

PRACTICA 2 - Limpieza y análisis de datos

Presentada por Jhon Harry Loaliza y Pablo Jesús Sánchez

Tipología y ciclo de vida del dato - Aula 3

18 dic 2020

Contenido:

1. Descripción del dataset

1. Importancia 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 objeto analizado se llama *Automobile.dats.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 aspirador que posee el motor del vehículo
- num-of-doors**: número de puertas
- body-style**: tipo de 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 de combustible 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 permitan reconocer y clasificar nuevos vehículos en cada uno de estos segmentos, lo cual también ayudará 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

import warnings
warnings.filterwarnings('ignore')

plt.style.use('ggplot')
```

```
In [2]: # Cargue de los datos
df = pd.read_csv('Automobile.dats.csv')
# Verifica que los datos no sean nulos
df.head()
```

height	height	0.00000
curb-weight	curb-weight	0.00000
engine-type	engine-type	0.00000
num-of-cylinders	num-of-cylinders	0.00000
engine-size	engine-size	0.00000
fuel-system	fuel-system	0.00000
bore	bore	1.95122
stroke	stroke	1.95122
compression-ratio	compression-ratio	0.00000
horsepower	horsepower	0.97561
peak-rpm	peak-rpm	0.97561
city-mpg	city-mpg	0.00000
highway-mpg	highway-mpg	0.00000
price	price	1.95122

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 numéricas usando **imputación iterativa** basada en algoritmo Bayesian ridge, que posee la librería *sklearn*.

```
# Cargamos la librería
from sklearn.experimental import enable_iterative_imputer
from sklearn.impute import IterativeImputer
```

5 rows x 24 columns

```
In [3]: # Dimensión 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   engine-size          205 non-null    int64
17   fuel-system          205 non-null    object
18   horsepower           205 non-null    object
19   stroke              205 non-null    object
20   compression-ratio    205 non-null    float64
21   horsepower           205 non-null    object
22   peak-rpm            205 non-null    object
23   city-mpg             205 non-null    int64
24   highway-mpg          205 non-null    int64
25   price               205 non-null    object
dtypes: float64(5), int64(5), object(16)
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 actuariales 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()
```

	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	rwd	front	88.6	168.8	64.1	...	130	mpfi	3.47	2.68	9.0
1	alfa-romero	gas	std	two	convertible	rwd	front	88.6	168.8	64.1	...	130	mpfi	3.47	2.68	9.0
2	alfa-romero	gas	std	two	hatchback	rwd	front	94.5	171.2	65.5	...	152	mpfi	2.68	3.47	9.0
3	audi	gas	std	four	sedan	4wd	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 "o" en forma incorrecta, por ejemplo en el caso de las variables *price* y *horsepower*, vemos que aparecen unos "0" 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()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 205 entries, 0 to 204
Data columns (total 15 columns):
 #   Column              Non-Null Count  Dtype
---  --
 0   make                 205 non-null    object
 1   fuel-type            205 non-null    object
 2   aspiration            205 non-null    object
 3   num-of-doors         205 non-null    object
 4   body-style           205 non-null    object
 5   drive-wheels         205 non-null    object
 6   engine-location      205 non-null    object
 7   engine-type          205 non-null    object
 8   num-of-cylinders     205 non-null    object
 9   fuel-system          205 non-null    object
10   horsepower           205 non-null    object
11   stroke              205 non-null    object
12   compression-ratio    205 non-null    object
13   peak-rpm            205 non-null    object
14   price               205 non-null    object
dtypes: object(15)
memory usage: 24.1 KB

None
```

En este caso se encuentra que, según la definición de las variables, las que tiene errado el tipo de variable son: *horsepower*, *peak-rpm*, *price*, *bore*, *stroke*, y entre las variables que son de tipo categórico pero que poseen el carácter "0" en sus registros, encontramos *num-of-doors*.

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 **cuantitativas** 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 perdido 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['stroke'] = pd.to_numeric(df['stroke'], errors='coerce')
df['bore'] = pd.to_numeric(df['bore'], 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.isnull().sum()
valores_perdidos = df.isnull().sum()

valores_perdidos_df
```

	Variable	% de perdidos
	make	0.00000
	fuel-type	0.00000
	aspiration	0.00000
	num-of-doors	0.00000
	body-style	0.00000
	drive-wheels	0.00000
	engine-location	0.00000
	wheel-base	0.00000
	length	0.00000
	width	0.00000
	height	0.00000
	curb-weight	0.00000
	engine-type	0.00000
	num-of-cylinders	0.00000
	engine-size	0.00000
	fuel-system	0.00000
	bore	1.95122
	stroke	1.95122
	compression-ratio	0.00000
	horsepower	0.97561
	peak-rpm	0.97561
	city-mpg	0.00000
	highway-mpg	0.00000
	price	1.95122

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 representativos 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 **datos imputación iterativa** basada en algoritmo Bayesiano *ridge*, que posee la librería *sklearn*.

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

In [11]: # Seleccionamos solo las variables numericas 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)

# Extraemos el array con los valores de las variables escogidas
data = data.values
data.shape
```

(205, 14)

```
In [12]: # Extraemos el rango de columnas a imputar para crear la lista de valores
ix = [i for i in range(data.shape[1]) if i != 14]
X = data[:, ix]
```

```
In [13]: # Cantidad de valores nulos antes de imputar
print("Datos perdidos antes de imputar: %d" % sum(np.isnan(X).flatten()))

# Establecemos el imputador, por defecto usa el algoritmo BayesianoRidge
imputer = IterativeImputer()

# Aplicamos el algoritmo a los datos
imputer.fit(X)

# Transformamos el dataset
Xtrans = imputer.transform(X)

# Total de perdidos después de la imputación
print("Datos perdidos después de imputar: %d" % sum(np.isnan(Xtrans).flatten()))

Datos perdidos antes de imputar: 16
Datos perdidos después de imputar: 0
```

Luego de imputar los datos en las variables numéricas, procedemos a convertir el resultado en *dataframe* y unirlo con las variables categoricas del dataset original.

```
In [14]: # Convertimos los valores de las variables sin nulos en un dataframe
df_num = pd.DataFrame(Xtrans, columns=lista_cols_num)

# Extraemos las variables categoricas del dataset original
df_str = df.select_dtypes(include = ["object"])

# Creamos un nuevo dataset con las variables imputadas y las variables categoricas
df_1 = pd.concat([df_str, df_num], axis=1)
df_1.shape
```

(205, 24)

Para la variable *num-of-doors* podemos tomar dos opciones, imputar los valores con la moda de la variable o al ser tan pocos registros nulos, podrían eliminarse, en este caso optamos por la opción de imputarlos usando la moda de la variable

```
In [15]: # Calculamos la moda de la variable y reemplazamos el valor nulo con esta
df['num-of-doors'] = df['num-of-doors'].mode() # El resultado es 'four'
df['num-of-doors'] = df['num-of-doors'].replace('?', 'four').astype('object')
```

```
In [16]: # Comprobamos que ha quedado OK el cambio
df['num-of-doors'].unique()
```

array(['two', 'four'], dtype=object)

```
In [17]: # Modificamos el nombre del nuevo dataset para mayor comodidad al momento de trabajar
df = df_1
```

3.2. Identificación y tratamiento de valores extremos

Para identificar si existen valores extremos o atípicos en los datos, usaremos visualizaciones como los boxplots y el criterio del rango intercuartílico para identificar los extremos de las variables numéricas.

Gráfico de boxplot para las variables numéricas

```
In [18]: plt.figure(figsize=(12,10))
sns.set(font_scale=4)
df.boxplot(vert=False)
```

Out [18]: <AxesSubplot>

Observando el gráfico de cajas para las variables numéricas, podría decirse que las variables *price*, *horsepower*, *engine-size* y *peak-rpm* poseen valores por fuera del nivel superior del gráfico de cajas ($Q3 + 1.5 \times IQR$), aunque también se advierte que no son muy alejados estos valores, salvo algunos valores de la variable *price*. Pasamos a analizar cada una de estas variables para ver su comportamiento y si estos valores atípicos pueden ser justificados.

Variable price

```
In [19]: # Extraemos los valores que están por encima o por debajo del límite del diagrama de cajas
Q1 = df['price'].quantile(0.25)
Q3 = df['price'].quantile(0.75)
IQR = Q3 - Q1

df[(df['price'] < Q1-1.5*IQR) | (df['price'] > Q3+1.5*IQR)][['price']]
```

```
Out [19]:   price
15    30760.000000
16    30760.000000
17    36880.000000
18    32250.000000
48    35550.000000
49    36000.000000
50    31600.000000
71    34184.000000
72    35056.000000
73    40360.000000
74    45400.000000
75    52528.000000
76    5000.000000
126   37028.000000
127   22004.000000
Name: price, dtype: float64
```

Revisemos los valores que están iguales o por encima de 32000 USD (el cual aparece como valor atípico anteriormente):

```
In [20]: df[(df['price'] >= 32000)].head(10)
```

	make	fuel-type	aspiration	num-of-doors	body-style	drive-wheels	engine-location	engine-type	num-of-cylinders	fuel-system	curb-weight	engine-size	bore	stroke	compression-ratio
16	bmw	gas	std	two	sedan	rwd	front	ohc	six	mpfi	...	3380.0	209.0	3.62	3.39
17	bmw	gas	std	four	sedan	rwd	front	ohc	six	mpfi	...	3505.0	209.0	3.62	3.39
48	jaguar	gas	std	four	sedan	rwd	front	dohc	six	mpfi	...	4066.0	258.0	3.63	4.17
49	jaguar	gas	std	two	sedan	rwd	front	ohcv	twelve	mpfi	...	3950.0	326.0	3.54	2.76
71	mercedes-benz	gas	std	four	sedan	rwd	front	ohcv	eight	mpfi	...	3740.0	234.0	3.46	3.10
72	mercedes-benz	gas	std	two	convertible	rwd	front	ohcv	eight	mpfi	...	3680.0	304.0	3.60	3.10
73	mercedes-benz	gas	std	four	sedan	rwd	front	ohcv	eight	mpfi	...	3950.0	326.0	3.80	3.35
74	mercedes-benz	gas	std	two	hatchtop	rwd	front	ohcv	eight	mpfi	...	3715.0	304.0	3.80	3.35
126	porche	gas	std	two	hatchtop	rwd	rear	ohf	six	mpfi	...	2756.0	194.0	3.74	2.90

10 rows x 24 columns

Como se observa en los datos, al revisar las marcas de estos vehículos con precios superiores al límite superior del *complot*, vemos que corresponde a fabricantes como Porsche, BMW, Mercedes-Benz o Jaguar, marcas que son reconocidas por construir vehículos de alta gama y que además poseen prestaciones superiores a los autos comunes, como mayor potencia, mayor tamaño de motor, por lo que los valores extremos de esta variable parecen razonables al ser marcas costosas y de renombre en el mercado.

Variable engine-size

```
In [21]: # Extraemos los valores que están por encima o por debajo del límite del diagrama de cajas
Q1 = df['engine-size'].quantile(0.25)
Q3 = df['engine-size'].quantile(0.75)
IQR = Q3 - Q1

df[(df['engine-size'] < Q1-1.5*IQR) | (df['engine-size'] > Q3+1.5*IQR)][['engine-size']]
```

```
Out [21]:   engine-size
15    209.0
16    209.0
17    209.0
47    258.0
48    258.0
49    234.0
72    234.0
73    308.0
74    304.0
Name: engine-size, dtype: float64
```

Analizamos el valor del tamaño del motor por encima de 300 pulgadas cúbicas (4.916 cm cúbicos):

```
In [22]: df[(df['engine-size'] >= 300)].head(10)
```

	make	fuel-type	aspiration	num-of-doors	body-style	drive-wheels	engine-location	engine-type	num-of-cylinders	fuel-system	curb-weight	engine-size	bore	stroke	compression-ratio
49	jaguar	gas	std	two	sedan	rwd	front	ohcv	twelve	mpfi	...	3950.0	326.0	3.54	2.76
73	mercedes-benz	gas	std	four	sedan	rwd	front	ohcv	six	mpfi	...	3900.0	308.0	3.80	3.35
74	mercedes-benz	gas	std	two	hatchtop	rwd	front	ohcv	eight	mpfi	...	3715.0	304.0	3.80	3.35

3 rows x 24 columns

Apreciando los valores extremos de la variable *engine-size*, los cuales son 3 que superan el valor de 300 pulgadas cúbicas, vemos que corresponden a marcas de vehículos importantes, que además poseen valores altos de *peak-rpm* y *compression-ratio*, esto dado por el tamaño del motor y por tener una mayor potencia, los cuales parecen valores razonables en la variable.

