

SECTION A

1. The primary use of NumPy library is: Ans.(b) Numerical computation
2. Which algorithm is suitable for binary classification: Ans.(c) Logistic Regression
3. A method to detect outliers is: Ans.(b) Z-score
4. Train-test split is mainly used for: Ans.(b) Avoiding Overfitting
5. An evaluation metric for Regression is: Ans.(d) Mean Squared Error

SECTION B

6. Supervised Learning-Supervised learning uses labeled data, where both input and output are known. The model learns by mapping inputs to outputs. Example: Email spam detection (spam vs not spam) Unsupervised Learning-Unsupervised learning works with unlabeled data and finds hidden patterns or structures. Example: Customer segmentation using clustering
7. Feature scaling is the process of bringing all input features to a similar scale so that no single feature dominates due to its large values. Standardization transforms data to have zero mean and unit variance, which helps algorithms like gradient descent and KNN converge faster.
8. Overfitting occurs when a model learns the training data too well, including noise, and performs poorly on unseen data, when the performance gap between training and testing gap is greater than 5% is considered Overfitting and less than 5% is considered Underfitting. Use techniques such as cross-validation, pruning (for trees), dropout, early stopping, and increasing training data.
9. Classification predicts discrete or categorical outputs. The goal is to assign inputs to predefined classes, and also use in categorical data. Regression predicts continuous numerical values, and use for numerical data. The goal is to estimate a quantity.
10. Problem definition and understanding business requirements. 2. Data collection from various sources. 3. Data cleaning and preprocessing. 4. Exploratory Data Analysis (EDA). 5. Feature engineering and selection. 6. Model building and training. 7. Model evaluation and validation. 8. Deployment and monitoring.

SECTION C

Task 1 - Data Preprocessing

```
In [407]:  
import numpy as np  
import pandas as pd  
import matplotlib.pyplot as plt  
import seaborn as sns  
  
Load dataset using pandas  
  
In [408]:  
ds = pd.read_csv('Used_Bikes.csv')  
ds.head()  
  
Out[408]:  
   bike_name    price      city  kms_driven  owner  age  power  brand  
0   TVS Star City Plus Dual Tone 110cc  35000.0  Ahmedabad  17654.0  First Owner  3.0  110.0    TVS  
1   Royal Enfield Classic 350cc  119900.0    Delhi  11000.0  First Owner  4.0  350.0  Royal Enfield  
2   Triumph Daytona 675R 600000.0    Delhi  110.0  First Owner  8.0  675.0   Triumph  
3   TVS Apache RTR 180cc  65000.0  Bangalore  16329.0  First Owner  4.0  180.0    TVS  
4   Yamaha FZ S V 2.0 150cc-Ltd. Edition  80000.0  Bangalore  10000.0  First Owner  3.0  150.0   Yamaha  
  
Detect and handle missing values  
  
In [409]:  
ds.duplicated().sum()  
  
Out[409]:  
np.int64(25324)  
  
In [410]:  
ds.drop_duplicates(inplace=True)  
  
In [411]:  
ds.duplicated().sum()  
  
Out[411]:  
np.int64(0)  
  
In [412]:  
ds.isnull()  
  
Out[412]:  
   bike_name    price      city  kms_driven  owner  age  power  brand  
0   False  False  False  False  False  False  False  False  
1   False  False  False  False  False  False  False  False  
2   False  False  False  False  False  False  False  False  
3   False  False  False  False  False  False  False  False  
4   False  False  False  False  False  False  False  False  
...  ...  ...  ...  ...  ...  ...  ...  
9362  False  False  False  False  False  False  False  False  
9369  False  False  False  False  False  False  False  False  
9370  False  False  False  False  False  False  False  False  
9371  False  False  False  False  False  False  False  False  
9372  False  False  False  False  False  False  False  False  
  
7324 rows × 8 columns  
  
In [413]:  
ds.isnull().sum()  
  
Out[413]:  
bike_name    0  
price        0  
city         0  
kms_driven  0  
owner        0  
age          0  
power        0  
brand        0  
dtype: int64  
  
In [414]:  
# if we have null values we can drop them using dropna() function  
ds1 = ds.dropna()  
ds1.head()  
  
Out[414]:  
   bike_name    price      city  kms_driven  owner  age  power  brand  
0   TVS Star City Plus Dual Tone 110cc  35000.0  Ahmedabad  17654.0  First Owner  3.0  110.0    TVS  
1   Royal Enfield Classic 350cc  119900.0    Delhi  11000.0  First Owner  4.0  350.0  Royal Enfield  
2   Triumph Daytona 675R 600000.0    Delhi  110.0  First Owner  8.0  675.0   Triumph  
3   TVS Apache RTR 180cc  65000.0  Bangalore  16329.0  First Owner  4.0  180.0    TVS  
4   Yamaha FZ S V 2.0 150cc-Ltd. Edition  80000.0  Bangalore  10000.0  First Owner  3.0  150.0   Yamaha  
  
In [415]:  
# now we fill the null values  
ds.fillna(ds.mean(numeric_only=True))  
ds.head()  
  
Out[415]:  
   bike_name    price      city  kms_driven  owner  age  power  brand  
0   TVS Star City Plus Dual Tone 110cc  35000.0  Ahmedabad  17654.0  1.3.0  110.0    TVS  
1   Royal Enfield Classic 350cc  119900.0    Delhi  11000.0  1.4.0  350.0  Royal Enfield  
2   Triumph Daytona 675R 600000.0    Delhi  110.0  1.8.0  675.0   Triumph  
3   TVS Apache RTR 180cc  65000.0  Bangalore  16329.0  1.4.0  180.0    TVS  
4   Yamaha FZ S V 2.0 150cc-Ltd. Edition  80000.0  Bangalore  10000.0  1.3.0  150.0   Yamaha  
  
Summary Insights : The dataset was successfully loaded using the pandas library, allowing structured data manipulation and analysis. Missing values were identified using null-value checks and handled appropriately by either removing incomplete records or imputing values using statistical measures such as mean or median. Numerical features were analyzed for outliers to ensure data quality and reliability. Handling missing values and outliers improved data quality and prepared the dataset for accurate and stable machine learning modeling. The Z-score method was used to identify values far from the mean, which is effective when data follows a normal distribution.
```

