

# R Notebook

This is an R [Markdown](#) Notebook. When you execute code within the notebook, the results appear beneath the code.

Try executing this chunk by clicking the *Run* button within the chunk or by placing your cursor inside it and pressing *Ctrl+Shift+Enter*.

Hide

Hide

```
install.packages("reshape2")
```

WARNING: Rtools is required to build R packages but is not currently installed. Please download and install the appropriate version of Rtools before proceeding:

<https://cran.rstudio.com/bin/windows/Rtools/>  
Installing package into 'C:/Users/Edgar/AppData/Local/R/win-library/4.3'  
(as 'lib' is unspecified)  
also installing the dependency 'plyr'

trying URL 'https://cran.rstudio.com/bin/windows/contrib/4.3/plyr\_1.8.8.zip'  
Content type 'application/zip' length 1161886 bytes (1.1 MB)  
downloaded 1.1 MB

trying URL 'https://cran.rstudio.com/bin/windows/contrib/4.3/reshape2\_1.4.4.zip'  
Content type 'application/zip' length 454195 bytes (443 KB)  
downloaded 443 KB

package 'plyr' successfully unpacked and MD5 sums checked  
package 'reshape2' successfully unpacked and MD5 sums checked

The downloaded binary packages are in  
C:/Users/Edgar/AppData/Local/Temp/Rtmpmfn20/downloaded\_packages

Hide

Hide

```
library(reshape2)
```

Attaching package: 'reshape2'

The following object is masked from 'package:tidyr':

smiths

Hide

Hide

```
library(gridExtra)
library(ggplot2)
library(tidyverse)
library(dplyr)
```

Add a new chunk by clicking the *Insert Chunk* button on the toolbar or by pressing *Ctrl+Alt+I*.

When you save the notebook, an HTML file containing the code and output will be saved alongside it (click the *Preview* button or press *Ctrl+Shift+K* to preview the HTML file).

The preview shows you a rendered HTML copy of the contents of the editor. Consequently, unlike *Knitr*, *Preview* does not run any R code chunks. Instead, the output of the chunk when it was last run in the editor is displayed.

## 3 DATA EXPLORATION

### 1 Extracto

cargar los datos del data set carIns\_final.RData la cual ya tiene descrtados los valores NA

1 usando el paquete de dplyr, responda las siguientes preguntas

a) Obtener el número de carros agrupados por bodystyle

Hide

Hide

```
load("C:/Users/Edgar/Documents/GitHub/DataMining_and_MachineLearning_EdgarV/data/HandsOn_Data/carIns_final.Rdata")
## 1s()

df <- carIns_final

num_carros_por_bodystyle <- df %>%
  group_by(bodyStyle) %>%
  count()
```

num\_carros\_por\_bodystyle

bodyStyle <fctr> n <int>

convertible 6

hardtop 8

hatchback 70

sedan 96

wagon 25

5 rows

Hide

Hide

NA

b) Obtener el número de carros agrupados por bodystyle y ademas por fuelType

Hide

Hide

```
num_carros_por_bodystyle_fueltype <- df %>%
  group_by(bodyStyle, fuelType) %>%
  count()
```

## NOS MUESTRA EL NUMERO DE VEHICULOS SI ESTE ES A GASOLINA O A DIESEL Y POR CADA DISEÑO

num\_carros\_por\_bodystyle\_fueltype

bodyStyle <fctr> fuelType <fctr> n <int>

convertible gas 6

hardtop diesel 1

hardtop gas 7

hatchback diesel 1

hatchback gas 69

sedan diesel 15

sedan gas 81

wagon diesel 3

wagon gas 22

9 rows

Hide

Hide

NA

c) Obtener la media y la desviación estandar del atributo cityMpg agrupado por bodyStyle y en orden ascendente

