

# Actividad 4. Explorando bases

Oscar Gutierrez

2024-08-13

Cargar dataset

```
M= read.csv("mc-donalds-menu.csv")
```

Seleccionar las variables

```
calorias= M$Calories  
proteinas = M$Protein
```

## Analisis de datos atipicos

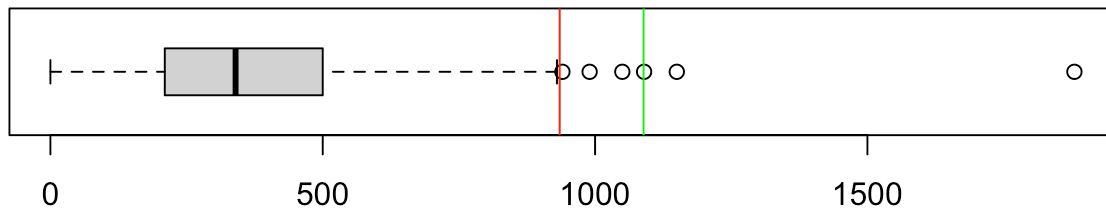
### Calorias

```
q1=quantile(calorias, 0.25)  
q3 = quantile(calorias, 0.75)  
ri=IQR(calorias)    #Rango intercuartílico de X  
par(mfrow=c(2,1))  #Matriz de gráficos de 2x1  
boxplot(calorias,horizontal=TRUE)  #y1=min en la escala del eje Y, y2=máx en la escala  
del eje Y  
abline(v=q3+1.5*ri, col="red")  #linea vertical en el límite de los datos atipicos o ext  
remos  
abline(v= mean(calorias)+ 3*sd(calorias), col="green") # linea vertical a 3 sd de la med  
ia  
X1= M[M$calorias<q3+1.5*ri ]  #En la matriz M, quita datos más allá de 1.5 rangos interc  
uartílicos arriba de q3 de la variable X  
summary(X1)
```

```
## < table of extent 0 x 0 >
```

```
summary(calorias)
```

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	0.0	210.0	340.0	368.3	500.0	1880.0



En este caso, se puede observar que hay varios valores atípicos, sin embargo, es común que haya ciertos alimentos que tengan un alto valor calórico por lo que no es necesario removerlos. Hay datos que están después de 3 desviaciones estándar de la media, por lo que estos serían valores extremos.

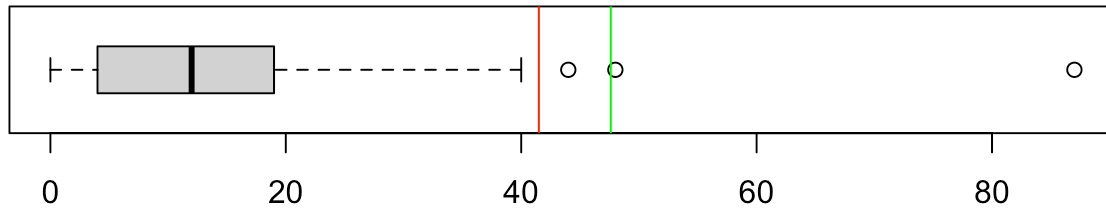
## Proteínas

```
q1=quantile(proteinas, 0.25)
q3 = quantile(proteinas, 0.75)
ri=IQR(proteinas)    #Rango intercuartílico de X
par(mfrow=c(2,1))   #Matriz de gráficos de 2x1
boxplot(proteinas,horizontal=TRUE)  #y1=min en la escala del eje Y, y2=máx en la escala del eje Y
abline(v=q3+1.5*ri, col="red")  #línea vertical en el límite de los datos atípicos o extremos
abline(v= mean(proteinas)+ 3*sd(proteinas), col="green") # línea vertical a 3 sd de la media
X1= M[M$proteinas<q3+1.5*ri ]  #En la matriz M, quita datos más allá de 1.5 rangos intercuartílicos arriba de q3 de la variable X
summary(X1)
```

```
## < table of extent 0 x 0 >
```

```
summary(proteinas)
```

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	0.00	4.00	12.00	13.34	19.00	87.00



Al igual que con las calorías, no es necesario remover los valores atípicos. En este caso también hay valores que se encuentran a 3 desviaciones estándar de la media, estos serían considerados extremos.

