# Capstone Project Summary Report

Sanchaita Biswas

# Capstone Project - 1: Predicting Property Prices in a Specific Location Using Machine Learning
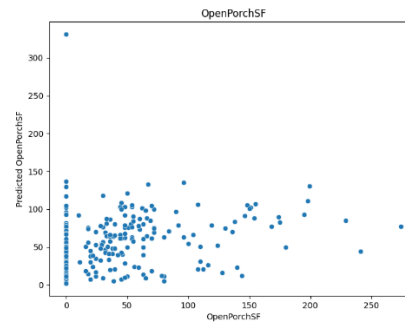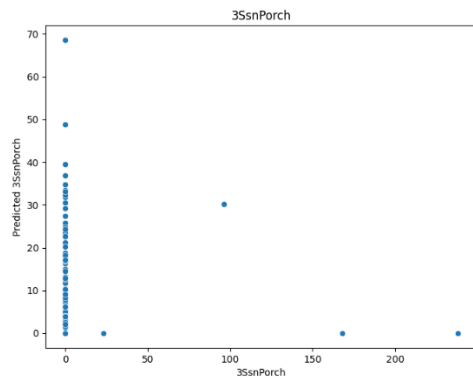
## 1. Introduction to the Project:

The real estate market is influenced by numerous factors, making property price prediction a complex yet crucial task for various stakeholders. This capstone project uses machine learning techniques to predict property prices in a specific location. The project aims to provide valuable insights for real estate agents, buyers, and sellers by utilizing advanced data analysis and modeling.

## 2. Objectives of the Project:

1. Collect and clean real estate data from a specific location.

2. Handle ordinal and nominal columns separately to enhance prediction accuracy.

3. Implement scaling, PCA, and fillna() techniques to handle missing data effectively.

4. Perform exploratory data analysis (EDA) to identify key variables influencing property prices.

5. Determine appropriate encoding techniques for ordinal and nominal variables based on model requirements.

6. Develop a machine learning model capable of accurately predicting property prices.

7. Evaluate model performance and compare with alternative algorithms.

8. Present project findings and insights clearly and concisely.

# 3. Flow Chart of Operations:





# 4. Python Codes:

```python
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score

# Load the dataset from Excel file
data = pd.read_excel('D:/DigiCrome/Project1/Datafields/NewModified_dataset.xlsx', sheet_name='Sheet1')

# Data Preprocessing - Handle missing values
data.dropna(inplace=True)  # Drop rows with missing values

# Assuming 'PropertyClass' is the target variable and other columns are features
X = data.drop('PropertyClass', axis=1)
y = data['PropertyClass']

# Split data into train and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Train a Random Forest Classifier
model = RandomForestClassifier()
model.fit(X_train, y_train)

# Model Evaluation
y_pred = model.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy)
```

```python
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.ensemble import RandomForestRegressor
from sklearn.metrics import accuracy_score
from sklearn.metrics import mean_squared_error

# Load the dataset
data = pd.read_csv("D:/DigiCrome/Project1/Property_data.csv")

# EDA (Explore the data)
print(data.head())
print(data.info())
print(data.describe())

# Check data types of all columns
print("Data types of all columns:")
print(data.dtypes)

# Load the new data set
dataOne = pd.read_excel("D:/DigiCrome/Project1/NewModified_dataset.xlsx")

X = dataOne.drop("PropertySize", axis=1)
y = dataOne["PropertySize"]
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Choose a model
model = RandomForestClassifier()

# Train the model
model.fit(X_train, y_train)

# Step 6: Model evaluation
y_pred = model.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy)
```
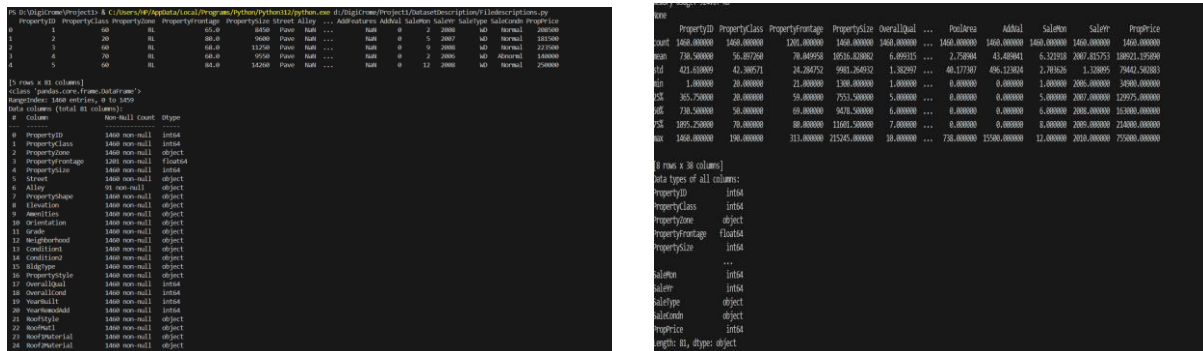
# 5. Screenshot of the Outputs:



# 6. Report on EDA:

Exploratory Data Analysis (EDA) is crucial in understanding the underlying patterns and relationships within the dataset, providing valuable insights into the factors influencing property prices. This section presents the key findings from our EDA process, accompanied by graphs and visualizations.

## 1. Distribution of Property Prices:

   - **Histogram**: The distribution of property prices reveals whether they are normally distributed or skewed.

   - **Boxplot:** Identifying outliers and understanding the spread of property prices across different quartiles.

## 2. Correlation Analysis:

   - **Heatmap:** Visualizing the correlation matrix to identify features strongly correlated with property prices.

- **Scatterplots:** Exploring relationships between property prices and other numerical variables, such as property size, number of bedrooms, and overall quality.

## 3. Categorical Variables Analysis:

 - **Bar Charts:** Examining the distribution of categorical variables like property class, zoning classification, and neighborhood about property prices.

 - **Boxplots:** Comparing property prices across different categories of categorical variables to identify any significant differences.

## 4. Time Trends Analysis:

 - **Line Plot:** Analyzing trends in property prices over time (years old) to understand market dynamics and seasonality effects.

## 5. Feature Engineering Insights:

 - **Derived Features:** Exploring newly created features based on domain knowledge or statistical techniques and their relationship with property prices.

# Key Insights:

1. **Property Size:** Larger properties tend to have higher prices, indicating a positive correlation between property size and price.

2. **Neighbourhood**: Certain neighborhoods exhibit higher property prices compared to others, highlighting the importance of location in determining property values.

3. **Overall Quality**: Properties with higher overall quality ratings command higher prices, emphasizing the significance of property condition.

4. **Time Trends**: Property prices may exhibit temporal trends, with fluctuations influenced by economic factors and market conditions.

5. **Categorical Variables**: Certain property classes, zoning classifications, and amenities may significantly impact property prices, warranting further investigation.

# Recommendations:

1. Focus on properties in neighbourhoods with historically high price premiums.

2. Prioritize properties with larger sizes and superior quality ratings for potentially higher returns.

3. Consider the impact of time trends on property prices when making investment decisions.

4. Investigate the influence of specific categorical variables on property prices to tailor marketing strategies accordingly.

In conclusion, the EDA process has provided valuable insights into the determinants of property prices, laying the groundwork for developing an accurate predictive model. Further analysis and feature selection based on these insights will enhance the model's performance and predictive capability.

# 7. Learning Outcomes:

- Proficiency in data collection, cleaning, and preprocessing techniques.

- Understanding of feature engineering methods to enhance model performance.

- Familiarity with various machine learning algorithms and their application in property price prediction.

- Skills in model evaluation and comparison to determine the best-performing approach.

- Experience in presenting project findings and insights effectively.

# 8. Conclusion:

This capstone project endeavours to develop a robust machine-learning model for predicting property prices in a specific location. By following a structured methodology encompassing data collection, cleaning, analysis, modeling, and evaluation, the project aims to deliver accurate predictions and valuable insights for real estate stakeholders.

# 9. Citations:

- Dean De Cock. (n.d.). Ames, Iowa: Alternative to the Boston Housing Data as an End of Semester Regression Project. Retrieved from http://jse.amstat.org/v19n3/decock.pdf

- Scikit-learn: Machine Learning in Python. (n.d.). Retrieved from https://scikit-learn.org/stable/