Hide

Hide

```
media_st_cityMpg_por_bodystyle <- df %>%
  group_by(bodyStyle)%>%
  summarise(cityMpg_mean = mean(cityMpg), cityMpg_st = sd(cityMpg))%>%
  arrange(cityMpg_mean, cityMpg_st)
```

media\_st\_cityMpg\_por\_bodystyle

bodyStyle <fctr> cityMpg\_mean <dbl> cityMpg\_st <dbl>

convertible 20.50000 3.391165

hardtop 21.62500 5.423165

wagon 24.04000 4.217819

sedan 25.32292 6.599035

hatchback 26.31429 7.169870

5 rows

Hide

Hide

NA

d) Agrupelos por bodyStyle los atributos de cityMpg y highWayMpg, obtener la media, la mediana, desviación estandar y el rango inter-cuartil

Hide

Hide

```
resumen_cityMpg_highwayMpg <- df %>%
  group_by(bodyStyle)%>%
  summarise(cityMpg_mean=mean(cityMpg), highwayMpg_mean=mean(highwayMpg),
    cityMpg_mediam=median(cityMpg), highwayMpg_mediam=median(highwayMpg),
    cityMpg_sd=sd(cityMpg), highwayMpg_sd=sd(highwayMpg),
    cityMpg_rec=IQR(cityMpg), highwayMpg_rec=IQR(highwayMpg))
```

resumen\_cityMpg\_highwayMpg

bodyStyle <fctr> cityMpg\_mean <dbl> highwayMpg\_mean <dbl> cityMpg\_mediam <dbl> highwayMpg\_mediam <dbl> cityMpg\_sd <dbl>

convertible 20.50000 26.00000 21 27.0 3.391165

hardtop 21.62500 27.25000 23 27.5 5.423165

hatchback 26.31429 32.17143 26 31.5 7.169870

sedan 25.32292 30.83333 25 30.5 6.599035

wagon 24.04000 28.72000 24 29.0 4.217819

5 rows | 1-6 of 9 columns

Hide

Hide

NA

### 2. VISUALIZACIÓN

2.- Usando el paquete ggplot2, crear el grafico que usted encuentre adecuado para responder las siguientes preguntas

e) Mostrar la relación entre cityMpg y highwayMpg

Hide

Hide

```
ggplot(df, aes(x = cityMpg, y = highwayMpg)) +
  geom_point() + labs(title="Relación consumo combustible ciudad vs carretera")
```

Relación consumo combustible ciudad vs carretera



f) Mostrar la distribución de carros por su diseño

Hide

Hide

```
ggplot(num_carros_por_bodystyle, aes(x = bodyStyle, y = n)) +
  geom_bar(stat = "identity") + labs(title="Distribución de vehiculos por el diseño de su carrocería")
```

Distribución de vehiculos por el diseño de su carrocería



g) Mostrar la distribución de carros por su precio (Sugerencia: establecer el ancho de la barra en 5000)

Hide

Hide

```
carros_por_precio <- df%>%
  group_by(price)

ggplot(carros_por_precio, aes(x = price)) +
  geom_histogram(binwidth = 5000) + labs(title="Distribución de vehiculos por precio")
```

Distribución de vehiculos por precio



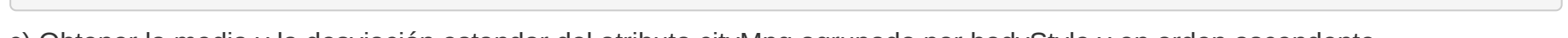
h) agregar la información de la estimación de densidad el grafico anterior

Hide

Hide

```
ggplot(carros_por_precio, aes(x = price)) +
  geom_histogram(aes(y = after_stat(density)), binwidth = 5000, color = "purple", fill = "lightblue") + labs(title="Distribución de vehiculos por precio vs densidad")
```

Distribución de vehiculos por precio y densidad



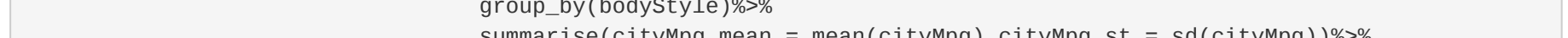
i) comprobar (visualmente) si es aceptable considerar ese precio para conseguir una distribución normal

Hide

Hide

```
carros_por_precio <- rnorm(103)
ggplot(data.frame(carros_por_precio), aes(sample = carros_por_precio)) +
  geom_qq() + labs(title = "distribución normal de una muestra de datos de precio") +
  geom_qq_line()
```

distribución normal de una muestra de datos de precio



j) Mostrar la distribución del precio por el atributo marca (Sugerencia: usar Boxplots y la función coord\_flip())

