

```

1 !pip install pandas
2 !pip install numpy
3 !pip install matplotlib

```

Requirement already satisfied: pandas in /usr/local/lib/python3.11/dist-packages (2.2.2)
 Requirement already satisfied: numpy>=1.23.2 in /usr/local/lib/python3.11/dist-packages (from pandas) (2.0.2)
 Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.11/dist-packages (from pandas) (2.8.2)
 Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.11/dist-packages (from pandas) (2025.1)
 Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.11/dist-packages (from pandas) (2025.1)
 Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.11/dist-packages (from python-dateutil>=2.8.2->pandas) (1.16.0)
 Requirement already satisfied: matplotlib in /usr/local/lib/python3.11/dist-packages (3.10.0)
 Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (1.3.0)
 Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (0.12.1)
 Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (4.55.3)
 Requirement already satisfied: kiwisolver>=1.3.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (1.4.7)
 Requirement already satisfied: numpy>=1.23 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (2.0.2)
 Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (24.2)
 Requirement already satisfied: pillow>=8 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (11.1.0)
 Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (3.2.1)
 Requirement already satisfied: python-dateutil>=2.7 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (2.8.2)
 Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.11/dist-packages (from python-dateutil>=2.7->matplotlib) (1.16.0)

```

1 from google.colab import drive
2
3
4 drive.mount('/content/drive')
5
6
7 file_path = '/content/drive/My Drive/iaa/precos_carros_brasil.csv'
8
9

```

Mounted at /content/drive

```

1 import pandas as pd
2 import numpy as np
3 #a. Carregue a base de dados media_precos_carros_brasil.csv
4 def extractData(dataBaseFile):
5     db = pd.read_csv(dataBaseFile)
6     return db
7
8 def treatData(db):
9     # d Crie duas categorias, para separar colunas numéricas e categóricas. Imprima o resumo
10    # de informações das variáveis numéricas e categóricas (estatística descritiva dos dados)
11
12    colunas_numericas = db.select_dtypes(include=[np.number]).columns
13    colunas_categoricas = db.select_dtypes(exclude=[np.number]).columns
14
15    # d.1 -> colunas numéricas
16    print("Resumo estatístico das colunas numéricas:")
17    print(db[colunas_numericas].describe()) # Estatísticas descritivas das variáveis numéricas
18
19    # d.2 -> categóricas
20    print("\nResumo das colunas categóricas:")
21    print(db[colunas_categoricas].describe()) # Resumo das variáveis categóricas
22
23    return db
24
25 def identifyNan(value):
26     value = np.nan
27     return value
28
29 def patternColumn(db):
30     db = pd.DataFrame(db)
31     colunas = db.columns
32     print(colunas)
33
34
35     colunas_lista = list(colunas)
36
37    # b. Verifique se há valores faltantes nos dados. Caso haja, escolha uma tratativa para resolver o problema de
38    # b.1 -> se a coluna ou linha inteira for composta por valor
39    db_cleaned = db.dropna(axis=1, how='all')
40    db_cleaned = db.dropna(axis=0, how='all')
41
42
43    # b.2 -> Atribua um novo valor aos dados Nan presentes em colunas e linhas nas regras de acordo com a linha de

```

```
43     # D.T. -> Atribui um novo valor aos dados nan presente em colunas e linhas nao nulas de acordo com o tipo de c
44     for coluna_nome in colunas_lista:
45         for index, value in db[coluna_nome].items():
46             if pd.isna(value):
47                 if coluna_nome in ("year_of_reference", "year_model", "avg_price_brl"):
48                     value = 0
49             else:
50                 value = "Não possui"
51
52     return db_cleaned
53
54
```

Double-click (or enter) to edit

```
1 db = extractData(file_path)
2 #c. Verifique se há dados duplicados nos dados
3 duplicates = db.duplicated().sum()
4 print(f"Número de linhas duplicadas: {duplicates}")
5 db = patternColum(db)
6 treatData(db)
7
```

<ipython-input-3-f34971953ffd>:5: DtypeWarning: Columns (1,2,3,4,5,6,7,8) have mixed types. Specify dtype option c

db = pd.read_csv(dataBaseFile)

Número de linhas duplicadas: 65246

Index(['year_of_reference', 'month_of_reference', 'fipecode', 'authentication', 'brand', 'model', 'fuel', 'gear', 'engine_size', 'year_model', 'avg_price_brl'], dtype='object')

Resumo estatístico das colunas numéricas:

	year_of_reference	year_model	avg_price_brl
count	202297.000000	202297.000000	202297.000000
mean	2021.564694	2011.271527	52756.909153
std	0.571903	6.376234	51628.677716
min	2021.000000	2000.000000	6647.000000
25%	2021.000000	2006.000000	22855.000000
50%	2022.000000	2012.000000	38027.000000
75%	2022.000000	2016.000000	64064.000000
max	2023.000000	2023.000000	979358.000000

