



Universidad Rafael Landívar
Segundo Ciclo 2021
Facultad de Ingeniería
Análisis de Datos
Ing. Dan Stanly Bolaños

Documentación Técnica Examen Final

Andrés Gálvez [1024718]
Alexander Villatoro [1182118]
Luis Chuta [1320016]
Sergio Lara [1044418]

Ciudad de Guatemala, Martes 23 de Noviembre de 2021

Contenido

Descripción de objetos	3
Objetos visuales de PowerBI:.....	3
Clustering	3
Modelos de Clasificación:	5
Árbol de decisión:	6
Regresión Lineal	7
Lineas de Tiempo:	9

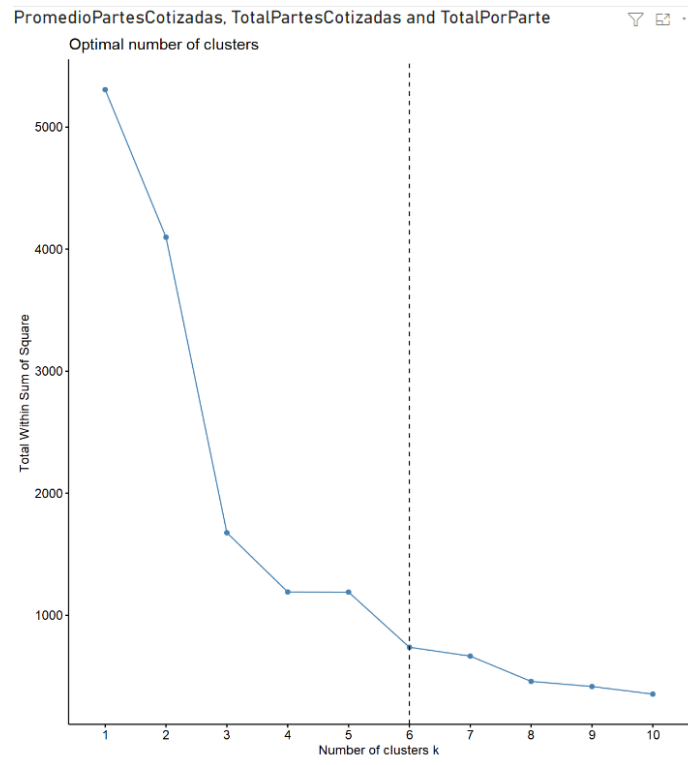
Descripción de objetos

Objetos visuales de PowerBI:

Para la presentación de los gráficos requeridos por el enunciado del proyecto, estos se hicieron en base a los laboratorios hechos en clase. Así que, para presentar, los resultados de nuestros modelos de predicción, clasificación, agrupamiento y líneas de tiempo para los cuales se crearon las siguientes graficas.

