

Ejercicios Tema 1. Regresión lineal simple

Máster en Ciencia de Datos. Módulo: Análisis exploratorio de datos

Ana Navarro Quiles

Curso 2022/2023

Objetivos:

Con estos ejercicios se pretende practicar los conceptos vistos en el Tema 1: Regresión Lineal Simple

- Cálculo e interpretación de la recta de regresión lineal.
- Coeficiente de determinación y coeficiente de correlación de Pearson.
- Modelo de regresión lineal simple (intervalos de confianza, contrastes de hipótesis, predicciones).
- Condiciones de validez del modelo.
- Transformaciones y alternativas no lineales.

Comandos útiles

A continuación indicamos algunos comandos útiles que complementan aquellos que ya hemos visto en la teoría. Para la realización de los ejercicios se recomienda revisar los comandos vistos durante el Tema 1.

Regresión lineal

Como ejemplo de aplicación vamos utilizar el banco de datos **deportistas** (que podemos encontrar en **datosTema1.RData**), empleando la variable suma de pliegues, *SumPliegues*, para explicar el porcentaje de grasa, *PrctGrasa*, que es la variable respuesta.

En particular, vamos a realizar la regresión lineal para el subconjunto de los hombres, *male*.

```
load('./data/datosTema1.Rdata')

#Selección subconjunto:
hombres <- subset(deportistas, Genero=='male')

#Ajuste modelo
lm_hombres <- lm(PrctGrasa ~ SumPliegues, data=hombres)

#Resumen Modelo
summary(lm_hombres)

#Cálculo de residuos
residuos<-residuals(lm_hombres)

#Representación del ajuste
plot(hombres$SumPliegues, hombres$PrctGrasa, col='BLUE')
abline(coef=coef(lm_hombres), col='RED')
```

```

#Extracción de coeficientes y sus Intervalos de Confianza
coef(lm_hombres)
confint(lm_hombres)

#Obtención de bandas de estimación:
minx<-range(hombres$SumPliegues)[1]; maxx<-range(hombres$SumPliegues)[2]
nuevos <- data.frame(list(SumPliegues = seq(minx,maxx,length=100)))
bandas_est<-predict(lm_hombres, newdata = nuevos, interval = "confidence")

#Representación gráfica:
plot(hombres$SumPliegues, hombres$PrctGrasa, col='BLUE')
abline(coef=coef(lm_hombres), col='RED')
lines(nuevos$SumPliegues,bandas_est[,2],col='BLACK')
lines(nuevos$SumPliegues,bandas_est[,3],col='BLACK')

#predicción x0=100:
predict100<-predict(lm_hombres, newdata = data.frame(SumPliegues=c(100)),
                    interval = "prediction")

#Obtención de bandas de predicción
minx<-range(hombres$SumPliegues)[1]; maxx<-range(hombres$SumPliegues)[2]
nuevos <- data.frame(list(SumPliegues = seq(minx,maxx,length=100)))
bandas_pred<-predict(lm_hombres, newdata = nuevos, interval = "prediction")

#Representación gráfica:
plot(hombres$SumPliegues, hombres$PrctGrasa, col='BLUE')
abline(coef=coef(lm_hombres), col='RED')
lines(nuevos$SumPliegues,bandas_pred[,2],col='BLACK')
lines(nuevos$SumPliegues,bandas_pred[,3],col='BLACK')

# diagnóstico linealidad y homocedasticidad
residuos <- residuals(lm_hombres)
predichos <- fitted.values(lm_hombres)
par(mfcol=c(1,2))
plot(predichos,residuos, col='BLUE',main = 'Gráfica de residuos')
abline(h=0,lty=2)

# diagnóstico normalidad residuos
qqnorm(residuos, col='BLUE')
qqline(residuos)

```

Alternativas no lineales. Modelos más flexibles

Como ejemplo de aplicación usamos los datos de la base **Boston** que está en el paquete de R **MASS**. Son datos sobre los suburbios de Boston, con variables como el precio medio de la vivienda *medv*, que vamos a utilizar como variable respuesta, y el estatus de la población *lstat* que vamos a utilizar como predictora.

```
library(MASS)
attach(Boston)

#ajuste knn vecinos
require(FNN) #libreria
xx <- seq(min(lstat),max(lstat),0.25) # puntos (ordenados)
new_lstat <- data.frame(list(lstat = xx))
reg_knn_80 <- knn.reg(lstat, new_lstat, y=medv, k=5) #ajuste para k=5

#representación gráfica knn
plot(lstat, medv,cex=.4)
lines(xx,reg_knn_80$pred,col='GREEN',lwd=2)

#ajuste loess
xx <- seq(min(lstat),max(lstat),0.25) # puntos (ordenados)
new_lstat <- data.frame(list(lstat = xx))
reg_loess_20 <- loess(medv ~ lstat, span = 0.20) #ajuste
pred_loess_20 <- predict(reg_loess_20, newdata = new_lstat, se = T) #prediccion

#representación gráfica loess
plot(lstat, medv,cex=.4)
lines(xx,pred_loess_20$fit,col='BLUE',lwd=2)
```

