
Explorando clusters de autos

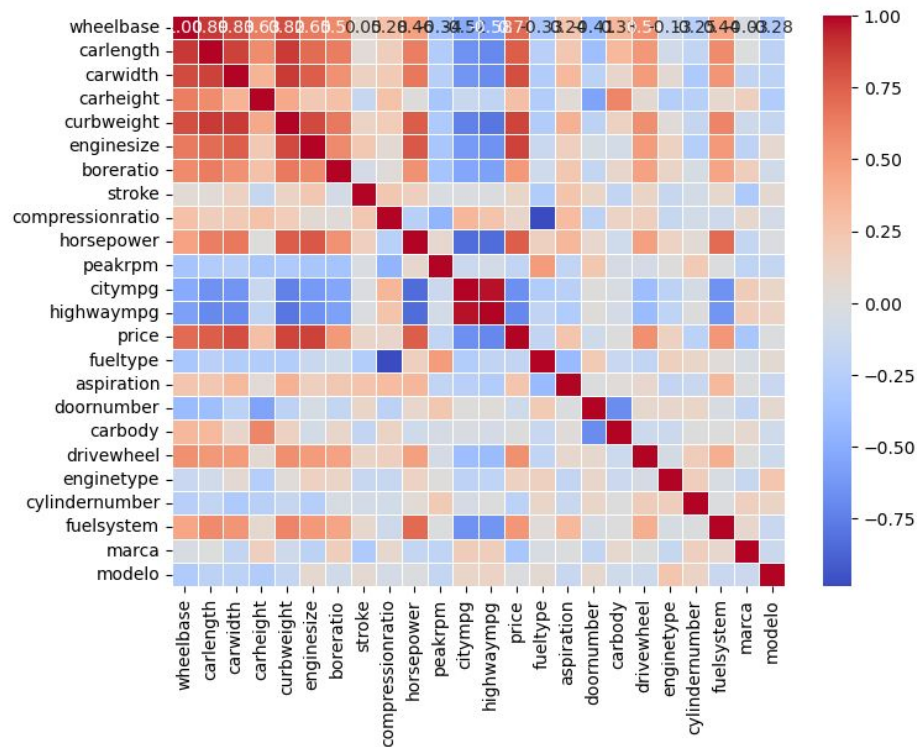
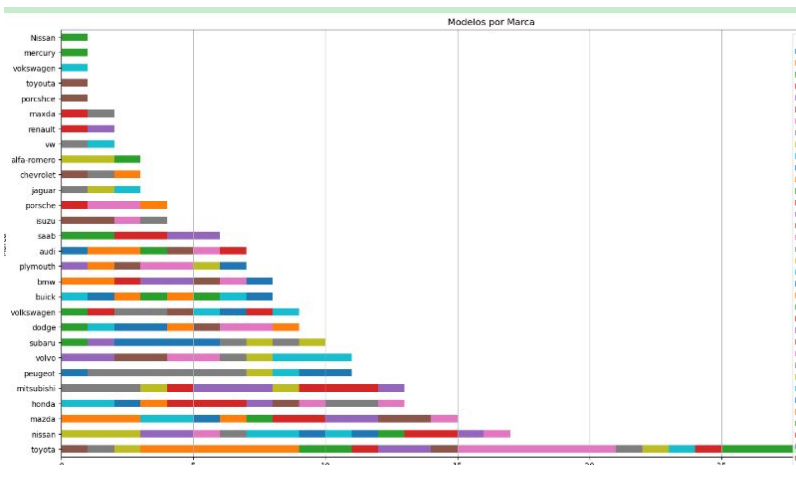
Estadística Multivariada

