

PRACTICA 2 - Limpieza y análisis de datos

Presentada por Jhon Harry Loaiza y Pablo Jesús Sánchez

Tipología y ciclo de vida del dato - Aula 3

18 dic 2020

Contenido:

1. Descripción del dataset

2. Integración y selección de los datos de interés

3. Limpieza de los datos

- ¿Los datos contienen ceros o elementos vacíos? ¿Cómo gestionarlos cada uno de estos casos?
- Identificación y tratamiento de valores extremos

4. Análisis de los datos

- Selección de los grupos de datos que se quieren analizar/comparar (planificación de los análisis a aplicar)
- Comprobación de la normalidad y homogeneidad de la varianza
- Aplicación de pruebas estadísticas para comparar los grupos de datos

5. Visualización de los datos

6. Machine Learning para agrupar en categorías los automóviles

7. Conclusión

1. Descripción del dataset

El conjunto de datos objetado análisis se llama *Automobile_data.csv*, el cual se ha obtenido en Kaggle y está conformado por 26 columnas y 205 filas, las cuales corresponden a diferentes características que posee un vehículo como el tamaño del motor, la potencia, el número de puertas, etc, y distintos vehículos por fabricante, respectivamente.

Como se indica en la página de Kaggle donde se aloja el dataset, el conjunto de datos corresponde a información que ha sido consolidada de distintos fuentes:

- 1) 1985 Model Import Car and Truck specifications, 1985 Ward's Automotive Yearbook
- 2) Manuales de automóviles personales, Oficina de servicios de seguros, 160 Water Street, Nueva York, NY 10038
- 3) Informe de colisión de seguros, Instituto de seguros para la seguridad en las carreteras, Watergate 600, Washington, DC 20037

La descripción de las variables del dataset es la siguiente:

- symboling**: índice de riesgo del vehículo en términos de equipamiento de seguridad
- normalized-losses**: promedio relativo de precio del pago de un vehículo asegurado anualmente
- make**: marca del fabricante del vehículo
- fuel-type**: tipo de combustible del motor
- aspiration**: tipo de aspiración que posee el motor del vehículo
- num-of-doors**: número de puertas
- body-style**: tipo del vehículo
- drive-wheels**: tipo de tracción en las ruedas del carro
- engine-location**: ubicación del motor en el vehículo
- wheel-base**: distancia entre los ejes del vehículo
- length**: longitud o largo del vehículo en centímetros
- width**: ancho del vehículo en centímetros
- height**: altura del vehículo en centímetros
- curb-weight**: peso del vehículo
- engine-type**: tipo de motor del vehículo
- num-of-cylinders**: número de cilindros del motor
- engine-size**: tamaño del motor
- fuel-system**: sistema de combustión del motor
- bore**: diámetro de los cilindros
- stroke**: distancia que posee el pistón
- compression-ratio**: índice de compresión del motor
- horsepower**: potencia del motor (caballos de fuerza)
- peak-rpm**: revoluciones por minuto del motor
- city-mpg**: consumo de combustible del vehículo en millas (en ciudad)
- highway-mpg**: consumo vehículo en millas (en autopista)
- price**: precio del vehículo en USD

Este dataset en principio podría ser usado para predecir el valor de un vehículo dados sus características como largo, ancho, el tipo de motor que usa, su rendimiento, el tipo de tracción, entre otros. En nuestro caso, se pretende realizar un análisis de clusters para encontrar los tipos de vehículos que pueden existir en el conjunto de datos, usando solo las variables inherentes al vehículo.

1.1 Importancia del dataset

El objetivo con este dataset es construir a partir de este, un modelo de clustering, con el fin de clasificar los vehículos según sus características y discriminarlos por grupos homogéneos, que permita reconocer y clasificar nuevos vehículos en cada uno de estos segmentos, lo cual también ayudaría a los clientes a tomar una decisión de compra por ejemplo, al conocer ciertas características, podría saber de antemano cual es el vehículo que mejor se ajusta a sus necesidades dados los segmentos generados. Por lo cual esta clasificación de los vehículos permite a los concesionarios y al cliente tener una mejor perspectiva de lo que se necesita para un vehículo y poder tener una mejor manera de comprar.

2. Integración y selección de los datos de interés

En este paso procedemos a importar las librerías necesarias para el análisis de los datos además de cargar los datos

