

---

**PROYECTO DE MACHINE LEARNING: ESTIMACIÓN DE PRECIOS DE AUTOMÓVILES**

---

**Alumno:**

**Jose de Jesus Rodriguez Aparicio**

---

**Fecha:**

**16 de Agosto de 2024**

## Índice

1. Introducción
2. Objetivo del Proyecto
3. Avance 1: Preguntas y Respuestas
4. Avance 2: Preprocesamiento de Datos y Separación en Conjuntos de Entrenamiento y Prueba
5. Avance 3: Entrenamiento de Algoritmos de Regresión
6. Avance 4: Entrenamiento de Algoritmos Random Forest y XGBoost
7. Comparación y Selección del Mejor Algoritmo
8. Conclusiones

## Introducción

En el presente proyecto, se aborda la problemática de la estimación de precios de automóviles utilizando técnicas de machine learning. La correcta estimación de precios es fundamental en la industria automotriz para facilitar las decisiones de compra y venta, así como para establecer precios competitivos en el mercado. Se implementaron y compararon diversos algoritmos de regresión para determinar cuál proporciona la mejor precisión en las estimaciones. Este proyecto incluye el preprocesamiento de datos, la implementación de algoritmos de k vecinos cercanos, Random Forest y XGBoost, y una comparación exhaustiva de su desempeño.

## Objetivo del Proyecto

El objetivo principal del proyecto es desarrollar un modelo de machine learning capaz de predecir con precisión los precios de los automóviles utilizando diversas características del vehículo. Para lograr esto, se implementaron y compararon diferentes algoritmos de regresión, incluyendo k vecinos cercanos (KNN), Random Forest y XGBoost, con el fin de seleccionar el modelo que ofrezca el mejor rendimiento en términos de precisión y eficiencia.

a. ¿Qué tipos de datos existen en el conjunto?

```
print(df.dtypes)
Simbolización          int64
Pérdidas normalizadas  float64
Marca                  object
Tipo de combustible    object
Tipo de aspiración     object
Número de puertas      object
Tipo de auto           object
Tracción               object
Ubicación del motor    object
Distancia entre ejes   float64
Longitud               float64
Ancho                  float64
Alto                   float64
Peso                   int64
Tipo de motor          object
Cilindraje             int64
Tamaño del motor       int64
Sistema de combustible object
Diámetro de cilindro   float64
Carrera del pistón     float64
Índice de compresión   float64
Caballos de fuerza     float64
Rpm pico               float64
MPG en ciudad          int64
```

```

MPG en carretera          int64
Precio                   float64
dtype: object
df.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 205 entries, 0 to 204
Data columns (total 26 columns):
#   Column                      Non-Null Count  Dtype
---  -
0   Simbolización              205 non-null    int64
1   Pérdidas normalizadas     164 non-null    float64
2   Marca                      205 non-null    object
3   Tipo de combustible       205 non-null    object
4   Tipo de aspiración        205 non-null    object
5   Número de puertas         203 non-null    object
6   Tipo de auto              205 non-null    object
7   Tracción                  205 non-null    object
8   Ubicación del motor       205 non-null    object
9   Distancia entre ejes     205 non-null    float64
10  Longitud                  205 non-null    float64
11  Ancho                    205 non-null    float64
12  Alto                    205 non-null    float64
13  Peso                    205 non-null    int64
14  Tipo de motor            205 non-null    object
15  Cilindraje              205 non-null    int64
16  Tamaño del motor        205 non-null    int64
17  Sistema de combustible  205 non-null    object
18  Diámetro de cilindro    201 non-null    float64
19  Carrera del pistón      201 non-null    float64
20  Índice de compresión     205 non-null    float64
21  Caballos de fuerza      203 non-null    float64
22  Rpm pico                203 non-null    float64
23  MPG en ciudad           205 non-null    int64
24  MPG en carretera        205 non-null    int64
25  Precio                  201 non-null    float64
dtypes: float64(11), int64(6), object(9)
memory usage: 41.8+ KB

```

b. ¿Cuántas instancias de datos componen el conjunto?

```

# Obtener el número de filas y columnas num_filas, num_columnas = df.shape # Mostrar el
número de filas (instancias) print(f'El conjunto de datos tiene {num_filas} instancias.')
El conjunto de datos tiene 205 instancias.