## Pruebas de normalidad

###Prueba Anderson Darling para las calorías h0: Los datos siguen una distribución normal h1: Los datos no siguen una distribución normal

```
library(nortest)
ad.test(calorias)
```

```
##
## Anderson-Darling normality test
##
## data: calorias
## A = 2.5088, p-value = 2.369e-06
```

El valor p es muy pequeño, por lo cual los datos no siguen una distribución normal

# Prueba Anderson Darling para proteínas

$H_0$ : Los datos siguen una distribución normal  $H_1$ : Los datos no siguen una distribución normal

```
library(nortest)
ad.test(proteinas)
```

```
##
## Anderson-Darling normality test
##
## data:  proteínas
## A = 4.7515, p-value = 8.515e-12
```

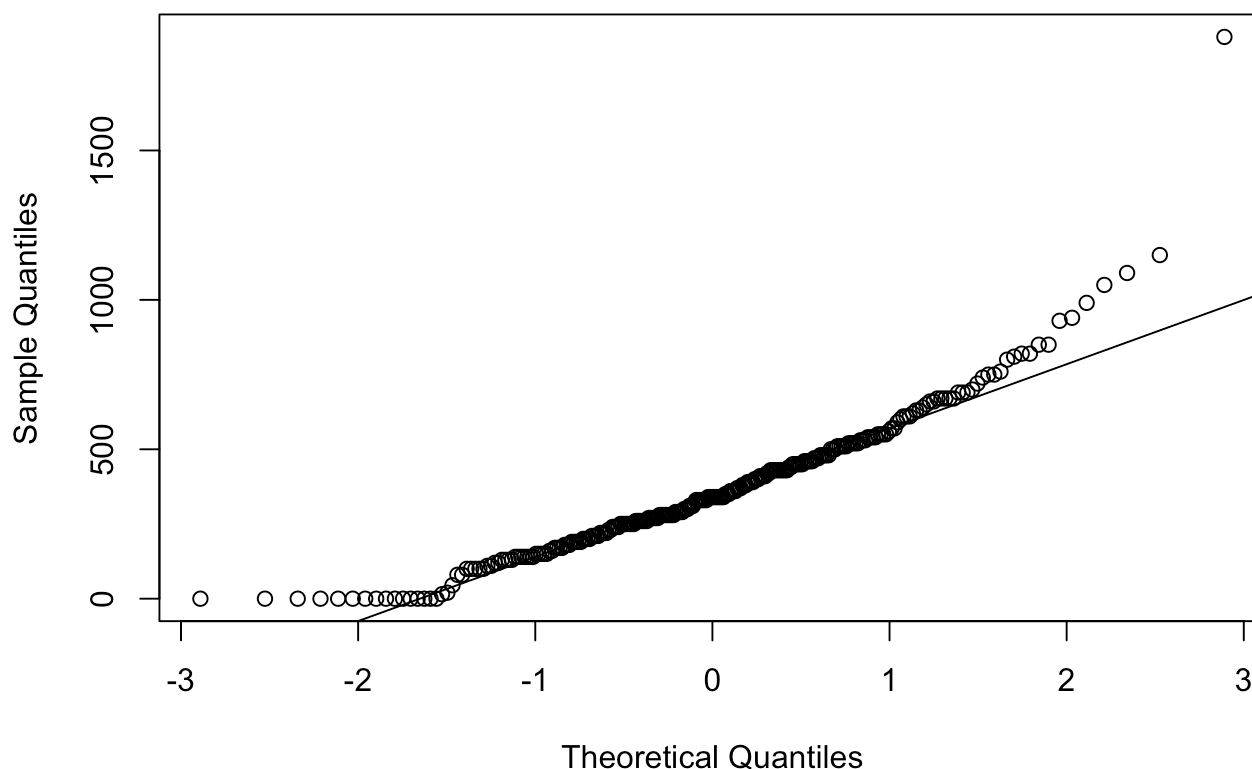
El valor p es muy pequeño, por lo cual las proteínas no siguen una distribución normal.