Task 2 - Model Building

```
Train-test split  
  
In [416]:  
x = pd.get_dummies(x, drop_first=True)  
  
In [417]:  
ds.fillna(ds.mean(numeric_only=True), inplace=True)  
  
In [418]:  
dt = {'First Owner':1, 'Second Owner':2, 'Third Owner':3, 'Fourth Owner Or More':4}  
  
Out[418]:  
{'First Owner': 1,  
 'Second Owner': 2,  
 'Third Owner': 3,  
 'Fourth Owner Or More': 4}  
  
In [419]:  
ds['owner']=ds['owner'].map(dt)  
  
In [420]:  
ds.head()  
  
Out[420]:  
   bike_name    price      city  kms_driven  owner  age  power  brand  
0   TVS Star City Plus Dual Tone 110cc  35000.0  Ahmedabad  17654.0  1.3.0  110.0    TVS  
1   Royal Enfield Classic 350cc  119900.0    Delhi  11000.0  1.4.0  350.0  Royal Enfield  
2   Triumph Daytona 675R 600000.0    Delhi  110.0  1.8.0  675.0   Triumph  
3   TVS Apache RTR 180cc  65000.0  Bangalore  16329.0  1.4.0  180.0    TVS  
4   Yamaha FZ S V 2.0 150cc-Ltd. Edition  80000.0  Bangalore  10000.0  1.3.0  150.0   Yamaha  
  
In [421]:  
brands=['TVS', 'Yamaha', 'Bajaj', 'Hero', 'Royal Enfield', 'Honda', 'Suzuki', 'TVKTM', 'KTM', 'Kawasaki', 'Harley-Davidson', 'Kawasaki', 'Hyosung', 'Honda', 'Benelli', 'Triumph', 'Ducati', 'BMW']  
ds['brand'].isin(brands)  
  
Out[421]:  
   bike_name    price      city  kms_driven  owner  age  power  brand  
0   TVS Star City Plus Dual Tone 110cc  35000.0  Ahmedabad  17654.0  1.3.0  110.0    TVS  
1   Royal Enfield Classic 350cc  119900.0    Delhi  11000.0  1.4.0  350.0  Royal Enfield  
2   Triumph Daytona 675R 600000.0    Delhi  110.0  1.8.0  675.0   Triumph  
3   TVS Apache RTR 180cc  65000.0  Bangalore  16329.0  1.4.0  180.0    TVS  
4   Yamaha FZ S V 2.0 150cc-Ltd. Edition  80000.0  Bangalore  10000.0  1.3.0  150.0   Yamaha  
  
In [422]:  
brand_ls = list(ds['brand'].value_counts().head(16).keys())  
brand_ls  
  
Out[422]:  
['Bajaj',  
 'Royal Enfield',  
 'Hero',  
 'Honda',  
 'TVKTM',  
 'KTM',  
 'Suzuki',  
 'Harley-Davidson',  
 'Kawasaki',  
 'Hyosung',  
 'Honda',  
 'Benelli',  
 'Triumph',  
 'Ducati',  
 'BMW']  
  
In [423]:  
ds2 = ds[ds['brand'].isin(brand_ls)]  
ds2.head()  
  
Out[423]:  
   bike_name    price      city  kms_driven  owner  age  power  brand  
0   TVS Star City Plus Dual Tone 110cc  35000.0  Ahmedabad  17654.0  1.3.0  110.0    TVS  
1   Royal Enfield Classic 350cc  119900.0    Delhi  11000.0  1.4.0  350.0  Royal Enfield  
2   Triumph Daytona 675R 600000.0    Delhi  110.0  1.8.0  675.0   Triumph  
3   TVS Apache RTR 180cc  65000.0  Bangalore  16329.0  1.4.0  180.0    TVS  
4   Yamaha FZ S V 2.0 150cc-Ltd. Edition  80000.0  Bangalore  10000.0  1.3.0  150.0   Yamaha  
  
In [424]:  
brand_dict = {brand: 1, brand_ls[i]: i+1 for i in range(len(brand_ls))}  
brand_dict  
  
Out[424]:  
{'Bajaj': 1,  
 'Royal Enfield': 2,  
 'Hero': 3,  
 'Honda': 4,  
 'TVKTM': 5,  
 'KTM': 6,  
 'Suzuki': 7,  
 'Harley-Davidson': 8,  
 'Kawasaki': 9,  
 'Hyosung': 10,  
 'Honda': 11,  