Resumo das colunas categóricas:

	month_of_reference	fipecode	authentication	brand \
count	202297	202297	202297	202297
unique	12	2091	202295	6
top	January	001216-5	3r6c277cnqcb	Fiat
freq	24260	425	2	44962

	model	fuel	gear	engine_size
count	202297	202297	202297	202297
unique	2112	3	2	29
top	Palio Week. Adv/Adv TRYON 1.8 mpi Flex	Gasoline	manual	1,6
freq	425	168685	161885	47420

	year_of_reference	month_of_reference	fipecode	authentication	brand	model	fuel	gear	eng
0	2021.0	January	004001-0	cfzlcztzfwrcp	GM - Chevrolet	Corsa Wind 1.0 MPFI / EFI 2p	Gasoline	manual	
1	2021.0	January	004001-0	cdqwxwpw3y2p	GM - Chevrolet	Corsa Wind 1.0 MPFI / EFI 2p	Gasoline	manual	
2	2021.0	January	004001-0	cb1t3xwwj1xp	GM - Chevrolet	Corsa Wind 1.0 MPFI / EFI 2p	Gasoline	manual	
3	2021.0	January	004001-0	cb9gct6j65r0	GM - Chevrolet	Corsa Wind 1.0 MPFI / EFI 2p	Alcohol	manual	
4	2021.0	January	004003-7	g15wg0gbz1fx	GM - Chevrolet	Corsa Pick-Up GL/ Champ	Gasoline	manual	

...							1.6 MPFI / EFI	...	
...							Saveiro Robust 1.6 Total Flex 16V	Gasoline	manual
202292	2023.0	January	005538-7	ccv3mvxn	sz0dqw	VW - VolksWagen			
...							Gol Last Edition 1.0 Flex 12V 5p	Gasoline	manual
202293	2023.0	January	005539-5	chmwfg3l	5hbp	VW - VolksWagen			
...							Gol Last Edition 1.0 Flex 12V 5p	Gasoline	manual
202294	2023.0	January	005539-5	cdj27srtc	vcddqw	VW - VolksWagen			
...							Polo Track 1.0 Flex 12V 5p	Gasoline	manual
202295	2023.0	January	005540-9	9w64fg6	dhhqp	VW - VolksWagen			
...							Polo Track 1.0 Flex 12V 5p	Gasoline	manual
202296	2023.0	January	005540-9	7hbnjmi	j9z5dqw	VW - VolksWagen			

202297 rows × 11 columns

```
1 #c imprima a contagem de valores por modelo (model) e marca do carro (brand)
2
3 def contar_modelo_e_marca(db):
4
5     if 'model' in db.columns and 'brand' in db.columns:
6
7         print("\nContagem de valores por modelo (model):")
8         print(db['model'].value_counts())
9         print("\nContagem de valores por marca (brand):")
10        print(db['brand'].value_counts())
11
12    else:
13        print("As colunas 'model' e/ou 'brand' não estão presentes no DataFrame.")
14
```

