

```

1 install.packages("dplyr")
2 install.packages("MASS")
3 install.packages("e1071")
4 install.packages("nortest")
5 install.packages("VGAM")

```

Installing package into ‘/usr/local/lib/R/site-library’
(as ‘lib’ is unspecified)

Installing package into ‘/usr/local/lib/R/site-library’
(as ‘lib’ is unspecified)

Installing package into ‘/usr/local/lib/R/site-library’
(as ‘lib’ is unspecified)

also installing the dependency ‘proxy’

Installing package into ‘/usr/local/lib/R/site-library’
(as ‘lib’ is unspecified)

Installing package into ‘/usr/local/lib/R/site-library’
(as ‘lib’ is unspecified)

```

1 library(dplyr)
2 library(ggplot2)
3 library(MASS)
4 library(e1071)
5 library(nortest)

```

No se ha podido completar el guardado automático. Este archivo se ha actualizado de forma remota o en otra pestaña.

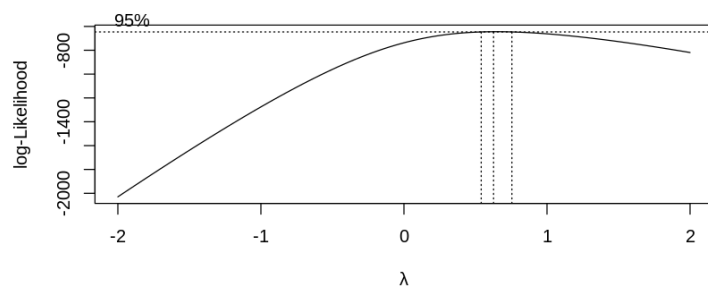
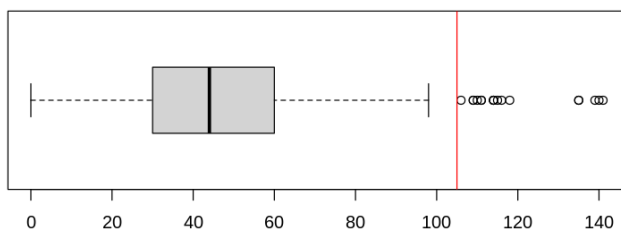
[Mostrar diferencias](#)

```

9 {
10  result = ((data + 1)^lda - 1)/lda
11  return(result) # Return the result
12 }
13
14 # LOADING A DATASET
15 data = read.csv("mc-donalds-menu.csv")
16 cal = data$Carbohydrates
17
18 q1c = quantile(cal, 0.25) # First quartile of the "Calories" variable
19 q3c = quantile(cal, 0.75) # Third quartile of the "Calories" variable
20 ri_c = IQR(cal) # Interquartile range of "Calories" variable
21
22 par(mfrow=c(2,1)) # Create a 2x1 grid for plots
23 boxplot(cal, horizontal=TRUE)
24 abline(v=q3c + 1.5*ri_c, col="red")
25
26 # BOX-COX TRANSFORMATION
27 bc = boxcox((cal+1)~1)
28 lda = bc$x[which.max(bc$y)]
29 cat("Optimal lambda: ",lda,"\n")
30 cat("Best transformation equation: ((x-1)^",lda,"-",lda,")/",lda)

```

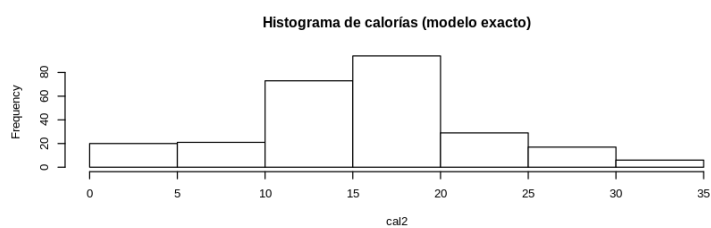
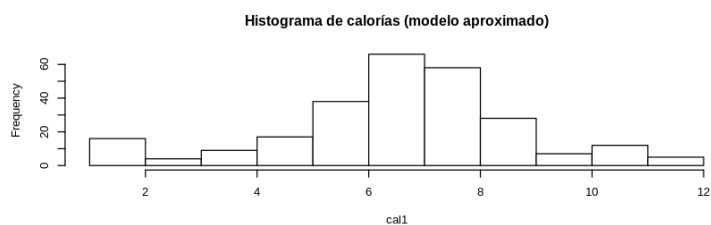
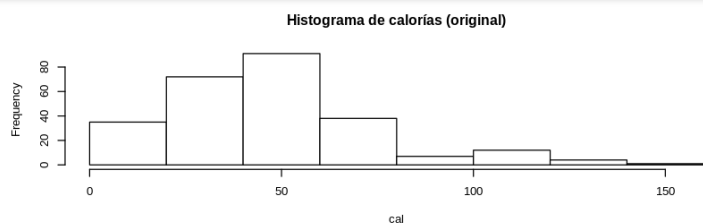
Optimal lambda: 0.6262626
Best transformation equation: $((x-1)^{0.6262626} - 0.6262626) / 0.6262626$



```
1 # HISTOGRAM COMPARATION
2 cal1 = sqrt(cal + 1)
3 cal2 = BoxCox(cal, lda)
4
5 par(mfrow=c(3,1))
6 hist(cal, col=0, main="Histograma de calorías (original)")
7 hist(cal1, col=0, main="Histograma de calorías (modelo aproximado)")
8 hist(cal2, col=0, main="Histograma de calorías (modelo exacto)")
```