Variable horsepower

```
In [23]: # Extraemos los valores que están por encima o por debajo del límite del diagrama de cajas
Q1 = df['horsepower'].quantile(0.25)
Q3 = df['horsepower'].quantile(0.75)
IQR = Q3 - Q1

df[(df['horsepower'] < Q1-1.5*IQR) | (df['horsepower'] > Q3+1.5*IQR)][['horsepower']]
```

```
Out [23]:   horsepower
49    262.0
126    20.0
127    20.0
128    20.0
129    288.0
Name: horsepower, dtype: float64
```

Out [24]: df[(df['horsepower'] >= 210)]

	make	fuel-type	aspiration	num-of-doors	body-style	drive-wheels	engine-location	engine-type	num-of-cylinders	fuel-system	curb-weight	engine-size	bore	stroke	compression-ratio
49	jaguar	gas	std	two	sedan	rwd	front	ohcv	twelve	mpfi	...	3950.0	326.0	3.54	2.76
129	porche	gas	std	two	hatchback	rwd	front	ohcv	eight	mpfi	...	3366.0	203.0	3.94	3.11

2 rows x 24 columns

Revisando los valores extremos de la variable *horsepower*, se puede ver que los valores están por encima de

Las hipótesis a contrastar son:

$$H_0: \mu_1 = \mu_2$$
$$H_1: \mu_1 > \mu_2$$

Donde μ_1 es la media de potencia de motores a Diesel y μ_2 es la media de potencia de motores a Gasolina.

```
In [71]: # Calculamos la media y la desviación para los dos grupos
diesel_mot_media = round(diesel_mot.mean(),1)
diesel_mot_dev = round(diesel_mot.std(),1)
gasolina_mot_media = round(gasolina_mot.mean(),1)
gasolina_mot_dev = round(gasolina_mot.std(),1)
print("Media potencia motor Diesel:",diesel_mot_media)
print("Desviación potencia motor Diesel:",diesel_mot_dev)
print("Media potencia motor Gasolina:",gasolina_mot_media)
print("Desviación potencia motor Gasolina:",gasolina_mot_dev)

# Calculamos el estadístico de contraste
ttest_pval = ttest_ind(diesel_mot, gasolina_mot, axis=0, equal_var=False)

# Ya que el comando ttest_ind genera únicamente el pvalor para un contraste de dos colas,
# se calcula el de una sola cola a partir de este
pval_unacola = 1-pval/2

print("p-value",pval_unacola)
if pval_unacola < 0.05:
    print("Se rechaza la hipótesis nula")
else:
    print("Se acepta la hipótesis nula")
```

Media potencia motor Diesel: 84.5
Desviación potencia motor Diesel: 26.0
Media potencia motor Gasolina: 106.5
Desviación potencia motor Gasolina: 40.2
p-value 0.99126930502611
Se acepta la hipótesis nula

Dado el resultado del contraste de hipótesis, se determina que no existe evidencia significativa de que la potencia del motor diesel sea mayor a la de un motor a gasolina.

5. Visualización de los datos

Para el caso de visualizar los datos, se revisará la distribución de las variables continuas y su posible relación entre ellas, además de visualizar esta información respecto a las variables categóricas.

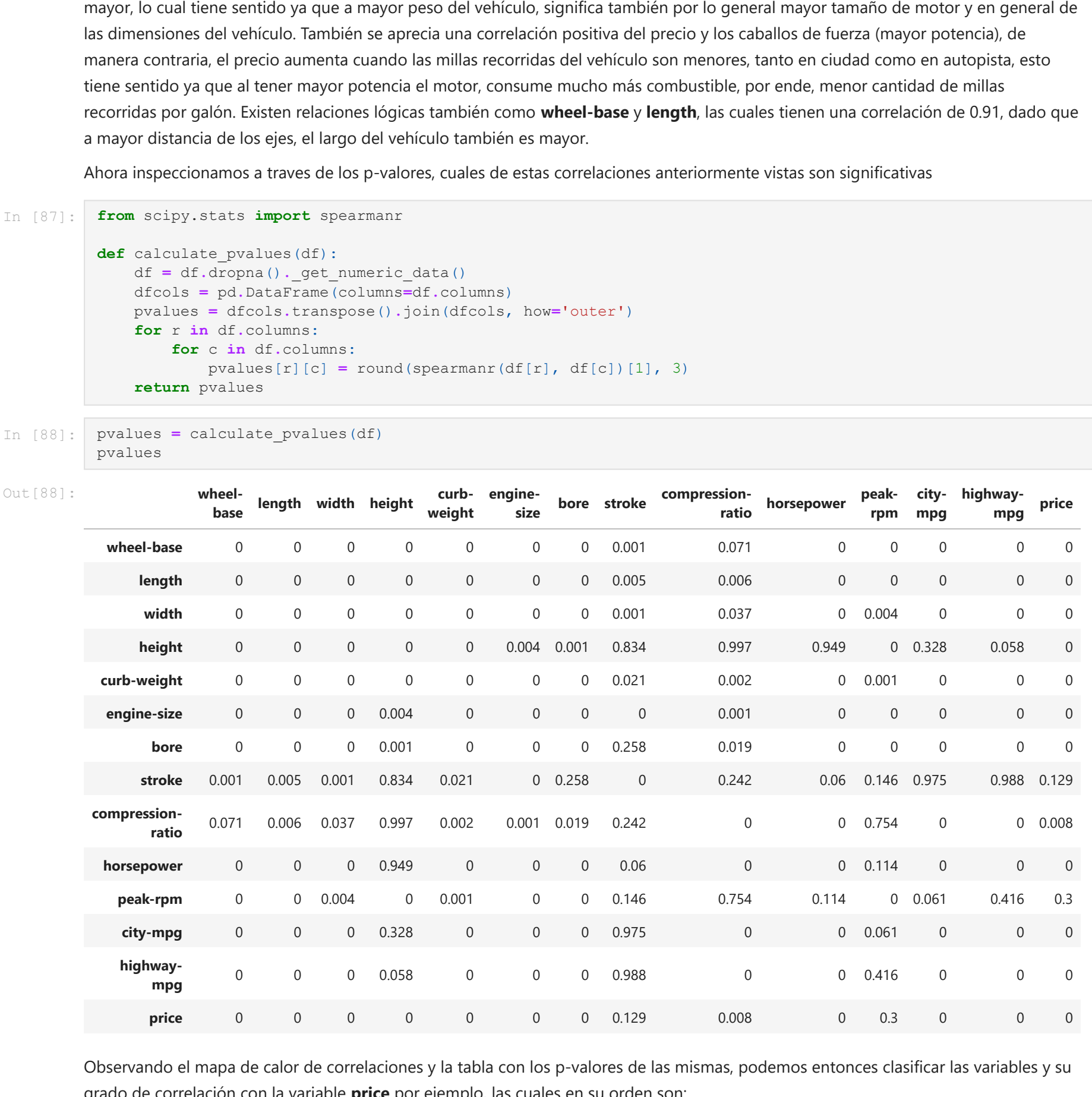
Histogramas y gráficos de distribución

Para los gráficos y visualizaciones, partiremos analizando las variables continuas, observando la correlación existente entre estas, esto con el fin de hacer transformaciones o hacer ajustes en caso de ser necesario.

Evaluaremos en primera instancia la matriz de correlaciones de las variables numéricas, para esto se usará la correlación de spearman, esta se escoge debido a que es un contraste no paramétrico, apropiado para el caso dado que algunas variables no se distribuyen de manera normal. Se procede a construir la matriz de correlaciones.

```
In [36]: #Calculamos la correlación de las variables
correlation_df = df.corr(method='spearman')
```

```
In [38]: # Mapa de correlaciones para las variables continuas
sns.set_theme(style='ticks')
plt.figure(figsize=(12,10))
sns.heatmap(correlation_df, annot=True, cmap='viridis', annot_kws={"size": 12})
plt.show()
```



De este mapa de calor, se puede ver que existe múltiples correlaciones importantes en las variables, muchas de esta por la naturaleza de las mismas, como el caso por ejemplo de la variable **length** del carro con el **ancho**, la cual es una correlación del 89%, lo mismo podemos apreciar de las variables **city-mpg** vs **horsepower**, con correlación de -0.91, la cual tiene sentido, a mayor potencia, menores millas recorridas en ciudad (y también en highway). Otra relación es que precio del automóvil crece cuando el peso vacío (**curb-weight**) es menor, lo cual tiene sentido ya que a menor peso del vehículo, significa también por lo general mayor tamaño de motor y en general de las dimensiones del vehículo. También se aprecia una correlación positiva del precio y los caballos de fuerza (mayor potencia), de manera contraria, el precio aumenta cuando las millas recorridas del vehículo son menores, tanto en ciudad como en autopista, esto tiene sentido ya que al tener mayor potencia el motor, consume mucho más combustible, por ende, menor cantidad de millas recorridas por galón. Existen relaciones lógicas también como **wheel-base** vs **length**, las cuales tienen una correlación de 0.91, dado que a mayor distancia de los ejes, el largo del vehículo también es mayor.

Ahora inspeccionamos a través de los p-valores, cuales de estas correlaciones anteriormente vistas son significativas

```
In [87]: from scipy.stats import spearmanr

def calculate_pvalues(df):
    df = df.dropna().get_numeric_data()
    dfcols = pd.DataFrame(columns=df.columns)
    pvalues = dfcols.transpose().join(dfcols, how='outer')
    for r in df.columns:
        for c in df.columns:
            pvalues[c][r] = round(spearmanr(df[r], df[c])[1], 3)
    return pvalues

pvalues = calculate_pvalues(df)
pvalues
```

```
Out[88]:
```