Miguel Zamorano

Índice

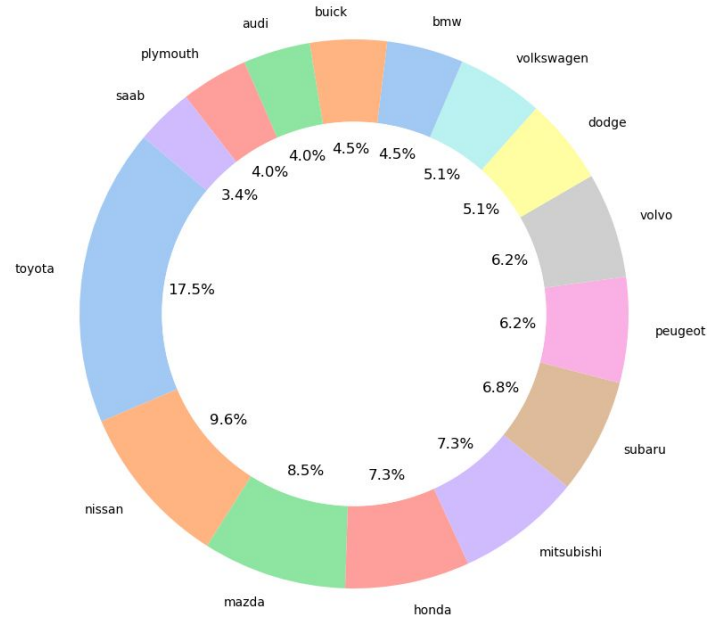
1. EDA
 2. Preprocesamiento
 3. Limpieza
 4. Prueba de Kolmogorov
 5. Clasificación
 - a. SVM
 - b. RandomForest
 6. Clustering
 - a. kmeans
 - b. Aglo.Jerárquico simple
 - c. Aglo.Jerárquico completo
 - d. Aglo.Jerárquico centroide
 7. Clustering con Reducción de dimensionalidad
 - a. PCA y kmeans
 8. Análisis de Calidad
 - a. Tabla Siluetas
 9. Interpretación de Resultados
 - a. **Características correspondientes a cada cluster**
 10. Conclusiones
-

post limpieza: 25 variables numéricas y 177 registros



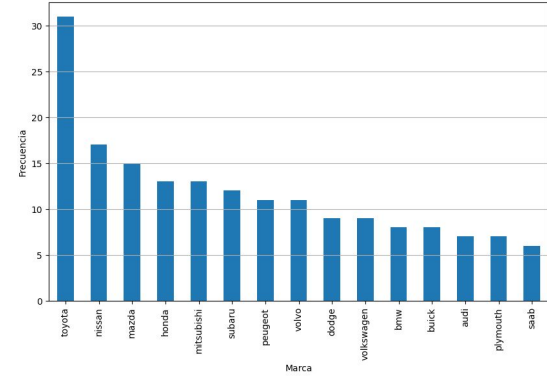
Post Limpieza

Distribución de Marcas



toyota	31
nissan	17
mazda	15
honda	13
mitsubishi	13
subaru	12
peugeot	11
volvo	11
dodge	9
volkswagen	9
bmw	8
buick	8
audi	7
plymouth	7
saab	6

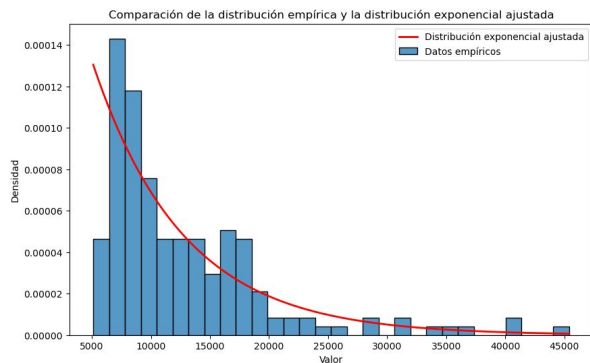
Frecuencia de Marcas



15 marcas de 28

177 registros de 205

Prueba de Bondad:Kolmogorov-Smirnov



```
# Realizamos prueba de Kolmogorov-Smirnov
```

```
ks_statistic, p_value = kstest(p_marca_mayor, 'powerlaw', args=params)
```

```
print(f'Estadístico KS: {ks_statistic}')
```