No se ha podido completar el guardado automático. Este archivo se ha actualizado de forma remota o en otra pestaña.

[Mostrar diferencias](#)



```
1 # STATSTICAL SUMMARY
```

```

1 # STATISTICAL SUMMARY
2 # Normality test (Anderson-Darling)
3 D0 = ad.test(cal)
4 D1 = ad.test(cal1)
5 D2 = ad.test(cal2)
6
7 # Summary
8 m0 = round(c(as.numeric(summary(cal)),kurtosis(cal),skewness(cal),D0$p.value),5)
9 m1 = round(c(as.numeric(summary(cal1)),kurtosis(cal1),skewness(cal1),D1$p.value),5)
10 m2 = round(c(as.numeric(summary(cal2)),kurtosis(cal2),skewness(cal2),D2$p.value),5)
11
12 # Print results
13 m = as.data.frame(rbind(m0,m1,m2))
14 row.names(m) = c("Original","Primer modelo","Segundo Modelo")
15 names(m) = c("Minimo","Q1","Mediana","Media","Q3","Máximo","Curtosis","Sesgo","Valor p")
16 m
17 print(m$'Valor p')

```

A data.frame: 3 × 9

	Minimo	Q1	Mediana	Media	Q3	Máximo	Curtosis	Sesgo	Valor
	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
Original	0	30.00000	44.00000	47.34615	60.00000	141.00000	1.32408	0.90220	
Primer modelo	1	5.56776	6.70820	6.58324	7.81025	11.91638	0.90923	-0.49396	
Segundo Modelo	0	12.11923	15.72485	15.66877	19.36021	33.97793	0.63820	-0.08250	

[1] 0 0 0

```

1 # REMOVING OUTLIERS AND ZEROS
2 threshold_c = a3c + 1.5*ri_c

```

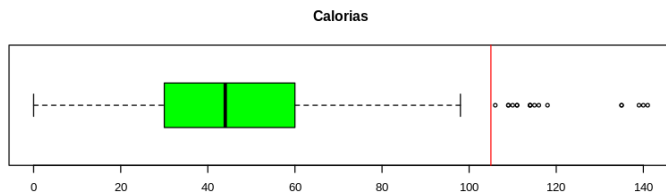
No se ha podido completar el guardado automático. Este archivo se ha actualizado de forma remota o en otra pestaña.

[Mostrar diferencias](#)

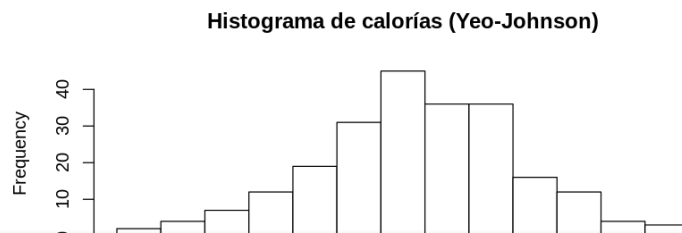
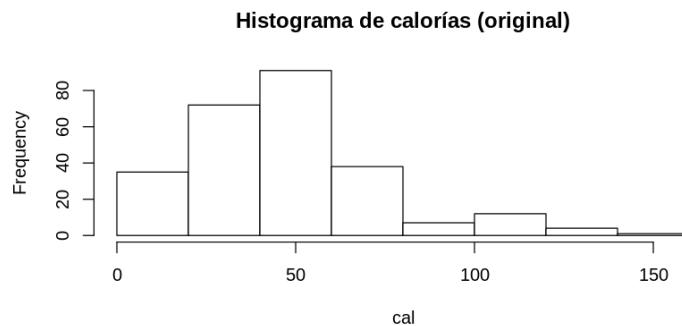
```

5
6 par(mfrow=c(3,1)) # Create a 2x1 grid for plots
7 boxplot(cal, horizontal = TRUE, col="green", main="Calorias")
8 abline(v=threshold_c, col="red")
9 boxplot(cal_NoOutliers, horizontal = TRUE, col="yellow", main="Calorias sin outliers")
10 abline(v=threshold_c, col="red")
11 hist(cal_NoOutliers, col=0, main="Histograma de calorías (no outliers)")

