

Practica 2 Limpieza y Validacion de Datos

Carlos A Gutierrez Cruz 10 de junio de 2018

```
{r setup, include=FALSE} knitr::opts_chunk$set(echo = TRUE)
```

Presentaci3n

En esta pr3ctica se elabora un caso pr3ctico orientado a aprender a identificar los datos relevantes para un proyecto anal3tico y usar las herramientas de integraci3n, limpieza, validaci3n y an3lisis de las mismas.

Competencias

En esta pr3ctica se desarrollan las siguientes competencias del M3ster de Data Science:

- Capacidad de analizar un problema en el nivel de abstracci3n adecuado a cada situaci3n y aplicar las habilidades y conocimientos adquiridos para abordarlo y resolverlo.
- Capacidad para aplicar las t3cnicas espec3ficas de tratamiento de datos (integraci3n, transformaci3n, limpieza y validaci3n) para su posterior an3lisis.

Objetivos

Los objetivos concretos de esta pr3ctica son:

- Aprender a aplicar los conocimientos adquiridos y su capacidad de resoluci3n de problemas en entornos nuevos o poco conocidos dentro de contextos m3s amplios o multidisciplinarios.
- Saber identificar los datos relevantes y los tratamientos necesarios (integraci3n, limpieza y validaci3n) para llevar a cabo un proyecto anal3tico.
- Aprender a analizar los datos adecuadamente para abordar la informaci3n contenida en los datos.
- Identificar la mejor representaci3n de los resultados para aportar conclusiones sobre el problema planteado en el proceso anal3tico.
- Actuar con los principios 3ticos y legales relacionados con la manipulaci3n de datos en funci3n del 3mbito de aplicaci3n.
- Desarrollar las habilidades de aprendizaje que les permitan continuar

Desarrollo de la Pr3ctica

1. Descripci3n del dataset. ¿Por qu3 es importante y qu3 pregunta/problema pretende responder?

Para el desarrollo de la pr3ctica se ha escogido un conjunto de datos (dataset) del sitio de Kaggle con la siguiente direcci3n. “<https://www.kaggle.com/checoalejandro/autos-consumo-gasolina-mexico/data>”

La intenci3n de haber escogido este dataset fu3 primero escoger un dataset de algo relacionado con M3xico que es mi pa3s de origen y en segundo porque es conocido que la Ciudad de M3xico es una de las ciudades m3s pobladas del mundo y que uno de sus principales problemas es la contaminaci3n originados principalmente por los veh3culos automotores que circulan diariamente en la ciudad.

Bajo el supuesto de que el dataset es una muestra aleatoria y representativa sobre el consumo de gasolina de los veh3culos automotores de modelos 2011-2018 en M3xico puede ejemplificar perfectamente el uso de t3cnicas con programaci3n en R para el an3lisis de datos.

El dataset se compone de una muestra de 4617 observaciones y 18 variables que a continuación se enlistan y describen.

- Marca - Marca fabricante del vehículo
- Submarca - Sublinea
- Modelo - Año del vehículo
- Trans. - Tipo de transmisión (manual, automática, CVT)
- Comb. - Tipo de combustible (Diesel, Gasolina)
- Cilindros - Número de cilindros del vehículo (4, 6, 8)
- Potencia (HP) - Potencia del motor en caballos de fuerza (HP)
- Tamaño (L) - Tamaño del motor (1.5, 1.8, 2.0)
- Categoría - Auto compacto, lujo, deportivo, SUV
- R. Ciudad (km/l) - Rendimiento en Ciudad kilómetros por litro
- R. Carr. (km/l) - Rendimiento en Carretera kilómetros por litro
- R. Comb. (km/l) - Rendimiento combinado (ciudad/carretera)
- R. Ajust. (km/l) - Rendimiento ajustado
- CO2(g/km) - Cantidad de emisión de Dioxido de Carbono
- NOx (g/1000km) - Cantidad de Oxígeno de Nitrogeno
- Calificación Gas Ef. Inv. - Categorización de uso de efectivo de combustible
- Calificación Contam. Aire - Categorización sobre el nivel de contaminación.

Dado la información anterior es interesante plantearse algunas preguntas como:

Si los vehículos cubren las normas de emisión planteadas El nivel de emisión de los vehículos tiene relación con el modelo El rendimiento en el consumo de combustible esta relacionado con el modelo o tipo de vehículo.

Se pueden hacer muchas preguntas al respecto, pero lo importante del ejercicio es mostrar la utilidad que tienen las técnicas y herramientas dentro del procesamiento y análisis de los datos.