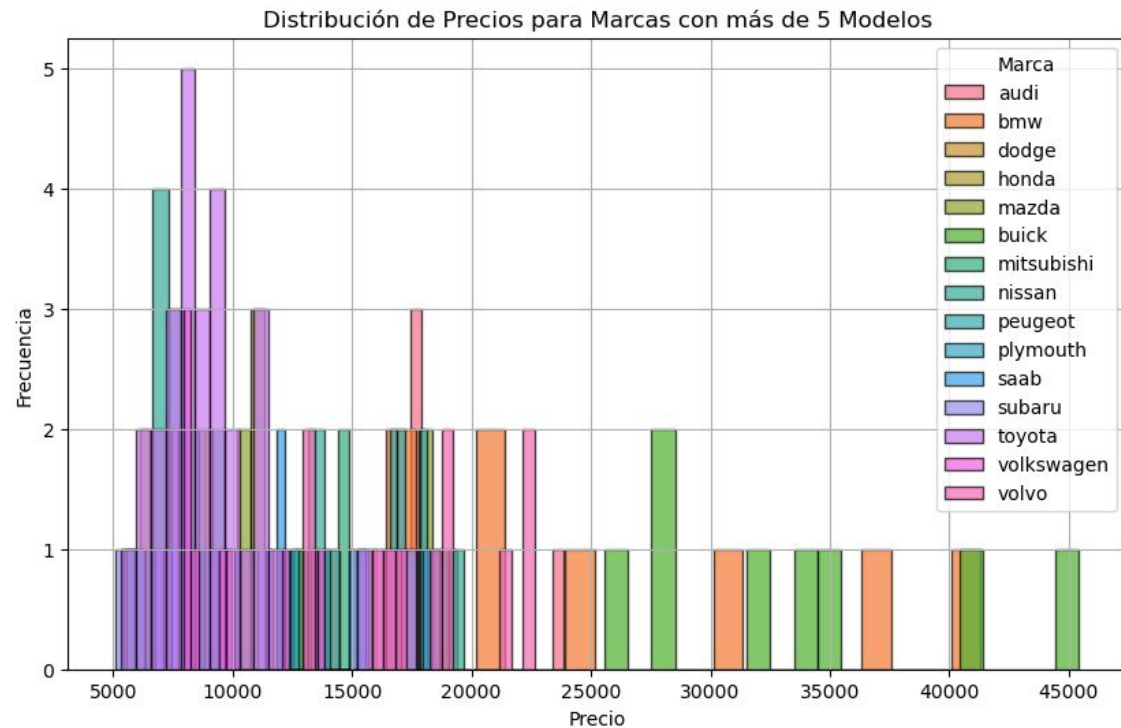
```
print(f'Valor p: {p_value}')
```

```
Estadístico KS: 0.2669693040214083
```

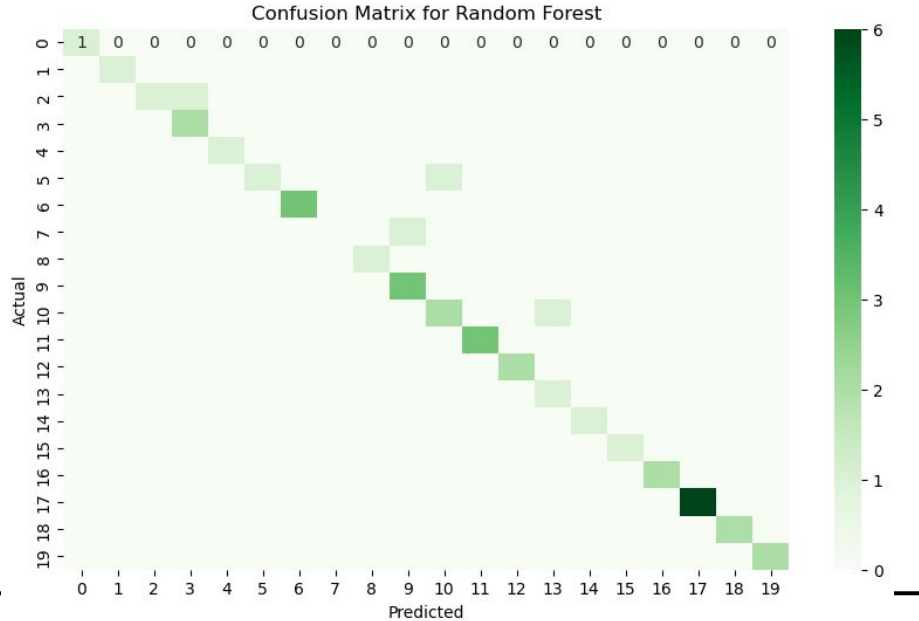
```
Valor p: 1.2543809951898195e-11
```

Rechazamos la hipótesis nula: los datos de la variable precio no siguen una distribución exponencial porque valor P es menor que 0.05

Distribución Precio con Marcas



	Model	Accuracy	Precision	Recall	F1-Score
0	SVM	0.8	0.770833	0.8	0.774405
1	Random Forest	0.9	0.902083	0.9	0.887619



Clustering

1. Clustering

- a. Aglo.Jerárquico simple
- b. Aglo.Jerárquico completo
- c. Aglo.Jerárquico centroide
- d. kmeans