```



```
1 cal3 = yeo.johnson(cal_NoOutliers, lambda = lda)
2 par(mfrow=c(2,1))
3 hist(cal, col=0, main="Histograma de calorías (original)")
4 hist(cal3, col=0, main="Histograma de calorías (Yeo-Johnson)")
```



No se ha podido completar el guardado automático. Este archivo se ha actualizado de forma remota o en otra pestaña.
[Mostrar diferencias](#)

```
1 # STATISTICAL SUMMARY
2 # Normality test (Anderson-Darling)
3 D0 = ad.test(cal)
4 D3 = ad.test(cal3)
5
6 # Summary
7 m0 = round(c(as.numeric(summary(cal)),kurtosis(cal),skewness(cal),D0$p.value),5)
8 m3 = round(c(as.numeric(summary(cal3)),kurtosis(cal3),skewness(cal3),D3$p.value),5)
9
10 # Print results
11 M = as.data.frame(rbind(m0,m3))
12 row.names(M) = c("Original","Modelo Yeo-Johnson")
13 names(M) = c("Mínimo","Q1","Mediana","Media","Q3","Máximo","Curtosis","Sesgo","Valor p")
14 M
15 print(M2$`Valor p`)
```

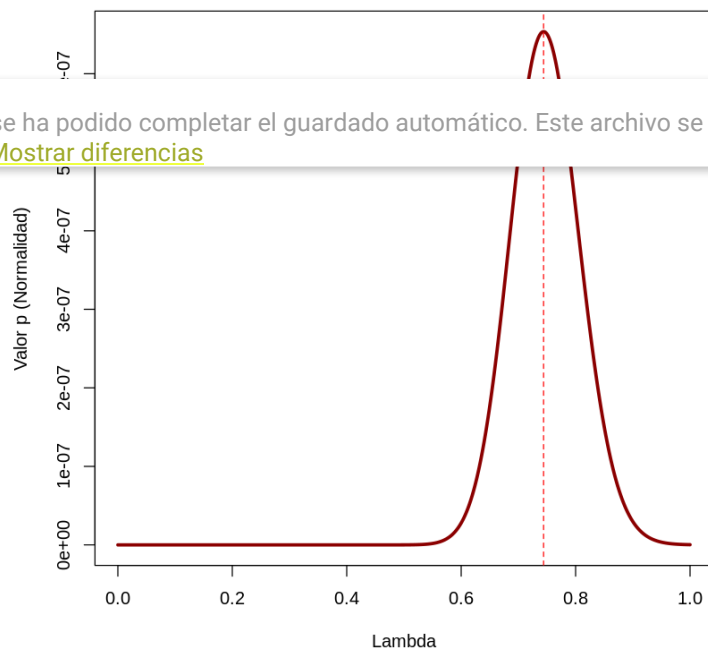
A data.frame: 2 × 9

	Minimo	Q1	Mediana	Media	Q3	Máximo	Curtosis	Sesgo	Val
	✓dh1\	✓dh1\	✓dh1\	✓dh1\	✓dh1\	✓dh1\	✓dh1\	✓dh1\	✓

```
1 lp = seq(0, 1, 0.001) # Proposed lambda values
2 nlp = length(lp)
3 n = length(cal)
4 D = matrix(as.numeric(NA), ncol = 2, nrow = nlp)
5
6 for (i in 1:nlp) {
7   d = yeo.johnson(cal, lambda = lp[i])
8   p = ad.test(d)
9   D[i,] = c(lp[i], p$p.value)
10 }
11
12 N = as.data.frame(D)
13 names(N) <- c("Lambda", "Valor-p")
14 plot(N$Lambda, N$`Valor-p`, type = "l", col = "darkred", lwd = 3,
15 xlab = "Lambda", ylab = "Valor p (Normalidad)", main="Lamda Optimization")
16
17 G = data.frame(subset(N, N$`Valor-p`==max(N$`Valor-p`)))
18 lda2 = G$Lambda
19 abline(v = lda2, col = "red", lty = 2)
20 cat("Optimal lambda: ", lda2, "\n")
```

