

ARTIFICIAL INTELLIGENCE

**(Real-World Problem Solving Using Artificial Intelligence
Approaches)**

Project on: “London Housing Price Prediction Model”

Table of Contents

ABSTRACT.....	3
INTRODUCTION.....	4
PROBLEM STATEMENT.....	5
Project Aim and Objectives.....	6
AI APPROACH FOR PREDICTING LONDON HOUSING PRICES DATA.....	7
PROPOSED AI APPROACH	7
Justification for the Chosen AI Approach.....	7
LIBRARY & DATASET IMPORT.....	8
Exploratory Data Analysis (EDA)	9
DATA CLEANING.....	11
Data Visualization	13
Data preprocessing	19
X, Y split.....	19
Model Selection	20
Machine Learning Models.....	20
Linear Regression	20
Random Forest.....	20
Elastic Net	21
Neural Networks	21
Model Comparison.....	22
SUMMARY	23
Conclusion.....	26
Recommendation.....	26
Future Directions	26
References	28
Appendix	29

ABSTRACT

In the dynamic intersection of machine learning and real-world applications, our study evaluates the predictive capabilities of four distinct models—Linear Regression, Random Forest Regressor, Elastic Net, and an Artificial Neural Network (ANN)—in the context of socioeconomics, financial dynamics, and health outcomes. This appraisal employs metrics, including Mean Absolute Error (MAE), Mean Squared Error (MSE), Root Mean Squared Error (RMSE), R2 Score, and Cross-Validation RMSE, to analyze and compare model performance.

The Linear Regression model, serving as a baseline, exhibits moderate predictive ability, with limitations apparent in capturing intricate relationships within socioeconomic, financial, or health datasets. The Random Forest Regressor and Elastic Net strike a balance between interpretability and accuracy, making them promising candidates for applications where understanding model decisions is crucial.

Exploring the landscape of financial dynamics, the Elastic Net, with its regularization techniques, demonstrates competitive performance, emphasizing the importance of balancing model complexity in datasets with numerous features. In contrast, the Random Forest Regressor shows potential in capturing feature importance, guiding the selection of key factors influencing financial outcomes.

Venturing into health-related predictions, the Artificial Neural Network (ANN), while presenting challenges in its current state, holds promise for capturing complex patterns within intricate health datasets. However, optimization through careful architecture design, hyperparameter tuning, and data preprocessing is imperative for unlocking its true potential.

Recommendations stemming from this study include hyperparameter tuning for the Random Forest Regressor and Elastic Net, feature engineering to uncover latent patterns, and comprehensive data exploration for robust predictions. Further, model ensemble approaches and iterative cycles of refinement are advocated to navigate the complexities of socioeconomic, financial, or health-related data.

This study contributes to the broader discourse on machine learning applications in real-world scenarios, providing insights into model strengths and limitations within the context of socioeconomics, financial dynamics, and health outcomes. The findings underscore the importance of a nuanced approach to model selection and refinement to ensure optimal predictions in domains crucial to societal well-being.

INTRODUCTION

The main objective of this project is to develop and evaluate machine learning models for predicting house prices in London, leveraging a comprehensive dataset (London Housing Prices). Four distinct models were employed: Linear Regression, Random Forest Regressor, Elastic Net, and a Neural Network. The models underwent thorough evaluation using various metrics, providing insights into their performance and suitability for the given task. The model will consider various property features to deduce accurate predictions of housing prices. A comparison of these algorithms' performance will be conducted, and the most accurate result will be chosen.

PROBLEM STATEMENT

There are many factors that can affect the value of a house property (e.g., location, size, condition, number of bedrooms, bathroom etc.), these factors can change quite substantially from one property to another. The housing market itself is quite a volatile industry, and is quite dependent on demand and supply fluctuations, not to even mention economic factors such as interest rates & inflation, so it is quite a challenge to predict the price variation over time.

It is also quite challenging to predict housing prices due to the limited data that is available, most datasets contain a limited number of features related to each property, such is why feature engineering is quite important. As a result, it is quite difficult to accurately predict property prices that consider all the factors that influence them.

The London housing dataset contains different house related attributes for properties located in London.

Project Aim and Objectives

Aim:

This project's main aim is to compare the predictive performance of machine learning models and develop a predictive model for London Price Housing prices.

Objectives:

- Model Evaluation
- Regularization Techniques:
- Comprehensive Model Comparison
- Model Comparison and Selection
- Identification of Best-Performing Model
- Recommendations for Model Refinement
- Future Work Considerations

AI APPROACH FOR PREDICTING LONDON HOUSING PRICES DATA

Real estate brokers, potential purchasers, and investors face difficulties in predicting housing costs. This is a major challenge today. There is a need for accurate and unambiguous estimates of property values, which is essential for making effective decisions.

PROPOSED AI APPROACH

- Linear Regression
- Random Forest
- Elastic Net
- Artificial Neural Networks

I will use data from the London Housing Prices dataset to train a machine learning model. The model will include a comprehensive range of property-specific features, neighborhood characteristics, and market trends to provide accurate and informed predictions of housing prices. I propose employing four models; linear regression, random forest, Elastic Net and Neural Networks for predicting housing prices.

Justification for the Chosen AI Approach

The proposed AI approach, employing machine learning algorithms and a range of features, offers the most efficacious and precise method for predicting housing prices. The choice of linear regression, random forest, Elastic net, and Neural network algorithms provides a balance between simplicity and complexity, which allows me to investigate both approaches and choose the one that best suits the data and the problem at hand. Furthermore, the continuous updating and improvement of the model based on new data will ensure its adaptability to changing market conditions.

LIBRARY & DATASET IMPORT

I will be importing the necessary libraries for this project.

Importing the necessary libraries

```
In [52]: import numpy as np
import tensorflow as tf
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.preprocessing import StandardScaler
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense
from sklearn.metrics import r2_score, mean_absolute_error, mean_squared_error

from sklearn.linear_model import LinearRegression
from sklearn.ensemble import RandomForestRegressor
from sklearn.linear_model import ElasticNet
```

Loading the dataset

```
In [2]: df1 = pd.read_csv("london.csv")
```


In [5]: df1

Out[5]:

	Unnamed: 0	Property Name	Price	House Type	Area in sq ft	No. of Bedrooms	No. of Bathrooms	No. of Receptions	Location	City/County	Postal Code
0	0	Queens Road	1675000	House	2716	5	5	5	Wimbledon	London	SW19 8NY
1	1	Seward Street	650000	Flat / Apartment	814	2	2	2	Clerkenwell	London	EC1V 3PA
2	2	Hotham Road	735000	Flat / Apartment	761	2	2	2	Putney	London	SW15 1QL
3	3	Festing Road	1765000	House	1986	4	4	4	Putney	London	SW15 1LP
4	4	Spencer Walk	675000	Flat / Apartment	700	2	2	2	Putney	London	SW15 1PL
...
3475	3475	One Lillie Square	3350000	New development	1410	3	3	3	NaN	Lillie Square	SW6 1UE
3476	3476	St. James's Street	5275000	Flat / Apartment	1749	3	3	3	St James's	London	SW1A 1JT
3477	3477	Ingram Avenue	5995000	House	4435	6	6	6	Hampstead Garden Suburb	London	NW11 6TG
3478	3478	Cork Street	6300000	New development	1506	3	3	3	Mayfair	London	W1S 3AR
3479	3479	Courtenay Avenue	8850000	House	5395	6	6	6	Highgate	London	N6 4LP

