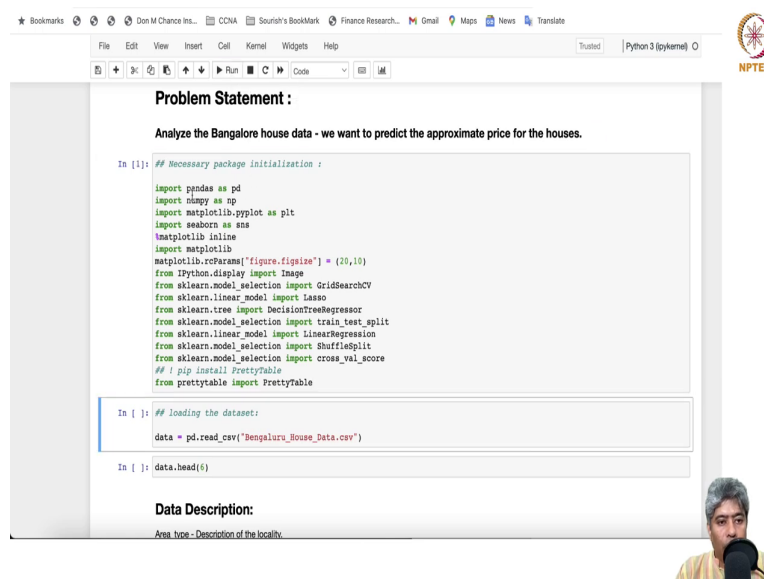**Predictive Analytics - Regression and Classification**
**Prof. Sourish Das**
**Department of Mathematics**
**Chennai Mathematical Institute**

**Lecture - 54**
**Hands on with Python: Analysis of Bangalore House Price Data**

Hello all. In this video, we are going to do some hands-on. We are going to analyze Bangalore House Price Data using Python. I have shared the Bangalore house price data in the NPTEL platform. You can also get this Bangalore house price data in Kaggle or in internet somewhere. In some GitHub repository also it is available.

(Refer Slide Time: 00:54)



So, in this notebook I prepared the problem statement is we will use this Bangalore house price data and we will try to approximate or estimate a possible price of any house that are kind of part of the Bangalore you know in the vicinity of approximately we will try to

estimate a possible price of a possible house. So, first we are going to load this bunch of Python packages.

(Refer Slide Time: 01:28)



And, then we are going to call the Bangalore house price data. So, the first I will print this dataset. It has 9 columns; area type, availability, location, size, society, total square feet, bath, a balcony and price. I have also provided the description of each column, area type refers to description of the locality.

(Refer Slide Time: 01:04)



In fact, I can try this also. So, yeah, you can just put.

(Refer Slide Time: 02:07)



So, I will do that beautification later and location stands for where it is located in the Bangalore, size stands for BHK or bedrooms like whether it is a 3 BHK or 2 BHK back pay. Society stands for which society looks like there are some in a values. Bath stands for number of bathrooms, the apartment of the house has, balcony is the number of balcony and price stands for the value of the price in rupees lakhs. Data dot shape will give you the shape of the data. It is it has 13,320 rows and 9 columns.

(Refer Slide Time: 03:07)



Data columns, data dot columns will give you the all the columns that are there and data in dot info will give you basic information about each column. Area type, availability location, size, society, square feet, these are objects. Type data, bath, balcony price, these are float. Now, total square feet cannot be object. If you look into the data, total square feets are numeric.

(Refer Slide Time: 03:46)



There must be some character that was required that came in the data and that resulted the made it object. We will see how to solve that problem. If we say data dot describes, it will give you basic summary of the float variables which is bath, balcony and price. So, like there are 13,247 instances where the bathroom, bath variable is available average 2.69, minimum 1 bathroom.

There are most of the median is 2, most of the values are here. Most of the apartment or house has 2 bathroom and maximum 40. There could be very big house or mansions which has a very. And, the price is average price is 112 lakhs and median is 72 lakhs whereas, maxima is 3600 lakhs so, near about 36 crore.

