



HOUSING: PRICE PREDICTION



Submitted by:
Deepak kr. Singh

ACKNOWLEDGMENT

This includes mentioning of all the references, research papers, data sources, professionals and other resources that helped you and guided you in completion of the project.

- I would like to thank FlipRobo Technologies for providing me this opportunity and guidance throughout the project and all the steps that are implemented.
- I have primarily referred to various articles scattered across various websites for the purpose of getting an idea on Housing project.
- I would like to thank the technical support team also for helping me out and reaching out to me on clearing all my doubts as early as possible.
- I would like to thank my project SME M/S Sapna Verma for providing the flexibility in time and also for giving us guidance in creating the project.
- I have referred to various articles in Towards Data Science and Kaggle

For more please visit : [DS0003/Housing-Project: HOUSING: PRICE PREDICTION \(github.com\)](https://github.com/DS0003/Housing-Project)



INTRODUCTION

● Business Problem Framing

- Houses are one of the necessary need of each and every person around the globe and therefore housing and real estate market is one of the markets which is one of the major contributors in the world's economy. It is a very large market and there are various companies working in the domain. Data science comes as a very important tool to solve problems in the domain to help the companies increase their overall revenue, profits, improving their marketing strategies and focusing on changing trends in house sales and purchases. Predictive modelling, Market mix modelling, recommendation systems are some of the machine learning techniques used for achieving the business goals for housing companies. Our problem is related to one such housing company.
- A US-based housing company named Surprise Housing has decided to enter the Australian market. The company uses data analytics to purchase houses at a price below their actual values and flip them at a higher price. For the same purpose, the company has collected a data set from the sale of houses in Australia. The data is provided in CSV
- The company is looking at prospective properties to buy houses to enter the market. You are required to build a model using Machine Learning in order to predict the actual value of the prospective properties and decide whether to invest in them or not. For this company wants to know:
 - Which variables are important to predict the price of variable?
 - How do these variables describe the price of the house?

Business Goal:

- You are required to model the price of houses with the available independent variables. This model will then be used by the management to understand how exactly the prices vary with the variables. They can accordingly manipulate the strategy of the firm and concentrate on areas that will yield high returns. Further, the model will be a good way for the management to understand the pricing dynamics of a new market.

● Conceptual Background of the Domain Problem

Estimating the sale prices of houses is one of the basic projects to have on your Data Science CV. By finishing this article, you will be able to predict continuous variables using various types of linear regression algorithm: (Linear regression is an algorithm used to predict values that are continuous in nature. It became more popular because it is the best algorithm to start with if you are a newbie to ML.)

Technical Requirements:

- Data contains 1460 entries each having 81 variables.
- Data contains Null values. Data treatment using domain knowledge, understanding.
- Extensive EDA has to be performed to gain relationships of important variable and price.
- Data contains numerical as well as categorical variable, handle them accordingly.
- Machine Learning models, apply regularization and determine values of Hyper Parameters.

- You need to find important features which affect the price positively or negatively.
- Two datasets are being provided (test.csv, train.csv), train model on train.csv dataset and predict on test.csv file.

• Review of Literature

Housing is one out of the 3 basic needs for human survival hence research about housing and all that relates to it can never be over emphasized. In this section of this research work I will examine different papers or research work that have been published previously as regards housing and prices in the past and the research gap noticed which forms the basis for my own research work.

In this study, we will use a housing dataset presented by De Cock (2011). This dataset describes the sales of residential units in Ames, Iowa starting from 2006 until 2010. The dataset contains a large number of variables that are involved in determining a house price.

The performance of the model were evaluated using mean absolute percentage error (MAPE) performance metric. MAPE was calculated using this formula:

$$MAPE = \frac{\sum_{i=1}^r (abs(y_i - \hat{y}_i) / y_i)}{r} \times 100$$

where \hat{y}_i is the predicted stock price on day i , y_i is the actual stock price on day, i , and r is the number of trading days.

Steps :

- ✓ Importing the required packages into our python environment
- ✓ Importing the house price data and do some EDA on it
- ✓ Data Visualization on the house price data
- ✓ Feature Selection & Data Split
- ✓ Modelling the data using the algorithms
- ✓ Evaluating the built model using the evaluation metrics

Finally, we conclude which model is best suitable for the given case by evaluating each of them using the evaluation metrics provided by the scikit-learn package.

• Motivation for the Problem Undertaken

- As we know that housing and real estate market is one of the markets which is one of the major contributors in the world's economy. The model is use by the management to understand how exactly the prices vary with the variables.
- Analytics data will help company to purchase houses at a price below their actual values and flip them at a higher price. For the same purpose, the company has collected a data set from the sale of houses in Australia.
- Record of low home loan rates and deficiency of stock is helpful to keep the US real market solid, home costs have been flooding month-by-month and breaking new records.

- Cost are ascending as there is a lot of capital uninvolved, just as exceptionally modest home loan rates, as well as new gaint are also getting involved and for longer period of time home costs have been set in mid-single digits.
- As a less number house owners in different reasons; due to developing expenses and land financial backers gathering , starter houses lodging supply is by and by at its most minimal level.
- As houses are one of the necessary need of each and every person around the globe and therefore housing and real estate market is needed new people are migrating so there is requirement of new houses as per their requirements.

Analytical Problem Framing

- Mathematical/ Analytical Modeling of the Problem

- We are going to work with the house price dataset that contains various features and information about the house and its sale price. Using the `'read_csv'` function provided by the Pandas package, we can import the data into our python environment. After importing the data, we can use the `'head'` function to get a glimpse of our dataset.
- A US-based housing company named Surprise Housing has decided to enter the Australian market,
- Housing is one of the fundamental essential of every living thing hence the reason for continuous research in this sector. This project simply examines a dataset, which consists 1460 observations and 80 features that contribute to the sale price of the houses. Dataset was cleaned and transformed and some explorations were done on it to answer some basic questions that anybody would like to ask about housing. Feature engineering was performed on the transformed data using Principal Component Analysis (PCA) and encoding and this is to ensure our dataset is ready in the right form with the right variables to be used in the algorithms, which results in improved model accuracy. Different ensemble algorithms were used on the dataset in this project. The overall result of this project shows that the most important variables that determine the price of a house being sold.

