



e-ISSN: 2319-8753 | p-ISSN: 2347-6710

IJRSET

International Journal of Innovative Research in
SCIENCE | ENGINEERING | TECHNOLOGY

INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY

Volume 12, Special Issue 2, April 2023

ISSN INTERNATIONAL
STANDARD
SERIAL
NUMBER
INDIA

9940 572 462

6381 907 438

ijirset@gmail.com

www.ijirset.com

Prediction of House Price Using SMLT

Mrs Sruthi Nath¹, Nithi Sneha², Vandana S³

Assistant Professor, Department of Computer Science Engineering, Velammal Engineering College, Chennai,
Tamil Nadu, India¹

U.G. Students, Department of Computer Science and Engineering, Velammal Engineering College, Chennai,
Tamil Nadu, India²⁻⁴

ABSTRACT : Generally, predicting how the house price will perform is one of the most difficult things to do. It can be described as one of the most critical process to predict that. This is a very complex task and has uncertainties. To prevent this problem in One of the most interesting (or perhaps most profitable) time series data using machine learning techniques. Hence, house price prediction has become an important research area. The aim is to predict machine learning based techniques for house price prediction results in error based calculation. The analysis of dataset by supervised machine learning technique(SMLT) to capture several information's like, variable identification, uni-variate analysis, bi-variant and multi-variant analysis, missing value treatments and analyze the data validation, data cleaning/preparing and data visualization will be done on the entire given dataset. To propose a machine learning-based method to accurately predict the house price Index value by prediction results in the form of house price increase or stable state best regression from comparing supervise machine learning algorithms. Additionally, to compare and discuss the performance of various machine learning algorithms. dataset with evaluation classification report, identify the confusion matrix and to categorizing data from priority and the result shows that the effectiveness of the proposed machine learning algorithm technique can be compared with best accuracy MAE,MSE,R².

I. INTRODUCTION

I.

Data science is an interdisciplinary field that uses scientific methods, processes, algorithms and systems to extract knowledge and insights from structured and unstructured data, and apply knowledge and actionable insights from data across a broad range of application domains. Data science can be defined as a blend of mathematics, business acumen, tools, algorithms and machine learning techniques, all of which help us in finding out the hidden insights or patterns from raw data which can be of major use in the formation of big business decisions.

II. RELATED WORK

Urban housing price is widely accepted as an economic indicator which is of both business and research interest in urban computing. However, due to the complex nature of influencing factors and the sparse property of transaction records, to implement such a model is still challenging. To address these challenges, in this work, we study an effective and fine-grained model for urban subregion housing price predictions. Compared to existing works, our proposal improves the forecasting granularity from city-level to mile-level, with only publicly released transaction data. We employ a feature selection mechanism to select more relevant features. a fine-grained forecasting model, JGC MMN, for subregion spatiotemporal housing price prediction. In particular, we modify the structure of DenseNet and adopt the method of bagging by fusing the KF-based method to improve the accuracy.

III. METHODOLOGY

Data Pre-processing:

Validation techniques in machine learning are used to get the error rate of the Machine Learning (ML) model, which can be considered as close to the true error rate of the dataset. If the data volume is large enough to be representative of the population, you may not need the validation techniques. However, in real-world scenarios, to work with samples of data that may not be a true representative of the population of given dataset. To finding the missing value, duplicate value and description of data type whether it is float variable or integer. The sample of data used to provide an unbiased evaluation of a model fit Data Processing Test dataset Training dataset Regression ML Algorithm Model on the training dataset while tuning model hyper parameters. The evaluation becomes more biased as skill on the validation



dataset is incorporated into the model configuration. The validation set is used to evaluate a given model, but this is for frequent evaluation. It as machine learning engineers use this data to fine-tune the model hyper parameters. Data collection, data analysis, and the process of addressing data content, quality, and structure can add up to a time-consuming to-do list. During the process of data identification, it helps to understand your data and its properties; this knowledge will help you choose which algorithm to use to build your model.

