PERTEMUAN 4

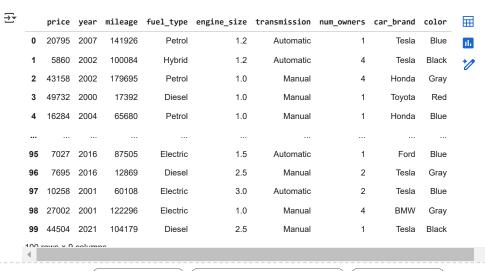
LATIHAN PRAKTIKUM FEATURE ENGINEERING

Nama : Kinanti AnggraeniNIM : 4112322010

Import Dataset

import pandas as pd

df = pd.read_csv('/content/used_cars_dataset.csv')
df



Cek missing values

```
# Cek jumlah missing values di setiap kolom
missing_values = df.isnull().sum()
```

Menampilkan kolom yang memiliki missing values beserta jumlahnya missing_values[missing_values > 0]



tidak terdapat missing value pada dataset ini sehingga tidak memerlukan imputasi

Ubah fitur kategorikal menjadi numerik

melihat masing-masing kategori unik setiap fitur

```
# Melakukan One-Hot Encoding
df_baru = pd.get_dummies(df, columns=['fuel_type', 'transmission', 'car_brand', 'color'], drop_first=True)
# Konversi nilai boolean ke integer (0 dan 1)
# nilai False setara dengan 0 dan True setara dengan 1.
df_baru = df_baru.astype(int)
# Menampilkan hasil One-Hot Encoding
print("\nData Setelah One-Hot Encoding:")
df baru.head()
     Data Setelah One-Hot Encoding:
         price
               year
                     mileage engine_size
                                           num_owners fuel_type_Electric fuel_type_Hybrid fuel_type_Petrol transmission_Manual car_brand_Ford car_brand_
      0 20795 2007
                                                                                           0
                                                                                                                                                   0
                       141926
         5860 2002
                       100084
                                         1
                                                     4
                                                                         0
                                                                                            1
                                                                                                              O
                                                                                                                                   n
                                                                                                                                                   n
                                                     4
                                                                         0
                                                                                            0
                                                                                                                                                   0
        43158 2002
                       179695
      3 49732 2000
                        17392
                                                                         0
                                                                                           0
                                                                                                              0
                                                                                                                                                   0
      4 16284 2004
                        65680
                                                                         0
                                                                                           0
                                                                                                                                                   0
 Langkah berikutnya: (Buat kode dengan df_baru)

    Lihat plot yang direkomendasikan

                                                                                   New interactive sheet
```

· Dampak dari One-Hot Encoding terhadap jumlah fitur dalam dataset

Dampak dari One-Hot Encoding adalah penambahan jumlah fitur dalam dataset. Setiap kategori unik dalam fitur yang di-encode akan menjadi kolom baru. Misalnya, fitur fuel_type dengan 4 kategori (Petrol, Hybrid, Diesel, Electric), setelah One-Hot Encoding, akan mendapatkan 3 kolom baru (karena menggunakan drop_first=True untuk menghindari multikolinearitas).

Jika memiliki fitur dengan n kategori, maka One-Hot Encoding akan menambah n-1 kolom baru ke dataset.

• Efektivitas One-Hot Encoding jika terdapat banyak kategori unik

Jika dataset memiliki 100 kategori unik dalam satu fitur, One-Hot Encoding mungkin tidak efektif karena akan menghasilkan 99 kolom baru. Ini dapat menyebabkan Curse of Dimensionality (Dengan banyaknya fitur, model dapat menjadi lebih kompleks dan sulit untuk dilatih) dan juga Overfitting (Model mungkin belajar noise dalam data daripada pola yang sebenarnya)

scaling pada Fitur Numerik

```
from sklearn.preprocessing import StandardScaler
# Memilih fitur numerik untuk scaling
features_to_scale = ['mileage', 'engine_size', 'num_owners']
# Inisialisasi StandardScaler
scaler = StandardScaler()
# Melakukan scaling
df_baru[features_to_scale] = scaler.fit_transform(df[features_to_scale])
# Menampilkan hasil setelah scaling
\verb|print("\nData Setelah Scaling dengan StandardScaler:")|\\
print(df_baru[features_to_scale].head())
₹
     Data Setelah Scaling dengan StandardScaler:
                 engine_size
         mileage
                              num owners
                                -1.189562
       0.697751
                    -0.948936
     1 -0.078852
                    -0.948936
                                 1.306023
     2 1.398758
                    -1.225191
                                 1.306023
     3 -1.613648
                    -1.225191
                                -1.189562
     4 -0.717404
                    -1.225191
                                -1.189562
```

Perbedaan Hasil antara StandardScaler, MinMaxScaler, dan RobustScaler

StandardScaler:

Mengubah data sehingga memiliki rata-rata 0 dan deviasi standar 1. Cocok untuk data yang terdistribusi normal. Sensitif terhadap outlier.