Clustering

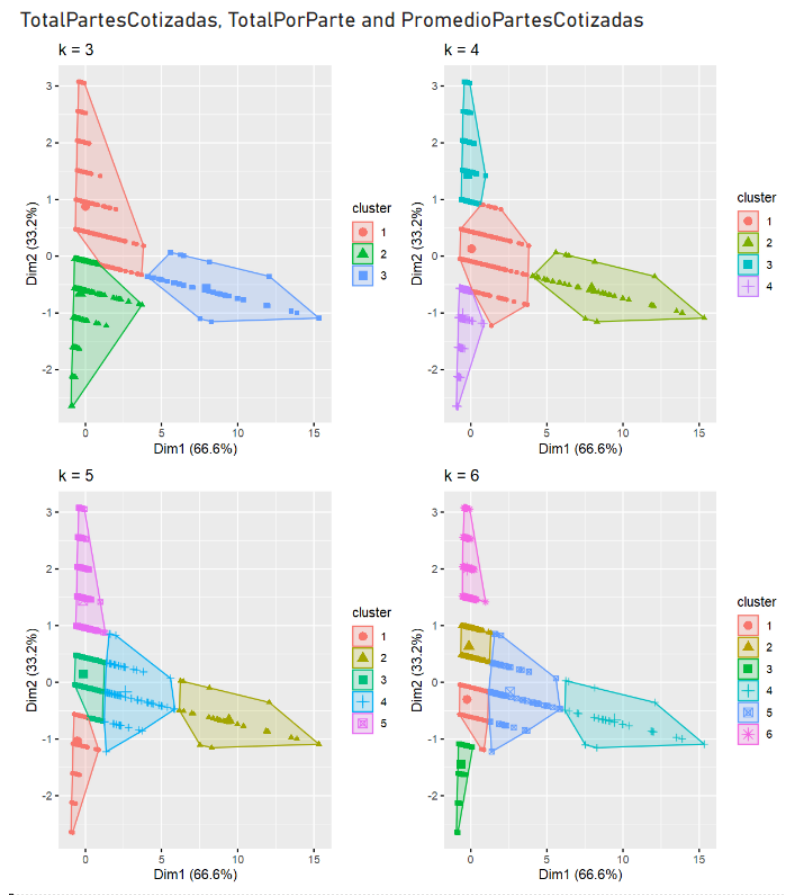
Gráfica WSS



Código

```
1 # The following code to create a dataframe and remove duplicated rows is always executed and acts as a preamble for your script.
2
3 # dataset <- data.frame(PromedioPartesCotizadas, TotalPartesCotizadas, TotalPorParte)
4 # dataset <- unique(dataset)
5
6 # Paste or type your script code here:
7
8
9 library(ggplot2)
10 library(dplyr)
11 library(DBI)
12 library(odbc)
13
14 rownames(dataset) <- dataset$Nombre
15 Repuestos <- scale(dataset)
16 Repuestos <- na.omit(Repuestos) #quitar los registros que estan en null
17
18
19 #install.packages("factoextra") #paquete para graficar
20 library(factoextra) #libreria de paquete instalado
21
22
23 distancia <- get_dist(Repuestos)
24 set.seed(123) #seed permitira fijar un id para generar valores random
25
26 fviz_nbclust(Repuestos, kmeans, method = "wss") +
27   geom_vline(xintercept = 6, linetype = 2)
```

Gráfica de K-means:



Codigo K-means

```
# The following code to create a dataframe and remove duplicated rows is always executed and acts as a preamble for your script:

# dataset <- data.frame(PromedioPartesCotizadas, TotalPartesCotizadas, TotalPorParte)
# dataset <- unique(dataset)

# Paste or type your script code here:

library(ggplot2)
library(dplyr)
library(DBI)
library(odbc)

rownames(dataset) <- dataset$Nombre
Repuestos <- scale(dataset)
Repuestos <- na.omit(Repuestos) #quitar los registros que estan en null

#install.packages("factoextra") #paquete para graficar
library(factoextra) #libreria de paquete instalado

distancia <- get_dist(Repuestos)
set.seed(123) #seed permitira fijar un id para generar valores random

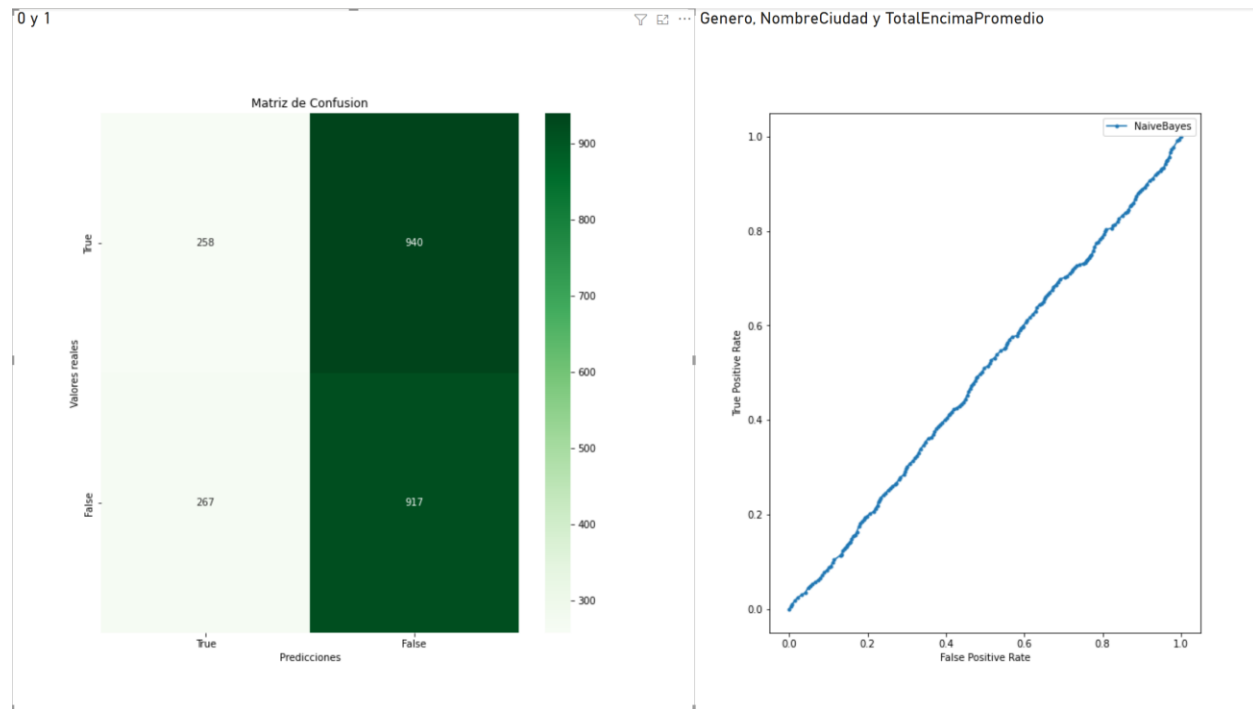
clusterk3 <- kmeans(Repuestos, 3, nstart = 25)
clusterk4 <- kmeans(Repuestos, 4, nstart = 25)
clusterk5 <- kmeans(Repuestos, 5, nstart = 25)
clusterk6 <- kmeans(Repuestos, 6, nstart = 25)
clusterk7 <- kmeans(Repuestos, 7, nstart = 25)

#grafica1 <- fviz_cluster(clusterk2, geom = "point", data = Repuestos) + ggtitle("k = 2")
grafica2 <- fviz_cluster(clusterk3, geom = "point", data = Repuestos) + ggtitle("k = 3")
grafica3 <- fviz_cluster(clusterk4, geom = "point", data = Repuestos) + ggtitle("k = 4")
grafica4 <- fviz_cluster(clusterk5, geom = "point", data = Repuestos) + ggtitle("k = 5")
grafica5 <- fviz_cluster(clusterk6, geom = "point", data = Repuestos) + ggtitle("k = 6")
#grafica6 <- fviz_cluster(clusterk7, geom = "point", data = Repuestos) + ggtitle("k = 7")

library(gridExtra)
grid.arrange(grafica2, grafica3, grafica4, grafica5, nrow = 2)
```

Modelos de Clasificación:

Gráfica:



Código Heatmap:

```

Editor de scripts de Python
⚠ Las filas duplicadas se quitarán de los datos.
1 # El código siguiente, que crea un dataframe y quita las filas duplicadas, siempre se ejecuta y actúa como un preámbulo del script:
2
3 # dataset = pandas.DataFrame(0, 1)
4 # dataset = dataset.drop_duplicates()
5
6 # Pegue o escriba aquí el código de script:
7 import seaborn as sns
8 import matplotlib.pyplot as plt
9 ax = plt.subplot()
10 sns.heatmap(dataset, annot=True, fmt='g', ax=ax, cmap='Greens'); #annot=True to annotate cells, fmt='g' to disable scientific notation
11
12 ax.set_xlabel('Predicciones'); ax.set_ylabel('Valores reales');
13 ax.set_title('Matriz de Confusion');
14 ax.xaxis.set_ticklabels(['True', 'False']); ax.yaxis.set_ticklabels(['True', 'False']);
15
16 plt.show()

