

HOUSE PRICE PREDICTION

The real estate industry is a significant sector of the economy, and accurately predicting house prices is crucial for various stakeholders, including homebuyers, sellers, and real estate investors. Machine learning models can play a pivotal role in providing more precise and data-driven predictions of house prices.

1.PROBELM STATEMENT

The aim of this project is to develop a machine learning model that predicts house prices based on a set of features related to the properties. Specifically, the goal is to create a model that can accurately estimate the selling price of houses, taking into account factors such as location, size, number of bedrooms and bathrooms, neighborhood characteristics, and other relevant attributes.

2.DATA

Data is collected from kaggle

```
In [41]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline

In [42]: !pip install seaborn

Requirement already satisfied: seaborn in c:\users\dinesh\desktop\jio_al\env\lib\site-packages (0.12.2)
Requirement already satisfied: numpy<1.24.0,>=1.17 in c:\users\dinesh\desktop\jio_al\env\lib\site-packages (from seaborn) (1.24.4)
Requirement already satisfied: pandas<=25 in c:\users\dinesh\desktop\jio_al\env\lib\site-packages (from seaborn) (2.0.3)
Requirement already satisfied: matplotlib<3.8.1,>=3.1 in c:\users\dinesh\desktop\jio_al\env\lib\site-packages (from seaborn) (3.7.2)
Requirement already satisfied: contourpy<1.0.1 in c:\users\dinesh\desktop\jio_al\env\lib\site-packages (from matplotlib<3.8.1,>=3.1->seaborn) (1.0.6)
Requirement already satisfied: fonttools<4.22.0 in c:\users\dinesh\desktop\jio_al\env\lib\site-packages (from matplotlib<3.8.1,>=3.1->seaborn) (4.26.0)
Requirement already satisfied: kiwisolver<3.2 in c:\users\dinesh\desktop\jio_al\env\lib\site-packages (from matplotlib<3.8.1,>=3.1->seaborn) (1.4.4)
Requirement already satisfied: packaging<20.8 in c:\users\dinesh\desktop\jio_al\env\lib\site-packages (from matplotlib<3.8.1,>=3.1->seaborn) (23.1)
Requirement already satisfied: pillow<=2.0 in c:\users\dinesh\desktop\jio_al\env\lib\site-packages (from matplotlib<3.8.1,>=3.1->seaborn) (9.4.0)
Requirement already satisfied: pyparsing<3.1,>=2.3.1 in c:\users\dinesh\desktop\jio_al\env\lib\site-packages (from matplotlib<3.8.1,>=3.1->seaborn) (3.0.9)
Requirement already satisfied: python-dateutil<=2.7 in c:\users\dinesh\desktop\jio_al\env\lib\site-packages (from matplotlib<3.8.1,>=3.1->seaborn) (2.8.2)
Requirement already satisfied: pytz<=2020.1 in c:\users\dinesh\desktop\jio_al\env\lib\site-packages (from pandas<=25->seaborn) (2022.7)
Requirement already satisfied: tzdata<=2022.1 in c:\users\dinesh\desktop\jio_al\env\lib\site-packages (from pandas<=25->seaborn) (2023.3)
Requirement already satisfied: six<=1.5 in c:\users\dinesh\desktop\jio_al\env\lib\site-packages (from python-dateutil<=2.7->matplotlib<3.8.1,>=3.1->seaborn) (1.16.0)

In [43]: import seaborn as sns

In [44]: df = pd.read_csv("house-prices - ML Dataset.csv")

In [45]: df.head()
df

Out[45]:
```