2. Integración y selección de los datos de interés a analizar.

Procederemos a cargar el archivo “Consumo Gasolina Autos Ene 2018.csv” y haremos un análisis exploratorio para observar el tipo de información contenida que pueda ser de utilidad.

```
dir()
setwd("D:/UOC/02 Tipología y ciclo de vida de los datos/Practica 2 Limpieza y Validacion de Datos/Pract")
# Carga de datos
data <- read.csv("Consumo Gasolina Autos Ene 2018.csv", header=TRUE)

attach(data)

summary(data)

# Tipo de datos de cada variable
sapply(data, function(x) class(x))

table (data[,6])
```

Dada la información mostrada, no sería útil considerar los vehículos con uso de combustible diesel ya que por el número de observaciones que se tienen no son significantes, así mismo nos centraremos en los autos para lo cual no consideraremos en este trabajo las camionetas o vehículos utilitarios ya que además de estar sujetos a otros niveles ambientales, tienen un rendimiento de combustible mucho menor a los autos a fin de evitar mucha variabilidad y dispersión de los datos durante el análisis.

```
datos2 <- subset(data, Comb=="Gasolina" & Categoría != "CAMIONETAS DE USO MULTIPLE (SUV)")
attach(datos2)
summary(datos2)
```

3. Limpieza de los datos.

3.1. ¿Los datos contienen ceros o elementos vacíos? ¿Cómo gestionarlos cada uno de estos casos?

las siguientes instrucciones en R nos permitan conocer si existen ceros o elementos vacíos.

```
sapply(datos2, function(x) sum(is.na(x)))
```

Lo cual nos permite apreciar que no existen dentro de las variables de la nuevo dataset(datos2) elementos vacíos.

Si hubieran existido elementos vacíos dentro de variables se pudieran haber tratado de la siguiente forma.

- Eliminar los registros si fueran pocas las observaciones con elementos vacíos
- Completar la información con el promedio de las categorías de la variable si es que no hay mucha variabilidad en los datos.
- Imputación de valores basado en valores próximos

3.2. Identificación y tratamiento de valores extremos.

Los siguientes diagramas en R para analizar variables de interés nos ayudan a identificar outliers o valores extremos.

```
#diagramas de boxplot
```

```
par(mfrow=c(1,1))
```

```
boxplot(CO2.g.km.~ Modelo, data=datos2, main="Emision de CO2 por modelo", xlab="Modelo", ylab="CO2 g/km")
```

```
boxplot(NOx..g.1000km.~ Modelo, data=datos2,main="Emision de NOx por modelo", xlab="Modelo", ylab="NOx g/km")
```

```
boxplot(R..Comb...km.l. ~ Cilindros ,data=datos2,main="Rendimiento Combustible Combinado (Ciudad/Carr.) por Cilindros")
```

```
boxplot(R..Comb...km.l. ~ Modelo, data=datos2,main="Rendimiento Combustible Combinado (Ciudad/Carr.) por Modelo")
```

```
#grabamos los datos a utilizar para el analisis
```

```
write.csv(datos2, "Consumo Gasolina Autos Ene 2018_Nuevo.csv")
```

4. Análisis de los datos.

4.1. Selección de los grupos de datos que se quieren analizar/comparar (planificación de los análisis a aplicar).

Las variables o grupos de datos para analizar y comparar serán:

- Modelo del vehículo
- Cilindros
- Rendimientos de combustible
- Niveles de emisión

```
#Resumen de variables de interes
```

```
#Agrupación por Modelo
```

```
table(Modelo)
```

```
#Agrupación por Num. cilindros
```

```
table(datos2$Cilindros)
```

```
#Resumen de medidas de tendencia central
summary(datos2$R..Comb...km.l.)
summary(datos2$CO2.g.km.)
summary(datos2$NOx..g.1000km.)

#Grupos por tipo num de cilindros
datos2.hasta4 <- subset(datos2, Cilindros <= 4)
datos2.mayor4 <- subset(datos2, Cilindros > 4)
```

4.2. Comprobaci3n de la normalidad y homogeneidad de la varianza.

Para la comprobacion de la normalidad existen dos pruebas importantes Kolmogorov-Smirnov y Shapiro-Wilk. Dado que el numero de observaciones es mayor a 50 usaremos la de Kolmogorov-Smirnov

A continuacin se muestra un ejemplo las variables relacionadas con los niveles de emision CO2 y NOx

#Prueba de Normalidad

```
ks.test(CO2.g.km., "pnorm", mean=mean(CO2.g.km., sd=sd(CO2.g.km.)))
```

```
ks.test(NOx..g.1000km. , "pnorm", mean=mean(NOx..g.1000km.), sd=sd(NOx..g.1000km.))
```

#Prueba de F para igualdad de varianzas

```
var.test(datos2.hasta4$CO2.g.km., datos2.mayor4$CO2.g.km.)
```

El resultado de p-value en la prueba de Kolmogorov-Smirnov es tan pequea±o nos ayuda a determinar que se rechaza la hip3tesis alternativa y se acepta la hip3tesis nula de que las observaciones provienen de una distribuci3n normal.