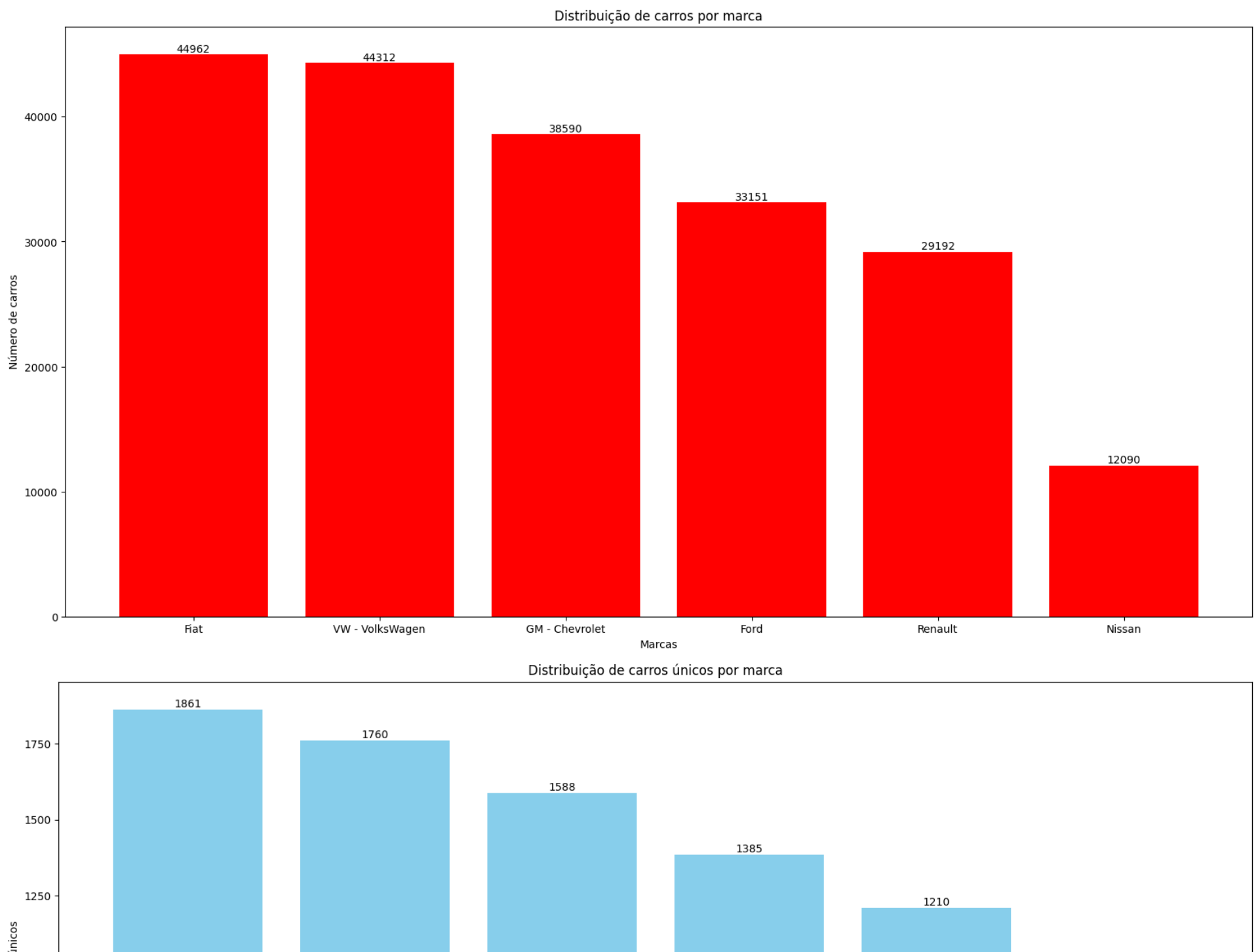
```
1 contar_modelo_e_marca(db)
```

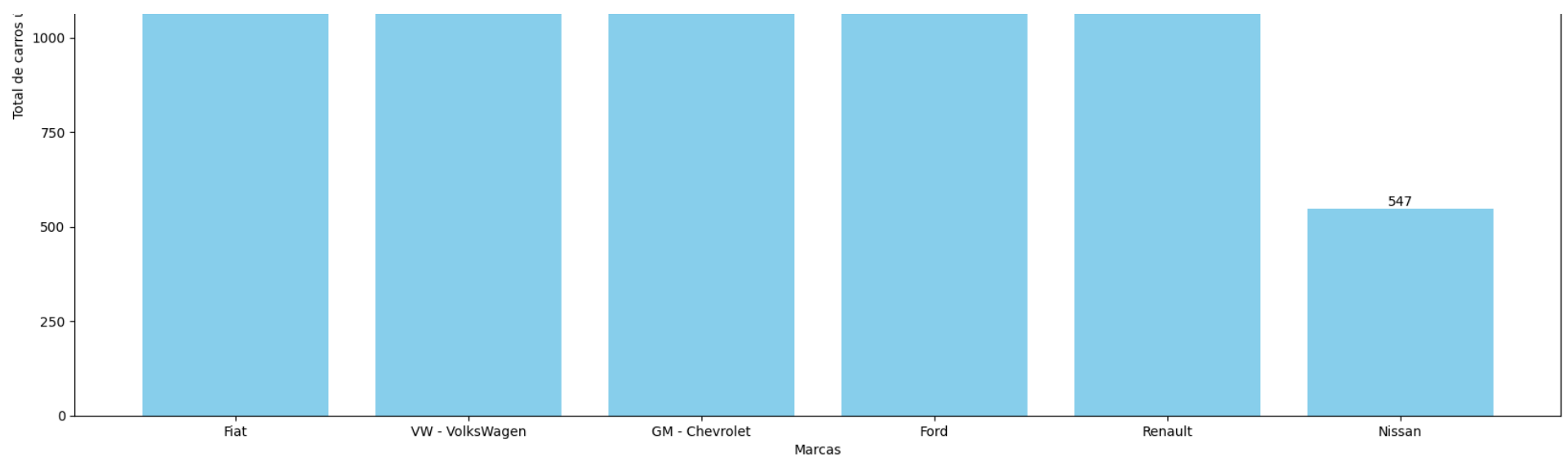
```
Contagem de valores por modelo (model):
model
Palio Week. Adv/Adv TRYON 1.8 mpi Flex    425
Focus 1.6 S/SE/SE Plus Flex 8V/16V 5p    425
Focus 2.0 16V/SE/SE Plus Flex 5p Aut.    400
Saveiro 1.6 Mi/ 1.6 Mi Total Flex 8V     400
Doblo Adv/Adv TRYON/LOCKER 1.8 Flex      375
...
STEPWAY Zen Flex 1.0 12V Mec.             2
Polo Track 1.0 Flex 12V 5p                2
Saveiro Robust 1.6 Total Flex 16V         2
KICKS Active 1.6 16V Flex Aut.            2
PULSE ABARTH 1.3 Turbo 16V Flex Aut.      2
Name: count, Length: 2112, dtype: int64
```

```
Contagem de valores por marca (brand):
brand
Fiat            44962
VW - Volkswagen 44312
GM - Chevrolet  38590
Ford            33151
Renault         29192
Nissan           12090
Name: count, dtype: int64
```

```
1 # a. Gere um gráfico da distribuição da quantidade de carros por marca
2
3 import matplotlib.pyplot as plt
4
5 def plot_car_count_by_brand(db):
6
7     plt.figure(figsize=(20, 10))
8     brand_counts = db['brand'].value_counts()
9     graphic = plt.bar(brand_counts.index, brand_counts.values, color='red')
10
11     plt.bar_label(graphic, fontsize=10)
12     plt.title('Distribuição de carros por marca')
13     plt.xlabel('Marcas')
14     plt.ylabel('Número de carros')
15     plt.xticks(rotation=0)
16
17     plt.show()
18
19