| | Home | Price | SqFt | Bedrooms | Bathrooms | Offers | Brick | Neighborhood |
|-----|------|--------|------|----------|-----------|--------|-------|--------------|
| 0 | 1 | 114300 | 1790 | 2 | 2 | 2 | No | East |
| 1 | 2 | 114200 | 2030 | 4 | 2 | 3 | No | East |
| 2 | 3 | 114800 | 1740 | 3 | 2 | 1 | No | East |
| 3 | 4 | 94700 | 1980 | 3 | 2 | 3 | No | East |
| 4 | 5 | 119800 | 2130 | 3 | 3 | 3 | No | East |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 123 | 124 | 119700 | 1900 | 3 | 3 | 3 | Yes | East |
| 124 | 125 | 147900 | 2180 | 4 | 3 | 3 | Yes | East |
| 125 | 126 | 113500 | 2070 | 2 | 2 | 2 | No | North |
| 126 | 127 | 149900 | 2020 | 3 | 3 | 1 | No | West |
| 127 | 128 | 124600 | 2250 | 3 | 3 | 4 | No | North |

128 rows x 8 columns

```
In [46]: df.isna().sum()
Out[46]:
Home      0
Price      0
Sqft       0
Bedrooms  0
Bathrooms  0
Offers     0
Brick      0
Neighborhood
dtype: int64

In [47]: df.duplicated().sum()
Out[47]: 0

In [48]: df.shape
Out[48]: (128, 8)

In [49]: df.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 128 entries, 0 to 127
Data columns (total 8 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   Home        128 non-null    int64
 1   Price       128 non-null    int64
 2   Sqft        128 non-null    int64
 3   Bedrooms    128 non-null    int64
 4   Bathrooms   128 non-null    int64
 5   Offers      128 non-null    int64
 6   Brick       128 non-null    object
 7   Neighborhood 128 non-null    object
dtypes: int64(6), object(2)
memory usage: 8.1+ KB
```

Plotting a scatter plot for finding outliers

```
In [50]: plt.figure(figsize=(10,6))
plt.scatter(df['Sqft'],df['Price'], alpha = 0.5)
plt.grid(True)
plt.title("SCATTER PLOT")
plt.xlabel("Sqft")
plt.ylabel("Prices")
plt.show()

SCATTER PLOT
```

Finding outliers using IQR and dropping them.

```
In [51]: Q1 = df['Price'].quantile(0.25)
Q3 = df['Price'].quantile(0.75)
IQR = Q3 - Q1

lower_bound = Q1 - 1.5*IQR
upper_bound = Q3 + 1.5*IQR
df = df[(df['Price']>= lower_bound)&& (df['Price']<=upper_bound)]
df

Out[51]:
```

| | Home | Price | SqFt | Bedrooms | Bathrooms | Offers | Brick | Neighborhood |
|-----|------|--------|------|----------|-----------|--------|-------|--------------|
| 0 | 1 | 114300 | 1790 | 2 | 2 | 2 | No | East |
| 1 | 2 | 114200 | 2030 | 4 | 2 | 3 | No | East |
| 2 | 3 | 114800 | 1740 | 3 | 2 | 1 | No | East |
| 3 | 4 | 94700 | 1980 | 3 | 2 | 3 | No | East |
| 4 | 5 | 119800 | 2130 | 3 | 3 | 3 | No | East |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 123 | 124 | 119700 | 1900 | 3 | 3 | 3 | Yes | East |
| 124 | 125 | 147900 | 2180 | 4 | 3 | 3 | Yes | East |
| 125 | 126 | 113500 | 2070 | 2 | 2 | 2 | No | North |
| 126 | 127 | 149900 | 2020 | 3 | 3 | 1 | No | West |
| 127 | 128 | 124600 | 2250 | 3 | 3 | 4 | No | North |

127 rows x 8 columns

```
In [52]: plt.figure(figsize=(10,6))
plt.scatter(df['Sqft'],df['Price'], alpha = 8.5)
plt.grid(True)
plt.title("SCATTER PLOT")
plt.xlabel("Sqft")
plt.ylabel("Prices")
plt.show()

SCATTER PLOT
```