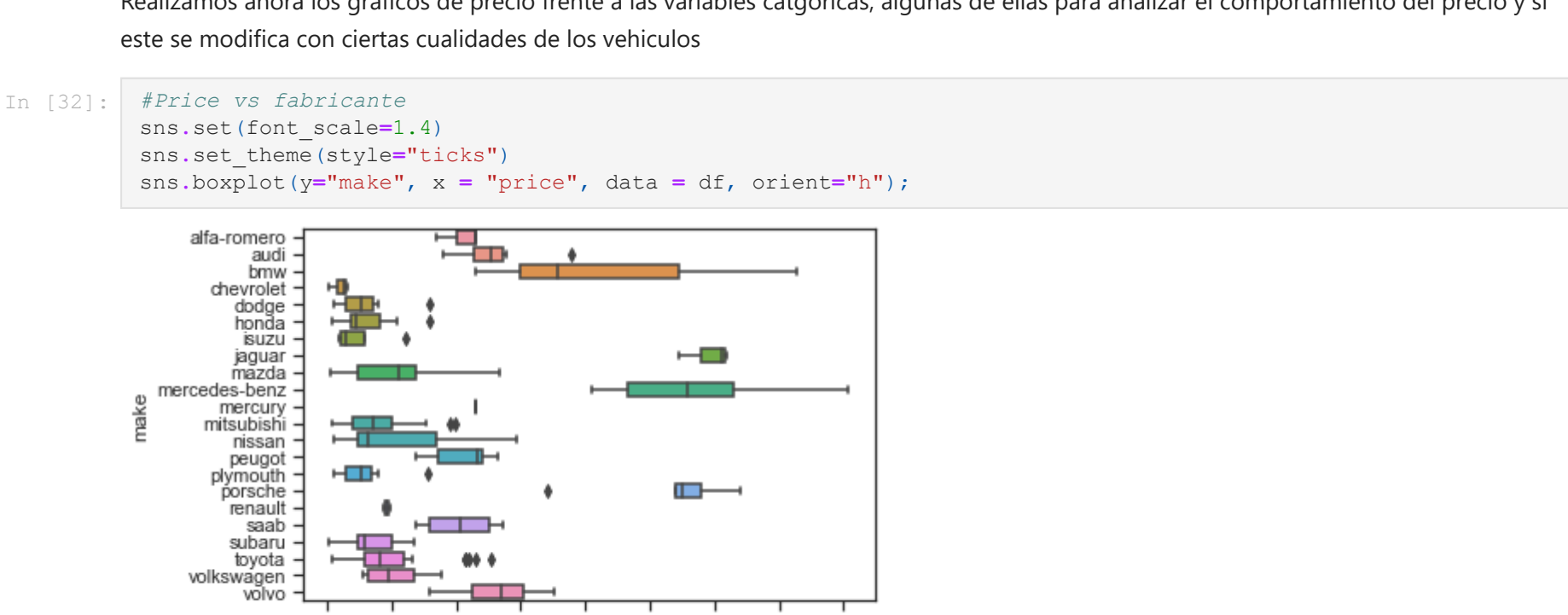
	wheel-base	length	width	height	curb-weight	engine-size	bore	stroke	compression-ratio	horsepower	peak- rpm	city- mpg	highway- mpg	price
wheel-base	0	0	0	0	0	0	0	0	0.001	0.071	0	0	0	0
length	0	0	0	0	0	0	0	0	0.005	0.006	0	0	0	0
width	0	0	0	0	0	0	0	0	0.001	0.037	0	0.004	0	0
height	0	0	0	0	0	0.004	0.001	0.834	0.997	0.949	0	0.328	0.058	0
curb-weight	0	0	0	0	0	0	0	0.021	0.002	0	0.001	0	0	0
engine-size	0	0	0	0	0.004	0	0	0	0.001	0	0	0	0	0
bore	0	0	0	0.001	0	0	0	0.258	0	0.019	0	0	0	0
stroke	0.001	0.005	0.001	0.834	0.021	0	0.258	0	0.242	0.06	0.146	0.975	0.988	0.029
compression-ratio	0.071	0.006	0.037	0.997	0.002	0.001	0.019	0.242	0	0	0.754	0	0	0.008
horsepower	0	0	0	0.949	0	0	0	0.06	0	0	0.114	0	0	0
peak-rpm	0	0	0.004	0	0.001	0	0	0.146	0.754	0.114	0	0.061	0.416	0.3
city-mpg	0	0	0	0.328	0	0	0	0.975	0	0	0.061	0	0	0
highway-mpg	0	0	0	0.058	0	0	0	0.988	0	0	0.416	0	0	0
price	0	0	0	0	0	0	0.029	0.008	0	0.3	0	0	0	0

Observando el mapa de calor de correlaciones y la tabla con los p-valores de las mismas, podemos entonces clasificar las variables y su grado de correlación con la variable **price** por ejemplo, las cuales en su orden son:

Variables con correlación ALTA

- curb-weight (0.92)
- horsepower (0.85)
- city-mpg (0.84)
- highway-mpg (0.83)
- engine-size (0.83)
- width (0.82)
- length (0.81)

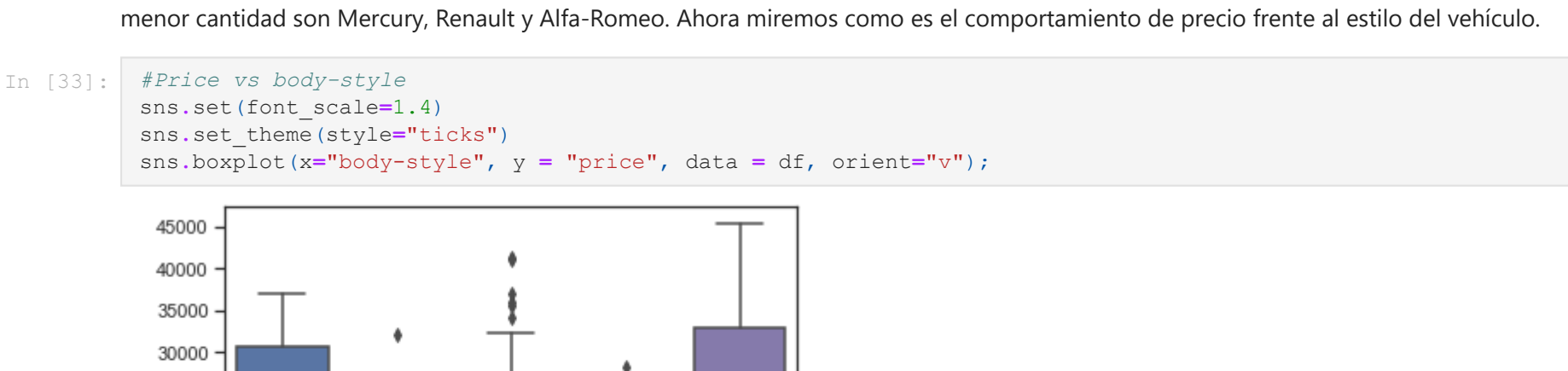
Construimos una matriz de diagramas de dispersión con las variables con alta correlación entre si y apreciar mejor su relación de manera visual



De los diagramas de dispersión, podemos ver unas relaciones fuertes entre las variables **city-mpg** y **highway-mpg**, siendo está de 0.97, y tiene sentido que las millas por galón en ciudad estén relacionadas con las millas por galón en autopista. También se aprecia una relación importante en las variables de **length** y **width**, con 0.89, la cual también tiene lógica ya que un carro entre mayor sea su largo, muy probablemente su ancho también lo sea.

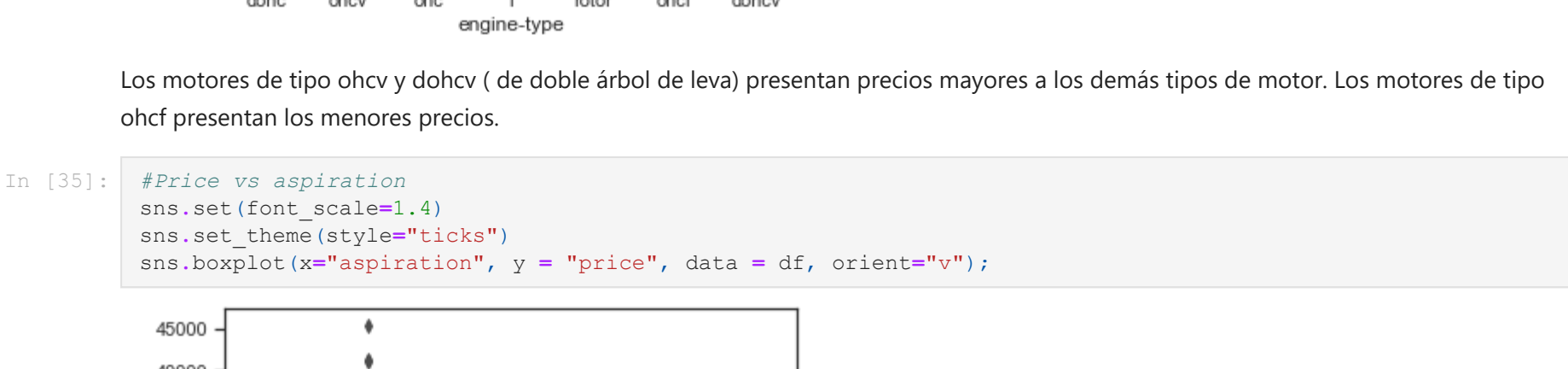
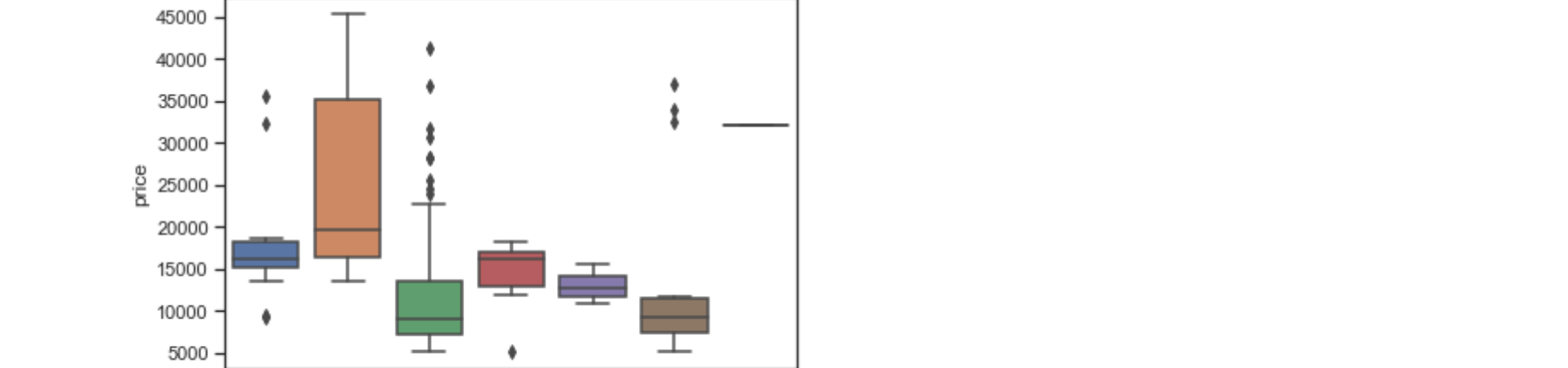
Realizamos ahora los gráficos de precio frente a las variables categóricas, algunas de ellas para analizar el comportamiento del precio y si este se modifica con ciertas cualidades de los vehículos

```
In [27]: #Price vs fabricante
sns.set(font_scale=1.4)
sns.set_theme(style='ticks')
sns.boxplot(x="make", y = "price", data = df, orient="h")
```

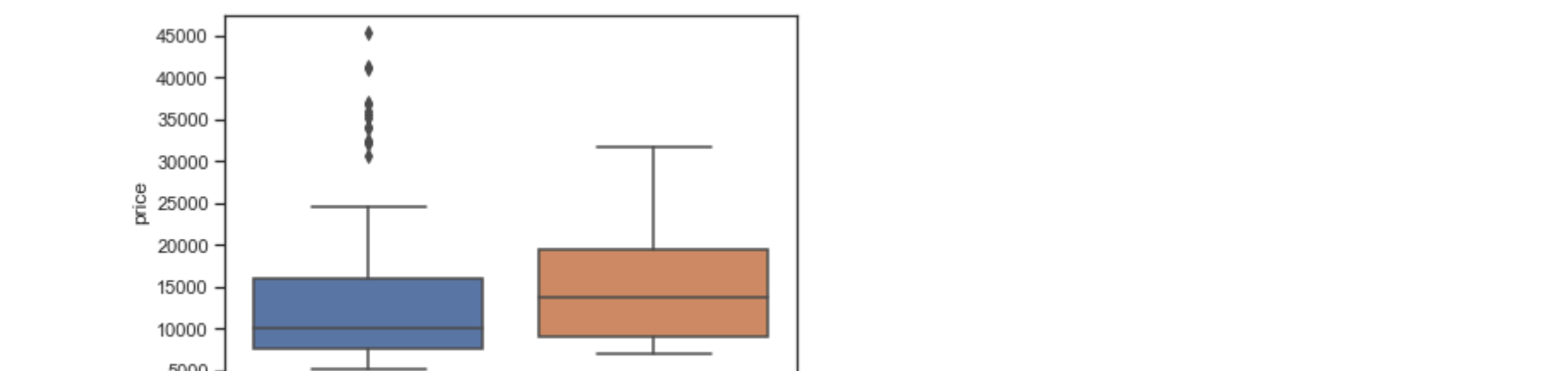
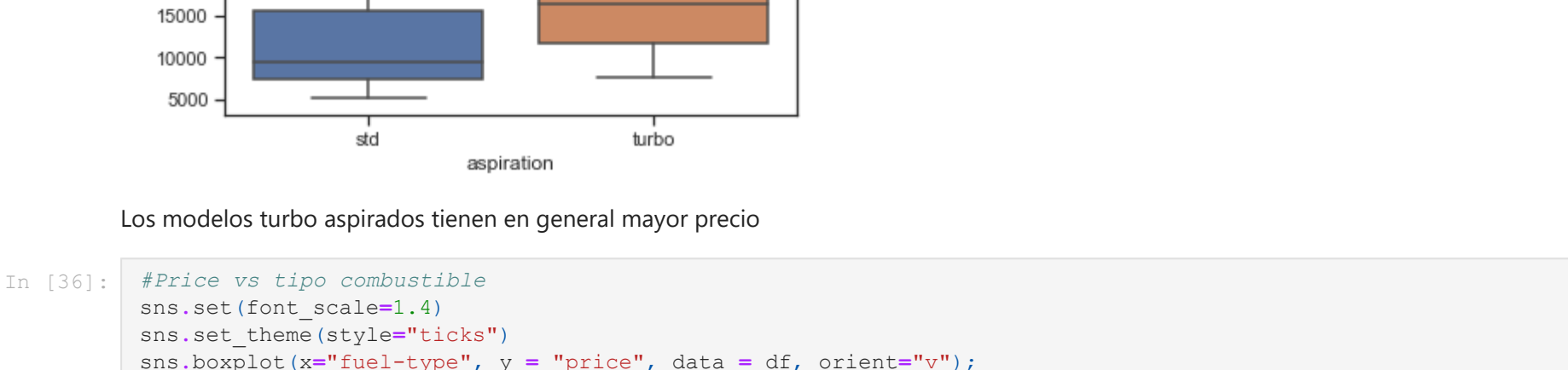


Como se vio en el análisis de los datos extremos de precio, la marca influye en el valor, se aprecian valores altos en marcas como Porsche, Mercedes Benz, BMW, Jaguar y algunas modelos de Audi. En contraparte, en las marcas con precios bajos encontramos la Chevrolet, Dodge, Plymouth, Honda e Isuzu.