20 def plot_car_unique_count_by_brand(db):
21
22     plt.figure(figsize=(20, 10))
23     unique_db = db.drop_duplicates(subset=['brand', 'model', 'year_model', 'engine_size'])
24     unique_values_brand = unique_db['brand'].value_counts()
25     graphic = plt.bar(unique_values_brand.index, unique_values_brand.values, color='skyblue')
26
27     plt.bar_label(graphic, fontsize=10)
28     plt.title('Distribuição de carros únicos por marca')
29     plt.xlabel('Marcas')
30     plt.ylabel('Total de carros únicos')
31     plt.xticks(rotation=0)
32
33     plt.show()
34
```

```
1 plot_car_count_by_brand(db)
2 plot_car_unique_count_by_brand(db)
```





```

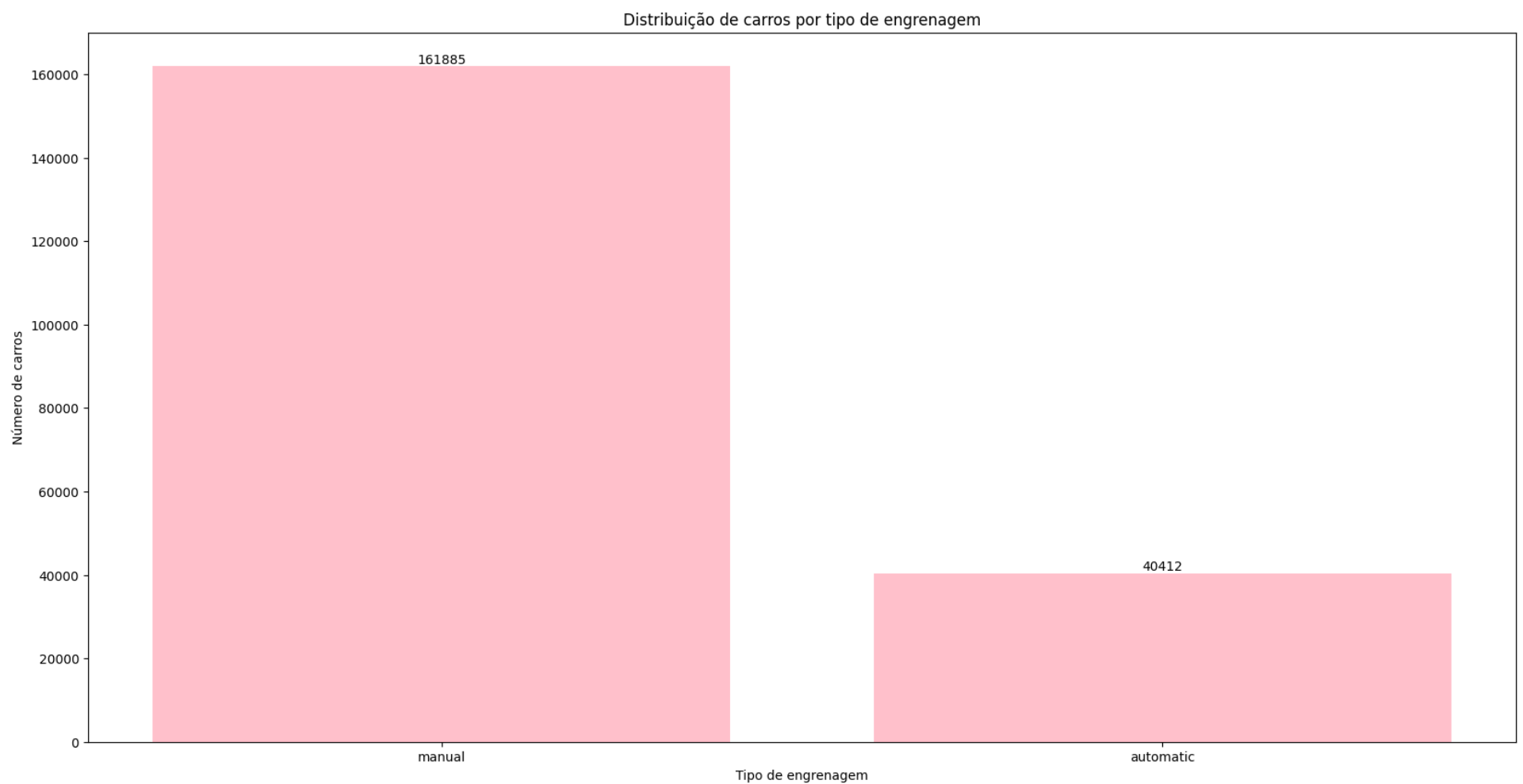
1 # b. Plot the distribution of the number of cars by type of gear
2 def plot_car_count_by_gear_type(db):
3
4     plt.figure(figsize=(20, 10))
5     gear_count = db['gear'].value_counts()
6     graphic = plt.bar(gear_count.index, gear_count.values, color='pink')
7
8     plt.bar_label(graphic, fontsize=10)
9     plt.title('Distribuição de carros por tipo de engrenagem')
10    plt.xlabel('Tipo de engrenagem')
11    plt.ylabel('Número de carros')
12    plt.xticks(rotation=0)
13
14    plt.show()
15

```

```

1 plot_car_count_by_gear_type(db)

```



```

1 #c. Gere um gráfico da evolução da média de preço dos carros ao longo dos meses de 2022 (variável de tempo no eixo
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6