(Refer Slide Time: 05:00)



So, there data dot null, if there are many null values are there and if you do some, then if you see that society has maximum number of null values and balcony also has quite a few null values. So, probably it is better to drop society. And we did some basic searching and we found that in this particular row, total square feet was recorded as 34.46 cube times meter. So, it is a weird way of recorded.

(Refer Slide Time: 05:42)



So, we may have to drop this recording and performing group by operation. So, if this will be like area types, these are the count of the area types. There are four kind of area types, Built-up Area, Carpet Area, Plot Area and Super built Area. So, now, Super built Area 8790 instances belongs to Super built Area. Majority of the instances belong to Super built Area. 2000 cases and 2004 and cases belongs to Plot Area or Built-up Area and there are only 87 cases where it is given as a Carpet Area.

(Refer Slide Time: 06:34)



So, if you run this. So, there are four cases of area types.

(Refer Slide Time: 06:40)



Let us do some visualization. So, there is a bar plot for the area type.

(Refer Slide Time: 06:48)



 If you do, we can run this. This is the bathroom has a bar plot.

(Refer Slide Time: 06:56)



This is the bar plot for the bathroom and you can see there are about to almost 7000 instances had 2 bathroom. So, maximum cases there are 2 bathroom.

(Refer Slide Time: 07:12)

(Refer Slide Time: 07:15)



And, again this is a pie chart for area type.

(Refer Slide Time: 07:21)



And, balcony there are 0, 1, 2, 3. Majority are either have 1 balcony or 2 balcony. There are few cases there are 3 balcony and quite a few cases there are no balcony. Now, if you count, you can count how many null values are there for each predictor. The society as we found, society as the maximum number of null values, then balcony followed by bath and size.

Price, total square feet, availability, area type does not have any balcony, any null values. We are dropping here area types, society, balcony and availability. Because, if you run the area types, if you actually run this, copy this and.

(Refer Slide Time: 08:50)



So, if you do availability so, you will see that there are 81 different type of availability. Most of the availability goes to ready to move and they stop them at just given a number date without any particular mention. So, it isit is not going to create any help in the model final analysis. And, area type is also not going to available and society has too many null values.

So, that is why we are going to drop these cases. And after we drop, we have now 13,320 cases still there and after dropping these columns. Now, we have 5 columns and now if we drop the NAs, we have still 13,246 now cases. Now, there is no null value in the data set. So, that is a good start to have.

(Refer Slide Time: 09:58)



So, now, it is a full data, there is no missing data. Now, let us do, we have to may have to do some engineering or data cleaning, not only data engineering, we have to do some and data cleaning ok. Now, first if you run data size, there are you will see that there are 2 BHK, 3 BHK, 4 BHK cases and if you run this, ok. So, there are 2 BHK and then ok, maybe there are 6 alright.

(Refer Slide Time: 10:44)



So, these are the this is the data set that we have. Ok. And, if we drop plot the BHK cases.

(Refer Slide Time: 10:55)



So, these are the different type of 1, 2, 3, 4 different kinds of BHKs and they are plot. So, most of the instances are either 2 BHK house or 3 BHK or 4 BHK houses. There are some 1 BHK and 5 BHK and then other cases are very rare. So, the unique values that we have here.

(Refer Slide Time: 11:25)



And, the BHK versus price box plot, these are like side by side box plot, these are like single value cases.

(Refer Slide Time: 11:33)



So, they just put up a dash kind of thing. Now, if we just look into number of unique value in the total square fit, it has 2067 and there are just too many unique values. So, th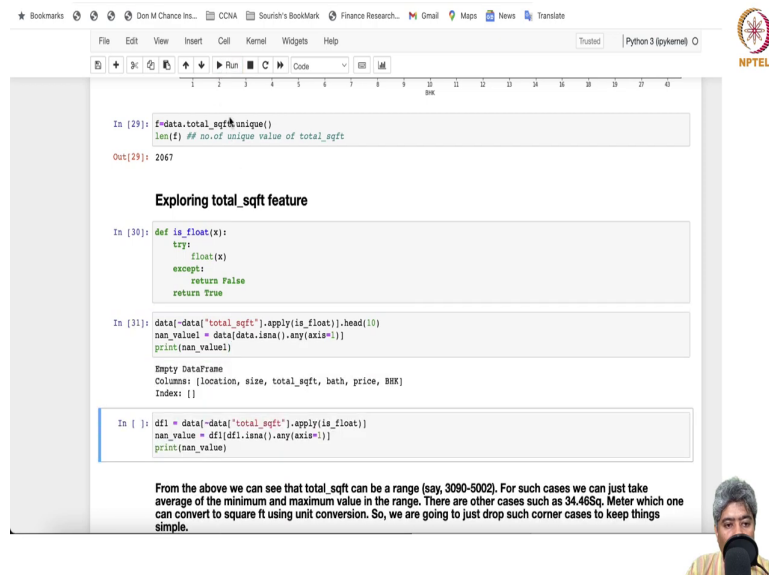at we have to, there are some character values. So, we have to first check which are the float and which are the not float. So, we can create definition.