Keywords: Housing price, Principal Component Analysis, encoding, Ensemble Algorithms and Feature Engineering.

- Data Sources and their formats

- The dataset has been provided by the FlipRobo technologies only for academic use, not for any commercial.
- The dataset describe data related to housing with 1460 records.
- The dataset is in csv. Format which contains train and test data.
- This dataset is to use simply examines data, which consists 1460 observations and 80 features for model predication.
- The dataset is in both numerical as well as categorical data.

- **Data Preprocessing Done**

The dataset received from FlipRobo technologies, data describe about housing for client US based company which wishes to get in Australian market based on dataset model.

Data Pre-processing and Transformation:

The extracted CSV file was loaded then, Structure of the data is investigated as against the meta data provided. It was discovered that some variables were numeric instead of factors so that was converted to factors.

The ID column in the dataset was dropped as it isn't necessary for the prediction

Missing values in the numerical variables were replaced by the mean of the column while for factor variables it was replaced with 'No_' and this is because the explanation in the metadata says where there is a missing value it is because that feature doesn't exist for that house.

Specificity was done whereby all the levels in the factor columns were correctly represented as it is in the meta data

Standard deviation of the dependent variable was calculated and it showed that the independent variables aren't too far away from the mean of the sale price. It shows that the values in the dataset are normally distributed.

There is no noticeable outlier to be worried about as checked; Multicollinearity was used to check for the correlation between the independent variables.

Important package required

Our primary packages for this project are going to be pandas for data processing, NumPy to work with arrays, matplotlib & seaborn for data visualizations, and finally scikit-learn for building and evaluating our ML model.

Importing required packages

```
import numpy as np
import pandas as pd

import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.tree import DecisionTreeRegressor
from sklearn.ensemble import RandomForestRegressor
from sklearn.ensemble import ExtraTreesRegressor
from sklearn.neighbors import KNeighborsRegressor
from sklearn.experimental import enable_hist_gradient_boosting
from sklearn.ensemble import HistGradientBoostingRegressor
from sklearn.ensemble import GradientBoostingRegressor

from sklearn.metrics import mean_squared_error, mean_absolute_error
from statsmodels.stats.outliers_influence import variance_inflation_factor

from sklearn.model_selection import train_test_split
from scipy.stats import zscore
from sklearn.metrics import r2_score

import warnings
warnings.filterwarnings("ignore")
```

Importing csv dataset :

```
# Load both "train" and "test" dataset:

df_train = pd.read_csv("E:\project flip robo\Project-Housing splitted\housing_train.csv")
df_test = pd.read_csv("E:\project flip robo\Project-Housing splitted\housing_test.csv")
```


• Data Inputs- Logic- Output Relationships

Exploratory Data Analysis (EDA)

- This section shows the exploration done on the dataset, which is what motivated the use of the algorithm. The following are the questions explored in this project and for the sake of writing I will only show some of the visuals here while I will provide the codes that shows the full visualization of all the questions explored.
- Is there a significant relationship between sale price and building's age? It was used to check for this and we can see that there is a relationship between how much old the building is and how much it was sold for.
- What is the average sale price based on overall condition of the house, year it was built, condition1 - proximity to social amenities and sale condition? For overall condition we have levels 1-9 with 1 been the lowest and 9 been the highest and the average of each level is shown. For the others the result is in the code provided.
- What is the sale price distribution based on the overall quality of the house?
- What categories of house (based on age built) have the highest sale price? Here it was obvious that houses built below 50 years have sales price higher than \$700,000
- Sale Price versus Month it was sold. From here we saw that house price increases more during winter than autumn, spring and summer.
- What sale type has the highest sale price? There are 9 different types of sales type that was considered against the sale price.
- Price distribution and season. The bar charts shows there is higher percentage of people buy houses across all seasons at less than 200k.
- At what price will people buy more even with garage attached. Density of garage type has a high peak at claim size about 160k\$. It tells us that people are liable to buy houses at that point regardless of the sale price as long as a garage is attached to the house
- Sales per seasons. The probability of people buying houses is higher in summer and spring is more than autumn and winter.

Fill the missing values for both "Traning" and "Test" dataset

```
# Fill the columns with mean as its continous data
```

```
df_train["LotFrontage"].fillna(df_train["LotFrontage"].mean(), inplace=True)  
df_train["MasVnrArea"].fillna(df_train["MasVnrArea"].mean(), inplace=True)
```

```
# Fill the columns with mode as its categorical data
```

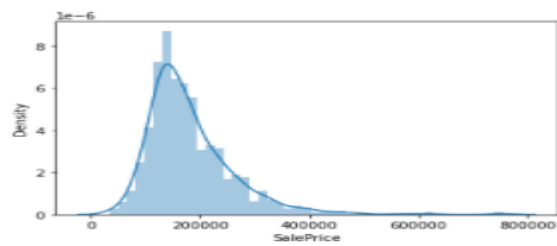
```
df_train["MasVnrType"].fillna(df_train["MasVnrType"].value_counts().index[0], inplace=True)  
df_train["BsmtQual"].fillna(df_train["BsmtQual"].value_counts().index[0], inplace=True)  
df_train["BsmtCond"].fillna(df_train["BsmtCond"].value_counts().index[0], inplace=True)  
df_train["BsmtExposure"].fillna(df_train["BsmtExposure"].value_counts().index[0], inplace=True)  
df_train["BsmtFinType1"].fillna(df_train["BsmtFinType1"].value_counts().index[0], inplace=True)  
df_train["BsmtFinType2"].fillna(df_train["BsmtFinType2"].value_counts().index[0], inplace=True)  
df_train["FireplaceQu"].fillna(df_train["FireplaceQu"].value_counts().index[0], inplace=True)  
df_train["GarageType"].fillna(df_train["GarageType"].value_counts().index[0], inplace=True)  
df_train["GarageYrBlt"].fillna(df_train["GarageYrBlt"].value_counts().index[0], inplace=True)  
df_train["GarageFinish"].fillna(df_train["GarageFinish"].value_counts().index[0], inplace=True)  
df_train["GarageQual"].fillna(df_train["GarageQual"].value_counts().index[0], inplace=True)  
df_train["GarageCond"].fillna(df_train["GarageCond"].value_counts().index[0], inplace=True)
```