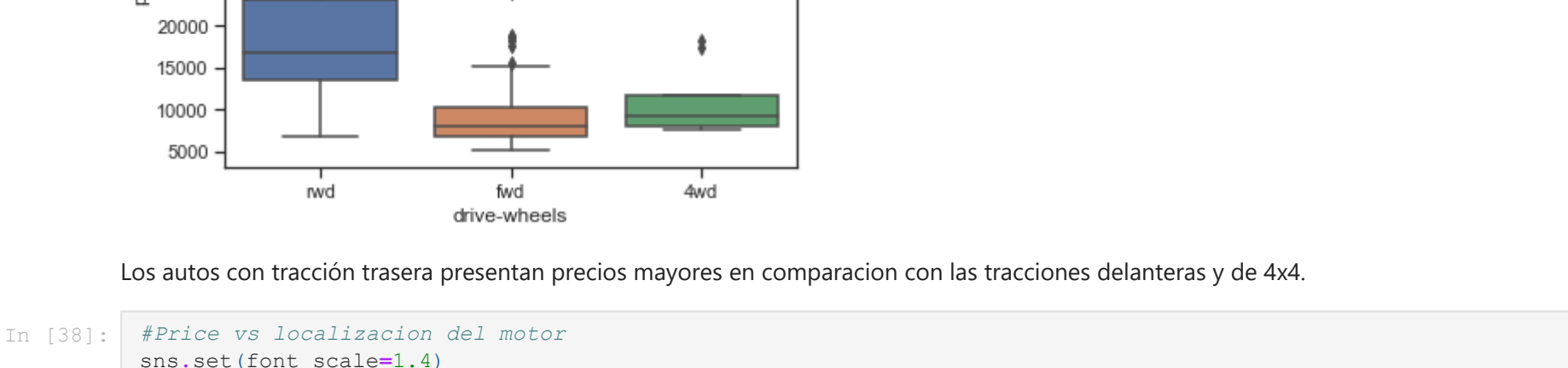
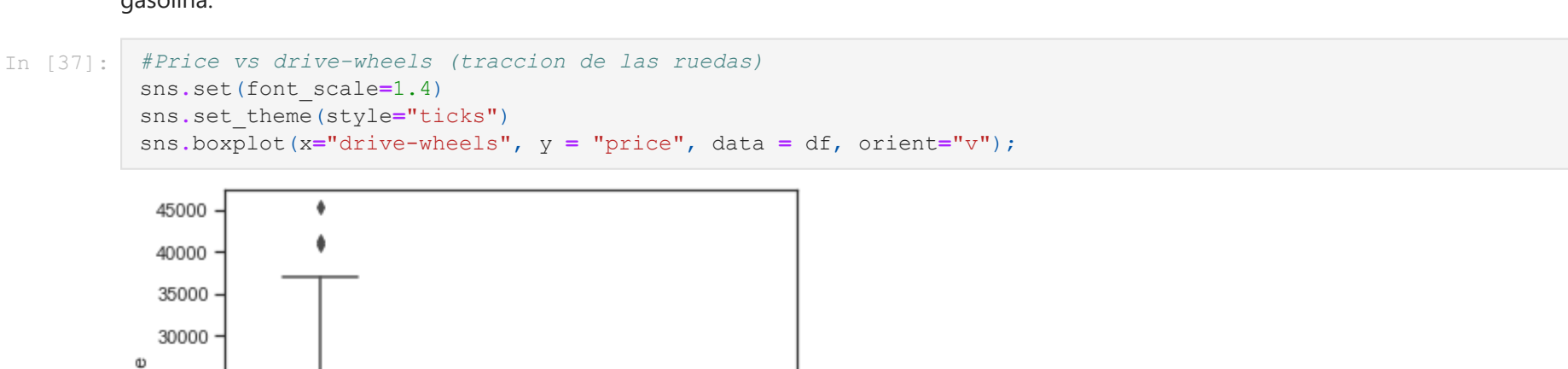
Vemos como se distribuyen los vehículos en cantidad por cada uno de los fabricantes en el dataset



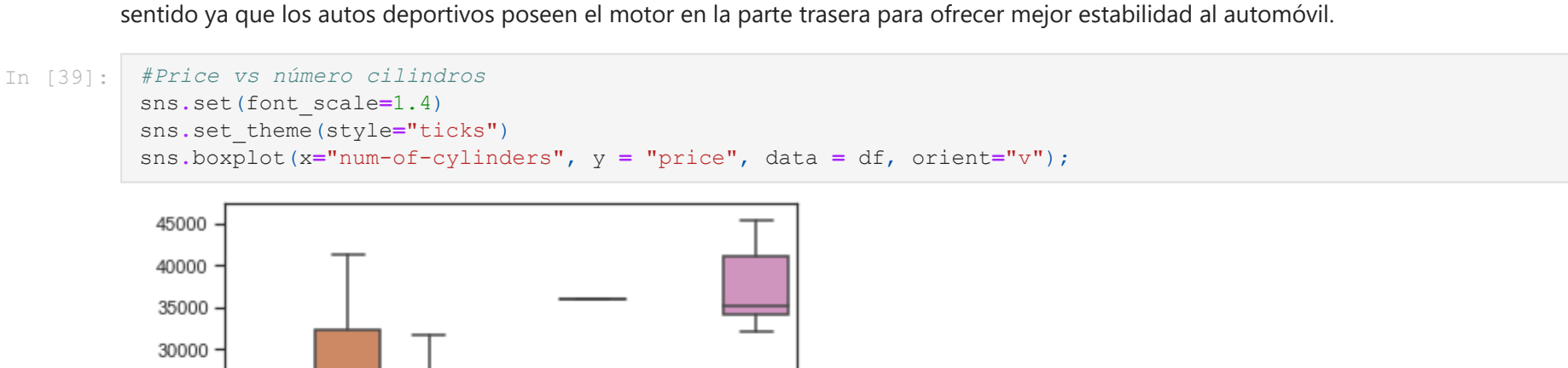
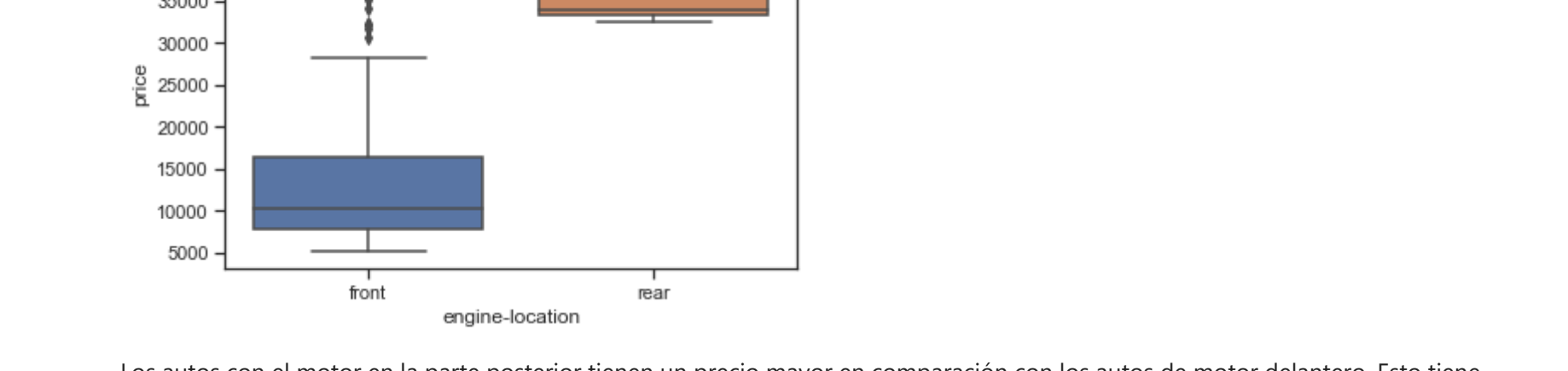
Como se puede observar, la mayor cantidad de registros corresponde a autos de la marca Toyota, Nissan y Mazda, mientras que los de menor cantidad son Mercury, Renault y Alfa-Romeo. Ahora miremos como es el comportamiento de precio frente al estilo del vehículo.



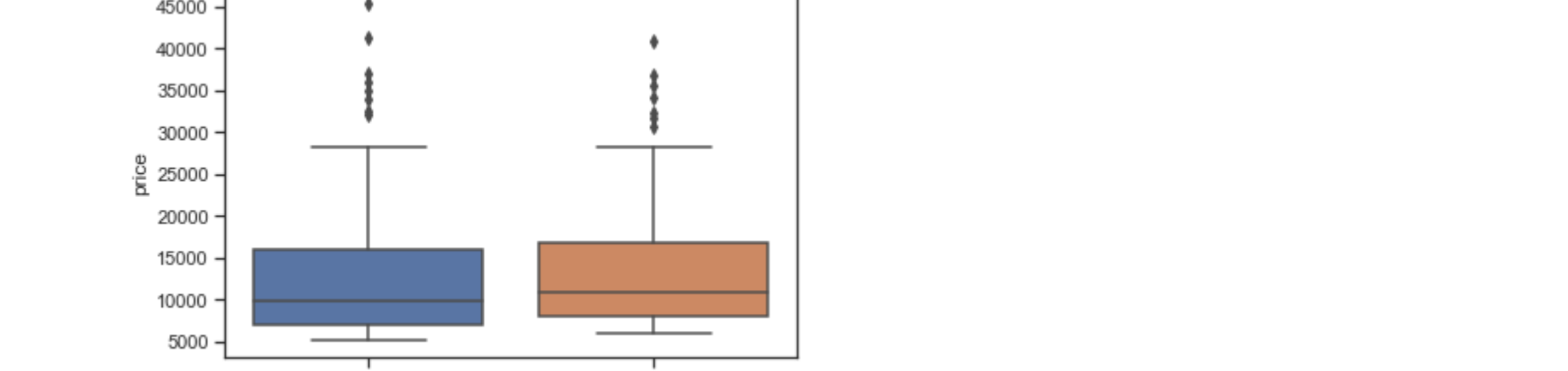
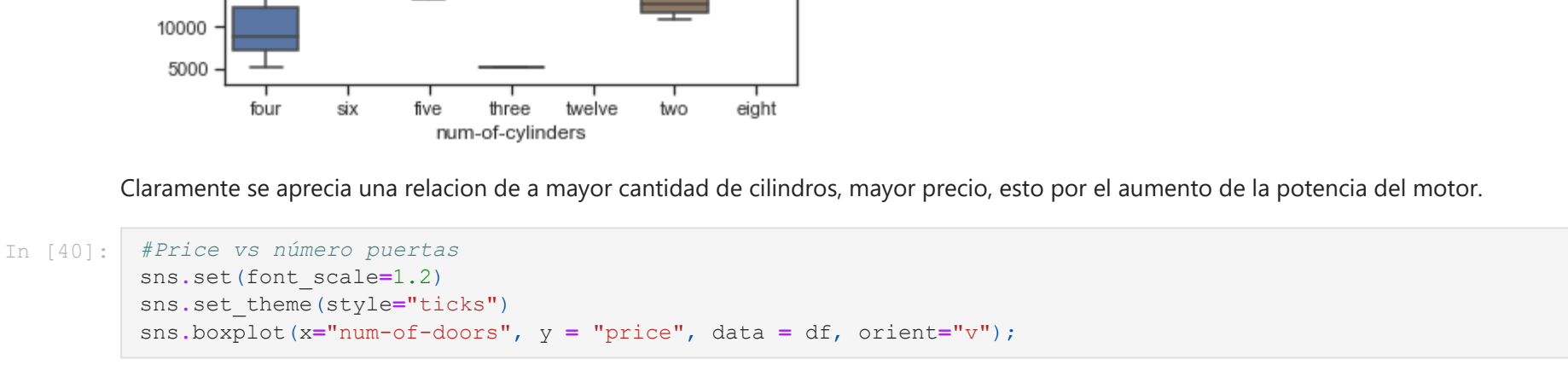
El precio de los autos convertibles se muestra mayor al resto de tipos de vehículos, aunque el tipo Hardtop o techo rígido presenta algunos valores altos de precio.



