**MITA CAPSTONE PROJECT**

**Boston House Price Prediction**

# RUTGERS
## THE STATE UNIVERSITY
## OF NEW JERSEY

**Program: Masters in Information Technology And Analytics**

**Semester – Spring 2021**

**Submitted By: Akshata Chavan**

**RUID: 190001210**

**Under Guidence of: Prof. Michail Xyntarakis**

**INDEX:**

# 1. Introduction:

### 1.1 Context
This Capstone Project has been done as part of Masters in Information Technology and Analytics at Rutgers Business School, Newark and mentored by Prof. Michail Xyntarakis. The Project was completed in the duration of three months including deciding project topic, research, learning, development and reporting.

### 1.2 Motivation
The project represent the interest in the field of Machine Learning, Data analytics and Visualization. The skills gained in the curriculum over the past semesters have been applied practically which helped in enhancing the existing theoretical knowledge.

### 1.3 Idea
Residing in Boston currently, this topic immensely interested me to know better about the city and its residential facts. Also, when the project was completed I was able to verify the results in the real world scenario which is accurate to great extent.

### 1.4 Objective
The main objective is to evaluate the performance of various models used and predict the price of the houses in the suburbs of Boston, Massachusetts. The project is created in two parts using Machine Learning Algorithms in Python (Jupyter Notebook) and through Tableau Visualizations. In the initial part of the project, models trained and tested on the dataset can be seen as a good fit and thus could be used to make predictions about the monetary value of the houses. In the later part, various charts and representations have been used to get a better idea of various factors that effect the house price in Boston. When an individual plans to buy a house he considers various aspects like price, the number of rooms, the distance from major locations, educational institutes, crime rate, other residents in that area etc. Here, multiple factors are used to predict the price of the house. This model and visualizations can be very useful to individuals who are new to the city and want to explore/buy a house.

## 2. Dataset

### 2.1 Getting the dataset
The Dataset obtained is the corrected version of original Boston dataset. This is the updated dataset and contains the Harrison and Rubinfield data corrected for minor errors in the latitude and longitude. Also, the town column has been added in the dataset. The corrected dataset was in a txt format and had to be converted to xlsx format.



Figure 2.1: The original dataset

### 2.2. Dataset Description
The dataset used in this project contains 506 rows and 20 columns. The columns in the dataset are as follows:

CRIM - per capita crime rate by town

ZN - proportion of residential land zoned for lots over 25,000 sq. ft.

INDUS - proportion of non-retail business acres per town

CHAS - Charles River dummy variable (= 1 if tract bounds river; 0 otherwise)

NOX - nitric oxides concentration (parts per 10 million)

RM - average number of rooms per dwelling

AGE - proportion of owner-occupied units built prior to 1940

DIS - weighted distances to five Boston employment centers

RAD - index of accessibility to radial highways

TAX - full-value property-tax rate per 10,000usd

PTRATIO - pupil-teacher ratio by town

B - 1000(Bk - 0.63)^2 - where Bk is the proportion of blacks by town

LSTAT - % lower status of the population



Figure 2.2. : Viewing the dataset

4

# 3. Data Exploration and Cleaning

Deriving information about the dataset through exploration and checking for unnecessary columns, null values, repeated values etc.



Figure 3.1: Exploring the datatype information



Figure 3.2: Checking for null values



Figure 3.3: Dropping unnecessary columns

# 4. Exploratory Data Analysis

Exploratory Data Analysis (EDA) is important part of any data analysis. EDA helps to generate questions about the data, visualize, transform and summarize the data. It makes easier to find patterns, spot anomalies and check assumptions. EDA was originally developed by American Mathematician John Turkey in the 1970s and these techniques are widely used in data discovery process till date. The specific statistical techniques to perform with EDA tools include Clustering and dimension reduction, Univariate visualization, Bivariate visualizations and summary statistics, Multivariate visualizations, K-means Clustering etc.

The EDA performed on thr Boston Dataset can be seen in the figures below.