Drop the column:

```
df_train.drop(columns = ["Alley", "PoolQC", "Fence", "MiscFeature"], axis=1, inplace=True)  
df_test.drop(columns = ["Alley", "PoolQC", "Fence", "MiscFeature"], axis=1, inplace=True)
```

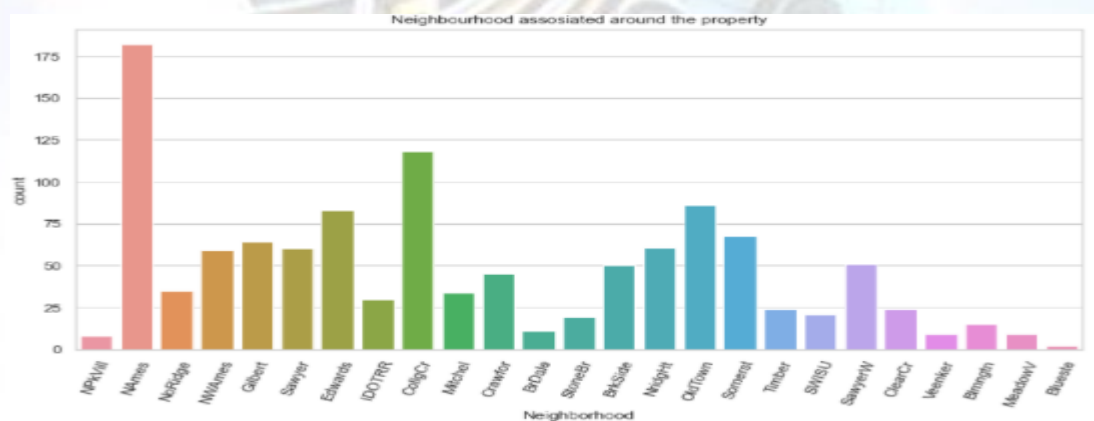
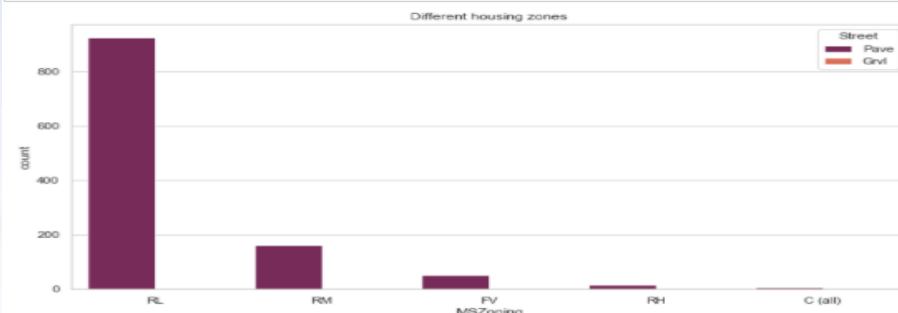
Data Visualization:

```
sns.distplot(df_train['SalePrice'])
<AxesSubplot:xlabel='SalePrice', ylabel='Density'>
```



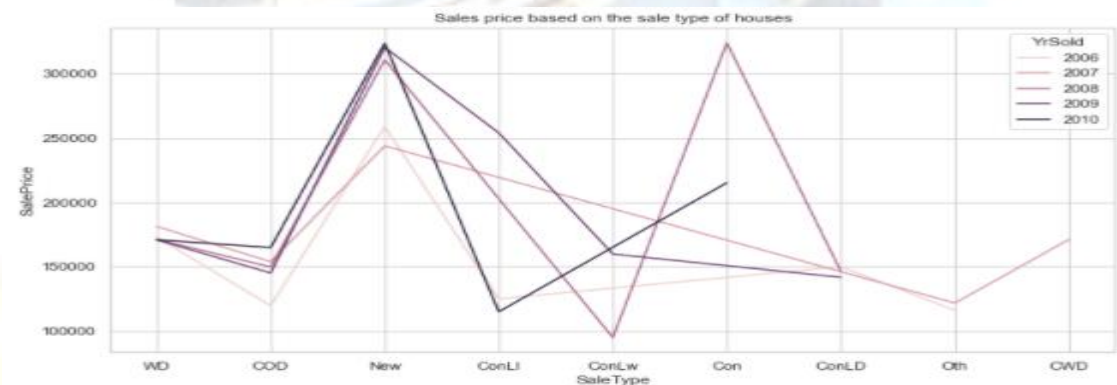
Plot below can identify the zoning area of the properties / apartments that are up for sale.

```
plt.figure(figsize=(12, 6))
sns.set_theme(style="whitegrid")
ax = sns.countplot("MSZoning", data=df_train, hue="Street", palette="rocket").set(title='Different housing zones')
plt.show()
```

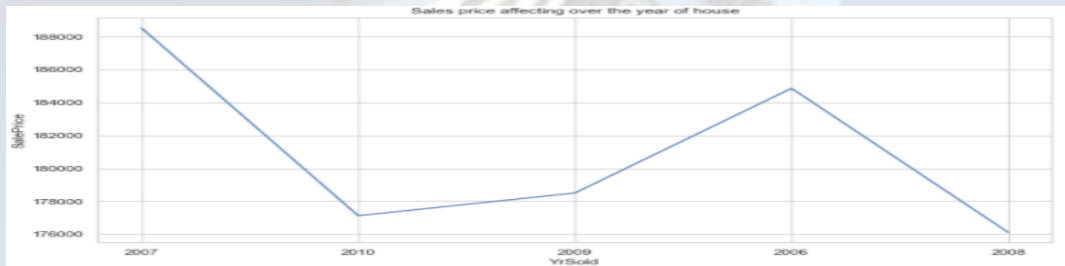


Let's observe if there is any relation between sales price of the properties over the years.