2. Clustering con Reducción de dimensionalidad

- a. PCA y kmeans

3. Análisis de Calidad

- a. Tabla Siluetas

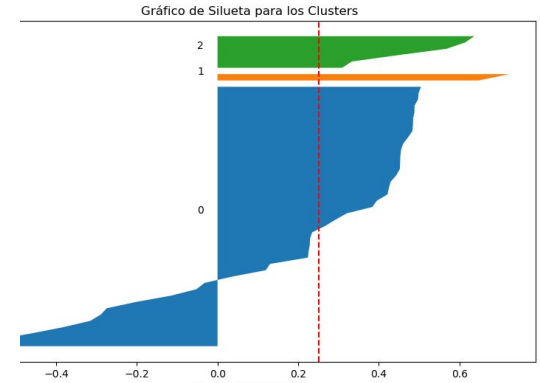
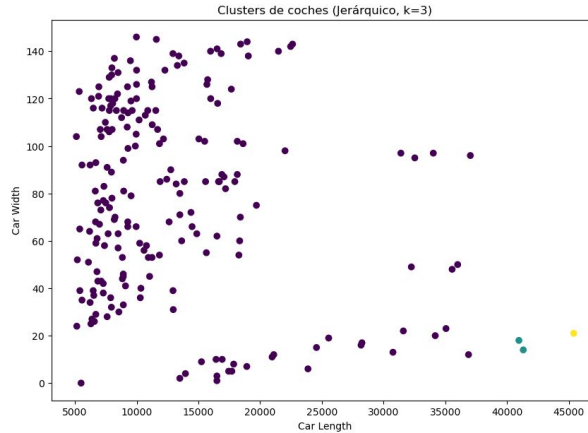
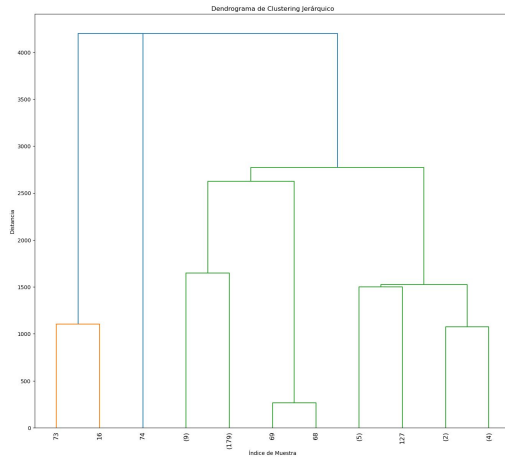
4. Interpretación de Resultados

- a. **Características correspondientes a cada cluster**
-

Resultados

Método	Silhouette Coefficient
kmeans	0.643946205
Aglo.Jerárquico simple	0.25
Aglo.Jerárquico completo	0.42
Aglo Jerárquico centroide	0.42
PCA y kmeans	0.7243985953663576

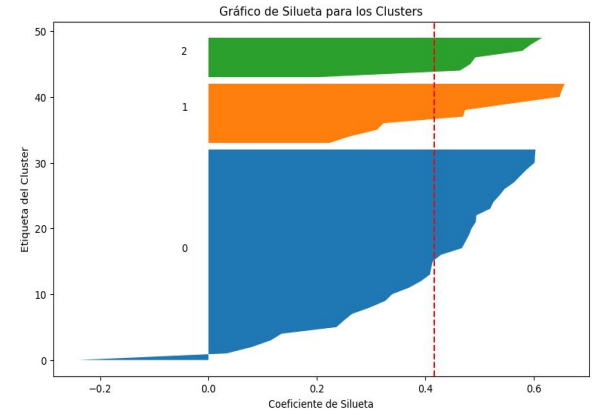
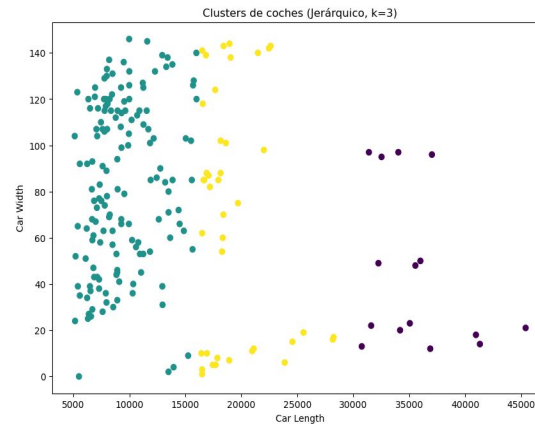
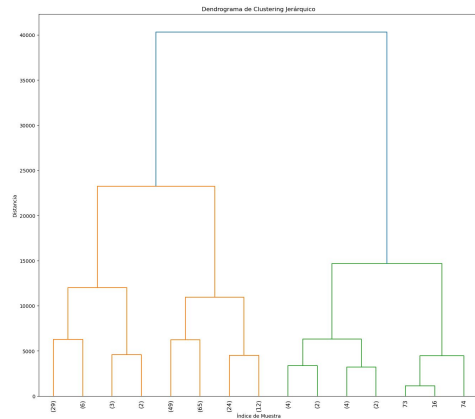
Aglomerativo Jerárquicos simple