Figure 4.1. : Box and Whisker plot

The above figure is a box plot for various features in the dataset. It is used here to depict the data through quartiles. The lines extending are whiskers and indicate the variability outside the upper and lower quartiles. It can tell you if your data is symmetrical, how tightly your data is grouped, and if and how your data is skewed.

```
In [19]:  #Create DIST plots
          fig, ax = plt.subplots(ncols=7, nrows=2, figsize=(20, 10))
          index = 0
          ax = ax.flatten()

          for col, value in df2.items():
              sns.distplot(value, ax=ax[index])
              index += 1
          plt.tight_layout(pad=0.5, w_pad=0.7, h_pad=5.0)
```

C:\Users\chava\anaconda3\lib\site-packages\seaborn\distributions.py:369: UserWarning: Default bandwidth for data is 0; sk
ipping density estimation.
  warnings.warn(msg, UserWarning)

*Min-Max Normalization*

Figure 4.2. Dist plot

The above figure is a seaborn dist plot. It consist mainly of a histogram with line on it. Dist plot is a univariate distribution of observations.



```
In [20]:  cols = ['CRIM', 'ZN', 'TAX', 'B']
          for col in cols:
              # find minimum and maximum of that column
              minimum = min(df2[col])
              maximum = max(df2[col])
              df2[col] = (df2[col] - minimum) / (maximum - minimum)

In [21]:  fig, ax = plt.subplots(ncols=7, nrows=2, figsize=(20, 10))
          index = 0
          ax = ax.flatten()

          for col, value in df2.items():
              sns.distplot(value, ax=ax[index])
              index += 1
          plt.tight_layout(pad=0.5, w_pad=0.7, h_pad=5.0)
```

C:\Users\chava\anaconda3\lib\site-packages\seaborn\distributions.py:369: UserWarning: Default bandwidth for data is 0; skip
ping density estimation.
  warnings.warn(msg, UserWarning)

Figure 4.3: Dist plot after Min-Max Scaler

7

The min-max scaler is a preprocessing technique in the sklearn.preprocessing library. It is used to transform features by scaling each feature to a given range. The formula for transformation is given by

X_std = (X - X.min (axis=0)) / (X.max(axis=0) - X.min(axis=0))
X_scaled = X_std * (max - min) + min
Where min, max = feature_range.

The syntax in sklearn is
*class* sklearn.preprocessing.**MinMaxScaler**(*feature_range=0, 1*, *\**, *copy=True*, *clip=False)*

```
In [22]:  # standardization
          from sklearn import preprocessing
          scalar = preprocessing.StandardScaler()

          # fit our data
          scaled_cols = scalar.fit_transform(df2[cols])
          scaled_cols = pd.DataFrame(scaled_cols, columns=cols)
          scaled_cols.head()
```
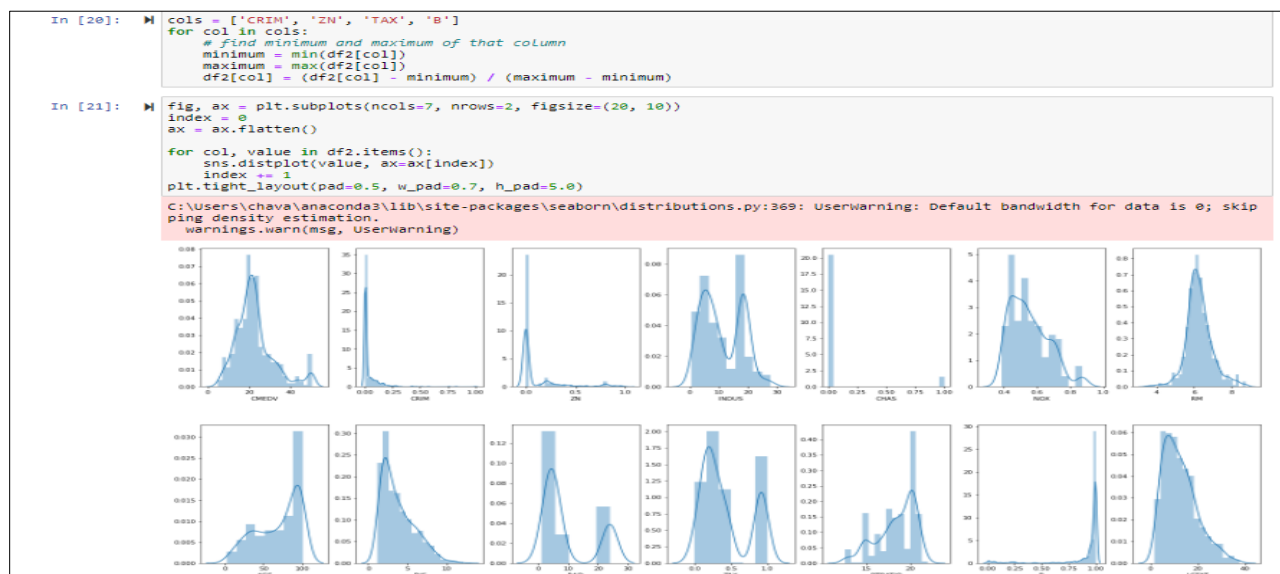
Out[22]:

|   | CRIM | ZN | TAX | B |
|---|------|----|-----|---|
| 0 | -0.419782 | 0.284830 | -0.666608 | 0.441052 |
| 1 | -0.417339 | -0.487722 | -0.987329 | 0.441052 |
| 2 | -0.417342 | -0.487722 | -0.987329 | 0.396427 |
| 3 | -0.416750 | -0.487722 | -1.106115 | 0.416163 |
| 4 | -0.412482 | -0.487722 | -1.106115 | 0.441052 |

```
In [23]:  for col in cols:
              df2[col] = scaled_cols[col]
```

```
In [24]:  fig, ax = plt.subplots(ncols=7, nrows=2, figsize=(20, 10))
          index = 0
          ax = ax.flatten()

          for col, value in df2.items():
              sns.distplot(value, ax=ax[index])
              index += 1
          plt.tight_layout(pad=0.5, w_pad=0.7, h_pad=5.0)

          C:\Users\chava\anaconda3\lib\site-packages\seaborn\distributions.py:369: UserWarning: Default bandwidth for data is 0; skip
          ping density estimation.
            warnings.warn(msg, UserWarning)
```
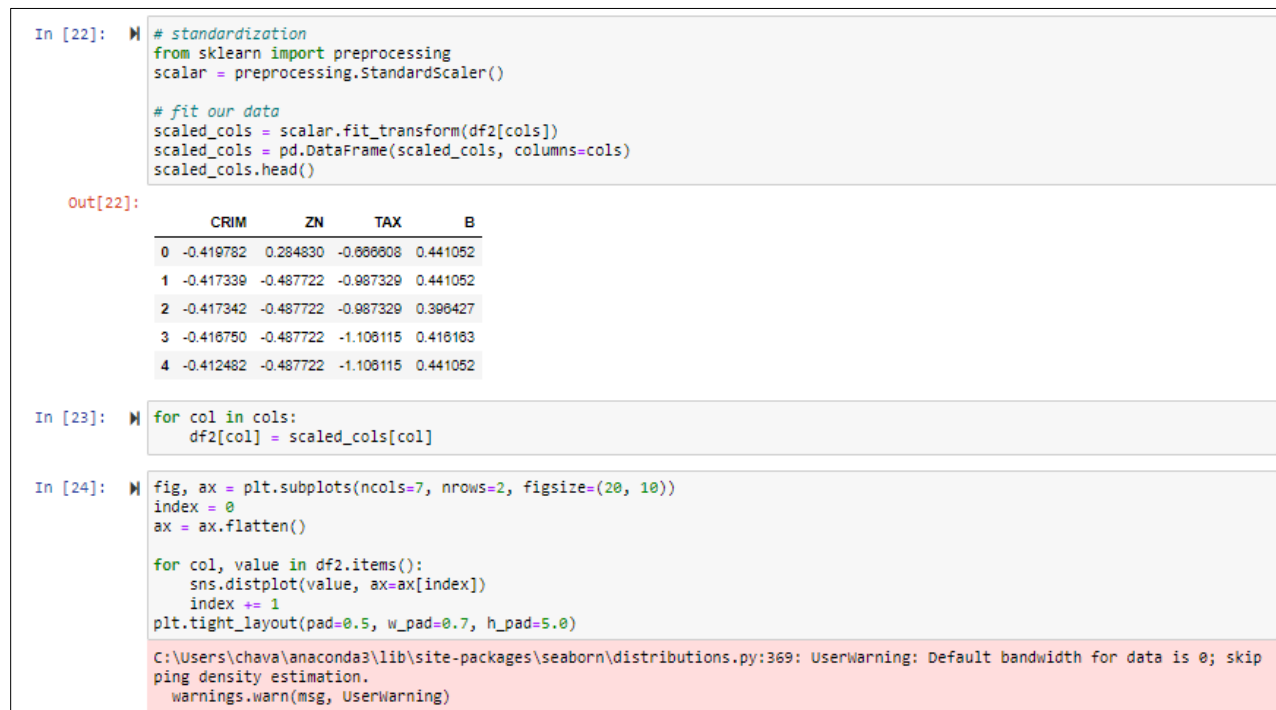
Figure 4.4 : Standard scaler technique application