Changing Categorical values to numerical values

```
In [53]: dfe = pd.get_dummies(dfo,columns=['Brick','Neighborhood'])

In [54]: dfe.drofe.astype(int)
dfe

Out[54]:
```

| | Home | Price | SqFt | Bedrooms | Bathrooms | Offers | Brick_No | Brick_Yes | Neighborhood_East | Neighborhood_North | Neighborhood_West |
|-----|------|--------|------|----------|-----------|--------|----------|-----------|-------------------|--------------------|-------------------|
| 0 | 1 | 114300 | 1790 | 2 | 2 | 2 | 1 | 0 | 1 | 0 | 0 |
| 1 | 2 | 114200 | 2030 | 4 | 2 | 3 | 1 | 0 | 1 | 0 | 0 |
| 2 | 3 | 114800 | 1740 | 3 | 2 | 1 | 1 | 0 | 1 | 0 | 0 |
| 3 | 4 | 94700 | 1980 | 3 | 2 | 3 | 1 | 0 | 1 | 0 | 0 |
| 4 | 5 | 119800 | 2130 | 3 | 3 | 3 | 1 | 0 | 1 | 0 | 0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 123 | 124 | 119700 | 1900 | 3 | 3 | 3 | 0 | 1 | 1 | 0 | 0 |
| 124 | 125 | 147900 | 2180 | 4 | 3 | 3 | 0 | 1 | 1 | 0 | 0 |
| 125 | 126 | 113500 | 2070 | 2 | 2 | 2 | 1 | 0 | 0 | 1 | 0 |
| 126 | 127 | 149900 | 2020 | 3 | 3 | 1 | 1 | 0 | 0 | 0 | 1 |
| 127 | 128 | 124600 | 2250 | 3 | 3 | 4 | 1 | 0 | 0 | 1 | 0 |

127 rows x 11 columns

Applying "SCALING" on numerical values

```
In [55]: from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()
scaled_features = scaler.fit_transform(dfe)

scaled_df = pd.DataFrame(scaled_features, columns=dfe.columns)

In [56]: scaled_df

Out[56]:
```

| | Home | Price | SqFt | Bedrooms | Bathrooms | Offers | Brick_No | Brick_Yes | Neighborhood_East | Neighborhood_North | Neighborhood_West |
|-----|----------|----------|----------|----------|-----------|--------|----------|-----------|-------------------|--------------------|-------------------|
| 0 | 0.000000 | 0.346026 | 0.298246 | 0.000000 | 0.0 | 0.2 | 1.0 | 0.0 | 1.0 | 0.0 | 0.0 |
| 1 | 0.000784 | 0.346889 | 0.508772 | 0.666667 | 0.0 | 0.4 | 1.0 | 0.0 | 1.0 | 0.0 | 0.0 |
| 2 | 0.015748 | 0.350480 | 0.254386 | 0.333333 | 0.0 | 0.0 | 1.0 | 0.0 | 1.0 | 0.0 | 0.0 |
| 3 | 0.023922 | 0.196319 | 0.464912 | 0.333333 | 0.0 | 0.4 | 1.0 | 0.0 | 1.0 | 0.0 | 0.0 |
| 4 | 0.031496 | 0.388044 | 0.596491 | 0.333333 | 0.5 | 0.4 | 1.0 | 0.0 | 1.0 | 0.0 | 0.0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 122 | 0.008504 | 0.388037 | 0.394737 | 0.333333 | 0.5 | 0.4 | 0.0 | 1.0 | 1.0 | 0.0 | 0.0 |
| 123 | 0.976378 | 0.604284 | 0.622807 | 0.666667 | 0.5 | 0.4 | 0.0 | 1.0 | 1.0 | 0.0 | 0.0 |
| 124 | 0.984262 | 0.340491 | 0.543800 | 0.000000 | 0.0 | 0.2 | 1.0 | 0.0 | 0.0 | 1.0 | 0.0 |
| 125 | 0.992126 | 0.619652 | 0.500000 | 0.333333 | 0.5 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 1.0 |
| 126 | 1.000000 | 0.425613 | 0.701754 | 0.333333 | 0.5 | 0.6 | 1.0 | 0.0 | 0.0 | 1.0 | 0.0 |

