

HOUSE PRICE
PREDICTION PROJECT
REPORT-

ACKNOWLEDGMENT

First and foremost, I would like to thank, Flip Robo Technologies to provide me a chance to work on this project. It was a great experience to work on this project under your guidance.

I would like to present my gratitude to the following websites:

- Zendesk
- Kaggle
- Datatrained Notes
- Sklearn.org
- Crazyegg

These websites were of great help and due to this, I was able to complete my project effectively and efficiently.

INTRODUCTION

- Business Problem Framing

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

- Conceptual Background of the Domain Problem Basic EDA concepts and regression algorithms must be known to work on this project. One should know what is Housing Price and how it is going to affect the real estate business. Why predicting the house prices is important and how can it be going to help the company?

- Review of Literature Houses are one of the necessary needs of each and every person around the globe and therefore housing and real estate market is one of the markets which is one of the major contributors in the world's economy. It is a very large market and there are various companies working in the domain. Data science comes as a very important tool to solve problems in the domain to help the companies increase their overall revenue, profits, improving their marketing strategies and focusing on changing trends in house sales and purchases. Predictive modelling, Market mix modelling, recommendation systems are some of the machine learning techniques used for achieving the business goals for housing companies. Our problem is related to one such housing company.

A US-based housing company named Surprise Housing has decided to enter the Australian market. The company uses data analytics to purchase houses at a price below their actual values and flip them at a higher price. For the same purpose, the company has collected a data set from the sale of houses in Australia. The data is provided in the CSV file below.

The company is looking at prospective properties to buy houses to enter the market. You are required to build a model using Machine Learning in order to predict the actual value of the prospective properties and decide whether to invest in them or not. For this company wants to know:

- Which variables are important to predict the price of variable?
- How do these variables describe the price of the house?

Analytical Problem Framing

- Data Sources and their formats

-The dataset is provided the internship organization in an excel format which contains the data in both in code sheet and categorical data. It contains 81 columns and 1168 rows. There are so many factors which can be used for the prediction of sale price of a house. It contains the factors on which the sale price of a house can depend. Dataset contain both numerical as well as categorical data.

- Libraries Used

I am using different libraries to explore the dataset.

1. Pandas – It is used to load and store the dataset. We can discuss the dataset with the pandas different attributes like .info, .columns, .shape.
2. Seaborn – It is used to plot the different types of plots like catplot, lineplot, countplot and more to have a better visualization of the dataset.
3. Matplotlib.pyplot – It helps to give a proper description to the plotted graph by seaborn and make our graph more informative.
4. Numpy – It is the library to perform the numerical analysis to the dataset.

- ***Load the Dataset***

- ***Checking the Attributes***

- First & last five rows of both the dataset
- Shape of the datasets
- Columns present in the datasets
- Brief info about the datasets
- Null values present in both the dataset

Now we have checked the attributes for the dataset and get a rough idea about the dataset like the no of rows & columns, datatype & null values in the dataset.

Dealing with the Null Values

In both the dataset null values are present, so we have to handled them for better model learning. As we have categorical & numerical data so we have to handled them accordingly. We also drop those rows who are having more than 50% null values.

EXPLORATORY DATA ANALYSIS

Which street house has higher price?

- House in Pave street have higher sale price

What type of Land Contour has higher sale price?

- HLS type houses has higher sale price

What Lot configuration is in higher demand?

- CulDSac followed by FR3 lot configuration are in higher demand.

Whose neighborhood increased the sale price?

- The one who has NoRidge & NridgHt in their neighbourhood has the high sale price.

Which house style has high sale price?

- The 2.5Fin has the highest sale price followed by 2Story and 1.5Unf has the lowest sale price.

How lot area affects the price?

- Most of the houses have low lot area, very little on the higher side & the price is very high for some of the houses

How land slope affects the price?

- Land slope doesn't affect the price much more; it is same almost for every type.

What type of building has high sale price?

- TwnhsE & IFarm type buildings are on higher side.

OverallQual Vs Price

- A higher overall quality grade means a higher sale price.