3480 rows x 11 columns

Exploratory Data Analysis (EDA)

1. To get the first five rows of the data

Data Exploration

In [3]: `#Let's take a look at the first 5 rows of the dataset
df1.head()`

Out[3]:

	Unnamed: 0	Property Name	Price	House Type	Area in sq ft	No. of Bedrooms	No. of Bathrooms	No. of Receptions	Location	City/County	Postal Code
0	0	Queens Road	1675000	House	2716	5	5	5	Wimbledon	London	SW19 8NY
1	1	Seward Street	650000	Flat / Apartment	814	2	2	2	Clerkenwell	London	EC1V 3PA
2	2	Hotham Road	735000	Flat / Apartment	761	2	2	2	Putney	London	SW15 1QL
3	3	Festing Road	1765000	House	1986	4	4	4	Putney	London	SW15 1LP
4	4	Spencer Walk	675000	Flat / Apartment	700	2	2	2	Putney	London	SW15 1PL

2. To get the last five rows of the data

```
In [4]: #Let's take a look at the last 5 rows of the dataset
df1.tail()
```

Unnamed: 0	Property Name	Price	House Type	Area in sq ft	No. of Bedrooms	No. of Bathrooms	No. of Receptions	Location	City/County	Postal Code
3475	One Lillie Square	3350000	New development	1410	3	3	3	NaN	Lillie Square	SW6 1UE
3476	St James's Street	5275000	Flat/ Apartment	1749	3	3	3	St James's	London	SW1A 1JT
3477	Ingram Avenue	5995000	House	4435	6	6	6	Hampstead Garden Suburb	London	NW11 6TG
3478	Cork Street	6300000	New development	1508	3	3	3	Mayfair	London	W1S 3AR
3479	Courtenay Avenue	8650000	House	5395	6	6	6	Highgate	London	N8 4LP

3. Checking the shape of the data and getting the summary of the Data Frame's information

```
In [27]: #checking the shape of the data
df1.shape
```

```
Out[27]: (3480, 11)
```

```
In [7]: #To get summary of the DataFrame's information
df1.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3480 entries, 0 to 3479
Data columns (total 11 columns):
#   Column              Non-Null Count  Dtype
---  --
0   Unnamed: 0          3480 non-null   int64
1   Property Name       3480 non-null   object
2   Price               3480 non-null   int64
3   House Type         3480 non-null   object
4   Area in sq ft      3480 non-null   int64
5   No. of Bedrooms     3480 non-null   int64
6   No. of Bathrooms    3480 non-null   int64
7   No. of Receptions   3480 non-null   int64
8   Location            2518 non-null   object
9   City/County         3480 non-null   object
10  Postal Code         3480 non-null   object
dtypes: int64(6), object(5)
memory usage: 299.2+ KB
```

4. To get the statistical summary of the dataset and counting the number of missing values of each column

```
In [6]: #getting the statistical summary of dataset
df1.describe().T
```

	count	mean	std	min	25%	50%	75%	max
Unnamed: 0	3480.0	1.739500e+03	1.004734e+03	0.0	889.75	1739.5	2809.25	3479.0
Price	3480.0	1.884173e+06	2.267283e+06	180000.0	750000.00	1220000.0	2150000.00	39750000.0
Area in sq ft	3480.0	1.712974e+03	1.384256e+03	274.0	834.00	1310.0	2157.25	15405.0
No. of Bedrooms	3480.0	3.103736e+00	1.517698e+00	0.0	2.00	3.0	4.00	10.0
No. of Bathrooms	3480.0	3.103736e+00	1.517698e+00	0.0	2.00	3.0	4.00	10.0
No. of Receptions	3480.0	3.103736e+00	1.517698e+00	0.0	2.00	3.0	4.00	10.0

```
In [11]: #To count the number of missing values in each column of the DataFrame
df1.isnull().sum()
```

```
Out[11]: Unnamed: 0          0
Property Name          0
Price                  0
House Type             0
Area in sq ft          0
No. of Bedrooms        0
No. of Bathrooms       0
No. of Receptions      0
Location              962
City/County            0
Postal Code            0
dtype: int64
```

5. To get the count of the number of unique values in each column

```
In [13]: #To count the number of unique values in each column of a DataFrame
df1.nunique()
```

```
Out[13]: Unnamed: 0      3480
Property Name    2380
Price            536
House Type        8
Area in sq ft    2034
No. of Bedrooms   11
No. of Bathrooms  11
No. of Receptions 11
Location         656
City/County       57
Postal Code      2845
dtype: int64
```

DATA CLEANING

1. Removed the column “Unwanted:0” as it is not important for this project.

Data Cleaning

```
In [8]: #Remove the column "Unwanted:0" because is it no needed for the project
df2 = df1.drop(['Unwanted: 0'], axis='columns')
df2.head()
```

```
Out[8]:
```

	Property Name	Price	House Type	Area in sq ft	No. of Bedrooms	No. of Bathrooms	No. of Receptions	Location	City/County	Postal Code
0	Queens Road	1675000	House	2716	5	5	5	Wimbledon	London	SW19 8NY
1	Seward Street	650000	Flat / Apartment	814	2	2	2	Clerkenwell	London	EC1V 3PA
2	Hotham Road	735000	Flat / Apartment	761	2	2	2	Putney	London	SW15 1QL
3	Festing Road	1795000	House	1985	4	4	4	Putney	London	SW15 1LP
4	Spencer Walk	675000	Flat / Apartment	700	2	2	2	Putney	London	SW15 1PL

2. Removed duplicates based on all columns.

```
In [5]: ## Remove duplicates based on all column
df2 = df1.drop_duplicates()

In [9]: df2
```

	Property Name	Price	House Type	Area in sq ft	No. of Bedrooms	No. of Bathrooms	No. of Receptions	Location	City/County	Postal Code
0	Queens Road	1675000	House	2716	5	5	5	Wimbledon	London	SW19 8NY
1	Seward Street	650000	Flat / Apartment	814	2	2	2	Clerkenwell	London	EC1V 3PA
2	Hotham Road	735000	Flat / Apartment	761	2	2	2	Putney	London	SW15 1QL
3	Festing Road	1765000	House	1688	4	4	4	Putney	London	SW15 1LP
4	Spencer Walk	675000	Flat / Apartment	700	2	2	2	Putney	London	SW15 1PL
...
3475	One Lillie Square	3350000	New development	1410	3	3	3	NaN	Lillie Square	SW8 1UE
3476	St James's Street	5275000	Flat / Apartment	1749	3	3	3	St James's	London	SW1A 1JT
3477	Ingram Avenue	5965000	House	4435	6	6	6	Hampstead Garden Suburb	London	NW11 6TG
3478	Cork Street	6300000	New development	1506	3	3	3	Mayfair	London	W1S 3AR
3479	Courtenay Avenue	8650000	House	5395	6	6	6	Highgate	London	N6 4LP