127 rows x 11 columns

Splitting data into training and testing sets

```
In [57]: x = scaled_df.drop('Price',axis = 1)
y = scaled_df['Price']

In [58]: from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.2)

In [59]: scaled_df.describe()

Out[59]:
```

| | Home | Price | SqFt | Bedrooms | Bathrooms | Offers | Brick_No | Brick_Yes | Neighborhood_East | Neighborhood_North | Neighborhood_West |
|-------|------------|------------|------------|------------|------------|------------|------------|------------|-------------------|--------------------|-------------------|
| count | 127.000000 | 127.000000 | 127.000000 | 127.000000 | 127.000000 | 127.000000 | 127.000000 | 127.000000 | 127.000000 | 127.000000 | 127.000000 |
| mean | 0.497854 | 0.464912 | 0.464912 | 0.333333 | 0.333333 | 0.333333 | 0.333333 | 0.333333 | 0.333333 | 0.333333 | 0.333333 |
| std | 0.291917 | 0.196319 | 0.196319 | 0.196319 | 0.196319 | 0.196319 | 0.196319 | 0.196319 | 0.196319 | 0.196319 | 0.196319 |
| min | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 |
| 25% | 0.248031 | 0.332326 | 0.377193 | 0.333333 | 0.000000 | 0.200000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 |
| 50% | 0.496063 | 0.434049 | 0.462456 | 0.333333 | 0.000000 | 0.400000 | 1.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 |
| 75% | 0.744094 | 0.602629 | 0.602629 | 0.333333 | 0.500000 | 0.400000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| max | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |

```
In [60]: scaled_df.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 127 entries, 0 to 126
Data columns (total 11 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   Home        127 non-null    int64
 1   Price       127 non-null    float64
 2   Sqft        127 non-null    float64
 3   Bedrooms    127 non-null    float64
 4   Bathrooms   127 non-null    float64
 5   Offers      127 non-null    float64
 6   Brick_No    127 non-null    float64
 7   Brick_Yes   127 non-null    float64
 8   Neighborhood_East 127 non-null    float64
 9   Neighborhood_North 127 non-null    float64
10  Neighborhood_West 127 non-null    float64
dtypes: float64(11)
memory usage: 11.0 KB

In [77]: plt.hist(scaled_df['Price'])
plt.title("Distribution of house prices")

Out[77]:
```

Dropping Home column as there is no use with Home field.

```
In [78]: scaled_df = scaled_df.drop('Home',axis=1)

In [79]: scaled_df

Out[79]:
```

| | Price | SqFt | Bedrooms | Bathrooms | Offers | Brick_No | Brick_Yes | Neighborhood_East | Neighborhood_North | Neighborhood_West |
|-----|----------|----------|----------|-----------|--------|----------|-----------|-------------------|--------------------|-------------------|
| 0 | 0.346026 | 0.298246 | 0.000000 | 0.0 | 0.2 | 1.0 | 0.0 | 1.0 | 0.0 | 0.0 |
| 1 | 0.346889 | 0.508772 | 0.666667 | 0.0 | 0.4 | 1.0 | 0.0 | 1.0 | 0.0 | 0.0 |
| 2 | 0.350480 | 0.254386 | 0.333333 | 0.0 | 0.0 | 1.0 | 0.0 | 1.0 | 0.0 | 0.0 |
| 3 | 0.196319 | 0.464912 | 0.333333 | 0.0 | 0.4 | 1.0 | 0.0 | 1.0 | 0.0 | 0.0 |
| 4 | 0.388044 | 0.596491 | 0.333333 | 0.5 | 0.4 | 1.0 | 0.0 | 1.0 | 0.0 | 0.0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 122 | 0.388037 | 0.394737 | 0.333333 | 0.5 | 0.4 | 0.0 | 1.0 | 1.0 | 0.0 | 0.0 |
| 123 | 0.604284 | 0.622807 | 0.666667 | 0.5 | 0.4 | 0.0 | 1.0 | 1.0 | 0.0 | 0.0 |
| 124 | 0.340491 | 0.543800 | 0.000000 | 0.0 | 0.2 | 1.0 | 0.0 | 0.0 | 1.0 | 0.0 |
| 125 | 0.619652 | 0.500000 | 0.333333 | 0.5 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 1.0 |
| 126 | 0.425613 | 0.701754 | 0.333333 | 0.5 | 0.6 | 1.0 | 0.0 | 0.0 | 1.0 | 0.0 |

