

~\OneDrive\Documents\regression.py

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.model_selection import train_test_split
6 from sklearn.preprocessing import StandardScaler, OneHotEncoder
7 from sklearn.compose import ColumnTransformer
8 from sklearn.linear_model import LinearRegression
9 from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score
10
11 # Data Loading
12 data = pd.read_csv('books_data.csv')
13
14 # Data Printing
15 print("Dataset Head:")
16 print(data.head())
17
18 # Data Preprocessing
19 # Separate features and target
20 X = data.drop('Price', axis=1) # Features (e.g., Title)
21 y = data['Price'] # Target (Price)
22
23 # Remove the currency symbol and convert Price to float
24 y = y.str.replace('£', '').astype(float)
25
26 # Encoding categorical data (Title)
27 preprocessor = ColumnTransformer(
28     transformers=[
29         ('title_encoder', OneHotEncoder(sparse_output=True), ['Title']) # One-hot encode the
30         Title column
31     ]
32 )
33 X = preprocessor.fit_transform(X)
34
35 # Splitting Data
36 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
37
38 # Feature Scaling
39 scaler = StandardScaler(with_mean=False) # Set with_mean=False for sparse matrices
40 X_train = scaler.fit_transform(X_train)
41 X_test = scaler.transform(X_test)
42
43 # Model Initialization
44 model1 = LinearRegression()
45 model2 = LinearRegression()
46
47 # Model Training
48 model1.fit(X_train, y_train)
```

```
48 model2.fit(X_train, y_train)
49
50 # Model Evaluation
51 y_pred1 = model1.predict(X_test)
52 y_pred2 = model2.predict(X_test)
53
54 rmse1 = np.sqrt(mean_squared_error(y_test, y_pred1))
55 rmse2 = np.sqrt(mean_squared_error(y_test, y_pred2))
56 mae1 = mean_absolute_error(y_test, y_pred1)
57 mae2 = mean_absolute_error(y_test, y_pred2)
58 r2_1 = r2_score(y_test, y_pred1)
59 r2_2 = r2_score(y_test, y_pred2)
60
61 print("\nModel 1 Metrics:")
62 print("RMSE:", rmse1, "MAE:", mae1, "R^2:", r2_1)
63 print("\nModel 2 Metrics:")
64 print("RMSE:", rmse2, "MAE:", mae2, "R^2:", r2_2)
65
66 # Model Comparison
67 if r2_1 > r2_2:
68     print("\nModel 1 performs better.")
69 else:
70     print("\nModel 2 performs better.")
71
72 # Visualization: Actual vs Predicted
73 plt.figure(figsize=(12, 6))
74 plt.scatter(y_test, y_pred1, alpha=0.7, label='Model 1 Predictions', color='blue')
75 plt.scatter(y_test, y_pred2, alpha=0.7, label='Model 2 Predictions', color='red')
76 plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], 'k--', lw=2)
77 plt.title('Actual vs Predicted')
78 plt.xlabel('Actual Prices')
79 plt.ylabel('Predicted Prices')
80 plt.legend()
81 plt.show()
82
83 # Visualization: Residuals
84 plt.figure(figsize=(12, 6))
85 sns.histplot(y_test - y_pred1, kde=True, color='blue', label='Model 1 Residuals', bins=30)
86 sns.histplot(y_test - y_pred2, kde=True, color='red', label='Model 2 Residuals', bins=30)
87 plt.title('Residual Distribution')
88 plt.xlabel('Residuals')
89 plt.ylabel('Frequency')
90 plt.legend()
91 plt.show()
```