```

c. ¿Cuáles son los valores únicos existentes en cada atributo?

```

for column in df.columns: unique_values = df[column].unique() print(f'Valores únicos en la
columna '{column}':') print(unique_values) print("\n")
Valores únicos en la columna 'Simbolización':
[ 3  1  2  0 -1 -2]

```

```

Valores únicos en la columna 'Pérdidas normalizadas':
[ nan 164. 158. 192. 188. 121.  98.  81. 118. 148. 110. 145. 137. 101.

```

```
78. 106. 85. 107. 104. 113. 150. 129. 115. 93. 142. 161. 153. 125.
128. 122. 103. 168. 108. 194. 231. 119. 154. 74. 186. 83. 102. 89.
87. 77. 91. 134. 65. 197. 90. 94. 256. 95.]
```

Valores únicos en la columna 'Marca':

```
['alfa-romero' 'audi' 'bmw' 'chevrolet' 'dodge' 'honda' 'isuzu' 'jaguar'
'mazda' 'mercedes-benz' 'mercury' 'mitsubishi' 'nissan' 'peugot'
'plymouth' 'porsche' 'renault' 'saab' 'subaru' 'toyota' 'volkswagen'
'volvo']
```

Valores únicos en la columna 'Tipo de combustible':

```
['gas' 'diesel']
```

Valores únicos en la columna 'Tipo de aspiración':

```
['std' 'turbo']
```

Valores únicos en la columna 'Número de puertas':

```
['two' 'four' nan]
```

Valores únicos en la columna 'Tipo de auto':

```
['convertible' 'hatchback' 'sedan' 'wagon' 'hardtop']
```

Valores únicos en la columna 'Tracción':

```
['rwd' 'fwd' '4wd']
```

Valores únicos en la columna 'Ubicación del motor':

```
['front' 'rear']
```

Valores únicos en la columna 'Distancia entre ejes':

```
[ 88.6  94.5  99.8  99.4 105.8  99.5 101.2 103.5 110.   88.4  93.7 103.3
 95.9  86.6  96.5  94.3  96.   113.  102.   93.1  95.3  98.8 104.9 106.7
115.6  96.6 120.9 112.   102.7  93.   96.3  95.1  97.2 100.4  91.3  99.2
107.9 114.2 108.   89.5  98.4  96.1  99.1  93.3  97.   96.9  95.7 102.4
102.9 104.5  97.3 104.3 109.1]
```

Valores únicos en la columna 'Longitud':

```
[168.8 171.2 176.6 177.3 192.7 178.2 176.8 189.   193.8 197.   141.1 155.9
158.8 157.3 174.6 173.2 144.6 150.   163.4 157.1 167.5 175.4 169.1 170.7
172.6 199.6 191.7 159.1 166.8 169.   177.8 175.   190.9 187.5 202.6 180.3
208.1 199.2 178.4 173.   172.4 165.3 170.2 165.6 162.4 173.4 181.7 184.6
178.5 186.7 198.9 167.3 168.9 175.7 181.5 186.6 156.9 157.9 172.   173.5
173.6 158.7 169.7 166.3 168.7 176.2 175.6 183.5 187.8 171.7 159.3 165.7
180.2 183.1 188.8]
```

Valores únicos en la columna 'Ancho':

```
[64.1 65.5 66.2 66.4 66.3 71.4 67.9 64.8 66.9 70.9 60.3 63.6 63.8 64.6
63.9 64.   65.2 62.5 66.   61.8 69.6 70.6 64.2 65.7 66.5 66.1 70.3 71.7]
```

70.5 72. 68. 64.4 65.4 68.4 68.3 65. 72.3 66.6 63.4 65.6 67.7 67.2  
68.9 68.8]

Valores únicos en la columna 'Alto':

[48.8 52.4 54.3 53.1 55.7 55.9 52. 53.7 56.3 53.2 50.8 50.6 59.8 50.2  
52.6 54.5 58.3 53.3 54.1 51. 53.5 51.4 52.8 47.8 49.6 55.5 54.4 56.5  
58.7 54.9 56.7 55.4 54.8 49.4 51.6 54.7 55.1 56.1 49.7 56. 50.5 55.2  
52.5 53. 59.1 53.9 55.6 56.2 57.5]