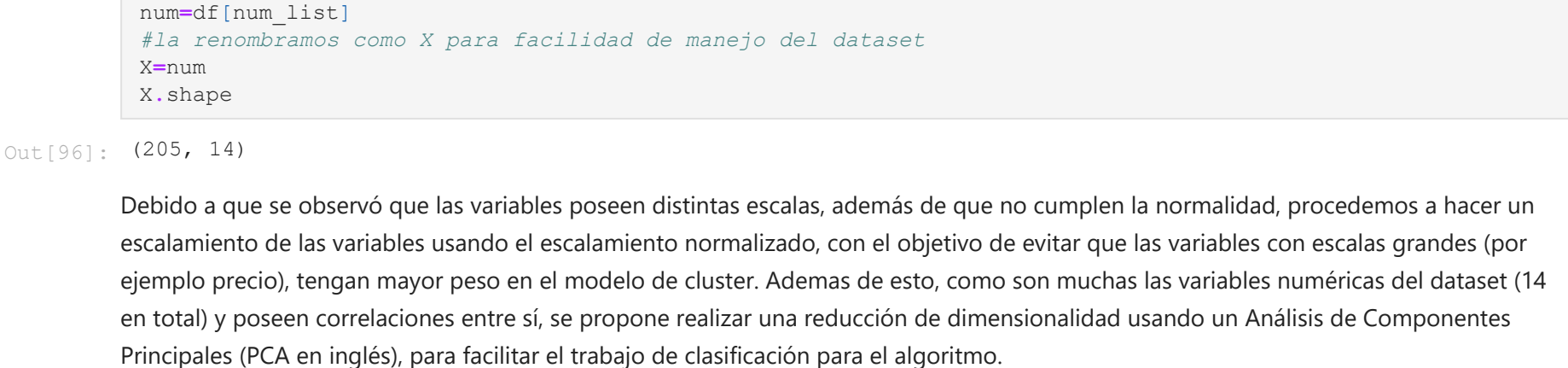
Los motores de tipo ohcv y dthcv (de doble árbol de leva) presentan precios mayores a los demás tipos de motor. Los motores de tipo ohc presentan los menores precios.



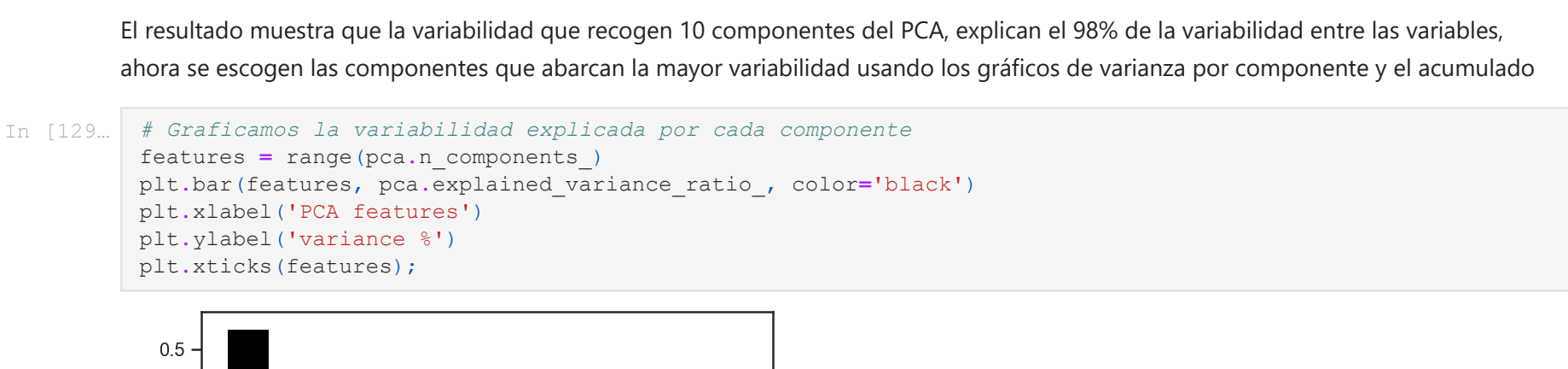
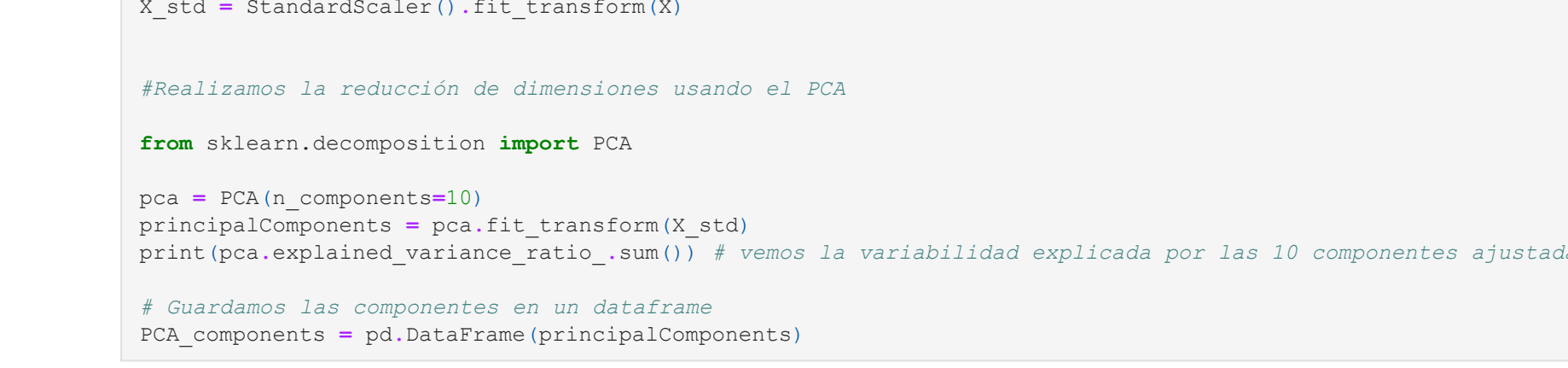
Los modelos turbo aspirados tienen en general mayor precio



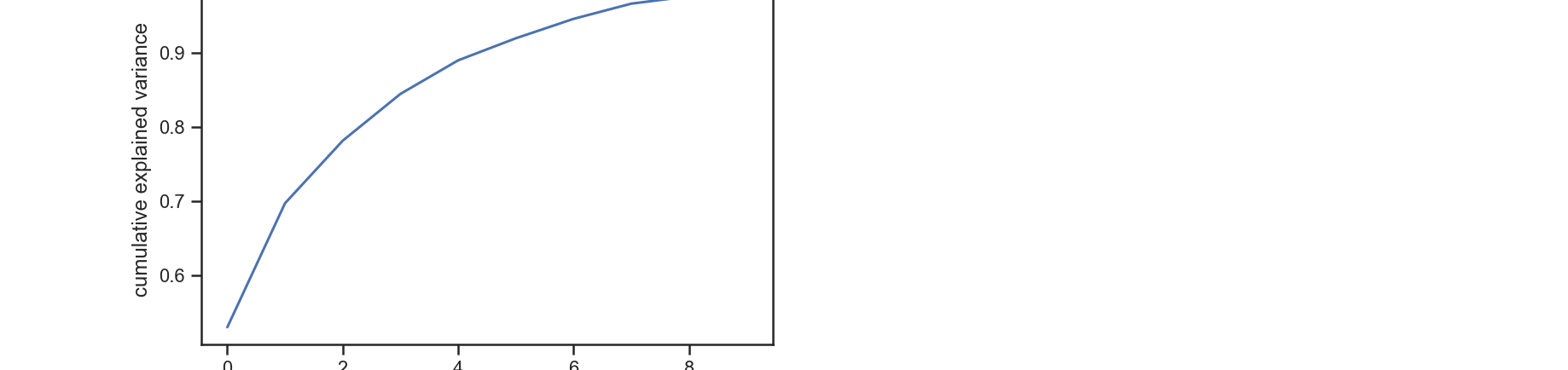
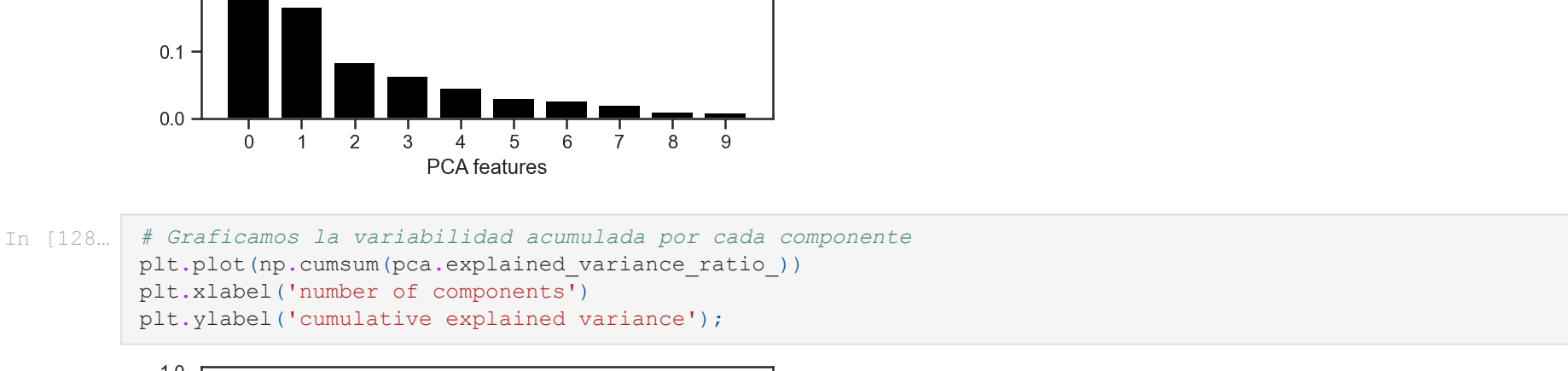
Los motores diesel presentan un valor superior a los motores de gasolina común, aunque hay una cantidad importante de motores atípicos en vehículos con gasolina corriente, correspondientes a autos de fabricantes de alta gama que potencian sus motores con gasolina.