Ejercicios propuestos

Ejercicio 1

Utilizando el banco de datos **deportistas**, considerad la variable respuesta *Peso* relacionandola con el predictor *PrctGrasa*

a) ¿Cuánto vale la pendiente de la recta? ¿Podemos afirmar que es positiva?

```
fit1<-lm(Peso~PrctGrasa, data=deportistas)
fit1_s<-summary(fit1)
fit1_s

##
## Call:
## lm(formula = Peso ~ PrctGrasa, data = deportistas)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -37.210  -8.482  -0.608   9.116  48.192
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  75.0137911   2.3625769   31.751  <2e-16 ***
## PrctGrasa    -0.0004346   0.1590772   -0.003    0.998
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 13.96 on 200 degrees of freedom
## Multiple R-squared:  3.732e-08, Adjusted R-squared:  -0.005
## F-statistic: 7.464e-06 on 1 and 200 DF,  p-value: 0.9978
```

Respuesta: la pendiente de la recta vale $-4.3459764 \times 10^{-4}$, negativa. Es prácticamente nula y su p-valor asociado al contraste es 0.9978229, por lo tanto la variable es no significativa. Si calculamos el intervalo de confianza para la pendiente obtenemos:

```
confint(fit1)

##              2.5 %      97.5 %
## (Intercept) 70.3550347 79.6725475
## PrctGrasa   -0.3141184  0.3132492
```

Como podemos apreciar, no podemos afirmar con un 95% de confianza que la pendiente sea no negativa.

b) Compara la varianza de la variable respuesta con la varianza de los residuos: ¿Qué porcentaje de la variabilidad inicial está explicado por la recta de mínimos cuadrados? ¿Qué porcentaje de la variabilidad inicial falta todavía por explicar?

```
# Comparación de las varianzas
var(deportistas$Peso)

## [1] 193.9112

var(fit1$residuals)

## [1] 193.9112

# Porcentaje de la variabilidad inicial explicada por la recta de mínimos cuadrados
fit1_s$r.squared*100
```

```
## [1] 3.731889e-06
```

```
# Porcentaje de la variabilidad NO explicada por el modelo  
100-fit1_s$r.squared*100
```

```
## [1] 100
```

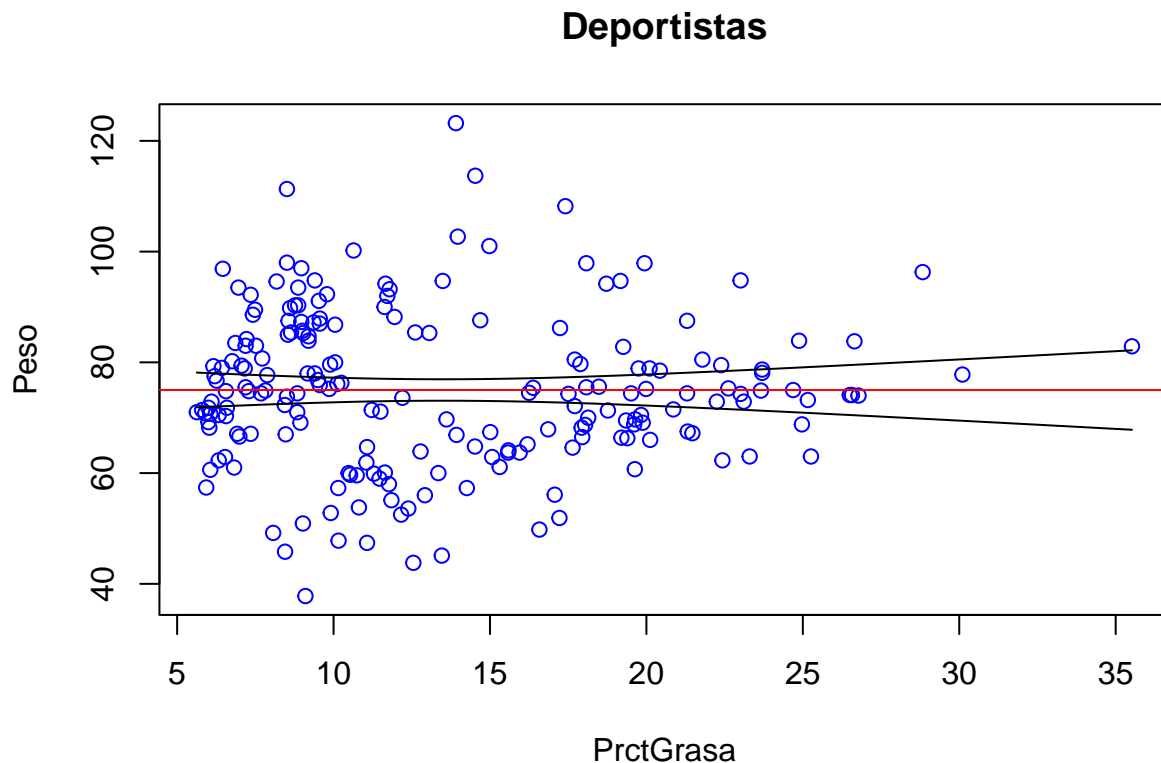
c) Obtén los intervalos de confianza al 95% sobre los parámetros de la recta.

```
confint(fit1)
```

```
##                2.5 %    97.5 %  
## (Intercept) 70.3550347 79.6725475  
## PrctGrasa   -0.3141184  0.3132492
```