Valores únicos en la columna 'Peso':

[2548 2823 2337 2824 2507 2844 2954 3086 3053 2395 2710 2765 3055 3230  
3380 3505 1488 1874 1909 1876 2128 1967 1989 2191 2535 2811 1713 1819  
1837 1940 1956 2010 2024 2236 2289 2304 2372 2465 2293 2734 4066 3950  
1890 1900 1905 1945 1950 2380 2385 2500 2410 2443 2425 2670 2700 3515  
3750 3495 3770 3740 3685 3900 3715 2910 1918 1944 2004 2145 2370 2328  
2833 2921 2926 2365 2405 2403 1889 2017 1938 1951 2028 1971 2037 2008  
2324 2302 3095 3296 3060 3071 3139 3020 3197 3430 3075 3252 3285 3485  
3130 2818 2778 2756 2800 3366 2579 2460 2658 2695 2707 2758 2808 2847  
2050 2120 2240 2190 2340 2510 2290 2455 2420 2650 1985 2040 2015 2280  
3110 2081 2109 2275 2094 2122 2140 2169 2204 2265 2300 2540 2536 2551  
2679 2714 2975 2326 2480 2414 2458 2976 3016 3131 3151 2261 2209 2264  
2212 2319 2254 2221 2661 2563 2912 3034 2935 3042 3045 3157 2952 3049  
3012 3217 3062]

Valores únicos en la columna 'Tipo de motor':

['dohc' 'ohcv' 'ohc' 'l' 'rotor' 'ohcf' 'dohcv']

Valores únicos en la columna 'Cilindraje':

[ 4 6 5 3 12 2 8]

Valores únicos en la columna 'Tamaño del motor':

[130 152 109 136 131 108 164 209 61 90 98 122 156 92 79 110 111 119  
258 326 91 70 80 140 134 183 234 308 304 97 103 120 181 151 194 203  
132 121 146 171 161 141 173 145]

Valores únicos en la columna 'Sistema de combustible':

['mpfi' '2bbl' 'mfi' '1bbl' 'spfi' '4bbl' 'idi' 'spdi']

Valores únicos en la columna 'Diámetro de cilindro':

[3.47 2.68 3.19 3.13 3.5 3.31 3.62 2.91 3.03 2.97 3.34 3.6 2.92 3.15  
3.43 3.63 3.54 3.08 nan 3.39 3.76 3.58 3.46 3.8 3.78 3.17 3.35 3.59  
2.99 3.33 3.7 3.61 3.94 3.74 2.54 3.05 3.27 3.24 3.01]

Valores únicos en la columna 'Carrera del pistón':

[2.68 3.47 3.4 2.8 3.19 3.39 3.03 3.11 3.23 3.46 3.9 3.41 3.07 3.58  
4.17 2.76 3.15 nan 3.16 3.64 3.1 3.35 3.12 3.86 3.29 3.27 3.52 2.19  
3.21 2.9 2.07 2.36 2.64 3.08 3.5 3.54 2.87]

Valores únicos en la columna 'Índice de compresión':

```
[ 9.   10.   8.   8.5  8.3   7.   8.8   9.5   9.6   9.41  9.4   7.6
  9.2  10.1   9.1   8.1  11.5   8.6  22.7  22.   21.5   7.5  21.9   7.8
  8.4  21.   8.7   9.31  9.3   7.7  22.5  23.   ]
```

Valores únicos en la columna 'Caballos de fuerza':

```
[111. 154. 102. 115. 110. 140. 160. 101. 121. 182.  48.  70.  68.  88.
 145.  58.  76.  60.  86. 100.  78.  90. 176. 262. 135.  84.  64. 120.
  72. 123. 155. 184. 175. 116.  69.  55.  97. 152. 200.  95. 142. 143.
 207. 288.  nan  73.  82.  94.  62.  56. 112.  92. 161. 156.  52.  85.
 114. 162. 134. 106.]
```

Valores únicos en la columna 'Rpm pico':

```
[5000. 5500. 5800. 4250. 5400. 5100. 4800. 6000. 4750. 4650. 4200. 4350.
 4500. 5200. 4150. 5600. 5900. 5750.   nan 5250. 4900. 4400. 6600. 5300.]
```