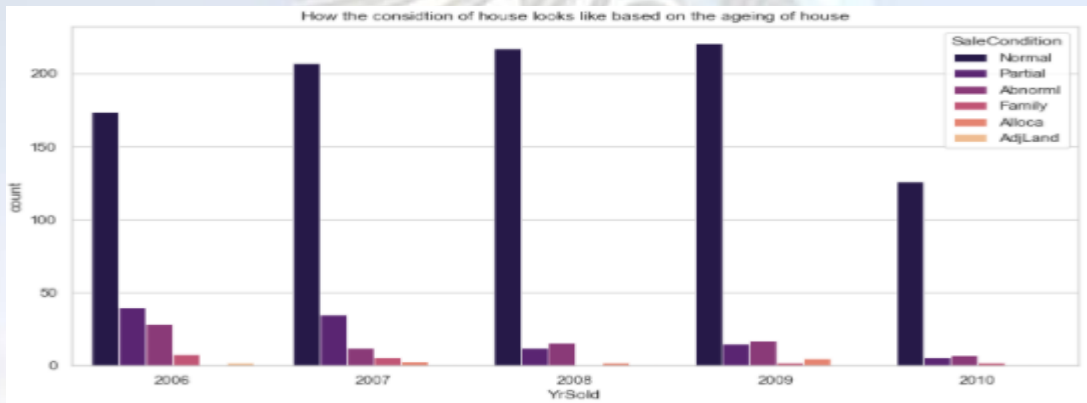
A such a downpayment is possible if the person looking forward to purchase has a steady income. Higher the downpayment, better the loan rates and since cost of living is so important. This could be an ideal option for people who cannot afford or do not wish to spend all at once.



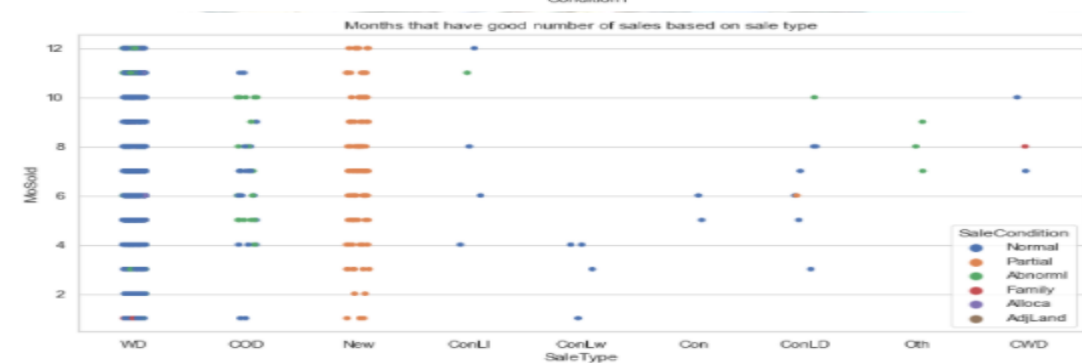
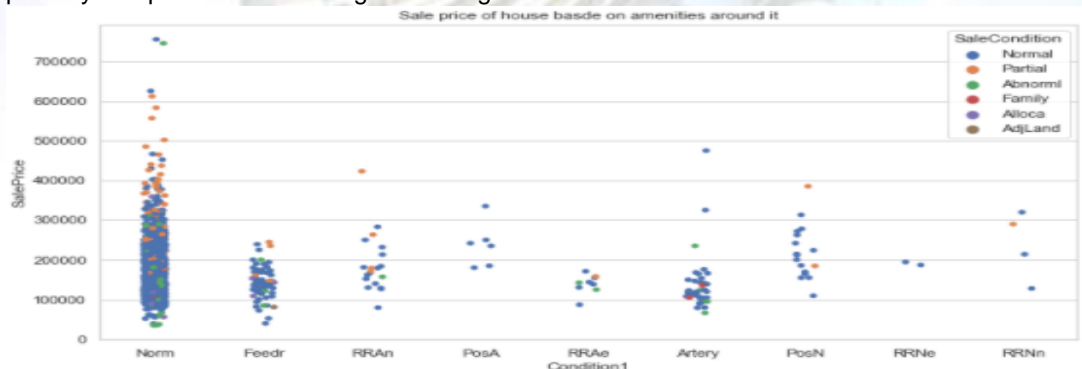
Plot shows us the years when properties were sold the highest.