```

```

7 def plot_avg_price_over_time(db):
8     db['year_of_reference'] = db['year_of_reference'].astype(str).str.strip()
9     db['year_of_reference'] = db['year_of_reference'].str[-4:]
10    db['year_of_reference'] = pd.to_numeric(db['year_of_reference'], errors='coerce')
11    db = db.dropna(subset=['year_of_reference'])
12    db['year_of_reference'] = db['year_of_reference'].astype(int)
13
14
15    db_2022 = db[db['year_of_reference'] == 22]
16
17    if db_2022.empty:
18        print("Nenhum dado encontrado para o ano de 2022.")
19        return
20
21    avg_price = db_2022.groupby('month_of_reference')['avg_price_brl'].mean().round(2).reset_index()
22
23    month_order = ['January', 'February', 'March', 'April', 'May', 'June', 'July',
24                  'August', 'September', 'October', 'November', 'December']
25
26    plt.figure(figsize=(20, 10))
27    sns.barplot(x='month_of_reference', y='avg_price_brl', data=avg_price, order=month_order)
28    plt.xticks(rotation=45)
29    plt.title('Preço Médio por Mês em 2022')
30    plt.xlabel('Mês de Referência')
31    plt.ylabel('Preço Médio (BRL)')
32    plt.show()
33
34

```

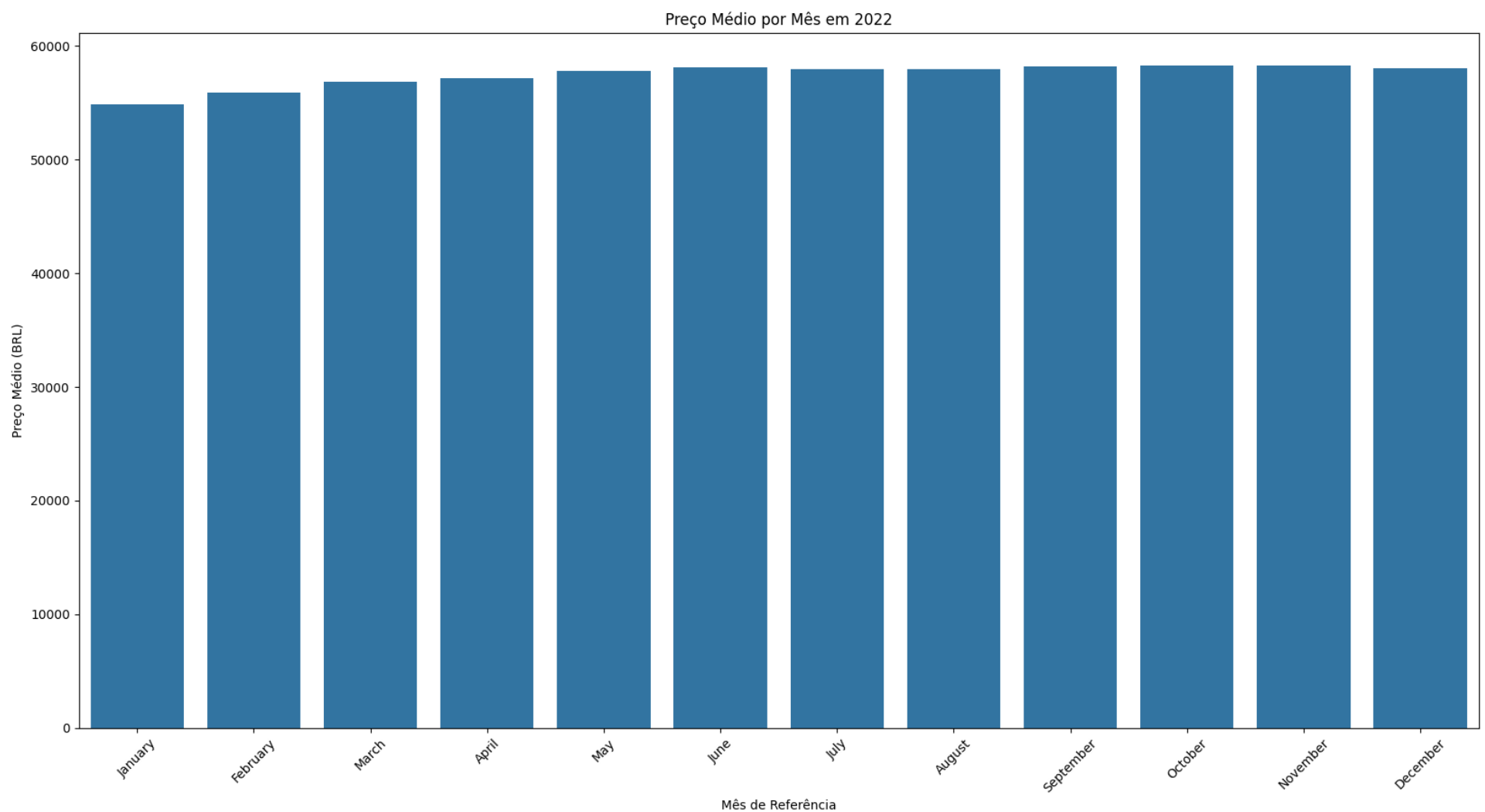
Double-click (or enter) to edit

```

1 print(np.dtype(db['year_of_reference']))
2
3
4 plot_avg_price_over_time(db)

