

FINAL YEAR PROJECT

COMPARISON OF PERFORMANCE OF SOME MACHINE LEARNING MODELS FOR HOUSE PRICE PREDICTION

PRANAB KUMAR PANDA

University Roll No. 7473U196013

College Roll No. 2019B006



**Institute of Mathematics & Applications
Andharua, Bhubaneswar- 751029
Odisha**

B. Sc. (Hons)

Mathematics & Computing

under the supervision of

INSTITUTE OF MATHEMATICS AND APPLICATIONS, BHUBANESWAR
Odisha, India

Institute of Mathematics & Applications

Andharua, Bhubaneswar- 751 029, Odisha.

Prof. Sudarsan Padhy

Guest Faculty

Institute of Mathematics and Applications

Date 04/08/2022

Supervisor's Certificate

This is to certify that the work presented in this project entitled "*COMPARISON OF PERFORMANCE OF SOME MACHINE LEARNING MODELS FOR HOUSE PRICE PREDICTION*" by *PRANAB KUMAR PANDA*, 7473U196013, is a record of original research/review work carried out by him/her under my supervision and guidance in partial fulfillment of the requirements of the B. Sc. (Hons) in Mathematics and Computing. Neither this project nor any part of it has been submitted for any degree or diploma to any institute or university in India or abroad.

Prof. Sudarsan Padhy

Dedicated

To Maa and Bapa

Declaration

I, *PRANAB KUMAR PANDA*, University Roll Number 7473U196013 hereby declare that this project entitled “*COMPARISON OF PERFORMANCE OF SOME MACHINE LEARNING MODELS FOR HOUSE PRICE PREDICTION*” represents my original work carried out as a B. Sc. (Hons) of IMA Bhubaneswar and, to the best of my knowledge, it is not a complete copy of previously published or written by another person, nor any material presented for award of any other degree or diploma of Institute of Mathematics and Applications, Bhubaneswar or any other institution. Any contribution made to this research by others, with whom I have worked at Institute of Mathematics and Applications, Bhubaneswar or elsewhere, is explicitly acknowledged in the dissertation. Works of other authors cited in this dissertation have been duly acknowledged under the section Bibliography.

Date 04/08/2022

Pranab Kumar Panda

Acknowledgment

This project is the document to support my candidature for the academic degree Bachelor of Science (B.Sc.), in which is detailed explanation of my research and findings on the topic COMPARISON OF PERFORMANCE OF SOME MACHINE LEARNING MODELS FOR HOUSE PRICE PREDICTION. Mentioned study was developed in Institute of Mathematics and Applications, Bhubaneswar at Utkal University, India. I would like to thank Prof. Sudarsan Padhy, who has supervised the academic and professional development of investigation of implementation of various Machine Learning Algorithms and comparing their accuracy. Also, I appreciate my family and personal friends; especially to Suman Sourav Biswal, Prateek Kumar for their enthusiastic support, constant motivation, and constant encouragement, and to my seniors Sampa Mondal and Swati Sucharita for their constructive comments that improved the presentation of these findings.

Abstract

This project proposes a performance comparison between some machine learning algorithms. The machine learning models used in this study are Multiple Linear regression, Ridge regression, Least Absolute Shrinkage and Selection Operator (Lasso), XGBoost, Random Forest, SVR(Support Vector Regression), ANN(Artificial Neural Network). Moreover, this study attempts to analyse the correlation between variables to determine the most important factors that affect house prices in a dataset taken from kaggle.com containing data of residential homes in Ames, Iowa.

The accuracy of the prediction is evaluated by checking the Mean Absolute Error (MAE), Mean Squared Error (MSE), RMSE (Root Mean Square Error), R^2 Square of all the applied model. The test is performed after applying the required pre-processing methods and splitting the data into two parts. However, one part will be used in the training and the other in the test phase. The correlation graphs show the variables' level of dependency. The empirical results show that 'EnclosedPorch', 'KitchenAbvGr', 'OverallCond' influence the house prices negatively whereas 'GarageCars', 'GarageArea', 'GrLivArea', 'OverallQual' impact the house prices positively to a great extent. This project shows that Random Forest Regressor gives the best prediction among other algorithms. At the end of the project, we will compare the accuracy of prediction by the machine learning models mentioned earlier by the help of bar graphs.

Keywords: Multiple Linear regression, Ridge regression, Lasso(Least Absolute Shrinkage and Selection Operator), XGBoost(Express Gradient Boosting), SVR(Support Vector Regression), ANN(Artificial Neural Network), EDA(Exploratory Data Analysis), RMSE(Root Mean Square Error)

Contents

1	Introduction	1
2	Background	2
2.1	Artificial Intelligence	2
2.2	Machine Learning	2
2.3	Models used in the project	3
2.3.1	Multiple Linear Regression	4
2.3.2	Ridge Regression	4
2.3.3	Lasso(Least Absolute Shrinkage and Selection Operator)	5
2.3.4	XGBoost(Express Gradient Boosting)	5
2.3.5	SVM(Support Vector Machine)	6
2.3.6	Random Forest Regressor	7
2.3.7	ANN(Artificial Neural Network)	8
3	Experiment	10
3.1	Dataset Used	10
3.2	Evaluation metrics	10
3.3	Computer Specification	12
3.4	Program Design	12
3.4.1	EDA (Exploratory Data Analysis)	12
3.4.2	Outliers	12

3.4.3 Train-Test Split of the dataset	14
4 Results and Discussion	16
5 Conclusion	34
5.1 Ethics	35
5.2 Future Work	35
A Features	36
B Python Code	39
C Github Link	49
Bibliography	50

Chapter 1

Introduction

“ Data science is the art of extracting meaningful insights from various data sources and deriving useful information from them! ”

In this project, we will compare some of the well-known machine learning models like Multiple Linear regression, Ridge regression, Lasso, XGBoost, SVM(Support Vector Machine), Random Forest regressor, ANN(Artificial Neural Network) after applying them on a dataset imported from the kaggle competition "House Prices - Advanced Regression Techniques". Here, our aim is to find the best model that will fit our problem (predicting the price of a house). We begin by EDA (Exploratory Data Analysis), where we try to visualize various parameters by analyzing the datasets and summarizing their main characteristics. It helps us to get an initial impression of the data we have before making any assumptions. In fact, it can help us to identify obvious errors, as well as understand patterns within the data, detect outliers or anomalous events, find interesting relations among the variables. One of the results that we get is that, 'EnclosedPorch', 'KitchenAbvGr', 'OverallCond' influence the house prices negatively whereas 'GarageCars', 'GarageArea', 'GrLivArea', 'OverallQual' impact the house prices positively to a great extent. Then, we prepare the data for modelling, by cleaning the data, processing the missing data, selecting the relevant variables, deducing some features, running statistical tests, defining the machine learning models and finally choosing the best price prediction model for the test set. For comparison between models we will use four common evaluation metrics i.e, Mean Absolute Error (MAE), Mean Squared Error (MSE), RMSE (Root Mean Square Error), R^2 Square. We have used Python programming language in Jupyter notebook for implementing all the tasks. Please find the github link and code attached in the end of this report to see the codes and the respective outputs.

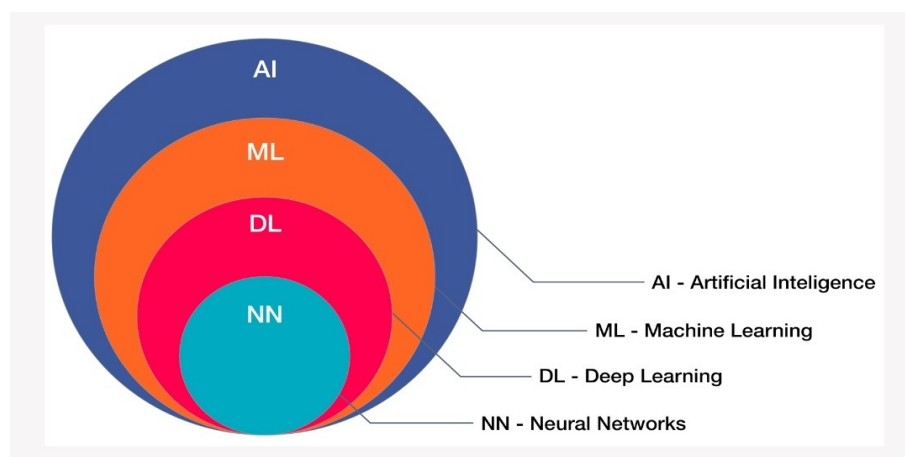
Chapter 2

Background

2.1 Artificial Intelligence

In the simplest terms, AI which stands for artificial intelligence refers to systems or machines that mimic human intelligence to perform tasks and can iteratively improve themselves based on the information they collect. AI is the ability of a machine to perform cognitive functions that we associate with human minds, such as perceiving, reasoning, learning, interacting with the environment, problem solving, and even exercising creativity. [1]

Five basic components of artificial intelligence include learning, reasoning, problem-solving, perception, and language understanding. [2]



2.2 Machine Learning

Machine learning is a subfield of AI. In a nutshell, we can say that ML makes computer learn from experience. With the use of machine learning (ML), software pro-

grammes can predict outcomes more accurately without having to be explicitly instructed to do so. In order to forecast new output values, machine learning algorithms use historical data as input.

MACHINE LEARNING IS OF BASICALLY THREE TYPES:-

1. SUPERVISED LEARNING
2. UNSUPERVISED LEARNING
3. REINFORCEMENT LEARNING

1. Supervised Learning

Supervised comes from the word supervise which means to observe and direct the execution of a task. In supervised learning, a supervisor is required for training the machine to learn from the data.

In supervised learning we do predictions and classification (image classification, audio classification) mainly.

2. Unsupervised Learning

Unsupervised learning is a type of machine learning that involves algorithms that train on unlabeled data. In this we can do clustering and recommendation. Clustering means grouping similar data in one group and unsimilar data in others. Here we can see that provide recommendations like TSNE, PCA (for dimensionality reduction), FACTORING etc.

3. Reinforcement Learning

Reinforcement learning helps us to make decision. It works on the goal and prescribed set of rules for accomplishing the goal. In Reinforcement learning we train robot (agent) to take decision (action) so that it will earn maximum reward. We teach the robot to take decision by giving them reward. Reward in general we can say it's a token of appreciation for doing correct things or to justify the performance.

2.3 Models used in the project

In this project the machine learning models that we are going to use, come under supervised learning. The models that we'll use are

1. Multiple Linear Regression
2. Ridge Regression
3. Lasso (Least Absolute Shrinkage and Selection Operator)
4. XGBoost (Express Gradient Boosting)
5. SVM (Support Vector Machine)

6. Random Forest Regressor
7. ANN(Artificial Neural Network)

2.3.1 Multiple Linear Regression

Multiple linear regression (MLR), or simply multiple regression, is a statistical technique that uses several explanatory variables to predict the outcome of a response variable. The goal of multiple linear regression is to model the linear relationship between the explanatory (independent) variables and response (dependent) variables.

It comes under supervised machine learning and is used to estimate the relationship between one dependent variable and more than one independent variables. Identifying the correlation and its cause-effect helps to make predictions by using these relations. To estimate these relationships, the prediction accuracy of the model is essential; the complexity of the model is of more interest. However, Multiple Linear Regression is prone to many problems such as multicollinearity, noises, and overfitting, which effect on the prediction accuracy.

Regularised regression plays a significant part in Multiple Linear Regression because it helps to reduce variance at the cost of introducing some bias, avoid the overfitting problem and solve ordinary least squares (OLS) problems. There are two types of regularisation techniques L1 norm (least absolute deviations) and L2 norm (least squares). L1 and L2 have different cost functions regarding model complexity .

2.3.2 Ridge Regression

The Ridge Regression is an L2-norm regularised regression technique that was introduced by Hoerl in 1962. It is an estimation procedure to manage collinearity without removing variables from the regression model. In multiple linear regression, the multicollinearity is a common problem that leads least square estimation to be unbiased, and its variances are far from the correct value. Therefore, by adding a degree of bias to the regression model, Ridge Regression reduces the standard errors, and it shrinks the least square coefficients towards the origin of the parameter space.