Importe de librerías y cargue de los datos

```
In [1]: import pandas as pd
pd.options.display.html.table_schema = True

import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import matplotlib.pyplot as plt
warnings.filterwarnings("ignore")
plt.style.use("ggplot")
```

```
In [2]: # Cargue de los datos
df = pd.read_csv('Automobile_data.csv')
# Verifica que no haya datos faltantes
df.head()
```

```
Out[2]:
```

	symboling	normalized-losses	make	fuel-type	aspiration	num-of-doors	body-style	drive-wheels	engine-location	wheel-base	length	width	engine-size	fuel-system	bore	stroke	compression-ratio	
0	alfa-romero	3	7	alfa-romero	gas	std	two	convertible	nwd	front	88.6	168.8	64.1	130	mpfi	3.47	2.68	9.0
1	alfa-romero	3	7	alfa-romero	gas	std	two	convertible	nwd	front	88.6	168.8	64.1	130	mpfi	3.47	2.68	9.0
2	alfa-romero	1	7	alfa-romero	gas	std	two	hatchback	nwd	front	94.5	171.2	65.5	152	mpfi	2.68	3.47	9.0
3	audi	2	164	audi	gas	std	four	sedan	fwd	front	99.8	176.6	66.2	109	mpfi	3.19	3.4	10.0
4	audi	2	164	audi	gas	std	four	sedan	4wd	front	99.4	176.6	66.4	136	mpfi	3.19	3.4	8.0

5 rows x 26 columns

```
In [3]: # Dimensiones de los datos
print("Dimensiones del dataset:" + str(df.shape))
# Información de los datos
print(df.info())

Dimensiones del dataset: (205, 26)
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 205 entries, 0 to 204
Data columns (total 26 columns):
 #   Column              Non-Null Count  Dtype
---  --
 0   symboling            205 non-null    int64
 1   normalized-losses    205 non-null    object
 2   make                 205 non-null    object
 3   fuel-type            205 non-null    object
 4   aspiration            205 non-null    object
 5   num-of-doors         205 non-null    object
 6   body-style           205 non-null    object
 7   drive-wheels         205 non-null    object
 8   engine-location      205 non-null    object
 9   wheel-base          205 non-null    float64
10   length               205 non-null    float64
11   width               205 non-null    float64
12   height              205 non-null    float64
13   curb-weight         205 non-null    int64
14   engine-type          205 non-null    object
15   engine-size          205 non-null    object
16   fuel-system          205 non-null    object
17   horsepower           205 non-null    float64
18   peak-rpm            205 non-null    object
19   city-mpg            205 non-null    object
20   highway-mpg         205 non-null    object
21   price               205 non-null    object
Dtypes: float64(5), int64(3), object(18)
memory usage: 41.84 KB
None
```

En general, los datos han sido cargados correctamente en la definición de la variable, salvo algunas como **price**, **horsepower**, entre otras que deberían aparecer como numéricas y no lo son.

Dado que nuestro objetivo es hacer cluster de los vehículos basados en sus características, podemos descartar las variables **symboling** y **normalized-losses**, ya que estas variables tienen que ver más con temas actitudinales o de asegurabilidad y no son necesarias para clasificar los vehículos según sus características para el objeto de análisis de los datos.

```
In [4]: # Eliminamos las variables 'symboling' y 'normalized-losses'
df.drop(["symboling", "normalized-losses"], axis = 1, inplace = True)
```

```
In [5]: # Comprobamos el dataset
df.head()
```

```
Out[5]:
```

	make	fuel-type	aspiration	num-of-doors	body-style	drive-wheels	engine-location	wheel-base	length	width	engine-size	fuel-system	bore	stroke	compression-ratio
0	alfa-romero	gas	std	two	convertible	nwd	front	88.6	168.8	64.1	130	mpfi	3.47	2.68	9.0
1	alfa-romero	gas	std	two	convertible	nwd	front	88.6	168.8	64.1	130	mpfi	3.47	2.68	9.0
2	alfa-romero	gas	std	two	hatchback	nwd	front	94.5	171.2	65.5	152	mpfi	2.68	3.47	9.0
3	audi	gas	std	four	sedan	fwd	front	99.8	176.6	66.2	109	mpfi	3.19	3.4	10.0
4	audi	gas	std	four	sedan	4wd	front	99.4	176.6	66.4	136	mpfi	3.19	3.4	8.0

5 rows x 24 columns

3. Limpieza de los datos

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

Como se vio en la descripción del dataset, los tipos de variables son diversos, y se pudo observar que existen algunas variables catalogadas en su tipo de forma incorrecta, por ejemplo en el caso de las variables **price** y **horsepower**, vemos que aparecen unos caracteres "?" en algunos registros, lo que indica que en este dataset, los valores nulos han sido marcados con este carácter. Se lleva a cabo una inspección en general por lo que se procede a verificar cuántas variables poseen este carácter y corregir así el tipo de dato de las mismas.

```
In [6]: # Seleccionamos las variables categoricas del dataset
df_str = df.select_dtypes(include = "object")
df_str.info()
```

```
Out[6]:
```

	make	fuel-type	aspiration	num-of-doors	body-style	drive-wheels	engine-location	wheel-base	length	width	engine-size	fuel-system	bore	stroke	compression-ratio
0	alfa-romero	gas	std	two	convertible	nwd	front	88.6	168.8	64.1	130	mpfi	3.47	2.68	9.0
1	alfa-romero	gas	std	two	convertible	nwd	front	88.6	168.8	64.1	130	mpfi	3.47	2.68	9.0
2	alfa-romero	gas	std	two	hatchback	nwd	front	94.5	171.2	65.5	152	mpfi	2.68	3.47	9.0
3	audi	gas	std	four	sedan	fwd	front	99.8	176.6	66.2	109	mpfi	3.19	3.4	10.0
4	audi	gas	std	four	sedan	4wd	front	99.4	176.6	66.4	136	mpfi	3.19	3.4	8.0

