

Escalamiento no metrico

2024-04-24

Librerias necesarias

```
library(car)

## Loading required package: carData

library(smacof)

## Loading required package: plotrix
## Loading required package: colorspace
## Loading required package: e1071
##
## Attaching package: 'smacof'
## The following object is masked from 'package:base':
##
##      transform

library(cluster)
library(lubridate)

##
## Attaching package: 'lubridate'
## The following objects are masked from 'package:base':
##
##      date, intersect, setdiff, union

source("utilerias/funciones.R")
```

Indicaciones:

Utiliza el siguiente analisis para despues hacer tu propio analisis, seleccionando la mejor dimension con los ratings(columnas 5:18) de la db RockHard del paquete smacof.

Sinapsis

El escalamiento multidimensional no métrico tiene por objetivo preservar las disimilaridades mientras se posicionan los objetos en una menor dimension. Se aplica principalmente sobre datos ordinales o ratings.

Pasos esenciales:

- 0.- Preprocesamiento(Datos, atípicos, escalamiento, transformaciones)
- 1.- Matriz de datos con escala ordinales(Definir rangos)

- 2.- Cálculo de disimilaridades
- 4.- Escalamiento no métrico (Regresión monótona)
- 5.- Mejoras

Carga de la informacion

Práctica sobre el Índice de Rezago Social en Mexico

Variables: Indices

```
data <- read.csv("Indica.csv")
```

Formato Correcto

```
rownames(data) <- data$Entidad.Federativa # Estableciendo como indice las entidades
data$Entidad.Federativa <- NULL # Estableciendo como indice las entidades

# Estableciendo de escalas ordinales
data$Grado.de.rezago.social <- factor(data$Grado.de.rezago.social, levels= c("Muy bajo", "Bajo", "Medio"
```

Selección de las columnas auxiliares y de análisis.

La información debe ser suministrada por el dueño de los datos. Por ejemplo definir el comportamiento del índice de rezago

```
auxiliares <- colnames(data[, c(2,15)])
analisis <- colnames(data[,3:12]) # Selección de columnas
columnas <- c(auxiliares, analisis)
datos <- data[, columnas] # Extracción
```

Escalas Iniciales

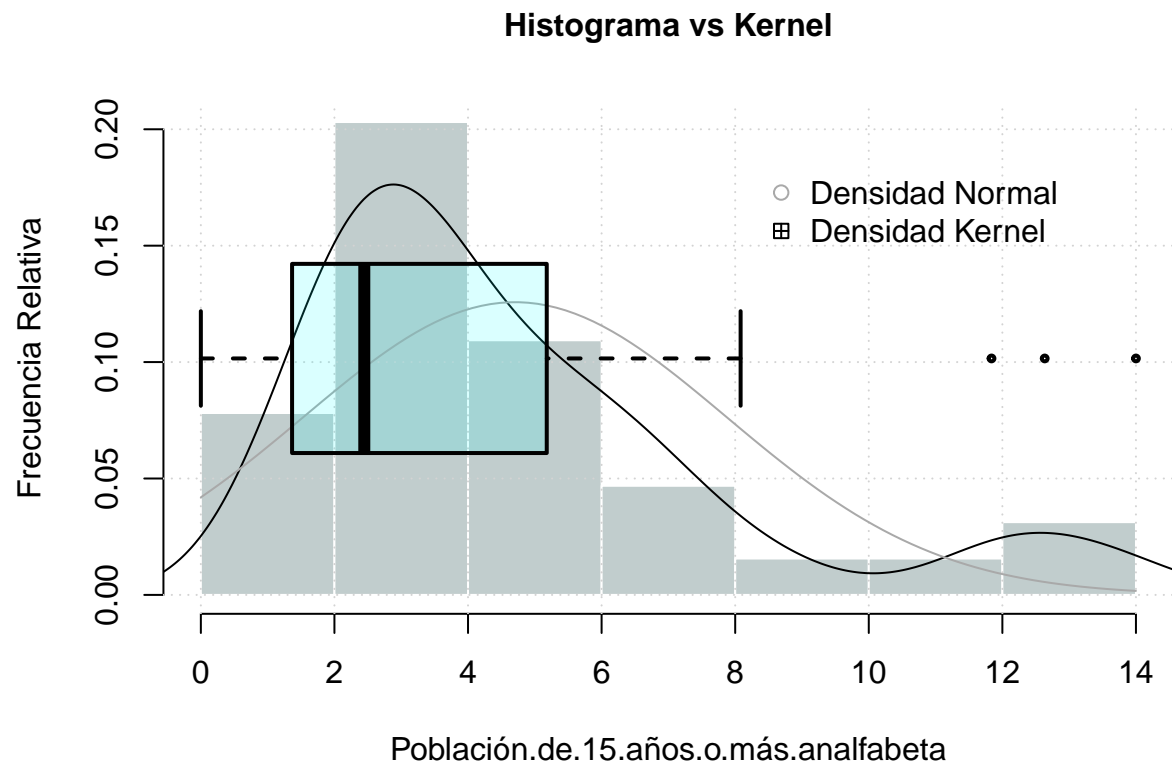
```
# Escalas
tipo <- sapply(datos, class)
continuas <- which(tipo == "numeric") # continuas
enteras <- which(tipo == "integer") # enteras
numericas <- names(c(continuas,enteras))

# Variables Categóricas
nominales <- which( tipo == "factor") # categóricas
ordinales <- which( sapply(datos, is.ordered) ) # ordinales
fecha <- which(tipo == "Date") # Fecha
categoricas <- names(c(nominales, ordinales, fecha))
```