A number of different data cleaning tasks using Python's Pandas library and specifically, it focus on probably the biggest data cleaning task, missing values and it able to more quickly clean data. It wants to spend less time cleaning data, and more time exploring and modeling.

Some of these sources are just simple random mistakes. Other times, there can be a deeper reason why data is missing. It's important to understand these different types of missing data from a statistics point of view. The type of missing data will influence how to deal with filling in the missing values and to detect missing values, and do some basic imputation and detailed statistical approach for dealing with missing data. Before, joint into code, it's important to understand the sources of missing data. Here are some typical reasons why data is missing:

- └ User forgot to fill in a field.
- └ Data was lost while transferring manually from a legacy database.
- └ There was a programming error.
- └ Users chose not to fill out a field tied to their beliefs about how the results would be used or interpreted.

Variable identification with Uni-variate, Bi-variate and Multi-variate analysis:

- └ import libraries for access and functional purpose and read the given dataset
- └ General Properties of Analyzing the given dataset
- └ Display the given dataset in the form of data frame
- └ show columns
- └ shape of the data frame
- └ To describe the data frame
- └ Checking data type and information about dataset
- └ Checking for duplicate data
- └ Checking Missing values of data frame
- └ Checking unique values of data frame
- └ Checking count values of data frame
- └ Rename and drop the given data frame
- └ To specify the type of values
- └ To create extra columns

Data Validation/ Cleaning/Preparing Process

Importing the library packages with loading given dataset. To analyzing the variable identification by data shape, data type and evaluating the missing values, duplicate values. A validation dataset is a sample of data held back from training your model that is used to give an estimate of model skill while tuning models and procedures that you can use to make the best use of validation and test datasets when evaluating your models. Data cleaning / preparing by rename the given dataset and drop the column etc. to analyze the uni-variate, bi-variate and multi-variate process. The steps and techniques for data cleaning will vary from dataset to dataset. The primary goal of data cleaning is to detect and remove errors and anomalies to increase the value of data in analytics and decision making.



```

In [1]: #import library packages
import pandas as p
import numpy as n

In [2]: #Load given dataset
data = p.read_csv("house.csv")

In [3]: import warnings
warnings.filterwarnings('ignore')

Before drop the given dataset

In [4]: data.head()
Out[4]:
   id      date  price  bedrooms  bathrooms  sqft_living  sqft_lot  floors  waterfront  view  ...  grade  sqft_above  sqft_baseament  yr_built
0  7129300520  20141013T000000  221900.0      3      1.00      1180      5650      1.0      0      0  ...    7      1180      0      1955
1  6414100192  20141209T000000  538000.0      3      2.25      2570      7242      2.0      0      0  ...    7      2170      400      1951
2  5631500400  20150225T000000  180000.0      2      1.00      770      10000     1.0      0      0  ...    6      770      0      1933
3  2487200875  20141209T000000  604000.0      4      3.00      1960      5000      1.0      0      0  ...    7      1050      910      1965
4  1054100610  20141013T000000  510000.0      2      2.00      1680      8080      1.0      0      0  ...    8      1680      0      1987
    
```

```

In [23]: print("Minimum bedrooms :", df.bedrooms.min())
print("Maximum bedrooms :", df.bedrooms.max())
Minimum bedrooms : 0
Maximum bedrooms : 33

In [24]: print("sorted:", sorted(df['bedrooms'].unique()))
sorted: [0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 33]

In [25]: p.Categorical(df['view']).describe()
Out[25]:
   counts  freqs
categories
0  19489  0.901726
1   332  0.015361
2   963  0.044557
3   510  0.023597
4   319  0.014760

In [26]: p.Categorical(df['bedrooms']).describe()
Out[26]:
   counts  freqs
categories
0     13  0.000601
    
```

Preparing the Dataset :

This dataset contains 15129 training 6213 testing records extracted, which were regression pattern:

- Predicting to the House price using regression pattern.