d) Dibuja el diagrama de dispersión, la recta de regresión y las bandas de confianza para la estimación al 95%.

```
minx<-range(deportistas$PrctGrasa)[1]  
maxx<-range(deportistas$PrctGrasa)[2]  
  
nuevos <- data.frame(list(PrctGrasa = seq(minx,maxx,length=100)))  
bandas_est<-predict(fit1, newdata = nuevos, interval = "confidence")  
  
plot(deportistas$PrctGrasa, deportistas$Peso, col='BLUE', xlab="PrctGrasa", ylab="Peso",  
      main="Deportistas")  
abline(coef=coef(fit1), col='RED')  
lines(nuevos$PrctGrasa,bandas_est[,2],col='BLACK')  
lines(nuevos$PrctGrasa,bandas_est[,3],col='BLACK')
```



- e) Si te parece adecuado estima el peso correspondiente a nuevos individuos con los siguientes porcentajes de grasa: 25, 50, 75%. Calcula sus respectivos intervalos de confianza al 95%.

No es adecuado hacer una predicción con este modelo ya que hemos obtenido un R^2 muy pequeño. No obstante, para practicar, la realizaremos igualmente.

```
predict25<-predict(fit1, newdata=data.frame(PrctGrasa=c(25)), interval = "prediction")
predict50<-predict(fit1, newdata=data.frame(PrctGrasa=c(50)), interval = "prediction")
predict75<-predict(fit1, newdata=data.frame(PrctGrasa=c(75)), interval = "prediction")
```

```
predict25
```

```
##          fit          lwr          upr
## 1 75.00293 47.17278 102.8331
```

```
predict50
```

```
##          fit          lwr          upr
## 1 74.99206 45.11636 104.8678
```

```
predict75
```

```
##          fit          lwr          upr
## 1 74.9812 41.3123 108.6501
```

Ejercicio 2

Repite el ejercicio anterior considerando *IMC* en lugar de peso y compara los resultados con los del ejercicio anterior.

- a) ¿Cuánto vale la pendiente de la recta? ¿Podemos afirmar que es positiva?

```
fit2<-lm(IMC~PrctGrasa, data=deportistas)
fit2_s<-summary(fit2)
fit2_s
```

```
##
## Call:
## lm(formula = IMC ~ PrctGrasa, data = deportistas)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -5.7340 -2.0182 -0.1511  1.4287 11.4292
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  21.78372    0.47728   45.64  < 2e-16 ***
## PrctGrasa     0.08678    0.03214    2.70  0.00752 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.82 on 200 degrees of freedom
## Multiple R-squared:  0.03518,    Adjusted R-squared:  0.03035
## F-statistic: 7.292 on 1 and 200 DF,  p-value: 0.007519
```

Respuesta: la pendiente de la recta vale 0.08678, positiva. Es prácticamente nula pero su p-valor asociado al contraste es 0.0075195, por lo tanto la variable es significativa y podemos descartar que sea nula.

- b) Compara la varianza de la variable respuesta con la varianza de los residuos: ¿Qué porcentaje de la

variabilidad inicial está explicado por la recta de mínimos cuadrados? ¿Qué porcentaje de la variabilidad inicial falta todavía por explicar?

```
# Comparación de las varianzas
```

```
var(deportistas$IMC)
```

```
## [1] 8.202111
```

```
var(fit2$residuals)
```

```
## [1] 7.913578
```

```
# Porcentaje de la variabilidad inicial explicada por la recta de mínimos cuadrados
```

```
fit2_s$r.squared*100
```

```
## [1] 3.517791
```

```
# Porcentaje de la variabilidad NO explicada por el modelo
```

```
100-fit2_s$r.squared*100
```

```
## [1] 96.48221
```

c) Obtén los intervalos de confianza al 95% sobre los parámetros de la recta.

```
confint(fit2)
```

```
##                2.5 %    97.5 %
```

```
## (Intercept) 20.84257560 22.724859
```

```
## PrctGrasa    0.02341092 0.150149
```

d) Dibuja el diagrama de dispersión, la recta de regresión y las bandas de confianza para la estimación al 95%.

```
minx<-range(deportistas$PrctGrasa)[1]
```

```
maxx<-range(deportistas$PrctGrasa)[2]
```

```
nuevos <- data.frame(list(PrctGrasa = seq(minx,maxx,length=100)))
```