Ridge formula is:

$$R = \text{Min}(\text{sum of squared residuals} + \alpha * \text{slope}^2) \quad (2.3.1)$$

When Least Squared Error determines the values of parameters, it minimises the sum of squared residuals. However, when Ridge determines the values of parameters, it reduces the sum of squared residuals. It adds a penalty term, where α determines the sever-

ity of the penalty and the length of the slope. In addition, increasing the α makes the slope asymptotically close to zero. Like Lasso, α is determined by applying the Cross-validation method. Therefore, Ridge helps to reduce variance by shrinking parameters and make the prediction less sensitive.

2.3.3 Lasso(Least Absolute Shrinkage and Selection Operator)

LASSO (least absolute shrinkage and selection operator) is a regression analysis method in machine learning that performs both variable selection and regularization in order to enhance the prediction accuracy and interpretability of the resulting statistical model. It was originally introduced in geophysics, and later by Robert Tibshirani, who coined the term in 1996. It is an L1-norm regularised regression technique that can also perform regularisation and feature selection.

Lasso introduces a bias term, but instead of squaring the slope like Ridge regression, the absolute value of the slope is added as a penalty term.

Lasso is defined as:

$$L = \text{Min}(\text{sum of squared residuals} + \alpha * |slope|) \quad (2.3.2)$$

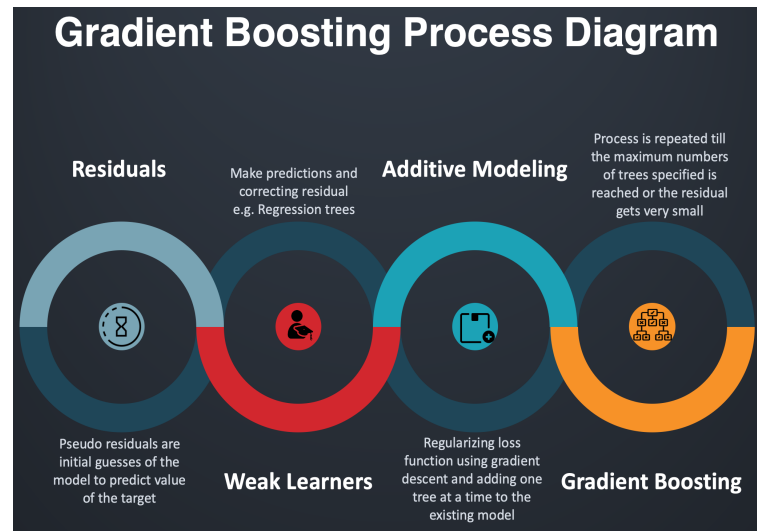
Where $\text{Min}(\text{sum of squared residuals})$ is the Least Squared Error, and $\alpha * |slope|$ is the penalty term. However, α is the tuning parameter which controls the strength of the penalty term. In other words, α the tuning parameter is the value of shrinkage. $|slope|$ the sum of the absolute value of the coefficients.

2.3.4 XGBoost(Express Gradient Boosting)

XGBoost is an optimized distributed gradient boosting library designed to be highly efficient, flexible and portable. It implements machine learning algorithms under the Gradient Boosting framework. XGBoost provides a parallel tree boosting (also known as GBDT, GBM) that solve many data science problems in a fast and accurate way. [9]

Thus it is a scalable Tree Boosting System. It is actually the implementation of gradient boosting trees designed for speed and performance. The gradient boosting decision tree algorithm is implemented in the XGBoost library. Boosting is an ensemble technique that adds new models to correct errors made by existing models. Models are added in a sequential order until no further advancements can be made. Gradient boosting is a method

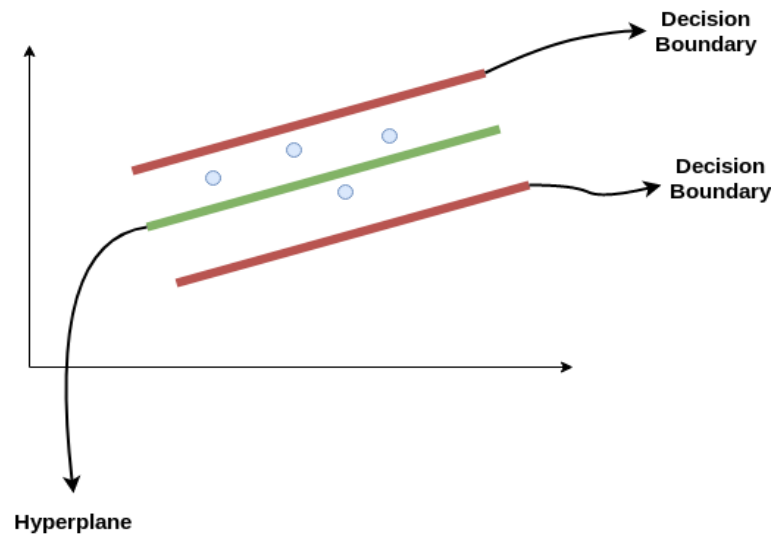
in which new models are created that predict the residuals or errors of prior models and are then combined to make the final prediction. It is called gradient boosting because it employs a gradient descent algorithm to reduce the loss while adding new models. [8]



2.3.5 SVM(Support Vector Machine)

It is a supervised machine learning algorithm used for both classification and regression. [6] It is very effective in high dimensional cases. In classification problem, the objective of SVM algorithm is to find a hyperplane in an N-dimensional space that distinctly classifies the data points. The dimension of the hyperplane depends upon the number of features. If the number of input features is two, then the hyperplane is just a line. If the number of input features is three, then the hyperplane becomes a 2-D plane. It becomes difficult to imagine when the number of features exceeds three.

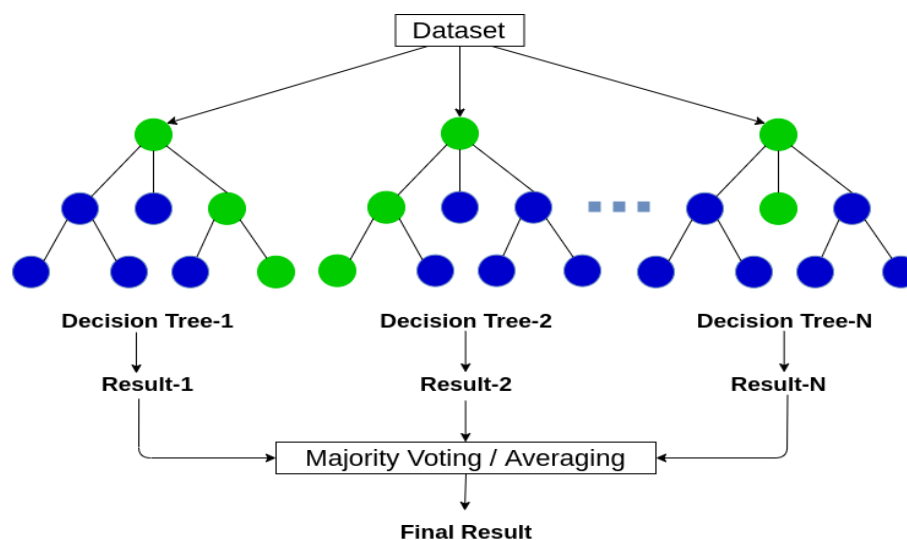
SVR(Support Vector Regression) is the implementation of SVM in regression problems. The problem of regression is to find a function that approximates mapping from an input domain to real numbers on the basis of a training sample



In the above figure the two red lines are the decision boundary and the green line is the hyperplane. The objective of SVR is to basically consider the points that are within the decision boundary line. Our best fit line is the hyperplane that has a maximum number of points. [7]

2.3.6 Random Forest Regressor

Random forest is an ensemble of decision trees. This is to say that many trees, constructed in a certain “random” way form a Random Forest. [5] In a random forest regressor each tree is created from a different sample of rows and at each node, a different sample of features is selected for splitting. Each of the trees makes its own individual prediction. These predictions are then averaged to produce a single result. The averaging makes a Random Forest better than a single Decision Tree hence improves its accuracy and reduces overfitting.

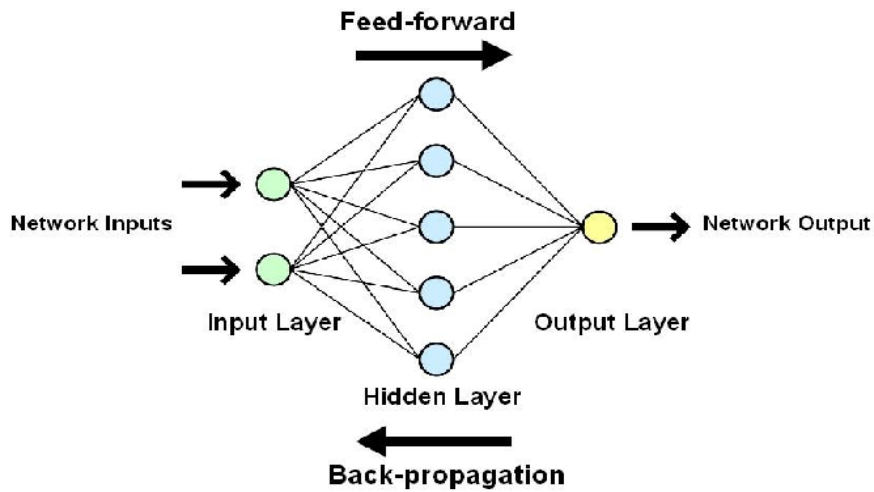


2.3.7 ANN(Artificial Neural Network)

Artificial neural network (ANN) is a simulation of the work of a biological brain. Just as the brain learns and evolves through the experiments that it faces through time to make decisions and predict the result of particular actions, ANN also tries to simulate the brain to learn the pattern in a given data to predict the output of that data.

ANN is based on a collection of connected elements or nodes called neurons. Neurons act as channels that take an input, process it, and then pass it to next neurons for further processing. This transaction or the process of transferring data between neurons is handled in various layers. A simple neural network consists of at least three layers that are input layer, one or more of hidden layers and an output layer. Each layer holds a set of neurons that takes input and process data and finally pass the output to other neurons in the next layer. This process is repetitive until the output layer has been reached, so that eventually, the result can be presented.

An ANN architecture is shown in the following figure



The data that is being held in each neuron is called activation. Activation value ranges from 0 to 1. As shown in the above figure, each neuron is linked to all neurons in the previous layer. Together, all activations from the first layer will decide if the activation will be triggered or not, which is done by taking all activations from the first layer and computing their weighted sum

$$w_1 a_1 + w_2 a_2 + w_3 a_3 \dots + w_n a_n \quad (2.3.3)$$

However, the output can be any number, although it should be only between 0 and 1. Thus, specifying the range of the output value to be within the accepted range. It can be done

by using the Sigmoid function that will put the output to be ranging from 0 to 1. Then the bias is added for inactivity to the equation so it can limit the activation to when it is meaningfully active

$$\sigma(w_1 a_1 + w_2 a_2 + w_3 a_3 \dots + w_n a_n - b) \quad (2.3.4)$$

Where a_i is activation, w_i presents the weight, b is the bias and σ is the sigmoid function.

Nevertheless, after getting the final activation, its predicted value needs to be compared with the actual value. The difference between these values is considered as an error, and it is calculated with the cost function. The cost function helps to detect the error percentage in the model, which needs to be reduced. Applying back-propagation on the model reduces the error percentage by running the procedures backwards to check on how the weights and bias are affecting the cost function.

Back-propagation is simply the process of reversing the whole activations transference among neurons. The method calculates the gradient of the cost function concerning the weight. It is performed in the training stage of the feed-forward for supervised learning

Chapter 3

Experiment

3.1 Dataset Used

The dataset used for this project is taken from kaggle. It is uploaded there in the name of "House Prices - Advanced Regression Techniques". [3] In the project, we have only taken the train.csv file from the site and divided it into two parts, one for training purpose and another one for prediction and testing purpose. The dataset is split into 3:1 ratio.