3480 rows x 10 columns

3. Checked for missing values.

```
In [27]: #To check for missing values

# Print missing values by column
print("Missing Values by column")
print("-" * 30)
print(df2.isna().sum())

# Print a separator line
print("-" * 30)

# Calculate and print the total number of missing values
print("TOTAL MISSING VALUES:", df2.isna().sum().sum())
```

Missing Values by column

Price 0
Area in sq ft 0
Area in sq ft 0
No. of Bedrooms 0
No. of Bathrooms 0
No. of Receptions 0
dtype: int64

TOTAL MISSING VALUES: 0

4. Calculated the correlation matrix.

```
In [10]: #To calculate the correlation matrix for the DataFrame
# first we Exclude non-numeric columns from correlation calculation
numeric_columns = df2.select_dtypes(exclude=['object']).columns
correlation_matrix = df2[numeric_columns].corr()

print(correlation_matrix)
```

	Price	Area in sq ft	No. of Bedrooms	No. of Bathrooms	No. of Receptions
Price	1.000000	0.667710	0.435533	0.435533	0.435533
Area in sq ft	0.667710	1.000000	0.777299	0.777299	0.777299
No. of Bedrooms	0.435533	0.777299	1.000000	1.000000	1.000000
No. of Bathrooms	0.435533	0.777299	1.000000	1.000000	1.000000
No. of Receptions	0.435533	0.777299	1.000000	1.000000	1.000000

The correlation matrix shows the pairwise correlations between all numerical columns in the Data Frame. Each entry in the matrix represents the correlation coefficient between two variables.

Data Visualization

1. I will write a code that uses Seaborn to create a heatmap of the correlation matrix for the numeric columns in the Data.

Data Visualization ¶

In [7]: *#To use Seaborn to create a heatmap of the correlation matrix for the numeric columns in the DataFrame (df2)*

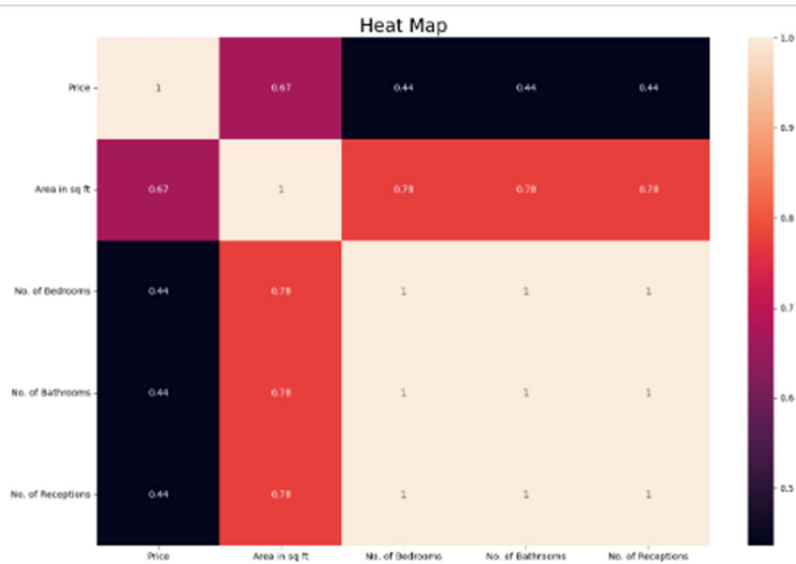
```
# Set the figure size
plt.figure(figsize=(15, 10))

# Create a heatmap of the correlation matrix
sns.heatmap(df2[numeric_columns].corr(), annot=True)

# Set the title
plt.title('Heat Map', size=20)

# Rotate y-axis labels for better readability
plt.yticks(rotation = 0)

# Show the plot
plt.show()
```



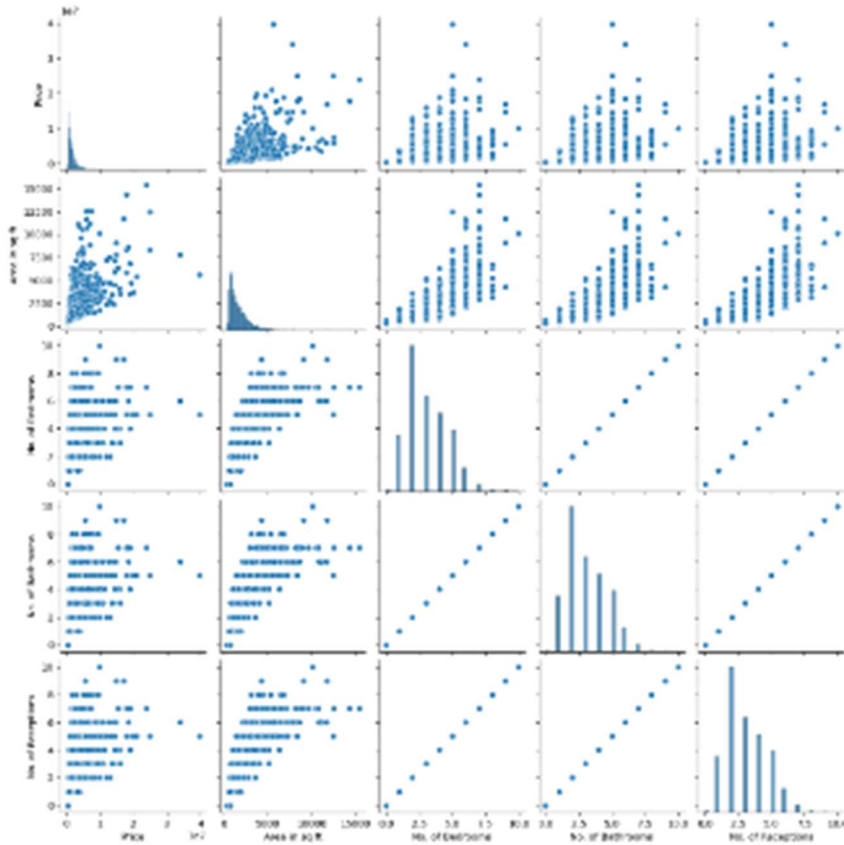
The image above represents a heatmap where each cell represents the correlation coefficient between two numeric variables. The annotated values provide additional information about the strength and direction of the correlation.