Exploratory Data Analysis of house Prediction

Multiple datasets from different sources would be combined to form a generalized dataset, and then different machine learning algorithms would be applied to extract patterns and to obtain results with MAE.



```

In [1]: #import library packages
import pandas as p
import matplotlib.pyplot as plt
import seaborn as s
import numpy as n

In [2]: data = p.read_csv("house.csv")

In [3]: import warnings
warnings.filterwarnings('ignore')

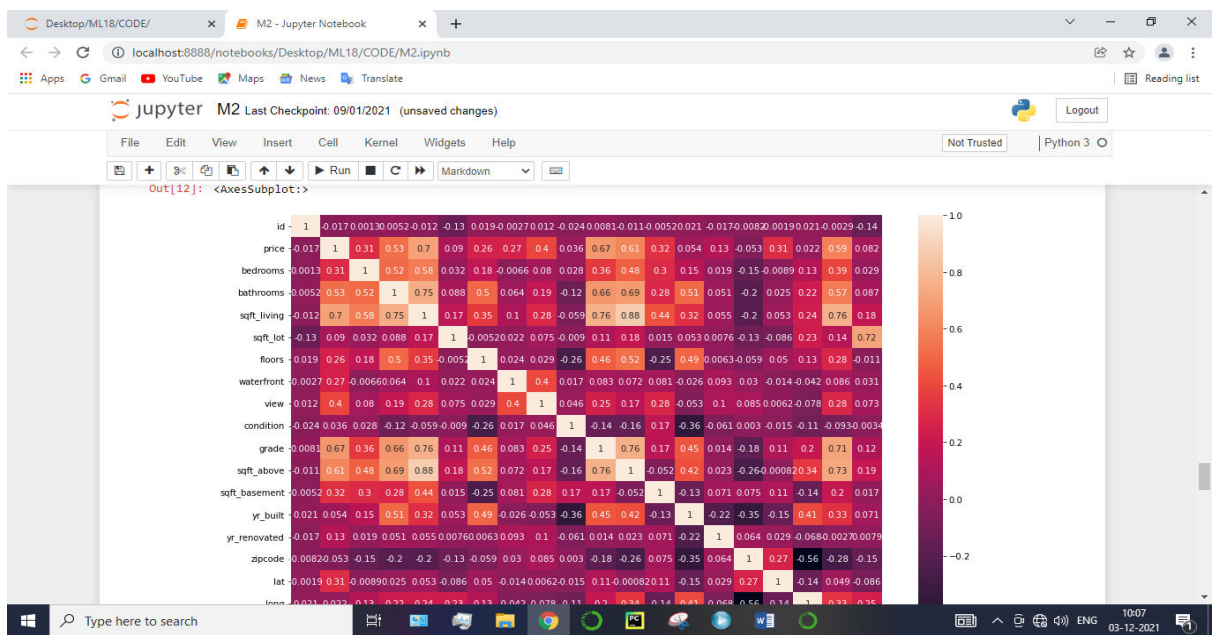
In [4]: df=data.dropna()

In [5]: df.columns

Out[5]: Index(['id', 'date', 'price', 'bedrooms', 'bathrooms', 'sqft_living',
             'sqft_lot', 'floors', 'waterfront', 'view', 'condition', 'grade',
             'sqft_above', 'sqft_baseament', 'yr_built', 'yr_renovated', 'zipcode',
             'lat', 'long', 'sqft_living15', 'sqft_lot15'],
            dtype='object')

In [6]: #Histogram Plot

```



Comparing Algorithm with prediction in the form of best accuracy result

It is important to compare the performance of multiple different machine learning algorithms consistently and it will discover to create a test harness to compare multiple different machine learning algorithms in Python with scikit-learn. It can use this test harness as a template on your own machine learning problems and add more and different algorithms to compare. Each model will have different performance characteristics. Using resampling methods like cross validation, you can get an estimate for how accurate each model may be on unseen data. It needs to be able to use these estimates to choose one or two best models from the suite of models that you have created. When have a new dataset, it is a good idea to visualize the data using different techniques in order to look at the data from different perspectives. The same idea applies to model selection. You should use a number of different ways of looking at the estimated accuracy of your machine learning algorithms in order to choose the one or two to finalize. A way to do this is to use different visualization methods to show the average accuracy, variance and other properties of the distribution of model accuracies.