It has 81 columns including "Id" that is being used to identify each house uniquely and the "Sale Price" that we will try to predict using various machine learning models. A screenshot of the dataset is attached below

	Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	LandContour	Utilities	...	PoolArea	PoolQC	Fence	MiscFeature	MiscVal
0	1	60	RL	65.0	8450	Pave	NaN	Reg	Lvl	AllPub	...	0	NaN	NaN	NaN	0
1	2	20	RL	80.0	9600	Pave	NaN	Reg	Lvl	AllPub	...	0	NaN	NaN	NaN	0
2	3	60	RL	68.0	11250	Pave	NaN	IR1	Lvl	AllPub	...	0	NaN	NaN	NaN	0
3	4	70	RL	60.0	9550	Pave	NaN	IR1	Lvl	AllPub	...	0	NaN	NaN	NaN	0
4	5	60	RL	84.0	14260	Pave	NaN	IR1	Lvl	AllPub	...	0	NaN	NaN	NaN	0

<https://www.kaggle.com/competitions/house-prices-advanced-regression-techniques/data>

3.2 Evaluation metrics

The prediction accuracy will be evaluated by four common evaluation metrics i.e, Mean Absolute Error (MAE), Mean Squared Error (MSE), RMSE (Root Mean Square Error), R^2 Square. R^2 will show if the model is overfitted, whereas MAE, MSE, RMSE show the error percentage between the actual and predicted data, which in this case, the house prices.

Mean Absolute Error (MAE) is the mean of the absolute value of the errors:

$$\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Mean Squared Error (MSE) is the mean of the squared errors:

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Root Mean Squared Error (RMSE) is the square root of the mean of the squared errors:

$$\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

NOTE: Here y_i represents the actual value (SalePrice) and \hat{y}_i represents the predicted value (SalePrice).

R^2 Square is the proportion of the variation in the dependent variable that is predictable from the independent variable(s). It is also known as the coefficient of determination.

If $R^2 = 0$, it indicates that the dependent variable cannot be predicted from the independent variable(s).

If $R^2 = 1$, it indicates the dependent variable can be predicted without error from the independent variable(s).

If R^2 lies in between 0 to 1, it indicates the extent to which the dependent variable is predictable.

Comparing these metrics:

MAE is the easiest to understand, because it's the average error.

MSE is more popular than MAE, because MSE "punishes" larger errors, which tends to be useful in the real world.

RMSE is even more popular than MSE, because RMSE is interpretable in the "y" units.

The first three are loss functions, so we want to minimize them. And more is the value of R^2 Square better and more accurate will be our predictions.

3.3 Computer Specification

The needed time to train the model depends on the capability of the used system during the experiment. Some libraries use GPU resources over the CPU to take a shorter time to train a model.

The specification of the computer system on which this code is being implemented is given below:

1. **OPERATING SYSTEM** Windows 11
2. **PROCESSOR** Intel i5 10th generation
3. **RAM** 8 GB
4. **GRAPHICS CARD** NVIDIA MX110

3.4 Program Design

The algorithms used in this study have different properties that will be used during the implementation. The experiment is done with the Jupyter Notebook using Python programming language. However, in all algorithms, the data is split into four variables, namely; `X_train`, `X_test`, `y_train` and `y_test`, by using `train_test_split` class from the library `sklearn.model_selection`. In addition, in all algorithms, the `train_test_split` class takes as parameters the independent variables, which is the data, the dependent variable, which is the `SalePrice`, `test_size = 0.25`, and `random_state = 0`.

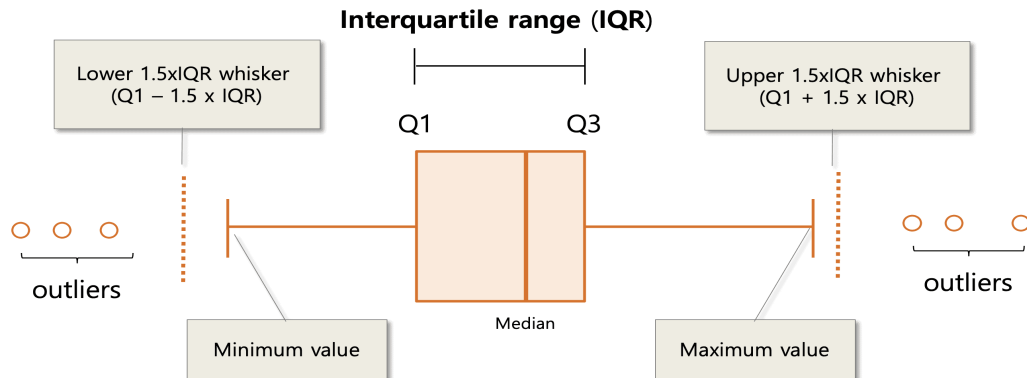
3.4.1 EDA (Exploratory Data Analysis)

The main motive of EDA is to get an overview of the numerical and categorical features of a dataset. In this project, python libraries like **pandas**, **matplotlib**, **seaborn** are imported and pre-defined functions are used for graphical representation of various features and to see the relationships in between them. Also we have used **heat map** to see the covariance between various numerical features. In the program we have used various graphs like scatter plot, box plot to give a pictorial description dataset and compare the accuracy of different models.

3.4.2 Outliers

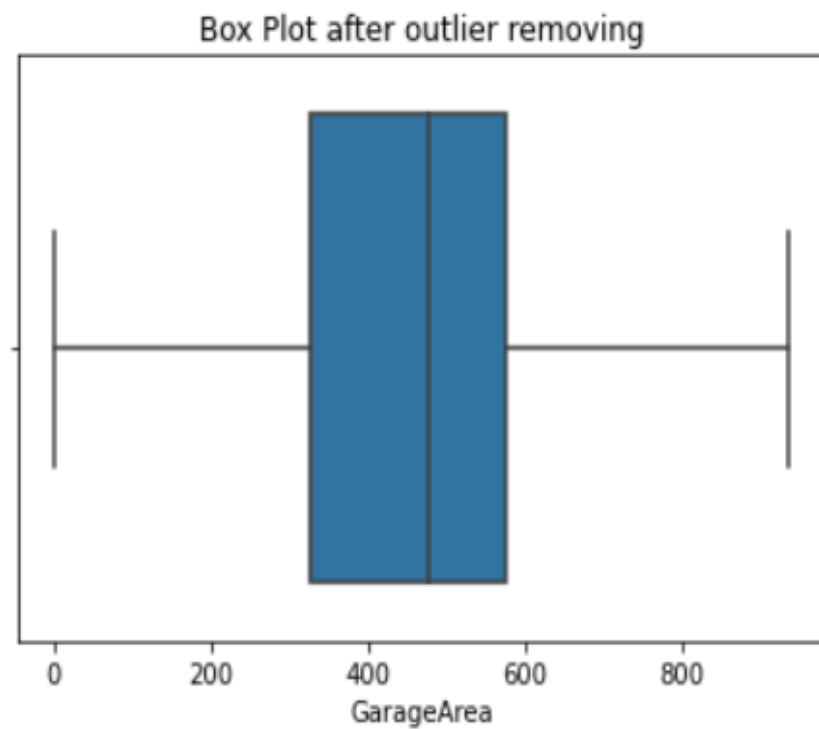
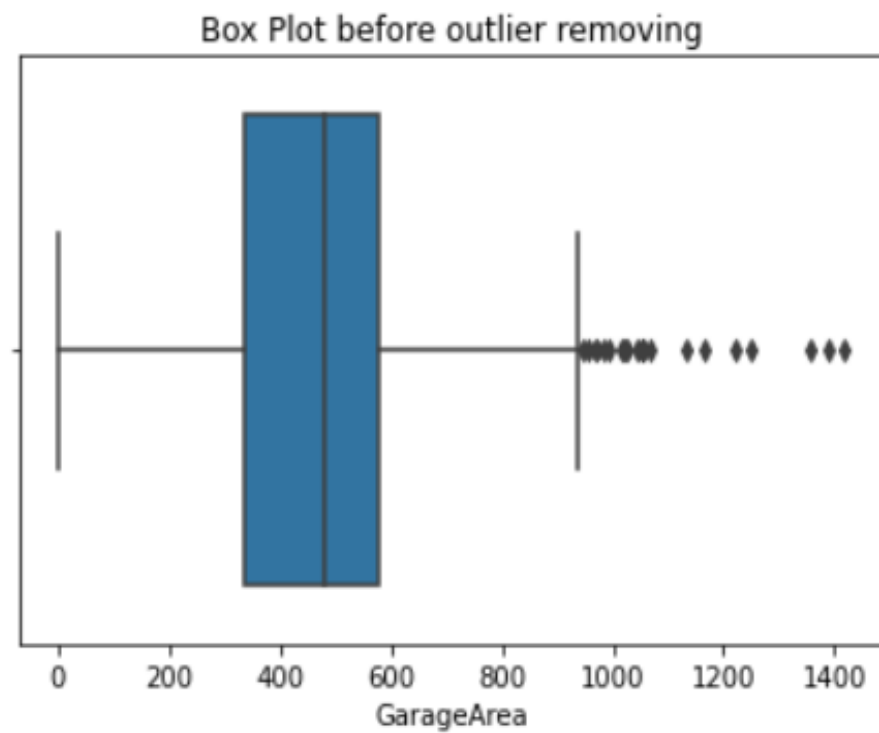
Outliers are those data points which differs significantly from other observations present in given dataset. It can occur because of variability in measurement and due to misinterpretation in filling data points. [4] They can be natural or because of entry errors, or due to measurement errors as well. In the project we have used **box plots** for visualising the outliers. Box plot is a graphical display for describing the distributions of the data.

Box plot uses the median and the lower and upper quartiles to showcase the distribution of data points .



```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5
6 ghar = pd.read_csv("../PROJECT/OTHERS/house-prices-advanced-
    regression-techniques/train.csv")
7 sns.boxplot(ghar['GarageArea'])
8 plt.title("Box Plot before removing outlier")
9 plt.show()
10
11 def drop_outliers(df, field_name):
12     iqr = 1.5 * (np.percentile(df[field_name], 75) - np.percentile(
13         df[field_name], 25))
14     df.drop(df[df[field_name] > (iqr + np.percentile(df[field_name], 75))].index, inplace=True)
15     df.drop(df[df[field_name] < (np.percentile(df[field_name], 25) - iqr)].index, inplace=True)
16
17 drop_outliers(ghar, 'GarageArea')
18 sns.boxplot(ghar['GarageArea'])
19 plt.title("Box Plot after outlier removing")
20 plt.show()
```

Listing 3.1: Python example



3.4.3 Train-Test Split of the dataset

The procedure of splitting a dataset involves taking a dataset and dividing it into two subsets. The first subset is used to fit the model and is referred to as the training dataset. The second subset is not used to train the model; instead, the input element of the dataset

is provided to the model, then predictions are made and compared to the expected values. This second dataset is referred to as the test dataset.

Train Dataset: Used to fit the machine learning model.

Test Dataset: Used to evaluate the fit machine learning model.

In the project we have used **sklearn.model_selection** library and **train_test_split** function to split the data into 3:1 ratio, which means we have used 75% of data for training the ML models and 25% for testing the accuracy of the predictions made by the ML models.

```
1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X, y,
    random_state=0, train_size = .75)
```

Listing 3.2: Program code for dataset split

Chapter 4

Results and Discussion

In the beginning, we got to see that the dataset on which we were going to work had 81 columns including 'Id' that is used for identifying houses and 'SalePrice' which depicts the price of a house. On evaluation, we get to know that the used dataset has **38 Numerical features** and **43 Categorical features**.

In [6]:

```
numerical_feats = ghar.dtypes[ghar.dtypes != "object"].index
print("Number of Numerical features: ", len(numerical_feats))

categorical_feats = ghar.dtypes[ghar.dtypes == "object"].index
print("Number of Categorical features: ", len(categorical_feats))
```

```
Number of Numerical features: 38
Number of Categorical features: 43
```

We then drop all the columns having categorical features, since we wish to work on numerical features to predict the Sale Price of house.