Optimal lambda: 0.744

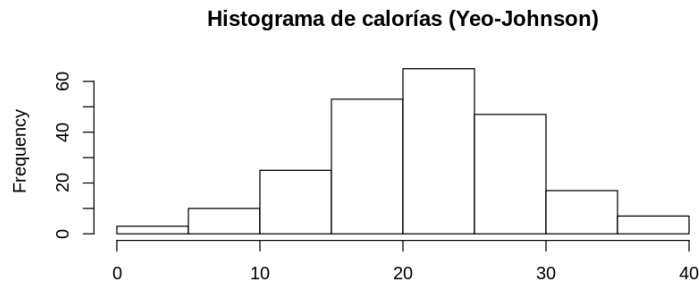
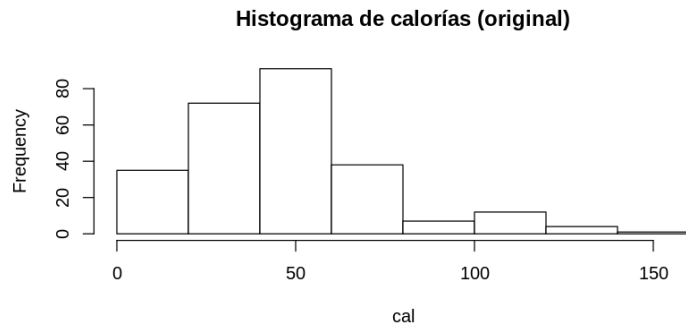
Lamda Optimization



No se ha podido completar el guardado automático. Este archivo se ha actualizado de forma remota o en otra pestaña.

[Mostrar diferencias](#)

```
1 cal3 = yeo.johnson(cal_NoOutliers, lambda = lda2)
2 par(mfrow=c(2,1))
3 hist(cal, col=0, main="Histograma de calorías (original)")
4 hist(cal3, col=0, main="Histograma de calorías (Yeo-Johnson)")
```



```
1 # STATISTICAL SUMMARY
2 # Normality test (Anderson-Darling)
3 D0 = ad.test(cal)
4 D3 = ad.test(cal3)
5
6 # Summary
7 m0 = round(c(as.numeric(summary(cal)),kurtosis(cal),skewness(cal),D0$p.value),5)
8 m3 = round(c(as.numeric(summary(cal3)),kurtosis(cal3),skewness(cal3),D3$p.value),5)
9
10 # Print results
```

No se ha podido completar el guardado automático. Este archivo se ha actualizado de forma remota o en otra pestaña.

[Mostrar diferencias](#)

```
13 names(M2) = c("Mínimo", "Q1", "Mediana", "Media", "Q3", "Máximo", "Curtosis", "Sesgo", "Valor p")
14 M2
15 print(M2$`Valor p`)
```

A data.frame: 2 × 9

	Minimo	Q1	Mediana	Media	Q3	Máximo	Curtosis	Sesgo	Valor p
	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
Original	0.00000	30.00000	44.00000	47.34615	60.00000	141.0000	1.32408	0.90220	0.82772
Modelo Yeo-Johnson 2	3.10695	16.77772	21.48125	21.56740	26.22469	39.6933	-0.01462	-0.03104	0.82772

[1] 0.00000 0.82772

La mejor transformación de los datos de acuerdo a las características de los modelos que encontraste:

- La mejor transformación de acuerdo los valores-p obtenidos para cada distribución fue la transformación Yeo-Johnson, para el caso de los carbohidratos de la base de datos.

Concluye sobre las ventajas y desventajas de los modelos de Box Cox y de Yeo Johnson:

Ventajas de modelo Cox-Box:

- Es muy simple

- Es fácil de interpretar

Desventajas de modelo Cox-Box:

- No acepta valores negativos
- No acepta valores ceros

Ventajas de modelo Yeo-Johnson:

- Acepta valores ceros y negativos
- Es un modelo robusto

Desventajas de modelo Yeo-Johnson:

- Es más complejo de programar
- Se debe estimar los parámetros cuidadosamente

Diferencias entre la transformación y el escalamiento de los datos:

- Cambios en la Distribución: Mientras que las transformaciones pueden hacer que una distribución no normal se aproxime a la normalidad, el escalamiento no cambia la forma de la distribución ni su normalidad.
- Una transformación puede convertir una relación lineal en una relación no lineal con sus variables y viceversa. El escalamiento no afecta la relación entre las variables.
- Las transformaciones como Yeo-Johnson o Box-Cox pueden manejar o no valores negativos y ceros (se requiere hacer un ajustamiento). El escalamiento no modifica la naturaleza de los valores negativos.

Cuándo Utilizar Cada Uno:

No se ha podido completar el guardado automático. Este archivo se ha actualizado de forma remota o en otra pestaña.

[Mostrar diferencias](#)

- Cuando se desee ajustar los datos a una escala determinada y se requieren mantener las relaciones relativas entre las variables, lo mejor será usar un escalamiento de datos.

✓ 0 s completado a las 20:57



No se ha podido completar el guardado automático. Este archivo se ha actualizado de forma remota o en otra pestaña.

[Mostrar diferencias](#)