## QQ plots

### QQ plot para calorías

```
qqnorm(calorias)
qqline(calorias)
```

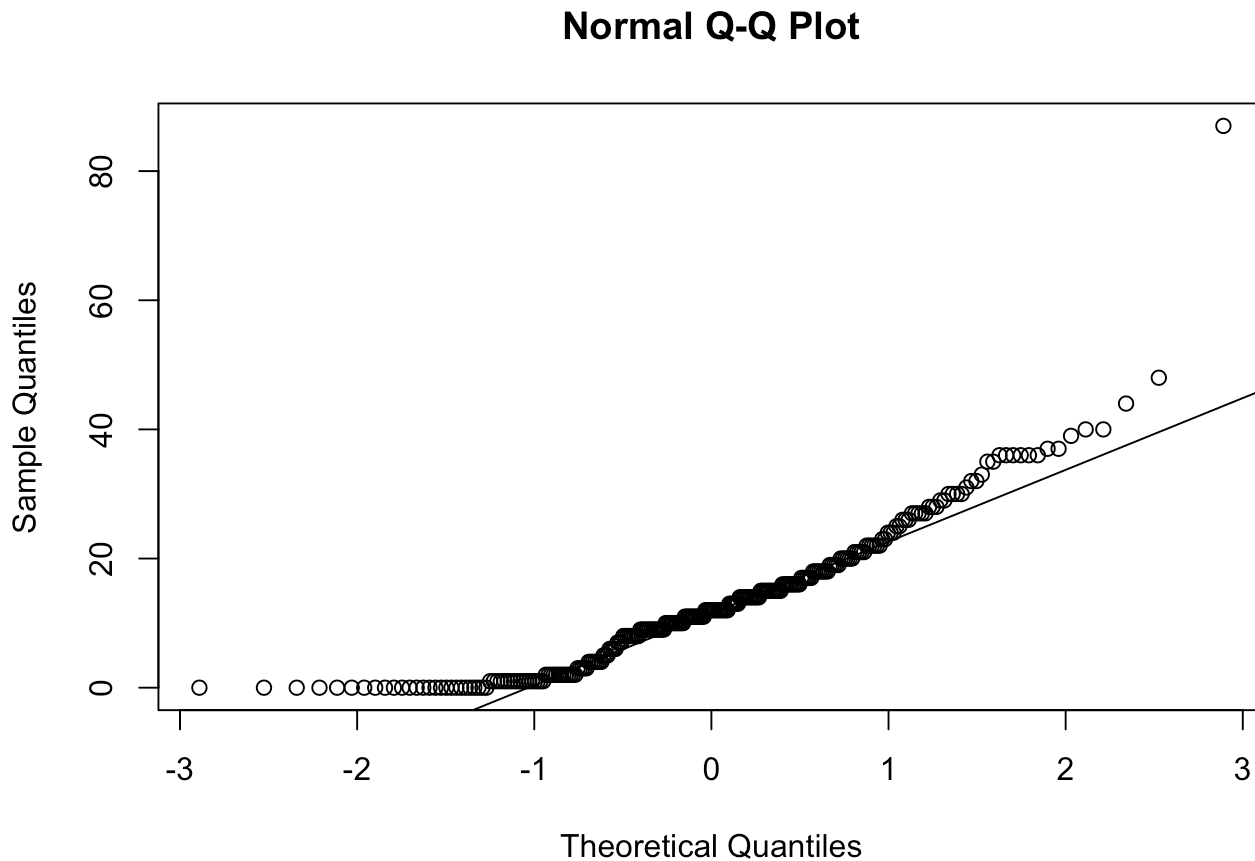
Normal Q-Q Plot



Los datos no siguen la normalidad en las colas.

## ###QQ plots para proteínas

```
qqnorm(proteinas)  
qqline(proteinas)
```



Los datos no cuentan con distribución normal debido a que que no siguen los cuantiles teoricos.

## Coeficientes de sesgo y curtosis

### Coeficientes para calorías

```
library(moments)  
cat("sesgo= ",skewness(calorias), "\ncurtosis=",kurtosis(calorias))
```

```
## sesgo= 1.444105  
## curtosis= 8.645274
```

### Coeficientes para proteínas

```
cat("sesgo= ",skewness(proteinas), "\ncurtosis=",kurtosis(proteinas))
```

```
## sesgo= 1.570794
## curtosis= 8.86355
```

Los coeficientes de sesgo son muy diferentes a 0 y los de curtosis son muy diferentes de 3, por lo que no concuerdan con una distribución normal.