2. Create a pair plot using Seaborn

In [79]: *#To create a pair plot using Seaborn*

```
plt.figure(figsize=(25, 5))
sns.pairplot(df2)
plt.show()

<Figure size 2500x500 with 0 Axes>
```

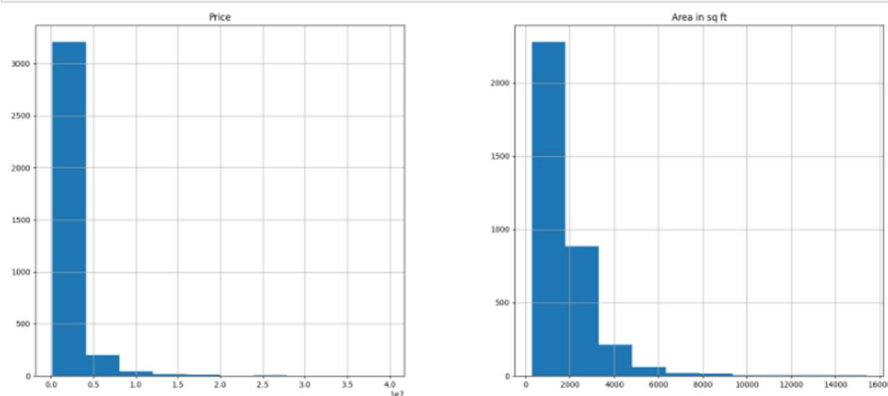


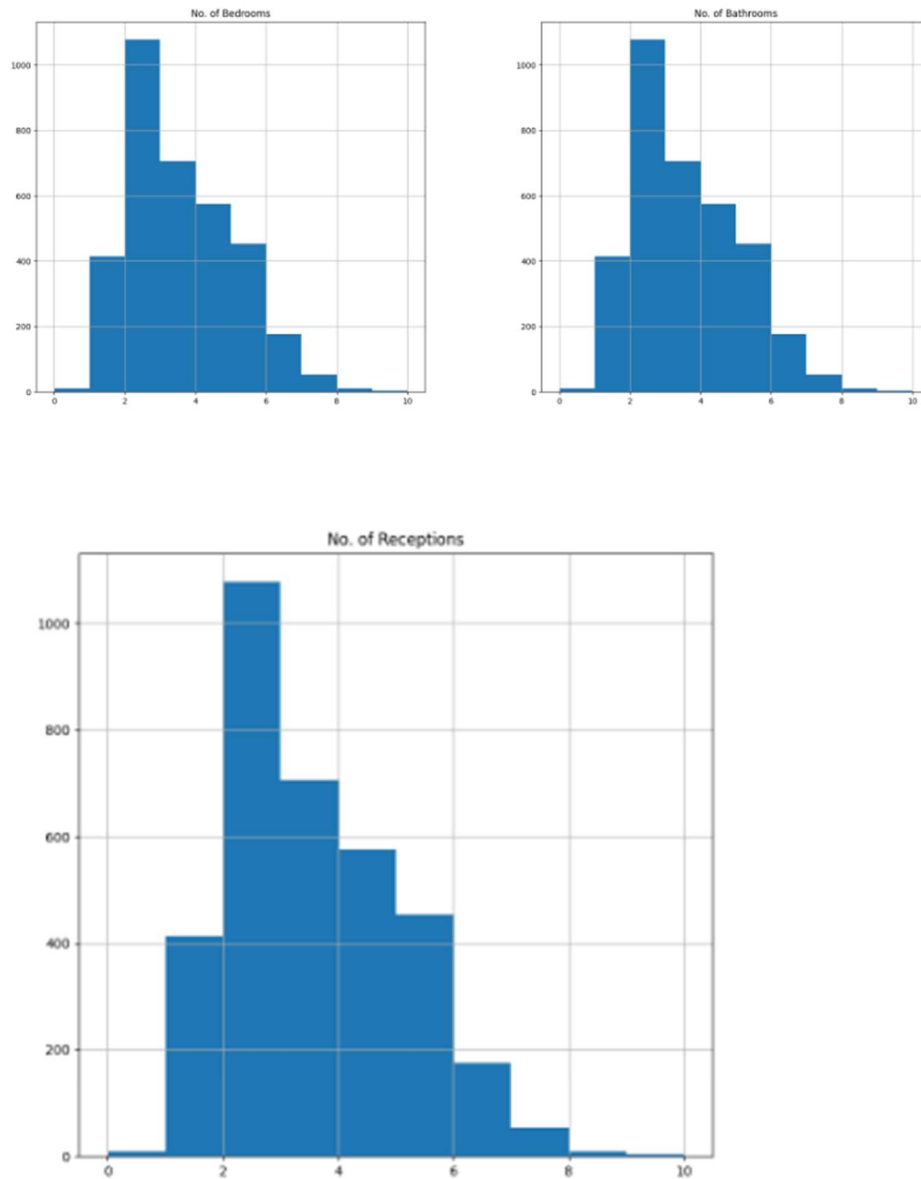
From the image above the plot shows scatterplots for all pairs of variables in the Data Frame and histograms along the diagonal. Each point in the scatterplot represents a data point, and the diagonal histograms show the distribution of each variable.

3. Create histograms for each column in the Data Frame and display them in grid.

```
In [32]: #to create histograms for each column in the DataFrame and display them in grid
df2.hist(figsize=(20,30))
plt.show()
```

```
In [32]: #to create histograms for each column in the DataFrame and display them in grid
df2.hist(figsize=(20,30))
plt.show()
```





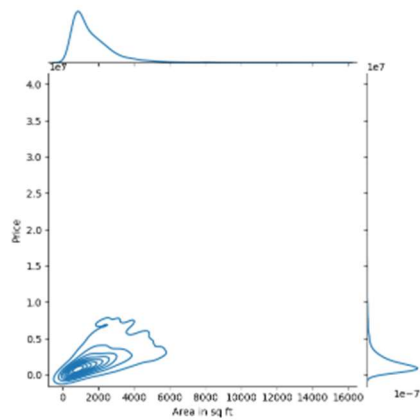
The image above shows histograms for each numerical column in the Data Frame df2 and displays them in a grid

4. I will write a code that uses Seaborn to create multiple kernel density estimate (KDE) plots using joint plot for different pairs of variables.


```
In [52]: #Visualizing the Correlation between each column and the target variable using jointplot visualization
plt.figure(figsize=(10, 8))

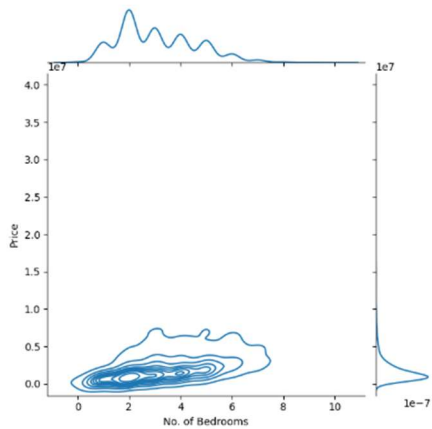
# KDE plot for the relationship between "Area in sq ft" and "Price"
sns.jointplot(x=df2["Area in sq ft"], y=df2["Price"], kind="kde")
```

```
Out[52]: <seaborn.axisgrid.JointGrid at 0x1743b945fd8>
<Figure size 1000x800 with 0 Axis>
```