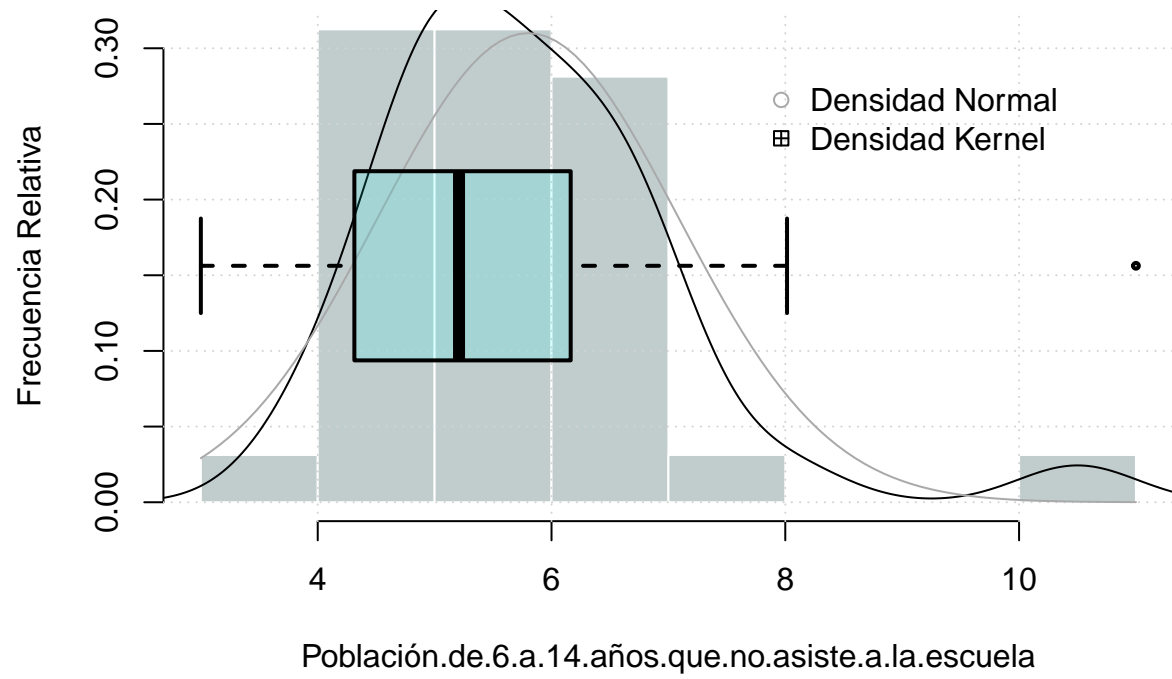
Descriptivos Multivariados

- Identificar Atípicos
- Problemas de escala
- Distribuciones

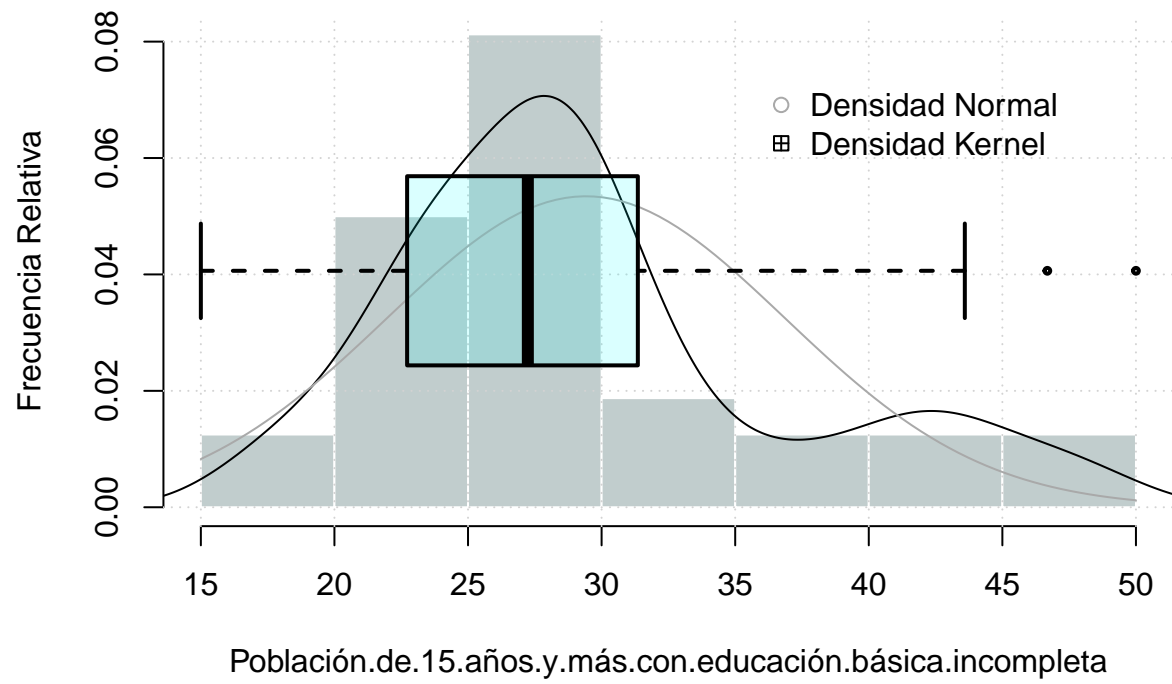
```
# Histogramas
multi.hist(datos[, numericas])
```