Now, we have 38 columns which means 37 features (including Id) to predict the house price. But we can't use them without pre-processing, since on checking we get to see that there are three columns (**LotFrontage**, **GarageYrBlt**, **MasVnrArea**) having some missing values under them(**NULL VALUES**).

```
total = ghar.isnull().sum().sort_values(ascending=False)
percent = (ghar.isnull().sum()/ghar.isnull().count()).sort_values(ascending=False)
missing_data = pd.concat([total, percent], axis=1, keys=['Total', 'Percent'])
missing_data.head(10)
```

	Total	Percent
LotFrontage	259	0.177397
GarageYrBlt	81	0.055479
MasVnrArea	8	0.005479
Id	0	0.000000
OpenPorchSF	0	0.000000
KitchenAbvGr	0	0.000000
TotRmsAbvGrd	0	0.000000
Fireplaces	0	0.000000
GarageCars	0	0.000000
GarageArea	0	0.000000

Numerical features having null values

As we can see here, 'LotFrontage' has 259 null missing row entries under it out of 1460 entries that is around 17.7 % missing entries. Similarly, 'GarageYrBlt', 'MasVnrArea' have 81, 8 missing entries respectively. So, we drop these three columns. So, finally we have 34 features to predict the Sale Price of a house.

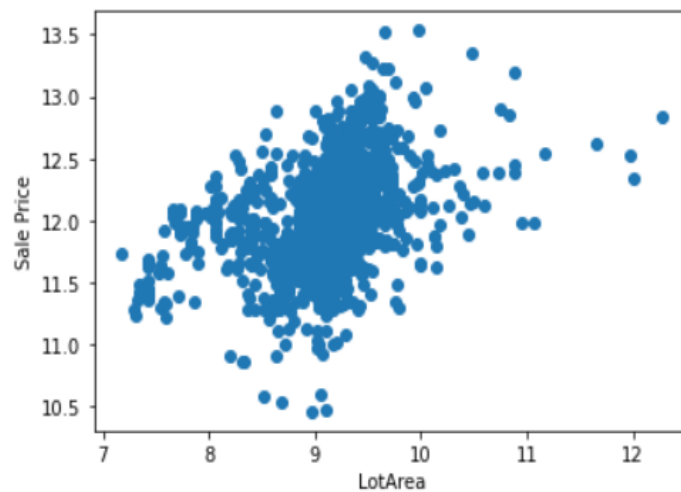
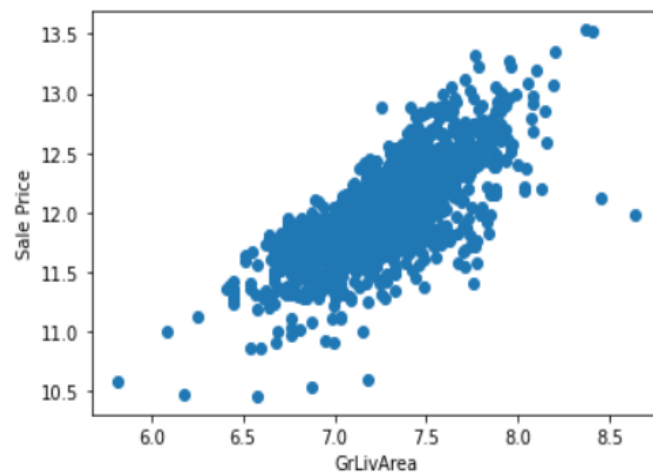
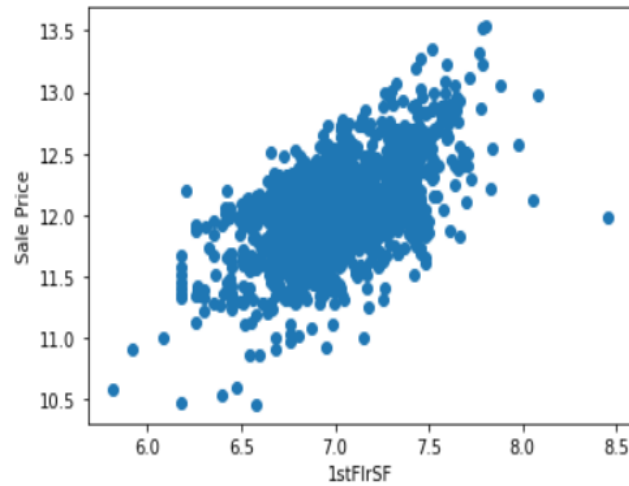
```
In [13]: ghar.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1460 entries, 0 to 1459
Data columns (total 35 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Id                     1460 non-null   int64
1   MSSubClass             1460 non-null   int64
2   LotArea                1460 non-null   int64
3   OverallQual            1460 non-null   int64
4   OverallCond            1460 non-null   int64
5   YearBuilt              1460 non-null   int64
6   YearRemodAdd           1460 non-null   int64
7   BsmtFinSF1             1460 non-null   int64
8   BsmtFinSF2             1460 non-null   int64
9   BsmtUnfSF              1460 non-null   int64
10  TotalBsmtSF            1460 non-null   int64
11  1stFlrSF               1460 non-null   int64
12  2ndFlrSF               1460 non-null   int64
13  LowQualFinSF           1460 non-null   int64
14  GrLivArea              1460 non-null   int64
15  BsmtFullBath           1460 non-null   int64
16  BsmtHalfBath           1460 non-null   int64
17  FullBath               1460 non-null   int64
18  HalfBath               1460 non-null   int64
19  BedroomAbvGr           1460 non-null   int64
20  KitchenAbvGr           1460 non-null   int64
21  TotRmsAbvGrd           1460 non-null   int64
22  Fireplaces             1460 non-null   int64
23  GarageCars             1460 non-null   int64
24  GarageArea             1460 non-null   int64
25  WoodDeckSF             1460 non-null   int64
26  OpenPorchSF            1460 non-null   int64
27  EnclosedPorch          1460 non-null   int64
28  3SsnPorch              1460 non-null   int64
29  ScreenPorch            1460 non-null   int64
30  PoolArea               1460 non-null   int64
31  MiscVal                1460 non-null   int64
32  MoSold                 1460 non-null   int64
33  YrSold                 1460 non-null   int64
34  SalePrice              1460 non-null   int64
dtypes: int64(35)
memory usage: 399.3 KB
```

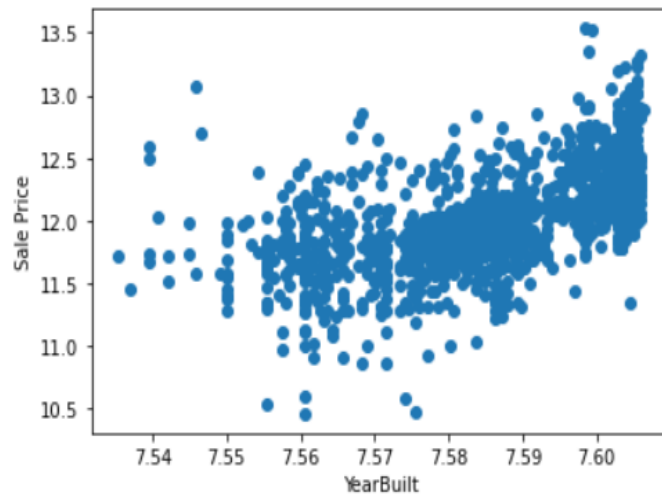
Features in the final dataset after dropping unwanted features

The code and its output given above shows the features that we have retained for our work. Now we can see we have 1460 non-null entries.

During EDA, we plot several graphs to see the co-relations in between Sale Price and other features. The following scatter plots represent the relationship of '1stFlrSF', 'GrLivArea', 'LotArea', 'YearBuilt' with 'Sale Price'.

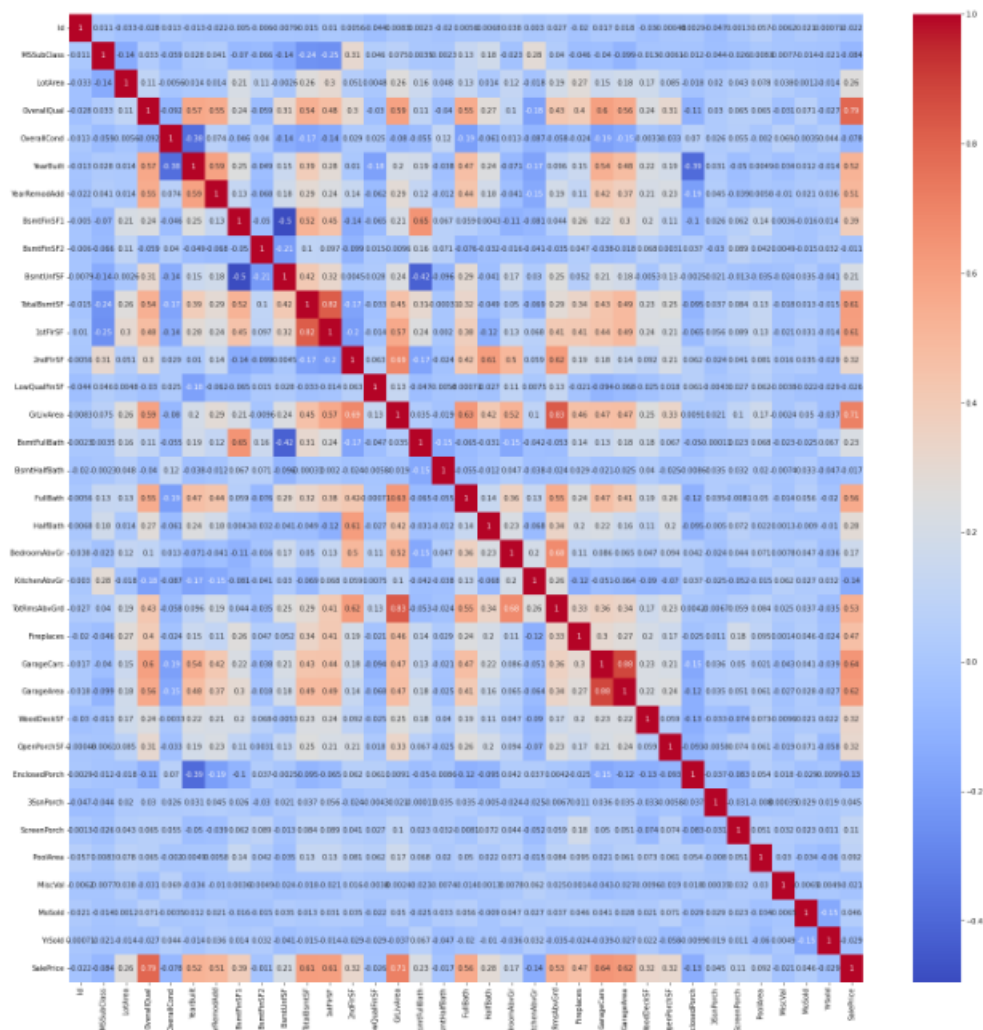


Results and Discussion



```
In [15]: plt.figure(figsize=(25,25))
sns.heatmap(ghar_corr(),annot=True,cmap='coolwarm')

Out[15]: <AxesSubplot>
```



Heat Map of features in the dataset

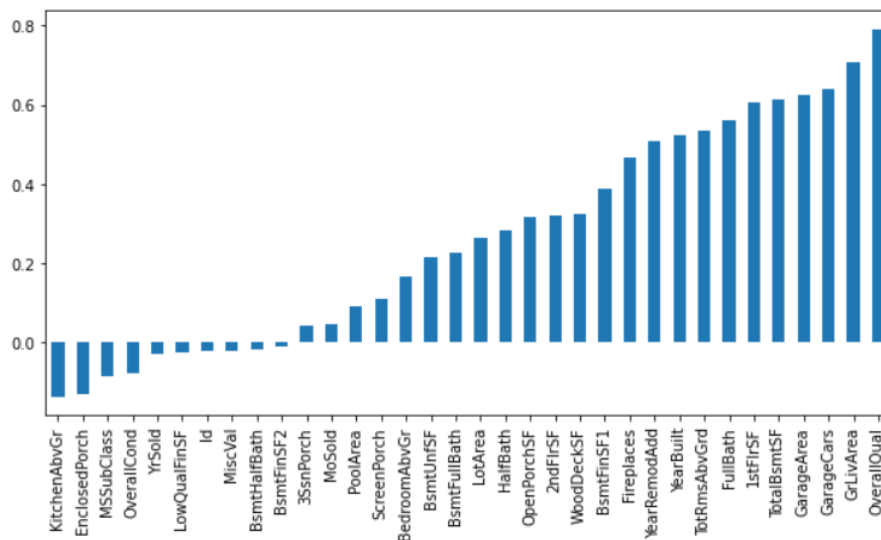
The above figure is known as heat-map, it can be seen as a **variance-covariance matrix**

between the features in the dataset. Deeper the red color more positive is the correlation between the features and deeper the blue color more will be the negative correlation between the features

Since we are interested in predicting the Sale Price of houses, we also plotted a **bar graph** to see the dependence of Sale Price with other features. The following graph depicts the same.