Data Wrangling

In this section of the report will load in the data, check for cleanliness, and then trim and clean given dataset for analysis. Make sure that the document steps carefully and justify for cleaning decisions.

Data collection

The data set collected for predicting given data is split into Training set and Test set. Generally, 7:3 ratios are applied to split the Training set and Test set. The Data Model which was created using machine learning algorithms are applied on the Training set and based on the test result MAE, Test set process is done.

Building the regression model

The predicting the house price problem, decision tree regression model is effective because of the following reasons: It provides better results in regression problem.

- It is strong in preprocessing outliers, irrelevant variables, and a mix of continuous, regression and discrete variables.
- It produces out of bag estimate error which has proven to be unbiased in many tests and it is relatively easy to tune with.

The key to a fair comparison of machine learning algorithms is ensuring that each algorithm is evaluated on the same way on the same data and it can achieve this by following the process shown in the diagram below.

ALGORITHM AND TECHNIQUE

Algorithm Explanation

In machine learning and statistics, classification is a supervised learning approach in which the computer program learns from the data input given to it and then uses this learning to classify new observation. This data set may simply be bi-class (like identifying whether the person is male or female or that the mail is spam or non-spam) or it may be multi-class too. Some examples of classification problems are: speech recognition, handwriting recognition, bio metric identification, document classification etc. In Supervised Learning, algorithm learn from labeled data. After understanding the data, the algorithm determines which label should be given to new data based on pattern and associating the patterns to the unlabeled new data.

Random Forest Regression:

Random Forest Regression is a supervised learning algorithm that uses **ensemble learning** method for regression.

Ensemble learning method is a technique that combines predictions from multiple machine learning algorithms to make a more accurate prediction than a single model.

we have a basic understanding of how the Random Forest Regression model works, we can assess its performance on a real-world dataset. Similar to my previous posts, I will be using data on House Sales in King County, USA.

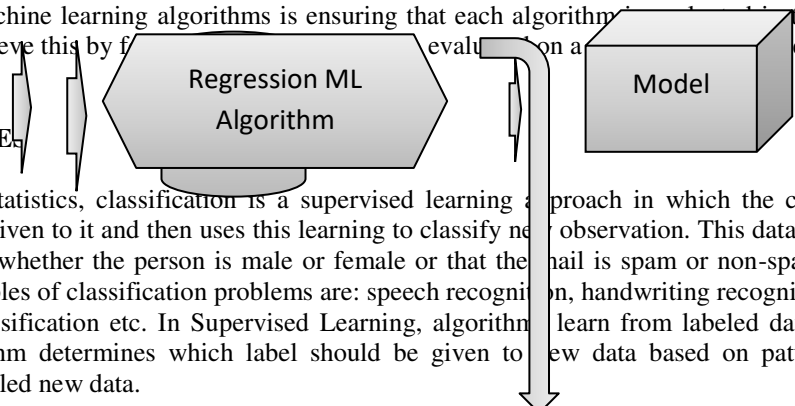
After importing the libraries, importing the dataset, addressing null values, and dropping any necessary columns, we are ready to create our Random Forest Regression model

From the sklearn package containing ensemble learning, we import the class **Random Forest Regressor**, create an instance of it, and assign it to a variable.

The parameter *n_estimators* creates *n* number of trees in your random forest, where *n* is the number you pass in. We passed in 10.

The *.fit()* function allows us to train the model, adjusting weights according to the data values in order to achieve better accuracy. After training, our model is ready to make predictions, which is called by the *.predict()* method.

