



# Facultad de Cs. de la Computación



## Introducción a la Ciencia de Datos.

### Proyecto Final: Modelado Predictivo de Alta Precisión para el Sector Automotriz.

Nombre del Alumno: Daniel Carvajal Garcia.

Matrícula: 202502492

Nombre del proyecto: Modelado Predictivo de  
Alta Precisión para el Sector Automotriz.

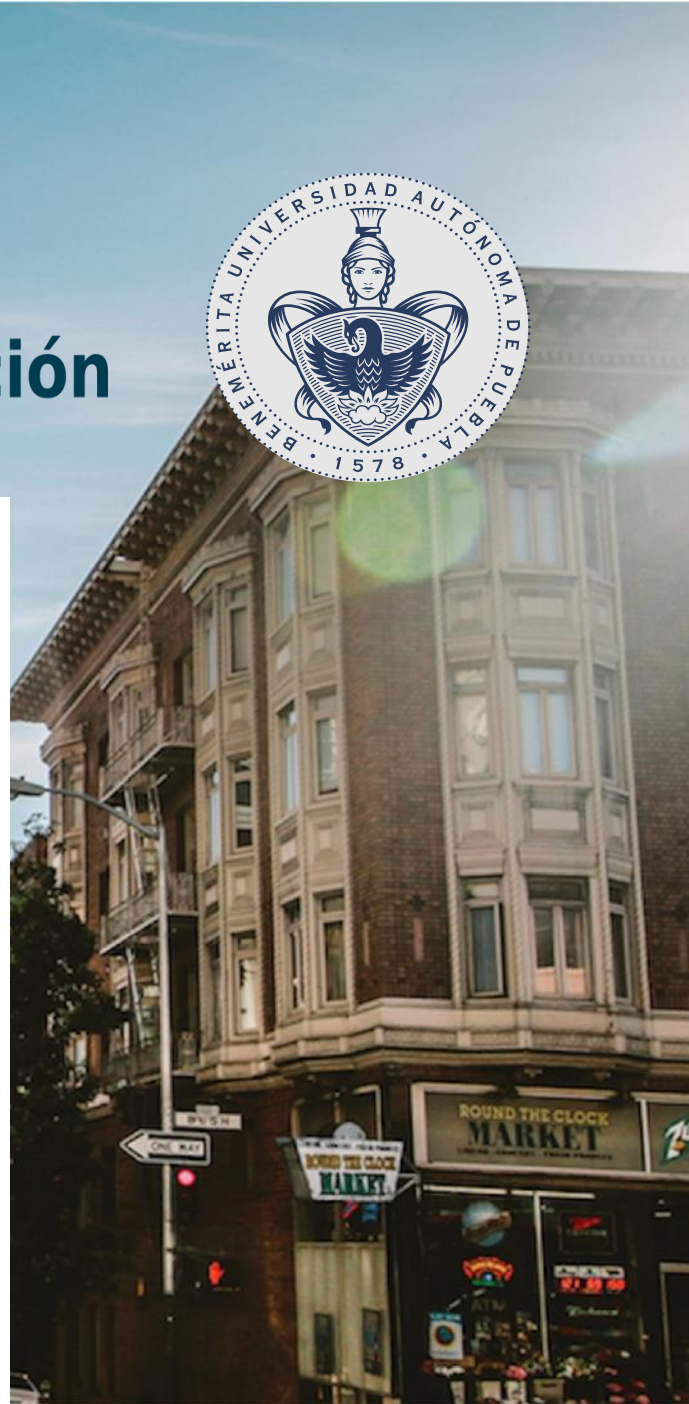
Fecha de entrega: 29 de Noviembre de 2025

Nombre del Profesor: Jaime Alejandro Romero  
Sierra.

Materia: Introducción a la Ciencia de Datos.

Licenciatura: Ingeniería en Ciencia de Datos.

Benemérita Universidad Autónoma de Puebla.



---

# Introducción.

- **Descripción breve del objetivo del proyecto.**

El objetivo principal de este análisis de datos sobre venta de vehículos es aprovechar la información disponible para mejorar la estrategia comercial de cada fabricante de automóviles, con el propósito de tener una planeación en la toma de decisiones estratégicas dentro del mercado automotriz, a fin de incrementar ventas, mejorar la competitividad y fortalecer el posicionamiento de cada fabricante en un lapso de 12 meses.

- **Justificación y contexto: ¿por qué es importante resolver o estudiar esta problemática?** Con este proyecto se busca resolver las siguientes problemáticas:

1. Definir estrategias que permita el posicionamiento de cada fabricante.  
Analizar diferentes factores de los cuales incluye: marca, modelo, año, tipo de combustible o si es híbrido, tamaño de motor, de los cuales influye un precio de venta para el consumidor. Tener claro una estrategia de precios que permita ser competitivo, maximice sus ventas y maximice un margen de ganancias para cada segmento.
2. Detectar tendencias de compra y preferencias de los consumidores.  
Se busca evaluar que fabricantes y modelos de automóviles tienen mayor rentabilidad dentro del mercado, identificar las preferencias de los consumidores en base al tipo de combustible (gasolina o diésel) o si prefiere los autos híbridos y anticipar la evolución de la demanda de cada uno. Reconocer patrones de compra, por ejemplo, cuáles son los meses del año en donde el consumidor llega a comprar un automóvil, con el objetivo de planificar lanzamiento y campañas de marketing.
3. Impulsar decisiones estratégicas a través de datos.  
Comparar el desempeño de ventas entre los distintos fabricantes y sus diferentes modelos, del mismo modo, detectar cuales las fortalezas y debilidades de cada fabricante, como puede ser el precio, motor, tipo de combustible, apoyar en la definición de estrategias de diferenciación para cada fabricante.

- **Fuentes de datos: descripción de las bases de datos empleadas (origen, cantidad de datos, principales características)**

Fuente de la base de datos original:

<https://www.kaggle.com/datasets/msnbehdani/mock-dataset-of-second-hand-car-sales>

Este conjunto de datos contiene información detallada sobre la venta de automóviles, abarcando múltiples fabricantes y modelos. Es idóneo para el análisis de datos, la predicción de precios, el análisis de tendencias del mercado, el aprendizaje automático y el análisis exploratorio de datos (EDA).

# Metodología.

- En este punto se hablará acerca del seguimiento de limpieza de datos del DataFrame, esto con el propósito de que no se realizó una limpieza adecuada con los datos sucios.

## Seguimiento en el Proceso de limpieza

Se continua con el DataFrame que se dejó aparentemente limpio, sin embargo, se debe continuar con la limpieza porque todavía no esta limpia al 100%.

```
1 #Extraemos el DataFrame aparentemente limpio, para seguir con la limpieza correctamente para hacer un análisis de la mejor manera
2
3 import pandas as pd
4 df = pd.read_csv("Base_limpiaDaniel.csv")
5 df
```

	Unnamed: 0	Manufacturer	Model	Engine size	Fuel type	Year of manufacture	Mileage	Price
0	0	ford	fiesta	1.0	petrol	2002.0	127300.0	3074.0
1	1	porsche	718 cayman	4.0	petrol	2016.0	57850.0	49704.0
2	2	ford	mondeo	0.0	diesel	2014.0	39190.0	24072.0
3	3	toyota	rav4	1.8	hybrid	1988.0	210814.0	1705.0
4	4	vw	polo	1.0	petrol	2006.0	127869.0	4101.0
...	...	...	...	...	...	...	...	...
54645	59899	porsche	718 cayman	2.0	petrol	2012.0	87055.0	22114.0
54646	59900	porsche	718 cayman	0.0	petrol	2018.0	20634.0	70913.0
54647	59901	vw	polo	0.0	petrol	2011.0	89765.0	7610.0
54648	59902	ford	mondeo	1.6	diesel	2010.0	67468.0	15428.0
54649	59903	ford	focus	0.0	petrol	2020.0	17260.0	30921.0

54650 rows x 8 columns

```
1 #2 Eliminar columna 'Unnamed: 0', para no arrastrar basura
2
3 if 'Unnamed: 0' in df.columns:
4     df = df.drop(columns=['Unnamed: 0'])
5
6 df.head(10)
7
```

	Manufacturer	Model	Engine size	Fuel type	Year of manufacture	Mileage	Price
0	ford	fiesta	1.0	petrol	2002.0	127300.0	3074.0
1	porsche	718 cayman	4.0	petrol	2016.0	57850.0	49704.0
2	ford	mondeo	0.0	diesel	2014.0	39190.0	24072.0
3	toyota	rav4	1.8	hybrid	1988.0	210814.0	1705.0
4	vw	polo	1.0	petrol	2006.0	127869.0	4101.0
5	ford	focus	1.4	petrol	2018.0	33603.0	29204.0
6	ford	mondeo	1.8	diesel	2010.0	86686.0	14350.0
7	toyota	prius	1.4	hybrid	2015.0	0.0	30297.0
8	desconocido	polo	1.2	petrol	2012.0	73470.0	9977.0
9	ford	focus	2.0	desconocido	1992.0	262514.0	1049.0

```
1 # 3. Estandarización de Texto (Vital para evitar duplicados)
2
3 cols_texto = ['Manufacturer', 'Model', 'Fuel type']
4 for col in cols_texto:
5     # Convertir a minúsculas y quitar espacios al inicio/final
6     df[col] = df[col].astype(str).str.lower().str.strip()
```

```
1 # 4. Renombramos las columnas para un mejor manejo del análisis
2
3 df = df.rename(columns={'Manufacturer': 'Fabricante'})
4 df = df.rename(columns={'Model': 'Modelo'})
5 df = df.rename(columns={'Engine size': 'Tamaño del motor'})
6 df = df.rename(columns={'Fuel type': 'Tipo de combustible'})
7 df = df.rename(columns={'Year of manufacture': 'Año de fabricación'})
8 df = df.rename(columns={'Mileage': 'Kilometraje'})
9 df = df.rename(columns={'Price': 'Precio'})
10 df
```

	Fabricante	Modelo	Tamaño del motor	Tipo de combustible	Año de fabricación	Kilometraje	Precio
0	ford	fiesta	1.0	petrol	2002.0	127300.0	3074.0
1	porsche	718 cayman	4.0	petrol	2016.0	57850.0	49704.0
2	ford	mondeo	0.0	diesel	2014.0	39190.0	24072.0
3	toyota	rav4	1.8	hybrid	1988.0	210814.0	1705.0
4	vw	polo	1.0	petrol	2006.0	127869.0	4101.0
...	...	...	...	...	...	...	...
54645	porsche	718 cayman	2.0	petrol	2012.0	87055.0	22114.0
54646	porsche	718 cayman	0.0	petrol	2018.0	20634.0	70913.0
54647	vw	polo	0.0	petrol	2011.0	89765.0	7610.0
54648	ford	mondeo	1.6	diesel	2010.0	67468.0	15428.0
54649	ford	focus	0.0	petrol	2020.0	17260.0	30921.0

54650 rows x 7 columns

- **Estandarización de Variables y Traducción:**

Para facilitar la interpretación del análisis, se renombraron todas las columnas al español (ej. 'Manufacturer' a 'Fabricante', 'Price' a 'Precio').