La prueba de F tambien nos ayuda a determinar la igualdad de las varianzas para los grupos de emisiones para automoviles de menos de 4 cilindros con los mayores a 4 cilindros

Lo cual tambi3n puede ejemplificarse de forma grafica

```
par(mfrow=c(2,2))
qqnorm(CO2.g.km., main = "Q-Q Norm Emisiones CO2")
qqline(CO2.g.km., col="red")
hist(CO2.g.km., main = "Histograma Emisiones CO2")
```

4.3. Aplicaci3n de pruebas estadísticas para comparar los grupos de datos.

En funci3n de los datos y el objetivo del estudio, aplicar pruebas de contraste de hip3tesis, correlaciones, regresiones, etc.

A continuacion vamos a comprobar que los vehiculos de menos de 4 cilindros cumplen con la norma de niveles de emision de Oxido de Nitrogeno NOx. La cual es de 21 g/1000km. Para ello nos apoyaremos de la prueba de t donde las hipotesis a comprobar son H0 La emision de NOx cumple al menos con la norma Ha La emisi

onde Nox es mayor a la norma

#Prueba de t para una muestra con intervalo de confianza al 95%

```
t.test(datos2.hasta4$NOx..g.1000km., alternative = "greater", mu=21)
```

```
t.test(datos2.mayor4$NOx..g.1000km., alternative = "greater", mu=21)
```

Estadísticamente se puede decir que los vehiculos con mayor a 4 cilindros cumplen con la norma.

Por ultimo realizaremos una relacion si el rendimiento de combustible tiene relacion con el tamaño del modelo vehiculo y numero de cilindros. Para ellos revisaremos los coeficientes de correlacion y modelos de regresion

```
cor(datos2$Modelo, datos2$R..Comb...km.l.)

cor(datos2$Cilindros, datos2$R..Comb...km.l.)

pairs(~datos2$Cilindros+datos2$Modelo+datos2$R..Comb...km.l.)

summary(lm(datos2$R..Comb...km.l.~datos2$Modelo+datos2$Cilindros))
```

Conforme al resultado de la funcion de regresion `lm()`, El rendimiento de combustible de los vehiculos tiene una significativa relacion con las varibales modelo y numero de cilindros. Ademas el coeficiente R^2 nos muestra que cerca del 67% de los datos quedarian explicados por el modelo.

5. Representación de los resultados a partir de tablas y gráficas.

```
ryx <-lm(datos2$R..Comb...km.l.~datos2$Tamaño..L.)

linea <- ryx$coefficients[1]+(ryx$coefficients[2]*datos2$Tamaño..L.)
plot(datos2$Tamaño..L.,datos2$R..Comb...km.l., xlab = "Tamaño", ylab="Rendimiento km/l")

lines(linea, col=4)
```

La grafica muestra el modelo lineal para la relación que hay entre el rendimiento de comustible dependiendo el tamaño del motor (1.5litros, 2.0 litros, etc)

6. Resolución del problema. A partir de los resultados obtenidos, ¿cuáles son las conclusiones? ¿Los resultados permiten responder al problema?

El caso práctico nos permitió mostrar la aplicación práctica de las técnicas y herramientas para el procesamiento y analisis de datos.

Basado en el dataset nos permitio ejemplificar que no existía una diferencia en cuanto a los niveles de emisión CO_2 o No_x independientemente del número de cilindros del motor

También nos permitio ver estadísticamente que el rendimientos de combustible de los vehiculos tomados de la muestra tienen una relación significativa con el modelo y numero de cilindros de la maquina.

Por último, también se puede decir que estadísticamente, los vehiculos con numero de cilindros mayor a 4 cumplieron con la norma en el nivel de emisión, lo que hace suponer que a pesar de que son maquinas mas grande nos hace suponer que trabajan a mayor estabilidad que los motores con menos cilindros.

Dado lo anterior siempre se puede decir que este tipo de conclusiones estan basadas en la muestra utilizada y para otro tipo de prueba similar siempre es importante llevar acabo una planeación y diseño para la recolección de los datos sobre el problema a resolver o analizar.