la puntuación de silueta promedio es: 0.25

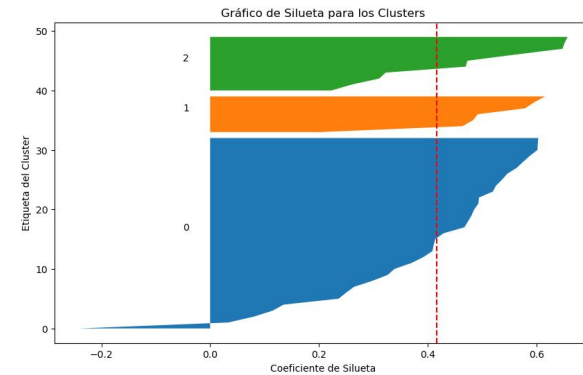
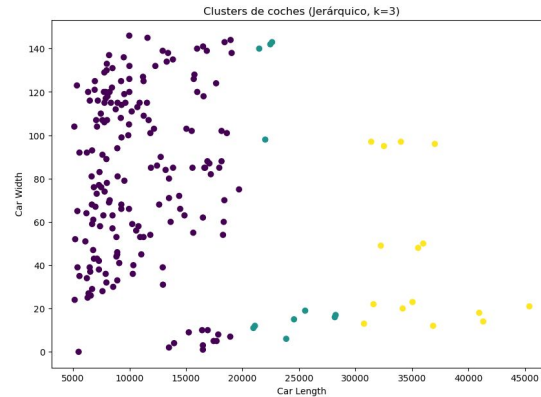
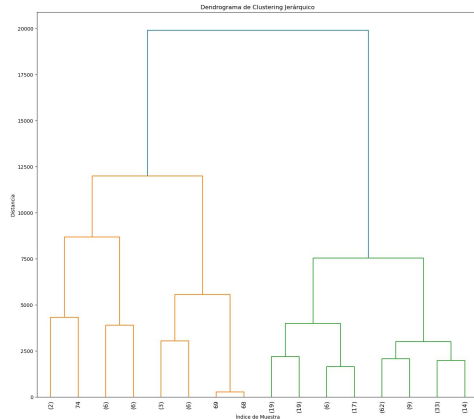
Aglomerativo Jerárquicos completo