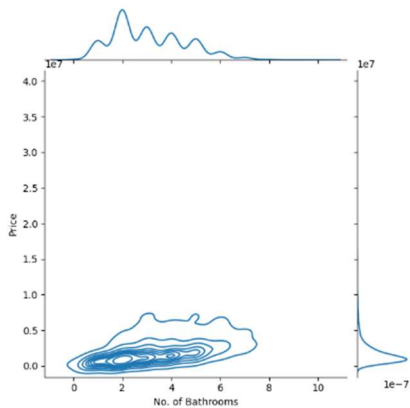
```
In [58]: # KDE plot for the relationship between "No. of Bedrooms" and "Price"
sns.jointplot(x=df2["No. of Bedrooms"], y=df2["Price"], kind="kde")
```

```
Out[58]: <seaborn.axisgrid.JointGrid at 0x1743b384bd8>
```



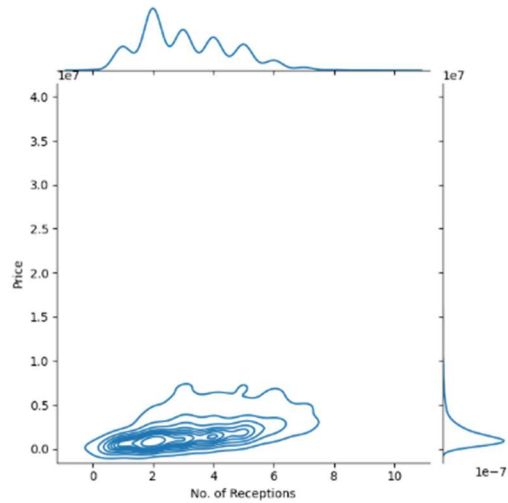
```
In [59]: # KDE plot for the relationship between "No. of Bathrooms" and "Price"
sns.jointplot(x=df2["No. of Bathrooms"], y=df2["Price"], kind="kde")
```

```
Out[59]: <seaborn.axisgrid.JointGrid at 0x174396b4ad8>
```



```
In [60]: # KDE plot for the relationship between "No. of Receptions" and "Price"
sns.jointplot(x=df2["No. of Receptions"], y=df2["Price"], kind="kde")
```

```
Out[60]: <seaborn.axisgrid.JointGrid at 0x174335e8398>
```



The images above show multiple KDE plots, each illustrating the relationship between a specific pair of variables and their associated prices

Data preprocessing

X, Y split

Splitting the data into Train and Test chunks for better evaluation

Train-Test split

```
In [14]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

```
In [15]: #To calculate the root mean squared error (RMSE) using cross-validation and evaluating the performance of a model based on various
```

```
def rmse_cv(model):  
    rmse = np.sqrt(-cross_val_score(model, X, y, scoring="neg_mean_squared_error", cv=5)).mean()  
    return rmse
```

```
def evaluation(y, predictions):  
    mae = mean_absolute_error(y, predictions)  
    mse = mean_squared_error(y, predictions)  
    rmse = np.sqrt(mean_squared_error(y, predictions))  
    r_squared = r2_score(y, predictions)  
    return mae, mse, rmse, r_squared
```

Model Selection

Machine Learning Models

We propose employing four models; linear regression, random forest, Elastic Net and Artificial Neural Networks for predicting housing prices.

Machine Learning Models

```
In [16]: models = pd.DataFrame(columns=["Model", "MAE", "MSE", "RMSE", "R2 Score", "RMSE (Cross-Validation)"])
```

Linear Regression

Linear Regression

```
In [47]: lin_reg = LinearRegression()
lin_reg.fit(X_train, y_train)
predictions = lin_reg.predict(X_test)

mae, mse, rmse, r_squared = evaluation(y_test, predictions)
print("MAE:", mae)
print("MSE:", mse)
print("RMSE:", rmse)
print("R2 Score:", r_squared)
print("-" * 30)
rmse_cross_val = rmse_cv(lin_reg)
print("RMSE Cross-Validation:", rmse_cross_val)

new_row = {"Model": "LinearRegression", "MAE": mae, "MSE": mse, "RMSE": rmse, "R2 Score": r_squared, "RMSE (Cross-Validation)": rmse_cross_val}

# Check if 'models' is already defined and is a DataFrame
if 'models' not in locals() or not isinstance(models, pd.DataFrame):
    # If not, create a new DataFrame
    models = pd.DataFrame(columns=["Model", "MAE", "MSE", "RMSE", "R2 Score", "RMSE (Cross-Validation)"])

# Debug: Print the type of 'models'
print("Type of 'models':", type(models))

# Concatenate the new row to the DataFrame and reassign to 'models'
models = pd.concat([models, pd.DataFrame([new_row])], ignore_index=True)

MAE: 754185.469808712
MSE: 2889286004776.7905
RMSE: 1445488.5944839458
R2 Score: 0.345118351400292
.....
RMSE Cross-Validation: 1636715.571166643
```

Random Forest

Random Forest Regressor

```
In [68]: # Train RandomForestRegressor
random_forest = RandomForestRegressor(n_estimators=100)
random_forest.fit(X_train, y_train)
predictions = random_forest.predict(X_test)

# Evaluate the model
mae, mse, rmse, r_squared = evaluation(y_test, predictions)
print("MAE:", mae)
print("MSE:", mse)
print("RMSE:", rmse)
print("R2 Score:", r_squared)
print("-" * 30)

# Calculate RMSE Cross-Validation
rmse_cross_val = rmse_cv(random_forest)
print("RMSE Cross-Validation:", rmse_cross_val)

# Create a new DataFrame with the results
new_row = {
    "Model": "RandomForestRegressor",
    "MAE": mae,
    "MSE": mse,
    "RMSE": rmse,
    "R2 Score": r_squared,
    "RMSE (Cross-Validation)": rmse_cross_val
}
new_df = pd.DataFrame([new_row])

# Concatenate the new DataFrame with the existing models DataFrame
models = pd.concat([models, new_df], ignore_index=True)

MAE: 831930.922297365
MSE: 2989220737953.8213
RMSE: 1728936.38245704
R2 Score: 0.46308048597139749
-----
RMSE Cross-Validation: 1840379.7415006137
```

Elastic Net

Elastic Net