 'Benelli': 12,  
 'Triumph': 13,  
 'Ducati': 14,  
 'BMW': 15}  
  
In [425]:  
brand_ls1 = {brand: 1, brand_ls[i]: i+1 for i in range(len(brand_ls))}  
brand_ls1  
  
Out[425]:  
{'Bajaj': 1,  
 'Royal Enfield': 2,  
 'Hero': 3,  
 'Honda': 4,  
 'TVS': 5,  
 'TVKTM': 6,  
 'KTM': 7,  
 'Suzuki': 8,  
 'Harley-Davidson': 9,  
 'Kawasaki': 10,  
 'Hyosung': 11,  
 'Honda': 12,  
 'Benelli': 13,  
 'Triumph': 14,  
 'Ducati': 15,  
 'BMW': 16}  
  
In [426]:  
brand_ls2 = pd.Series(ds['brand'].map(brand_ls1))  
brand_ls2  
  
Out[426]:  
   bike_name    price      city  kms_driven  owner  age  power  brand  
0   TVS Star City Plus Dual Tone 110cc  35000.0  Ahmedabad  17654.0  1.3.0  110.0    TVS  
1   Royal Enfield Classic 350cc  119900.0    Delhi  11000.0  1.4.0  350.0  Royal Enfield  
2   Triumph Daytona 675R 600000.0    Delhi  110.0  1.8.0  675.0   Triumph  
3   TVS Apache RTR 180cc  65000.0  Bangalore  16329.0  1.4.0  180.0    TVS  
4   Yamaha FZ S V 2.0 150cc-Ltd. Edition  80000.0  Bangalore  10000.0  1.3.0  150.0   Yamaha  
  
In [427]:  
ds2.drop(['bike_name', 'city'], axis=1, inplace=True)  
  
Out[427]:  
/var/folders/lc/vzgpk1x9w0j360t2zjm_w8000gn/T/ipykernel_34389/2965310615.py:1: SettingWithCopyWarning:  
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy  
ds2.drop(['bike_name', 'city'], axis=1, inplace=True)  
  
In [428]:  
ds2.head()  
  
Out[428]:  
   price  kms_driven  owner  age  power  brand  
0  35000.0  17654.0  1.3.0  110.0    T  
1  119900.0  11000.0  1.4.0  350.0  Royal Enfield  
2  600000.0  110.0  1.8.0  675.0   Triumph  
3  65000.0  16329.0  1.4.0  180.0    T  
4  80000.0  10000.0  1.3.0  150.0   Yamaha  
  
In [429]:  
x = ds2.drop(['price'], axis=1)  
y = ds2[['price']]  
  
Out[429]:  
   price  
0  35000.0  
1  119900.0  
2  600000.0  
3  65000.0  
4  80000.0  
...  ...  
9362  48587.0  
9369  60000.0  
9370  3430.0  
9371  21300.0  
9372  7127.0  
  
7307 rows × 5 columns  
  
In [431]:  
x  
Out[431]:  
   kms_driven  owner  age  power  brand  
0  17654.0  1.3.0  110.0    T  
1  11000.0  1.4.0  350.0  Royal Enfield  
2  110.0  1.8.0  675.0   Triumph  
3  16329.0  1.4.0  180.0    T  
4  10000.0  1.3.0  150.0   Yamaha  
  
In [432]:  
y  
Out[432]:  
   price  
0  35000.0  
1  119900.0  
2  600000.0  
3  65000.0  
4  80000.0  
...  ...  
9362  48587.0  
9369  60000.0  
9370  3430.0  
9371  21300.0  
9372  7127.0  
  
7307 rows × 5 columns  
  
Train Linear Regression model  
  
In [433]:  
from sklearn.model_selection import train_test_split  
  
In [434]:  
x_train, x_test, y_train, y_test = train_test_split(x,y, test_size=0.2)  
  
In [435]:  
x_train  
Out[435]:  
   price  kms_driven  owner  age  power  brand  
0  35000.0  
1  119900.0  
2  600000.0  
3  65000.0  
4  80000.0  
...  ...  
9362  48587.0  
9369  60000.0  
9370  3430.0  
9371  21300.0  
9372  7127.0  
  
5645 rows × 5 columns  
  
In [436]:  
y_train  
Out[436]:  
