**FLIP ROBO**

# HOUSING: PRICE PREDICTION

Submitted by:

*SAYAN PAUL*

# ACKNOWLEDGMENT

# INTRODUCTION

- ## Business Problem Framing

  We were 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.

- ## Review of Literature

  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 called 'Surprise Housing'.

- ## Motivation for the Problem Undertaken
  1. Developing an understanding of the real estate market.
  2. Improving my machine learning skills.

# Analytical Problem Framing

- ## Data Sources and their formats

  We were provided the training data from Flip_Robo Technologies in csv format.
  Test data was also provided by the same in csv format.

- ## Data Preprocessing Done
  1. Dropping columns having very high number of null values.
  2. Replacing other NaN values with mean or mode depending on the data type.

     | Column having NaN values | Replacement of NaN values. |
     |---|---|
     | LotFrontage | mean |
     | MasVnrType<br>GarageYrBlt | median |
     | MasVnrArea<br>GarageCond<br>GarageQual<br>GarageFinish<br>GarageType<br>BsmtQual<br>BsmtCond<br>BsmtExposure | mean |

  3. Label encoding
  4. Skewness revomal using power transformer.
  5. Outlier removal
  6. Droppping of columns having very low correlation with target.

**7.** Multi collinearity removal.

- Data Inputs- Logic- Output Relationships
  Random forest Regressor was used here because it was giving the best results.
- State the set of assumptions (if any) related to the problem under consideration
    1. During outlier removal, we kept the absolute value of z less than 3.5 instead of 3.0 to avoid huge loss of data.
    2. We have created two datasets for out analysis→df_rev and df_irev. df_rev doesnot contain the columns which has low correlation with the target. df_irev is created after viewing the visualizations and only the absolutely insensible columns are removed.
    3. df_rev is mainly used for linear,lasso ,ridge regression. While df_irev is used for KNNeighbors and RandomForest Regression.
- Hardware and Software Requirements and Tools Used

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import datetime
import warnings
warnings.filterwarnings('ignore')
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import StandardScaler,power_transform
from sklearn.preprocessing import PowerTransformer
from sklearn.decomposition import PCA
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression,Lasso,Ridge
from sklearn.neighbors import KNeighborsRegressor
from sklearn.ensemble import RandomForestRegressor
from sklearn.preprocessing import PolynomialFeatures
from sklearn.model_selection import cross_val_score
from sklearn.metrics import mean_squared_error,r2_score,mean_absolute_error
from sklearn.model_selection import GridSearchCV
```

# Model/s Development and Evaluation

- ## Identification of possible problem-solving approaches (methods)

    The algorithms tried for solving this problem are:-

    1. Linear Regression
    2. Ridge Regression
    3. Lasso Regression
    4. KNeighbors Regression
    5. Random Forest Regression

- ## Testing of Identified Approaches (Algorithms)

    Hyperparameter tuning was performed on the above above algorithms and the best results are as shown below.

| Algorithms | Best r2_score |
|---|---|
| Linear Regression | 0.9183461118175837 |
| Ridge Regression | 0.9013675441601351 |
| Lasso Regression | 0.9013903564592171 |
| KNN Regression | 0.9041579347151206 |
| Random Forest Regression | 0.927325610260045 |

- ## Run and Evaluate selected models

```python
1   #final LinearRegression model
2   #testing with df_rev
3   x=df_irev.drop(columns=['SalePrice'],axis=1)
4   y=df_irev['SalePrice'].values
5   pca_reg=PCA(n_components=63)
6   xpca=pca_reg.fit_transform(x)
7   xpca=pd.DataFrame(xpca)
8
9   x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.20,random_state=12)
10  lr=LinearRegression()
11  lr.fit(x_train,y_train)
12  y_pred=lr.predict(x_test)
13  print('r2_score is:',r2_score(y_test,y_pred))
14  print('root_mean_sqaured_error:',mean_squared_error(y_test,y_pred)**0.5)
15  print('mean_absolute_error:',mean_absolute_error(y_test,y_pred))
```