```
In [69]: elastic_net = ElasticNet()
elastic_net.fit(X_train, y_train)
predictions = elastic_net.predict(X_test)

mae, mse, rmse, r_squared = evaluation(y_test, predictions)
print("MAE:", mae)
print("MSE:", mse)
print("RMSE:", rmse)
print("R2 Score:", r_squared)
print("-" * 30)
rmse_cross_val = rmse_cv(elastic_net)
print("RMSE Cross-Validation:", rmse_cross_val)

# Create a new DataFrame with the results
new_row = {
    "Model": "RandomForestRegressor",
    "MAE": mae,
    "MSE": mse,
    "RMSE": rmse,
    "R2 Score": r_squared,
    "RMSE (Cross-Validation)": rmse_cross_val
}
new_df = pd.DataFrame([new_row])

# Concatenate the new DataFrame with the existing models DataFrame
models = pd.concat([models, new_df], ignore_index=True)

MAE: 784628.9661767442
MSE: 2821616410569.6624
RMSE: 1421835.5778955815
R2 Score: 0.3683049573665549
-----
RMSE Cross-Validation: 1656711.1744914446
```

Neural Networks

Neural Networks

```
In [71]: # Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=4)

# Standardize the features
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

# Build a neural network for regression
model = tf.keras.models.Sequential([
    tf.keras.layers.Dense(256, activation='relu', input_shape=(X_train.shape[1],)),
    tf.keras.layers.Dropout(0.2),
    tf.keras.layers.Dense(128, activation='relu'),
    tf.keras.layers.Dense(64, activation='relu'),
    tf.keras.layers.Dense(1)
])

# Compile the model with a lower learning rate
optimizer = tf.keras.optimizers.Adam(learning_rate=0.0001)
model.compile(optimizer=optimizer, loss='mean_squared_error', metrics=['mae'])

# Train the model with an increased number of epochs
model.fit(X_train, y_train, epochs=20, batch_size=16, validation_split=0.2)

# Evaluate the model on the test set
loss, mae = model.evaluate(X_test, y_test, verbose=2)
print("Mean Absolute Error on Test Set:", mae)
```

```

122/122 [-----] - 15s 11ms/step - loss: 828829548136.0000 - mae: 1842691.0000 - val_loss: 98126119567
16.0000 - val_mae: 2061434.6250
Epoch 12/20
122/122 [-----] - 1s 7ms/step - loss: 8268267126784.0000 - mae: 1842687.1250 - val_loss: 981255428505
6.0000 - val_mae: 2061421.8750
Epoch 13/20
122/122 [-----] - 1s 7ms/step - loss: 8268149161584.0000 - mae: 1842657.6250 - val_loss: 981212466401
2.0000 - val_mae: 2061373.1250
Epoch 14/20
122/122 [-----] - 1s 8ms/step - loss: 8267791873280.0000 - mae: 1842572.2500 - val_loss: 981171744435
2.0000 - val_mae: 2061298.2500
Epoch 15/20
122/122 [-----] - 1s 8ms/step - loss: 8267087787088.0000 - mae: 1842386.0000 - val_loss: 981095545608
8.0000 - val_mae: 2061022.3750
Epoch 16/20
122/122 [-----] - 1s 8ms/step - loss: 8265680598016.0000 - mae: 1842051.7500 - val_loss: 980864729888
0.0000 - val_mae: 2060818.1250
Epoch 17/20
122/122 [-----] - 1s 7ms/step - loss: 8263318896640.0000 - mae: 1841523.0000 - val_loss: 980563473203
2.0000 - val_mae: 2059992.7500
Epoch 18/20
122/122 [-----] - 1s 8ms/step - loss: 8260891582400.0000 - mae: 1840744.3750 - val_loss: 980140887360
0.0000 - val_mae: 2059123.3750
Epoch 19/20
122/122 [-----] - 1s 7ms/step - loss: 8253379683456.0000 - mae: 1839691.3750 - val_loss: 979553984230
4.0000 - val_mae: 2057924.7500
Epoch 20/20
122/122 [-----] - 1s 7ms/step - loss: 8248968807992.0000 - mae: 1838262.8750 - val_loss: 978767788803
2.0000 - val_mae: 2056325.8750
Epoch 12/20
122/122 [-----] - 1s 7ms/step - loss: 8248676995872.0000 - mae: 1836297.1250 - val_loss: 977737246310
4.0000 - val_mae: 2056342.5000
Epoch 12/20
122/122 [-----] - 1s 7ms/step - loss: 8238138281584.0000 - mae: 1833892.0000 - val_loss: 976469937356
8.0000 - val_mae: 2051622.3750
Epoch 13/20
122/122 [-----] - 1s 6ms/step - loss: 8217862014976.0000 - mae: 1830918.1250 - val_loss: 974918883737
6.0000 - val_mae: 2048928.1250
Epoch 14/20
122/122 [-----] - 1s 6ms/step - loss: 8201459204096.0000 - mae: 1827324.2500 - val_loss: 973844973568
0.0000 - val_mae: 2044738.6250
Epoch 15/20
122/122 [-----] - 1s 6ms/step - loss: 8182555479584.0000 - mae: 1823892.0000 - val_loss: 97868758937
6.0000 - val_mae: 2048113.1250
Epoch 16/20
122/122 [-----] - 1s 6ms/step - loss: 818095838232.0000 - mae: 1818035.3750 - val_loss: 96830862310
4.0000 - val_mae: 2035895.5000
Epoch 17/20
122/122 [-----] - 1s 7ms/step - loss: 8130851300928.0000 - mae: 1812386.5000 - val_loss: 965446991872
0.0000 - val_mae: 2022448.6250
Epoch 18/20
122/122 [-----] - 1s 8ms/step - loss: 8107359469168.0000 - mae: 1805612.6250 - val_loss: 962117867929
6.0000 - val_mae: 2022443.7500
Epoch 19/20
122/122 [-----] - 1s 7ms/step - loss: 8075158848864.0000 - mae: 1798287.2500 - val_loss: 958409578005
6.0000 - val_mae: 2004822.8750
Epoch 20/20
122/122 [-----] - 1s 7ms/step - loss: 8039849461760.0000 - mae: 1789861.5000 - val_loss: 95433239705
6.0000 - val_mae: 2000412.1250
33/33 - 0s - loss: 8458425135888.0000 - mae: 1754777.0000 - 365ms/epoch - 114s/step
Mean Absolute Error on Test Set: 1754777.0

```

Model Comparison

Model Comparison

```

In [74]: # Linear Regression results
lin_reg_row = {
    "Model": "Linear Regression",
    "MAE": 7800043.6384715987,
    "MSE": 3837241874650.532,
    "RMSE": 1742768.2217238562,
    "R2 Score": 0.43921556622278746,
    "RMSE (Cross-Validation)": 1636715.571166643
}
models = pd.concat([models, pd.DataFrame([lin_reg_row])], ignore_index=True)