Hide

Hide

```
carros_por_precio1 <- df%>%
  group_by(price)

ggplot(carros_por_precio1, aes(x = make, y = price)) +
  geom_boxplot() +
  coord_flip()
```

make



k) Mostrar la distribución de precios por elatributo nDoors(numero de puertas) (Sugerencia: usar histogramas)

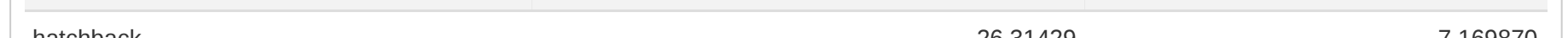
Hide

Hide

```
histograma1 <- ggplot(new_carro_por_precio, aes(x = price)) +
  geom_histogram(fill = "steelblue", color = "white", bins = 30) +
  labs(x = "Atributo 1", y = "Frecuencia") +
  ggtitle("Histograma de Precio")

histograma2 <- ggplot(new_carro_por_precio, aes(x = nDoors)) +
  geom_histogram(fill = "green", color = "white", bins = 30) +
  labs(x = "Atributo 2", y = "Frecuencia") +
  ggtitle("Histograma de Atributo 2")
grid.arrange(histograma1, histograma2, ncol = 3)
```

Histograma de Precio



Histograma de Atributo 2



N ##### l) Mostrar la distribución del precio por bodyStyle y atributos nDoors, (Sugerencia: usar Histogramas)

Hide

Hide

```
histograma1 <- ggplot(new_carro_por_precio, aes(x = price)) +
  geom_histogram(fill = "steelblue", color = "white", bins = 30) +
  labs(x = "Precio", y = "count") +
  ggtitle("Histograma de Precio")

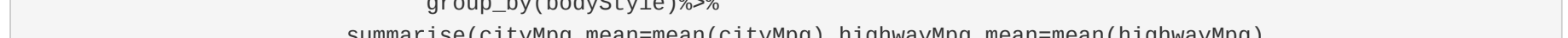
histograma2 <- ggplot(new_carro_por_precio, aes(x = nDoors)) +
  geom_histogram(fill = "green", color = "white", bins = 30) +
  labs(x = "num Puertas", y = "count") +
  ggtitle("Histograma de nDoors")

histograma3 <- ggplot(new_carro_por_precio, aes(x = bodyStyle)) +
  geom_histogram(fill = "red", color = "white", bins = 30) +
  labs(x = "Diseño", y = "count") +
  ggtitle("Histograma de bodyStyle")
grid.arrange(histograma1, histograma2, histograma3, ncol = 3)
```

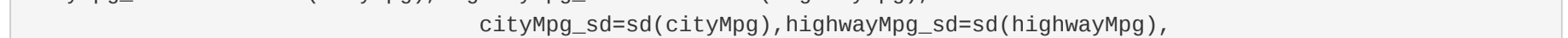
Histograma de Precio



Histograma de nDoors



Histograma de bodySt



m) agregar el parámetro free\_y a la función de la faceta en el gráfico anterior

Hide

Hide

```
carros <- new_carro_por_precio %>%
  group_by(bodyStyle, nDoors)%>%
  select(price)
```

Adding missing grouping variables: 'bodyStyle', 'nDoors'

Hide

Hide

carros

bodyStyle <fctr> nDoors <int> price <dbl>

1 2 13495

1 2 16500

3 2 16500

4 1 13950

4 1 17450

4 2 15250

4 1 17710

5 1 18920

4 1 23875

3 2 5572

1-10 of 205 rows

Previous 1 2 3 4 5 6 21 Next

Hide

Hide

```
carros_melt <- reshape2::melt(carros)
```

No id variables; using all as measure variables

Hide

Hide

```
histograma_final <- ggplot(carros_melt, aes(x = value)) +
  geom_histogram(fill = "steelblue", color = "white", bins = 30) +
  facet_wrap(~ variable, scales = "free_y") +
  labs(x = "Valor", y = "Frecuencia") +
  ggtitle("Histograma de los Atributos")

histograma_final
```

Histograma de los Atributos

