:3] "Sex" "Bwt" "Hwt"
## $ call      : language lda(formula = Sex ~ ., data = cats, CV = TRUE)
## $ xlevels   : Named list()
```

Vamos a calcular el error de predicción:

```
table(cats$Sex,cats.ldaJ$class)
```

```
##
##      F  M
## F 31 16
## M 14 83
```

```
errorJ<- mean(cats$Sex!=cats.ldaJ$class)
cat("Error Jackknife=",100*errorJ ,"% \n")
```

```
## Error Jackknife= 20.83333 %
```

```
#Cómo se generan muestras bootstrap de conjuntos #de datos. Por ejemplo
datos=cats[c(1:5,140:144),]
datos
```

```
##      Sex Bwt  Hwt
## 1      F 2.0  7.0
## 2      F 2.0  7.4
## 3      F 2.0  9.5
## 4      F 2.1  7.2
## 5      F 2.1  7.3
## 140     M 3.7 11.0
## 141     M 3.8 14.8
## 142     M 3.8 16.8
## 143     M 3.9 14.4
## 144     M 3.9 20.5
```

Vamos a crear muestras bootstrap

```
#Estimaciones Bootstrap:
#####
B<- 2000 # Construyo 2000 muestras bootstrap
errorboot<-numeric(B)
error00B<- numeric(B)
n<- nrow(cats)
indin<- 1:n
for (b in 1:B)
{
  if (b%%500==0)
    cat("Muestra bootstrap número ",b,"\n") #Generar muestra bootstrap de los índices
  indiB<- sample(indin,rep=T) # muestra de entrenamiento
```



```

# los datos que no esten en la muestra de entrenamiento serán los datost
#Obtener los ?ndices no incluidos en imuestrab
indi00B<-setdiff(indin,indiB) # La diferencias (los que no han salido)
#Construir modelo lda sobre la muestra bootstrap
cats.lda.boot<-lda(Sex~.,cats[indiB,])
#Calcular tasa de error en la muestra original
errorboot[b]<- mean(cats$Sex!=predict(cats.lda.boot,cats)$class) # error
#Obtener predicciones OOB
predi00B<- predict(cats.lda.boot,cats[indi00B,])$class # predicciones sobre las observaciones no elegid
#Calcular la tasa de error OOB
error00B[b]<- mean(cats$Sex[indi00B]!=predi00B)
}

## Muestra bootstrap número 500
## Muestra bootstrap número 1000
## Muestra bootstrap número 1500
## Muestra bootstrap número 2000

errorB<- mean(errorboot) #no recomendable
errorB

## [1] 0.2224826

error00B<- mean(error00B) # media de todos los errores, ligera variante de  $Err_b^{1}$ 
error00B

## [1] 0.2357634

error632B<-0.368*erroremp+0.632*error00B # Estimación del error bootstrap
error632B

## [1] 0.220558

#Calcular cada elemento Lij #directamente:
matrizL<- matrix(NA,n,n)
for (i in 1:n)
  for (j in 1:n)
    matrizL[i,j]<- (cats[i,$Sex!=predict(cats.lda,cats[j,])$class)
print(Noinf<- mean(matrizL))

## [1] 0.4300733

#0 bien, en un problema de clasificación, #con error 0-1,
# Noinf se puede calcular #de forma más eficiente:
p1<- mean(cats$Sex=="M")
q1<- mean(predict(cats.lda)$class=="M")
p1*(1-q1)+(1-p1)*q1

## [1] 0.4300733

Noinf

## [1] 0.4300733

erroremp

## [1] 0.1944444

error00B=min(error00B,Noinf)
error00B

```

```
## [1] 0.2357634
#Cálculo de tsr y w
(tsr<- (error00B-erroremp)/(Noinf-erroremp))

## [1] 0.1753561
(w<- 0.632/(1-0.368*tsr))

## [1] 0.675597
#Posible redefinición de tsr
if ( (error00B<= erroremp) | (Noinf <=erroremp) ) # en alguno de estos casos sedará la corrección
  tsr=0
tsr

## [1] 0.1753561
#(error632masB<-(1-w)*erroremp+w*error00B)
# #Definición general (la línea anterior no vale #si se redefinen tsr o error00B)
error632masB=error632B+
(error00B-erroremp)*(0.368*0.632*tsr)/(1-0.368*tsr)
error632masB

## [1] 0.2223594
op_ant = options(digits=4)
cat(" Error Empírico=\t",100*erroremp ,"% \n",
"Error Jackknife=\t",100*errorJ ,"% \n",
"Error 00B=\t\t", 100*error00B,"% \n",
"Error 0.632Boot=\t", 100*error632B,"% \n",
"Error 0.632+Boot=\t", 100*error632masB,"% \n")

## Error Empírico= 19.44 %
## Error Jackknife= 20.83 %
## Error 00B= 23.58 %
## Error 0.632Boot= 22.06 %
## Error 0.632+Boot= 22.24 %

options(op_ant)
```

## 8 Ejercicio 8

8. Estimar el error de predicción para los datos “Renta.txt”, siendo “rentsqm” (precio del alquiler por m2) la variable dependiente de un modelo de regresión lineal. Utilizar el criterio RECM.

### 8.1 Parte 1 Leer datos

ESTIMACION DEL ERROR DE PREDICCIÓN (REGRESION)

