

ML_Pipeline\main.py

```
1  # #imports
2  import pandas as pd
3  import numpy as np
4  import joblib
5
6
7
8  from sklearn.model_selection import train_test_split
9  from sklearn.pipeline import Pipeline
10 from sklearn.compose import ColumnTransformer
11 from sklearn.preprocessing import StandardScaler , OneHotEncoder
12 from sklearn.impute import SimpleImputer
13 from sklearn.ensemble import RandomForestRegressor
14 from sklearn.metrics import mean_absolute_error,mean_squared_error,r2_score,accuracy_score
15
16
17
18
19 # loading Data
20
21 data=pd.read_csv('C:\\Users\\morea\\Desktop\\ML\\ML_Pipeline\\CO2_emission.csv')
22 print(data.head())
23 data.info()
24 # create features and target variable
25 # Create features and target variable
26 x = data.drop(['CO2_Emissions'], axis=1) # Correct column name
27 y = data['CO2_Emissions'] # Correct column name
28
29 # Split categorical and numerical
30 numerical_cols = ['Model_Year', 'Engine_Size', 'Cylinders',
31                  'Fuel_Consumption_in_City(L/100 km)',
32                  'Fuel_Consumption_in_City_Hwy(L/100 km)',
33                  'Fuel_Consumption_comb(L/100km)', 'Smog_Level']
34 categorical_cols = ['Make', 'Model', 'Vehicle_Class', 'Transmission']
35
36 print(numerical_cols)
37 print(categorical_cols)
38 # Rest of your pipeline steps...
39
40
41 numerical_pipeline=Pipeline([
42     ('imputer',SimpleImputer(strategy='mean')),
43     ('scaler',StandardScaler())
44 ])
45
46 categorical_pipeline=Pipeline([
47     ('imputer',SimpleImputer(strategy='most_frequent')),
48     ('encoder',OneHotEncoder(handle_unknown='ignore')) #convert all strings into binary
```

```
49 ])  
50  
51 #join the pipelines together # combine the pipelines  
52 preprocessor=ColumnTransformer([  
53     ('num',numerical_pipeline,numerical_cols),  
54     ('cat',categorical_pipeline,categorical_cols)  
55 ])  
56 pipeline=Pipeline([  
57     ('preprocessor',preprocessor),  
58     ('model',RandomForestRegressor()),  
59 ])  
60  
61  
62 # split into train and test  
63 x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.2,random_state=42)  
64  
65 # train and predict  
66 pipeline.fit(x_train,y_train)  
67  
68 prediction=pipeline.predict(x_test)  
69 # evaluate accuract  
70 mse=mean_squared_error(y_test,prediction)  
71 rmse=np.sqrt(mse)  
72 r2=r2_score(y_test,prediction)  
73 print(mse,rmse,r2)  
74
```