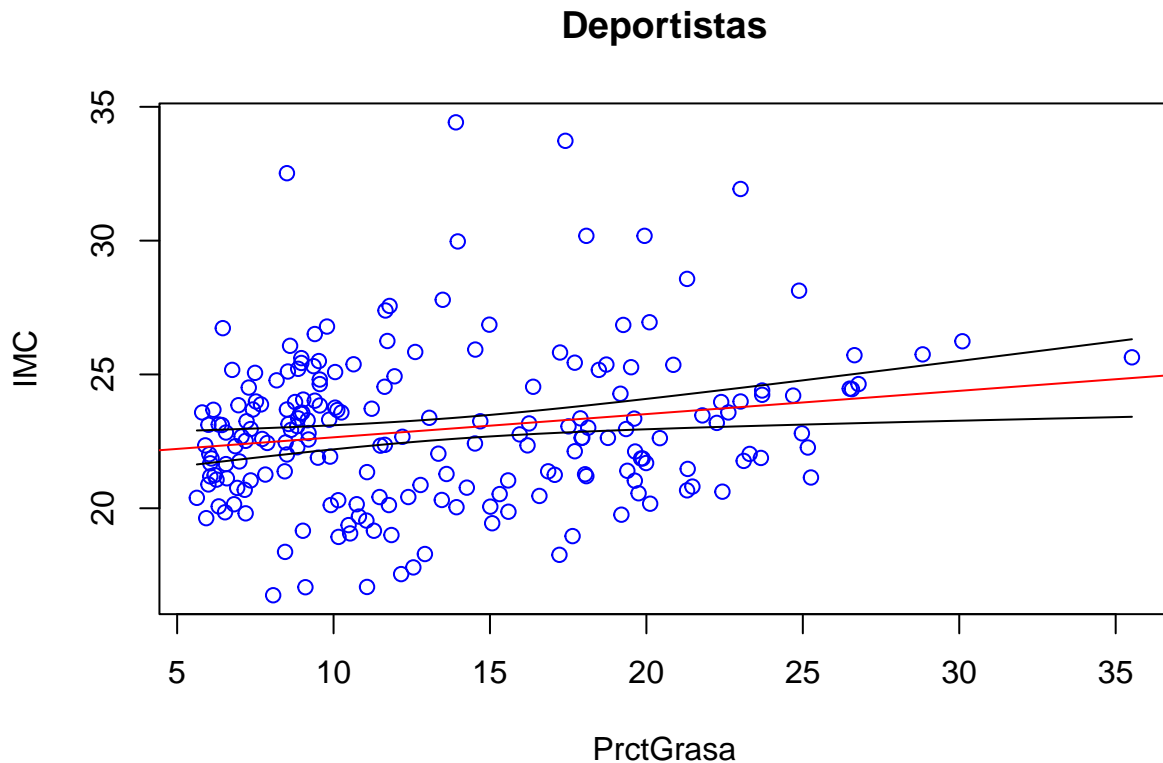
```
bandas_est<-predict(fit2, newdata = nuevos, interval = "confidence")
```

```
plot(deportistas$PrctGrasa, deportistas$IMC, col='BLUE', xlab="PrctGrasa", ylab="IMC",  
      main="Deportistas")
```

```
abline(coef=coef(fit2), col='RED')
```

```
lines(nuevos$PrctGrasa,bandas_est[,2],col='BLACK')
```

```
lines(nuevos$PrctGrasa,bandas_est[,3],col='BLACK')
```



- e) Si te parece adecuado estima el peso correspondiente a nuevos individuos con los siguientes porcentajes de grasa: 25, 50, 75%. Calcula sus respectivos intervalos de confianza al 95%.