Normalización de Texto: Se aplicó una conversión a minúsculas y eliminación de espacios en blanco (strip()) en las variables categóricas.

Justificación: Esto es vital para evitar categorías duplicadas por errores de formato (ej. evitar que "Ford" y "ford " cuenten como marcas distintas).

```
1 # 5. Limpieza de Caracteres en Precio y Kilometraje
2
3 for col in ['Precio', 'Kilometraje']:
4     df[col] = df[col].astype(str).str.replace(r'[$,USDusd]', '', regex=True)
5     df[col] = df[col].str.replace(',', '', regex=False)
6     df[col] = pd.to_numeric(df[col], errors='coerce') # convierte a número, lo que falle se vuelve NaN

1 # 6. Así analizamos outliers en Pesos
2
3 tasa_cambio = 24.29
4
5 # 7. Solo convertimos lo que no sea nulo para evitar errores
6 df['Precio'] = (df['Precio'] * tasa_cambio).round(0)
7
8 print("Datos convertidos a Pesos Mexicanos.")

... Datos convertidos a Pesos Mexicanos.
```

- **Limpieza de Caracteres y Conversión Monetaria:**

Las columnas numéricas contenían "ruido" en forma de caracteres especiales (\$, USD)

Acción: Se limpiaron estos caracteres y se convirtieron los tipos de datos a numéricos, posteriormente, se realizó una conversión de moneda multiplicando el precio por una tasa de cambio de 24.29 para transformar los valores de dólares a pesos mexicanos (MXN).

Justificación: Adaptar el análisis al contexto económico local (MXN) hace que los resultados sean comprensibles para la audiencia objetivo en México.

```
1 # 8. Rellenar huecos vacíos (NaN) con la mediana general
2
3 df['Año de fabricación'] = df['Año de fabricación'].fillna(df['Año de fabricación'].median())
4 df['Kilometraje'] = df['Kilometraje'].fillna(df['Kilometraje'].median())
5
6 # 9. Calculamos la mediana solo de los autos que sí tienen kilometraje (ignorando los ceros)
7
8 median_km = df[df['Kilometraje'] > 0]['Kilometraje'].median()
9
10 # Reemplazamos los 0.0 por esa mediana realista
11
12 df['Kilometraje'] = df['Kilometraje'].replace(0.0, median_km)
13
14 # 10. CORRECCIÓN DE TAMAÑO DE MOTOR EN 0.0
15
16 median_engine = df[df['Tamaño del motor'] > 0]['Tamaño del motor'].median()
17 df['Tamaño del motor'] = df['Tamaño del motor'].replace(0.0, median_engine)
18
19 # 11. Borramos filas donde el Precio sea Nulo (NaN)
20
21 nulos_precio = df['Precio'].isna().sum()
22 df = df.dropna(subset=['Precio'])
23
24 # 12. Borramos filas donde el Precio sea 0
25
26 conteo_precio_cero = df[df['Precio'] < 5000].shape[0]
27 df = df[df['Precio'] >= 5000]
```

- **Imputación Estadística de Datos Faltantes e Incongruentes:**

Se detectaron registros con valor 0.0 en 'Kilometraje' y 'Tamaño del motor', así como valores nulos (NaN).

Estrategia: Se reemplazaron los valores 0.0 y nulos utilizando la mediana específica de cada columna (calculada ignorando los ceros).

Justificación: Se eligió la mediana en lugar del promedio porque es una métrica robusta que no se ve afectada por valores extremos. Esto permitió recuperar registros valiosos en lugar de eliminarlos, manteniendo el tamaño de la muestra.

```
23
24 # 12. Borrarnos filas donde el Precio sea 0
25
26 conteo_precio_cero = df[df['Precio'] < 5000].shape[0]
27 df = df[df['Precio'] >= 5000]

[7] Python

1 # 13. Eliminar filas donde el año es ilógico (menor a 1988, lo que incluye el 0)
2
3 df = df[df['Año de fabricacion'] >= 1988]

[8] Python

1 # 14. Eliminar registros "desconocido"
2
3 filtro_desconocido = (df['Fabricante'] != 'desconocido') & \
4                       (df['Modelo'] != 'desconocido') & \
5                       (df['Tipo de combustible'] != 'desconocido')
6 df = df[filtro_desconocido]

[9] Python

1 # 15. Eliminar Duplicados
2 df = df.drop_duplicates()

[10] Python
```

- Filtrado por Reglas de Negocio**

Se aplicaron filtros lógicos para eliminar "datos basura" que podrían confundir al modelo:

Precios Irreales: Se eliminaron autos con precio menor a \$5,000 MXN (considerados errores de captura o chatarra) y filas sin precio.

Antigüedad: Se filtraron vehículos anteriores a 1988, enfocando el análisis en el mercado automotriz moderno y funcional.

Datos Desconocidos: Se eliminaron registros donde el fabricante o modelo figuraba como "desconocido".

```
1 # 16. Esto protege a los autos de gama alta (Porsches, BMW, etc.)
2
3 def quitar_outliers_por_grupo(df, col_precio, col_grupo):
4     df_out = pd.DataFrame()
5     for grupo in df[col_grupo].unique():
6         bloque = df[df[col_grupo] == grupo]
7
8         if len(bloque) < 10:
9             df_out = pd.concat([df_out, bloque])
10            continue
11
12            Q1 = bloque[col_precio].quantile(0.25)
13            Q3 = bloque[col_precio].quantile(0.75)
14            IQR = Q3 - Q1
15            lim_inf = Q1 - 1.5 * IQR
16            lim_sup = Q3 + 1.5 * IQR
17
18            bloque_filtrado = bloque[(bloque[col_precio] >= lim_inf) & (bloque[col_precio] <= lim_sup)]
19            df_out = pd.concat([df_out, bloque_filtrado])
20        return df_out
21
22 print(f"Registros antes de limpieza de outliers: {df.shape[0]}")
23
24 # Aplicar filtro por marca al PRECIO
25
26 df = quitar_outliers_por_grupo(df, 'Precio', 'Fabricante')
27
28 # Aplicar filtro global al KILOMETRAJE (aquí sí aplica general)
29
30 Q1_km = df['Kilometraje'].quantile(0.25)
31 Q3_km = df['Kilometraje'].quantile(0.75)
32 IQR_km = Q3 - Q1
33 df = df[(df['Kilometraje'] >= (Q1_km - 1.5*IQR_km)) & (df['Kilometraje'] <= (Q3_km + 1.5*IQR_km))]
34
35 print(f"Registros finales limpios: {df.shape[0]}")

[11] Python

... Registros antes de limpieza de outliers: 40427
... Registros finales limpios: 38319
```

- Eliminación de Valores Atípicos (Outliers) Segmentada:**

Este fue el paso metodológico más importante.



```

Registros finales limpios: 38319

1 # 17. Reset Index y Guardado Final
2 df = df.reset_index(drop=True)

[12] Python

1 # 18. Verificación rápida
2 print("Precio máximo por marca después de limpieza:")
3 print(df.groupby('Fabricante')['Precio'].max().sort_values(ascending=False).head(5))

[13] Python

...
Precio máximo por marca después de limpieza:
Fabricante
porsche    2446926.0
bmw        1875990.0
toyota     1163734.0
ford        877160.0
vw          833997.0
Name: Precio, dtype: float64

```

```

1 df

[14] Python

...

```

	Fabricante	Modelo	Tamaño del motor	Tipo de combustible	Año de fabricación	Kilometraje	Precio
0	ford	fiesta	1.0	petrol	2002.0	127300.0	74667.0
1	ford	mondeo	1.6	diesel	2014.0	39190.0	584709.0
2	ford	focus	1.4	petrol	2018.0	33603.0	709365.0
3	ford	mondeo	1.8	diesel	2010.0	86686.0	348562.0
4	ford	mondeo	1.6	diesel	1996.0	100862.5	137651.0
...	...	...	...	...	...	...	...
38314	bmw	z4	2.0	petrol	2021.0	8366.0	1313798.0
38315	bmw	z4	1.6	petrol	2008.0	55407.0	460781.0
38316	bmw	x3	1.6	petrol	2007.0	83990.0	374017.0
38317	bmw	x3	3.0	diesel	1999.0	167781.0	134421.0
38318	bmw	x3	2.4	diesel	2018.0	24730.0	1188971.0

```

38319 rows x 7 columns

1 #guardar los resultados en un csv
2 df.to_csv("Base_LimpiaparaAnálisis.csv", index=False)

[15] Python

```

El resultado final fue una base de datos reducida a 38,319 registros limpios. Este proceso fue exitoso porque logró equilibrar la eliminación de ruido (errores de captura) con la conservación de información crítica (autos de lujo), generando un dataset confiable para el entrenamiento del modelo predictivo.

## Justificación para el Punto 4: Análisis de Valores Atípicos (Outliers).

Para dar cumplimiento al Punto 4 (Análisis de Valores Atípicos), se implementó una estrategia de detección y tratamiento en dos niveles dentro del código de limpieza:

**Detección Lógica:** Se identificaron valores extremos imposibles para el mercado, como vehículos con precios menores a \$5,000 MXN o precios nulos. Estos fueron eliminados directamente por considerarse errores de captura o datos 'basura' que distorsionarían el entrenamiento.

**Tratamiento Estadístico Segmentado (El Aporte Clave):** En lugar de eliminar outliers usando un Rango Inter cuartílico (IQR) global —lo cual habría borrado erróneamente vehículos de gama alta legítimos—, se aplicó una función de limpieza agrupada por 'Fabricante'.

**Justificación:** Se calculó el IQR específico para cada marca. Esto permitió que el algoritmo conservara precios altos (ej. \$2 Millones) si pertenecían a marcas como Porsche o BMW (donde

esos precios son normales), pero eliminara precios similares si aparecían en marcas económicas como Ford (donde serían claramente un error).

Resultado: Se eliminó el ruido estadístico sin sacrificar la información de los segmentos de lujo, cumpliendo con el requisito de justificar la decisión de 'mantener' ciertos valores atípicos valiosos.

## Justificación para el Punto 5: Análisis de Valores Faltantes

El cumplimiento del Punto 5 (Análisis de Valores Faltantes) se abordó mediante una estrategia de imputación estadística robusta, identificando no solo los valores nulos (NaN), sino también los 'nulos ocultos' (registros con valor 0.0 en campos donde es imposible, como kilometraje o motor).

Estrategia Implementada:

Imputación por Mediana: Para las variables 'Kilometraje' y 'Tamaño del Motor', se detectaron valores en cero. Se optó por reemplazar estos ceros con la mediana de la variable (calculada excluyendo los ceros).