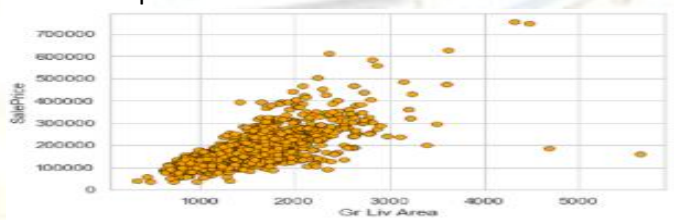
SalesCondition: Condition of sale



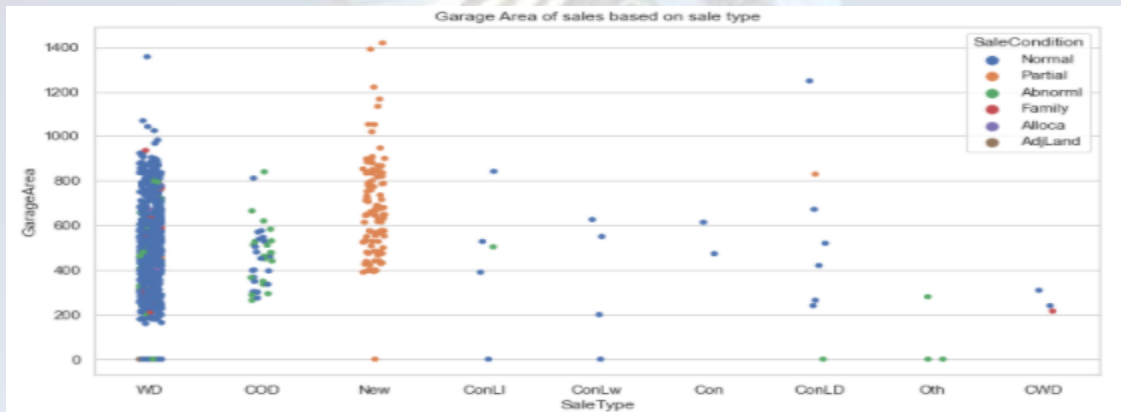
We can see that the roadways or streets seems to play an important role. Houses / properties that are in neighborhoods, condos, villa's etc fetch good amount of values and if it was remodelled or partially completed it seems to get even higher sale to some extent.



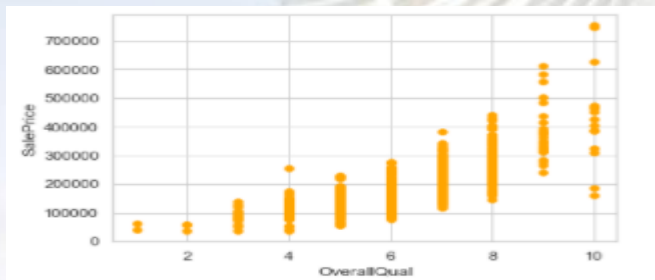
Relationship between GrLivArea and SalePrice



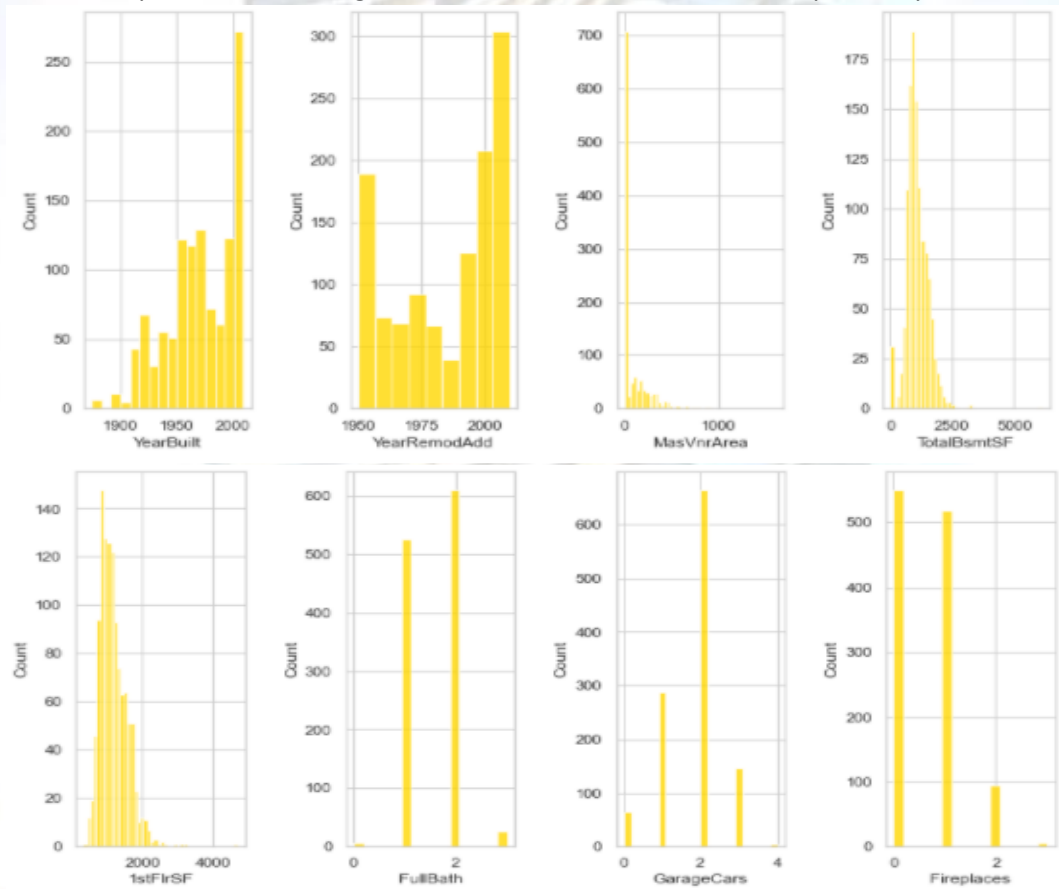
Garage area provided as per the sales type:



OverallQual: Rates the overall material and finish of the house



Relationship between the target variable and the variables that are positively correlated