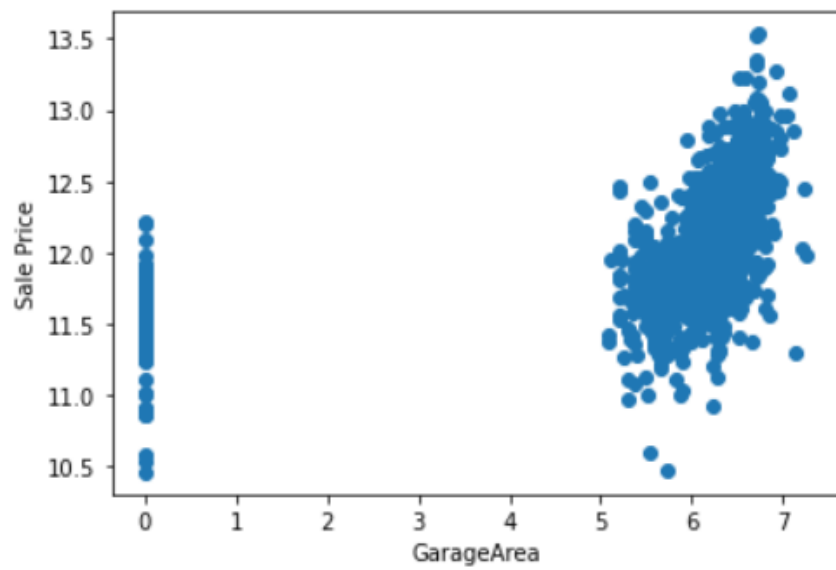
```
In [16]: plt.figure(figsize=(10,5))
ghar.corr()['SalePrice'].sort_values().drop('SalePrice').plot(kind='bar')

Out[16]: <AxesSubplot:>
```

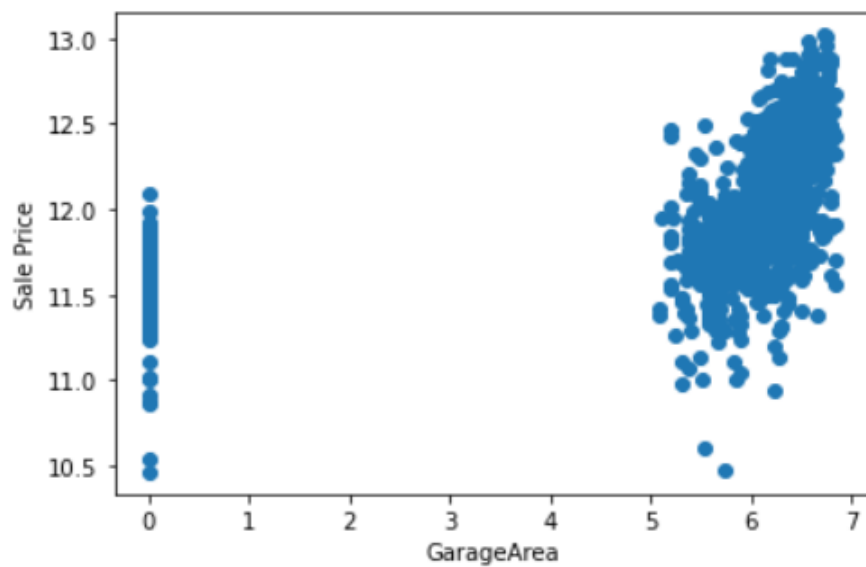


Bar graph of dependence of Sale Price upon features

Then we have removed outliers from certain features like 'GarageArea', 'GrLivArea', '1stFlrSF', 'TotalBsmtSF' and have visualised them by the help of box plots. Now, we have 1347 entries to work on with.



Before removing outliers



After removing outliers

Then, we split the dataset into two parts for training and testing ML models in the ratio of 3:1. Thus, 75% of data will be used for training the ML models and 25% for testing the accuracy of the predictions made by the ML models

Then we define some functions for our evaluation metrics using predefined functions from **sklearn** library. This helps us to avoid repetition of same block of codes again and again in the program.

evaluate() function is defined to evaluate all the four required metrics that we are going to use to compare the accuracy of ML models. **print_evaluate** function is defined to print the results.

```
1 import numpy as np
2 from sklearn import metrics
3
4 def print_evaluate(true, predicted):
5     mae = metrics.mean_absolute_error(true, predicted)
6     mse = metrics.mean_squared_error(true, predicted)
7     rmse = np.sqrt(metrics.mean_squared_error(true, predicted))
8     r2_square = metrics.r2_score(true, predicted)
9     print('MAE:', mae)
10    print('MSE:', mse)
11    print('RMSE:', rmse)
12    print('R2 Square', r2_square)
13    print('-----')
14
15 def evaluate(true, predicted):
16     mae = metrics.mean_absolute_error(true, predicted)
17     mse = metrics.mean_squared_error(true, predicted)
18     rmse = np.sqrt(metrics.mean_squared_error(true, predicted))
19     r2_square = metrics.r2_score(true, predicted)
20     return mae, mse, rmse, r2_square
```

Listing 4.1: Code for defining evaluation metrics

ML MODELS IMPLEMENTATION

We begin by **Multiple linear regression**. The code given below fits multiple linear regression by the help of training data that is X_train and y_train then, it is used to predict the Sale Price of the houses present in the testing dataset that is X_test. Finally by the help of our user-defined functions, we evaluate the accuracy of the multiple linear regression model and store it into a dataframe results_df so that we can compare all the results later, altogether.

```
1 from sklearn.linear_model import LinearRegression
2
3 lin_reg = LinearRegression(normalize=True)
4 lin_reg.fit(X_train,y_train)
5
6 test_pred = lin_reg.predict(X_test)
7 train_pred = lin_reg.predict(X_train)
8
9 print('Test set evaluation:\n-----')
10 print_evaluate(y_test, test_pred)
11 print('Train set evaluation:\n-----')
12 print_evaluate(y_train, train_pred)
13
14 results_df = pd.DataFrame(data=[["Linear Regression", *evaluate(
15     y_test, test_pred) ]],
16     columns=['Model', 'MAE', 'MSE', 'RMSE', 'R2 Square'])
```

Listing 4.2: Code for Multiple linear regression

The output that we got is as follows:

Test set evaluation:

MAE: 18031.09994805906
MSE: 654754851.4198682
RMSE: 25588.17796209547
R2 Square 0.8616161915916168

Train set evaluation:

MAE: 16594.941404038913
MSE: 508797862.3608397
RMSE: 22556.548103839817
R2 Square 0.8728078724387939

Then, we have applied **Ridge regression** after importing predefined functions from **sklearn.linear_model** library. Again we apply it to the training data and check the accuracy by applying our user-defined functions on testing data.

```
1 from sklearn.linear_model import Ridge, RidgeCV
2
3 alphas = np.geomspace(1e-9, 5, num=100)
4
5 ridgecv = RidgeCV(alphas = alphas, scoring = '
    neg_mean_squared_error', normalize = True)
6 ridgecv.fit(X_train, y_train)
7
8 ridge = Ridge(alpha = ridgecv.alpha_, normalize = True)
9 ridge.fit(X_train, y_train)
10
11 print('Ridge Regression:')
12 print("Alpha =", ridgecv.alpha_)
13
14 test_pred = ridge.predict(X_test)
15 train_pred = ridge.predict(X_train)
16
17 print('Test set evaluation:\n-----')
18 print_evaluate(y_test, test_pred)
19 print('Train set evaluation:\n-----')
20 print_evaluate(y_train, train_pred)
21
22 results_df_2 = pd.DataFrame(data=[["Ridge", *evaluate(y_test,
    test_pred)]],
23                             columns=['Model', 'MAE', 'MSE', 'RMSE', 'R2 Square'])
24 results_df = results_df.append(results_df_2, ignore_index=True)
```

Listing 4.3: Code for Ridge regression

The output that we got is as follows:

```
Ridge Regression:
Alpha = 0.03496578933827813
Test set evaluation:
-----
MAE: 17941.203909029748
MSE: 659197222.3203398
RMSE: 25674.836364042123
R2 Square 0.8606772872028419
-----
Train set evaluation:
-----
MAE: 16598.578018622953
MSE: 510812130.3043013
RMSE: 22601.1532958896
R2 Square 0.8723043345032797
-----
```

Then, we applied **LASSO** by importing predefined functions from **sklearn.linear_model** library. Again we apply it to the training data and check the accuracy by applying our user-defined functions on testing data.

```
1 from sklearn.linear_model import Lasso, LassoCV
2
3 lasso = Lasso(max_iter = 100000, normalize = True)
4
5 lassocv = LassoCV(alphas = None, cv = 10, max_iter = 100000,
6                   normalize = True)
7
8 lassocv.fit(X_train, y_train)
9
10 lasso.set_params(alpha=lassocv.alpha_)
11 lasso.fit(X_train, y_train)
12
13
14 test_pred = lasso.predict(X_test)
15 train_pred = lasso.predict(X_train)
16
17 print('Test set evaluation:\n_')
18 print_evaluate(y_test, test_pred)
19 print('Train set evaluation:\n_')
20 print_evaluate(y_train, train_pred)
21
22 print("Alpha =", lassocv.alpha_)
23
24 results_df_3 = pd.DataFrame(data=[["Lasso", *evaluate(y_test,
25                                     test_pred)]],
26                             columns=['Model', 'MAE', 'MSE', 'RMSE',
27                                     'R2 Square'])
28 results_df = results_df.append(results_df_3, ignore_index=True)
```

Listing 4.4: Code of LASSO

The output that we got is as follows:

Test set evaluation:

MAE: 17924.41488047078
MSE: 654162960.59647
RMSE: 25576.609638426864
R2 Square 0.8617412889561131

Train set evaluation:

MAE: 16622.69823507287
MSE: 513233502.30069155
RMSE: 22654.657408592422
R2 Square 0.871699026425905

Alpha = 11.219410925022853

Then, we applied **Random Forest Regressor** by importing predefined functions from **sklearn.linear_model** library. We apply it to the training data and check the accuracy by applying our user-defined functions on testing data.

```
1 from sklearn.ensemble import RandomForestRegressor
2
3 rf_reg = RandomForestRegressor(n_estimators=1000)
4 rf_reg.fit(X_train, y_train)
5
6 test_pred = rf_reg.predict(X_test)
7 train_pred = rf_reg.predict(X_train)
8
9 print('Test set evaluation:\n-----')
10 print_evaluate(y_test, test_pred)
11
12 print('Train set evaluation:\n-----')
13 print_evaluate(y_train, train_pred)
14
15 results_df_5 = pd.DataFrame(data=[["Random Forest Regressor", *
16     evaluate(y_test, test_pred)]],
17     columns=['Model', 'MAE', 'MSE', 'RMSE', 'R2 Square'])
18 results_df = results_df.append(results_df_5, ignore_index=True)
```

Listing 4.6: Code of Random Forest Regressor

The output that we got is as follows:

Test set evaluation:

MAE: 15986.833023738873
MSE: 556869679.7667505
RMSE: 23598.086358150962
R2 Square 0.8823044275179215