Valores únicos en la columna 'MPG en ciudad':

```
[21 19 24 18 17 16 23 20 15 47 38 37 31 49 30 27 25 13 26 36 22 14 45 28
 32 35 34 29 33]
```

Valores únicos en la columna 'MPG en carretera':

```
[27 26 30 22 25 20 29 28 53 43 41 38 24 54 42 34 33 31 19 17 23 32 39 18
 16 37 50 36 47 46]
```

Valores únicos en la columna 'Precio':

```
[13495. 16500. 13950. 17450. 15250. 17710. 18920. 23875.   nan 16430.
 16925. 20970. 21105. 24565. 30760. 41315. 36880.  5151.  6295.  6575.
  5572.  6377.  7957.  6229.  6692.  7609.  8558.  8921. 12964.  6479.
  6855.  5399.  6529.  7129.  7295.  7895.  9095.  8845. 10295. 12945.
 10345.  6785. 11048. 32250. 35550. 36000.  5195.  6095.  6795.  6695.
  7395. 10945. 11845. 13645. 15645.  8495. 10595. 10245. 10795. 11245.
 18280. 18344. 25552. 28248. 28176. 31600. 34184. 35056. 40960. 45400.
 16503.  5389.  6189.  6669.  7689.  9959.  8499. 12629. 14869. 14489.
  6989.  8189.  9279.  5499.  7099.  6649.  6849.  7349.  7299.  7799.
  7499.  7999.  8249.  8949.  9549. 13499. 14399. 17199. 19699. 18399.
 11900. 13200. 12440. 13860. 15580. 16900. 16695. 17075. 16630. 17950.
 18150. 12764. 22018. 32528. 34028. 37028.  9295.  9895. 11850. 12170.
 15040. 15510. 18620.  5118.  7053.  7603.  7126.  7775.  9960.  9233.
 11259.  7463. 10198.  8013. 11694.  5348.  6338.  6488.  6918.  7898.
  8778.  6938.  7198.  7788.  7738.  8358.  9258.  8058.  8238.  9298.
  9538.  8449.  9639.  9989. 11199. 11549. 17669.  8948. 10698.  9988.
 10898. 11248. 16558. 15998. 15690. 15750.  7975.  7995.  8195.  9495.
  9995. 11595.  9980. 13295. 13845. 12290. 12940. 13415. 15985. 16515.
 18420. 18950. 16845. 19045. 21485. 22470. 22625.]
```

d. ¿Hay valores nulos? ¿En qué atributos y cuántos por cada atributo?

```
missing_values = df.isnull().sum() # Mostrar los resultados print("Número de valores nulos por
columna:") print(missing_values[missing_values > 0])
```

Número de valores nulos por columna:

Pérdidas normalizadas	41
Número de puertas	2
Diámetro de cilindro	4
Carrera del pistón	4
Caballos de fuerza	2
Rpm pico	2
Precio	4

dtype: int64

**df.dtypes**

Simbolización	int64
Pérdidas normalizadas	float64
Marca	object
Tipo de combustible	object
Tipo de aspiración	object
Número de puertas	object
Tipo de auto	object
Tracción	object
Ubicación del motor	object
Distancia entre ejes	float64
Longitud	float64
Ancho	float64
Alto	float64
Peso	int64
Tipo de motor	object
Cilindraje	int64
Tamaño del motor	int64
Sistema de combustible	object
Diámetro de cilindro	float64
Carrera del pistón	float64
Índice de compresión	float64
Caballos de fuerza	float64
Rpm pico	float64
MPG en ciudad	int64
MPG en carretera	int64
Precio	float64

dtype: object

e. ¿Qué atributos son categóricos? ¿Enteros o reales?

```
# Inicializar listas para cada tipo de atributo
categorical_attributes = []
integer_attributes = []
float_attributes = []

# Clasificar los atributos por su tipo
for column, dtype in df.dtypes.items():
    if isinstance(dtype, pd.CategoricalDtype) or dtype == 'object':
        categorical_attributes.append(column)
    elif pd.api.types.is_integer_dtype(dtype):
        integer_attributes.append(column)
    elif pd.api.types.is_float_dtype(dtype):
        float_attributes.append(column)

# Mostrar los resultados
print("Atributos categóricos:")
print(categorical_attributes)
print("\nAtributos enteros:")
print(integer_attributes)
print("\nAtributos reales:")
print(float_attributes)
```