5 rows x 24 columns

3.2. Limpieza de los datos

Procedemos ahora a limpiar estas variables y reemplazar los registros faltantes haciendo uso del algoritmo de *regresión Bayesiana Ridge* para el caso de las variables **num-cylinders** y así permitir realizar una regresión combinando los observaciones de las demás variables del conjunto de datos que acompañan al dato nulo, realizando una regresión ridge (que evita la multicolinealidad) y a partir de este, predice el valor del dato nulo en la variable en cuestión.

Para la variable **num-of-doors**, imputaremos el valor predicho con la moda de la misma variable.

Imputación de valores

```
In [7]: # Primero se quitan los caracteres "?" y se cambian por nulos
```

```
In [8]: # Poramos la conversión de las variables con el carácter "?" a numéricas, para convertirlos estos registros
df["peak-rpm"] = pd.to_numeric(df["peak-rpm"], errors="coerce")
df["price"] = pd.to_numeric(df["price"], errors="coerce")
df["horsepower"] = pd.to_numeric(df["horsepower"], errors="coerce")
df["horsepower"] = pd.to_numeric(df["horsepower"], errors="coerce")
df["bore"] = pd.to_numeric(df["bore"], errors="coerce")
df["stroke"] = pd.to_numeric(df["stroke"], errors="coerce")
```

Observamos los campos que poseen valores nulos y el porcentaje que representa este en la variable

```
In [9]: prop_perdidos = df.isnull().sum() * 100 / len(df)
valores_perdidos_df = pd.DataFrame(prop_perdidos)
valores_perdidos_df
```

```
Out[9]:
```

	Variable	% de perdidos
	make	0.000000
	fuel-type	0.000000
	aspiration	0.000000
	num-of-doors	0.000000
	body-style	0.000000
	drive-wheels	0.000000
	engine-location	0.000000
	wheel-base	0.000000
	length	0.000000
	width	0.000000
	height	0.000000
	curb-weight	0.000000
	engine-type	0.000000
	num-of-cylinders	0.000000
	engine-size	0.000000
	fuel-system	0.000000
	bore	1.951222
	stroke	1.951222
	compression-ratio	0.000000
	horsepower	0.975611
	peak-rpm	0.975611
	city-mpg	0.000000
	highway-mpg	0.000000
	price	1.951222

Viendo los resultados, y dado que el porcentaje de valores perdidos en las variables es inferior al 5%, podríamos optar por eliminar los registros sin más y no habría problema alguno en perder representatividad en el conjunto de datos, pero en este caso se prefiere imputar estos datos, para tener la mayor cantidad de información de los vehículos. Ahora vamos a imputar los valores en las variables usando **imputación iterativa** basada en el algoritmo Bayesiano ridge, que posee la librería **sklearn**.

```
In [10]: # Cargamos la librería
from sklearn.experimental import enable_iterative_imputer
from sklearn.impute import IterativeImputer
```

```
In [11]: # Seleccionamos solo las variables numéricas para realizar la imputación
data = df.select_dtypes(include = "number")
# Guardamos los nombres de las columnas para usarlas posteriormente
lista_cols_num = list(data.columns.values)
```

```
Out[11]:
```

	make	fuel-type	aspiration	num-of-doors	body-style	drive-wheels	engine-location	wheel-base	length	width	engine-size	fuel-system	bore	stroke	compression-ratio
0	alfa-romero	gas	std	two	convertible	nwd	front	88.6	168.8	64.1	130	mpfi	3.47	2.68	9.0
1	alfa-romero	gas	std	two	convertible	nwd	front	88.6	168.8	64.1	130	mpfi	3.47	2.68	9.0
2	alfa-romero	gas	std	two	hatchback	nwd	front	94.5	171.2	65.5	152	mpfi	2.68	3.47	9.0
3	audi	gas	std	four	sedan	fwd	front	99.8	176.6	66.2	109	mpfi	3.19	3.4	10.0
4	audi	gas	std	four	sedan	4wd	front	99.4	176.6	66.4	136	mpfi	3.19	3.4	8.0

5 rows x 24 columns

3.3. Limpieza de los datos

Procedemos ahora a limpiar estas variables y reemplazar los registros faltantes haciendo uso del algoritmo de *regresión Bayesiana Ridge* para el caso de las variables **num-cylinders** y así permitir realizar una regresión combinando los observaciones de las demás variables del conjunto de datos que acompañan al dato nulo, realizando una regresión ridge (que evita la multicolinealidad) y a partir de este, predice el valor del dato nulo en la variable en cuestión.

Para la variable **num-of-doors**, imputaremos el valor predicho con la moda de la misma variable.