Para `n_clusters = 3`, la puntuación de silueta promedio es: 0.42



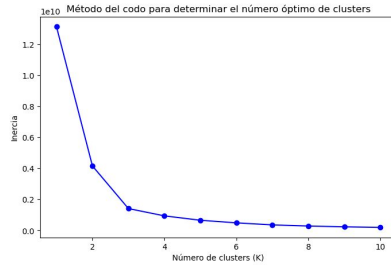
Aglomerativo Jerárquicos centroide

Para `n_clusters = 3`, la puntuación de silueta promedio es: 0.42

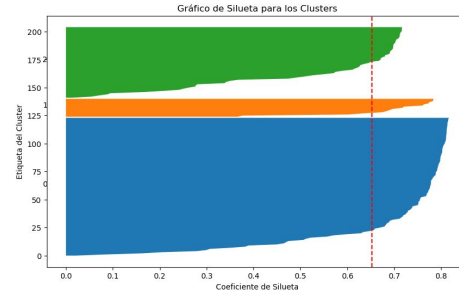


kmeans k=3

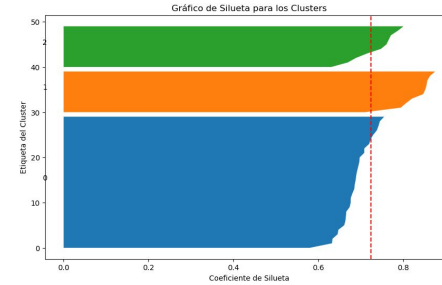
Resultados M.del Codo



Sin PCA 0.643946205

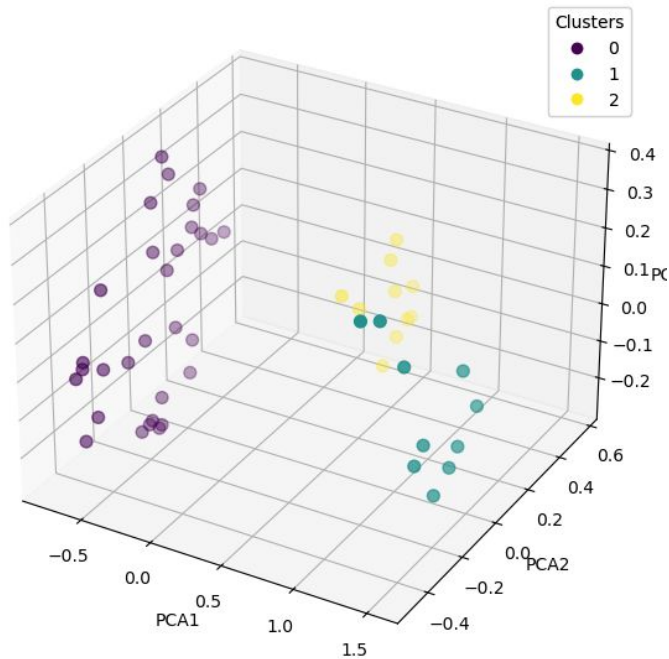


Con PCA.7243985953663576



kmeans con PCA

Visualización de Clustering en 3D



PCA nos encontró 2 componentes Principales

Determinación del Número Óptimo de Componentes Principales

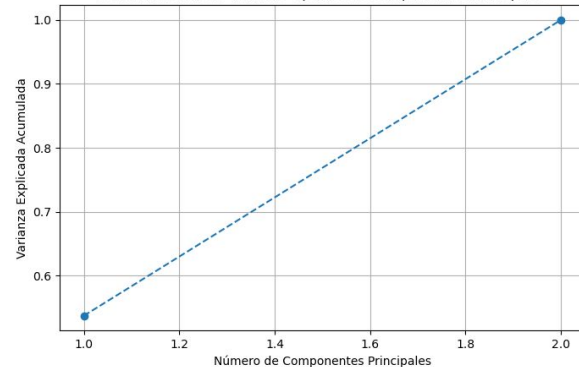
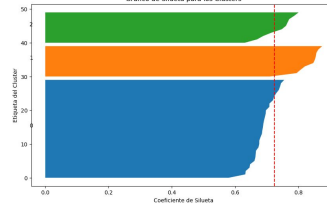
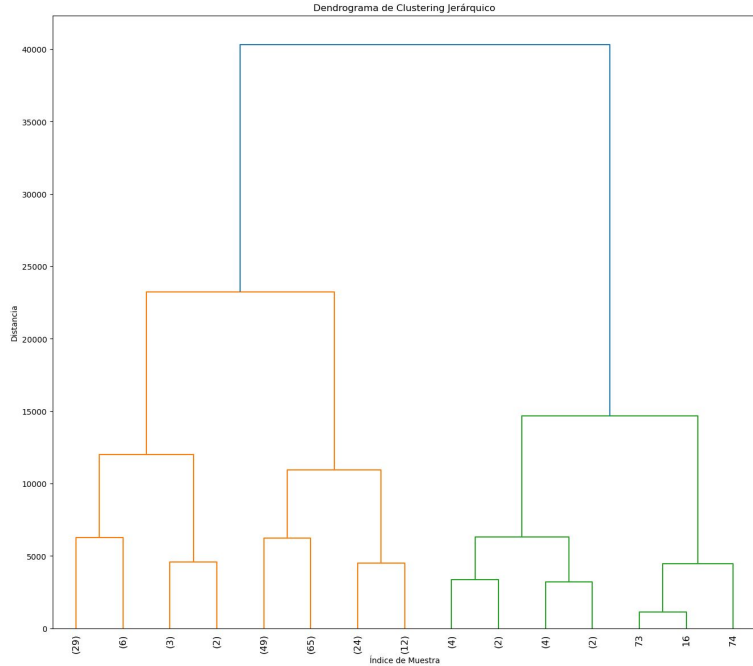