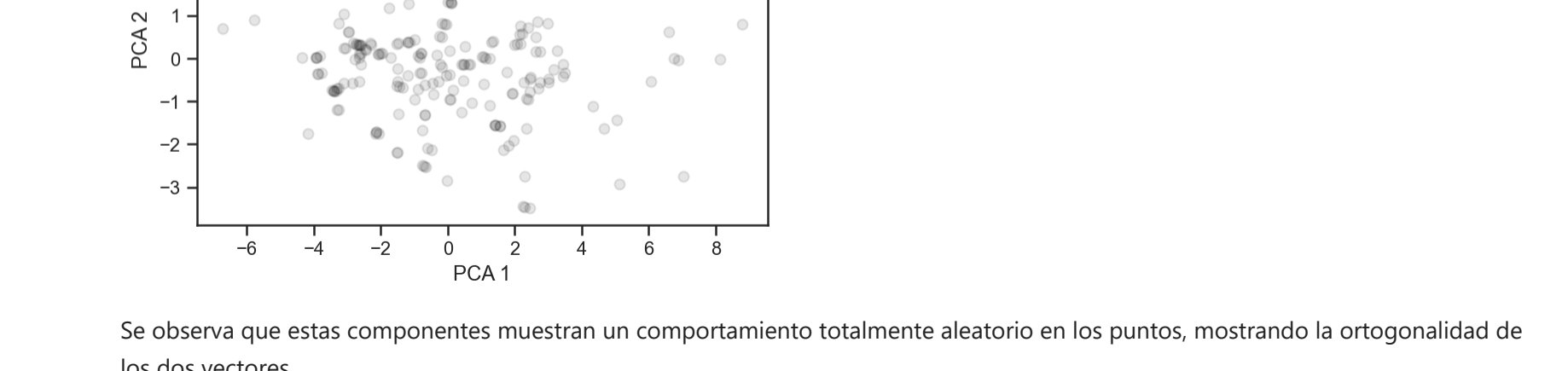
Los autos con tracción trasera presentan precios mayores en comparación con las tracciones delanteras y de 4x4.



Los autos con el motor en la parte posterior tienen un precio mayor en comparación con los autos de motor delantero. Esto tiene sentido ya que los autos deportivos poseen el motor en la parte trasera para ofrecer mejor estabilidad al automóvil.



Claramente se aprecia una relación de a mayor cantidad de cilindros, mayor precio, esto por el aumento de la potencia del motor.



El número de puertas no afecta el valor del vehículo

6. Machine Learning para modelos en categorías los automóviles

En esta sección se procede a construir el modelo de clustering para las variables del dataset, que permita clasificar los vehículos según sus características. Procedemos a preparar los datos para realizar el modelo planteado

```
In [96]: #Preparación de los datos: particionamos los datos en variables numericas en un conjunto y las categoricas en otro
str_list=[] # lista con las variables categoricas
for colname,colvalue in df.items():
    if type(colvalue)!=str:
        str_list.append(colname)
# obtenemos las numericas por defecto
num_list=df.columns.difference(str_list)

#Realizamos las variables numericas del dataset original
num=df[num_list]
#la renombramos como X para facilidad de manejo del dataset
X=num
X.shape
```

```
Out[96]: (205, 14)
```

Debido a que a se observó que las variables poseen distintas escalas, además de que no cumplen la normalidad, procedemos a hacer un escalado de las variables usando el escalamiento normalizado, con el objetivo de evitar que las variables con escalas grandes (por ejemplo precio), tengan mayor peso en el modelo de cluster. Además de esto, como son muchas las variables numéricas del dataset (14 en total) y poseen correlaciones entre si, se propone realizar una reducción de dimensionalidad usando un Análisis de Componentes Principales (PCA en inglés), para facilitar el trabajo de clasificación para el algoritmo.

```
In [121]: #cargamos el metodo de escalamiento, en este caso el escalamiento normalizado
from sklearn.preprocessing import StandardScaler

X_std = StandardScaler().fit_transform(X)

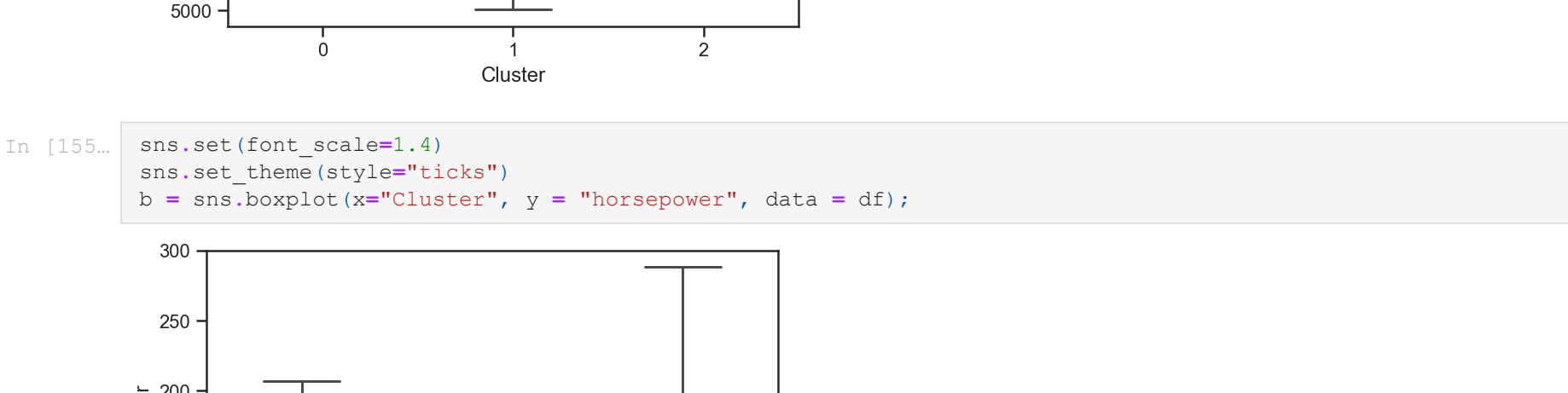
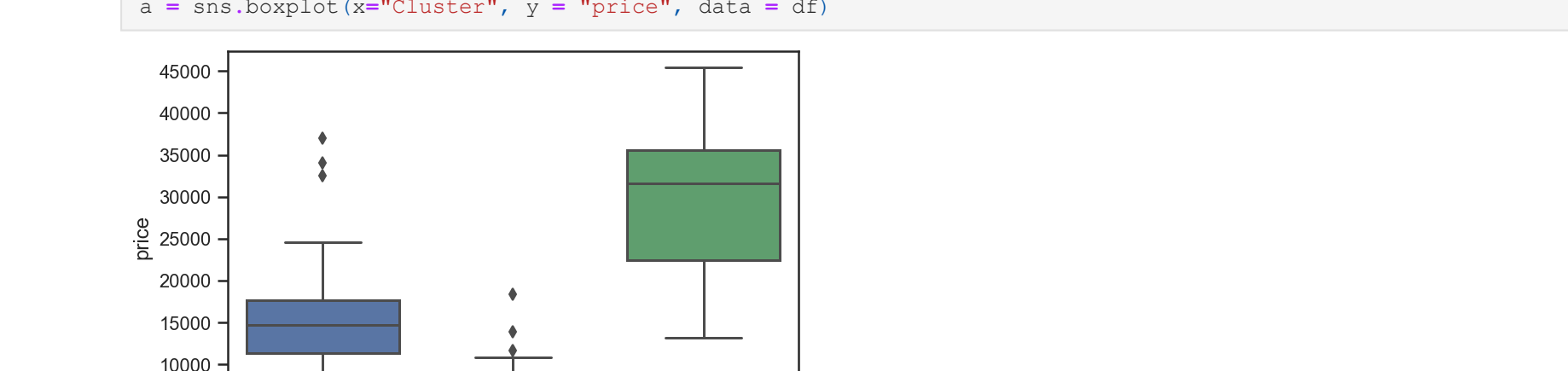
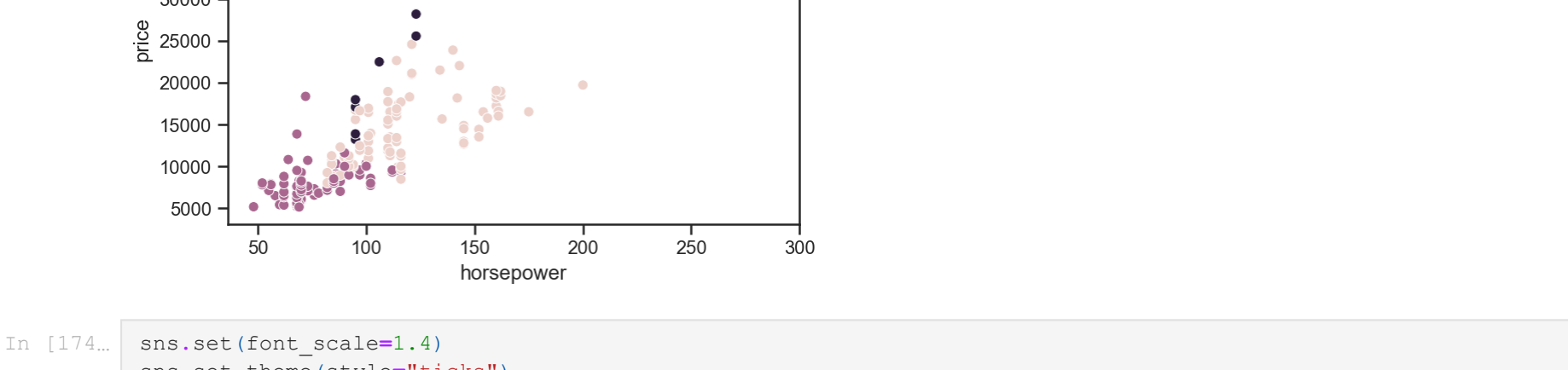
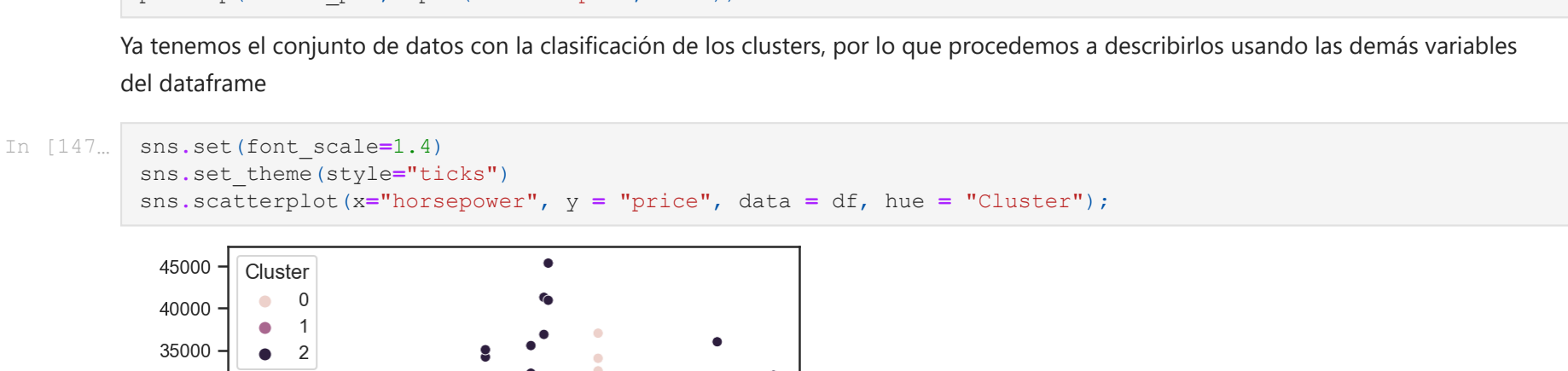
#Realizamos la reducción de dimensiones usando el PCA
from sklearn.decomposition import PCA

pca = PCA(n_components=10)
principalComponents = pca.fit_transform(X_std)
principalComponents.shape # Vemos la variabilidad explicada por las 10 componentes ajustadas
```

```
# Guardamos las componentes en un dataframe
PCA_components = pd.DataFrame(principalComponents)

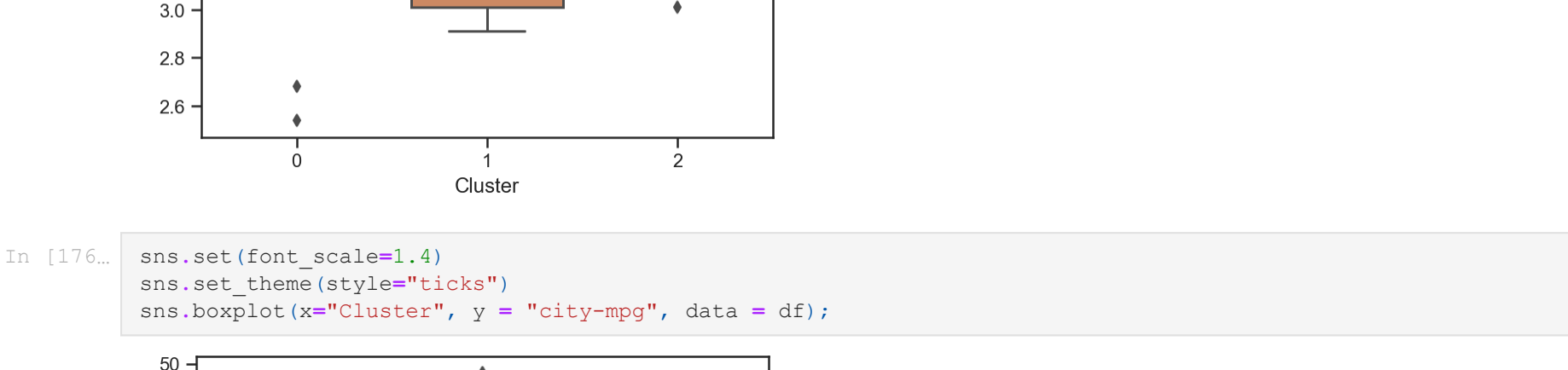
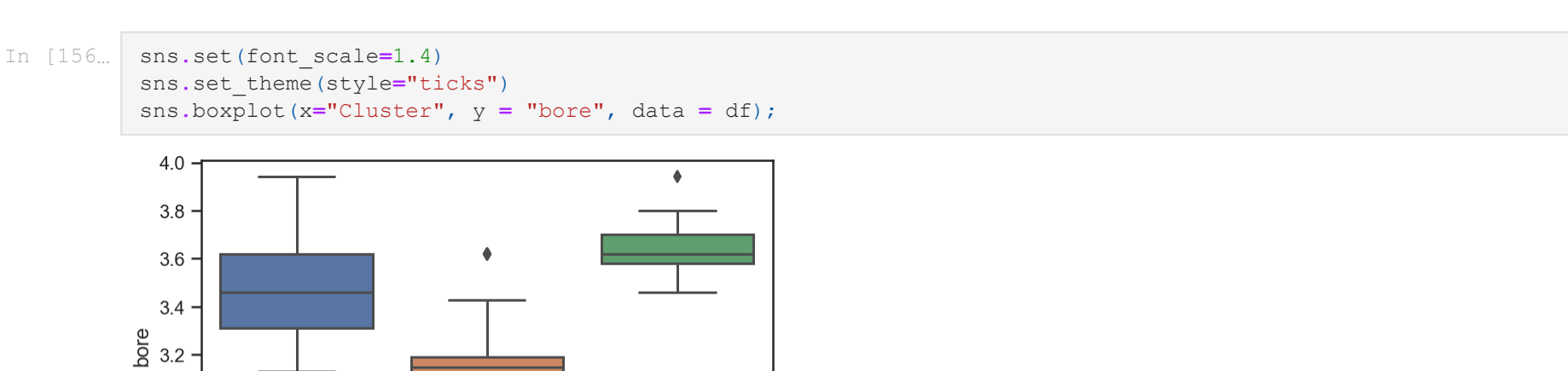
0.983939560762749
```