Overall Condition Vs Sale Price

- Whose overall condition is around 5 touches the higher side of price

Year built Vs Sale Price

- Newer houses sale price is high as compared to old houses but there are also some old houses whose sale price is high

YearRemodAdd Vs Sale Price - We almost have a distributed data in this but yes newer one has higher price than older one.

Which roof style & material increases the price of a house?

- A house with Gable roof style and made of WdShngl shown up a with higher sale price

Which exterior contribute more towards sale price?

- ImStucc & Stone followed by CementBd exterior sale price is high compare to others

Basement condition Vs Price

- Whose condition is GD or TA then obviously getting the higher sale price

Total Rms above ground Vs Sale price

- With Grade 9, 10 & 11 your sale price is going to be good

Year sold Vs Price

- Whatever be the year sale price doesn't affect too much

Is saletype affect the sale price?

- Yes, if your sale type is Con or New then definitely you are going to get a good price

Label Encoding & Correlation

As we have some categorical data we have to encoded those columns for machine learning model. We will use Label Encoder from `sklearn.preprocessing`. We will describe the statistical summary of the dataset and find the correlation of each column.

Removing the Outliers

We have some outliers present in the dataset, so let's handle them also. As the outliers in the dataset will affect our ML model. We need to remove all the outliers present in the dataset. There is something called zscore which indicates how many standard deviations away an element is from the mean. We consider the points as outliers whose zscore is above 3 or less than -3.

So we need to remove all such points from our dataset. Using the threshold, we have removed all the points where the zscore is greater than 3. Now the total number of rows after removing the outliers are 721.

MODEL BUILDING

We will import important libraries for the building the ML model and defining the different models for our easiness.

Regression Algorithms

We have use five different regression algorithms to find the best model for our problem:

1. Linear Regression- from `sklearn.linear_model` import `LinearRegression`
2. Decision Tree Regressor- from `sklearn.tree` import `DecisionTreeRegressor`
3. Support Vector Regressor- from `sklearn.svm` import `SVR`
4. Kneighbor Regressor- from `sklearn.neighbors` import `KNeighborsRegressor`
5. Random Forest Regressor- from `sklearn.ensemble` import `RandomForestRegressor`

MODEL ACCURACY

Linear Regression 0.878165525517784

Decision Tree Regressor 0.7819521438284093

Support Vector Regressor -0.044480459746993

Kneighbor Regressor 0.6839964008801616

Random Forest Regressor 0.8815066126350098

Hence, we are getting the best accuracy score through the Random Forest Classifier Model. We will go ahead with this to find the cross val score and hypermeter tuning.

Cross Val Score & Hypermeter Tuning

Cross-validation provides information about how well a classifier generalizes, specifically the range of expected errors of the classifier. Cross Val Score tells how the model is generalized at a particular cross validation.

At CV=3 we get the best results i.e. the Random Forest Classifier more generalized at cv=3, so we calculate the hyper parameters at this value. We will find which parameters of random forest classifier are the best for our model. We will do this using Grid Search CV method & also calculate the accuracy score at those best parameters.

Saving the Model

Saving the best model – Random Forest Classifier in this case for future predictions.

Hence up to some good extensions our model predicted so well.

Now, what our model predict for test dataset?

With the best model that we have saved earlier, let's predict the sale price of the houses.

CONCLUSION

Conclusion of the Study: - The results of this study suggest following outputs which might be useful for the company to enter into the Australian Market:

- There are lot of things that is going to decide the sale price of a house. As we see above in our visualizations, a lot of things affect the price like neighbourhood, Quality condition, basement, living area, roof type, building type and many more. One needs to analyse every aspect to have good hands on the prediction of the price.
- With the machine learning it become easier to predict the price but yes it is not 100% accurate, it provides an idea and accordingly we can analyse the market and prepare the strategies to grab the opportunities.
- Learning Outcomes of the Study in respect of Data Science
- I got to know the different factors required for the price prediction of a house.
- It was fun to deal with this project and learn how we can use our saved model to predict the price for given dataset.
- It was difficult to handle so much columns simultaneously but yes every difficulty learns the new things to us.

THANK YOU