Train set evaluation:

MAE: 5893.06822079208
MSE: 72930242.0443673
RMSE: 8539.920494030803
R2 Square 0.9817684913098197

Then, we have applied **Support Vector Machine (Support Vector Regression**, since here it is used for regression) after importing predefined functions from **sklearn.linear_model** library. Again we apply it to the training data and check the accuracy by applying our user-defined functions on testing data.

```
1 from sklearn.svm import SVR
2
3 svm_reg = SVR(kernel='rbf', C=1000000, epsilon=0.001)
4 svm_reg.fit(X_train, y_train)
5
6 test_pred = svm_reg.predict(X_test)
7 train_pred = svm_reg.predict(X_train)
8
9 print('Test set evaluation:\n_ _ _ _ _')
10 print_evaluate(y_test, test_pred)
11
12 print('Train set evaluation:\n_ _ _ _ _')
13 print_evaluate(y_train, train_pred)
14
15 results_df_3 = pd.DataFrame(data=[["SVM Regressor", *evaluate(
16     y_test, test_pred)]],
17     columns=['Model', 'MAE', 'MSE', 'RMSE', 'R2 Square'])
18 results_df = results_df.append(results_df_3, ignore_index=True)
```

Listing 4.7: Code of SVM

The output that we got is as follows:

Test set evaluation:

MAE: 23228.650822315394
MSE: 1136136310.904815
RMSE: 33706.62117306947
R2 Square 0.7598752124812613

Train set evaluation:

MAE: 21057.707781303285
MSE: 913885321.5834868
RMSE: 30230.536243730225
R2 Square 0.7715418499208946

Finally, we used Artificial Neural Network

```
1 from tensorflow.keras.models import Sequential
2 from tensorflow.keras.layers import Input, Dense, Activation,
  Dropout
3 from tensorflow.keras.optimizers import Adam
4 X_train = np.array(X_train)
5 X_test = np.array(X_test)
6 y_train = np.array(y_train)
7 y_test = np.array(y_test)
8
9 model = Sequential()
10
11 model.add(Dense(X_train.shape[1], activation='relu'))
12 model.add(Dense(32, activation='relu'))
13 # model.add(Dropout(0.2))
14
15 model.add(Dense(64, activation='relu'))
16 # model.add(Dropout(0.2))
17
18 model.add(Dense(128, activation='relu'))
19 model.add(Dropout(0.2))
20 model.add(Dense(1))
21
22 model.compile(optimizer=Adam(0.001), loss='mse')
23
24 r = model.fit(X_train, y_train,
25               validation_data=(X_test, y_test),
26               batch_size=128,
27               epochs=50)
28
29 test_pred = model.predict(X_test)
30 train_pred = model.predict(X_train)
31
32 print('Test set evaluation:\n_____')
33 print_evaluate(y_test, test_pred)
34
35 print('Train set evaluation:\n_____')
36 print_evaluate(y_train, train_pred)
37
38 results_df_4 = pd.DataFrame(data=[["Artificial Neural Network", *
39                                   evaluate(y_test, test_pred)]],
40                             columns=['Model', 'MAE', 'MSE', 'RMSE',
41                                     'R2 Square'])
42 results_df = results_df.append(results_df_4, ignore_index=True)
```

Listing 4.8: Code of ANN

The output that we got is as follows:

Test set evaluation:

MAE: 24946.108772255193
MSE: 1232085008.6104596
RMSE: 35101.068482461604
R2 Square 0.7395962543772652

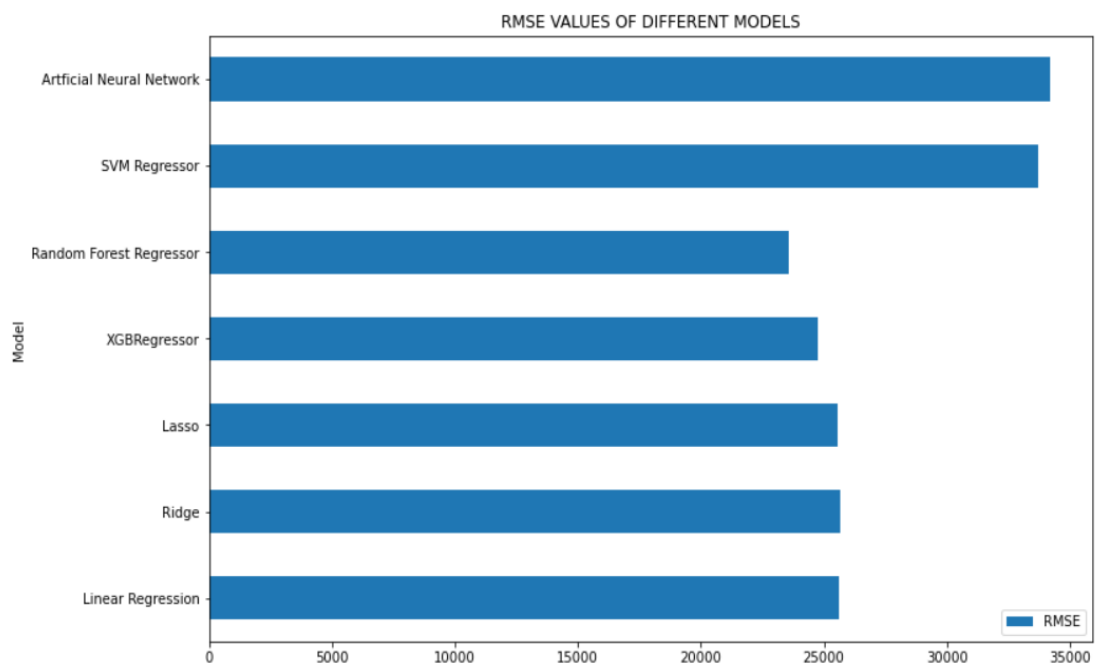
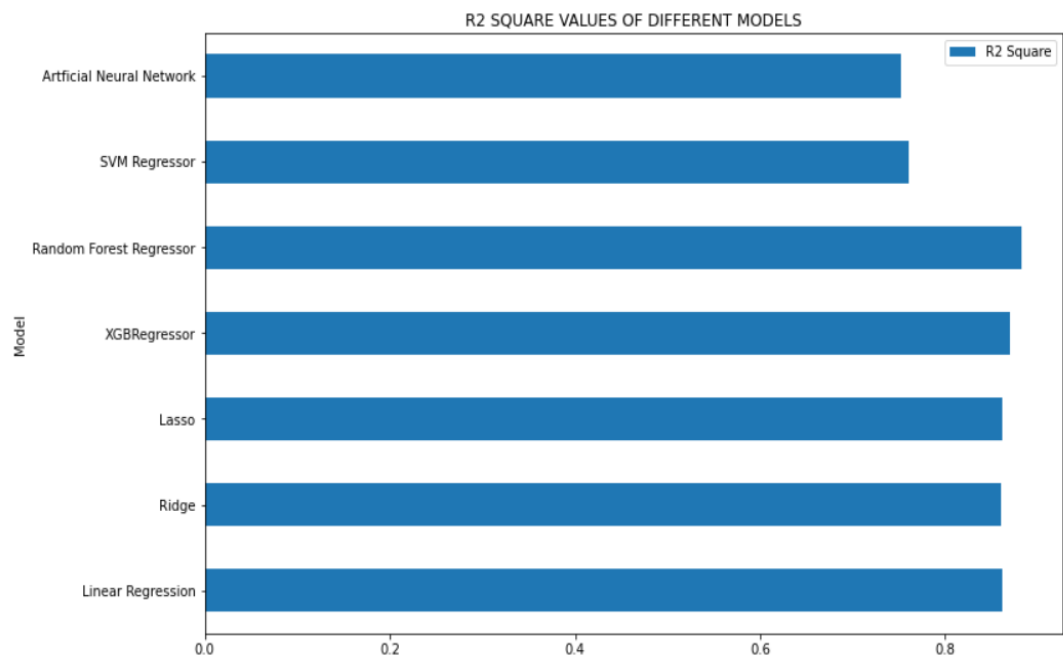
Train set evaluation:

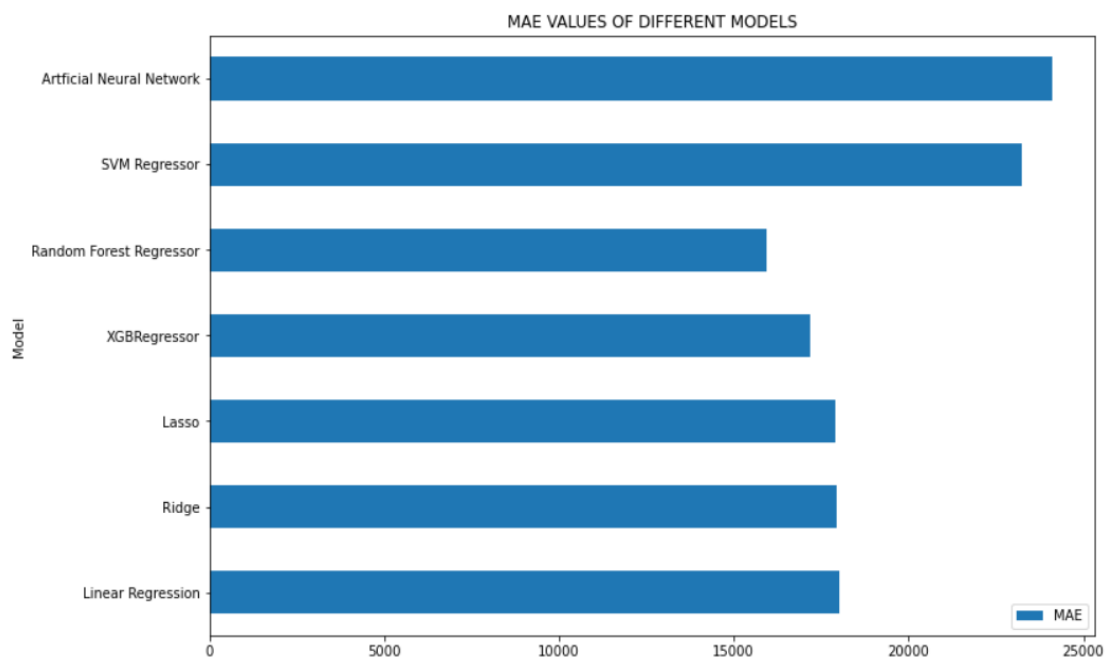
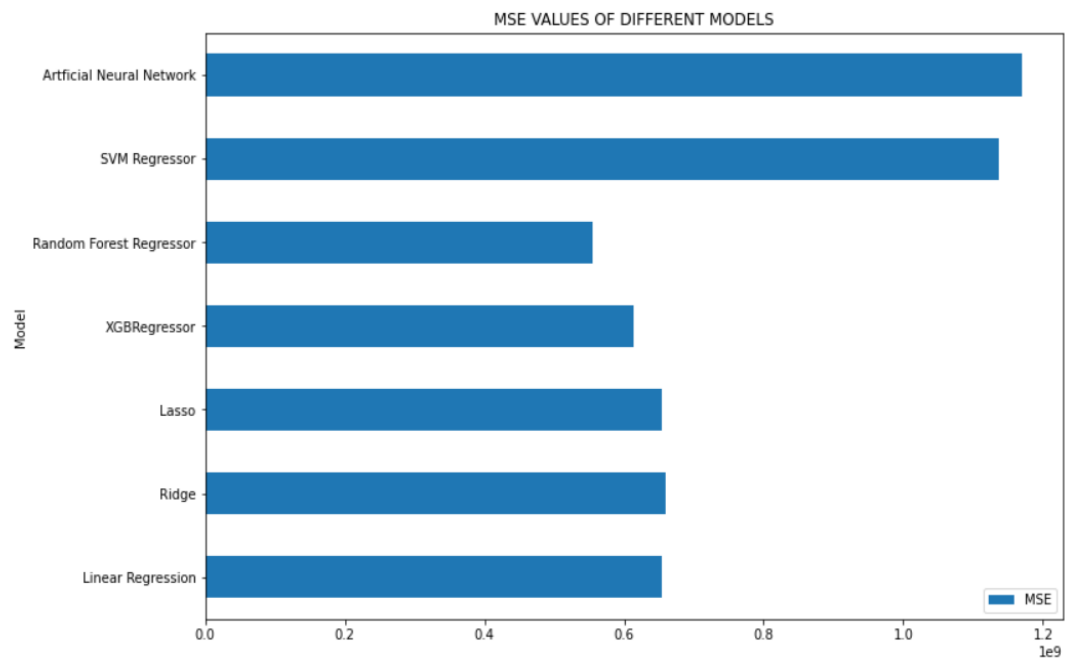
MAE: 22945.703991336635
MSE: 985160327.001907
RMSE: 31387.263770547233
R2 Square 0.7537241265149027