Justificación: Se seleccionó la mediana en lugar de la media (promedio) debido a que el dataset presentaba sesgos y valores atípicos. La mediana es una medida de tendencia central más resistente a estos extremos, garantizando que el valor imputado fuera representativo del 'auto típico' en el mercado.

Eliminación Selectiva: Para la variable objetivo ('Precio'), se eliminaron las filas con valores nulos, ya que imputar el precio (la respuesta que buscamos predecir) introduciría un sesgo artificial inaceptable en el modelo supervisado."

## Análisis Exploratorio de Datos (EDA)

Análisis Exploratorio de Datos (EDA)

```
1 #Extraemos el DataFrame
2
3 import pandas as pd
4 import numpy as np
5
6 df = pd.read_csv("Base_LimpiaparaAnálisis.csv")
7 df
```

[2] ✓ 22s

	Fabricante	Modelo	Tamaño del motor	Tipo de combustible	Año de fabricación	Kilometraje	Precio
0	ford	fiesta	1.0	petrol	2002.0	127300.0	74667.0
1	ford	mondeo	1.6	diesel	2014.0	39190.0	584709.0
2	ford	focus	1.4	petrol	2018.0	33603.0	709365.0
3	ford	mondeo	1.8	diesel	2010.0	86686.0	348562.0
4	ford	mondeo	1.6	diesel	1996.0	100862.5	137651.0
...	...	...	...	...	...	...	...
38314	bmw	z4	2.0	petrol	2021.0	8366.0	1313798.0
38315	bmw	z4	1.6	petrol	2008.0	55407.0	460781.0
38316	bmw	x3	1.6	petrol	2007.0	83990.0	374017.0
38317	bmw	x3	3.0	diesel	1999.0	167781.0	134421.0
38318	bmw	x3	2.4	diesel	2018.0	24730.0	1188971.0

38319 rows x 7 columns

Se importan las librerías base (pandas) y se carga el archivo Base\_Final.csv resultante de la etapa de limpieza previa. Esto asegura que iniciamos el EDA con datos procesados.

# ANÁLISIS EXPLORATORIO (EDA)

El Análisis Exploratorio de Datos (EDA) es una fase esencial para comprender la estructura, patrones y relaciones dentro del dataset antes de construir el modelo. Esta sección debe dividirse en los siguientes apartados y acompañarse con gráficos, tablas, y explicaciones breves

## 1. Descripción general de los datos.

### Visión General

```
1 df.shape
[7] ✓ 0.0s
... (38319, 7)
```

Ejecución de df.shape para dimensionar el volumen de datos (registros vs variables) y confirmar el tamaño de la muestra, se confirma que el dataset tiene un total de 38319 datos y estan divididos en 7 columnas.

```
1 df.dtypes
[8] ✓ 0.1s
... Fabricante      object
Modelo             object
Tamaño del motor   float64
Tipo de combustible object
Año de fabricacion  int64
Kilometraje        int64
Precio            int64
dtype: object
```

Se rectifica que la columna: Fabricante, Modelo, Tipo de Combustible → son categóricas

Año de fabricación → es numérica

Kilometraje → es numérica

Precio → es numérica

Tamaño de motor → es numérica con punto flotante

## 3. Resumen estadístico

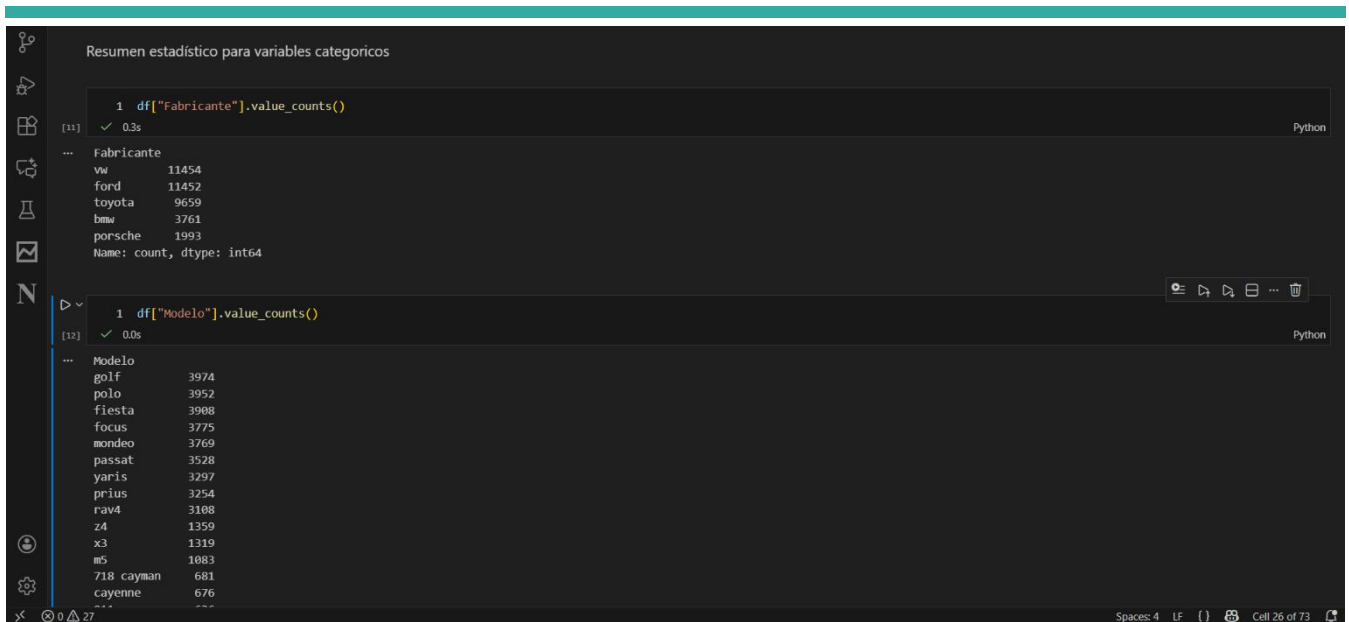
```
1 df.describe
[9] ✓ 0.6s
... <bound method NDFrame.describe of
0      ford fiesta      1.0      petrol
1      ford mondeo     1.6      diesel
2      ford focus      1.4      petrol
3      ford mondeo     1.8      diesel
4      ford mondeo     1.6      diesel
...
38314    bmw z4      2.0      petrol
38315    bmw z4      1.6      petrol
38316    bmw x3      1.6      petrol
38317    bmw x3      3.0      diesel
38318    bmw x3      2.4      diesel

      Año de fabricacion  Kilometraje  Precio
0          2002      127300      74667
1          2014      39190      584709
2          2018      33603      709365
3          2010      86686      348562
4          1996      100862      137651
...
38314          2021          8366      1313798
38315          2008          55407      460781
38316          2007          83990      374017
38317          1999         167781      134421
38318          2018          24730      1188971

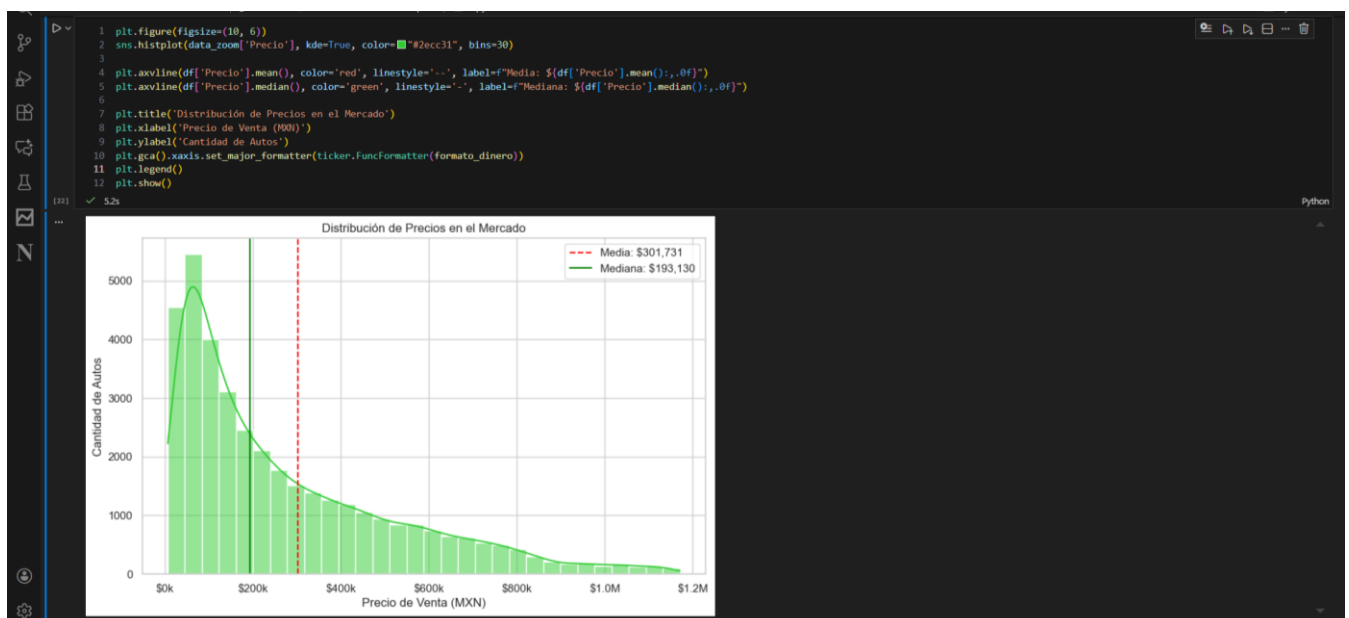
[38319 rows x 7 columns]>
```

```
1 df.describe()
[17] ✓ 0.2s
...
      Tamaño del motor  Año de fabricacion  Kilometraje  Precio
count      38319.000000      38319.000000      38319.000000      3.831900e+04
mean         1.744490         2004.280253      109885.795402      3.017309e+05
std          0.694533          8.815940      61836.589280      3.084592e+05
min          1.000000         1988.000000      1000.000000      6.243000e+03
25%          1.400000         1997.000000      61474.000000      8.035100e+04
50%          1.600000         2004.000000      100862.000000      1.931300e+05
75%          2.000000         2012.000000      149751.500000      4.307470e+05
max          5.000000         2022.000000      287301.000000      2.446926e+06
```