```
#1. Leer los datos
```

```
datos<-read.table("datos/Renta.txt",header=T)
dim(datos)
```

```
## [1] 118 6
```

```
summary(datos)
```

```
##      rentsqm      yearc      locat      bath
## Min.   : 5.146   Min.   :1918   Min.   :1.000   Min.   :0.00000
## 1st Qu.: 8.533   1st Qu.:1939   1st Qu.:1.000   1st Qu.:0.00000
## Median : 9.344   Median :1959   Median :2.000   Median :0.00000
## Mean   : 9.396   Mean   :1957   Mean   :1.771   Mean   :0.04237
## 3rd Qu.:10.405   3rd Qu.:1971   3rd Qu.:2.000   3rd Qu.:0.00000
## Max.   :12.613   Max.   :1995   Max.   :3.000   Max.   :1.00000
##      kitchen      cheating
## Min.   :0.0000   Min.   :0.0000
## 1st Qu.:0.0000   1st Qu.:1.0000
## Median :0.0000   Median :1.0000
## Mean   :0.0339   Mean   :0.8983
## 3rd Qu.:0.0000   3rd Qu.:1.0000
## Max.   :1.0000   Max.   :1.0000
```

### 8.2 Parte 2 (Modelo de regresión lineal múltiple)

```
modelo=lm(rentsqm~.,data=datos)
#Calcular MSE y RMSE empírico
error_emp=mean(residuals(modelo)^2)
RMSE_emp=sqrt(error_emp)
```

### 8.3 Parte 3 (Estimaciones Jackknife y bootstrap)

```
#3.1. Jackknife
```

```
#Se pueden calcular con cv.lm
```

```
#o bien recorriendo los n modelos #cada uno se construye dejando fuera
```

```
##el caso i, donde se aplica el
```

```
###modelo para calcular prediJ[i]
```

```
n=nrow(datos)
```

```
prediJ = numeric(n)
```

```
for(i in 1:n){
  modelo.i = lm(rentsqm~.,data=datos[-i,])
  prediJ[i]<-predict(modelo.i,datos[i,])
}
resi_J=datos$rentsqm - prediJ
RMSE_J<-sqrt( mean(resi_J^2) )
```

## 8.4 Parte 4 (generamos las muestras bootstrap)

```
#3.2. Bootstrap:
#Se pueden calcular los estimadores de ECM (MSE)
# y al final visualizar su raíz cuadrada

B<-2000
errorboot<-numeric(B)
errorOOB<- numeric(B)
indin<-1:n

for(b in 1:B){
  if (b%%500==0) cat("Muestra bootstrap número ",b,"\n")
  #Generar muestra bootstrap de los índices
  indiB<- sample(indin,rep=T)
  #Obtener los índices no incluidos en imuestrab
  indiOOB<-setdiff(indin,indiB)
  #Construir modelo lda sobre la muestra bootstrap
  modelo.boot<-lm(rentsqm~.,data=datos[indiB,])
  #Calcular ECM en la muestra original
  suppressWarnings({ errorboot[b]<- mean((datos$rentsqm[indiB]-predict(modelo.boot,datos))^2)})
  #Obtener predicciones OOB
  suppressWarnings({ prediOOB<- predict(modelo.boot,datos[indiOOB,]) })
  #Calcular ECM OOB
  errorOOB[b]<- mean((datos$rentsqm[indiOOB]-prediOOB)^2)
}

## Muestra bootstrap número 500
## Muestra bootstrap número 1000
## Muestra bootstrap número 1500
## Muestra bootstrap número 2000

Y los errores:

errorB<- mean(errorboot)
errorOOB<- mean(errorOOB)
error632B<-0.368*error_emp+0.632*errorOOB

#Calcular cada elemento L_ij
#Al no ser un problema de clasificación, #se debe calcular directamente:
matrizL<- matrix(NA,n,n)
for (i in 1:n)
  for (j in 1:n)
    matrizL[i,j]<- (datos$rentsqm[i]-predict(modelo,datos[j,]))^2 #error cuad.
print(Noinf<- mean(matrizL))

## [1] 3.164884

#Redefinición de errorOOB
##(por si hiciera falta)
errorOOB #Noinf< error_emp ?

## [1] 0.9459989
Noinf

## [1] 3.164884
```

```

error00B=min(error00B,Noinf)
error00B

## [1] 0.9459989
#Cálculo de tsr y w
(tsr<- (error00B-error_emp)/(Noinf-error_emp))

## [1] 0.06005573
(w<- 0.632/(1-0.368*tsr))

## [1] 0.6462832
#Posible redefinición de tsr
if ( (error00B<= error_emp) | (Noinf <= error_emp) )
tsr=0
tsr

## [1] 0.06005573
#(error632masB<-(1-w)*erroremp+w*error00B) #Definición general (la línea anterior no vale #si se redefi

error632masB=error632B+
(error00B-error_emp)*(0.368*0.632*tsr)/(1-0.368*tsr)
error632masB

## [1] 0.8958522
cat(" RMSE Empírico=\t",sqrt(RMSE_emp), "\n",
" RMSE Jackknife=\t", RMSE_J, "\n",
" RMSE 00B=\t\t", sqrt(error00B),"\n",
" RMSE 0.632Boot=\t", sqrt(error632B),"\n",
" RMSE 0.632+Boot=\t", sqrt(error632masB),"\n")

## RMSE Empírico= 0.9469887
## RMSE Jackknife= 0.9483709
## RMSE 00B= 0.9726248
## RMSE 0.632Boot= 0.9454244
## RMSE 0.632+Boot= 0.9464947

```