No es adecuado hacer una predicción con este modelo ya que hemos obtenido un R^2 muy pequeño. No obstante, para practicar, la realizaremos igualmente.

```
predict25<-predict(fit2, newdata=data.frame(PrctGrasa=c(25)), interval = "prediction")
predict50<-predict(fit2, newdata=data.frame(PrctGrasa=c(50)), interval = "prediction")
predict75<-predict(fit2, newdata=data.frame(PrctGrasa=c(75)), interval = "prediction")
```

```
predict25
```

```
##          fit          lwr          upr
## 1 23.95322 18.33109 29.57534
```

```
predict50
```

```
##          fit          lwr          upr
## 1 26.12271 20.08736 32.15807
```

```
predict75
```

```
##          fit          lwr          upr
## 1 28.29221 21.49057 35.09386
```

Ejercicio 3

Utilizando el banco de datos **deportistas.cs**, considerad la variable respuesta *PrctGrasa* relacionándola con el predictor *MCMagra*.

- Obtén la recta mínimos cuadrados utilizando todos los datos, sin tener en cuenta el *sexo*.
- Evalúa el efecto del *sexo* sobre *PrctGrasa*.
- Obtén ahora una recta para los *hombres* y otra para las *mujeres*.
- Dibuja en la misma gráfica las tres rectas y comenta los resultados.

Ejercicio 4

En la base **deportistas**,

- Evalúa mediante regresión lineal si el *PrctGrasa* explica los resultados analíticos: *Hematocrito*, *Ferritina* y *Hemoglobina*.
- Evalúa mediante regresión lineal si *Peso* y *Altura* explican el *PrctGrasa*.
- Evalúa la relación entre *IMC* y las variables *SumPliegues* y *PrctGrasa*.

Ejercicio 5

Utilizando el modelo $Y = 25 + 2X + \epsilon$, siendo ϵ Normal con media 0 varianza $\sigma^2 = 4$, simula $N = 10000$ muestras de tamaño $n = 50$. Para ello, utiliza valores de X simulados de una Uniforme definida en el intervalo $(0, 5)$. A continuación, para cada una de las muestras simuladas, obtén el intervalo de confianza al 95% sobre la pendiente de la recta. ¿Qué porcentaje de intervalos no contienen al verdadero valor de la pendiente?

```
set.seed(1)
N<-10000
n<-50
min_x<-0
max_x<-5
acierto<-0

for (i in 1:N){
  error<-rnorm(n, mean=0, sd=2)
  x<-runif(n, min=min_x, max=max_x)
  y<-25+2*x+error
  fit<-lm(y~x, data=data.frame(x=x, y=y))
  intervalo<-confint(fit)["x",]
  if(2>intervalo[1] & 2<intervalo[2]){
    acierto=acierto+1
  }
}
por1<-100-acierto/N*100
por1
```

```
## [1] 4.97
```

Vuelve a calcular ese porcentaje, pero ahora simulando ϵ de forma que $\epsilon/8$ sea t-Student con 4 grados de libertad. ¿Cómo afecta la falta de normalidad a la fiabilidad de ese intervalo?

```
N<-10000
n<-50
df<-4
min_x<-0
max_x<-5
acierto<-0

for (i in 1:N){
```

```

error<-8*rt(n, df)
x<-runif(n, min=min_x, max=max_x)
y<-25+2*x+error
fit<-lm(y~x, data=data.frame(x=x, y=y))
intervalo<-confint(fit)["x",]
if(2>intervalo[1] & 2<intervalo[2]){
  acierto=acierto+1
}
}
por2<-100-acierto/N*100
por2

```

```
## [1] 4.73
```

Ejercicio 6

Utilizando el banco de datos **Auto**, en el paquete de R **ISLR**, se desea explicar el consumo de carburante, variable *mpg*, a partir de la potencia del motor, variable *horsepower*.

- Dibuja el diagrama de dispersión y la recta de mínimos cuadrados.
- ¿Hay relación entre esas dos variables? ¿Cómo de fuerte es esa relación? ¿Podemos afirmar si es positiva o negativa?
- ¿Qué consumo se espera si potencia del motor es 75? Proporciona el intervalo de confianza y el de predicción para esa potencia de motor.
- Analiza gráficamente los residuos y comenta los resultados.

Ejercicio 7

El banco de datos **cerebros** es un banco de datos famoso. En él se recogen los pesos del cuerpo y del cerebro de diversos animales. Vamos a explicar el peso del cerebro (en g) *cerebro* a partir del peso del cuerpo (en Kg) *cuerpo*.