Imputación de valores

```
In [12]: # Primero se quitan los caracteres "?" y se cambian por nulos
```

```
In [13]: # Poramos la conversión de las variables con el carácter "?" a numéricas, para convertirlos estos registros
df["peak-rpm"] = pd.to_numeric(df["peak-rpm"], errors="coerce")
df["price"] = pd.to_numeric(df["price"], errors="coerce")
df["horsepower"] = pd.to_numeric(df["horsepower"], errors="coerce")
df["horsepower"] = pd.to_numeric(df["horsepower"], errors="coerce")
df["bore"] = pd.to_numeric(df["bore"], errors="coerce")
df["stroke"] = pd.to_numeric(df["stroke"], errors="coerce")
```

Observamos los campos que poseen valores nulos y el porcentaje que representa este en la variable

```
In [14]: prop_perdidos = df.isnull().sum() * 100 / len(df)
valores_perdidos_df = pd.DataFrame(prop_perdidos)
valores_perdidos_df
```

```
Out[14]:
```

	Variable	% de perdidos
	make	0.000000
	fuel-type	0.000000
	aspiration	0.000000
	num-of-doors	0.000000
	body-style	0.000000
	drive-wheels	0.000000
	engine-location	0.000000
	wheel-base	0.000000
	length	0.000000
	width	0.000000
	height	0.000000
	curb-weight	0.000000
	engine-type	0.000000
	num-of-cylinders	0.000000
	engine-size	0.000000
	fuel-system	0.000000
	bore	1.951222
	stroke	1.951222
	compression-ratio	0.000000
	horsepower	0.975611
	peak-rpm	0.975611
	city-mpg	0.000000
	highway-mpg	0.000000
	price	1.951222

Viendo los resultados, y dado que el porcentaje de valores perdidos en las variables es inferior al 5%, podríamos optar por eliminar los registros sin más y no habría problema alguno en perder representatividad en el conjunto de datos, pero en este caso se prefiere imputar estos datos, para tener la mayor cantidad de información de los vehículos. Ahora vamos a imputar los valores en las variables usando **imputación iterativa** basada en el algoritmo Bayesiano ridge, que posee la librería **sklearn**.

```
In [15]: # Cargamos la librería
from sklearn.experimental import enable_iterative_imputer
from sklearn.impute import IterativeImputer
```

```
In [16]: # Seleccionamos solo las variables numéricas para realizar la imputación
data = df.select_dtypes(include = "number")
# Guardamos los nombres de las columnas para usarlas posteriormente
lista_cols_num = list(data.columns.values)
```

```
Out[16]:
```

	make	fuel-type	aspiration	num-of-doors	body-style	drive-wheels	engine-location	wheel-base	length	width	engine-size	fuel-system	bore	stroke	compression-ratio
0	alfa-romero	gas	std	two	convertible	nwd	front	88.6	168.8	64.1	130	mpfi	3.47	2.68	9.0
1	alfa-romero	gas	std	two	convertible	nwd	front	88.6	168.8	64.1	130	mpfi	3.47	2.68	9.0
2	alfa-romero	gas	std	two	hatchback	nwd	front	94.5	171.2	65.5	152	mpfi	2.68	3.47	9.0
3	audi	gas	std	four	sedan	fwd	front	99.8	176.6	66.2	109	mpfi	3.19	3.4	10.0
4	audi	gas	std	four	sedan	4wd	front	99.4	176.6	66.4	136	mpfi	3.19	3.4	8.0

5 rows x 24 columns

3.3. Limpieza de los datos

Procedemos ahora a limpiar estas variables y reemplazar los registros faltantes haciendo uso del algoritmo de *regresión Bayesiana Ridge* para el caso de las variables **num-cylinders** y así permitir realizar una regresión combinando los observaciones de las demás variables del conjunto de datos que acompañan al dato nulo, realizando una regresión ridge (que evita la multicolinealidad) y a partir de este, predice el valor del dato nulo en la variable en cuestión.

Para la variable **num-of-doors**, imputaremos el valor predicho con la moda de la misma variable.

Imputación de valores

```
In [17]: # Primero se quitan los caracteres "?" y se cambian por nulos
```