Gráfico de Silueta para los Clusters

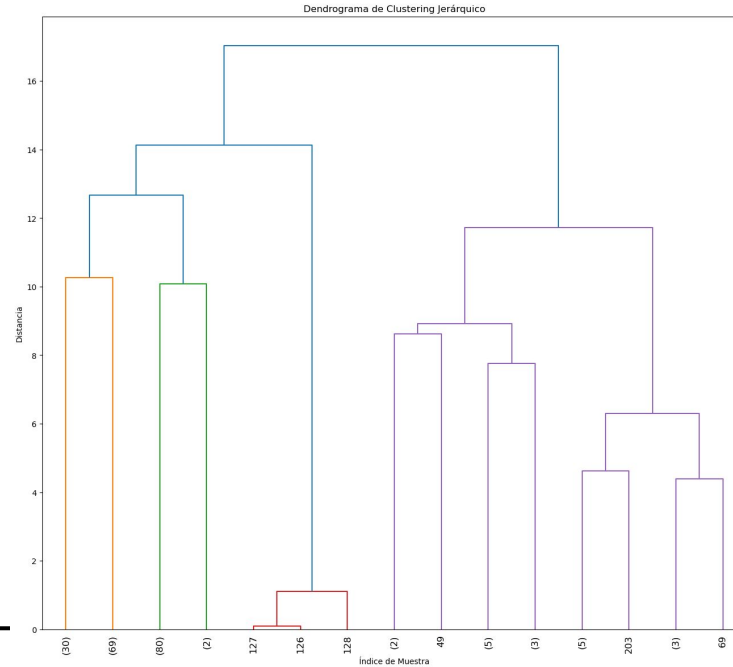


comparativo Aglo completo

Sin PCA



Con PCA

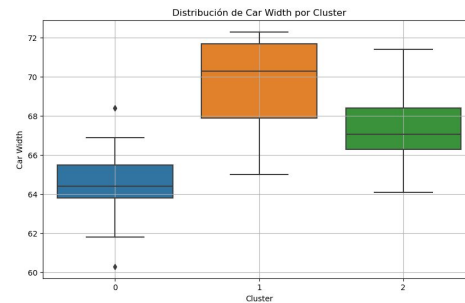
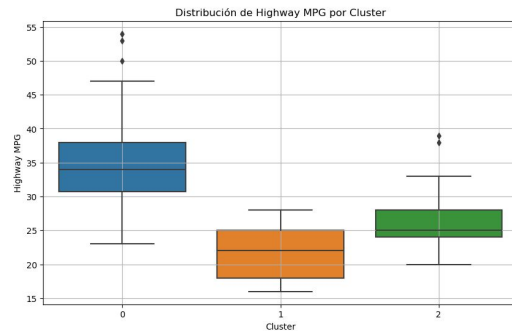
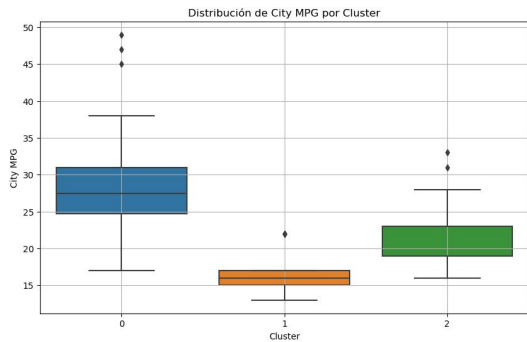
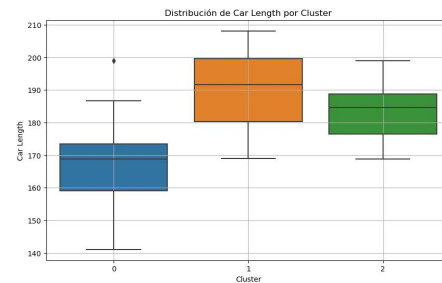
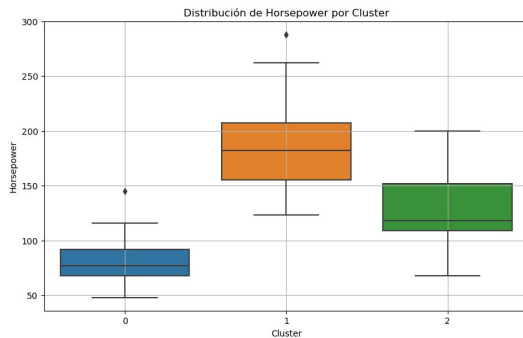
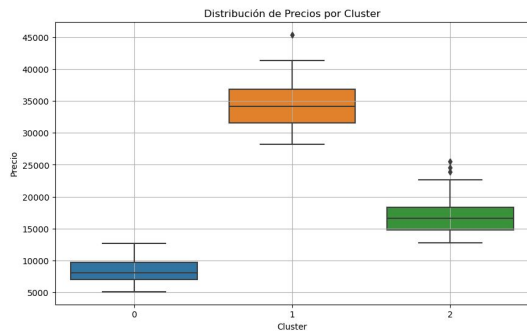


Resultados

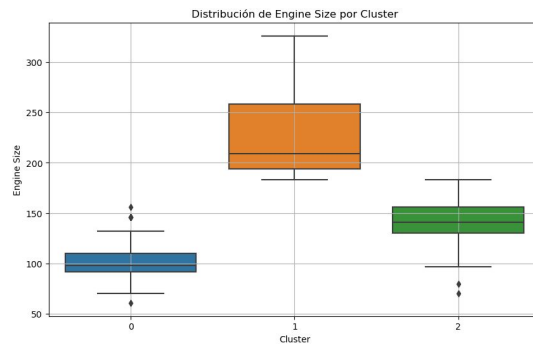
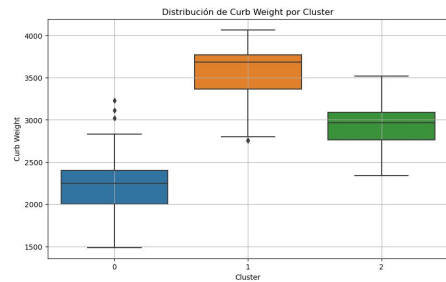
Método	Silhouette Coefficient
kmeans	0.643946205
Aglo.Jerárquico simple	0.25
Aglo.Jerárquico completo	0.42
Aglo Jerárquico centroide	0.42
PCA y kmeans	0.7243985953663576