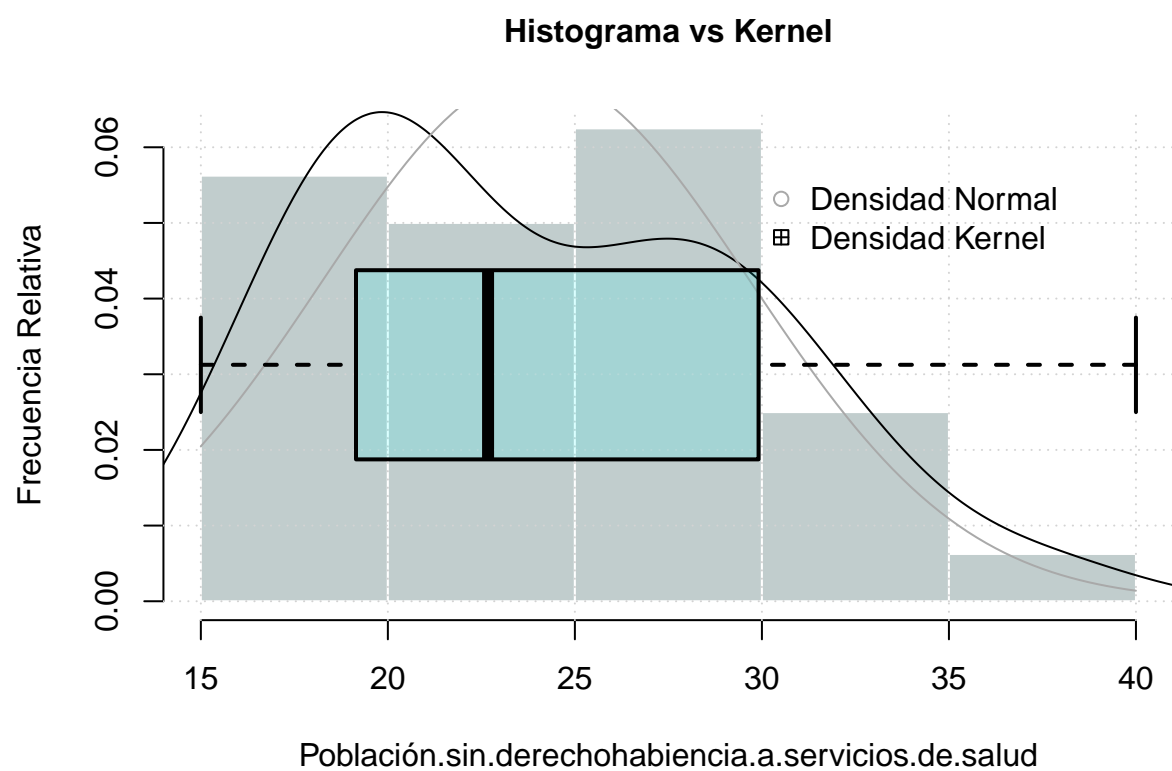


Histograma vs Kernel

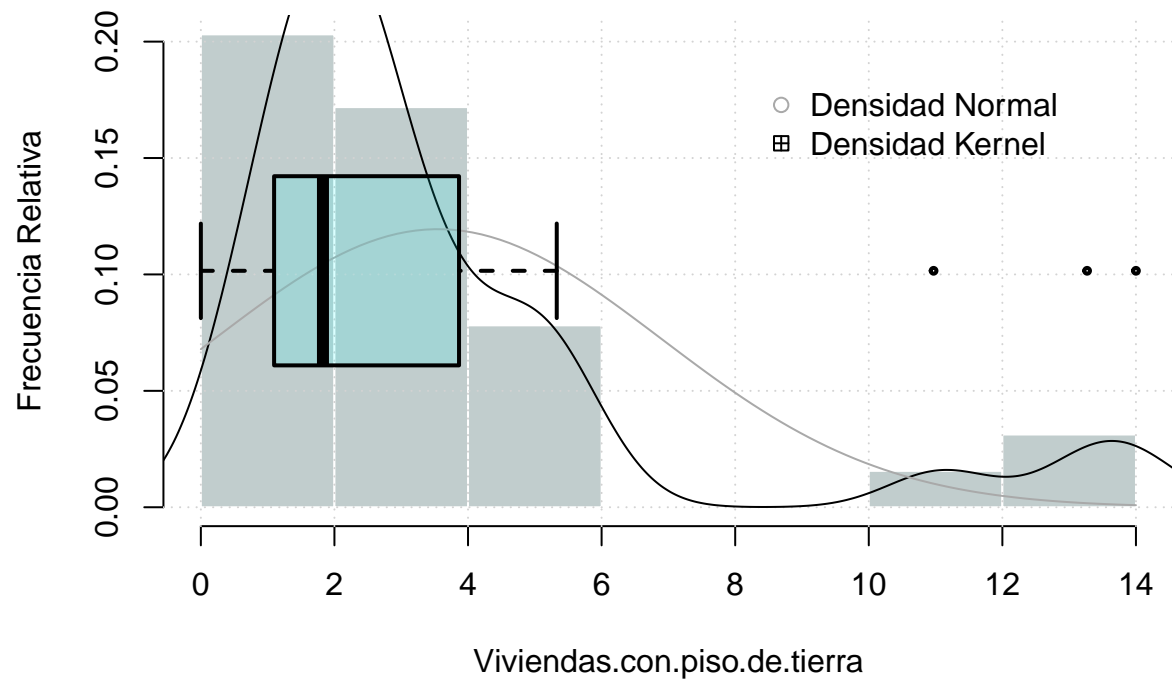


Histograma vs Kernel

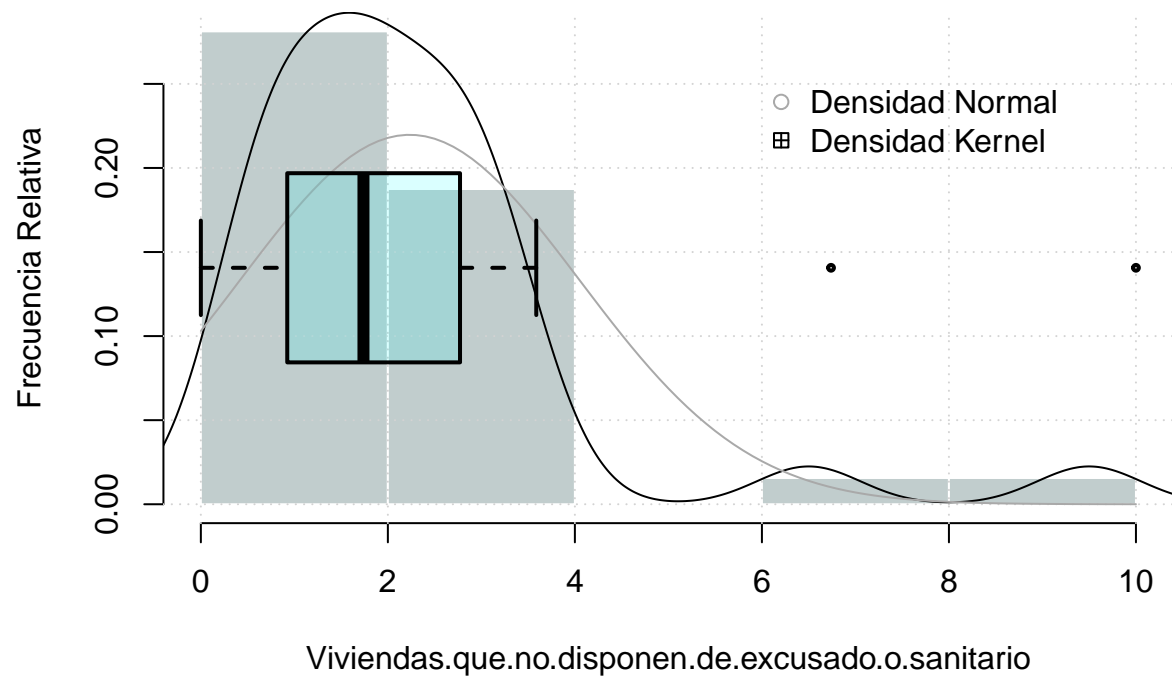




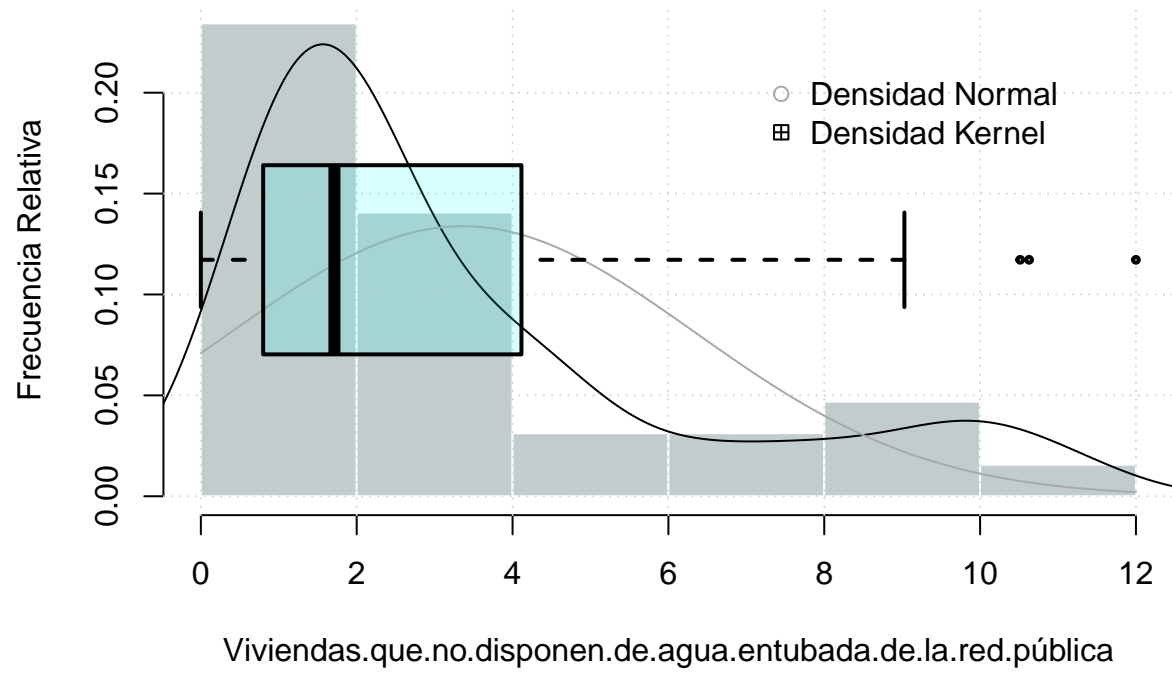
Histograma vs Kernel