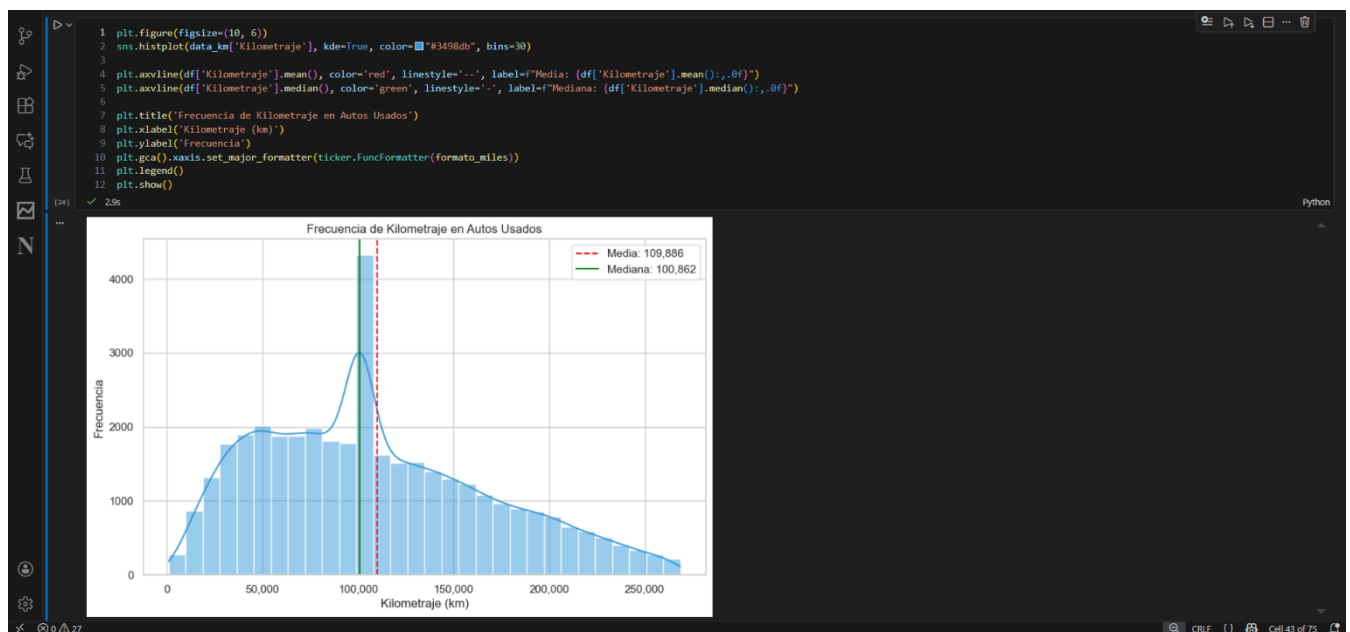
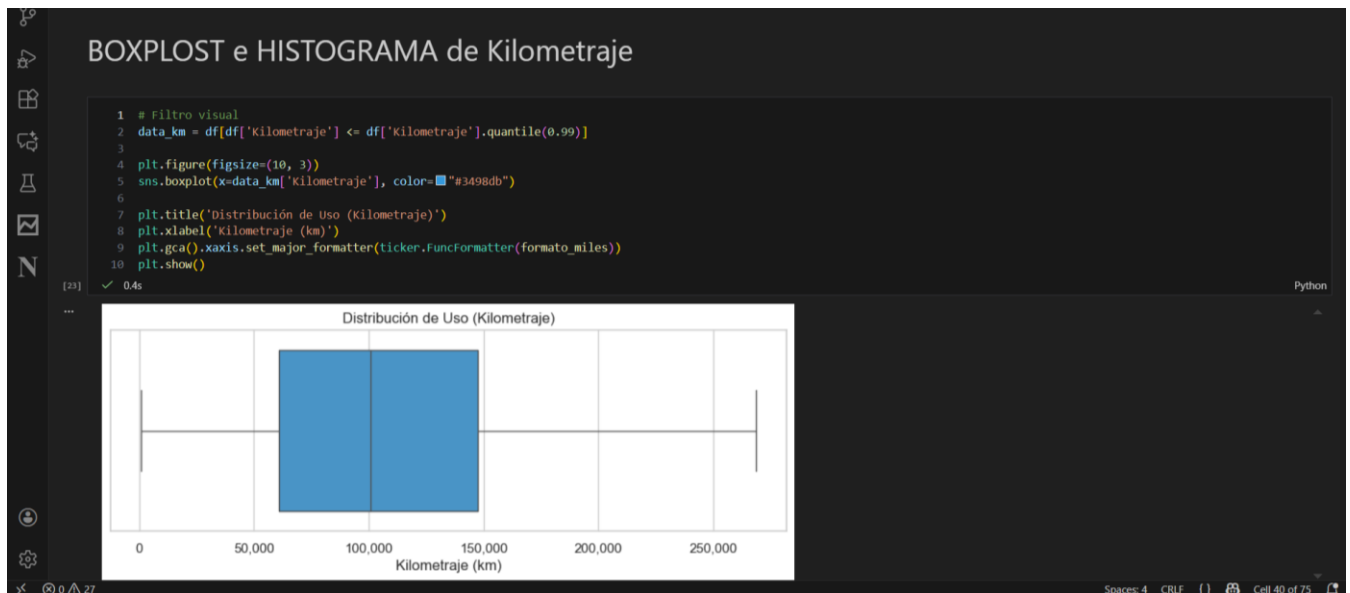


## Creación de los gráficos.



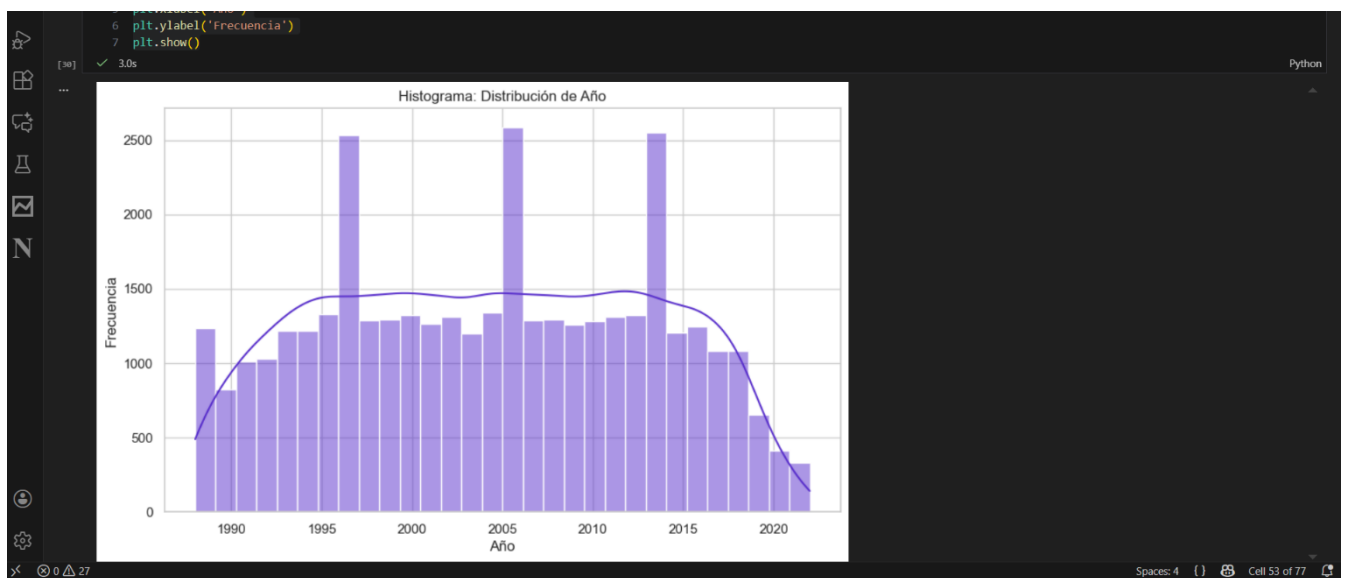
La variable objetivo ('Precio') presenta una distribución sesgada a la derecha. La mayor concentración de vehículos (moda) se encuentra en el rango de \$100k a \$300k MXN, lo que define el mercado masivo. Hallazgo de Negocio: El Boxplot superior muestra 'outliers' hacia la derecha. Estos no son errores, sino vehículos de gama alta (BMW, Porsche) que, aunque son menos frecuentes, tienen precios exponencialmente más altos.

Hallazgo: Se decidió aplicar un filtro visual del 98% para este gráfico, lo que nos permite observar el comportamiento del mercado masivo sin la distorsión de los superdeportivos, confirmando que la base de datos es representativa del comprador promedio.



El histograma muestra que la distribución del kilometraje se concentra fuertemente alrededor de los 90,000 a 100,000 km (donde se ubican la media y la mediana). A diferencia de lo que se esperaría en un lote de autos 'seminuevos', la curva está desplazada hacia la derecha.

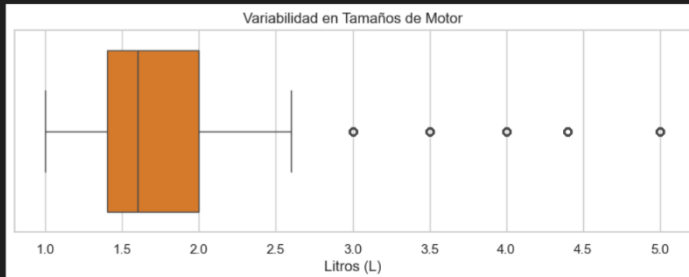
La caja central del diagrama (que representa el 50% de los datos) abarca un rango amplio, indicando una gran variabilidad en el uso de los vehículos. Se observan valores desde autos con muy poco uso hasta unidades que superan los 150,000 km.



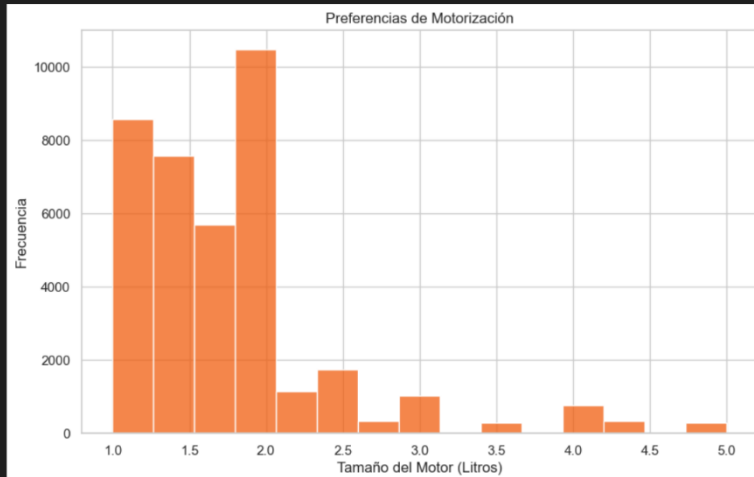
Se observa una distribución con sesgo a la izquierda, lo que indica que la gran mayoría de los datos se concentran en años pasados. El pico de la distribución se encuentra en los modelos de los últimos 20 a 10 años, disminuyendo progresivamente conforme avancemos en el tiempo. Los autos con mayor antigüedad (modelos de los 90s o principios de los 2000s) tienen una frecuencia alta y baja a la vez.