```

Código de Gráfica ROC-AUC

```

Editor de scripts de Python
⚠ Las filas duplicadas se quitarán de los datos.
6 # Pegue o escriba aquí el código de script:
7 import numpy as np
8 import matplotlib.pyplot as plt
9 import pandas as pd
10 import sklearn
11 from sklearn.preprocessing import LabelEncoder
12 import os
13 import pyodbc
14 conn = pyodbc.connect('DRIVER={ODBC Driver 17 for SQL Server};SERVER=DESKTOP-7ECAC8A;DATABASE=RepuestosWeb;Trusted_Connection=yes;')
15
16 query = "select * from VW_OrdenesEncimaPromedio ;"
17 dataset = pd.read_sql(query, conn)
18 print(dataset.head(26))
19 dataset.drop('ID_Parte', axis=1, inplace=True)
20 dataset.drop('ID_Categoria', axis=1, inplace=True)
21 dataset.drop('Total_Orden', axis=1, inplace=True)
22 X = dataset.iloc[:, 0:2].values
23 y = dataset.iloc[:, -1].values
24 X = dataset.iloc[:, 0:2].values
25 y = dataset.iloc[:, -1].values
26
27 le = LabelEncoder()
28 X[:,0] = le.fit_transform(X[:,0])

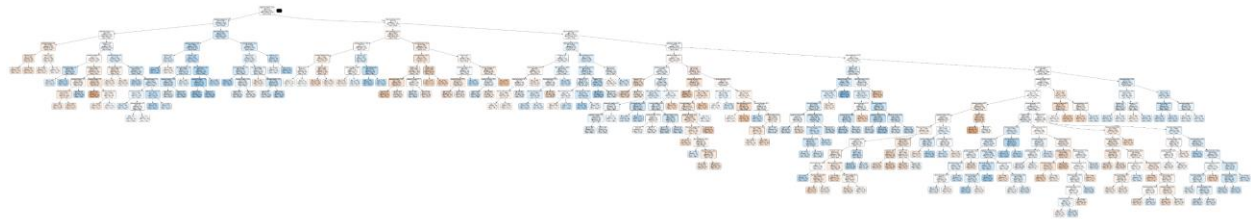
```

Árbol de decisión:

Gráfica:

NombreCiudad, Genero y TotalEncimaPromedio

▽ □ ...

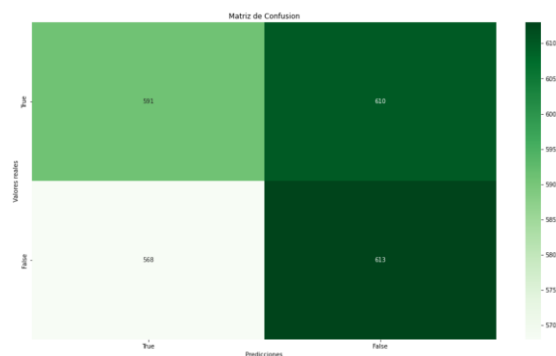


Código:

```

Editor de scripts de Python
⚠ Las filas duplicadas se quitarán de los datos.
6 # Pegue o escriba aquí el código de script:
7 #Almacenar las variables X y Y por separado
8 import numpy as np
9 import matplotlib.pyplot as plt
10 import pandas as pd
11 import sklearn
12 from sklearn.preprocessing import LabelEncoder
13 import os
14 import pyodbc
15
16
17 conn = pyodbc.connect('DRIVER={ODBC Driver 17 for SQL Server};SERVER=DESKTOP-7ECACBA;DATABASE=RepuestosWeb;Trusted_Connection=yes;')
18
19 query = "select * from VW_OrdenesEncimaPromedio ;"
20 df = pd.read_sql(query, conn)
21 print(df.head(26))
22
23
24 df.drop('ID_Parte',axis=1,inplace=True)
25 df.drop('ID_Categoria',axis=1,inplace=True)
26 df.drop('Total_Orden',axis=1,inplace=True)
27 df
28
  