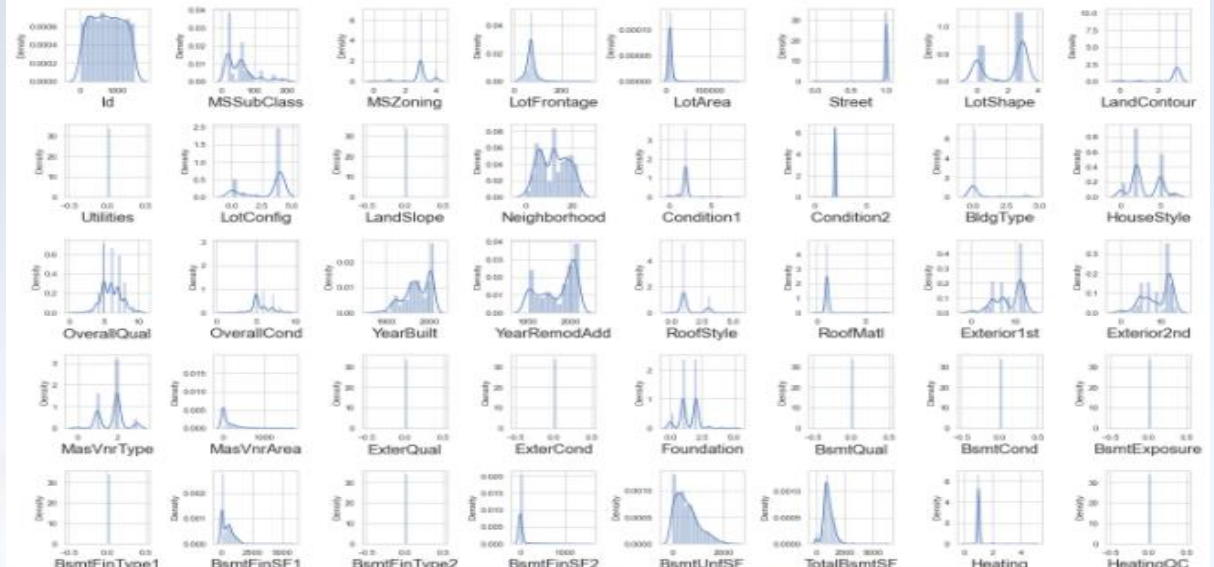
Apart from these customer prefer area located near to railway, highway and market. Apartments/plot near to road plays important role.

Distribution of dataset:

```
# Let us now see the distribution of the "Train dataset"
```

```
plt.figure(figsize=(20,25), facecolor="white")
plotnumber = 1
```

```
for column in df_train:
    if plotnumber <= 77:
        ax = plt.subplot(9,8, plotnumber)
        sns.distplot(df_train[column])
        plt.xlabel(column, fontsize=20)
        plotnumber+=1
plt.tight_layout()
```

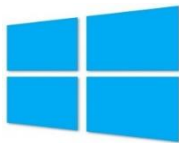


- State the set of assumptions (if any) related to the problem under consideration
- Hardware and Software Requirements and Tools Used

Windows 10 64bit

Anaconda 2021 / Python version – Python 3.9.5

Software: Jupyter notebook, Python, Panda library, numpy library, Matplotlib library, Seaborn library



Python: Python is a general-purpose, and high-level programming language which is best known for its efficiency and powerful functions. Its ease to use, which makes it more accessible. Python provides data scientists with an extensive amount of tools and packages to build machine learning models. One of its special features is that we can build various machine learning with less-code.

NumPy is a Python library used for working with arrays. It also has functions for working in domain of linear algebra, fourier transform, and matrices.

Matplotlib is a plotting library for the Python programming language and its numerical mathematics extension NumPy.

Seaborn is a library for making statistical graphics in Python. It builds on top of matplotlib and integrates closely with pandas data structures. Seaborn helps you explore and understand your data.

• Model/s Development and Evaluation

- Identification of possible problem-solving approaches (methods)
- The factors need to be found which can impact the housing price . This can be done by analysing the various factors and the store the respondent prefers. This will be done by checking each of the factors impacts the respondents decision making.

Encoding Dataset:

Using levelencoding :

```
# Using LabelEncoder:
from sklearn.preprocessing import LabelEncoder
encoder = LabelEncoder()

# Encode the training dataset
df_train.MSZoning = encoder.fit_transform(df_train.MSZoning)
df_train.Street = encoder.fit_transform(df_train.Street)
df_train.LotShape = encoder.fit_transform(df_train.LotShape)
df_train.LandContour = encoder.fit_transform(df_train.LandContour)
df_train.Utilities = encoder.fit_transform(df_train.Utilities)
df_train.LotConfig = encoder.fit_transform(df_train.LotConfig)
df_train.LandSlope = encoder.fit_transform(df_train.LandSlope)
df_train.Neighborhood = encoder.fit_transform(df_train.Neighborhood)
df_train.Condition1 = encoder.fit_transform(df_train.Condition1)
df_train.Condition2 = encoder.fit_transform(df_train.Condition2)
df_train.BldgType = encoder.fit_transform(df_train.BldgType)
df_train.HouseStyle = encoder.fit_transform(df_train.HouseStyle)
df_train.RoofStyle = encoder.fit_transform(df_train.RoofStyle)
df_train.RoofMatl = encoder.fit_transform(df_train.RoofMatl)
df_train.Exterior1st = encoder.fit_transform(df_train.Exterior1st)
df_train.Exterior2nd = encoder.fit_transform(df_train.Exterior2nd)
df_train.MasVnrType = encoder.fit_transform(df_train.MasVnrType)
df_train.ExterQual = encoder.fit_transform(df_train.ExterQual)
df_train.ExterCond = encoder.fit_transform(df_train.ExterCond)
df_train.Foundation = encoder.fit_transform(df_train.Foundation)
df_train.BsmtQual = encoder.fit_transform(df_train.BsmtQual)
df_train.BsmtCond = encoder.fit_transform(df_train.BsmtCond)
df_train.BsmtExposure = encoder.fit_transform(df_train.BsmtExposure)
df_train.BsmtFinType1 = encoder.fit_transform(df_train.BsmtFinType1)
df_train.BsmtFinType2 = encoder.fit_transform(df_train.BsmtFinType2)
```

Using train-test split

Train-Test split:

```
# train dataset with features only
x = df_train.drop(columns = ["SalePrice"], axis=1)
y = df_train["SalePrice"]

# test dataset with features only
x1 = df_test
```

Testing of Identified Approaches (Algorithms)

Below down all the algorithms used for the training and testing.