The standard scaler is a preprocessing technique in the sklearn.preprocessing library. It is used to standardize features by removing mean and scaling it to unit variance. The formula for transformation is given by

z = (x - u) / s

where u is the mean of the training samples or zero if with_mean=False, and s is the standard deviation of the training samples or one if with_std=False.

The syntax in sklearn is

*class* sklearn.preprocessing.StandardScaler(*, copy=True, with_mean=True, with_std=True)

```
"Coorelation Matrix"
```

```
In [25]:  corr = df.corr()
          plt.figure(figsize=(20,10))
          sns.heatmap(corr, annot=True, cmap='coolwarn')
```

```
Out[25]: <matplotlib.axes._subplots.AxesSubplot at 0x220ea51b130>
```



Figure 4.5. : Correlation Matrix

The correlation matrix is a table of correlation coefficients between variables. Every cell shows the correlation between two variables. A correlation matrix is used to summarize data, as an input into a more advanced analysis, and as a diagnostic for advanced analyses.

The common used of correlation matrix is to summarize large amount of data where the aim is to observe patterns, input into another analyses and as a diagnostic. Here, -1 indicates negative linear correlation, 0 indicates no linear correlation and 1 indicates positive linear correlation between two variables.

```
In [26]:  sns.regplot(y=df['CMEDV'],x=df['LSTAT'])
          Out[26]: <matplotlib.axes._subplots.AxesSubplot at 0x220e9fcf310>
```

```
In [27]:  sns.regplot(y=df['CMEDV'],x=df['RM'])
          Out[27]: <matplotlib.axes._subplots.AxesSubplot at 0x220e97d4f40>
```

```
In [28]:  sns.regplot(y=df['CMEDV'],x=df['CRIM'])
          Out[28]: <matplotlib.axes._subplots.AxesSubplot at 0x220e985b640>
```



Figure 4.6: Feature Observation

# 5. Models

## 5.1. Linear Regression

Linear regression is ML algorithm based on supervised learning. It performs regression tasks and target prediction values based on independent variables. It is used for finding out relationship between variables and forecasting.



Figure 5.1 : Linear Regression performed in the project



Figure 5.2 : Output of linear regression

10

## 5.2. Random Forest Regressor

Random forest regression is supervised learning algorithm using ensemble learning which combines multiple ML algorithms to makes accurate predictions. It is powerful and accurate and performs great on problems like nonlinear relationships. The disadvantage is that overfitting occurs easily.



Figure 5.3 : Random Forest Regression performed in the project



Figure 5.4 : Output of Random forest regression

## 5.3. Support Vector Regressor

The SVM Regressor works on the principle of SVM classification. It is adapted form of SVM when dependent variable is numerical and not categorical. The main benefit is that it is a non parametric technique. The output does not depend on distributions of underlying dependent and non dependent variables.



Figure 5.5: SVM Regression performed in the project



Figure 5.6 : Output of SVM regression

## 5.4. XGBoost Regressor

XGBoost is a tool for building supervised learning models. The results of the regression problems are continuous or real values. Some commonly used regression algorithms are Linear Regression and Decision Trees. Root-mean-squared error (RMSE) and mean-squared-error (MAE) are metrics involved in the regression. RMSE is square root of mean squared error. MAE is an absolute sum of actual and predicted differences.