A random forest is a meta estimator that fits a number of classifying decision trees on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting. The sub-sample size is controlled with the *max_samples* parameter if *Bootstrap = True* (default), otherwise the whole dataset is used to build each tree.





Module 4 : Performance measurements of RandomForestRegressor and DecisionTreeRegressor algorithms

```
In [1]: #import library packages
import pandas as p
import matplotlib.pyplot as plt
import seaborn as s
import numpy as n

In [2]: #Load given dataset
data = p.read_csv("house.csv")

In [3]: import warnings
warnings.filterwarnings('ignore')

In [4]: df=data.dropna()

In [5]: df.head()

Out[5]:
```

	id	date	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	waterfront	view	...	grade	sqft_above	sqft_basement	yr_built
0	7129300520	20141013T000000	221900.0	3	1.00	1180	5650	1.0	0	0	...	7	1180	0	1955
1	6414100192	20141209T000000	538000.0	3	2.25	2570	7242	2.0	0	0	...	7	2170	400	1951

RandomForestRegressor

```
In [14]: from sklearn.ensemble import RandomForestRegressor
RF= RandomForestRegressor()
RF.fit(X_train,y_train)

predictR = RF.predict(X_test)

MAE=(mean_absolute_error(y_test,predictR))
print("MEAN ABSOLUTE ERROR VALUE IS :",MAE)
print(" ")

MSE=(mean_squared_error(y_test,predictR))
print("MEAN SQUARED ERROR VALUE IS :",MSE)
print(" ")

MedianAE=(median_absolute_error(y_test,predictR))
print("MEDIAN ABSOLUTE ERROR VALUE IS :",MedianAE)
print(" ")

EVS=(explained_variance_score(y_test,predictR)*100)
print("ACCURACY RESULT OF RANDOM FOREST REGRESSOR IS :",EVS)
print(" ")

R2=(r2_score(y_test,predictR))
print("R2 SCORE VALUE IS :",R2)
print(" ")

MEAN ABSOLUTE ERROR VALUE IS : 86853.44266894333
```

Decision Tree Regression

Decision Tree - Regression. Decision tree builds regression or classification models in the form of a tree structure. It breaks down a dataset into smaller and smaller subsets while at the same time an associated decision tree is incrementally developed. The final result is a tree with decision nodes and leaf nodes.

Decision tree builds regression or classification models in the form of a tree structure. It breaks down a dataset into smaller and smaller subsets while at the same time an associated decision tree is incrementally developed.



The final result is a tree with decision nodes and leaf nodes. A decision node (e.g., Outlook) has two or more branches (e.g., Sunny, Overcast and Rainy), each representing values for the attribute tested. Leaf node (e.g., Hours Played) represents a decision on the numerical target. The topmost decision node in a tree which corresponds to the best predictor called root node. Decision trees can handle both categorical and numerical data.

Decision Tree is one of the most commonly used, practical approaches for supervised learning. It can be used to solve both Regression and Classification tasks with the latter being put more into practical application. It is a tree-structured classifier with three types of nodes.