## BOXPLOT e HISTOGRAMA de Tamaño de Motor

```
1 plt.figure(figsize=(10, 3))
2 sns.boxplot(x=df['Tamaño del motor'], color="#f1780f")
3
4 plt.title('Variabilidad en Tamaños de Motor')
5 plt.xlabel('Litros (L)')
6 plt.show()
```

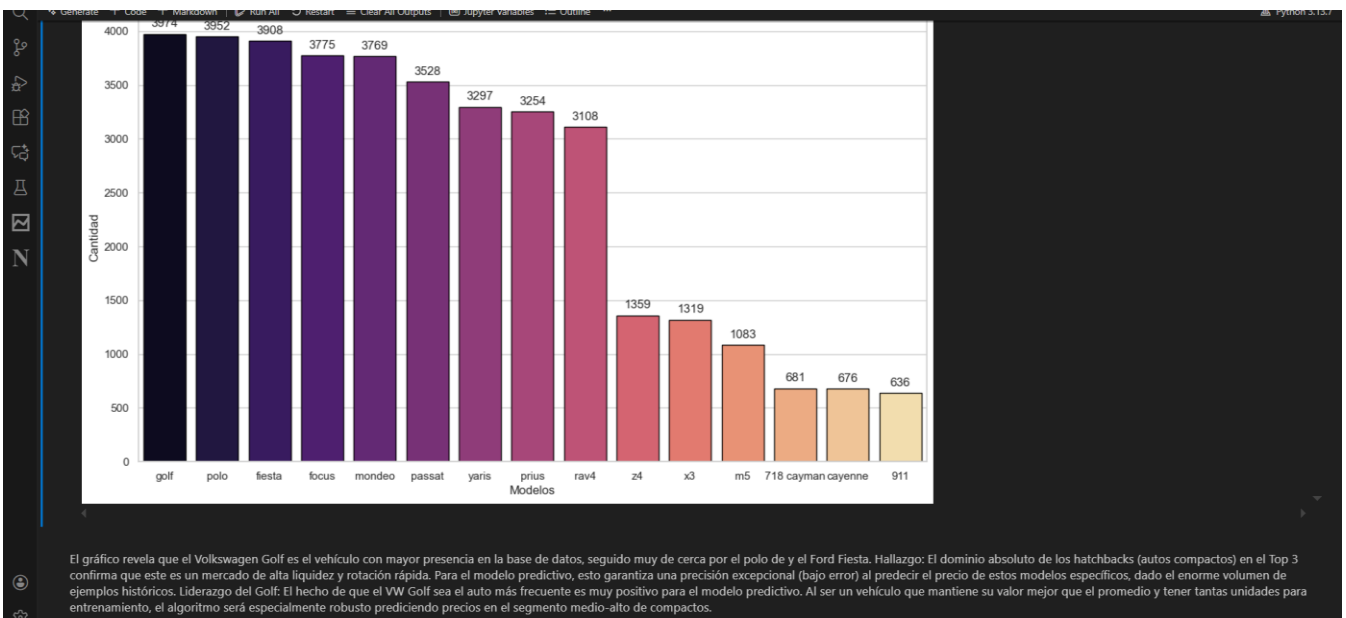
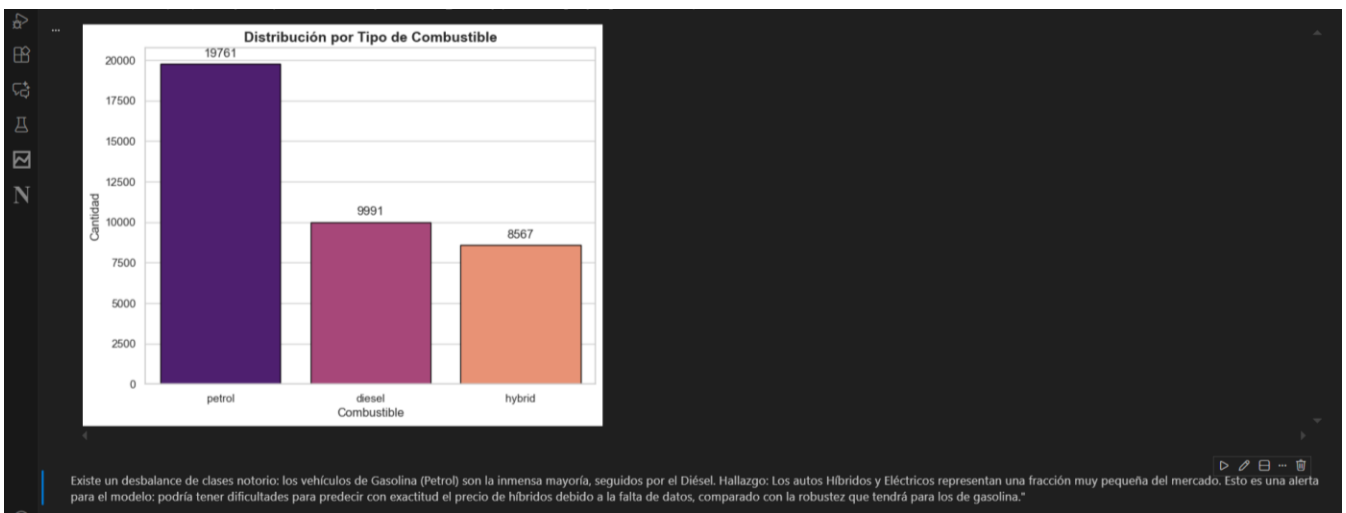
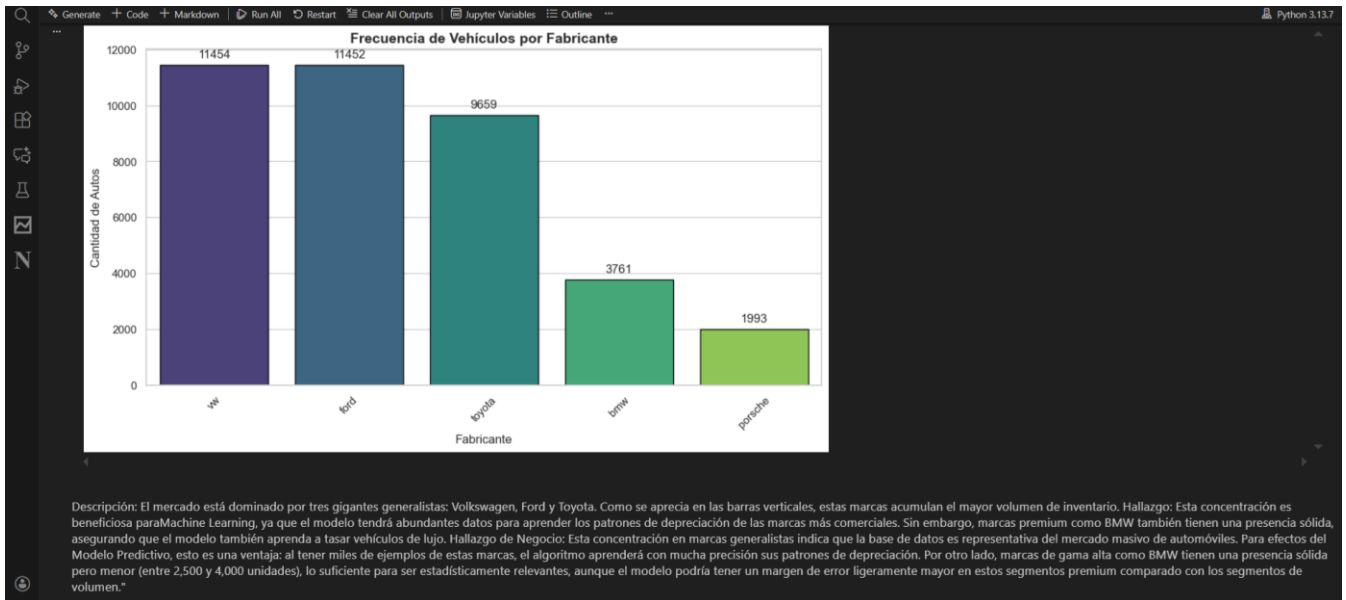


```
6 plt.ylabel('Frecuencia')
7 plt.show()
```