El resultado muestra que la variabilidad que abogan 10 componentes del PCA, explican el 98% de la variabilidad entre las variables, ahora se escogen las componentes que recorren la mayor variabilidad usando los graficos de varianza por componente y el acumulado



Dado que se sugiere trabajar con las componentes que por lo menos acumulen el 80% de la variabilidad total de los datos, se opta por usar las tres primeras componentes, las cuales acumulan cerca del 85% de la variabilidad.

Para efectos de este análisis, se grafican las dos primeras componentes y apreciar su comportamiento.



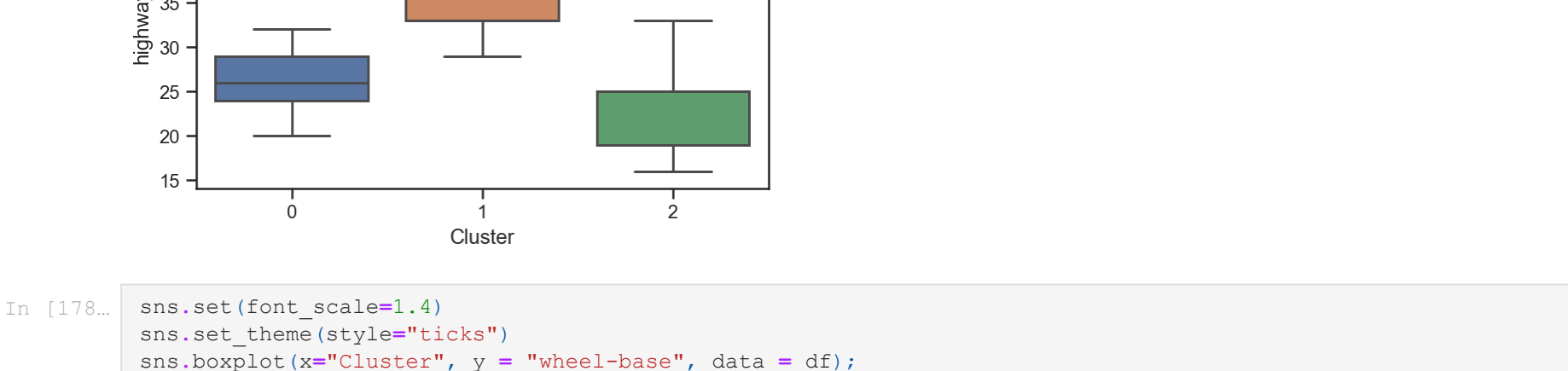
Se observa que estas componentes muestran un comportamiento totalmente aleatorio en los puntos, mostrando la ortogonalidad de los dos vectores.

Ahora vamos a realizar los clusters usando las 3 primeras componentes del PCA. Para encontrar el número óptimo de clusters, usaremos el método del codo, el cual, basado en la variabilidad intraclusters, permite definir la cantidad óptima de conglomerados a trabajar, para eso usaremos la librería k-me de Python.

```
In [139]: from sklearn import KMeansLocator
from sklearn.cluster import KMeans

wcss = []
max_clusters = 10
for i in range(1, max_clusters):
    kmeans_pca = KMeans(n_clusters=i, init='k-means++', random_state=42)
    kmeans_pca.fit(PCA_components.iloc[:,3:])
    wcss.append(kmeans_pca.inertia_)

# Encontramos el codo y numero optimo de clusters
n_clusters = KMeansLocator([for i in range(1, max_clusters)], wcss, curve='convex', direction='decreasing').
print("Optimal number of clusters", n_clusters)
```



Al hacer uso del método del codo, encontramos tres clusters optimos para nuestro dataset, procedemos entonces a construir el modelo

```
In [141]: kmeans_pca = KMeans(n_clusters=n_clusters, init='k-means++', random_state=42)
kmeans_pca.fit(PCA_components.iloc[:,3:])

# analizamos los cluster al dataset
df['Cluster'] = kmeans_pca.labels_
```

```
In [175]: df.head(3)
```

```
Out[175]:
```

	make	fuel-type	aspiration	num-of- doors	body- style	drive- wheels	engine- location	engine- type	num-of- cylinders	fuel- system	engine- size	bore	stroke	compression- ratio	horsepower
0	alfa-romero	gas	std	two	convertible	rwd	front	dthc	four	mpfi	1.300	3.47	2.68		9.0
1	alfa-romero	gas	std	two	convertible	rwd	front	dthc	four	mpfi	1.300	3.47	2.68		9.0
2	alfa-romero	gas	std	two	hatchback	rwd	front	ohcv	six	mpfi	1.520	2.68	3.47		9.0

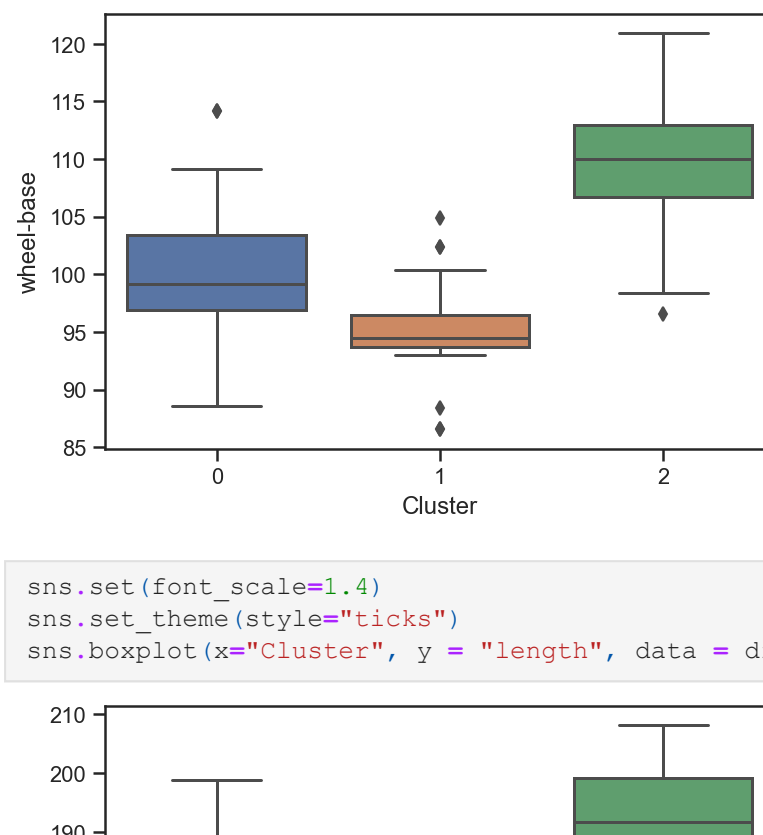
3 rows x 25 columns

```
In [202]: # Guardamos el modelo PCA para el modelo posterior uso
import pickle as pk

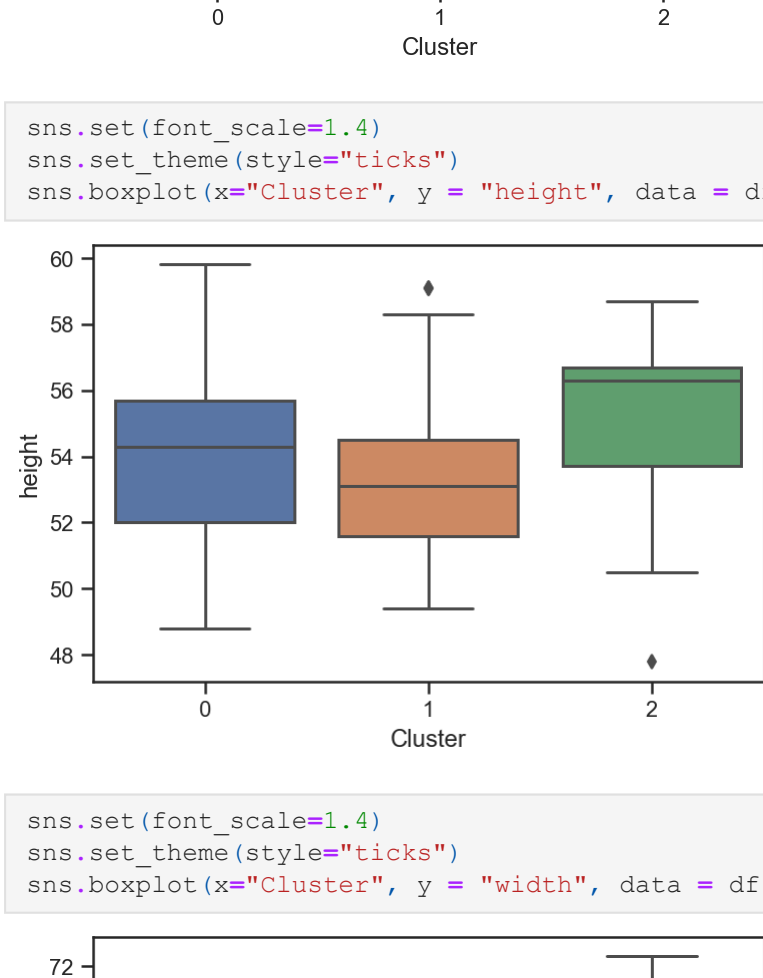
pk.dump(pca, open("pca.pkl", "wb"))

# Guardamos el modelo de kmedas para su posterior uso
pk.dump(kmeans_pca, open("kmeans.pkl", "wb"))
```