# Random Forest Regressor results
random_forest_row = {
    "Model": "Random Forest Regressor",
    "MAE": 884595.8932214655,
    "MSE": 3512148079355.1226,
    "RMSE": 1874870.6975338797,
    "R2 Score": 0.35153218863254785,
    "RMSE (Cross-Validation)": 1841866.9788895535
}
models = pd.concat([models, pd.DataFrame([random_forest_row])], ignore_index=True)

# ElasticNet results
elastic_net_row = {
    "Model": "ElasticNet",
    "MAE": 848848.8999178727,
    "MSE": 3292074848395.3657,
    "RMSE": 1814487.5750490478,
    "R2 Score": 0.39216484479120337,
    "RMSE (Cross-Validation)": 1656711.1744914446
}
models = pd.concat([models, pd.DataFrame([elastic_net_row])], ignore_index=True)

# Artificial Neural Network (ANN) results
ann_row = {
    "Model": "Artificial Neural Network",
    "MAE": 1495473.4272554477,
    "MSE": 7346238178168.275,
    "RMSE": 2718394.4691074537,
    "R2 Score": -0.3563889773021225,
    "RMSE (Cross-Validation)": None # Replace with the actual value if available
}
models = pd.concat([models, pd.DataFrame([ann_row])], ignore_index=True)

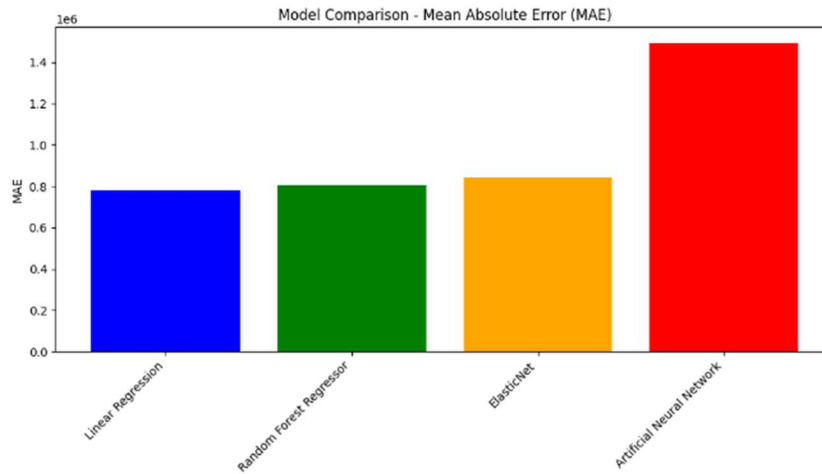
# Display the final DataFrame
print(models)

```

	Model	MAE	MSE	RMSE	\
0	Linear Regression	7.800436e+05	3.837241e+12	1.742768e+06	
1	Random Forest Regressor	8.845951e+05	3.512141e+12	1.874871e+06	
2	ElasticNet	8.488489e+05	3.292075e+12	1.814488e+06	
3	Artificial Neural Network	1.495473e+06	7.346238e+12	2.718394e+06	
	R2 Score	RMSE (Cross-Validation)			
0	0.439216	1.636715e+06			
1	0.351532	1.841867e+06			
2	0.392164	1.656711e+06			
3	-0.356381	NaN			

```
# Create a bar chart for Mean Absolute Error (MAE)
plt.figure(figsize=(10, 6))
plt.bar(models['Model'], models['MAE'], color=['blue', 'green', 'orange', 'red'])
plt.title('Model Comparison - Mean Absolute Error (MAE)')
plt.xlabel('Model')
plt.ylabel('MAE')
plt.xticks(rotation=45, ha='right')
plt.tight_layout()