```
r2_score is: 0.9183461118175837
root_mean_sqaured_error: 21779.904026224704
mean_absolute_error: 15966.113854410862
```
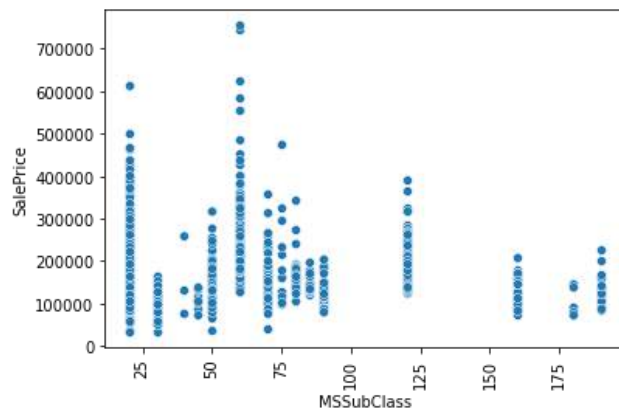
```python
1   #training the final Ridge regression model
2   x=df_rev.drop(columns=['SalePrice'],axis=1)
3   y=df_rev['SalePrice'].values
4
5   x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.20,random_state=27)
6
7   r=Ridge(alpha= 1, solver= 'auto')
8   r.fit(x_train,y_train)
9   y_pred=r.predict(x_test)
10  print('r2_score is:',r2_score(y_test,y_pred))
11  print('root_mean_sqaured_error:',mean_squared_error(y_test,y_pred)**0.5)
12  print('mean_absolute_error:',mean_absolute_error(y_test,y_pred))
```

```
r2_score is: 0.9013675441601351
root_mean_sqaured_error: 23115.716735518836
mean_absolute_error: 18527.42855509859
```

```python
1   #training the final Lasso regression model
2   x=df_rev.drop(columns=['SalePrice'],axis=1)
3   y=df_rev['SalePrice'].values
4
5   x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.20,random_state=27)
6
7   l=Lasso(alpha= 1, max_iter=1300,selection='random')
8   l.fit(x_train,y_train)
9   y_pred=l.predict(x_test)
10  print('r2_score is:',r2_score(y_test,y_pred))
11  print('root_mean_sqaured_error:',mean_squared_error(y_test,y_pred)**0.5)
12  print('mean_absolute_error:',mean_absolute_error(y_test,y_pred))
```

```
r2_score is: 0.9013903564592171
root_mean_sqaured_error: 23113.043410935774
mean_absolute_error: 18530.60519239076
```

```
1  #training the final KNNRegressor model
2  x=df_irev.drop(columns=['SalePrice'],axis=1)
3  y=df_irev['SalePrice'].values
4
5  scaler_knn = StandardScaler()
6  x_knc=scaler_knn.fit_transform(x)
7
8  x_train,x_test,y_train,y_test=train_test_split(x_knc,y,test_size=0.20,random_state=64)
9  knr=KNeighborsRegressor(algorithm='auto',n_neighbors=6,weights='distance',p=1)
10 knr.fit(x_train,y_train)
11 y_pred=knr.predict(x_test)
12 print('r2_score is:',r2_score(y_test,y_pred))
13 print('root_mean_sqaured_error:',mean_squared_error(y_test,y_pred)**0.5)
14 print('mean_absolute_error:',mean_absolute_error(y_test,y_pred))
```

```
r2_score is: 0.9041579347151206
root_mean_sqaured_error: 22032.875791292437
mean_absolute_error: 15482.53233372893
```

```
1  #training the final random forest regressor model.
2  x=df_irev.drop(columns=['SalePrice'],axis=1)
3  y=df_irev['SalePrice'].values
4
5  x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.20,random_state=45)
6  rfr=RandomForestRegressor(max_depth=15,n_estimators=90)
7  rfr.fit(x_train,y_train)
8  y_pred=rfr.predict(x_test)
9  print('r2_score is:',r2_score(y_test,y_pred))
10 print('root_mean_sqaured_error:',mean_squared_error(y_test,y_pred)**0.5)
11 print('mean_absolute_error:',mean_absolute_error(y_test,y_pred))
12
13 #random forest regressor is giving the best results.|
14 #So this is our final model.
```