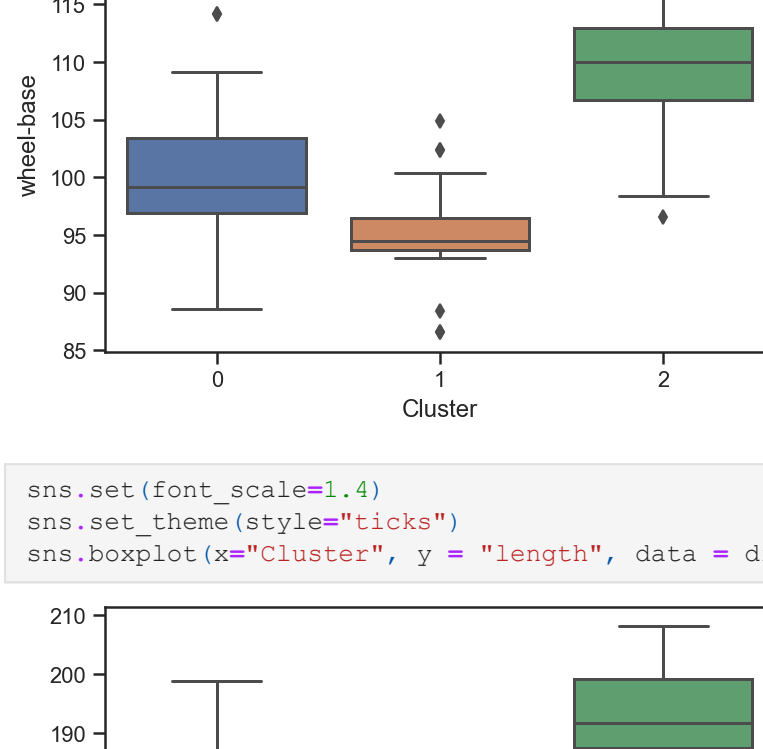
```
In [18]: # Poramos la conversión de las variables con el carácter "?" a numéricas, para convertirlos estos registros
df["peak-rpm"] = pd.to_numeric(df["peak-rpm"], errors="coerce")
df["price"] = pd.to_numeric(df["price"], errors="coerce")
df["horsepower"] = pd.to_numeric(df["horsepower"], errors="coerce")
df["horsepower"] = pd.to_numeric(df["horsepower"], errors="coerce")
df["bore"] = pd.to_numeric(df["bore"], errors="coerce")
df["stroke"] = pd.to_numeric(df["stroke"], errors="coerce")
```

Observamos los campos que poseen valores nulos y el porcentaje que representa este en la variable

```
In [19]: prop_perdidos = df.isnull().sum() * 100 / len(df)
valores_perdidos_df = pd.DataFrame(prop_perdidos)
valores_perdidos_df
```

```
Out[19]:
```

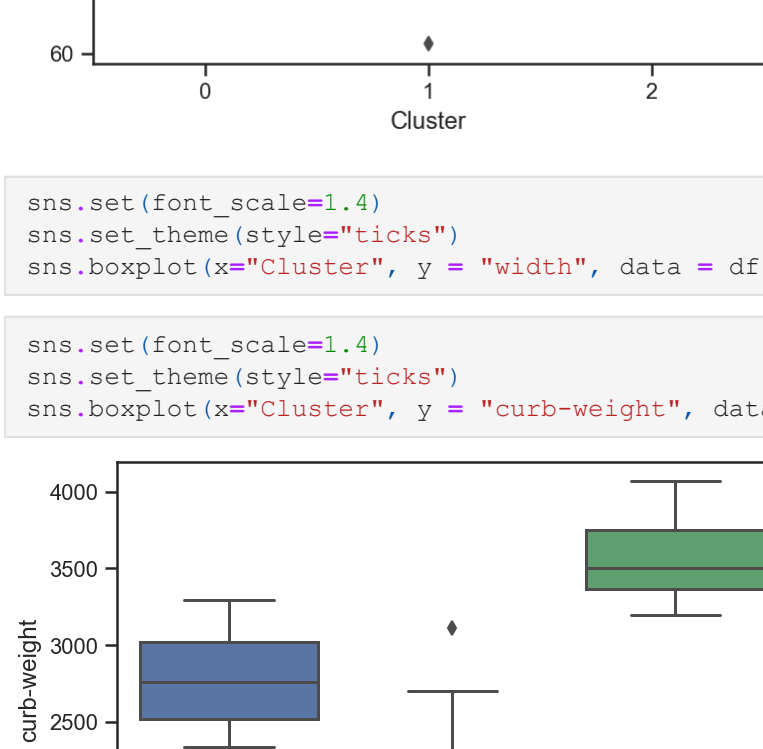
	Variable	% de perdidos
	make	0.000000
	fuel-type	0.000000
	aspiration	0.000000
	num-of-doors	0.000000
	body-style	0.000000
	drive-wheels	0.000000
	engine-location	0.000000
	wheel-base	0.000000
	length	0.000000
	width	0.000000
	height	0.000000
	curb-weight	0.000000
	engine-type	0.000000
	num-of-cylinders	0.000000
	engine-size	0.000000
	fuel-system	0.000000
	bore	1.951222
	stroke	1.951222
	compression-ratio	0.000000
	horsepower	0.975611
	peak-rpm	0.975611
	city-mpg	0.000000
	highway-mpg	0.000000