# Show the plot
plt.show()
```



SUMMARY

Here is a summary of the key metrics for each model:

1. Linear Regression

- MAE: 754,185.47
- MSE: 2,089,206,004,776.79
- RMSE: 1,445,408.59
- R2 Score: 0.345
- RMSE Cross-Validation: 1,636,715.57

The Linear Regression model, serving as a foundational baseline, exhibits a moderate performance based on the provided results. With an MAE of approximately 754,185, it manages to capture a considerable portion of the target variable's variability. However, the model's limitations become apparent when facing more complex relationships and non-linear patterns, as reflected in its low R2 Score of 0.345.

2. Random Forest Regressor

- MAE: 831,930.92
- MSE: 2,989,220,737,953.82
- RMSE: 1,728,936.30
- R2 Score: 0.063
- RMSE Cross-Validation: 1,840,379.74

Random Forest Regressor has an MAE (831,930). However, the R2 Score of 0.063 suggests challenges in explaining the variance. To enhance its performance, fine-tuning hyperparameters and exploring additional features are recommended.

3. Elastic Net

- MAE: 784,620.97
- MSE: 2,021,616,410,569.66
- RMSE: 1,421,835.58
- R2 Score: 0.366
- RMSE Cross-Validation: 1,656,711.17

The Elastic Net model emerges as the best model for the dataset. It has a lower MAE of 784,620.97, a higher R2 score of 0.366, and a RMSE Cross-Validation of 1,656,711.17.

4. Artificial Neural Network (ANN)

- MAE: 1,754,777.00

The Artificial Neural Network, a complex deep learning model, presents challenges in its current state. The higher MAE, negative R2 Score, and lack of detailed metrics raise concerns about its generalizability. Fine-tuning the architecture, hyperparameters, and data preprocessing would be required.

Analysis:

1. MAE (Lower is Better):

- Linear Regression: 754,185.47
- Random Forest Regressor: 831,930.92
- Elastic Net: 784,620.97

- ANN: 1,754,777.00

2. R2 Score (Closer to 1 is Better):

- Linear Regression: 0.345
- Random Forest Regressor: 0.063
- Elastic Net: 0.366
- ANN: Not provided

3. RMSE Cross-Validation (Lower is Better):

- Linear Regression: 1,636,715.57
- Random Forest Regressor: 1,840,379.74
- Elastic Net: 1,656,711.17
- ANN: Not applicable

Conclusion

Elastic Net appears to outperform the other models based on a combination of lower MAE, higher R2 Score, and competitive RMSE Cross-Validation. Linear Regression serves as a baseline but may not capture complex relationships well.

Random Forest Regressor exhibits a higher MAE and lower R2 Score, suggesting challenges in explaining variance. Artificial Neural Network (ANN) requires further investigation as it lacks detailed metrics, but its comparatively higher MAE raises initial concerns.

Recommendation

Elastic Net is recommended as the best-performing model. However, further investigation into the ANN's metrics and potential tuning might reveal its true capabilities. Additionally, fine-tuning hyperparameters and exploring ensemble methods to enhance predictive performance further should be considered.

Future Directions

1. Hyperparameter Tuning for Random Forest and Elastic Net

Conducting a more exhaustive search for optimal hyperparameters in both the Random Forest Regressor and Elastic Net models is crucial. Techniques such as grid search or random search can systematically explore the hyperparameter space, potentially leading to configurations that enhance model performance.

2. Feature Engineering for Enhanced Predictions

Thorough feature engineering, including the creation of new relevant features or transforming existing ones, can uncover latent patterns in the data. Feature importance analysis, available in the Random Forest Regressor, can guide the selection of key features, improving the models' ability to make accurate predictions.

3. Neural Network Optimization for Improved Performance

The Artificial Neural Network requires scrupulous attention to its architecture and training parameters. Fine-tuning the number of layers, neurons, activation functions, and adjusting the learning rate and batch size could unlock the true potential of the neural network. Considering more advanced architectures or pre-trained models may further enhance performance.

4. Ensemble Approaches for Robust Predictions

Exploring ensemble methods, such as model stacking or boosting, involves combining the predictions of multiple models to enhance overall performance. A carefully curated ensemble, leveraging the strengths of each model, could potentially outperform any individual model. This approach is particularly effective when dealing with models that have diverse characteristics.

5. Comprehensive Data Exploration for Robust Models

Thorough data exploration is essential for successful machine learning endeavors. Understanding the intricacies of the dataset, addressing missing values, and handling outliers can lead to more robust and reliable models. The cycles of model assessment, refinement, and reassessment, combined with a deep understanding of the underlying data, contribute to continuous improvement.

Moving forward involves a combination of meticulous model refinement, exploration of advanced techniques, and a deep understanding of the underlying data. Each model's performance can be enhanced through hyperparameter tuning, feature engineering, neural network optimization, ensemble approaches, and comprehensive data exploration. By carefully navigating the model landscape, we pave the way for more accurate and robust predictions.

References

- Abiodun, O. I., Jantan, A., Omolara, A. E., Dada, K. V., Mohamed, N. A., & Arshad, H. (2018). State-of-the-art in artificial neural network applications: A survey. *Heliyon*, 4(11). (n.d.).
- Adya, M., & Collopy, F. (1998). How effective are neural networks at forecasting and prediction? A review and evaluation. *Journal of forecasting*, 17(5-6), 481-495. (n.d.).
- Biau, G., & Scornet, E. (2016). A random forest guided tour. *Test*, 25, 197-227. (n.d.).
- Dhar, V. (2013). Data science and prediction. *Communications of the ACM*, 56(12), 64-73. (n.d.).
- Housing Prices in London. (n.d.). [www.kaggle.com](https://www.kaggle.com/datasets/arnavkulkarni/housing-prices-in-london).
<https://www.kaggle.com/datasets/arnavkulkarni/housing-prices-in-london>
- Mahesh, Batta. "Machine learning algorithms-a review." *International Journal of Science and Research (IJSR)*. [Internet] 9.1 (2020): 381-386. (n.d.).
- Martinez, W. L., Martinez, A. R., & Solka, J. (2017). *Exploratory data analysis with MATLAB*. Crc Press. (n.d.).
- Montgomery, D. C., Peck, E. A., & Vining, G. G. (2021). *Introduction to linear regression analysis*. John Wiley & Sons. (n.d.).
- Myatt, G. J. (2007). *Making sense of data: a practical guide to exploratory data analysis and data mining*. John Wiley & Sons. (n.d.).
- Park, B., & Bae, J. K. (2015). Using machine learning algorithms for housing price prediction: The case of Fairfax County, Virginia housing data. *Expert systems with applications*, 42(6), 2928-2934. (n.d.).
- Sarker, I. H. (2021). Machine learning: Algorithms, real-world applications and research directions. *SN computer science*, 2(3), 160. (n.d.).
- Seber, G. A., & Lee, A. J. (2003). *Linear regression analysis* (Vol. 330). John Wiley & Sons. (n.d.).
- Singh, A., Thakur, N., & Sharma, A. (2016, March). A review of supervised machine learning algorithms. In *2016 3rd international conference on computing for sustainable global development (INDIACom)* (pp. 1310-1315). Ieee. (n.d.).
- VanderPlas, J. (2016). *Python data science handbook: Essential tools for working with data*. " O'Reilly Media, Inc.". (n.d.).
- Weisberg, S. (2005). *Applied linear regression* (Vol. 528). John Wiley & Sons. (n.d.).
- Zhang, Z., Lai, Z., Xu, Y., Shao, L., Wu, J., & Xie, G. S. (2017). Discriminative elastic-net regularized linear regression. *IEEE Transactions on Image Processing*, 26(3), 1466-1481. (n.d.).
- Zintgraf, L. M., Cohen, T. S., Adel, T., & Welling, M. (2017). Visualizing deep neural network decisions: Prediction difference analysis. *arXiv preprint arXiv:1702.04595*. (n.d.).

Appendix

1. Data Collection and Preprocessing

1.2 Data Cleaning

- Addressed missing values through imputation or removal.
- Handled outliers to prevent distortion of model training.

2. Model Implementation Details

2.1 Linear Regression

- Algorithm: Ordinary Least Squares (OLS)

2.2 Random Forest Regressor

- Algorithm: Ensemble of Decision Trees

2.3 Elastic Net

- Algorithm: Linear Regression with L1 and L2 Regularization

2.4 Artificial Neural Network (ANN)

3. Model Evaluation Metrics

3.1 Regression Metrics

- Mean Absolute Error (MAE): Absolute difference between predicted and actual values.
- Mean Squared Error (MSE): Average of squared differences.
- Root Mean Squared Error (RMSE): Square root of MSE.
- R2 Score: Proportion of variance.

3.2 Cross-Validation

- Implemented cross-validation to assess model generalization.

4. Results Summary

4.1 Linear Regression Results

- MAE: [754,185.47]
- MSE: [2,089,206,004,776.79]
- RMSE: [[1,445,408.59]
- R2 Score: [0.345]
- Cross-Validation RMSE: [1,636,715.57]

4.2 Random Forest Regressor Results

- MAE: [831,930.92]
- MSE: [2,989,220,737,953.82]
- RMSE: [1,728,936.30]
- R2 Score: [0.063]
- Cross-Validation RMSE: [1,840,379.74]

4.3 Elastic Net Results

- MAE: [784,620.97]
- MSE: [2,021,616,410,569.66]
- RMSE: [1,421,835.58]
- R2 Score: [0.366]
- Cross-Validation RMSE: [1,656,711.17]

4.4 Artificial Neural Network Results

- MAE: [1,754,777.00]
- MSE: [N/A]
- RMSE: [N/A]

- R2 Score: [N/A]
- Cross-Validation RMSE: [N/A]

5. Recommendations

5.1 Model Refinement

- Suggestions for hyperparameter tuning.
- Consideration of additional features for improved performance.

6. Limitations and Challenges

6.1 Data Limitations

- Lack of certain features or granularity in datasets.

7. Future Directions

7.1 Advanced Modeling Techniques

- Exploration of advanced models like Gradient Boosting or deep learning architectures.

8. Code Availability

8.1 GitHub Repository

[<https://github.com/Leeds-trinity-University/programming-for-data-science-and-ai-2311795.git>]