- Ajusta el modelo y realiza el diagnóstico del modelo.

```

fitCerebro<-lm(cerebro~cuerpo, data=cerebros)
fitCerebro_s<-summary(fitCerebro)
fitCerebro_s

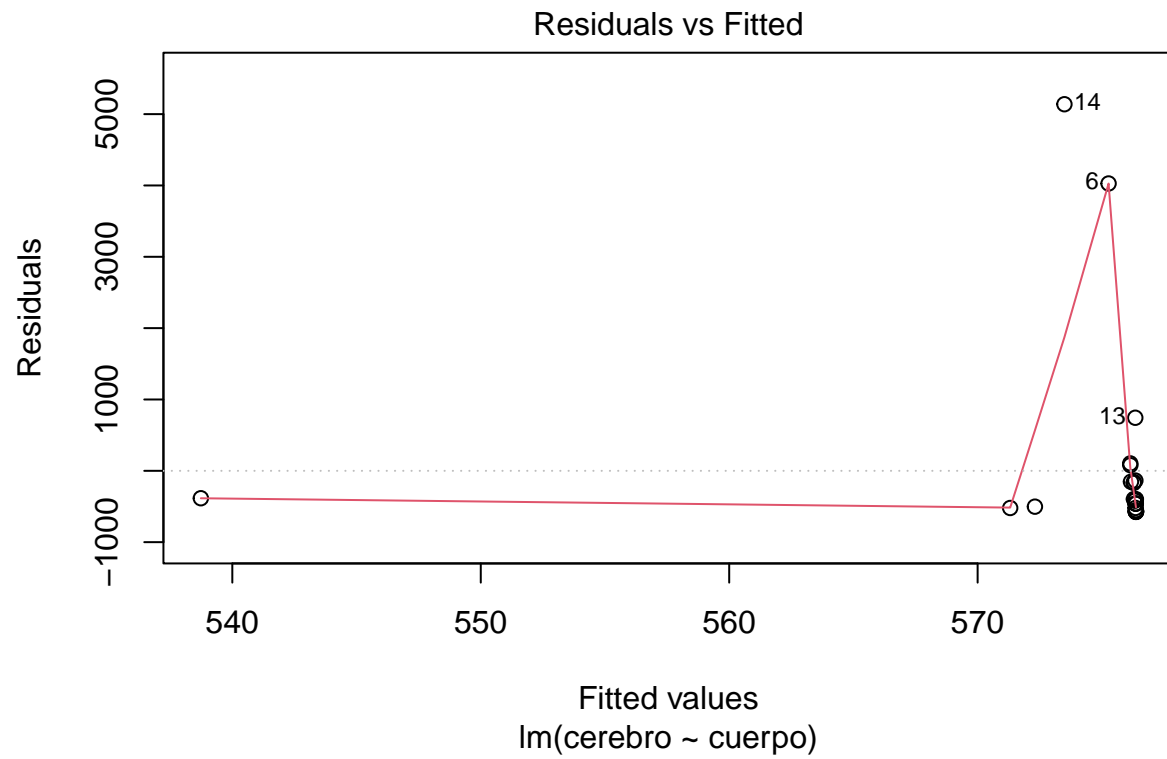
```