```

Heatmap del Árbol:



Código del mapa de calor:

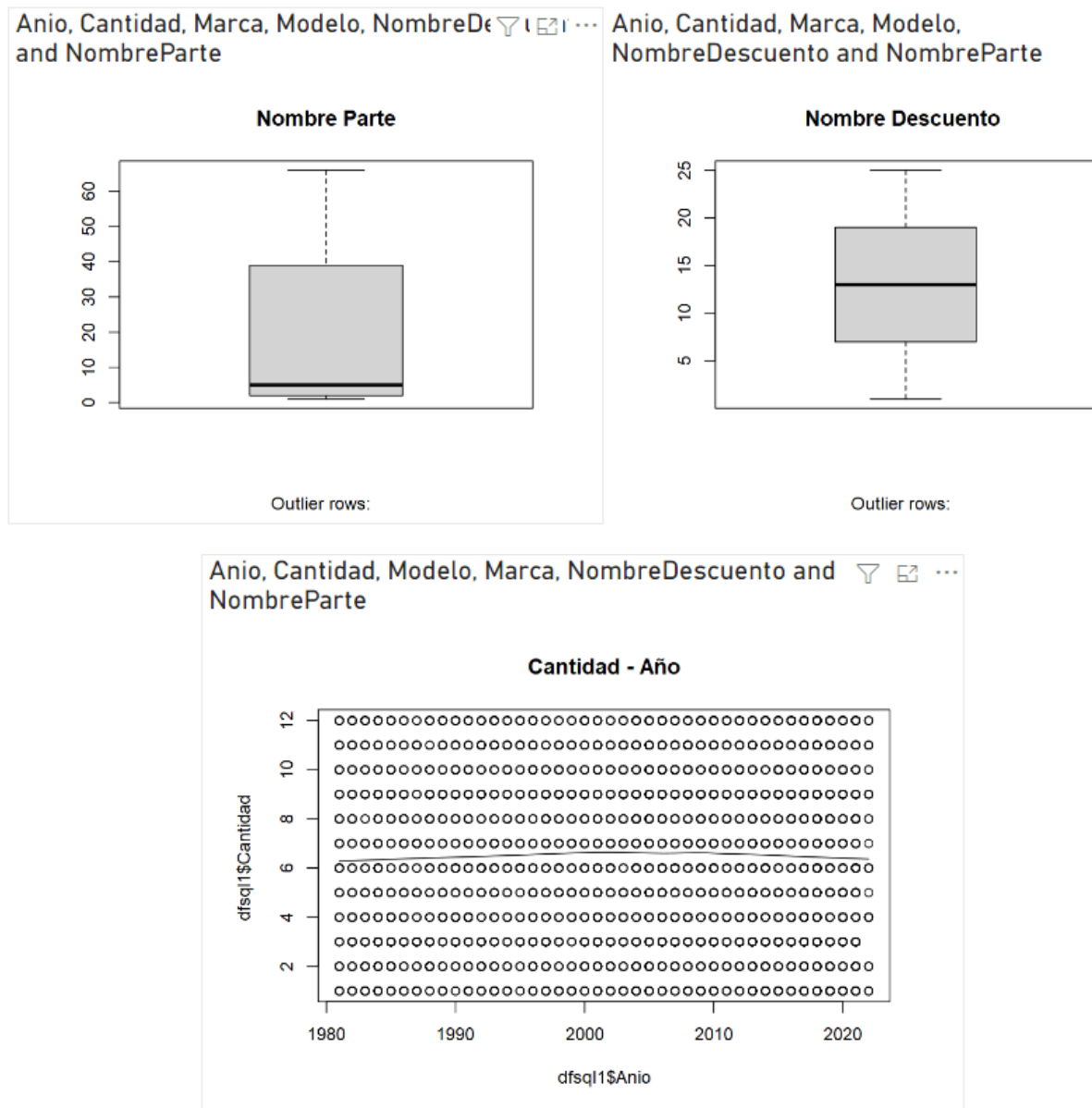
```

Editor de scripts de Python
Las filas duplicadas se quitarán de los datos.

6 # Pegue o escriba aquí el código de script:
7 import numpy as np
8 import matplotlib.pyplot as plt
9 import pandas as pd
10 import sklearn
11 from sklearn.preprocessing import LabelEncoder
12 import os
13 import pyodbc
14
15
16 conn = pyodbc.connect('DRIVER={ODBC Driver 17 for SQL Server};SERVER=DESKTOP-7ECACBA;DATABASE=RepuestosWeb;Trusted_Connection=yes;')
17
18 query = "select * from VM_OrdenesEncinaPromedio ;"
19 df = pd.read_sql(query, conn)
20 print(df.head(26))
21
22
23 df.drop('ID_Parte',axis=1,inplace=True)
24 df.drop('ID_Categoria',axis=1,inplace=True)
25 df.drop('Total_Orden',axis=1,inplace=True)
26 df
27
28
  