Ya tenemos el conjunto de datos con la clasificación de los clusters, por lo que procedemos a describirlos usando las demás variables del dataframe



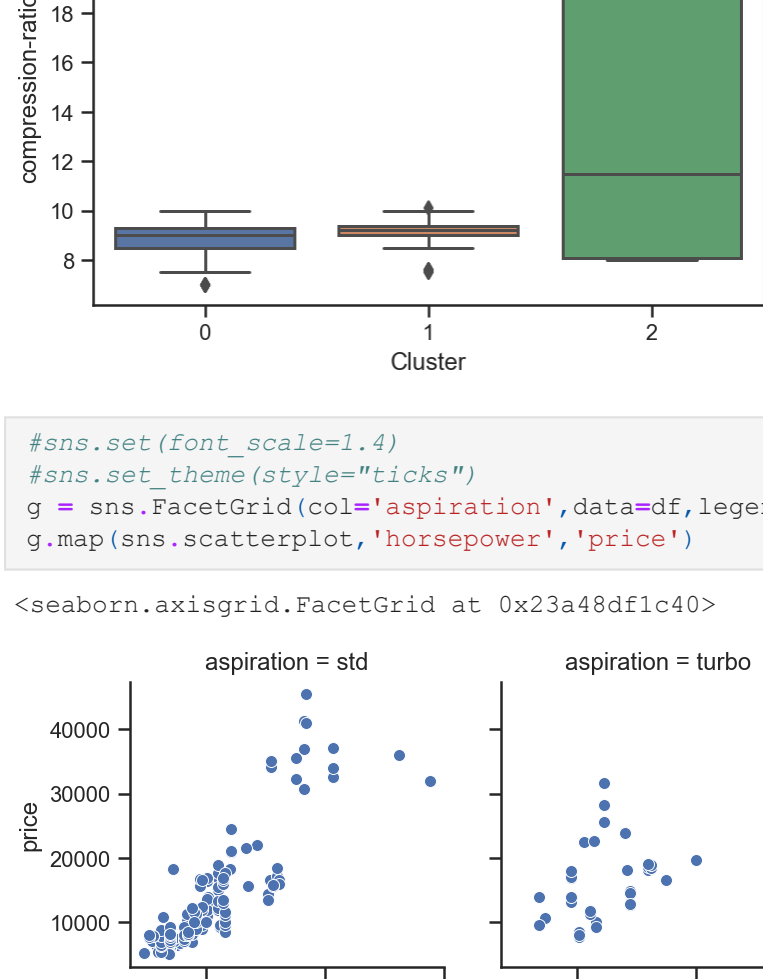
```
In [179]: sns.set(font_scale=1.4)
sns.set_theme(style="ticks")
sns.boxplot(x="Cluster", y = "length", data = df);
```



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

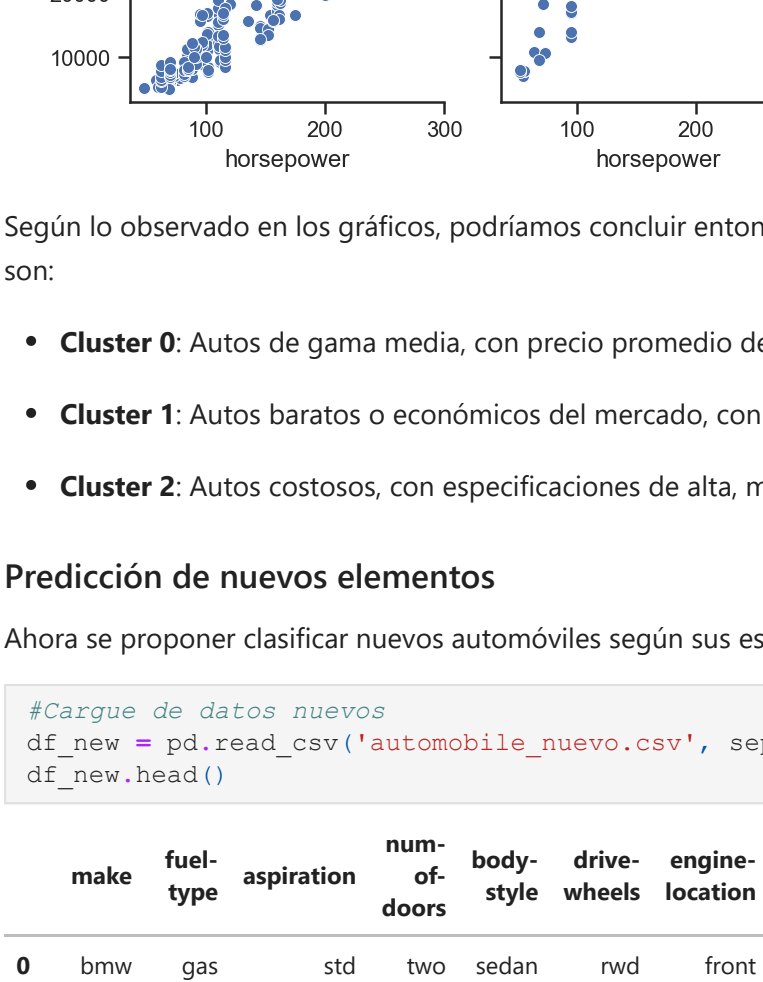


```
In [182]: sns.set(font_scale=1.4)
sns.set_theme(style="ticks")
sns.boxplot(x="Cluster", y = "width", data = df);
```

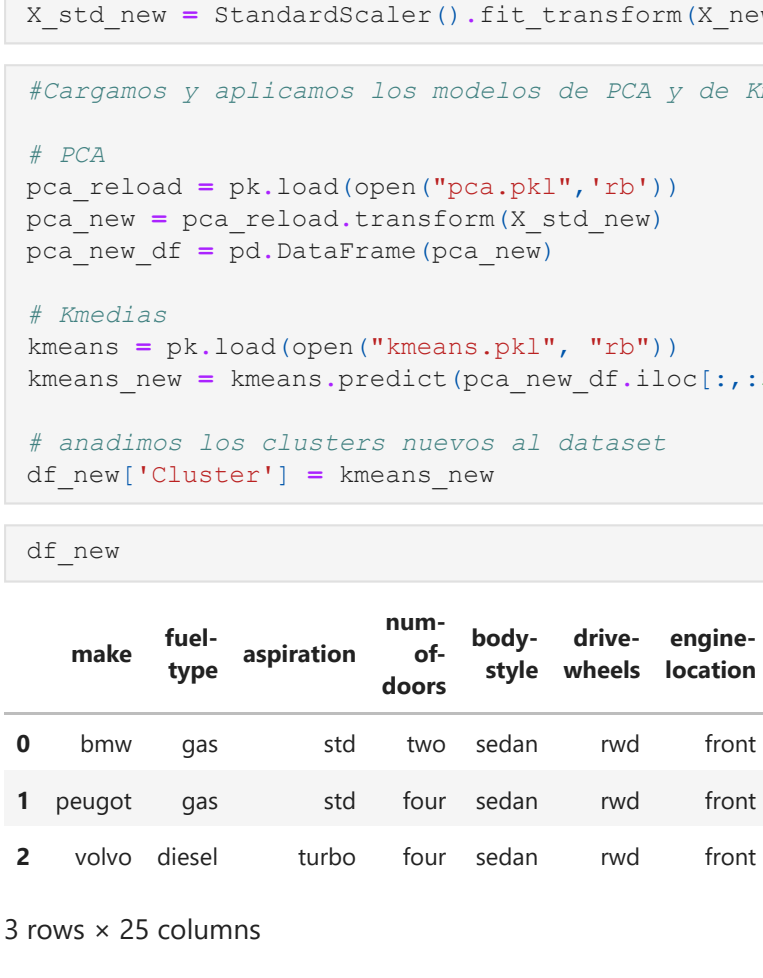


```
In [ ]: sns.set(font_scale=1.4)
sns.set_theme(style="ticks")
sns.boxplot(x="Cluster", y = "width", data = df);
```

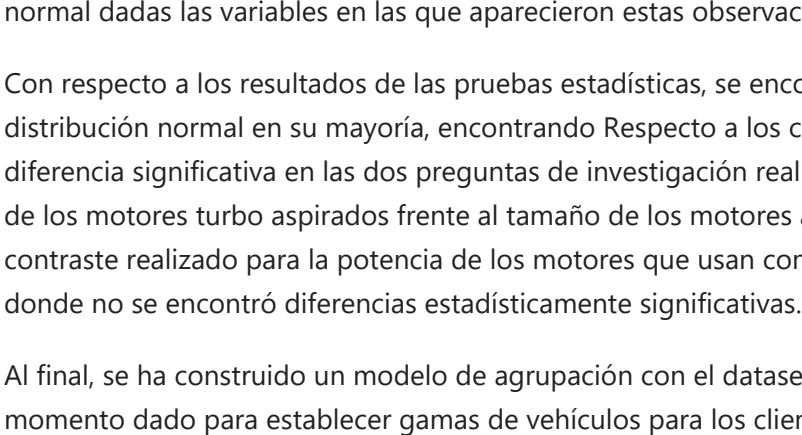
```
In [183]: sns.set(font_scale=1.4)
sns.set_theme(style="ticks")
sns.boxplot(x="Cluster", y = "curb-weight", data = df);
```



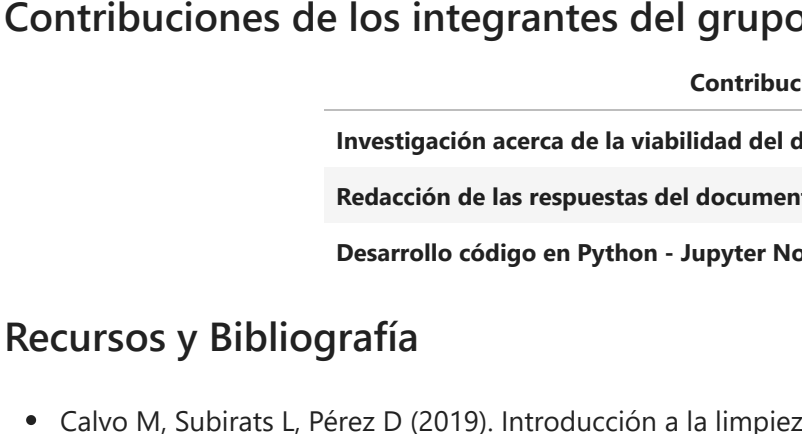
```
In [183]: sns.set(font_scale=1.4)
sns.set_theme(style="ticks")
sns.boxplot(x="Cluster", y = "compression-ratio", data = df);
```



```
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')
```



```
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')
```



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 propone clasificar nuevos automóviles según sus especificaciones en alguno de los clusters, usando el modelo ya construido

```
In [196]: #Carga de datos nuevos
df_new = pd.read_csv('automobile_nuevo.csv', sep = ';')
```

	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	...	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 × 16 columns

```
In [215]: # extruemos solo las variables numericas del nuevo grupo de datos
X_std_new = StandardScaler().fit_transform(X_std_new)
```

```
In [214]: #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 [215]: df_new
```

	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	...	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 × 16 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 atípicos 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 ganancias 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 evidenció 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.

Contribuciones de los integrantes del grupo en la práctica

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

Recursos y Bibliografía

- Calvo M, Subirats L, Pérez D (2019). Introducción a la limpieza y análisis de los datos. Editorial UOC.
- Wes McKinney (2012). Python for Data Analysis. O'Reilly Media, Inc.
- Gibergans J (2018). Contraste de varianzas. Editorial UOC.
- Rovira C (2018). Contraste de hipótesis. Editorial UOC.