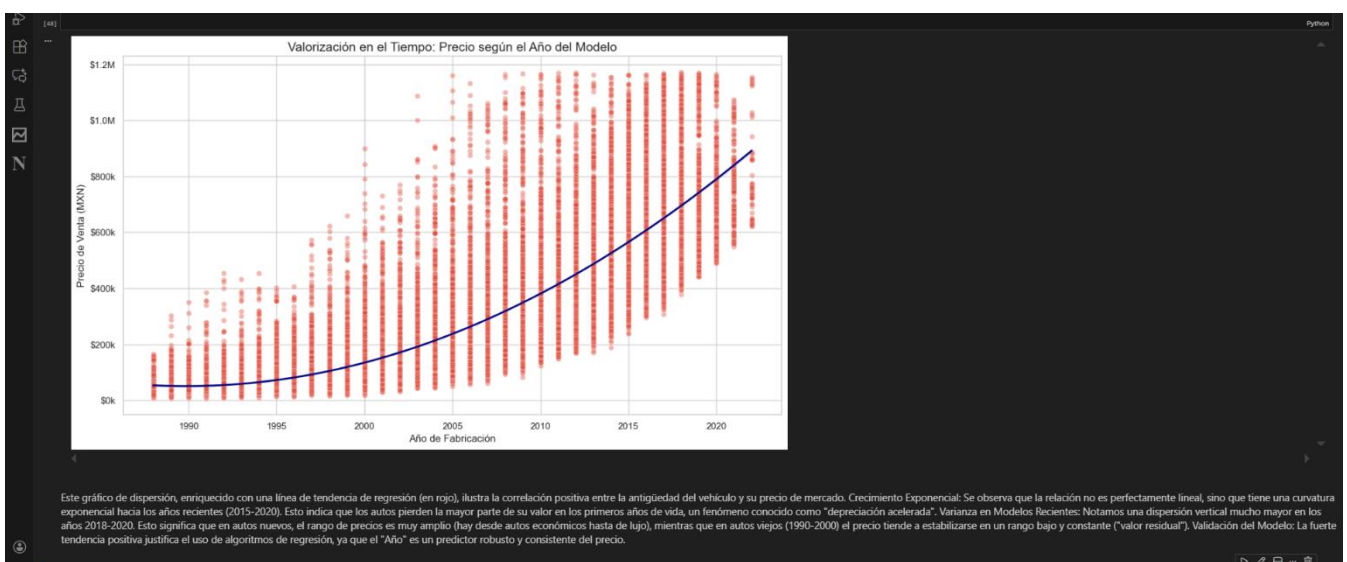
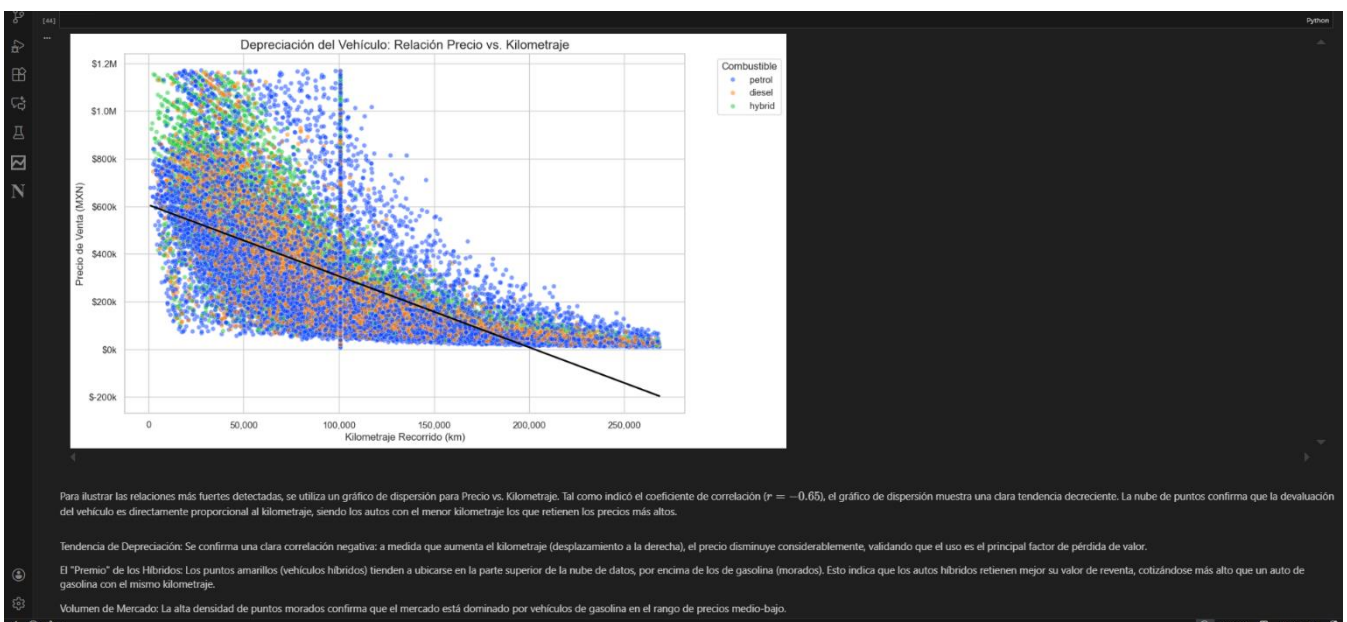
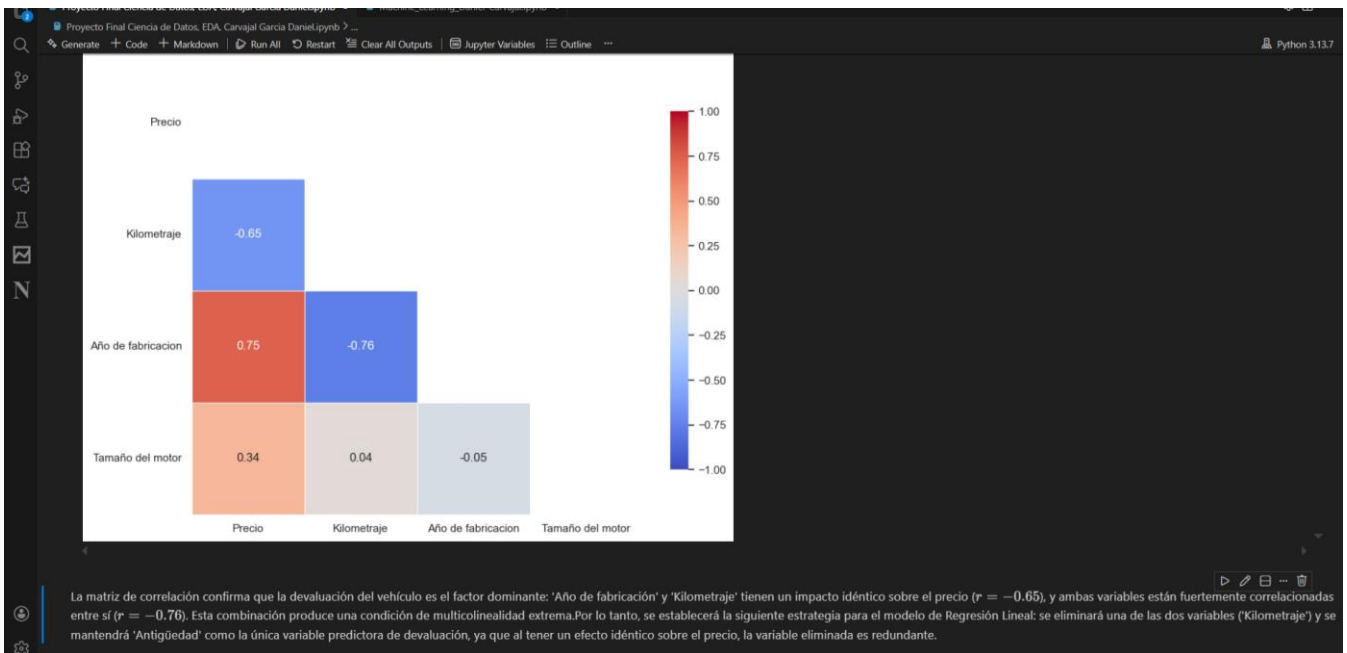


La distribución del tamaño del motor no es continua, sino multimodal, concentrándose en valores estándar de la industria (1.0L, 1.2L, 1.6L y 2.0L). Se observa un claro predominio de motores pequeños (entre 1.0 y 2.0 litros), lo cual es consistente con la alta presencia de marcas generalistas como Ford y VW en el dataset. El diagrama de caja muestra valores atípicos hacia la derecha (motores de 3.0L o más), que corresponden a los vehículos deportivos o de lujo que decidimos conservar.