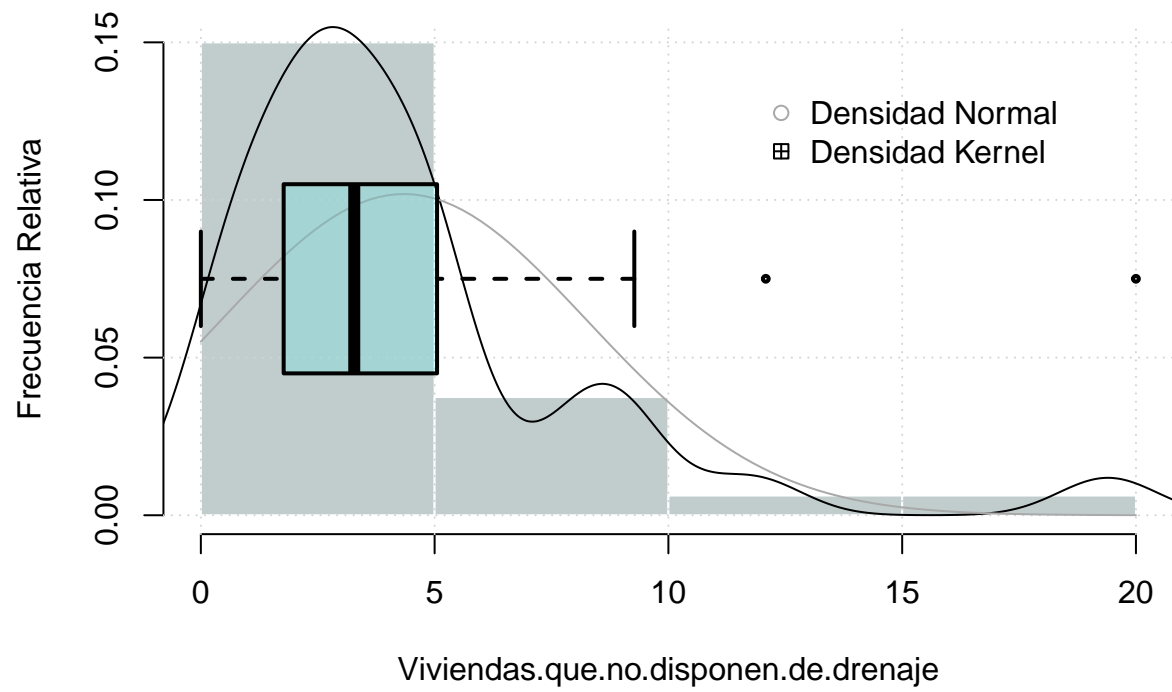
Histograma vs Kernel



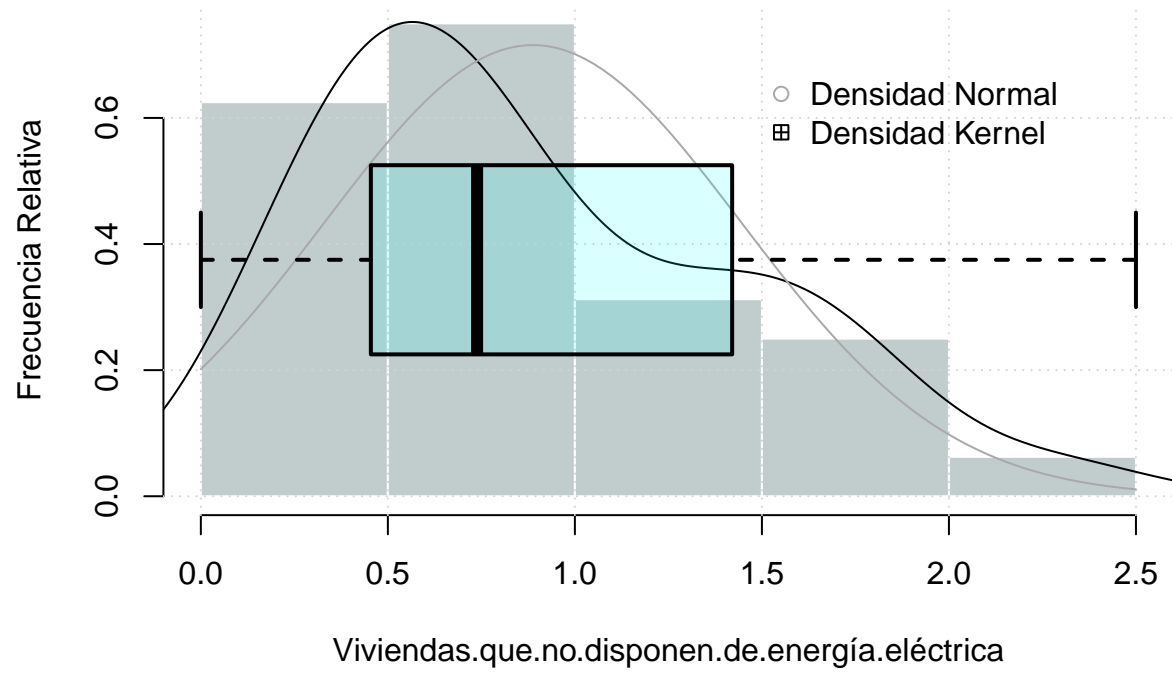
Histograma vs Kernel



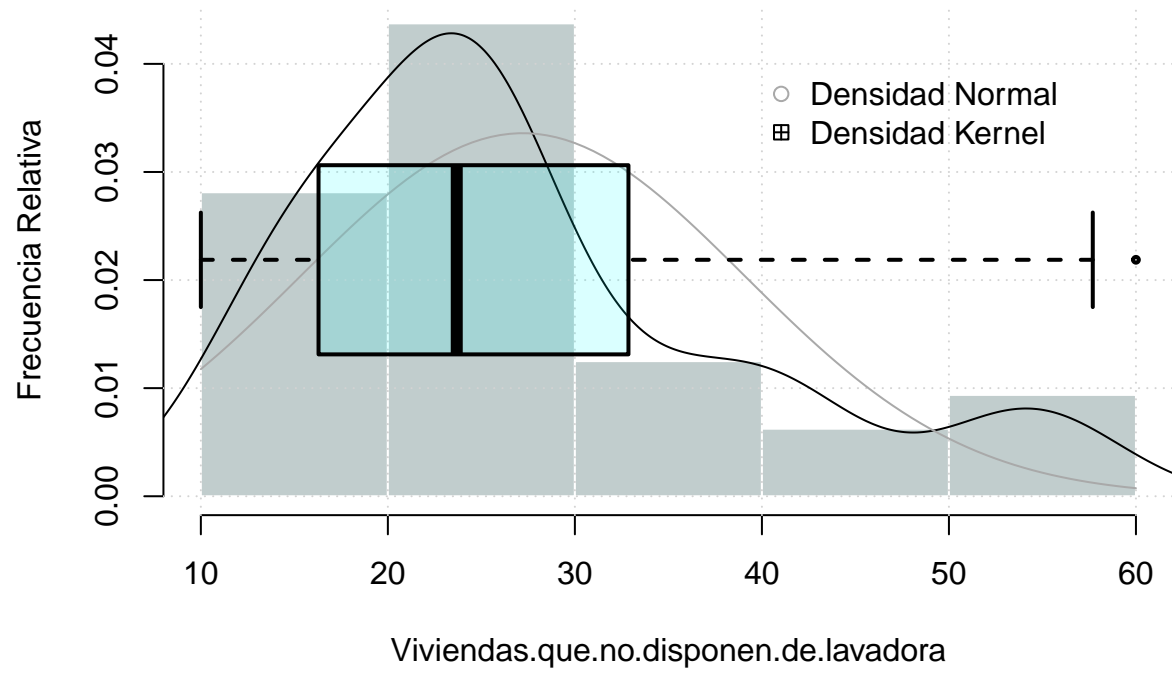
Histograma vs Kernel



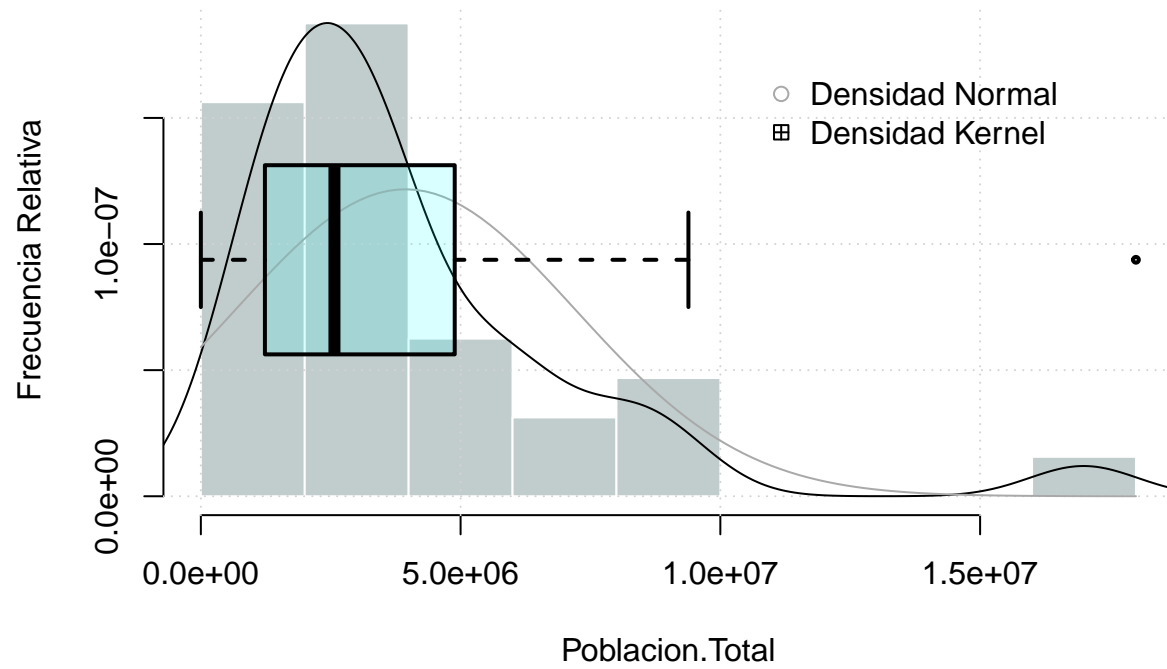