**Características
representativas
correspondientes a cada
cluster**

cluster vs columna n

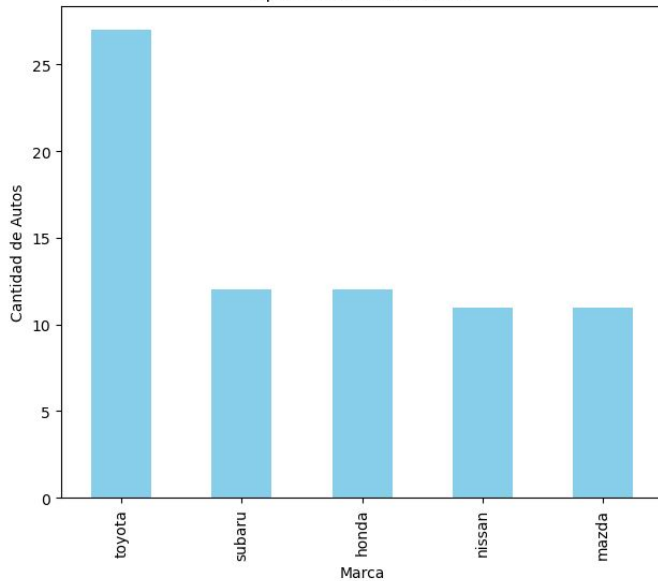


cluster vs columna n

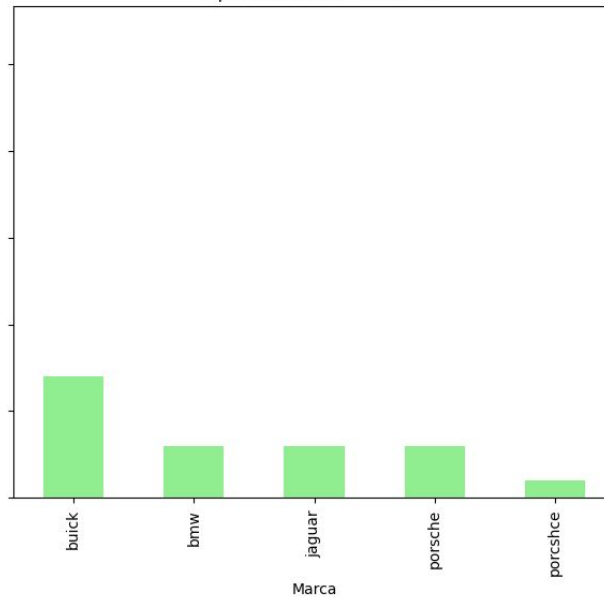


Top 5 Marcas por Cluster

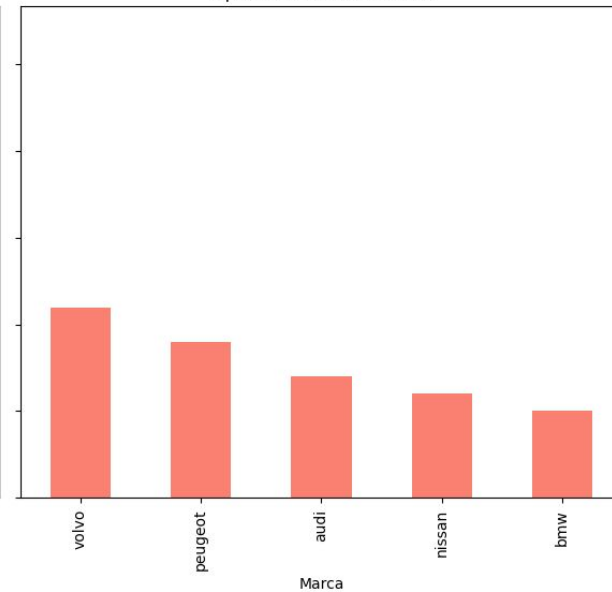
Top 5 Marcas en Cluster 0



Top 5 Marcas en Cluster 1



Top 5 Marcas en Cluster 2



Referencias

El dataset elegido fue obtenido de kaggle

<https://www.kaggle.com/code/davidcanorosillo/car-clustering>

titulado : car datasets