```
r2_score is: 0.927325610260045
root_mean_sqaured_error: 17792.291465361675
mean_absolute_error: 12929.348347531466
```

- **Key Metrics for success in solving problem under consideration**
    1. r2_score
    2. root_mean_squared_error
    3. mean_absolute_error

- **Visualizations**

STORY 1946 & NEWER ALL STYLES has higher price
STORY 1946 & NEWER has higher price.



Residential Low Density houses can be sold at higher prices.
Commercial(all) are generally sold at very low price.(Also present in very low numbers)
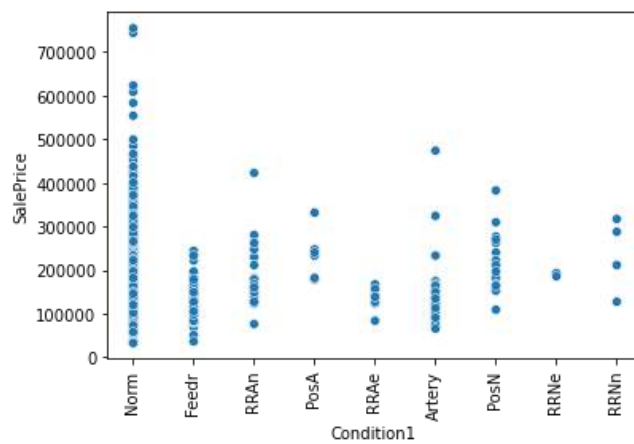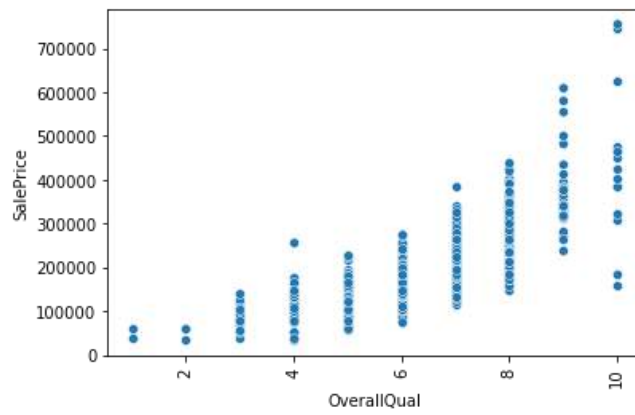


Price increases with LotFrontage
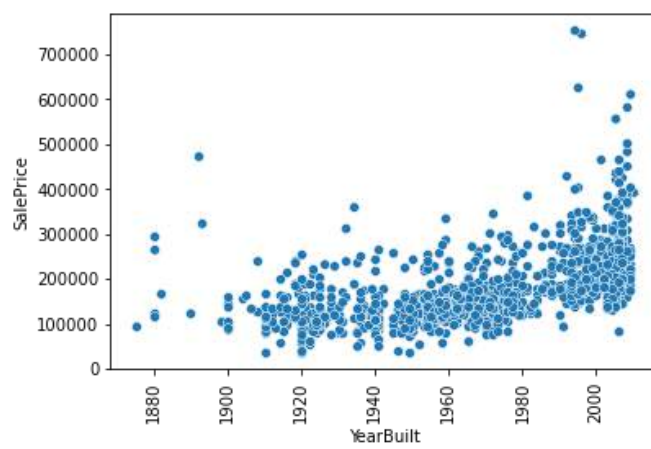
Highly Irregular LotShape has lower price.
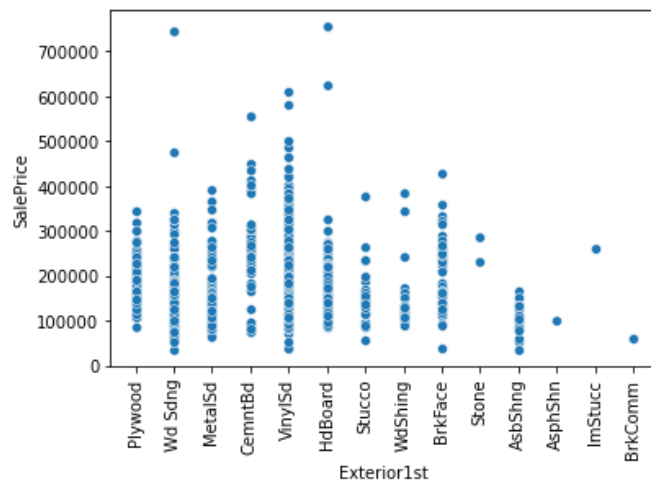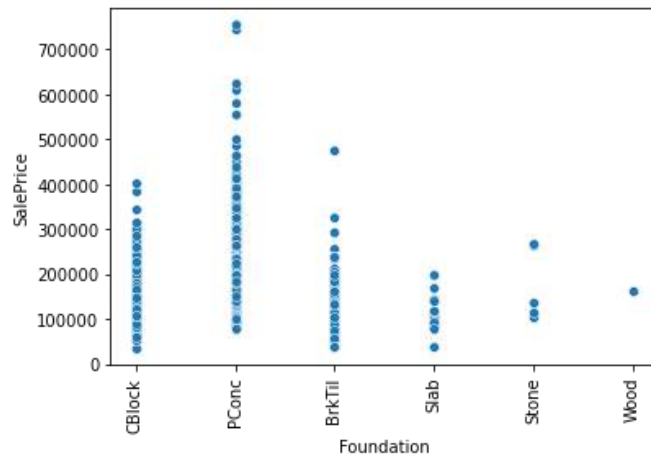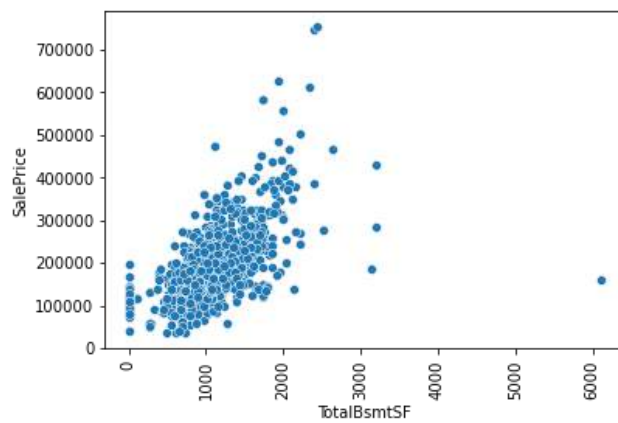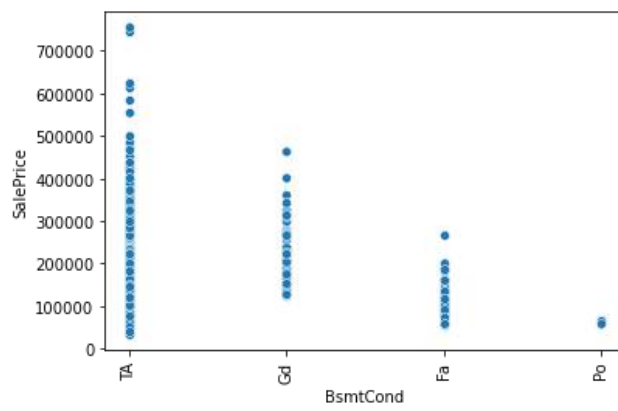


Level land is sold at higher price.

Price shows increase with overall quality.
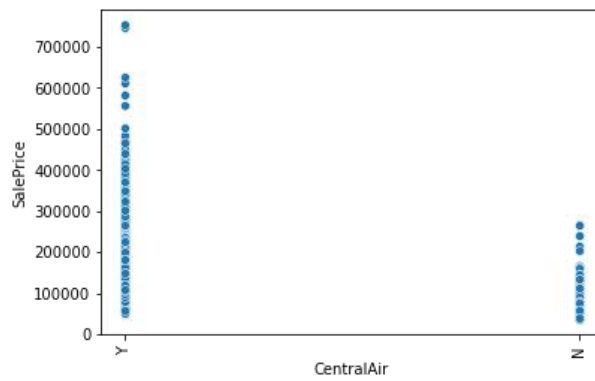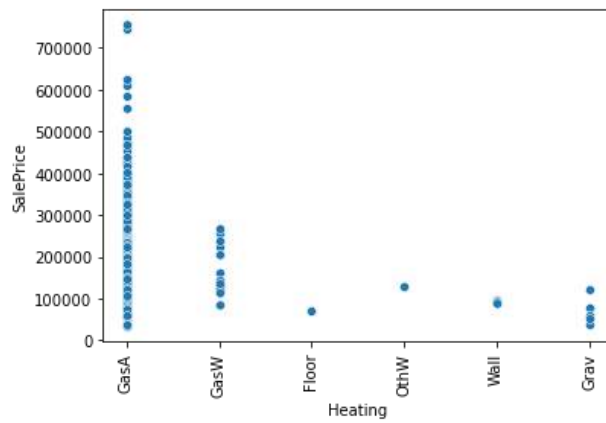