```

Regresión Lineal

Gráficas:



Codigo

```
dfsql <- na.omit(dataset)

dfsqlindex <- sample(1:nrow(dfsql), 0.8*nrow(dfsql))

dfsql1 <- dfsql[dfsqlindex,] #80%
dfsql2 <- dfsql[-dfsqlindex,] #20%

for (i in c("NombreParte", "NombreDescuento", "Marca", "Modelo")) {
  dfsql1[,i] = as.factor(dfsql1[,i])
}

for (i in c("NombreParte", "NombreDescuento", "Marca", "Modelo")) {
  dfsql1[,i] = as.numeric(dfsql1[,i])
}

lmHeight = lm(Cantidad~Anio+Marca+Modelo, data = dfsql1) #Create the linear regression
lmHeight2 = lm(Cantidad~NombreParte, data = dfsql1) #Create the linear regression
lmHeight3 = lm(Cantidad~NombreDescuento, data = dfsql1) #Create the linear regression

boxplot(dfsql1$NombreParte, main="Nombre Parte", sub=paste("Outlier rows: ", boxplot.stats(dfsql1$NombreParte)$out))
boxplot(dfsql1$NombreDescuento, main="Nombre Descuento ", sub=paste("Outlier rows: ", boxplot.stats(dfsql1$NombreDescuento)$out))
boxplot(dfsql1$Marca, main="Marca", sub=paste("Outlier rows: ", boxplot.stats(dfsql1$Marca)$out))
```

```
dfsql <- na.omit(dataset)

dfsqlindex <- sample(1:nrow(dfsql), 0.8*nrow(dfsql))

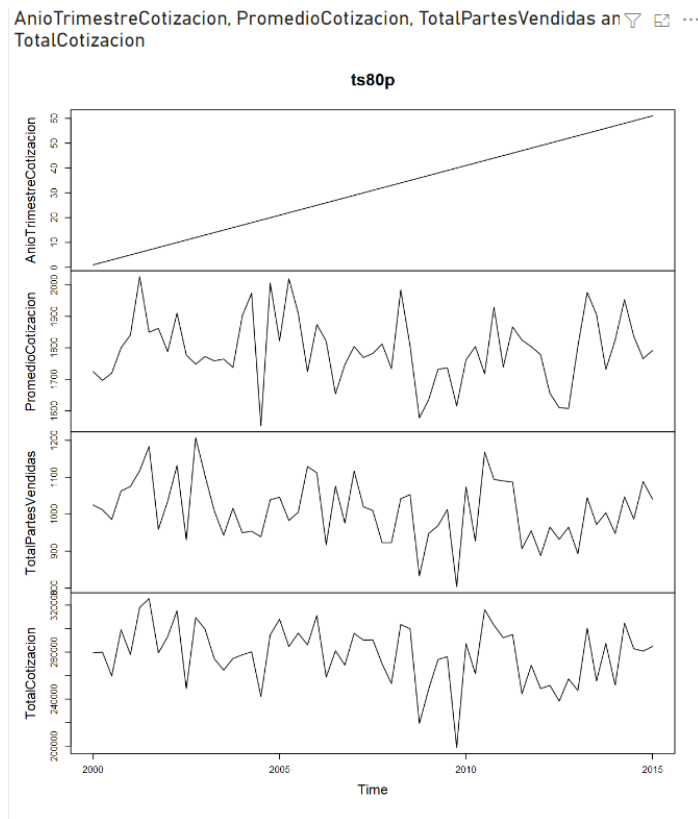
dfsql1 <- dfsql[dfsqlindex,] #80%
dfsql2 <- dfsql[-dfsqlindex,] #20%

for (i in c("NombreParte", "NombreDescuento", "Marca", "Modelo")) {
  dfsql1[,i] = as.factor(dfsql1[,i])
}

for (i in c("NombreParte", "NombreDescuento", "Marca", "Modelo")) {
  dfsql1[,i] = as.numeric(dfsql1[,i])
}

lmHeight = lm(Cantidad~Anio+Marca+Modelo, data = dfsql1) #Create the linear regression
lmHeight2 = lm(Cantidad~NombreParte, data = dfsql1) #Create the linear regression
lmHeight3 = lm(Cantidad~NombreDescuento, data = dfsql1) #Create the linear regression
scatter.smooth(y=dfsql1$Cantidad, x=dfsql1$Anio, main = "Cantidad - Año")
```