127 rows x 10 columns

Identifying correlations between features and the target variable

```
In [88]: corr_matrix = scaled_df.corr()
plt.subplots(figsize=(10,6))
ax = sns.heatmap(corr_matrix,annot = True,linewidth=0.5,fmt='.2f',cmap='vminbu')

Price      1.00    0.53    0.52    0.52    -0.34    -0.44    0.44    -0.13    -0.55    0.71
Sqft       0.53    1.00    0.47    0.52    0.34    -0.06    0.06    -0.28    0.23
Bedrooms   0.52    0.47    1.00    0.41    0.11    -0.03    0.03    -0.08    -0.36    0.46
Bathrooms  0.52    0.52    0.41    1.00    0.14    -0.16    0.16    0.01    -0.27    0.28
Offers     -0.34    0.34    0.11    0.14    1.00    0.15    -0.15    -0.01    -0.34
Brick_No   -0.44    -0.06    0.03    -0.16    0.15    1.00    -0.16    0.25    0.25    -0.10
Brick_Yes  0.44    0.06    0.03    0.16    -0.15    -1.00    1.00    0.16    -0.25    0.10
Neighborhood_East -0.13    0.06    -0.08    0.01    -0.01    -0.16    0.16    1.00    -0.54    -0.48
Neighborhood_North -0.55    -0.28    -0.36    -0.27    0.34    0.25    -0.25    -0.54    1.00    -0.48
Neighborhood_West 0.71    0.23    0.46    0.28    -0.34    -0.10    0.10    -0.48    -0.48    1.00

Price      Sqft      Bedrooms  Bathrooms  Offers  Brick_No  Brick_Yes  Neighborhood_East  Neighborhood_North  Neighborhood_West

In [81]: target_correlation = corr_matrix['Price']
target_correlation

Out[81]:
Price      1.000000
Sqft       0.531726
Bedrooms    0.518775
Bathrooms    0.518775
Offers      -0.335863
Brick_No     0.438237
Brick_Yes    0.438237
Neighborhood_East 0.562294
Neighborhood_North 0.138496
Neighborhood_West 0.562294
Name: Price, dtype: float64

Importing ML models

In [64]: from sklearn.ensemble import RandomForestRegressor
from sklearn.linear_model import LinearRegression
from sklearn import DecisionTreeRegressor

In [65]: models = ("Linear Regression",LinearRegression(),
                "DecisionTreeRegressor",DecisionTreeRegressor(),
                "Random Forest",RandomForestRegressor())

def fit_and_score(models,x_train,x_test,y_train,y_test):
    np.random.seed(42)
    model_scores = {}
    for name,model in models.items():
        model.fit(x_train,y_train)
        model_scores[name]=model.score(x_test,y_test)
    return model_scores

In [66]: fit_and_score(models,x_train,x_test,y_train,y_test)

Out[66]:
('Linear Regression', 0.837406921391103,
 'DecisionTreeRegressor', 0.615447253365919,
 'Random Forest', 0.735558682387638)

EVALUATION METRICS

In [67]: from sklearn.metrics import mean_absolute_error, r2_score, mean_squared_error

In [68]: def fit_and_score(models,x_train,x_test,y_train,y_test):
    np.random.seed(77)
    model_scores = {}
    for name,model in models.items():
        model.fit(x_train,y_train)
        score=model.score(x_test,y_test)
        y_pred = model.predict(x_test)
        mae = mean_absolute_error(y_test,y_pred)
        rmse = r2_score(y_test,y_pred)