Decision trees where the target variable can take continuous values (typically real numbers) are called regression trees. Decision trees are among the most popular machine learning algorithms given their intelligibility and simplicity.

```

DecisionTreeRegressor

In [15]: from sklearn.tree import DecisionTreeRegressor
DT= DecisionTreeRegressor()
DT.fit(X_train,y_train)

predictD = DT.predict(X_test)

MAE= (mean_absolute_error(y_test,predictD))
print('MEAN ABSOLUTE ERROR VALUE IS :',MAE)
print(" ")

MSE=(mean_squared_error(y_test,predictD))
print('MEAN SQUARED ERROR VALUE IS :',MSE)
print(" ")

MedianAE=(median_absolute_error(y_test,predictD))
print('MEDIAN ABSOLUTE ERROR VALUE IS :',MedianAE)
print(" ")

EVS=(explained_variance_score(y_test,predictD)*100)
print('ACCURACY RESULT OF DECISION TREE REGRESSOR IS :',EVS)
print(" ")

R2=(r2_score(y_test,predictD))
print('R2_SCORE VALUE IS :',R2)
print(" ")

MEAN ABSOLUTE ERROR VALUE IS : 120583.03789130167
    
```

```

print(" ")

EVS=(explained_variance_score(y_test,predictD)*100)
print('ACCURACY RESULT OF DECISION TREE REGRESSOR IS :',EVS)
print(" ")

R2=(r2_score(y_test,predictD))
print('R2_SCORE VALUE IS :',R2)
print(" ")

MEAN ABSOLUTE ERROR VALUE IS : 120583.03789130167

MEAN SQUARED ERROR VALUE IS : 46300104090.13283

MEDIAN ABSOLUTE ERROR VALUE IS : 65000.0

ACCURACY RESULT OF DECISION TREE REGRESSOR IS : 66.06391505761528

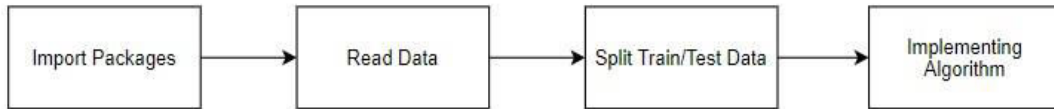
R2_SCORE VALUE IS : 0.6605934536942953

In [16]: import joblib
joblib.dump(RF,'rf1.pk1')

Out[16]: ['rf1.pk1']

In [ ]:
    
```


MODULE DIAGRAM:



GIVEN INPUT EXPECTED OUTPUT

Input : Data

Output : Getting Accuracy

II. IV. EXPERIMENTAL RESULTS





III. V. CONCLUSION

IV.

The analytical process started from data cleaning and processing, missing value, exploratory analysis and finally model building and evaluation. The best accuracy on public test set is higher accuracy score is will be find out. This application can help to find the House Price. We have overseen out how to set up a model that gives clients for an original best methodology with look at future dwelling esteem expectations. Straight previous suggest works bring been used inside our model, something appreciate that that future worth expectations will have an inclination towards even more reasonable qualities. We composed an approach with use in basically the same manner as significantly data as time licenses for our expectation framework, by embracing those thoughts from guaranteeing inclination supporting. These consolidate redesigns we didn't choose as a result of compelled length of the time. A genuine concern for the expectation structure may be the stacking time frame. Additionally, our informational collection takes more than one day ought to get ready. As gone against playing out the calculations consecutively, we may use different processors and equal the calculations in question, which may potentially diminish the planning time Furthermore expectation period. Incorporate even more functionalities under the model, we can give decisions for customer with select a region on the other hand district should deliver rather than entering in the rundown.

REFERENCES

1. John M. Clapp and Carmelo Giaccotto .” Evaluating House Price Forecasts” - 2002
2. Lynn Wu, Erik Brynjolfsson .” How Google Searches Foreshadow Housing Prices and Sales”- 2009
3. Vasilios Plakandaras, Rangan Gupta, Periklis Gogas and Theophilos Papadimitriou , “Forecasting the U.S. Real House Price Index” - 2013
4. Koen de Koning, Tatiana Filatova, Okmyung Bin,” Improved Methods for Predicting Property Prices in Hazard Prone Dynamic Markets” - 2016
5. Gaikwad Purva Chandrakant, Ganjave Pratiksha Namdev, Gorade Pooja Subhash, S. S. Gore, “Implementation of House Price Prediction Model Using Image Processing and Machine Learning” -2019



SJIF: 8.423



INTERNATIONAL
STANDARD
SERIAL
NUMBER
INTERNATIONAL CENTRE



INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY

 **9940 572 462**  **6381 907 438**  **ijirset@gmail.com**



www.ijirset.com

Scan to save the contact details