Histograma vs Kernel



Histograma vs Kernel

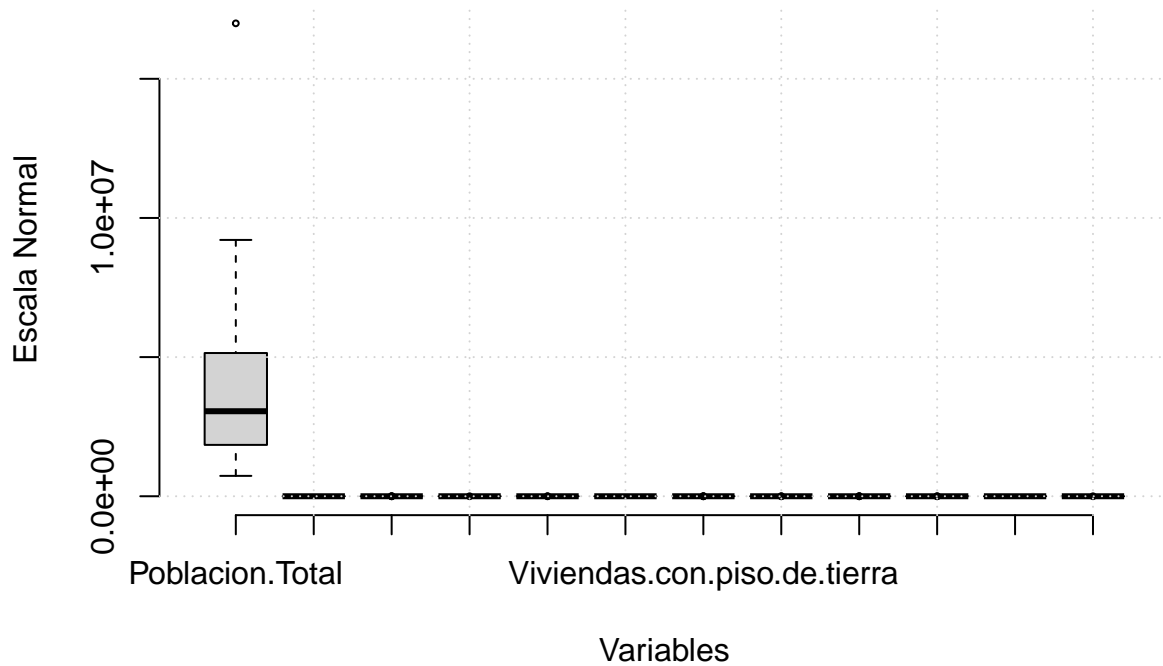


Histograma vs Kernel



```
# Boxplot
boxplot(datos, main="Caja y Bigotes",
        frame = FALSE, xlab="Variables", ylab= "Escala Normal", cex=0.4);grid()
```