        mse = np.sqrt(mean_squared_error(y_test,y_pred))

        print("Mean Absolute Error of ",name," :", mae)
        print("R-squared:",name," :", r2)
        print("Root Mean Squared Error:",name," :", rmse)

In [69]: fit_and_score(models,x_train,x_test,y_train,y_test)

Mean Absolute Error of Linear Regression : 0.06853863258512554
R-squared: Linear Regression : 0.817608893181103
Root Mean Squared Error of DecisionTreeRegressor : 0.2143279377894652
R-squared: DecisionTreeRegressor : 0.537441583256464
Root Mean Squared Error of DecisionTreeRegressor : 0.2359135963763352
Mean Absolute Error of Random Forest : 0.0865119599811232
R-squared: Random Forest : 0.748605114874011
Root Mean Squared Error of Random Forest : 0.3517464685875784

Considering Random Forest Regressor model

Applying hyperparameter tuning on Random Forest Regressor model, even though Linear regression has highest accuracy.Because there are no hyper parameters for Linear regression model, so it can't be tuned.

In [70]: from sklearn.model_selection import GridSearchCV
model = RandomForestRegressor()
model.fit(x_train,y_train)
y_pred = model.predict(x_test)
np.random.seed(42)
param_grid = {
    'n_estimators': [100, 200, 300],
    'max_depth': [None, 10, 20],
    'min_samples_split': [2, 5, 10]
}
scoring = "neg_mean_squared_error"
grid_search = GridSearchCV(model,param_grid,scoring=scoring,cv=5)
grid_search.fit(x_train,y_train)
best_params = grid_search.best_params_
best_model = grid_search.best_estimator_
test_predictions = best_model.predict(x_test)
rmse = mean_squared_error(y_test,test_predictions)
test_rmse = r2_score(y_test,test_predictions)

test_mae
Out[71]: 0.019428395321698682

test_R2
Out[72]: 0.7386756378995945

best_model.score(x_train,y_test)
Out[73]: 0.7386756378995945

Visualizing the predictions of the selected model against the actual house prices.

In [74]: plt.figure(figsize=(10,6))
plt.scatter(y_test,y_pred,c='blue',marker='o',edgecolor='k',label='Actual Prices')
plt.scatter(y_test,y_pred,c='red',marker='o',edgecolor='k',label='Predicted Prices')
plt.grid(True)
plt.plot([min(y_test),max(y_test)], [min(y_test),max(y_test)], linestyle='--',color='red')
plt.show()


```

Finding feature importance scores

```
In [75]: importances = best_model.feature_importances_
feature_names = x.columns
feature_importance_df = pd.DataFrame({'feature':feature_names,'Importance':importances})
feature_importance_df = feature_importance_df.sort_values(by='Importance',ascending=False)

In [76]: feature_importance_df

Out[76]:
```

| | Feature Importance |
|---|-----------------------------|
| 9 | Neighborhood_West 0.545651 |
| 1 | Sqft 0.225615 |
| 4 | Offers 0.050380 |
| 3 | Bathrooms 0.046732 |
| 0 | Home 0.045082 |
| 5 | Brick_No 0.030381 |
| 6 | Brick_Yes 0.027341 |
| 2 | Bedrooms 0.016039 |
| 8 | Neighborhood_North 0.005523 |
| 7 | Neighborhood_East 0.002476 |

Neighborhood_West has the highest feature importance

```
In [ ]:
In [ ]:
```