## Calculo de media, mediana y rango medio

### Calorias

```
cat("media=", mean(calorias), "\nmediana=", median(calorias), "\nrango medio=", (max(calorias)-min(calorias))/2)
```

```
## media= 368.2692
## mediana= 340
## rango medio= 940
```

### Proteinas

```
cat("media=", mean(proteinas), "\nmediana=", median(proteinas), "\nrango medio=", (max(proteinas)-min(proteinas))/2)
```

```
## media= 13.33846
## mediana= 12
## rango medio= 43.5
```

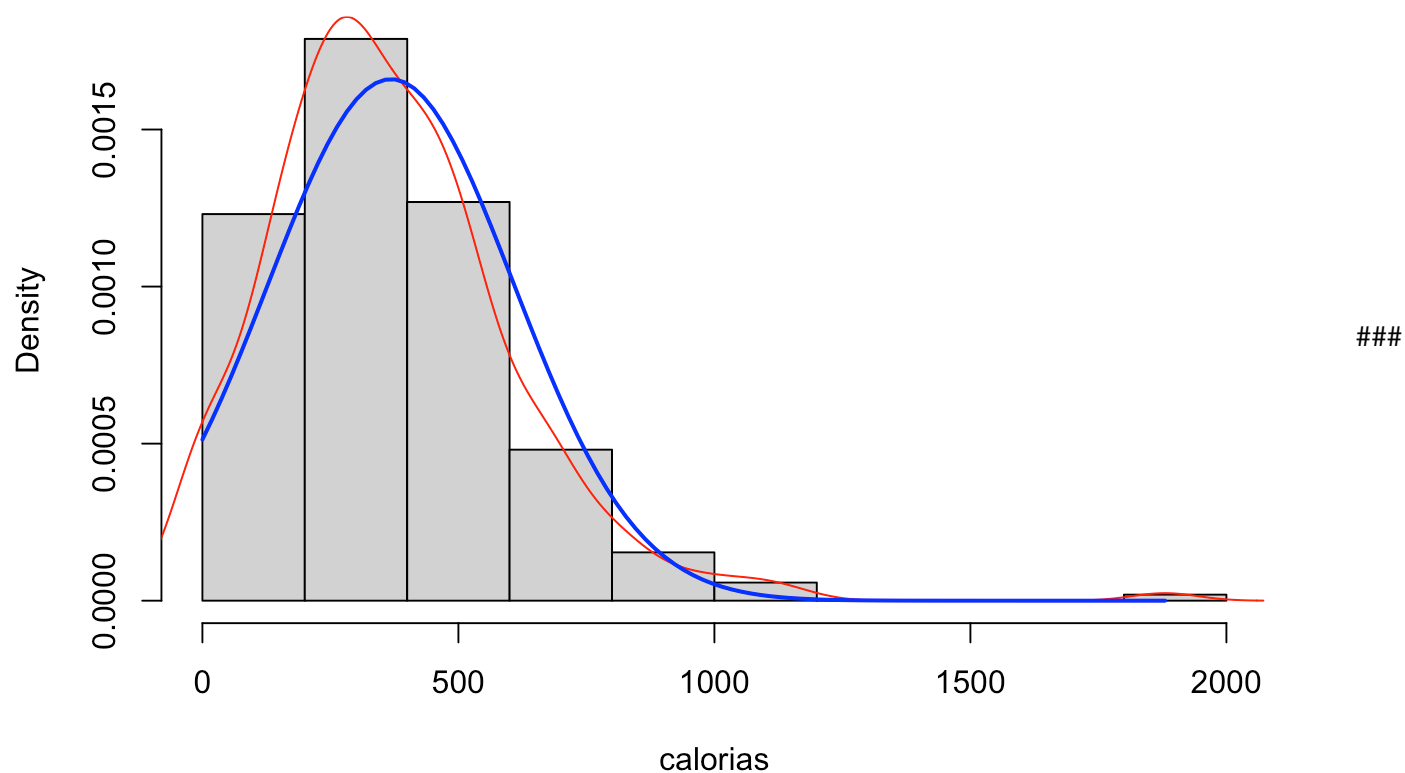
En una distribución normal estos 3 valores deberían ser iguales.