And, then where we have float and a. So, empty data frame and then ok. Let us run this. Hopefully, this is fine.

(Refer Slide Time: 12:27)

(Refer Slide Time: 12:30)



So, there are quite a few.

(Refer Slide Time: 12:38)



These are the cases where the total square feet you have NAs.

(Refer Slide Time: 12:40)



And so, these are the, we just created price per total square feet. So, price per square feet.

(Refer Slide Time: 12:59)



How much you are paying for each square feet?

(Refer Slide Time: 13:03)



And now, if you just plot, you can see this on the x axis we are plotting total square per feet versus price.

And we are here, we are going to, we are checking the categorical variable versus numerical variables, how many categorical, how many numerical, those are which are the. So, these are the numerical variables and these two are categorical, location and size. So, how many new unique locations we have? So, there are 1300 unique locations that we have.

(Refer Slide Time: 13:48)



And in the 1300 unique location, we count that and these is the situation, there are quite a few location where you have only single instances. So, if we do that, then there will be little bit of a problem because you see, there are 1300, about 1300 locations, unique location and if we create a dummy variable, if we do a one-hot one encoding and if we do create a dummy variable, for each location, it will create a dummy variable.

And there are instances we have one one-hot encoding and that will make the solution unreliable. In fact, probably you will not have a solution, you can have a solution, but you will not be able to calculate the standard error. So, you need at least 20 locations, a 20 instances for each for a locations, a solution to estimate. So, what we are doing here, we are saying that, ok, let us try to figure out which are the locations you have unique are less than 20 instances ok.

(Refer Slide Time: 15:10)



These are the locations where you have less than 20 instances, ok. And, let us drop those and we are discarding those locations. So, there are now 144 unique locations. After dropping 20 location, where we have now have 44 unique locations. And if you just say data dot shape, still you have 13,240, 7, 246 instances you have. So, this is a good thing.

(Refer Slide Time: 16:01)



Ploting the Scatter Chart for log-Total Square Feet Area Vs log-price

And now we are going to do some plotting, scattered chart for 2BHK and 3BHK properties. So, say for example, here in the Rajaji Nagar area, you can give a area name and for that area, you will get a 2BHK and 3BHK's price or say Whitefield, if you take Whitefield.

And you give another place, just a minute. So, let me just copy this here and instead of. So, now, let me just run this. So, in the Whitefield, these are the 2BHK prices and these are the 3BHK. The x axis we plot the total square feet and on the y axis, we have right plotting the price for each of each cases.

(Refer Slide Time: 17:28)



Now, we are plotting the chart for a log total, log scale in the total square feet versus log price.

(Refer Slide Time: 17:46)



So, if you plot that, they are somewhat in the much more linear in terms.

(Refer Slide Time: 17:51)



And if you just do the, do it for the Whitefield also. So, this is already for the Whitefield. Let me do it for Rajaji Nagar than. So, this was for the Whitefield. Let me do it for the Rajaji Nagar. Alright. Ok. So, this is less number of instances that you have.

(Refer Slide Time: 18:46)



So, price per square feet, if you see the price per square feet, ok this was, this is the minimum actually. Let me just first plot. This is the really high price per square feet. And, if you create the minimum of it, then the price per square feet is 267. So, the range is very big. That is bit of a weird, I find.