Caja y Bigotes

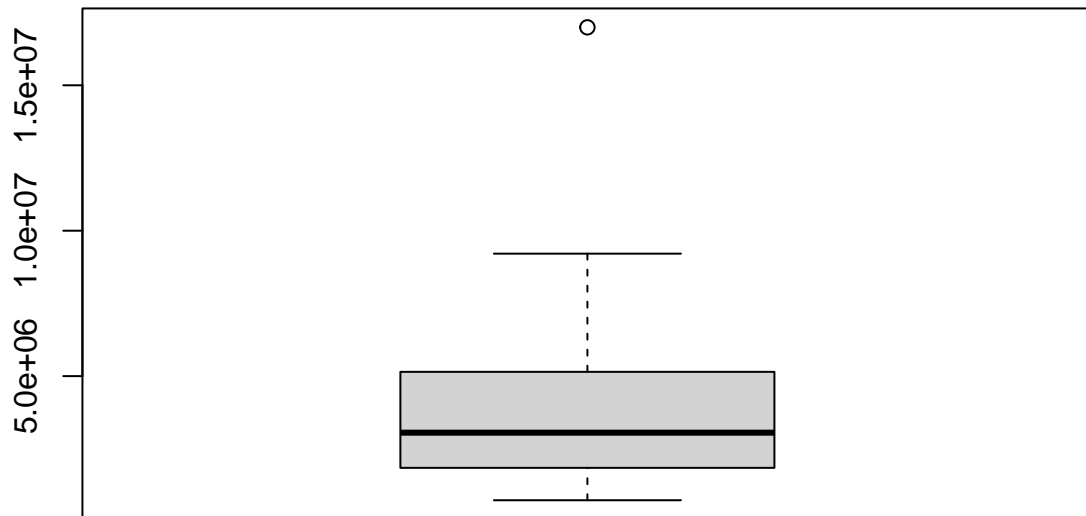


```
# Andrews ##CORREGIR!!!  
#andrews(df = datos, type=2, bty = "n", ylab="f(t)", xlab="t", lwd=1, main="Grafico Andrews" ); grid()
```

Eliminacion de datos atipicos

- Importante ver que la variable auxiliar ayuda a identificar observaciones que afecten el análisis.

```
outliers <- boxplot(datos$Poblacion.Total)$out
```



```
elementos <- which(datos$Poblacion.Total %in% outliers)

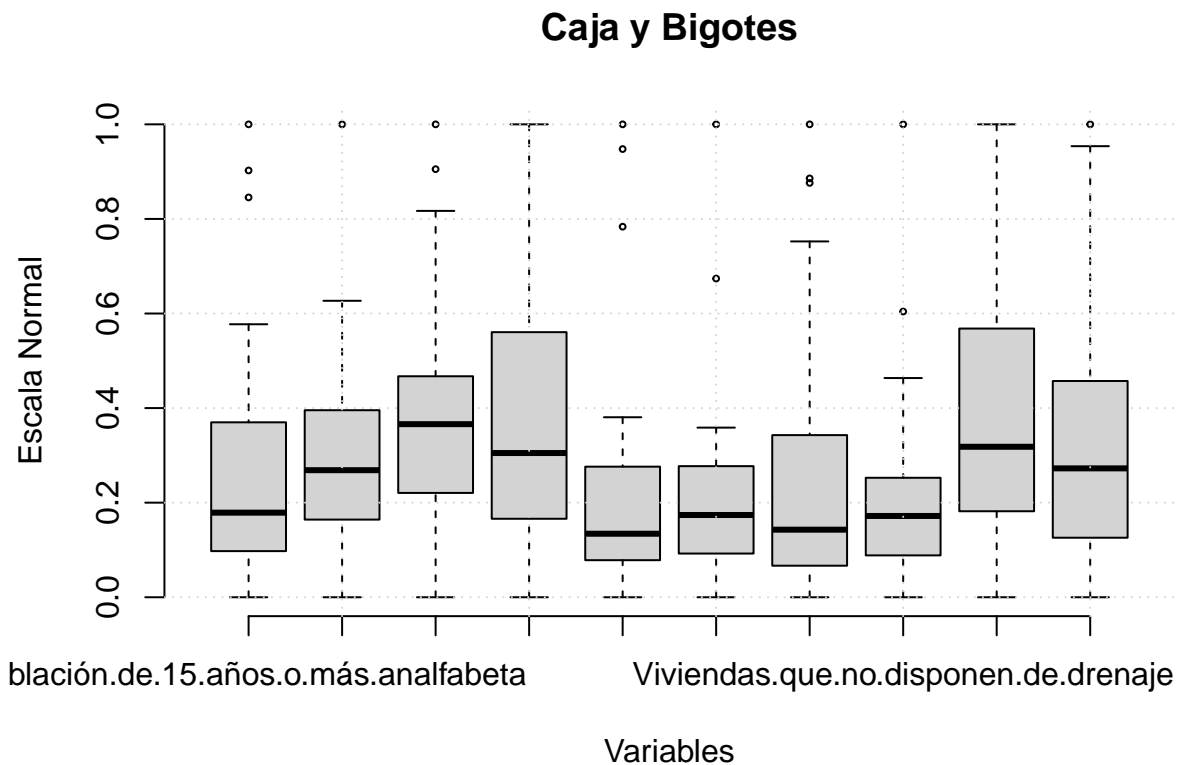
datos <- datos[-union(elementos,elementos), ]
```

Escalamiento

Importante que los índices se recodifican a una escala ordinal, pero primero se normalizan ya que se trata de un índice.

```
# Normalizacion
datos[, analisis] <- sapply(datos[, analisis], function(data){
  (data - min(data)) / (max(data) - min(data))})

# Boxplot
boxplot(datos[, analisis], main="Caja y Bigotes",
        frame = FALSE, xlab="Variables", ylab= "Escala Normal", cex=0.4);grid()
```



1 Matriz de datos con escala ordinales

Definicion de rangos para la escala de lickert

0-20 -> 1 21:40 -> 2 40:60 -> 3 61:80 -> 4 81:100 -> 5

```
# Transformacion a escala ordinal
datos[, analisis] <- datos[, analisis]*100
datos[, analisis] <- round(datos[, analisis])

for(indice in analisis){
  for(n in 1:nrow(datos)){
    datos[n,indice] = recode(datos[n,indice], "0:20=1; 21:40=2; 41:60=3; 61:80=4; 81:100=5")
  }
}

# Formato Correcto
for(indice in analisis){
  datos[, indice] <- factor(datos[, indice], order = TRUE)
}

# Redefinicion de Escalas
tipo <- sapply(datos, class)
continuas <- which(tipo == "numeric") # continuas
enteras <- which(tipo == "integer") # enteras
numericas <- names(c(continuas,enteras))
```



```
# Variables Categoricalas
nominales <- which( tipo == "factor") # categoricas
ordinales <- which( sapply(datos, is.ordered) ) # ordinales
fecha <- which(tipo == "Date") # Fecha
categoricas <- names(c(nominales, ordinales, fecha))
```

