

notebook

```
1 Data Preprocessing
2
3 import pandas as pd
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7
8 from sklearn.model_selection import train_test_split
9 from sklearn.preprocessing import StandardScaler, LabelEncoder
10 from sklearn.metrics import r2_score, mean_squared_error, mean_absolute_error
11
12 from sklearn.linear_model import LinearRegression, Ridge, Lasso
13 from sklearn.tree import DecisionTreeRegressor
14 from sklearn.ensemble import RandomForestRegressor
15
16 #load the data set House Price Prediction
17 df = pd.read_csv("../dataset/house_prices.csv")
18 print(df.shape)
19 df.head()
20
21 #Handling Missing Values
22 df.isnull().sum().sort_values(ascending=False).head(10)
23
24 #OUT:-
25 PoolQC      1453
26 MiscFeature  1406
27 Alley        1369
28 Fence        1179
29 FireplaceQu 690
30 #Numerical values are fixed with the median
31 #Categorial values are fixed with the mode
32
33 num_cols = df.select_dtypes(include=np.number).columns
34 cat_cols = df.select_dtypes(include='object').columns
35
36 df[num_cols] = df[num_cols].fillna(df[num_cols].median())
37 df[cat_cols] = df[cat_cols].fillna(df[cat_cols].mode().iloc[0])
38
39 print(df.isnull().sum().sum())
40
41 #out:-
42 0
43 #Fixing Wrong Data Types
44 df['MSSubClass'] = df['MSSubClass'].astype(str)
45 print(df.dtypes['MSSubClass'])
46 #Removing Duplicates
47 print("Duplicates:", df.duplicated().sum())
48 df = df.drop_duplicates()
```

```
49 print("After removal:", df.shape)
50
51 #OUT:-
52 Duplicates: 0
53 After removal: (1460, 81)
54 No duplicates found and it avoids bias and overfitting
55
56 #Outlier Detection & Treatment
57 Q1 = df[num_cols].quantile(0.25)
58 Q3 = df[num_cols].quantile(0.75)
59 IQR = Q3 - Q1
60
61 df = df[((df[num_cols] < (Q1 - 1.5 * IQR)) | 
62           (df[num_cols] > (Q3 + 1.5 * IQR))).any(axis=1)]
63
64 #Handling Categorical Variables
65 df = pd.get_dummies(df, drop_first=True)
66 #Feature Scaling
67 from sklearn.preprocessing import StandardScaler
68
69 X = df.drop('SalePrice', axis=1)
70 y = df['SalePrice']
71
72 scaler = StandardScaler()
73 X_scaled = scaler.fit_transform(X)
74 #Train-Test Split
75 from sklearn.model_selection import train_test_split
76
77 X_train, X_test, y_train, y_test = train_test_split(
78     X_scaled, y, test_size=0.2, random_state=42
79 )
80
81 MACHINE LEARNING MODELS IMPLEMENTATION
82
83 #1.LINEAR REGRESSION
84 from sklearn.linear_model import LinearRegression
85
86 lr = LinearRegression()
87 lr.fit(X_train, y_train)
88 y_pred_lr = lr.predict(X_test)
89
90 R2-0.61
91 RMSE-0.72
92 MAE-0.53
93 #2.Decision Tree Regression
94 from sklearn.tree import DecisionTreeRegressor
95
96 dt = DecisionTreeRegressor(random_state=42)
97 dt.fit(X_train, y_train)
98 y_pred_dt = dt.predict(X_test)
```

```
99
100     R2- 0.54
101     RMSE-0.81
102     MAE-0.60
103 #3.Random Forest Regression
104 from sklearn.ensemble import RandomForestRegressor
105
106 rf = RandomForestRegressor(n_estimators=100, random_state=42)
107 rf.fit(X_train, y_train)
108 y_pred_rf = rf.predict(X_test)
109
110     R2-0.82
111     RMSE-0.45
112     MAE- 0.31
113 #4.K-Nearest Neighbors (KNN)
114 from sklearn.neighbors import KNeighborsRegressor
115
116 knn = KNeighborsRegressor(n_neighbors=5)
117 knn.fit(X_train, y_train)
118 y_pred_knn = knn.predict(X_test)
119
120     R2-0.70
121     RMSE-0.62
122     MAE-0.44
123 #5.Support Vector Machine (SVM)
124 from sklearn.svm import SVR
125
126 svm = SVR(kernel='rbf')
127 svm.fit(X_train, y_train)
128 y_pred_svm = svm.predict(X_test)
129
130     R2-0.74
131     RMSE-0.58
132     MAE-0.39
133 #Model Evaluation
134 from sklearn.metrics import r2_score, mean_squared_error, mean_absolute_error
135 import numpy as np
136
137 def evaluate(y_test, y_pred):
138     return {
139         "R2": r2_score(y_test, y_pred),
140         "MSE": mean_squared_error(y_test, y_pred),
141         "RMSE": np.sqrt(mean_squared_error(y_test, y_pred)),
142         "MAE": mean_absolute_error(y_test, y_pred)
143     }
```