Using below algorithms for model building :

K-Neighbors Regressor

Decision Tree Regressor

Random Forest Regressor

Extra Trees Regressor

Run and Evaluate selected models

Models performance and accuracy is up-to-mark.

K-Neighbors Regressor:

```
from sklearn.neighbors import KNeighborsRegressor
k_neigh = KNeighborsRegressor()
k_neigh.fit(x_train,y_train)

y_pred = k_neigh.predict(x_test)

print("Adjusted R2 squared : ",k_neigh.score(x_train,y_train))
print("Mean Absolute Error (MAE): ", mean_absolute_error(y_test, y_pred))
print("Mean Squared Error (MSE): ",mean_squared_error(y_test, y_pred))
print("Root Mean Squared Error (RMSE): ",np.sqrt(mean_squared_error(y_test, y_pred)))

Adjusted R2 squared : 0.8643003104379634
Mean Absolute Error (MAE): 23744.676923076924
Mean Squared Error (MSE): 1688441073.2562962
Root Mean Squared Error (RMSE): 41090.64459528831
```


Decision Tree Regressor:

```
from sklearn.tree import DecisionTreeRegressor

dt_reg = DecisionTreeRegressor()
dt_reg.fit(x_train,y_train)

y_pred = dt_reg.predict(x_test)

print("Adjusted R2 squared : ",dt_reg.score(x_train,y_train))
print("Mean Absolute Error (MAE): ", mean_absolute_error(y_test, y_pred))
print("Mean Squared Error (MSE): ",mean_squared_error(y_test, y_pred))
print("Root Mean Squared Error (RMSE): ",np.sqrt(mean_squared_error(y_test, y_pred)))

Adjusted R2 squared : 0.9999949254390876
Mean Absolute Error (MAE): 29578.811965811965
Mean Squared Error (MSE): 2062475577.2222223
Root Mean Squared Error (RMSE): 45414.48642473261
```

Random Forest Regressor:

```
from sklearn.ensemble import RandomForestRegressor

rf_reg = RandomForestRegressor()
rf_reg.fit(x_train,y_train)

y_pred = rf_reg.predict(x_test)

print("Adjusted R2 squared : ",rf_reg.score(x_train,y_train))
print("Mean Absolute Error (MAE): ", mean_absolute_error(y_test, y_pred))
print("Mean Squared Error (MSE): ",mean_squared_error(y_test, y_pred))
print("Root Mean Squared Error (RMSE): ",np.sqrt(mean_squared_error(y_test, y_pred)))

Adjusted R2 squared : 0.9756971968676988
Mean Absolute Error (MAE): 21830.18027730294
Mean Squared Error (MSE): 1560035694.0942047
Root Mean Squared Error (RMSE): 39497.28717385796
```

Extra Trees Regressor:

```
from sklearn.ensemble import ExtraTreesRegressor

extra_reg = ExtraTreesRegressor()
extra_reg.fit(x_train,y_train)

y_pred = extra_reg.predict(x_test)

print("Adjusted R2 squared : ",extra_reg.score(x_train,y_train))
print("Mean Absolute Error (MAE): ", mean_absolute_error(y_test, y_pred))
print("Mean Squared Error (MSE): ",mean_squared_error(y_test, y_pred))
print("Root Mean Squared Error (RMSE): ",np.sqrt(mean_squared_error(y_test, y_pred)))

Adjusted R2 squared : 0.9999949254390876
Mean Absolute Error (MAE): 21995.66603988604
Mean Squared Error (MSE): 1426036354.5247543
Root Mean Squared Error (RMSE): 37762.896532506005
```

Cross-Validation:

Cross-validation

```
scr = cross_val_score(k_neigh, x, y, cv=5)
print("Cross Validation score of KNeighborsRegressor model is:", scr.mean())

scr = cross_val_score(dt_reg, x, y, cv=5)
print("Cross Validation score of DecisionTreeRegressor model is:", scr.mean())

scr = cross_val_score(rf_reg, x, y, cv=5)
print("Cross Validation score of RandomForestRegressor model is:", scr.mean())

scr = cross_val_score(extra_reg, x, y, cv=5)
print("Cross Validation score of ExtraTreesRegressor model is:", scr.mean())

Cross Validation score of KNeighborsRegressor model is: 0.696131368140551
Cross Validation score of DecisionTreeRegressor model is: 0.727306751561138
Cross Validation score of RandomForestRegressor model is: 0.8243606646421894
Cross Validation score of ExtraTreesRegressor model is: 0.8316932592995011
```

Hyper-parameter:

Hyper Parameter Tuning

```
from sklearn.model_selection import GridSearchCV
from pprint import pprint
pprint(rf_reg.get_params())

{'bootstrap': True,
 'ccp_alpha': 0.0,
 'criterion': 'mse',
 'max_depth': None,
 'max_features': 'auto',
 'max_leaf_nodes': None,
 'max_samples': None,
 'min_impurity_decrease': 0.0,
 'min_impurity_split': None,
 'min_samples_leaf': 1,
 'min_samples_split': 2,
 'min_weight_fraction_leaf': 0.0,
 'n_estimators': 100,
 'n_jobs': None,
 'oob_score': False,
 'random_state': None,
 'verbose': 0,
 'warm_start': False}
```

* Using hyper-parameter tuning re-building of new model; RandomForestRegressor is used for the new model creation

New-model:

New-model

```
rf_reg_tuned = RandomForestRegressor(n_jobs=1 )
rf_reg_tuned.fit(x_train,y_train)
rf_reg_tuned_ypred = rf_reg_tuned.predict(x_test)
r2_score_rf_reg_tuned = r2_score(y_test,rf_reg_tuned_ypred)

print("Adjusted R2 squared : ",rf_reg.score(x_train,y_train))
print("Mean Absolute Error (MAE): ", mean_absolute_error(y_test, y_pred))
print("Mean Squared Error (MSE): ",mean_squared_error(y_test, y_pred))
print("Root Mean Squared Error (RMSE): ",np.sqrt(mean_squared_error(y_test, y_pred)))

Adjusted R2 squared : 0.9756971968676988
Mean Absolute Error (MAE): 21995.66603988604
Mean Squared Error (MSE): 1426036354.5247543
Root Mean Squared Error (RMSE): 37762.896532506005
```

Key Metrics for success in solving problem under consideration

Using the sklearn.metrics I have calculated Adjusted R2 squared ,Mean Absolute Error (MAE),Mean Squared Error (MSE),Root Mean Squared Error (RMSE)

Visualizations

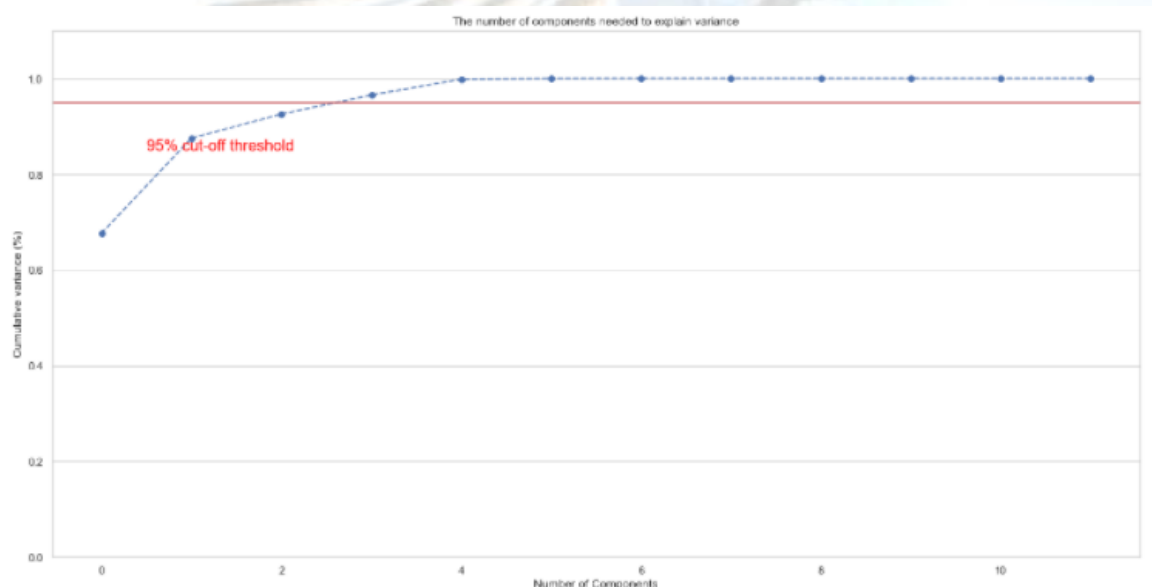
PCA

```
: pca = PCA().fit(x)
plt.rcParams["figure.figsize"] = (20,10)
fig, ax = plt.subplots()
y = np.cumsum(pca.explained_variance_ratio_)

plt.ylim(0.0,1.1)
plt.plot(y, marker='o', linestyle='--', color='b')

plt.xlabel('Number of Components')
plt.ylabel('Cumulative variance (%)')
plt.title('The number of components needed to explain variance')

plt.axhline(y=0.95, color='r', linestyle='--')
plt.text(0.5, 0.85, '95% cut-off threshold', color = 'red', fontsize=16)
ax.grid(axis='x')
plt.show()
```



Interpretation of the Results

- In this research, two experiments were performed, the first experiment was conducted using all the variables available in the dataset after pre-processing, while the second experiment was conducted using most important variables and the goal of this is to be able to improve the model's performance using fewer variables.
- There requirement of train and test and building of many models to get accuracy of the model.
- There are multiple of matric which decide the best fit model like as : R-squared ,RMSE value, etc.
- Database helped in making perfect model and will help in understanding Australian market.

CONCLUSION

In conclusion after reviewing the above papers the first thing I noticed with the dataset is that a lot of research has been done in the past on US housing sector hence I got data for other country, Australian houses. Also different algorithms have been used to predict housing prices but the most efficient one is regression of which my research work will perform on this dataset using different regression algorithms and present the result to know which gives the best accuracy using the R-squared and RMSE value.

Key Findings and Conclusions of the Study

Project findings:

It is seen that different machine learning algorithms have been used to predict housing prices; as well as the different steps taken for the visualization of the dataset and finding co-relation between them. I have use sklearn metrics and calculated Adjusted R2 squared ,Mean Absolute Error (MAE),Mean Squared Error (MSE),Root Mean Squared Error (RMSE)

Learning Outcomes of the Study in respect of Data Science

It find that dataset was quite interesting to handle as it is seen that dataset has both numerical as well as categorical data.

Sale Price versus Month it shows that house price increases more during winter than autumn, spring and summer.

Sales per seasons. The probability of people buying houses is higher in summer and spring is more than autumn and winter.

Geographical location or condition plays as vital role in house/property selection, as We can see that the roadways or streets seems to play an important role.Houses / properties that are in neighbourhoods, condos, villa's etc fetch good amount of values.

Overall developed area are on priority to any costumer rather than developing.

Limitations of this work and Scope for Future Work

About Real-estate industry: It is the huge industry and the data given is quite small but for getting a start and for the decision making the data is quite sufficient. In this industry there is always open opportunities to get start. There are different raw data are available for real estate and applying data science to it, will move this industry at a next level.

We can get the understanding of past, present and the future in real estate market. Today in fast moving world this could be a great investment for business. To beat the market there is requirement to understanding the value of data science in the field of real estate. Data can be driven and it will open opportunity for the investors.



By:Deepak Singh