Lineas de Tiempo: Grafica



Codigo

```
library(ggplot2)

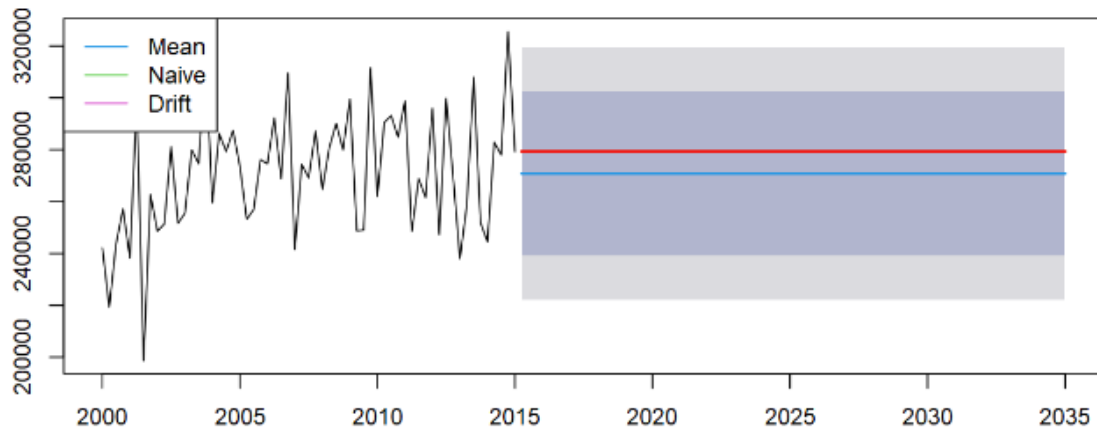
library(forecast)

ts1 <- ts(dataset, start = c(2000,1), frequency = 4)

# la funcion window nos permite dividir un objeto "ts" de un inicio a un fin
ts80p <- window(ts1, start = 2000, end = 2015)
plot(ts80p)
```