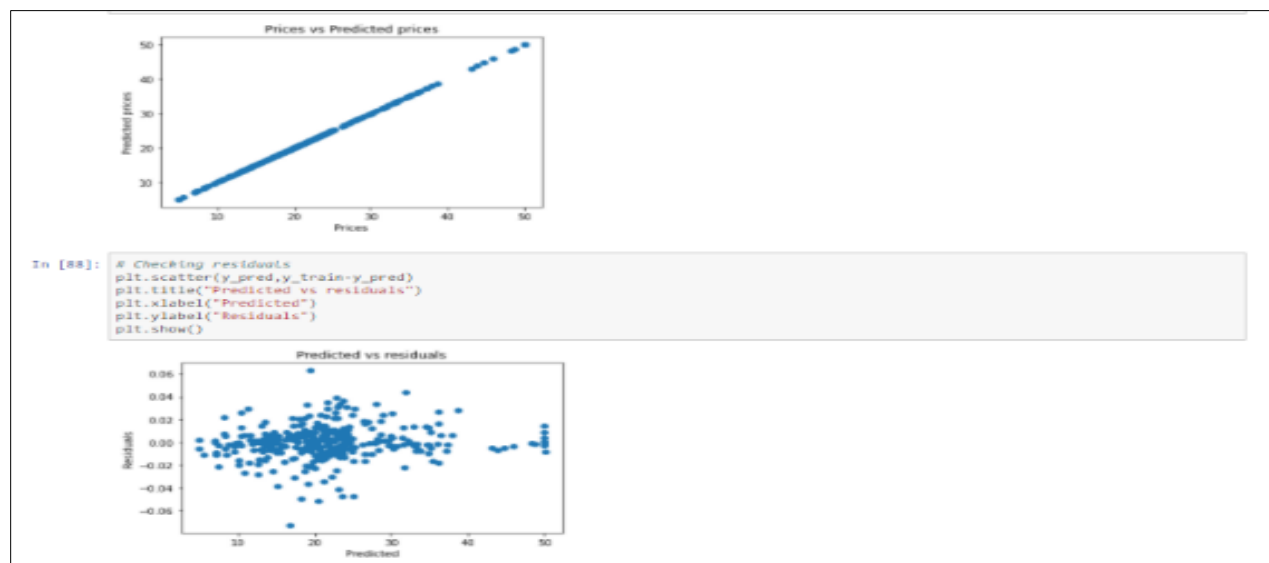Figure 5.7: XGBoost Regression performed in the project