Calculo de la matriz de Disimilaridad

- Como las variables son en escala ordinal, ent se utiliza distancia gower(mixtas).

```
gower_dist <- daisy(datos[, analisis], metric = "gower")
```

Escalamiento no métrico

- Métricas de ajuste: stress con valor entre [0,1] y entre mas pequeño mejor. Y rss; entre mas pequeño mejor. En este caso, a prueba y error se encontro que 7 es la mejor dimensión.

```
fit.datos <- smacofSym(gower_dist, type = "ordinal", ndim = 7)
fit.datos$stress
```

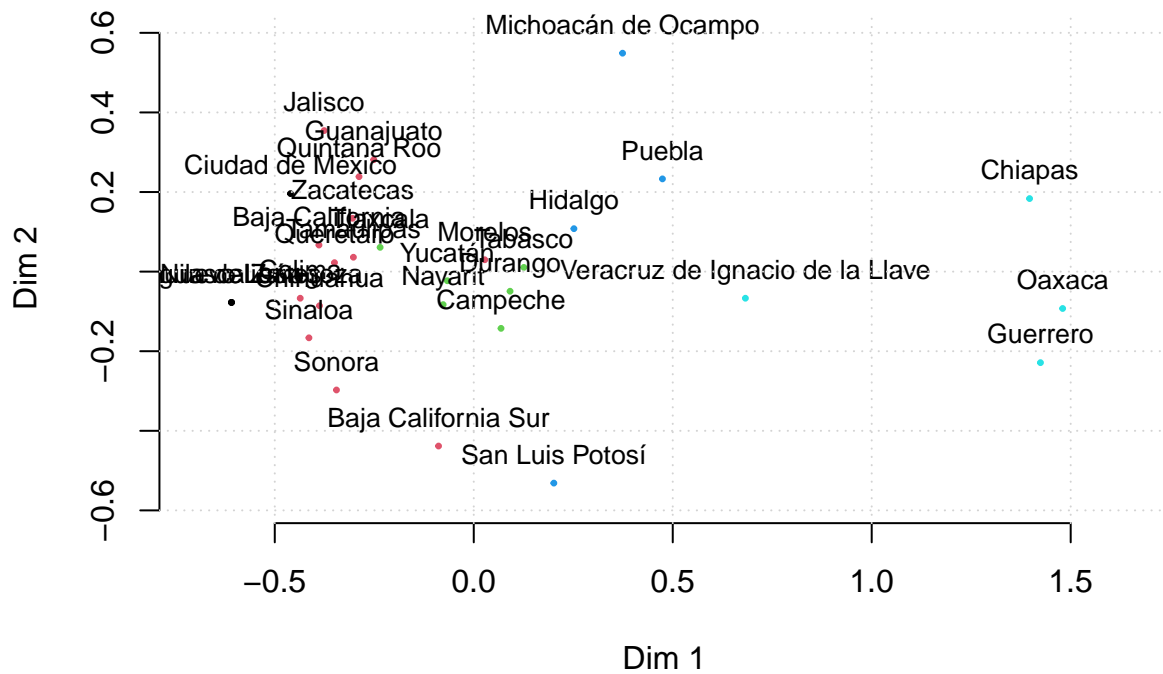
```
## [1] 0.019318
```

```
fit.datos$rss
```

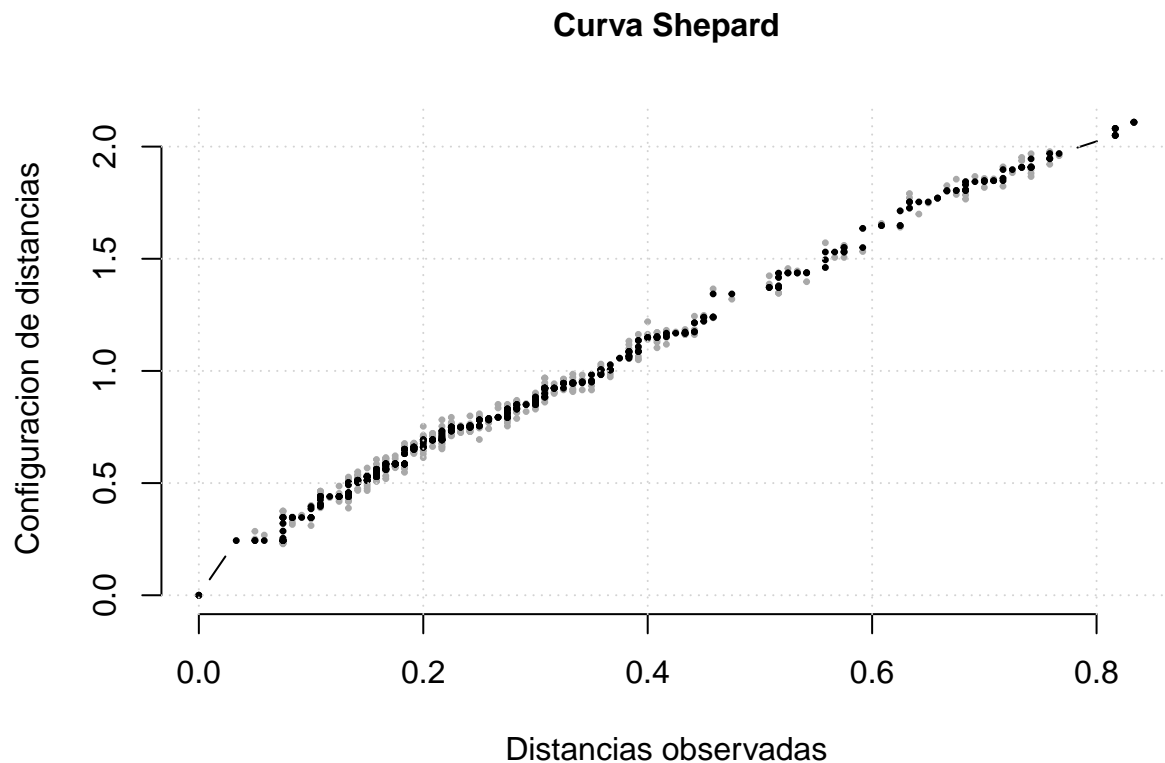
```
## [1] 0.1735311
```

```
# Dispersion
plot(fit.datos, plot.dim = c(1,2), main = "Escalamiento Multidimensional No metrico",
     xlab="Dim 1", ylab="Dim 2", cex=0.5, cex.main=1,
     bty = "n", col = datos$Grado.de.rezago.social );grid()
```

Escalamiento Multidimensional No metrico



```
# Curva Shape
plot(fit.datos, plot.type = "Shepard", main="Curva Shepard",
     xlab="Distancias observadas", ylab="Configuracion de distancias", cex=0.5, cex.main=1,
     col="skyblue", bty = "n");grid()
```



Mejoras

Para mejorar el ajuste, se puede intentar los siguiente:

1. Incrementar el numero de dimensiones(Capturar mayor variabilidad que implica menor rss)
2. Usar otra medida de disimilaridad
3. Usar otro algoritmo de optimizacion para el escalamiento
4. Problemas de preprocesamiento
5. Usar otro metodo como t-sne

Implementación de t-sne

```
from sklearn.manifold import TSNE
import seaborn as sns
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
```

Implementacion

```
# Datos
datos = r.datos
```

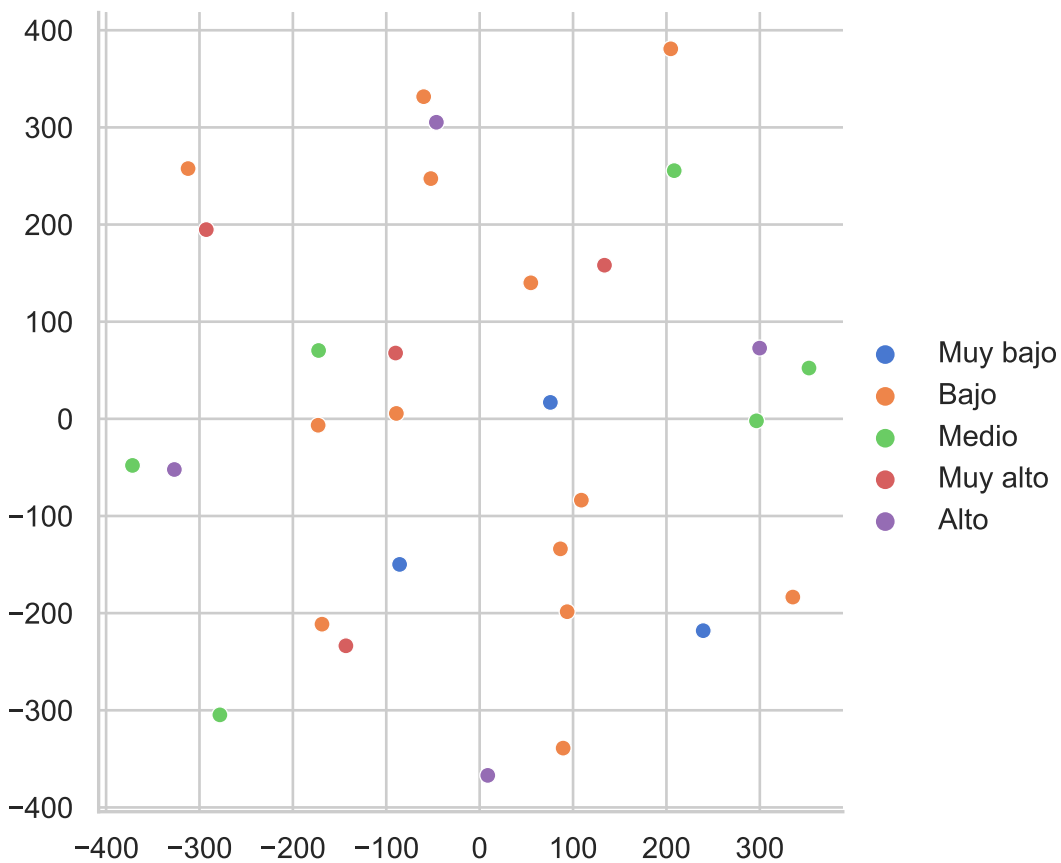
```
# Particion horizontal
x = np.array(datos[r.analisis])
y = np.array(datos[r.auxiliares[1]]) # Variable suplementaria
```

Ajuste

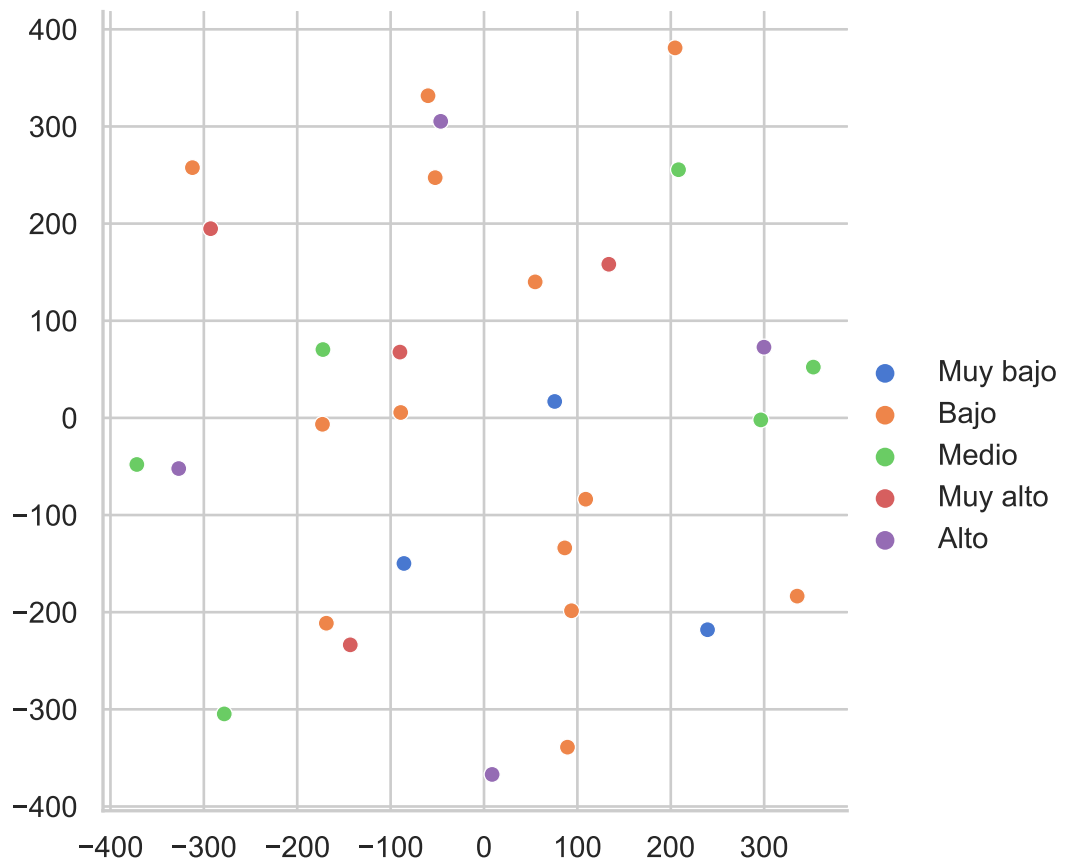
```
x_coord = TSNE(n_components = 3, perplexity = 30, n_iter = 4000).fit_transform(x)
```

Grafico

```
plt.clf()
sns.set(style="whitegrid")
sns.relplot(x=x_coord[:,0], y=x_coord[:,1], hue=y, palette="muted" )
```



```
plt.show()
```



```
exit
```

```
## Use exit() or Ctrl-Z plus Return to exit
```