MinMaxScaler:

Mengubah data ke dalam rentang [0, 1] (atau rentang lain yang ditentukan). Cocok untuk data yang tidak terdistribusi normal dan ketika kita ingin mempertahankan proporsi data. Tidak sensitif terhadap outlier, tetapi outlier dapat mempengaruhi rentang.

RobustScaler:

Menggunakan median dan interquartile range (IQR) untuk scaling. Cocok untuk data dengan banyak outlier. Mengurangi pengaruh outlier pada scaling.

• Dalam Kondisi Apa Lebih Baik Menggunakan MinMaxScaler dibandingkan StandardScaler?

lebih baik menggunakan MinMaxScaler dalam kondisi:

Data Tidak Terdistribusi Normal: Jika data tidak mengikuti distribusi normal, MinMaxScaler dapat memberikan hasil yang lebih baik.

Model yang Sensitif terhadap Skala: Beberapa algoritma, seperti K-Nearest Neighbors (KNN) dan Neural Networks, dapat bekerja lebih baik dengan data yang terletak dalam rentang yang sama.

Preservasi Proporsi Data: Jika kita ingin mempertahankan proporsi data asli dan tidak ingin outlier mempengaruhi skala, MinMaxScaler adalah pilihan yang baik.

Buat fitur baru car_age

```
# Menentukan tahun saat ini
current_year = 2025
# Membuat fitur baru car_age
df_baru['car_age'] = current_year - df['year']
# Menampilkan hasil setelah menambahkan fitur car_age
print("\nData Setelah Menambahkan Fitur car_age:")
print(df_baru[['year', 'car_age']].head())
₹
     Data Setelah Menambahkan Fitur car_age:
        year car_age
       2007
                   18
       2002
                   23
       2002
                   23
        2000
                   25
        2004
```

• Apakah fitur car_age lebih informatif dibandingkan fitur asli year?

fitur car_age lebih informatif dibandingkan year dalam analisis harga mobil bekas karena lebih langsung merepresentasikan usia mobil, yang berpengaruh terhadap nilai jual dan biaya perawatan. Selain itu, car_age menyederhanakan analisis dan mengurangi noise dengan memberikan konteks lebih jelas tentang lama pemakaian mobil.

Namun, fitur year tetap berguna untuk mengetahui tahun produksi mobil.

fitur baru mileage_per_year

```
# Membuat fitur baru mileage_per_year
df_baru['mileage_per_year'] = df_baru['mileage'] / df_baru['car_age']
# Menampilkan hasil setelah menambahkan fitur mileage_per_year
print("\nData Setelah Menambahkan Fitur mileage_per_year:")
print(df_baru[['mileage', 'car_age', 'mileage_per_year']].head())
₹
     Data Setelah Menambahkan Fitur mileage_per_year:
        mileage car_age mileage_per_year
                                  0.038764
       0.697751
                      18
     1 -0.078852
                      23
                                  -0.003428
     2 1.398758
                      23
                                  0.060816
     3 -1.613648
                      25
                                  -0.064546
     4 -0.717404
                                  -0.034162
```

• Mengapa mileage_per_year Bisa Menjadi Fitur yang Lebih Informatif Dibandingkan mileage Saja?

fitur mileage_per_year lebih informatif dibandingkan mileage saja karena memberikan konteks tentang seberapa intens mobil digunakan setiap tahunnya. Hal ini memudahkan perbandingan antar mobil dengan usia berbeda, menjadi indikator kondisi mobil, serta lebih relevan dalam penilaian harga mobil bekas. Dengan mileage_per_year, kita dapat memahami pola penggunaan mobil secara lebih akurat dibandingkan hanya melihat total jarak tempuhnya.

Menampilkan dataset terbaru

```
print("\nDataset Terbaru Setelah Menambahkan Fitur Baru:")
df_baru.head()
```



`u:						
wners	fuel_type_Electric	fuel_type_Hybrid	fuel_type_Petrol	transmission_Manual	car_brand_Ford	car_brand
39562	0	0	1	0	0	
)6023	0	1	0	0	0	
)6023	0	0	1	1	0	
39562	0	0	0	1	0	
39562	0	0	1	1	0	
4						

Langkah berikutnya: Buat kode dengan df_baru Lihat plot yang direkomendasikan New interactive sheet

Mulai coding atau <u>buat</u> kode dengan AI.