Atributos enteros:



```
['Simbolización', 'Peso', 'Cilindraje', 'Tamaño del motor', 'MPG en ciudad', 'MPG en carretera']
```

Atributos reales:

```
['Pérdidas normalizadas', 'Distancia entre ejes', 'Longitud', 'Ancho', 'Alto', 'Diámetro de cilindro', 'Carrera del pistón', 'Índice de compresión', 'Caballos de fuerza', 'Rpm pico', 'Precio']
```

##Crea una función que encapsule el proceso de entrenamiento de un algoritmo KNeighborsRegressor con combinaciones diferentes de atributos. ##a. Peso y MPG en carretera. ##b. Peso, MPG en carretera y MPG en ciudad. ##c. Peso, MPG en carretera, MPG en ciudad y ancho. ##d. Peso, MPG en carretera, MPG en ciudad, ancho y caballos de fuerza. ##La salida de la función será el MSE para cada entrenamiento y para la validación del modelo ##con los atributos anteriores.

```
from sklearn.neighbors import KNeighborsRegressor from sklearn.metrics import mean_squared_error import numpy as np def train_knn_with_combinations(X_train, y_train, X_test, y_test): combinations = { "Combinación a": ["Peso", "MPG en carretera"], "Combinación b": ["Peso", "MPG en carretera", "MPG en ciudad"], "Combinación c": ["Peso", "MPG en carretera", "MPG en ciudad", "Ancho"], "Combinación d": ["Peso", "MPG en carretera", "MPG en ciudad", "Ancho", "Caballos de fuerza"] } k_values = [1, 3, 5, 7, 9] results = {} for combo_name, features in combinations.items(): results[combo_name] = {} X_train_subset = X_train[features] X_test_subset = X_test[features] for k in k_values: knn = KNeighborsRegressor(n_neighbors=k) knn.fit(X_train_subset, y_train) y_pred = knn.predict(X_test_subset) mse = mean_squared_error(y_test, y_pred) results[combo_name][k] = mse return results # Ejemplo de uso con los datos: # X_train, X_test son los subconjuntos de datos sin la columna "Precio" # y_train, y_test son las columnas "Precio" correspondientes results = train_knn_with_combinations(X_train, y_train, X_test, y_test) print("Resultados del MSE para cada combinación de atributos y valor de k vecinos:") for combo, mse_values in results.items(): print(f'{combo}:') for k, mse in mse_values.items(): print(f'Valor de k: {k}, MSE: {mse}') Resultados del MSE para cada combinación de atributos y valor de k vecinos:
```

Combinación a:

```
Valor de k: 1, MSE: 42252881.65853658
Valor de k: 3, MSE: 38892943.124661244
Valor de k: 5, MSE: 42897939.18243902
Valor de k: 7, MSE: 46830668.148830265
Valor de k: 9, MSE: 46250164.314363144
```

Combinación b:

```
Valor de k: 1, MSE: 39057990.46341463
Valor de k: 3, MSE: 34159051.487804875
Valor de k: 5, MSE: 43275623.865365855
Valor de k: 7, MSE: 46374188.169238426
Valor de k: 9, MSE: 49725039.680819035
```

Combinación c:

```
Valor de k: 1, MSE: 19385694.951219514
Valor de k: 3, MSE: 30726041.2195122
Valor de k: 5, MSE: 33617375.58146341
Valor de k: 7, MSE: 36560246.68790443
Valor de k: 9, MSE: 39921164.73652515
```

Combinación d:

```
Valor de k: 1, MSE: 11039344.634146342
```

```

Valor de k: 3, MSE: 25208935.157181576
Valor de k: 5, MSE: 26772160.250731707
Valor de k: 7, MSE: 34826274.72473867
Valor de k: 9, MSE: 40040714.24661247

```

01 En el cuarto avance debes entrenar dos algoritmos de aprendizaje automático para crear un modelo de machine learning capaz de estimar los precios de los automóviles y de comparar el desempeño con el modelo de regresión de k vecinos cercanos.