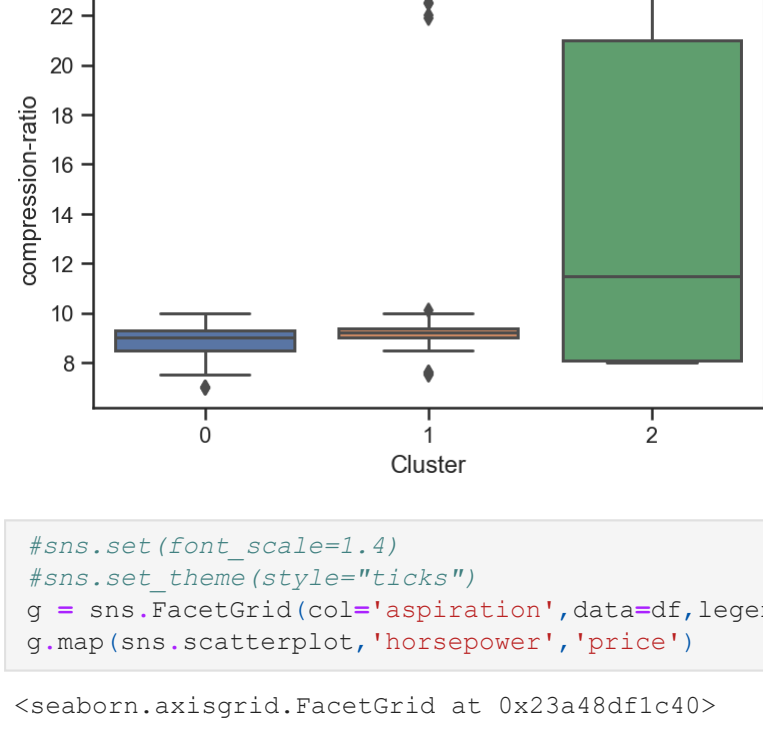
```
In [179]: sns.set(font_scale=1.4)
sns.set_theme(style="ticks")
sns.boxplot(x="Cluster", y = "length", data = df);
```



```
In [181]: sns.set(font_scale=1.4)
sns.set_theme(style="ticks")
sns.boxplot(x="Cluster", y = "height", data = df);
```

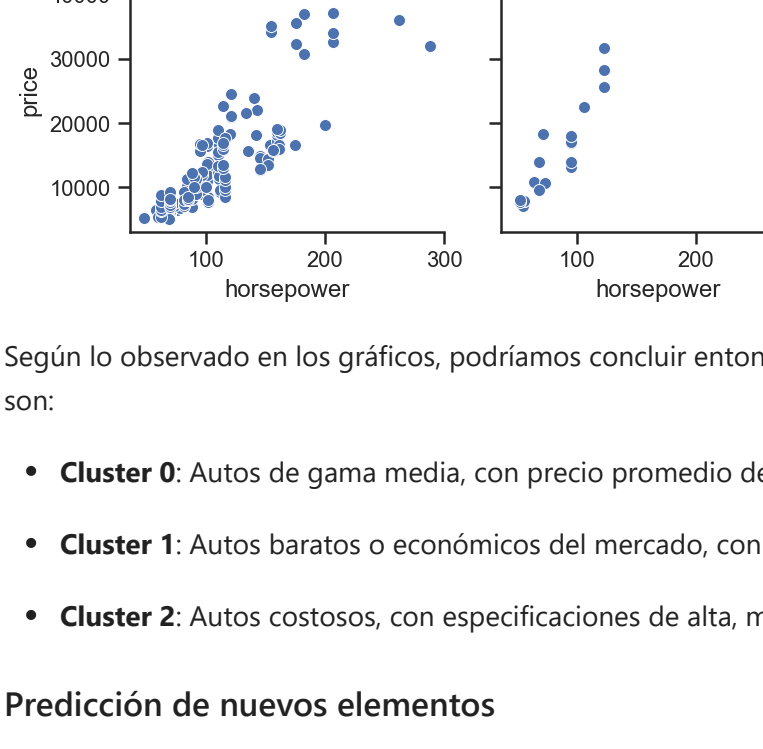


```
In [182]: sns.set(font_scale=1.4)
sns.set_theme(style="ticks")
sns.boxplot(x="Cluster", y = "width", data = df);
```



```
In [ ]: sns.set(font_scale=1.4)
sns.set_theme(style="ticks")
sns.boxplot(x="Cluster", y = "width", data = df);
```

```
In [183]: sns.set(font_scale=1.4)
sns.set_theme(style="ticks")
sns.boxplot(x="Cluster", y = "curb-weight", data = df);
```