Figure 5.8: Output of XGBoost regression

13

## 5.5. Model Evaluation

After evaluating all the four models to find the best model, it is concluded that the XGBoost Regressor is the best model with R-squared score of 88.8% and Cross validation score of 16.34.
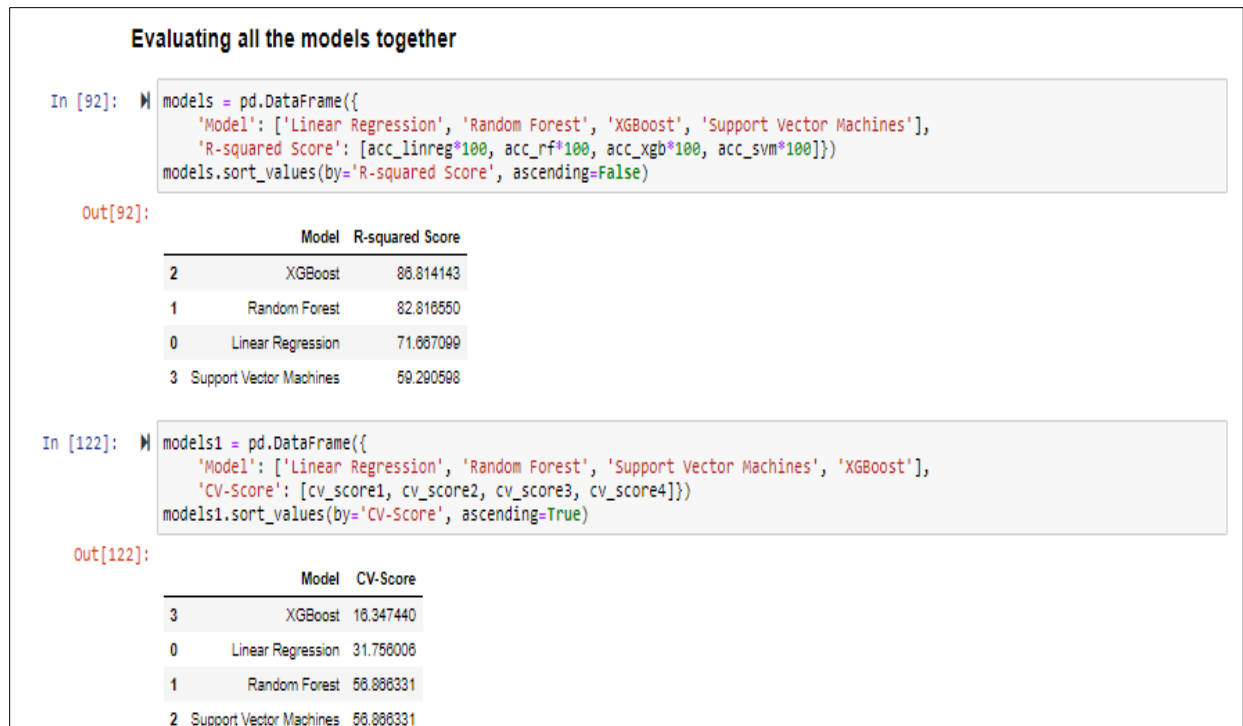


Figure 5.9: Results of model evaluation

# 6. Data Visualization using Tableau

The Tableau Story is the other half of this project where the price predictions that were made using the different models can be visualized. Here, the Tableau visualization will give an interactive demonstration of the factors and their direct or indirect effect on the price of the houses.

There are multiple charts used in the workbook for creating a story based on the Boston House price prediction.

There are four dashboards used to create a story. First dashboard consists of a map showing the areas with size varying marks with respect to the interest measure selected. Second visualizes the effect on the average measures due to different price range. Third dashboards consists of top areas with highest and lowest measures. All the dashboards are interactive and one can select the interest measure to visualize the necessary information.