After using these ML models, we finally pictorially represent the values of evaluation metrics that we obtained. If we compare them, we can see that **Random Forest Regressor** has the **maximum R^2 Square value** and **minimum Mean Absolute Error (MAE), Mean Squared Error (MSE), RMSE (Root Mean Square Error)**.

```
1 results_df.plot(x='Model', y='R2 Square', kind='barh', figsize=(12, 8), title="R2 SQUARE VALUES OF DIFFERENT MODELS")
2
3 results_df.plot(x='Model', y='RMSE', kind='barh', figsize=(12, 8), title="RMSE VALUES OF DIFFERENT MODELS")
4
5 results_df.plot(x='Model', y='MSE', kind='barh', figsize=(12, 8), title="MSE VALUES OF DIFFERENT MODELS")
6
7 results_df.plot(x='Model', y='MAE', kind='barh', figsize=(12, 8), title="MAE VALUES OF DIFFERENT MODELS")
```

Listing 4.9: Code for pictorial representation of evaluation metrics





From the above bar graphs, we can see that **Random Forest Regressor** performs better than other ML model used in the project.

Chapter 5

Conclusion

After running the python program, we compared the ML models that we had taken. Finally, the result we got is as follows:

	Model	MAE	MSE	RMSE	R2 Square
0	Linear Regression	18031.099948	6.547549e+08	25588.177962	0.861616
1	Ridge	17941.203909	6.591972e+08	25674.836364	0.860677
2	Lasso	17924.414880	6.541630e+08	25576.609638	0.861741
3	XGBRegressor	17178.648901	6.135903e+08	24770.755409	0.870316
4	Random Forest Regressor	15921.487059	5.546648e+08	23551.321735	0.882770
5	SVM Regressor	23228.650822	1.136136e+09	33706.621173	0.759875
6	Artificial Neural Network	24120.632789	1.170360e+09	34210.526267	0.752642

From the results, we are getting a MAE of about \$ 15,921 in the prediction by Random Forest regressor that is the least among all models used in this project. This means on an average the predictions of house price differ from the actual price by \$ 15,921. Whereas the highest MAE we got is of \$ 24,120 , that was from ANN which makes it the poorest ML model used for prediction, at least under these conditions that were taken in this project. Usually, ANN requires a lot of training data to get a good accuracy during testing.

The final results of this project showed that Random Forest regressor makes better prediction as compared to other used algorithms since it gives us the least MAE, MSE, RMSE and the highest R^2 Square values.

Although this study has shown that **Random forest Regressor** makes the best pre-

diction, one cannot guarantee that it will perform the same when used for other purposes than the ones that have been presented in this study. Infact, if we use different parameters for ML models like SVM, ANN, then their performance might increase too. Also, we can model the data, instead of deleting the features having some NULL entries, we could fill them up with the mean or mode of all the observed values so that we get more number of predictors. This can also boost the accuracy of some ML models.

5.1 Ethics

This project is taking into consideration the ethical part, where the dataset is downloaded from a public website called Kaggle. Kaggle allows users to find and publish data sets, explore and build models in a web-based data-science environment, and work with other data scientists. In this project, the algorithms are public and open source. In addition, the algorithms are trained and tested on the same public dataset.

5.2 Future Work

There is a great scope for improvement in accuracy of prediction. Future work on this project could be done by:

- Changing the way we have handled missing values. We can fill them up by mean or mode of observations present in that column. By this we will have more number of predictors which can also boost the accuracy of some ML models.
- Making use of categorical features which have not used in this project.
- Tweeking with the ANN model by using different parameters that can help it train the model better.
- The used pre-processing methods do help in the prediction accuracy. However, experimenting with different combinations of pre-processing methods to achieve better prediction accuracy.

Appendix A

Features

The list of all features present in the dataset is as follows:

FEATURE	DESCRIPTION
SalePrice	the property's sale price in dollars
MSSubClass	The building class
MSZoning	The general zoning classification
LotFrontage	Linear feet of street connected to property
LotArea	Lot size in square feet
Street	Type of road access
Alley	Type of alley access
LotShape	General shape of property
LandContour	Flatness of the property
Utilities	Type of utilities available
LotConfig	Lot configuration
LandSlope	Slope of property
Neighborhood	Physical locations within Ames city limits
Condition1	Proximity to main road or railroad
Condition2	Proximity to main road or railroad (if a second is present)
BldgType	Type of dwelling
HouseStyle	Style of dwelling
OverallQual	Overall material and finish quality
OverallCond	Overall condition rating
YearBuilt	Original construction date
YearRemodAdd	Remodel date
RoofStyle	Type of roof
RoofMatl	Roof material

FEATURE	DESCRIPTION
Exterior1st	Exterior covering on house
Exterior2nd	Exterior covering on house (if more than one material)
MasVnrType	Masonry veneer type
MasVnrArea	Masonry veneer area in square feet
ExterQual	Exterior material quality
ExterCond	Present condition of the material on the exterior
Foundation	Type of foundation
BsmtQual	Height of the basement
BsmtCond	General condition of the basement
BsmtExposure	Walkout or garden level basement walls
BsmtFinType1	Quality of basement finished area
BsmtFinSF1	Type 1 finished square feet
BsmtFinType2	Quality of second finished area (if present)
BsmtFinSF2	Type 2 finished square feet
BsmtUnfSF	Unfinished square feet of basement area
TotalBsmtSF	Total square feet of basement area
Heating	Type of heating
HeatingQC	Heating quality and condition
CentralAir	Central air conditioning
Electrical	Electrical system
1stFlrSF	First Floor square feet
2ndFlrSF	Second floor square feet
LowQualFinSF	Low quality finished square feet (all floors)
GrLivArea	Above grade (ground) living area square feet
BsmtFullBath	Basement full bathrooms
BsmtHalfBath	Basement half bathrooms
FullBath	Full bathrooms above grade
HalfBath	Half baths above grade
Bedroom	Number of bedrooms above basement level
Kitchen	Number of kitchens
KitchenQual	Kitchen quality
TotRmsAbvGrd	Total rooms above grade (does not include bathrooms)
Functional	Home functionality rating
Fireplaces	Number of fireplaces

FEATURE	DESCRIPTION
FireplaceQu	Fireplace quality
GarageType	Garage location
GarageYrBlt	Year garage was built
GarageFinish	Interior finish of the garage
GarageCars	Size of garage in car capacity
GarageArea	Size of garage in square feet
GarageQual	Garage quality
GarageCond	Garage condition
PavedDrive	Paved driveway
WoodDeckSF	Wood deck area in square feet
OpenPorchSF	Open porch area in square feet
EnclosedPorch	Enclosed porch area in square feet
3SsnPorch	Three season porch area in square feet
ScreenPorch	Screen porch area in square feet
PoolArea	Pool area in square feet
PoolQC	Pool quality
Fence	Fence quality
MiscFeature	Miscellaneous feature not covered in other categories
MiscVal	Value of miscellaneous feature
MoSold	Month Sold
YrSold	Year Sold
SaleType	Type of sale
SaleCondition	Condition of sale

Appendix B

Python Code

Below is the full python code that we have made to implement and evaluate the accuracy of various ML models.

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5
6 ghar = pd.read_csv(r"C:\Users\PRANAB\Desktop\B.Sc FINAL YR PROJECT\
    OTHERS\house-prices-advanced-regression-techniques\train.csv")
7 ghar.head()
8
9 ghar.info()
10
11 ghar.describe()
12
13 ghar.columns
14
15
16 numerical_feats = ghar.dtypes[ghar.dtypes != "object"].index
17 print("Number of Numerical features: ", len(numerical_feats))
18
19 categorical_feats = ghar.dtypes[ghar.dtypes == "object"].index
20 print("Number of Categorical features: ", len(categorical_feats))
21
22
23
24 print(ghar[numerical_feats].columns)
25 print("="*100)
26 print(ghar[categorical_feats].columns)
27
```



```
28
29 for i in categorical_feats:
30     ghar.drop(i, inplace=True, axis=1)
31
32 ghar.columns
33
34 ghar.shape
35
36 total = ghar.isnull().sum().sort_values(ascending=False)
37 percent = (ghar.isnull().sum()/ghar.isnull().count()).sort_values(
38     ascending=False)
39 missing_data = pd.concat([total, percent], axis=1, keys=['Total', '
40     Percent'])
41 missing_data.head(10)
42
43 ghar.drop(['LotFrontage', 'GarageYrBlt', 'MasVnrArea'], axis=1, inplace
44     =True)
45
46 ghar.info()
47
48 #EDA
49
50 features= []
51 for i in ghar.columns:
52     features.append(i)
53
54 for feature in features:
55     data=ghar.copy()
56     data['SalePrice']=np.log(data['SalePrice']+1)
57     data[feature]=np.log(data[feature]+1)
58     plt.scatter(data[feature], data['SalePrice'])
59     plt.xlabel(feature)
60     plt.ylabel('Sale Price')
61     plt.show()
62
63 plt.figure(figsize=(25,25))
64 sns.heatmap(ghar.corr(), annot=True, cmap='coolwarm')
65
66 plt.figure(figsize=(10,5))
67 ghar.corr()['SalePrice'].sort_values().drop('SalePrice').plot(kind=
68     'bar')
69
70 import warnings
71 warnings.filterwarnings("ignore")
72
73 #DATA PREPROCESSING AGAIN OUTLIER
```

```
71 sns.boxplot(ghar['GarageArea'])
72 plt.title("Box Plot before outlier removing")
73 plt.show()
74
75 def drop_outliers(df, field_name):
76     iqr = 1.5 * (np.percentile(df[field_name], 75) - np.percentile(
77         df[field_name], 25))
78     df.drop(df[df[field_name] > (iqr + np.percentile(df[field_name]
79         ], 75))].index, inplace=True)
80     df.drop(df[df[field_name] < (np.percentile(df[field_name], 25)
81         - iqr)].index, inplace=True)
82 drop_outliers(ghar, 'GarageArea')
83 sns.boxplot(ghar['GarageArea'])
84 plt.title("Box Plot after outlier removing")
85 plt.show()
86
87 sns.boxplot(ghar['GrLivArea'])
88 plt.title("Box Plot before outlier removing")
89 plt.show()
90
91 def drop_outliers(df, field_name):
92     iqr = 1.5 * (np.percentile(df[field_name], 75) - np.percentile(
93         df[field_name], 25))
94     df.drop(df[df[field_name] > (iqr + np.percentile(df[field_name]
95         ], 75))].index, inplace=True)
96     df.drop(df[df[field_name] < (np.percentile(df[field_name], 25)
97         - iqr)].index, inplace=True)
98 drop_outliers(ghar, 'GrLivArea')
99 sns.boxplot(ghar['GrLivArea'])
100 plt.title("Box Plot after outlier removing")
101 plt.show()
102
103 sns.boxplot(ghar['1stFlrSF'])
104 plt.title("Box Plot before outlier removing")
105 plt.show()
106
107 def drop_outliers(df, field_name):
108     iqr = 1.5 * (np.percentile(df[field_name], 75) - np.percentile(
109         df[field_name], 25))
110     df.drop(df[df[field_name] > (iqr + np.percentile(df[field_name]
111         ], 75))].index, inplace=True)
112     df.drop(df[df[field_name] < (np.percentile(df[field_name], 25)
113         - iqr)].index, inplace=True)
```