Grafica predicciones:

PromedioCotizacion, AnioTrimestreCotizacion, TotalCotizacion and TotalPartesVendidas



Codigo Predicciones:

```
library(ggplot2)

library(forecast)

library(ggplot2)
library(dplyr)
library(DBI)
library(odbc)

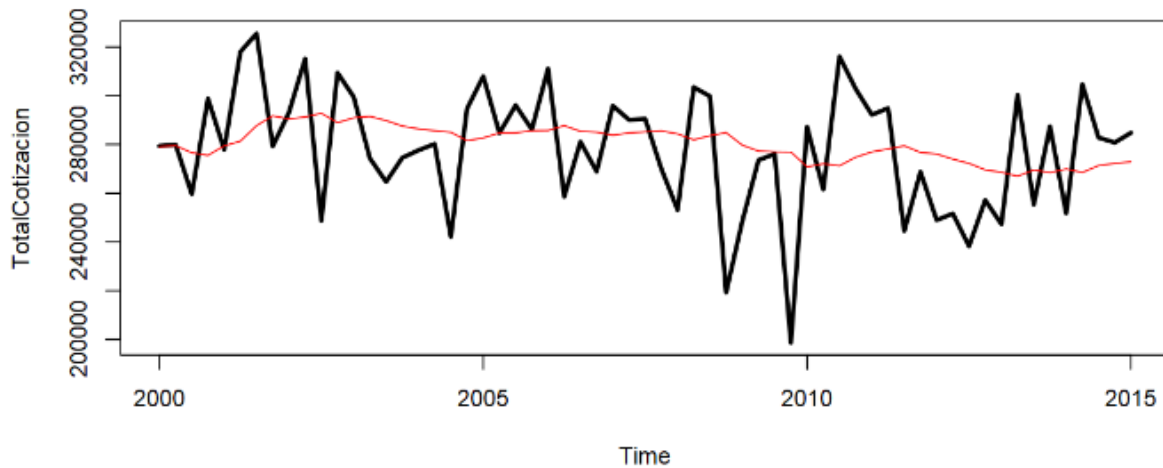
df1 <- dataset %>% select(TotalCotizacion)

ts1 <- ts(df1, start = c(2000,1), frequency = 4)
ts180p <- window(ts1, start = 2000, end = 2015)
meanmodel <- meanf(ts180p,h=80)
naivemodel <- naive(ts180p, h=80)
driftmodel <- rwf(ts180p, h=80)

plot(meanmodel, plot.conf = F, main = "")
lines(naivemodel$mean, col=123, lwd = 2)
lines(driftmodel$mean, col='red', lwd = 2)
legend("topleft",lty=1,col=c(4,123,22),
      legend=c("Mean","Naive","Drift"))
```

Gráfica predicción ajustando al valor real:

AnioTrimestreCotizacion, PromedioCotizacion, TotalCotizacion and
TotalPartesVendidas



Codigo predicción:

```
library(ggplot2)

library(forecast)

library(ggplot2)
library(dplyr)
library(DBI)
library(odbc)

df1 <- dataset %>% select(TotalCotizacion)

ts1 <- ts(df1, start = c(2000,1), frequency = 4)
ts180p <- window(ts1, start = 2000, end = 2015)

meanmodel3 <- meanf(ts180p,h=80)
naivemodel3 <- naive(ts180p, h=80)
driftmodel3 <- rwf(ts180p, h=80)

acf(ts180p, lag.max = 20)
pacf(ts180p, lag.max = 20)

auto.arima(ts180p)
arima(ts180p)

myar = auto.arima(ts180p, stepwise = F, approximation = F)

myar

# Ploteamos vs original
plot(ts180p, lwd = 3)
lines(myar$fitted, col = "red")
```