```

float64



```

1 #d. Gere um gráfico da distribuição da média de preço dos carros por marca e tipo de engrenagem
2 def plot_avg_price_by_brand_and_gear(db):
3     avg_price_by_brand_gear = db.groupby(['brand', 'gear'])['avg_price_brl'].mean().unstack()
4
5     plt.figure(figsize=(20, 10))
6     avg_price_by_brand_gear.plot(kind='bar', stacked=False)
7     plt.title('Média do preço do carro por engrenagem e marca')
8     plt.xlabel('Brand')

```

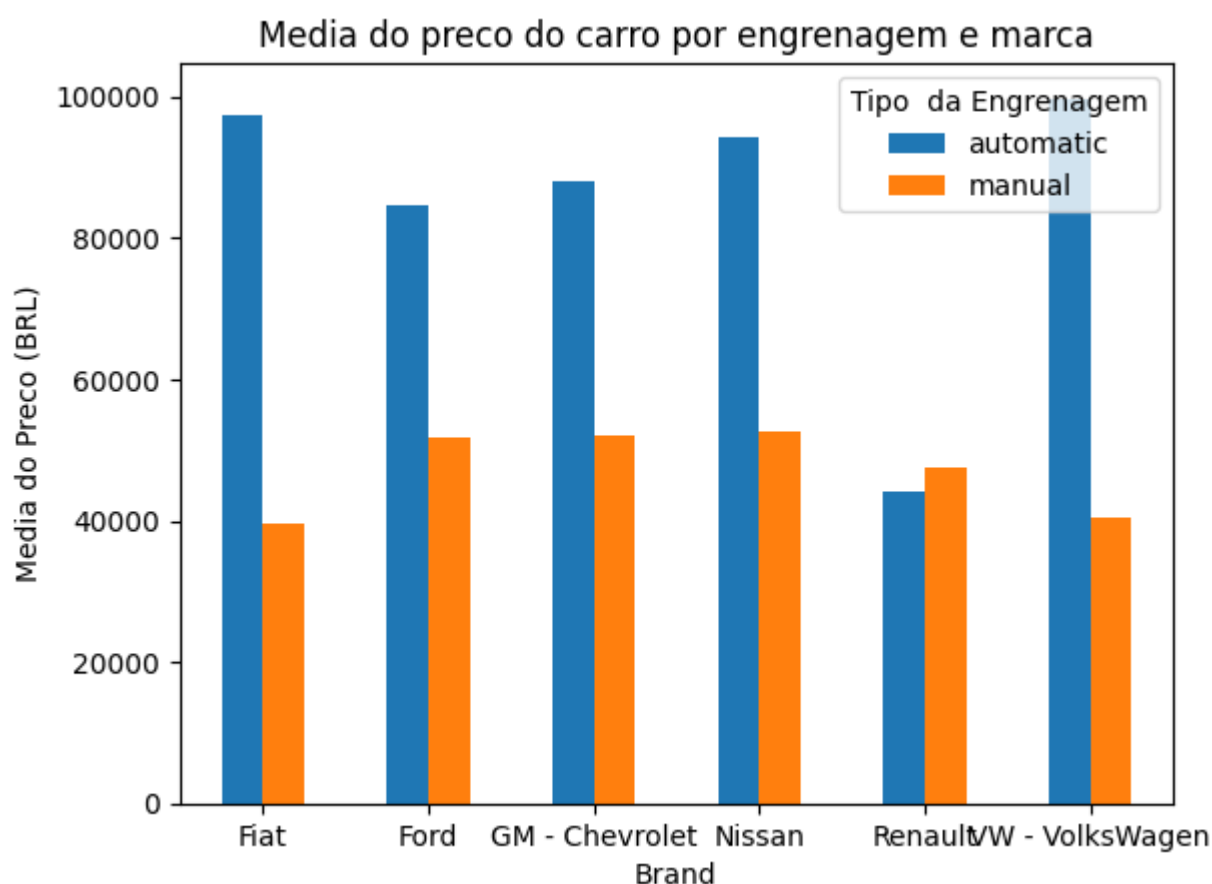
```

8 plt.xlabel('Brand')
9 plt.ylabel('Media do Preco (BRL)')
10 plt.xticks(rotation=0)
11 plt.legend(title='Tipo da Engrenagem')
12 plt.show()
13