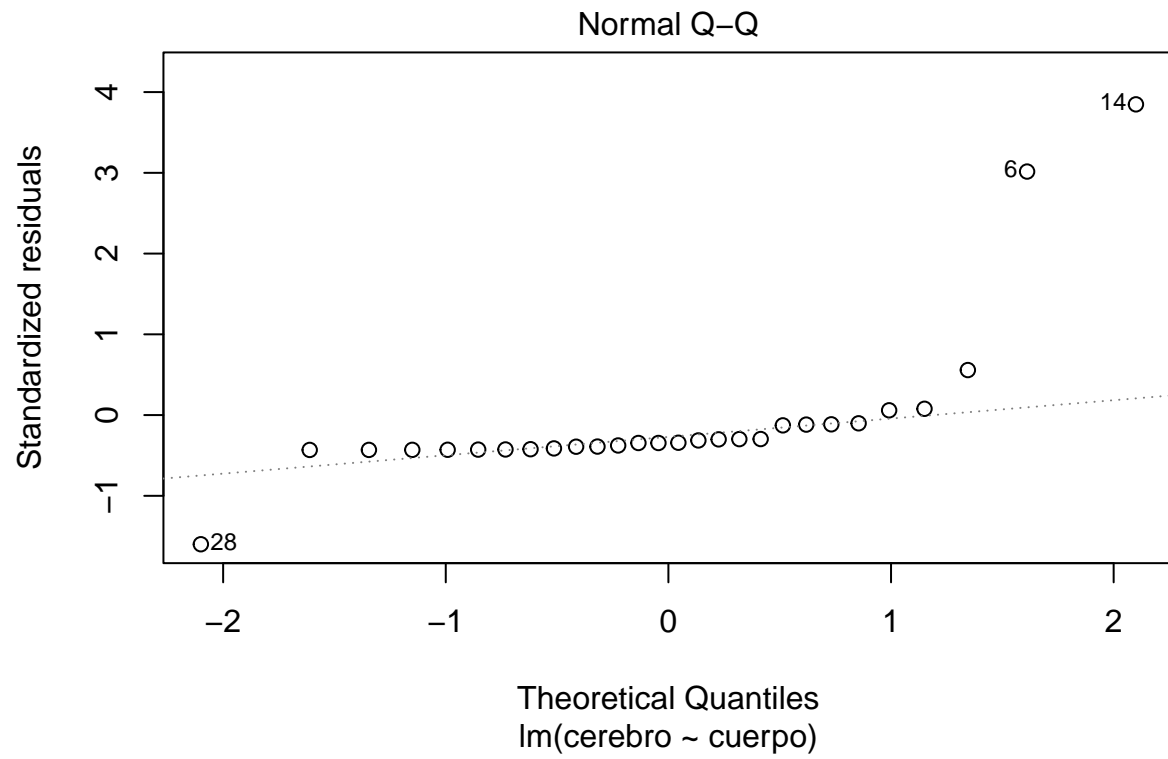
```

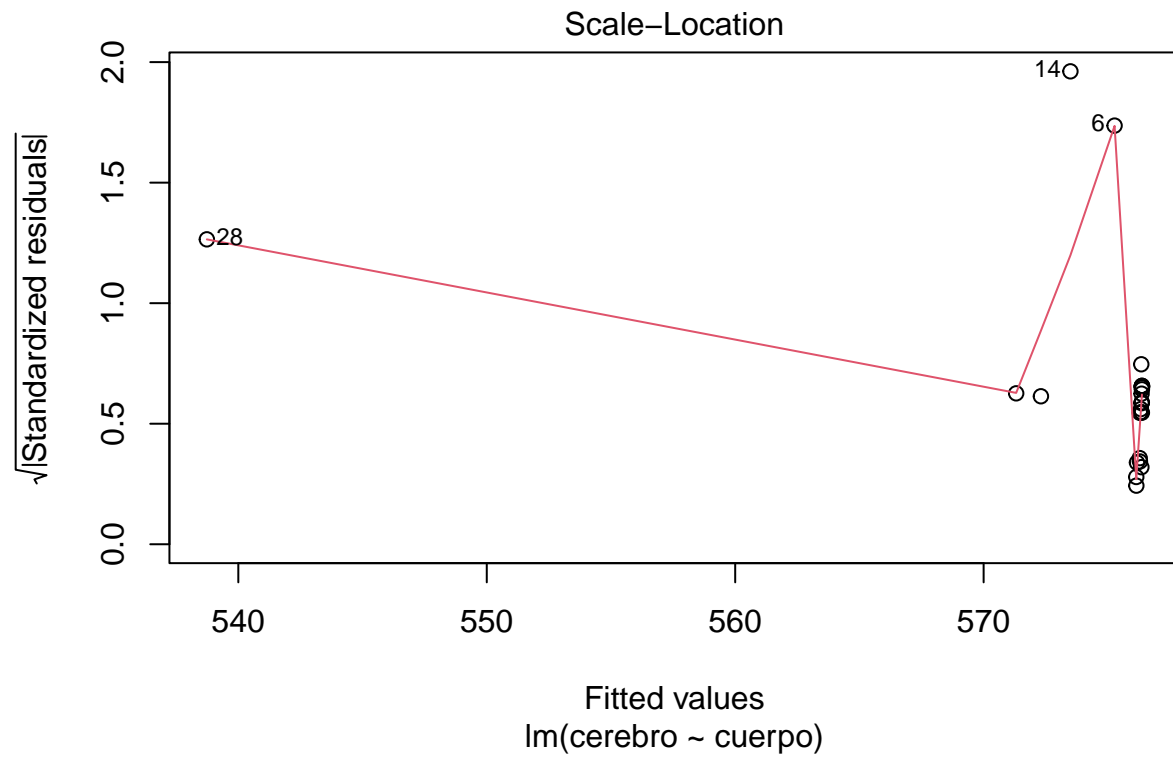
##
## Call:
## lm(formula = cerebro ~ cuerpo, data = cerebros)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -576.0  -554.1  -438.1  -156.3   5138.5
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  5.764e+02  2.659e+02   2.168  0.0395 *
## cuerpo      -4.326e-04  1.589e-02  -0.027  0.9785
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1360 on 26 degrees of freedom

```

```
## Multiple R-squared:  2.853e-05, Adjusted R-squared:  -0.03843
## F-statistic: 0.0007418 on 1 and 26 DF,  p-value: 0.9785
plot(fitCerebro)
```