Figure 6.1 : Tableau story point 1

Figure 6.2 : Tableau story point 2



Figure 6.3 : Tableau story point 3

Figure 6.4 : Tableau story point 4



Figure 6.5 : Tableau calculation used in most of the sheets

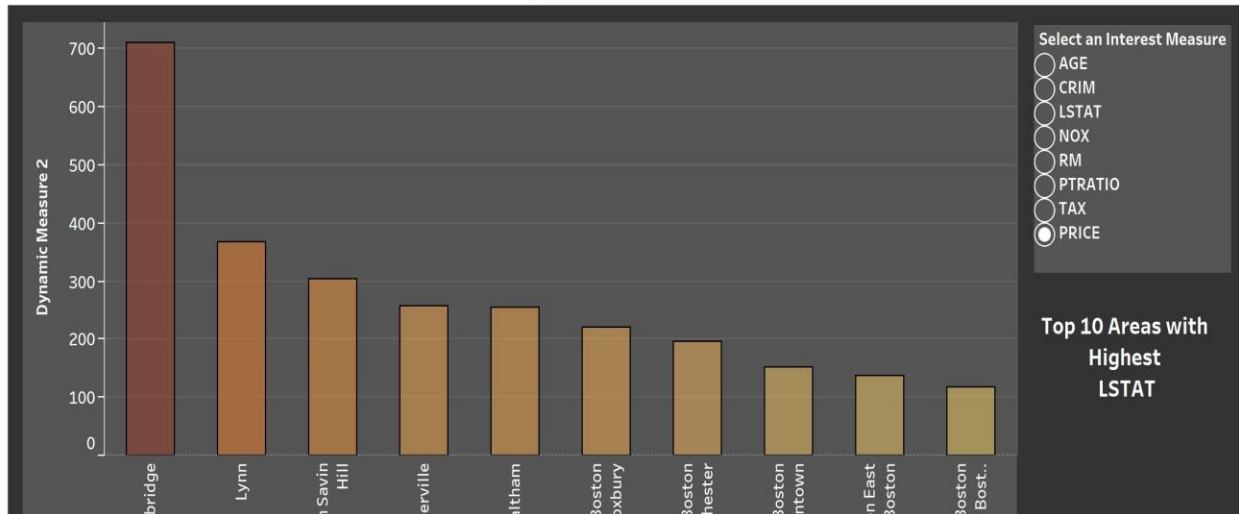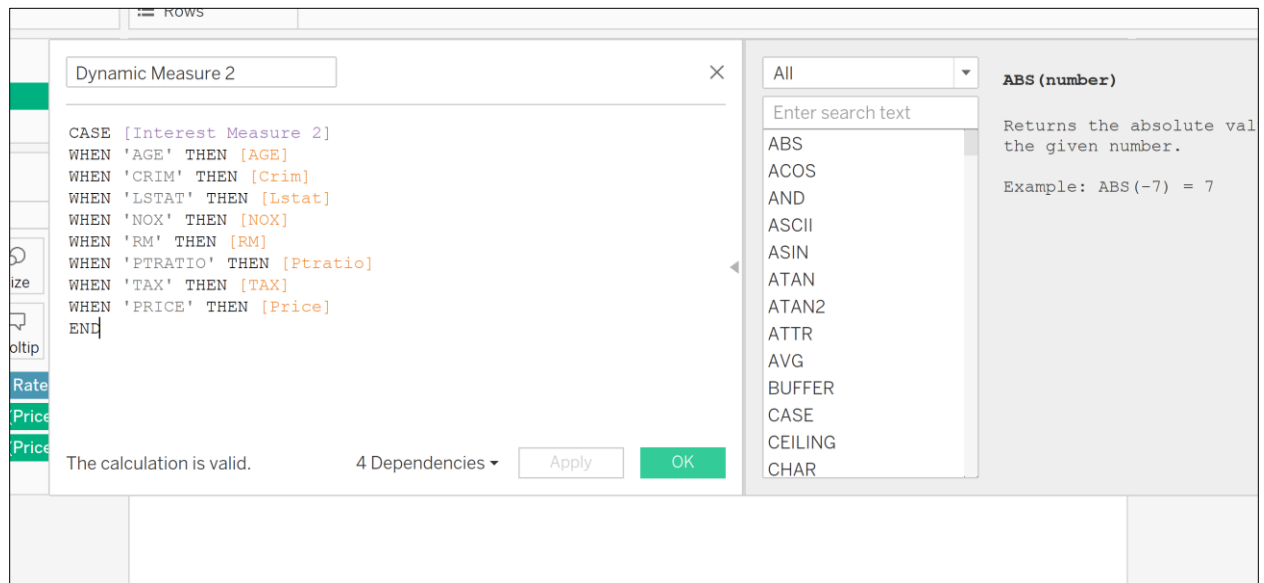# 7. Conclusion

- The XGBoost is the best model for predicting the monetary value of the house price based on various features with R-squared score of 88.8% and cross validation score of 16.34.

- The Price is comparatively low in the areas where the crime rate is high and vice versa in most cases.

- The highest crime rate is in Boston's Roxbury area but the price is not comparatively low.

- The price is high for the house with more number of rooms but that fails in case of some areas. This is due to other measures having more effect.

- The price is less where the Tax is high in most cases.

- The price is very low where the NOX is high.

- Pupil Teacher Ratio does not affect the price in all the areas.

- The price is low in area where the low status population is high.

# 8. References

- http://lib.stat.cmu.edu/datasets/boston_corrected.txt

- https://www.ibm.com/cloud/learn/exploratory-data-analysis

- https://pythonbasics.org/seaborn-distplot/

- https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.MinMaxScaler.html

- https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.StandardScaler.html

- https://www.geeksforgeeks.org/python-linear-regression-using-sklearn/

- https://levelup.gitconnected.com/random-forest-regression-209c0f354c84

- https://www.kdnuggets.com/2017/03/building-regression-models-support-vector-regression.html