   price  
0  35000.0  
1  119900.0  
2  600000.0  
3  65000.0  
4  80000.0  
...  ...  
9362  48587.0  
9369  60000.0  
9370  3430.0  
9371  21300.0  
9372  7127.0  
  
5645 rows × 1 columns  
  
In [437]:  
y_test  
Out[437]:  
   price  
0  27500.0  
1  70000.0  
2  110000.0  
3  140000.0  
4  120000.0  
...  ...  
9362  11735.0  
9369  139000.0  
9370  45000.0  
9371  61700.0  
9372  10000.0  
  
1462 rows × 1 columns  
  
Train Linear Regression model  
  
In [438]:  
from sklearn.linear_model import LinearRegression  
  
In [439]:  
lrr = LinearRegression()  
  
In [440]:  
lrr.fit(x_train, y_train)  
Out[440]:  
LinearRegression()  
Parameters  
  
Predictions & Evaluation  
  
In [441]:  
from sklearn.metrics import mean_squared_error, r2_score  
  
In [442]:  
y_pred = lrr.predict(x_test)  
y  
Out[442]:  
   price  
0  27500.0  
1  70000.0  
2  110000.0  
3  140000.0  
4  120000.0  
...  ...  
9362  11735.0  
9369  139000.0  
9370  45000.0  
9371  61700.0  
9372  10000.0  
  
1462 rows × 1 columns  
  
In [443]:  
y_test  
Out[443]:  
   price  
0  27500.0  
1  70000.0  
2  110000.0  
3  140000.0  
4  120000.0  
...  ...  
9362  11735.0  
9369  139000.0  
9370  45000.0  
9371  61700.0  
9372  10000.0  
  
1462 rows × 1 columns  
  
Train Linear Regression model  
  
In [444]:  
from sklearn.linear_model import LinearRegression  
  
In [445]:  
lrr = LinearRegression()  
  
In [446]:  
lrr.fit(x_train, y_train)  
Out[446]:  
LinearRegression()  
Parameters  
  
Predictions & Evaluation  
  
In [447]:  
from sklearn.metrics import mean_squared_error, r2_score  
  
In [448]:  
y_pred = lrr.predict(x_test)  
y  
Out[448]:  
   price  
0  27500.0  
1  70000.0  
2  110000.0  
3  140000.0  
4  120000.0  
...  ...  
9362  11735.0  
9369  139000.0  
9370  45000.0  
9371  61700.0  
9372  10000.0  
  
1462 rows × 1 columns  
  
In [449]:  
y_test  
Out[449]:  
   price  
0  27500.0  
1  70000.0  
2  110000.0  
3  140000.0  
4  120000.0  
...  ...  
9362  11735.0  
9369  139000.0  
9370  45000.0  
9371  61700.0  
9372  10000.0  
  
1462 rows × 1 columns  
  
Train Linear Regression model  
  
In [450]:  
from sklearn.linear_model import LinearRegression  
  
In [451]:  
lrr = LinearRegression()  
  
In [452]:  
lrr.fit(x_train, y_train)  
Out[452]:  
LinearRegression()  
Parameters  
  
Predictions & Evaluation  
  
In [453]:  
from sklearn.metrics import mean_squared_error, r2_score  
  
In [454]:  
y_pred = lrr.predict(x_test)  
y  
Out[454]:  
   price  
0  27500.0  
1  70000.0  
2  110000.0  
3  140000.0  
4  120000.0  
...  ...  
9362  11735.0  
9369  139000.0  
9370  45000.0  
9371  61700.0  
9372  10000.0  
  
1462 rows × 1 columns  
  
In [455]:  
y_test  
Out[455]:  
   price  
0  27500.0  
1  70000.0  
2  110000.0  
3  140000.0  
4  120000.0  
...  ...  
9362  11735.0  
9369  139000.0  
9370  45000.0  
9371  61700.0  
9372  10000.0  
  
1462 rows × 1 columns
```