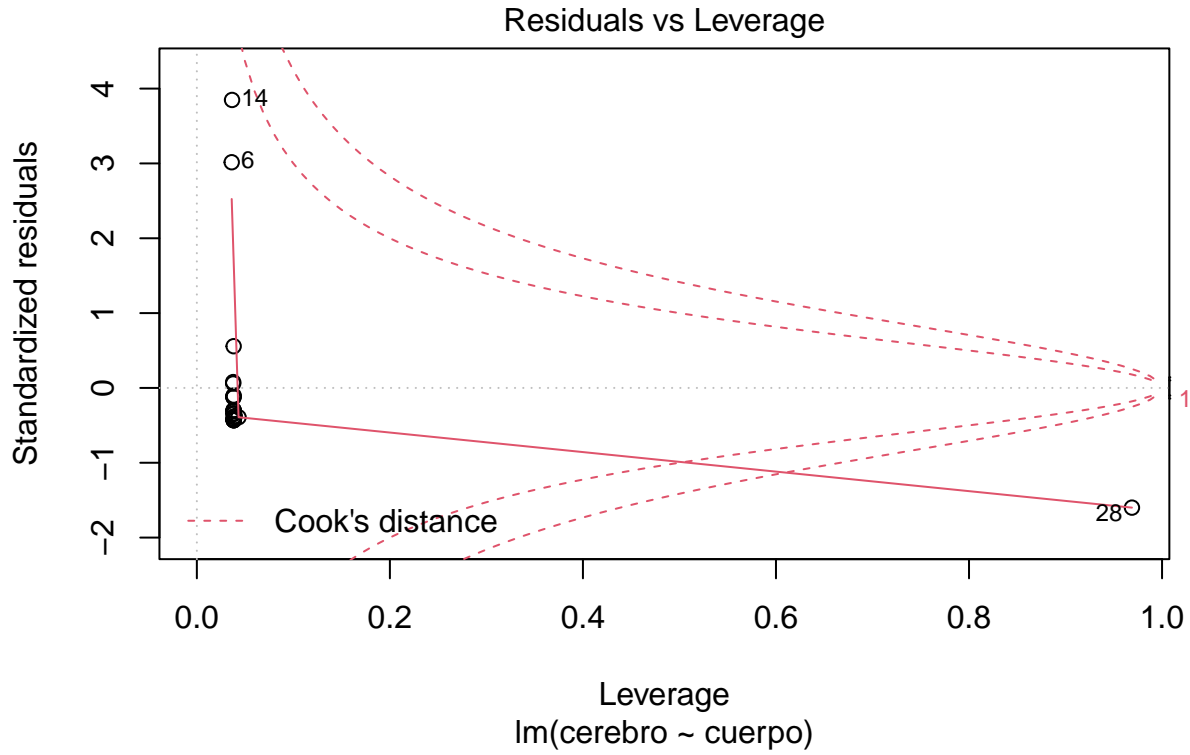






```
## Warning in sqrt(crit * p * (1 - hh)/hh): Se han producido NaNs
```

```
## Warning in sqrt(crit * p * (1 - hh)/hh): Se han producido NaNs
```



Respuesta: lo primero que vemos en el modelo es que el cuerpo no ha resultado significativo, además obtenemos un R^2 bajísimo y los residuos no siguen una distribución normal.

- b) Retira los datos que no pertenecen a la misma población que el resto y re-analiza.

Ejercicio 8

Utilizando los datos de mamíferos, del banco **cerebros**, y las variables en escala logarítmica, dibuja el diagrama de puntos con la recta de mínimos cuadrados. A continuación, analiza gráficamente los residuos ¿crees que el modelo lineal sería adecuado?

Suponiendo adecuado el modelo lineal, contesta a las siguientes preguntas:

- ¿Cuánto vale la pendiente de la recta? ¿Podemos afirmar que es positiva?
- Compara la varianza de la variable respuesta con la varianza de los residuos: ¿Qué porcentaje de la variabilidad inicial está explicado por la recta de mínimos cuadrados? ¿Qué porcentaje de la variabilidad inicial falta todavía por explicar?
- Obtén los intervalos de confianza al 90% sobre los parámetros de la recta.
- Estima el valor de la recta de regresión en el punto $\log(\text{cuerpo}) = 3$ y calcula su intervalo de confianza al 95%. Dibuja el diagrama de dispersión, la recta de regresión y las bandas de confianza al 95% sobre la estimación de la recta.
- Obtén la predicción puntual y por intervalos (al 95%) de un nuevo mamífero con $\log(\text{cuerpo}) = 6$. Añade a la gráfica anterior las bandas de predicción.

Ejercicio 9

El banco de datos **Advertising.csv** relaciona las ventas de ciertos productos con la inversión en publicidad, considerando diversos medios: televisión, radio y periódicos. Aquí vamos a estudiar la variable respuesta *sales* relacionándola con el predictor *TV*.

- a) Obtén un ajuste mediante el método *KNN*, decidiendo el valor de k que consideres adecuado.
- b) Obtén un ajuste mediante el método *loess*, decidiendo el valor de *span* que consideres adecuado.
- c) Dibuja, en la misma gráfica, los dos ajustes anteriores junto con la recta de mínimos cuadrados