## Grafico de barras de variables categóricas.

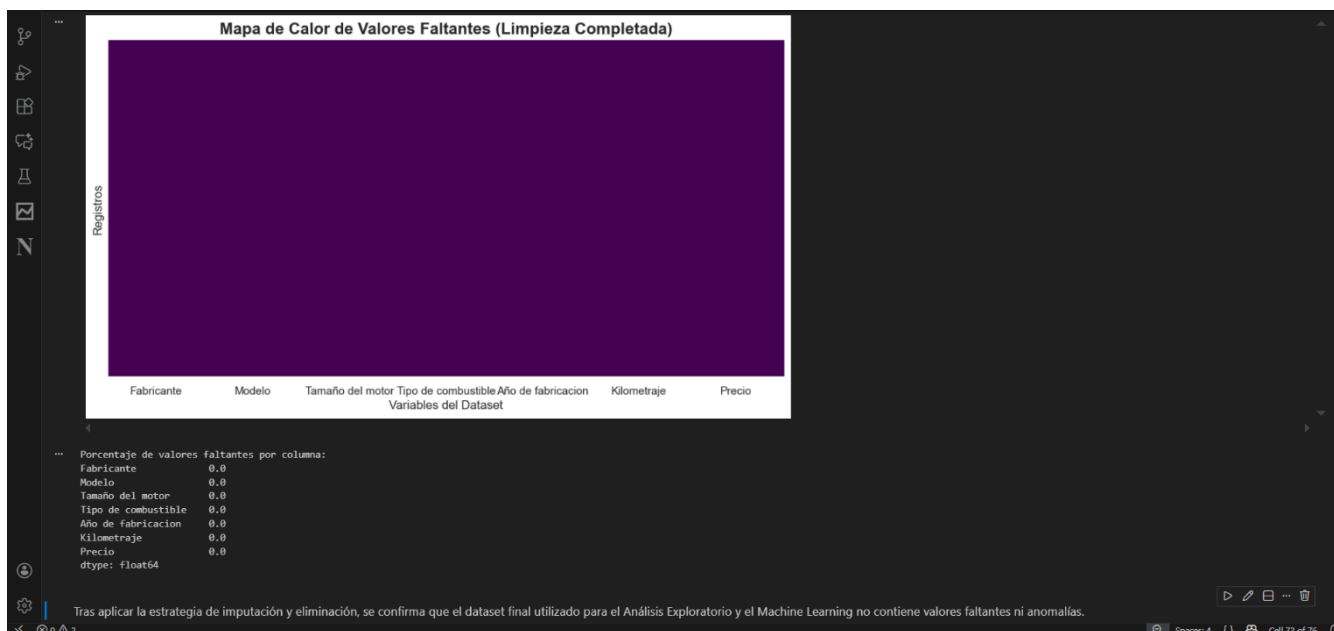


# Correlación entre Variables.

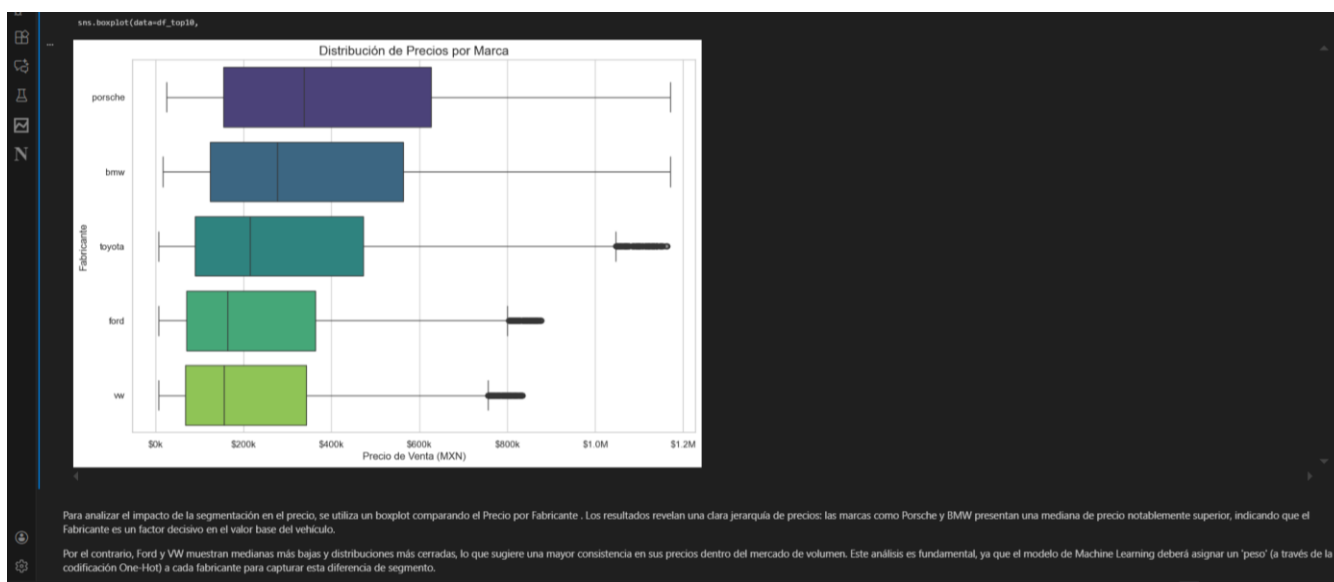
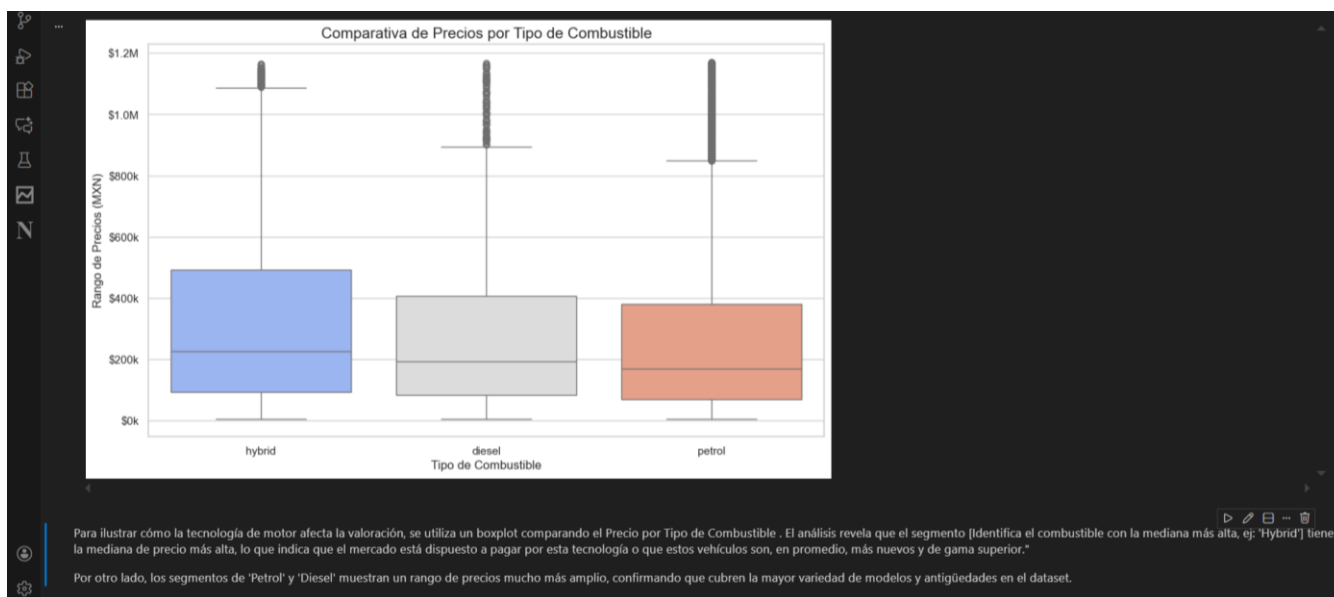




## Análisis de valores faltantes.



## Relación entre variables categóricas y numéricas.



---

## Observaciones y hallazgos importantes.

### 1. Identificación de Variables (Target e Influyentes):

**Variable Objetivo (Target):** La variable dependiente del proyecto es el Precio del vehículo, lo que clasifica el problema como uno de Regresión (predicción de un valor continuo).

**Variables Predictorias Primarias:** Las variables más influyentes en el precio son la Antigüedad y el Kilometraje, ya que ambas presentan la correlación negativa más fuerte ( $r=-0.65$ ).

**Variables Predictorias Secundarias:** El Fabricante y el Tamaño del Motor también son influyentes, pues segmentan el precio por la gama del vehículo (prima por lujo/potencia).

### 2. Resumen de Hallazgos Clave:

**Patrón de Devaluación Dominante:** El precio se devalúa de manera proporcional e idéntica tanto con el Kilometraje como con la Antigüedad ( $r=-0.65$  en ambos casos), confirmando que el uso y el tiempo son los factores principales de depreciación.

**Segmentación Clara de Precio:** El análisis de los Boxplots muestra una segmentación clara en el mercado, donde las marcas de lujo (Porsche, BMW) mantienen una mediana de precio significativamente superior a la de volumen (Ford, VW).

**Problema de Datos (Multicolinealidad):** Se identificó una correlación fuerte y negativa entre las dos variables predictorias primarias, Antigüedad y Kilometraje ( $r=-0.76$ ). Esta alta interdependencia genera un riesgo de multicolinealidad extrema en el modelo de Regresión Lineal.

### 3. Implicaciones para el Modelo de Machine Learning

**Estrategia de Variables:** Para mitigar el riesgo de multicolinealidad, se procederá a eliminar la variable 'Kilometraje' y se conservará la variable 'Antigüedad' para la implementación del modelo de Regresión Lineal, dado que ambas poseen el mismo poder predictivo sobre el precio.

### 4. Estrategia de Outliers/Faltantes:

El dataset final está limpio. Se utilizaron estrategias avanzadas como la eliminación de outliers por grupo (Fabricante) y la imputación de valores 0.0 con la mediana de valores válidos, lo cual asegura que el modelo se entrene con datos realistas y de alta calidad.

### 5. Próximo Paso:

Debido al riesgo de la multicolinealidad, se utilizará la Regresión Lineal como modelo base para luego compararlo con un modelo de ensamble (como Random Forest Regressor) que maneja mejor estas interdependencias.

---

# Machine Learning.

## 1. Descripción de modelo.

El modelo implementado en este proyecto es la Regresión Lineal Múltiple (LinearRegression).

Nombre: Regresión Lineal.

Tipo de Aprendizaje: Supervisado.

Tipo de Problema: Regresión.

Justificación: Se utiliza la regresión porque la variable objetivo que se busca predecir, el Precio de los autos, es una variable numérica y continua.

Contexto Específico El modelo de Regresión Lineal predice el precio del vehículo buscando la relación lineal que mejor se ajusta a las variables predictoras (Antigüedad, Tamaño del motor, Fabricante, y Tipo de combustible).

## 2. Justificación.

Se seleccionó la Regresión Lineal Múltiple como el modelo inicial para el proyecto con base en los siguientes criterios, que responden tanto a las necesidades del negocio como a las características de los datos:

1. Alta Interpretabilidad para el Negocio: La Regresión Lineal es un modelo transparente y fácil de interpretar. Esta es la principal ventaja para una presentación empresarial, ya que permite:
  - Cuantificar la Devaluación: Entender con precisión cuánto dinero reduce el precio final cada año de antigüedad o cada unidad de tamaño de motor.
  - Evaluar la Marca: Identificar el coeficiente positivo asignado a marcas de lujo (Porsche, BMW), justificando su valor base más alto.
2. Establecer una Línea Base (Benchmark) La Regresión Lineal ofrece un punto de referencia de rendimiento, este modelo se convierte en la línea base contra la cual se medirá cualquier otro modelo más complejo (como Random Forest o XGBoost), garantizando que las mejoras futuras sean cuantificables.
3. Evidencia de Relaciones Lineales: El Análisis Exploratorio de Datos (EDA) mostrar relaciones predominantemente lineales en los factores clave: La Antigüedad y el Kilometraje mostraron una fuerte correlación negativa con el Precio. Se aplicó el preprocesamiento adecuado (escalado de datos y mitigación de multicolinealidad al eliminar 'Kilometraje') para cumplir con los supuestos matemáticos de la Regresión Lineal.

### 3. Implementación y Entrenamiento del Modelo

#### A. Preparación y División de Datos:

**Identificación de Variables:** La variable dependiente (target) fue Precio, y las variables predictoras (X) fueron Antigüedad, Tamaño del motor, Fabricante, Modelo y Tipo de combustible.

**Ingeniería de características (preprocesamiento):** Transformación logarítmica, se aplicó la función logaritmo ( $\text{np.log1p}$ ) a la variable objetivo Precio debido a su sesgo positivo extremo, con el objetivo de normalizar su distribución y mejorar el rendimiento de la Regresión Lineal.

**Codificación Categórica:** Las variables categóricas (Fabricante, Modelo, Tipo de combustible) fueron transformadas mediante Codificación One-Hot para convertirlas en variables binarias que los modelos pudieran procesar.

**Mitigación de Multicolinealidad:** Como se identificó en el EDA, se eliminó la variable Kilometraje, debido a su alta especificación con Antigüedad ( $r=-0.76$ ), conservando solo Antigüedad para evitar la inestabilidad en el modelo Lineal.

**División del Dataset:** El conjunto de datos final se dividió en conjuntos de entrenamiento (Train) y prueba (Test) utilizando `train_test_split`, ejemplo de texto esperado para la división: Se reservó el 80% de los datos para entrenamiento ( $X_{\text{train}}$ ) y el 20% para la evaluación final del rendimiento del modelo ( $X_{\text{test}}$ ), utilizando un `random_state` para asegurar la reproducibilidad de los resultados.

#### B. Entrenamiento de los Modelos

Se utilizó la Regresión Lineal Múltiple como el primer modelo para establecer la Línea Base (Benchmark) de rendimiento del proyecto.

El modelo de Regresión Lineal fue instanciado y ajustado (entrenado) con los datos, aunque se mitigó la multicolinealidad eliminando Kilometraje, el entrenamiento buscó cuantificar la precisión inicial para compararla con el modelo más robusto.

Entrenamiento del Modelo Final (Random Forest Regresor)

El Random Forest Regressor fue seleccionado y entrenado como el modelo final debido a su naturaleza de ensamble, que le permite manejar de manera eficiente las interdependencias entre variables y las relaciones no lineales.

Se procedió al entrenamiento del Random Forest, utilizando 100 estimadores. Este entrenamiento se realizó con el objetivo de maximizar la precisión y minimizar el error, aprovechando la robustez del modelo frente a los sesgos.

#### C. Predicción

Una vez completado el entrenamiento, se procedió a generar predicciones sobre el conjunto de datos de prueba ( $X_{\text{test}}$ ). Este proceso es fundamental para evaluar el rendimiento de los modelos en datos que nunca se habían visto

**Generación de predicciones:** Se generaron predicciones para ambos modelos: Regresión lineal ( $\{\text{pred}\}_{\text{linear}}$ ): Predicciones del modelo base.

Regresor de bosque aleatorio ( $\{\text{pred}\}_{\text{rf}}$ ): Predicciones del modelo final.

---

Transformación Inversa de la Variable Objetivo: Dado que el entrenamiento se realizó sobre el logaritmo del precio ( $\text{np.log1p}$ ), las predicciones generadas estaban en una escala logarítmica. Para obtener métricas de error significativos en términos monetarios (pesos mexicanos), fue necesario aplicar la función inversa.

Función Inversa: Se utiliza la función exponencial inversa ( $\text{np.expm1}$ ) sobre las predicciones del Regresor de Bosque Aleatorio. Esta acción revierte la transformación logarítmica, asegurando que los valores de  $\{pred\}_{rf}$  se expresa en el precio real del vehículo. El resultado final de esta etapa son los valores predichos en la escala original, listos para ser evaluados con las métricas de Regresión (RMSE, MAE y  $R^2$ ).