## Histograma y distribución teórica

### Grafica para calorias

```
hist(calorias, freq=FALSE)
lines(density(calorias), col="red")
curve(dnorm(x, mean=mean(calorias), sd=sd(calorias)), from=min(calorias), to= max(calorias), add=TRUE, col="blue", lwd=2)
```

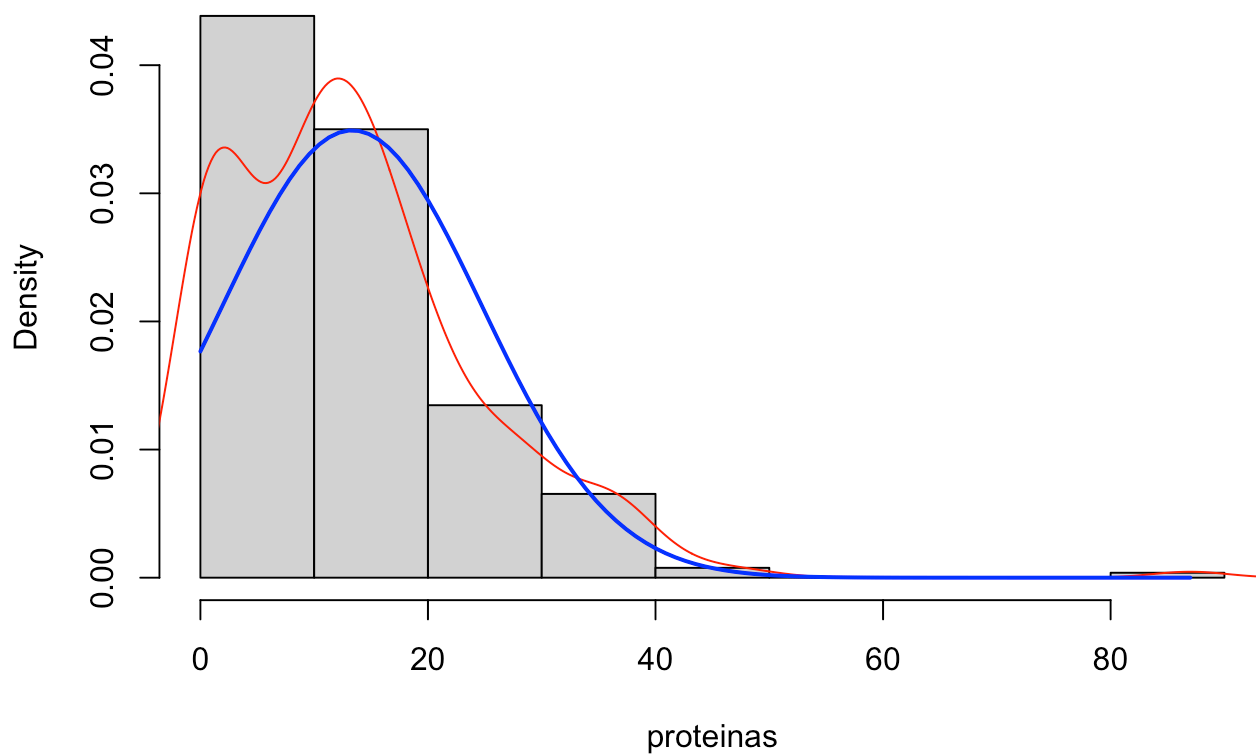
## Histogram of calorías



Grafica para proteínas

```
hist(proteinas,freq=FALSE)
lines(density(proteinas),col="red")
curve(dnorm(x,mean=mean(proteinas),sd=sd(proteinas)), from=min(proteinas), to= max(proteinas), add=TRUE, col="blue",lwd=2)
```

## Histogram of proteínas



Como se puede observar, las distribuciones de probabilidad no concuerdan con una distribución normal teorica. Los datos atipicos pueden influir en la normalidad debido a que pueden generar sesgo.