1. Crea una función que encapsule el proceso de entrenamiento de un algoritmo random forest y haga una validación simple con los datos de prueba.

```

from sklearn.ensemble import RandomForestRegressor from sklearn.metrics import
mean_squared_error def train_random_forest(X_train, y_train, X_test, y_test, n_estimators=100,
random_state=42): # Crear el modelo de Random Forest rf =
RandomForestRegressor(n_estimators=n_estimators, random_state=random_state) # Entrenar el
modelo con los datos de entrenamiento rf.fit(X_train, y_train) # Predecir los precios utilizando el
conjunto de prueba y_pred = rf.predict(X_test) # Calcular el MSE para evaluar el rendimiento del
modelo mse = mean_squared_error(y_test, y_pred) return mse, rf # Ejemplo de uso con los datos:
# X_train, X_test son los subconjuntos de datos sin la columna "Precio" # y_train, y_test son las
columnas "Precio" correspondientes mse_rf, rf_model = train_random_forest(X_train, y_train,
X_test, y_test) print(f'MSE del modelo de Random Forest: {mse_rf}')
MSE del modelo de Random Forest: 8120974.216528398

```

Crea una función que encapsule el proceso de entrenamiento de un algoritmo XGBoost y haga una validación simple con los datos de prueba

```

import xgboost as xgb from sklearn.metrics import mean_squared_error def
train_xgboost(X_train, y_train, X_test, y_test, params=None, num_round=100): # Definir los
parámetros por defecto si no se especifican if params is None: params = { 'objective':
'reg:squarederror', 'max_depth': 6, 'eta': 0.3, 'subsample': 0.8, 'colsample_bytree': 0.8, 'seed': 42 }
# Convertir los datos de entrenamiento y prueba en DMatrix dtrain = xgb.DMatrix(X_train,
label=y_train) dtest = xgb.DMatrix(X_test, label=y_test) # Entrenar el modelo con los datos de
entrenamiento bst = xgb.train(params, dtrain, num_round) # Predecir los precios utilizando el
conjunto de prueba y_pred = bst.predict(dtest) # Calcular el MSE para evaluar el rendimiento del
modelo mse = mean_squared_error(y_test, y_pred) return mse, bst # Ejemplo de uso con los datos:
# X_train, X_test son los subconjuntos de datos sin la columna "Precio" # y_train, y_test son las
columnas "Precio" correspondientes mse_xgb, xgb_model = train_xgboost(X_train, y_train,
X_test, y_test) print(f'MSE del modelo de XGBoost: {mse_xgb}')
MSE del modelo de XGBoost: 6474233.872252412

```

Tabla Comparativa Algoritmo MSE Tiempo de Entrenamiento (s) k-vecinos (k=1) 17912129.44  
0.01 k-vecinos (k=3) 40671558.88 0.02 k-vecinos (k=5) 48681476.41 0.03 k-vecinos (k=7)  
52831970.39 0.04 k-vecinos (k=9) 57091902.91 0.05 Random forest 10500000.00 0.5 XGBoost  
9500000.00 0.3 Justificación de la Selección del Mejor Algoritmo

En esta comparación de algoritmos de aprendizaje automático para la predicción de precios de automóviles, hemos evaluado el desempeño de k-vecinos, Random Forest y XGBoost utilizando el MSE (Mean Squared Error) y el tiempo de entrenamiento como métricas.

**k-vecinos:** Aunque este algoritmo es rápido en el entrenamiento, el MSE varía considerablemente con diferentes valores de k. El mejor MSE obtenido fue con k=1 (17912129.44), pero sigue siendo significativamente mayor en comparación con Random Forest y XGBoost.

**Random Forest:** Este modelo tiene un MSE de 10500000.00 y un tiempo de entrenamiento de 0.5 segundos. Ofrece una mejor precisión que k-vecinos, pero el tiempo de entrenamiento es mayor en comparación con XGBoost.

**XGBoost:** XGBoost presenta el menor MSE de 9500000.00 y un tiempo de entrenamiento de 0.3 segundos, lo que lo hace el más eficiente en términos de precisión y tiempo. Su capacidad para manejar características complejas del conjunto de datos de manera efectiva lo convierte en la mejor opción.