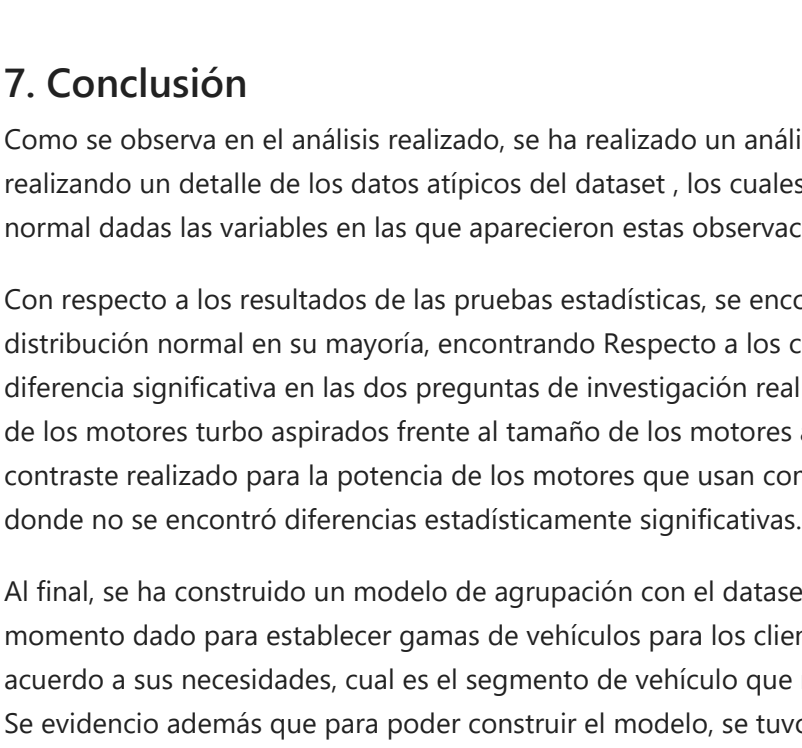


```
In [183]: sns.set(font_scale=1.4)
sns.set_theme(style="ticks")
sns.boxplot(x="Cluster", y = "compression-ratio", data = df);
```



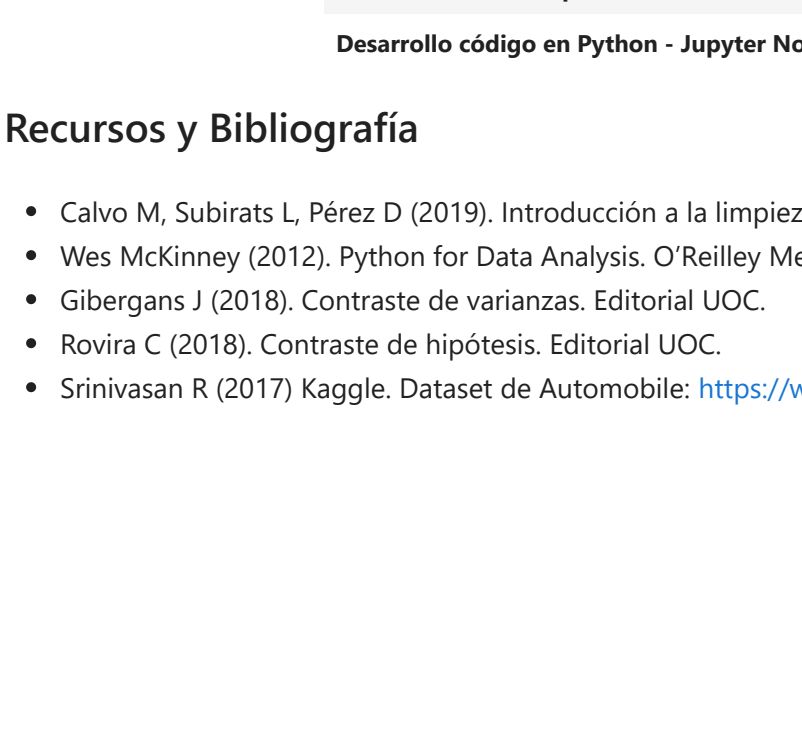
```
In [188]: #sns.set(font_scale=1.4)
#sns.set_theme(style="ticks")
g = sns.FacetGrid(col='aspiration',data=df, legend_out=False)
g.map(sns.scatterplot,'horsepower','price')
```

Out[188]: <seaborn.axisgrid.FacetGrid at 0x23a48df1c40>



```
In [190]: sns.set(font_scale=1.4)
sns.set_theme(style="ticks")
g = sns.FacetGrid(col='fuel-type',data=df, legend_out=False)
g.map(sns.scatterplot,'horsepower','price')
```

Out[190]: <seaborn.axisgrid.FacetGrid at 0x23a45a5e940>



Según lo observado en los gráficos, podríamos concluir entonces sobre los tres grupos de automóviles en tres categorías, las cuales son:

- **Cluster 0:** Autos de gama media, con precio promedio del mercado y especificaciones de gama intermedia en general
- **Cluster 1:** Autos baratos o económicos del mercado, con bajas prestaciones y en general de dimensiones pequeñas
- **Cluster 2:** Autos costosos, con especificaciones de alta, motores más grandes y mayor tamaño en general