```

```
1 plot_avg_price_by_brand_and_gear(db)
```

<Figure size 2000x1000 with 0 Axes>



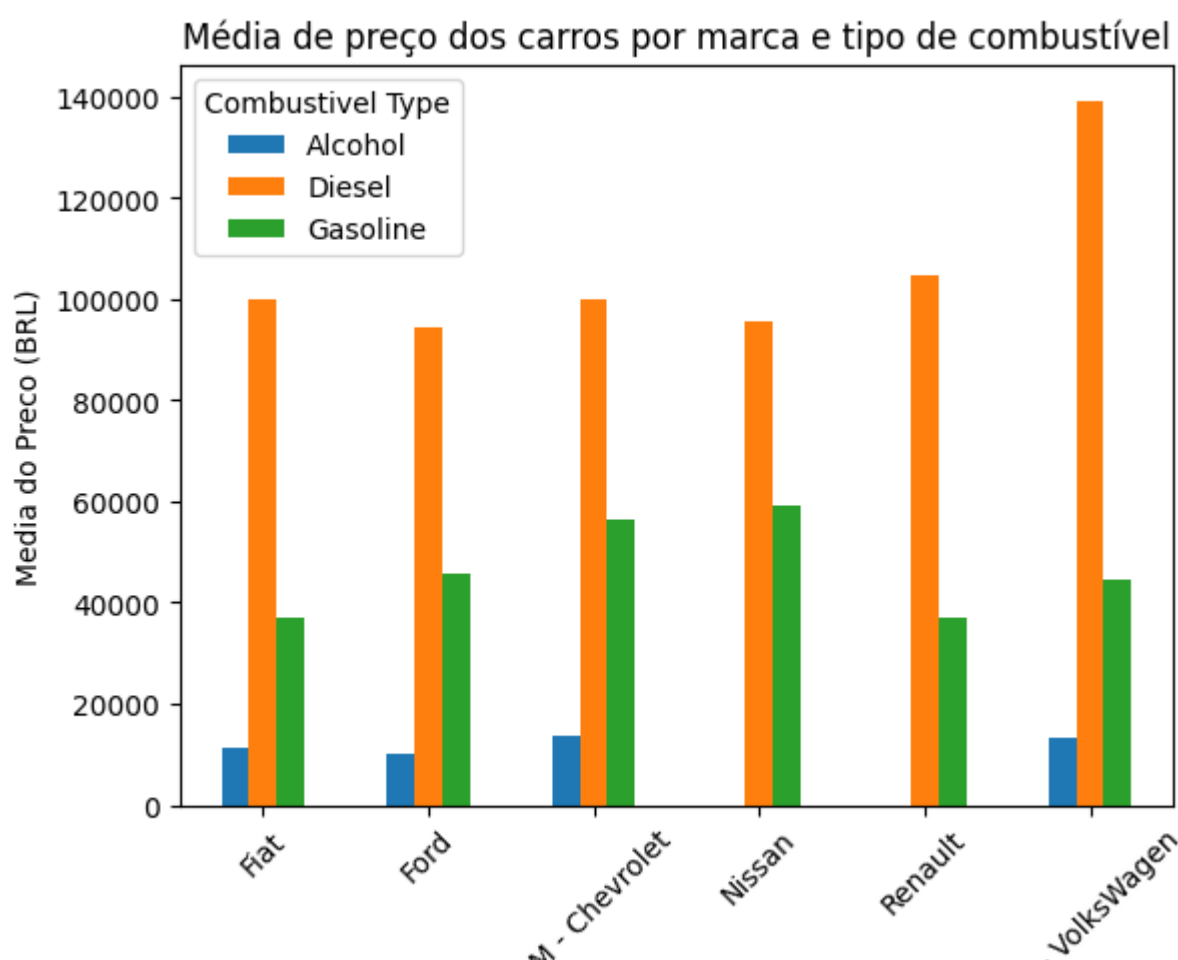
```

1 #f Gere um gráfico da distribuição da média de preço dos carros por marca e tipo de combustível
2
3 def plot_avg_price_by_brand_and_fuel(db):
4     avg_price_by_brand_fuel = db.groupby(['brand', 'fuel'])['avg_price_brl'].mean().unstack()
5
6     plt.figure(figsize=(20, 10))
7     avg_price_by_brand_fuel.plot(kind='bar', stacked=False)
8     plt.title('Média de preço dos carros por marca e tipo de combustível')
9     plt.xlabel('Marca')
10    plt.ylabel('Media do Preco (BRL)')
11    plt.xticks(rotation=45)
12    plt.legend(title='Combustivel Type')
13    plt.show()