Older buildings as sold as lower price.

Poured concrete foundation buildings has higher price.





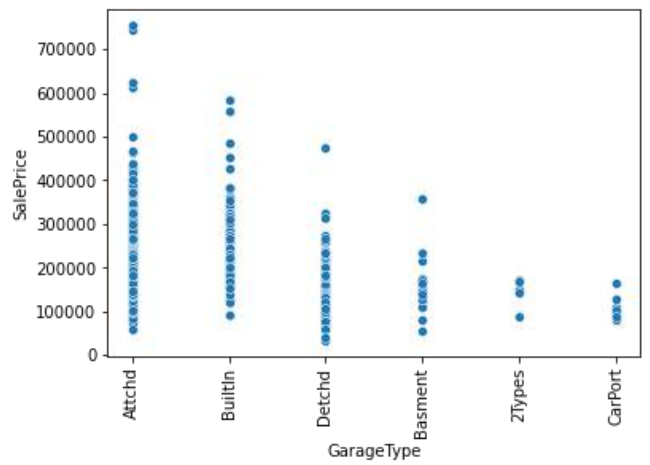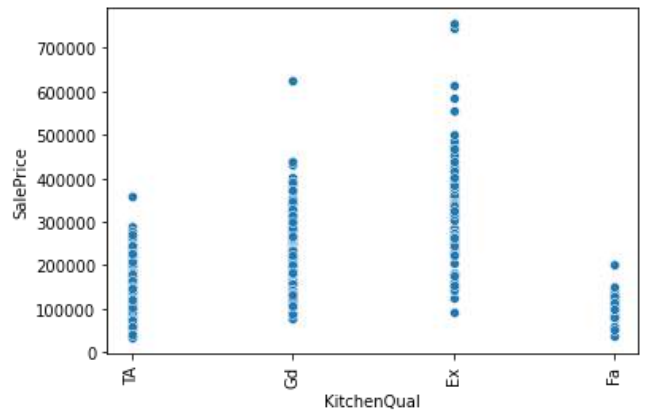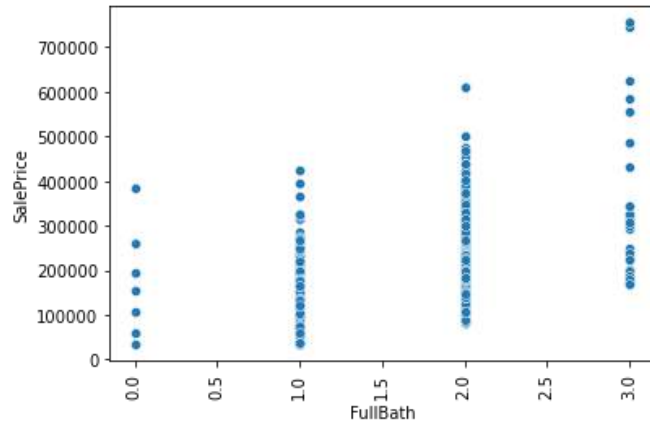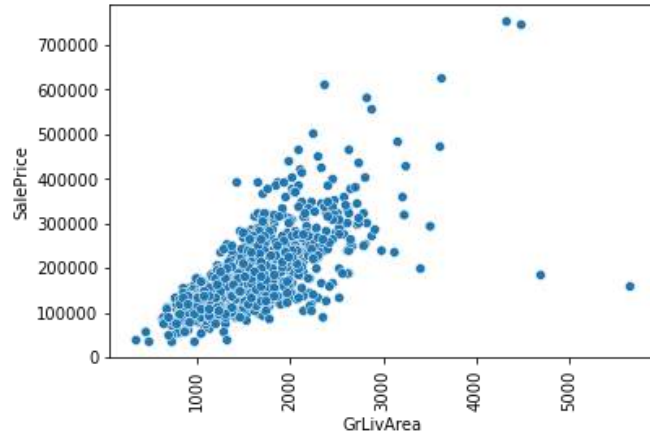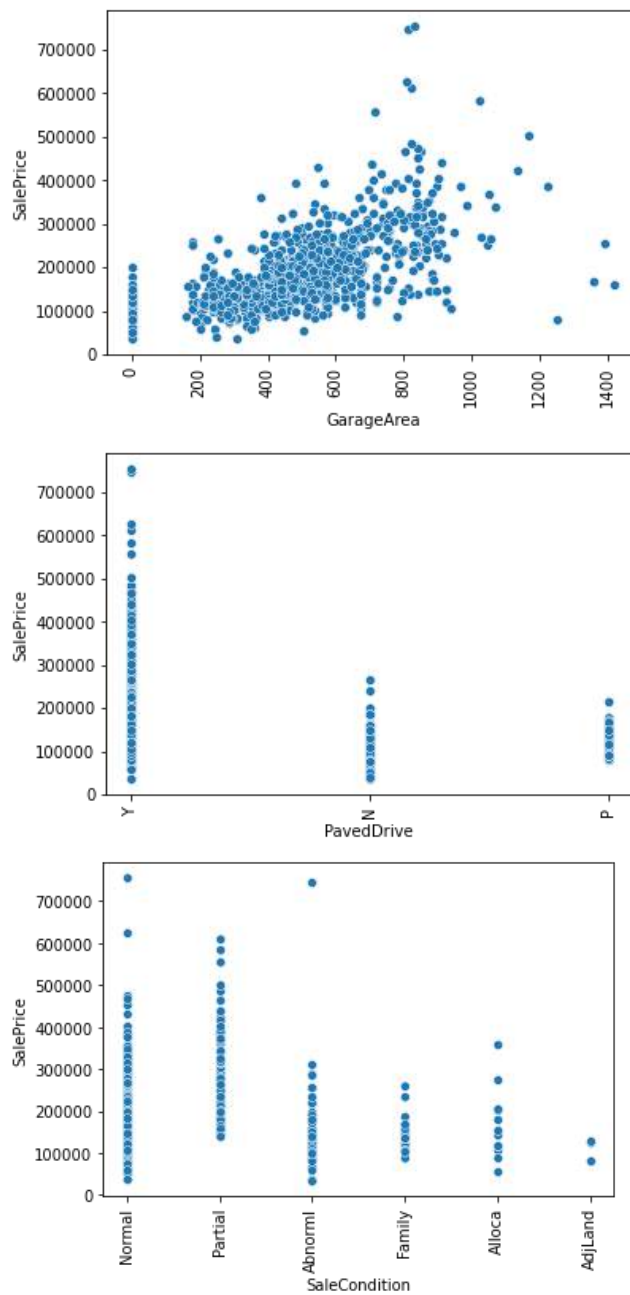Price is increasing with total basement finished surface area.

Buildings with central air conditioner is sold at higher price.

- ## Interpretation of the Results

  OverallQual,GrLivArea,GarageArea,TotalBsmtSF,1stFlrSF,FullBath,YearBuilt etc has very high correlation with SalePrice.

# CONCLUSION

- ## Learning Outcomes of the Study in respect of Data Science

  Dropping columns having low correlation with the target is not always the only way. Sometimes those correlation attributes are necessary to give good results with KNNRegressor and DecisionTree regressor, RandomForest Regressor.

- ## Limitations of this work and Scope for Future Work

  More hyper parameter tunings can be tried with the algorithms.