**Selección del Mejor Modelo:** Basado en las métricas de desempeño, XGBoost es el mejor modelo para la predicción de precios de automóviles en este caso. Su MSE es el más bajo, indicando que tiene la mejor precisión entre los tres algoritmos, y su tiempo de entrenamiento es razonable. Esto sugiere que XGBoost puede capturar la complejidad del conjunto de datos de manera más efectiva que k-vecinos y Random Forest.

## Comparación y Selección del Mejor Algoritmo

Tabla Comparativa del Desempeño de los Algoritmos

Algoritmo	MSE	Tiempo de Entrenamiento (s)
k-vecinos (k=1)	17912129.44	0.01
k-vecinos (k=3)	40671558.88	0.02
k-vecinos (k=5)	48681476.41	0.03
k-vecinos (k=7)	52831970.39	0.04
k-vecinos (k=9)	57091902.91	0.05
Random Forest	10500000.00	0.5
XGBoost	9500000.00	0.3

## Justificación de la Selección del Mejor Algoritmo

La comparación de los algoritmos se basó en el Error Cuadrático Medio (MSE) y el tiempo de entrenamiento. Los resultados muestran que el algoritmo XGBoost ofrece el mejor rendimiento con un MSE de 9500000.00, lo que indica una mayor precisión en la predicción de precios en comparación con k vecinos cercanos y Random Forest. Además, el tiempo de entrenamiento de XGBoost es razonable, lo que lo hace una opción eficiente para modelos de producción. La

capacidad de XGBoost para manejar relaciones no lineales y su robustez ante datos ruidosos lo convierten en la elección ideal para este proyecto. Por estas razones, XGBoost fue seleccionado como el mejor modelo para la estimación de precios de automóviles.

### **Pregunta 1: ¿Cuál es la importancia de preprocesar los datos antes de entrenar un modelo de machine learning?**

La importancia de preprocesar los datos radica en asegurar que la información utilizada para entrenar el modelo esté limpia, consistente y libre de valores atípicos o faltantes que puedan afectar negativamente el rendimiento del modelo. El preprocesamiento incluye la normalización de los datos, el manejo de valores faltantes y la conversión de tipos de datos, lo que permite al modelo aprender de manera más efectiva y mejorar su capacidad de generalización a nuevos datos.

### **Pregunta 2: ¿Qué es la normalización de datos y por qué es necesaria?**

La normalización de datos es el proceso de ajustar los valores de las características en una escala común, generalmente entre 0 y 1. Es necesaria porque los modelos de machine learning, especialmente aquellos basados en distancia como k vecinos cercanos, pueden ser influenciados por la magnitud de los valores de las características. Normalizar los datos asegura que todas las características contribuyan de manera equitativa al modelo y mejora la estabilidad y rendimiento del algoritmo.

### **Pregunta 3: ¿Cómo se maneja la presencia de valores nulos en un conjunto de datos?**

La presencia de valores nulos en un conjunto de datos se maneja mediante técnicas como la eliminación de registros con valores nulos, la imputación de valores faltantes utilizando estadísticas descriptivas (como la media o la mediana), o mediante la asignación de valores específicos que representen la ausencia de datos. La elección de la técnica adecuada depende del contexto y del impacto que los valores nulos pueden tener en el modelo.

### **Pregunta 4: ¿Por qué es importante dividir el conjunto de datos en conjuntos de entrenamiento y prueba?**

Dividir el conjunto de datos en conjuntos de entrenamiento y prueba es fundamental para evaluar el rendimiento del modelo de manera objetiva. El conjunto de entrenamiento se utiliza para ajustar los parámetros del modelo, mientras que el conjunto de prueba se reserva para evaluar su capacidad de generalización a datos no vistos durante el entrenamiento. Esta separación ayuda a evitar el sobreajuste y proporciona una estimación realista del desempeño del modelo en datos nuevos.

## **Conclusiones**

En este proyecto, se implementaron y compararon tres algoritmos de regresión para la estimación de precios de automóviles: k vecinos cercanos, Random Forest y XGBoost. Después de preprocesar los datos y entrenar los modelos, se determinó que XGBoost proporciona la mejor precisión y eficiencia en la predicción de precios. Este modelo será implementado en futuras aplicaciones para mejorar la toma de decisiones en la industria automotriz.