```

```
1 plot_avg_price_by_brand_and_fuel(db)
```

<Figure size 2000x1000 with 0 Axes>



G*

VW -

Marca

```
1 def convert_to_numeric(db):
2     db['year_of_reference'] = db['year_of_reference'].astype(float)
3     db['fuel_numeric'] = db['fuel'].replace({
4         'Alcohol': 1,
5         'Gasoline': 2,
6         'Diesel': 3
7     })
8     db['fuel_numeric'] = db['fuel_numeric'].astype(float)
9     db['gear_numeric'] = db['gear'].replace({
10     'automatic': 1,
11     'manual': 2,
12     })
13     db['gear_numeric'] = db['gear_numeric'].astype(float)
14     db['engine_size'] = db['engine_size'].str.replace(',', '.')
15
16
17     return db

1 #Parte 3
2 import pandas as pd
3 import numpy as np
4 from sklearn.model_selection import train_test_split
5 from sklearn.ensemble import RandomForestRegressor
6 from xgboost import XGBRegressor
7 from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score
8 from sklearn.preprocessing import LabelEncoder
9
10
11 def machine_learning_model(db):
12
13     categorical_columns = db.select_dtypes(exclude=[np.number]).columns
14
15     for col in categorical_columns:
16         if db[col].dtype == 'object':
17             le = LabelEncoder()
18             db[col] = le.fit_transform(db[col].fillna("missing"))
19
20     #a. Escolha as variáveis numéricas (modelos de Regressão) para serem as variáveis independentes do modelo.
21     X = db.drop(columns=["avg_price_brl"])
22     y = db["avg_price_brl"]
23
24     #b. Crie partições contendo 75% dos dados para treino e 25% para teste
25     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_state=42)
26
27     #c. Treine modelos RandomForest (biblioteca RandomForestRegressor) e XGBoost
28     #(biblioteca XGBRegressor)
29     #c. 1 -> Random Forest
30     rf_model = RandomForestRegressor(n_estimators=100, random_state=42)
31     rf_model.fit(X_train, y_train)
32
33     #c. 2 -> XGBRegressor
34     xgb_model = XGBRegressor(n_estimators=100, learning_rate=0.1, random_state=42)
35     xgb_model.fit(X_train, y_train)
36
37     #d. Grave os valores preditos em variáveis criadas
38     rf_predictions = rf_model.predict(X_test)
39     xgb_predictions = xgb_model.predict(X_test)
40
41
42     rf_mse = mean_squared_error(y_test, rf_predictions)
43     rf_mae = mean_absolute_error(y_test, rf_predictions)
44     rf_r2 = r2_score(y_test, rf_predictions)
45
46     xgb_mse = mean_squared_error(y_test, xgb_predictions)
47     xgb_mae = mean_absolute_error(y_test, xgb_predictions)
48     xgb_r2 = r2_score(y_test, xgb_predictions)
49
50
51     rf_importance = rf_model.feature_importances_
52     xgb_importance = xgb_model.feature_importances_
53
54     #e. Realize a análise de importância das variáveis para estimar a variável target, para cada
55     #modelo treinad
```



```

56
57     feature_names = X.columns
58     print("Importância das variáveis - Random Forest:")
59     for feature, importance in zip(feature_names, rf_importance):
60         print(f"{feature}: {importance}")
61
62     print("\nImportância das variáveis - XGBoost:")
63     for feature, importance in zip(feature_names, xgb_importance):
64         print(f"{feature}: {importance}")
65     #g. Escolha o melhor modelo com base nas métricas de avaliação MSE, MAE e R
66
67     if rf_r2 > xgb_r2:
68         print("\nModelo escolhido: RandomForestRegressor")
69         best_model = "RandomForestRegressor"
70         best_mse = rf_mse
71         best_mae = rf_mae
72         best_r2 = rf_r2
73     else:
74         print("\nModelo escolhido: XGBRegressor")
75         best_model = "XGBRegressor"
76         best_mse = xgb_mse
77         best_mae = xgb_mae
78         best_r2 = xgb_r2
79
80
81     print(f"\n{best_model} - Resultados:")
82     print(f"MSE: {best_mse}")
83     print(f"MAE: {best_mae}")
84     print(f"R²: {best_r2}")
85
86
87     if best_model == "RandomForestRegressor":
88         print("\nO modelo RandomForestRegressor teve um desempenho melhor baseado nas métricas de MSE, MAE e R².")
89     else:
90         print("\nO modelo XGBRegressor teve um desempenho melhor baseado nas métricas de MSE, MAE e R².")
91
92     return best_model, best_mse, best_mae, best_r2
93
94
95
96 db = extractData(file_path)
97 db = patternColum(db)
98 db = convert_to_numeric(db)
99 best_model, best_mse, best_mae, best_r2 = machine_learning_model(db)
100

```

```

<ipython-input-3-f34971953ffd>:5: DtypeWarning: Columns (1,2,3,4,5,6,7,8) have mixed types. Specify dtype option c
db = pd.read_csv(dataBaseFile)
Index(['year_of_reference', 'month_of_reference', 'fiipe_code',
      'authentication', 'brand', 'model', 'fuel', 'gear', 'engine_size',
      'year_model', 'avg_price_brl'],
      dtype='object')

```

```

<ipython-input-17-6ca10f1cf4c2>:3: FutureWarning: Downcasting behavior in `replace` is deprecated and will be remc
db['fuel_numeric'] = db['fuel'].replace({

```

```

<ipython-input-17-6ca10f1cf4c2>:9: FutureWarning: Downcasting behavior in `replace` is deprecated and will be remc
db['gear_numeric'] = db['gear'].replace({

```

```

Importância das variáveis - Random Forest:
year_of_reference: 0.00015000201751926782
month_of_reference: 1.0824179597380953e-05
fiipe_code: 0.006806591324575154
authentication: 0.2612201294919923
brand: 0.001867939703511724
model: 0.006721040761126544
fuel: 0.00011101886176718157
gear: 0.008111987335865854
engine_size: 0.3536233409122243
year_model: 0.2907043482230164
fuel_numeric: 0.06318122432247504
gear_numeric: 0.0074915528663289516

```

```

Importância das variáveis - XGBoost:
year_of_reference: 0.0010102057131007314
month_of_reference: 0.00011561973951756954
fiipe_code: 0.02325967513024807
authentication: 0.04566916823387146
brand: 0.001510525238700211
model: 0.004453557543456554
fuel: 0.006287244614213705
gear: 0.0005535298259928823
engine_size: 0.24163946509361267
year_model: 0.15439309179782867
fuel_numeric: 0.5211079120635986

```

```
gear_numeric: 0.0
```

```
Modelo escolhido: RandomForestRegressor
```

```
RandomForestRegressor - Resultados:
```

```
MSE: 7203821.98810656
```

```
MAE: 76.91745526445872
```

```
R²: 0.9972845770158916
```

```
O modelo RandomForestRegressor teve um desempenho melhor baseado nas métricas de MSE, MAE e R².
```