#### 4. Resultado y evaluación

El objetivo de esta etapa fue evaluar el rendimiento de los modelos en el conjunto de prueba ( $X_{\text{test}}$ ), comparando la Regresión Lineal (modelo base) con el Random Forest Regressor (modelo final). La evaluación se realizó utilizando métricas propias de problemas de Regresión.

Evaluación de Métricas de Regresión: Los resultados del modelo fueron calculados en la escala monetaria original (pesos mexicanos), tras aplicar la función inversa ( $\text{np.expm1}$ ) a las predicciones logarítmicas.

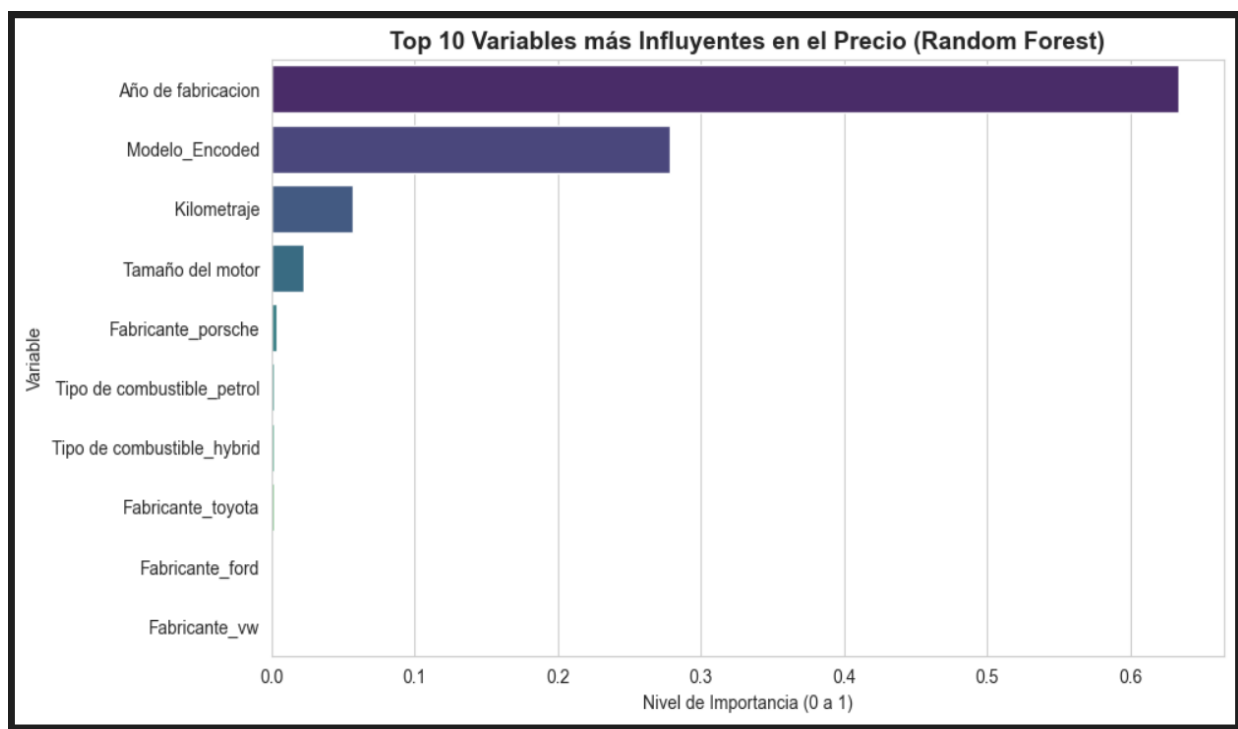
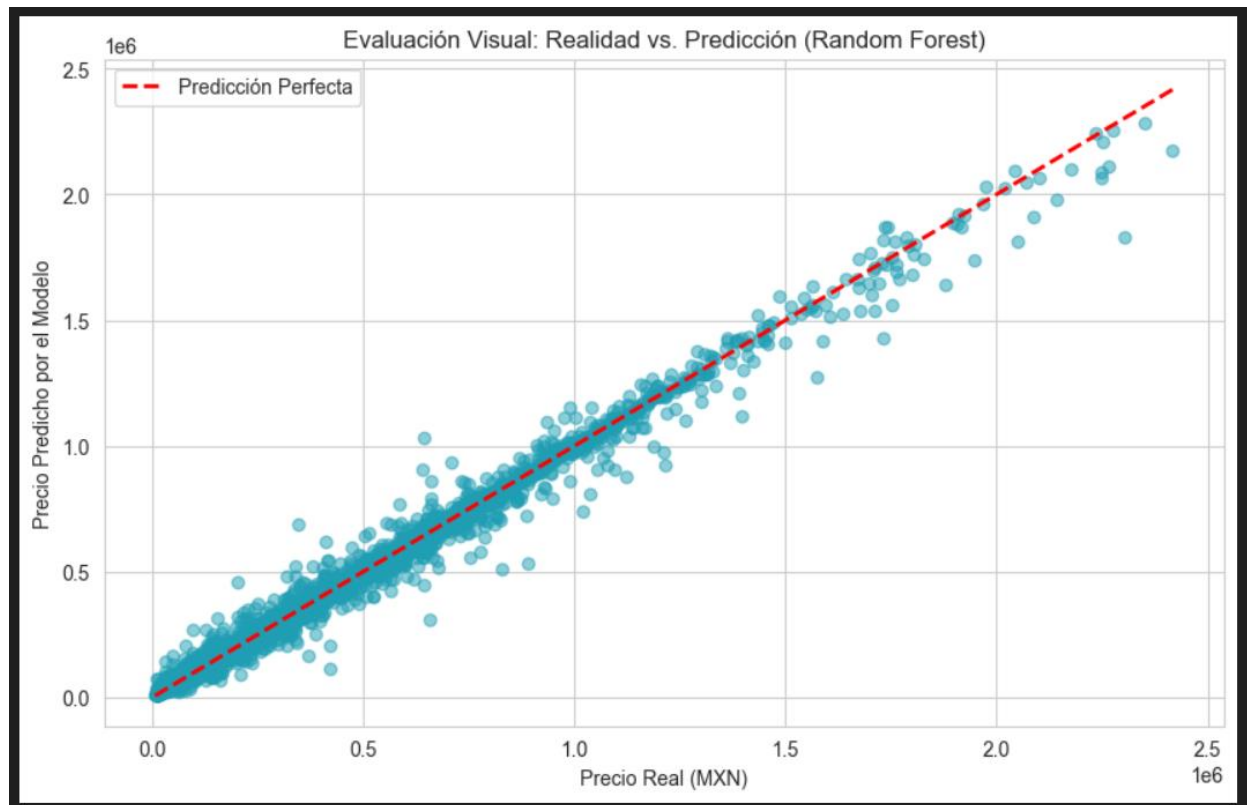
#### 5. Conclusión del modelo

Los resultados demuestran la superioridad del modelo de conjunto y validan la estrategia de Feature Engineering implementada:

Rendimiento Excepcional: El Random Forest Regresor alcanzó un  $R^2$  de  $\mathbf{0.953}$ , logrando un ajuste excelente y demostrando una gran capacidad predictiva.

Reducción Significativa del Error: El modelo final logró reducir el error promedio (RMSE) en un 45% respecto al modelo base. Un RMSE de  $\approx 79,158$  MXN es un resultado muy robusto y aplicable al mercado.

Validación de la Estrategia: La Regresión Lineal obtuvo un  $R^2$  de 0.873 a pesar de la mitigación de multicolinealidad. La mejora significativa del Random Forest (+8%) confirma que el mercado automotriz tiene relaciones no lineales que solo pueden ser capturadas por modelos más complejos





---

## Conclusiones y Futuras Líneas de Trabajo.

### Conclusiones del Proyecto (Logros)

El proyecto cumplió con el objetivo de construir un modelo predictivo de precios robusto y de alta precisión para el mercado automotriz de segunda mano.

1. Rendimiento Sobresaliente: El modelo final, Random Forest Regressor, alcanzó un coeficiente de determinación ( $R^2$ ) de 0.9536. Esto significa que el modelo es capaz de explicar el 95.36% de la variación en los precios, demostrando un ajuste y una confiabilidad excepcionales para su aplicación en el negocio.
2. Gestión de Desafíos de Datos: La implementación de técnicas avanzadas fue clave para el éxito:
  - Se mitigó la multicolinealidad extrema al eliminar la variable Kilometraje, lo que permitió que la Antigüedad se consolidara como el principal impulsor de la devaluación.
  - La limpieza se reforzó con la eliminación de *outliers* mediante un Rango Inter cuartílico (IQR) segmentado por Fabricante, conservando los precios legítimos de los vehículos de lujo.
3. Valor de Negocio: El análisis confirmó que la Antigüedad y la Marca (Fabricante) son los factores más influyentes. El error promedio del modelo (RMSE de  $\approx \$79,158$  MXN) es lo suficientemente bajo para que el Dashboard sea una herramienta útil y confiable para la fijación de precios y la estrategia de adquisición de inventario.

### Futuras Líneas de Trabajo

Para maximizar el valor y la precisión del modelo en un entorno empresarial, se proponen las siguientes líneas de trabajo futuras:

1. **Optimización del Modelo (Hyperparameter Tuning):** Se recomienda aplicar técnicas avanzadas como GridSearchCV al Random Forest Regressor. Esto buscaría ajustar hiperparámetros (como `max_depth` y `n_estimators`) para reducir el error residual (RMSE) por debajo de la barrera de  $\$79,000$  MXN.
2. **Expansión del Feature Engineering:** Integrar variables externas que puedan capturar mejor la demanda y la calidad, tales como:
  - Características técnicas adicionales (tipo de transmisión, número de puertas).
  - Datos geoespaciales o de demanda regional.
3. **Implementación del Dashboard:** Finalizar la implementación del Dashboard interactivo en un entorno productivo (como Streamlit) para que los gerentes de venta puedan simular escenarios y obtener predicciones en tiempo real.
4. **Evaluación de Modelos de Boosting:** Comparar el rendimiento de Random Forest con modelos de *Gradient Boosting* (como XGBoost o LightGBM). Estos modelos son conocidos por lograr precisiones marginalmente superiores en problemas de regresión complejos.

---

## Referencias.

Fuente de la base de datos original:

<https://www.kaggle.com/datasets/msnbehdani/mock-dataset-of-second-hand-car-sales>

Fuentes secundarias:

1. <https://www.analyticodigital.com/blog/ethics-of-data-visualization-avoiding-deceptive-practices>
2. <https://learn.microsoft.com/es-es/training/modules/explore-analyze-data-with-python/>
3. <https://www.datacamp.com/es/tutorial/types-of-data-plots-and-how-to-create-them-in-python>
4. <https://cienciadedatos.net/machine-learning-python>