```
109 drop_outliers(ghar, '1stFlrSF')
110 sns.boxplot(ghar['1stFlrSF'])
111 plt.title("Box Plot after outlier removing")
112 plt.show()
113
114
115
116 sns.boxplot(ghar['TotalBsmtSF'])
117 plt.title("Box Plot before outlier removing")
118 plt.show()
119
120 def drop_outliers(df, field_name):
121     iqr = 1.5 * (np.percentile(df[field_name], 75) - np.percentile(
122         df[field_name], 25))
123     df.drop(df[df[field_name] > (iqr + np.percentile(df[field_name]
124         ], 75))].index, inplace=True)
125     df.drop(df[df[field_name] < (np.percentile(df[field_name], 25)
126         - iqr)].index, inplace=True)
127 drop_outliers(ghar, 'TotalBsmtSF')
128 sns.boxplot(ghar['TotalBsmtSF'])
129 plt.title("Box Plot after outlier removing")
130 plt.show()
131
132
133
134 ghar.info()
135
136 #AGAIN VISUALISATION
137
138 plt.figure(figsize=(10,5))
139 ghar.corr()['SalePrice'].sort_values().drop('SalePrice').plot(kind=
140     'bar')
141
142
143
144 X = ghar.loc[:, ghar.columns!='SalePrice']
145
146 y = ghar['SalePrice']
147
148
149 from sklearn.model_selection import train_test_split
150
151 X_train, X_test, y_train, y_test = train_test_split(X, y,
152     random_state=0, train_size = .75)
153
154
155 from sklearn import metrics
156
157
158 def print_evaluate(true, predicted):
159     mae = metrics.mean_absolute_error(true, predicted)
```

```

151     mse = metrics.mean_squared_error(true, predicted)
152     rmse = np.sqrt(metrics.mean_squared_error(true, predicted))
153     r2_square = metrics.r2_score(true, predicted)
154     print('MAE:', mae)
155     print('MSE:', mse)
156     print('RMSE:', rmse)
157     print('R2 Square', r2_square)
158     print('-----')
159
160 def evaluate(true, predicted):
161     mae = metrics.mean_absolute_error(true, predicted)
162     mse = metrics.mean_squared_error(true, predicted)
163     rmse = np.sqrt(metrics.mean_squared_error(true, predicted))
164     r2_square = metrics.r2_score(true, predicted)
165     return mae, mse, rmse, r2_square
166
167
168 #MODELS
169 #LINEAR REGRESSION
170 from sklearn.linear_model import LinearRegression
171
172 lin_reg = LinearRegression(normalize=True)
173 lin_reg.fit(X_train,y_train)
174
175 test_pred = lin_reg.predict(X_test)
176 train_pred = lin_reg.predict(X_train)
177
178 print('Test set evaluation:\n-----')
179
180 print_evaluate(y_test, test_pred)
181 print('Train set evaluation:\n-----')
182
183 print_evaluate(y_train, train_pred)
184
185 results_df = pd.DataFrame(data=[["Linear Regression", *evaluate(
186     y_test, test_pred) ]],
187     columns=['Model', 'MAE', 'MSE', 'RMSE', '
188     R2 Square'])
189
190 #RIDGE
191 from sklearn.linear_model import Ridge, RidgeCV
192
193 alphas = np.geomspace(1e-9, 5, num=100)
194
195 ridgecv = RidgeCV(alphas = alphas, scoring = '
196     neg_mean_squared_error', normalize = True)

```

```
193 ridgecv.fit(X_train, y_train)
194
195 ridge = Ridge(alpha = ridgecv.alpha_, normalize = True)
196 ridge.fit(X_train, y_train)
197
198 print('Ridge Regression:')
199 print("Alpha =", ridgecv.alpha_)
200
201 test_pred = ridge.predict(X_test)
202 train_pred = ridge.predict(X_train)
203
204 print('Test set evaluation:\n_-----',
      )
205 print_evaluate(y_test, test_pred)
206 print('Train set evaluation:\n_-----',
      )
207 print_evaluate(y_train, train_pred)
208
209
210
211 results_df_2 = pd.DataFrame(data=[["Ridge", *evaluate(y_test,
      test_pred)]],
212                             columns=['Model', 'MAE', 'MSE', 'RMSE',
      'R2 Square'])
213 results_df = results_df.append(results_df_2, ignore_index=True)
214
215
216
217 #LASSO
218 from sklearn.linear_model import Lasso, LassoCV
219
220 lasso = Lasso(max_iter = 100000, normalize = True)
221
222 lassocv = LassoCV(alphas = None, cv = 10, max_iter = 100000,
      normalize = True)
223 lassocv.fit(X_train, y_train)
224
225 lasso.set_params(alpha=lassocv.alpha_)
226 lasso.fit(X_train, y_train)
227
228
229 test_pred = lasso.predict(X_test)
230 train_pred = lasso.predict(X_train)
231
232 print('Test set evaluation:\n_-----',
      )
233 print_evaluate(y_test, test_pred)
```

```
234 print('Train set evaluation:\n_-----')
235 print_evaluate(y_train, train_pred)
236
237
238 print("Alpha =", lasso_cv.alpha_)
239
240 results_df_3 = pd.DataFrame(data=[["Lasso", *evaluate(y_test,
241                                     test_pred)]],
242                             columns=['Model', 'MAE', 'MSE', 'RMSE',
243                                     'R2 Square'])
244 results_df = results_df.append(results_df_3, ignore_index=True)
245
246 #XGB
247 from numpy import loadtxt
248 from xgboost import XGBRegressor
249
250 # fit model on training data
251 model = XGBRegressor()
252 model.fit(X_train, y_train)
253 # make predictions for test data
254 #test_pred = model.predict(X_test)
255 #predictions = [round(value) for value in y_pred]
256 # evaluate predictions
257
258 test_pred = model.predict(X_test)
259 train_pred = model.predict(X_train)
260
261 print('Test set evaluation:\n_-----')
262 )
263 print_evaluate(y_test, test_pred)
264
265 print('Train set evaluation:\n_-----')
266 )
267 print_evaluate(y_train, train_pred)
268
269 results_df_4 = pd.DataFrame(data=[["XGBRegressor", *evaluate(y_test,
270                                     test_pred)]],
271                             columns=['Model', 'MAE', 'MSE', 'RMSE',
272                                     'R2 Square'])
273 results_df = results_df.append(results_df_4, ignore_index=True)
274
275 #RANDOM FOREST REGRESSOR
276 from sklearn.ensemble import RandomForestRegressor
277
```

```
274 rf_reg = RandomForestRegressor(n_estimators=1000)
275 rf_reg.fit(X_train, y_train)
276
277 test_pred = rf_reg.predict(X_test)
278 train_pred = rf_reg.predict(X_train)
279
280 print('Test set evaluation:\n-----',
      )
281 print_evaluate(y_test, test_pred)
282
283 print('Train set evaluation:\n-----',
      )
284 print_evaluate(y_train, train_pred)
285
286 results_df_5 = pd.DataFrame(data=[["Random Forest Regressor", *
      evaluate(y_test, test_pred)]],
287                             columns=['Model', 'MAE', 'MSE', 'RMSE',
      'R2 Square'])
288 results_df = results_df.append(results_df_5, ignore_index=True)
289
290 #SVM
291 from sklearn.svm import SVR
292
293 svm_reg = SVR(kernel='rbf', C=1000000, epsilon=0.001)
294 svm_reg.fit(X_train, y_train)
295
296 test_pred = svm_reg.predict(X_test)
297 train_pred = svm_reg.predict(X_train)
298
299 print('Test set evaluation:\n-----',
      )
300 print_evaluate(y_test, test_pred)
301
302 print('Train set evaluation:\n-----',
      )
303 print_evaluate(y_train, train_pred)
304
305 results_df_6 = pd.DataFrame(data=[["SVM Regressor", *evaluate(
      y_test, test_pred)]],
306                             columns=['Model', 'MAE', 'MSE', 'RMSE',
      'R2 Square'])
307 results_df = results_df.append(results_df_6, ignore_index=True)
308
309 #ANN
310
311 from tensorflow.keras.models import Sequential
```

```
312 from tensorflow.keras.layers import Input, Dense, Activation,
    Dropout
313 from tensorflow.keras.optimizers import Adam
314 X_train = np.array(X_train)
315 X_test = np.array(X_test)
316 y_train = np.array(y_train)
317 y_test = np.array(y_test)
318
319 model = Sequential()
320
321 model.add(Dense(X_train.shape[1], activation='relu'))
322 model.add(Dense(32, activation='relu'))
323 # model.add(Dropout(0.2))
324
325 model.add(Dense(64, activation='relu'))
326 # model.add(Dropout(0.2))
327
328 model.add(Dense(128, activation='relu'))
329 model.add(Dropout(0.2))
330 model.add(Dense(1))
331
332 model.compile(optimizer=Adam(0.001), loss='mse')
333
334 r = model.fit(X_train, y_train,
335               validation_data=(X_test, y_test),
336               batch_size=128,
337               epochs=150)
338
339 test_pred = model.predict(X_test)
340 train_pred = model.predict(X_train)
341
342 print('Test set evaluation:\n-----',
343       )
344 print_evaluate(y_test, test_pred)
345
346 print('Train set evaluation:\n-----',
347       )
348 print_evaluate(y_train, train_pred)
349
350 results_df_7 = pd.DataFrame(data=[["Artificial Neural Network", *
351                                   evaluate(y_test, test_pred)]],
352                             columns=['Model', 'MAE', 'MSE', 'RMSE',
353                                     'R2 Square'])
354 results_df = results_df.append(results_df_7, ignore_index=True)
355
356 #Models comparision
```



```
354 results_df.plot(x='Model', y='R2 Square', kind='barh', figsize=(12,
    8), title="R2 SQUARE VALUES OF DIFFERENT MODELS")
355
356 results_df.plot(x='Model', y='RMSE', kind='barh', figsize=(12, 8),
    title="RMSE VALUES OF DIFFERENT MODELS")
357
358 results_df.plot(x='Model', y='MSE', kind='barh', figsize=(12, 8),
    title="MSE VALUES OF DIFFERENT MODELS")
359
360 results_df.plot(x='Model', y='MAE', kind='barh', figsize=(12, 8),
    title="MAE VALUES OF DIFFERENT MODELS")
361
362
363
364 results_df
```

Listing B.1: Complete Python code

Appendix C

Github Link

The link to GitHub repository to access the source code <https://github.com/git-pkp/MLpro1>

Bibliography

- [1] ai explained "what is ai? learn about artificial intelligence". <https://www.oracle.com/in/artificial-intelligence/what-is-ai/>.
- [2] Artificial Intelligence (AI) Explained "what is ai? learn about artificial intelligence". <https://www.oracle.com/in/artificial-intelligence/what-is-ai/>.
- [3] Dataset "house prices - advanced regression techniques". <https://www.kaggle.com/competitions/house-prices-advanced-regression-techniques/data>.
- [4] Outlier "what are outliers and how to deal with them?". <https://medium.com/analytics-vidhya/how-to-remove-outliers-for-machine-learning-24620c4657e8>.
- [5] RFR "random forest regression: When does it fail and why?". <https://neptune.ai/blog/random-forest-regression-when-does-it-fail-and-why>.
- [6] SVR. [https://www.geeksforgeeks.org/support-vector-machine-algorithm/#:~:text=Support%20Vector%20Machine\(SVM\)%20is,distinctly%20classifies%20the%20data%20points](https://www.geeksforgeeks.org/support-vector-machine-algorithm/#:~:text=Support%20Vector%20Machine(SVM)%20is,distinctly%20classifies%20the%20data%20points).
- [7] SVR. <https://www.analyticsvidhya.com/blog/2020/03/support-vector-regression-tutorial-for-machine-learning/>.
- [8] XGB "xgboost: A deep dive into boosting". <https://dzone.com/articles/xgboost-a-deep-dive-into-boosting>.
- [9] XGBoost Documentation. <https://xgboost.readthedocs.io/en/stable/>.