Predicción de nuevos elementos

Ahora se proponer clasificar nuevos automóviles según sus especificaciones en alguno de los clusters, usando el modelo ya construido.

```
In [186]: #Cargue de datos nuevos
df_new = pd.read_csv('automobile_nuevo.csv', sep = ';')
df_new.head()
```

	make	fuel-type	aspiration	num-of-doors	body-style	drive-wheels	engine-location	wheel-base	length	width	engine-size	fuel-system	bore	stroke	compression-ratio	horsepower
0	bmw	gas	std	two	sedan	rwd	front	101.2	176.8	64.8	...	108	mpfi	3.50	2.80	8.8
1	peugot	gas	std	four	sedan	rwd	front	107.9	186.7	68.4	...	120	mpfi	3.46	3.19	8.4
2	volvo	diesel	turbo	four	sedan	rwd	front	109.1	188.8	68.9	...	145	idi	3.01	3.40	23.0

3 rows × 24 columns

```
In [225]: # extraemos solo las variables numericas del nuevo grupo de datos
X_new=df_new[num_list]
X_std_new = StandardScaler().fit_transform(X_new)
```

```
In [234]: #Cargamos y aplicamos los modelos de PCA y de Kmedias realizados anteriormente

# PCA
pca_reload = pk.load(open("pca.pkl", "rb"))
pca_new = pca_reload.transform(X_std_new)
pca_new_df = pd.DataFrame(pca_new)

# Kmedias
kmeans = pk.load(open("kmeans.pkl", "rb"))
kmeans_new = kmeans.predict(pca_new_df.iloc[:, :3])

# añadimos los clusters nuevos al dataset
df_new['Cluster'] = kmeans_new
```

In [235]: df_new

	make	fuel-type	aspiration	num-of-doors	body-style	drive-wheels	engine-location	wheel-base	length	width	fuel-system	bore	stroke	compression-ratio	horsepower	
0	bmw	gas	std	two	sedan	rwd	front	101.2	176.8	64.8	...	mpfi	3.50	2.80	8.8	101
1	peugot	gas	std	four	sedan	rwd	front	107.9	186.7	68.4	...	mpfi	3.46	3.19	8.4	97
2	volvo	diesel	turbo	four	sedan	rwd	front	109.1	188.8	68.9	...	idi	3.01	3.40	23.0	132

3 rows × 25 columns

Los nuevos autos aparecen clasificados en el cluster 1 el primero y los dos restantes en el cluster 0, indicando que son vehículos de gama media el primero y los otros dos son autos de gama baja o con prestaciones de rango bajo.

7. Conclusión

Como se observa en el análisis realizado, se ha realizado un análisis del dataset, primero realizando una limpieza de las variables, realizando un detalle de los datos típicos del dataset, los cuales al final se han decidido mantener estos registros, dato que parecían normal dadas las variables en las que aparecieron estas observaciones anómalas.

Con respecto a los resultados de las pruebas estadísticas, se encontró que las variables continuas evaluadas no presentaron una distribución normal en su mayoría, encontrando Respecto a los contrastes de hipótesis establecidos, pudimos ver que no hubo diferencia significativa en las dos preguntas de investigación realizadas, las cuales arrojaron que no existe una diferencia en el tamaño de los motores turbo aspirados frente al tamaño de los motores aspirados convencionalmente; este resultado también se reflejó en el contraste realizado para la potencia de los motores que usan combustible diesel frente a la potencia de los motores que usan gasolina, donde no se encontró diferencias estadísticamente significativas.

Al final, se ha construido un modelo de agrupación con el dataset, para encontrar segmentos de autos que pudiesen servir en un momento dado para establecer gamas de vehículos para los clientes en un concesionario o incluso para el cliente poder conocer, de acuerdo a sus necesidades, cual es el segmento de vehículo que mejor se adapta a su necesidad puntual al momento de adquirir uno. Se evidencio además que para poder construir el modelo, se tuvo que realizar una serie de transformaciones en los datos, como una estandarización (para evitar que variables con escalas grandes afectaran el modelo) y un análisis de componentes principales, los cuales fueron importantes para obtener un modelo más consistente, el cual produce resultados interesantes y que permite clasificar nuevos autos, como lo vimos en el caso de la predicción, entregando el resultado de manera sencilla y eficiente, lo cual sería muy útil a futuro para un distribuidor de vehículos poder segmentarlos usando este tipo de modelos, basado en las especificaciones de los vehículos.

ENLACE AL GITHUB

https://github.com/JhonHarry/automobile_data_kaggle_limpieza

Contribuciones de los integrantes del grupo en la práctica

Contribuciones	Firma
Investigación acerca de la viabilidad del dataset para la práctica	JHLG, PJSV
Redacción de las respuestas del documento final en markdown y conclusiones	JHLG, PJSV
Desarrollo código en Python - Jupyter Notebook, dataset final	JHLG, PJSV

Recursos y Bibliografía

- Calvo M, Subirats L, Pérez D (2019). Introducción a la limpieza y análisis de los datos. Editorial UOC.
- Wes McKinney (2018). Python for Data Analysis. O'Reilly Media, Inc.
- Gibergans J (2018). Contraste de varianzas. Editorial UOC.
- Rovina C (2018). Contraste de hipótesis. Editorial UOC.
- Srinivasan R (2017) Kaggle. Dataset de Automobile: <https://www.kaggle.com/toramky/automobile-dataset>