(Refer Slide Time: 19:27)

(Refer Slide Time: 19:30)

(Refer Slide Time: 19:36)



And, the histogram that you are plotting here for price per square feet and then bathroom for histogram. So, there are number of bedrooms. These are the number of bedrooms and, sorry, number of bathrooms. So, these are the cases where you have number of bedroom is 9 and number of bathroom and number of bedroom is number of bathroom is more than the number of bedroom.

So, these are like really, really big houses with, you know, lot of. So, this could be mansion or something. These are like really big places. So, a really big outlier with all these things so, we can drop this probably.

(Refer Slide Time: 21:04)



These are not really regular cases for sure.

(Refer Slide Time: 21:08)



So, now we are going to do some one-hot encoding for the location. This is very important. So, this gives us one-hot encoding for all the locations.

(Refer Slide Time: 21:31)



5 rows for 144 columns.

For each columns, we get the first five rows or all the 144 columns. Ok. And at the end, we get the others. Ok, alright. And, then concatenating with the data frames together, we have to add them, concatenate them. So, we just concatenate with the data frames.

(Refer Slide Time: 21:54)



And, then we have to drop the location because now we have for each location, we have done the one-hot encoding, correct? And, we can drop the NAs. Now, we take the, from the, we can drop the price and rest of the, we will become the X matrix. We can drop the size and drop the price per square feet.

(Refer Slide Time: 22:31)



So, this will be our final X. So, total square feet, bathroom, BHK and the locations. This will be the our thing. So, X shape will be 13200 and 146 columns. So, we have so many columns and yeah. And, there is no NA. Ok. So, now we are going to split the data into train and test.

(Refer Slide Time: 23:09)



So, we are going to call train test split. Now, if let me just go on the top. Ok. In the, from the sklearn package, scikit-learn package, model selection module, we imported the train and test split here. Ok. So, let us run that. And from there, we are importing the linear model, we linear regression from the linear model of sklearn, we call it regression and then we are called fitting the regression dot fit. Ok.

(Refer Slide Time: 24:15)



Now so, the variance score is 0.45.

(Refer Slide Time: 24:25)



And, then if we just plot the residual errors for train data versus test data. So, they are somewhat overlap each other. So, which is a good news, which is a good news.

And then if we just plot this, the model, if we just fit the model, OLS model. So, adjusted R square is 61. So, 61 percent gets, can be explained by the data.

(Refer Slide Time: 25:05)



So, the total square feet has a very positive coefficient as the and you know looks, this p value are very small for bathroom and BHK. And, then there are locations for which some p values sometimes are small, small. For example, there is a base line and against the base line. This Akshaya Nagar has a coefficient which is negative 52 and the p value is small. So, that means, it is small enough and the price that you pay is statistically significantly smaller than the base line.

(Refer Slide Time: 26:07)



So, we can look for places where it is bit high or bit low. So, for example, if you see this place has a positive coefficient, but the p value is really not large and the coefficient confidence interval also contains the 0. So that means, is its sort of a base line you know expected according to the expected price, you do not really expect much.

(Refer Slide Time: 26:41)



So, if you interestingly, Frazer Town is a place which has coefficient positive. That means for Frazer Town, you pay a premium and looks like p value is small and the coefficient confidence interval for the coefficient is also not including 0. So, I do not live in Bangalore. So, the people of Bangalore can tell me whether it makes sense that Frazer Town has a house price which is bit expensives looks like it is expensive place to be.

(Refer Slide Time: 27:26)



Let us see how it is. For example, Jigani is a place which has a coefficient value negative 68. That means, house price is lower than the expected average prices and then coefficient is also negative, quite low..

(Refer Slide Time: 28:01)



KR Puram, Kadugodi, these are all places where you have prices to be little on the lower side.

Now, here is one place, Koramangala looks like it is a bit expensive place.

(Refer Slide Time: 28:43)



Because, the coefficient is positive 47.112, p value is quite small and the confidence interval for the coefficient is completely right side of the 0. Does not include 0. So, the people from Bangalore tell me if Koramangala is a place where housing price is more than expected. At least the model is saying it is bit expensive place ok. There is one place looks really expensive Rajaji Nagar.

The coefficient is very high, 155.1818. And, then p value is really really small and interesting and the coefficient is completely off. One confidence in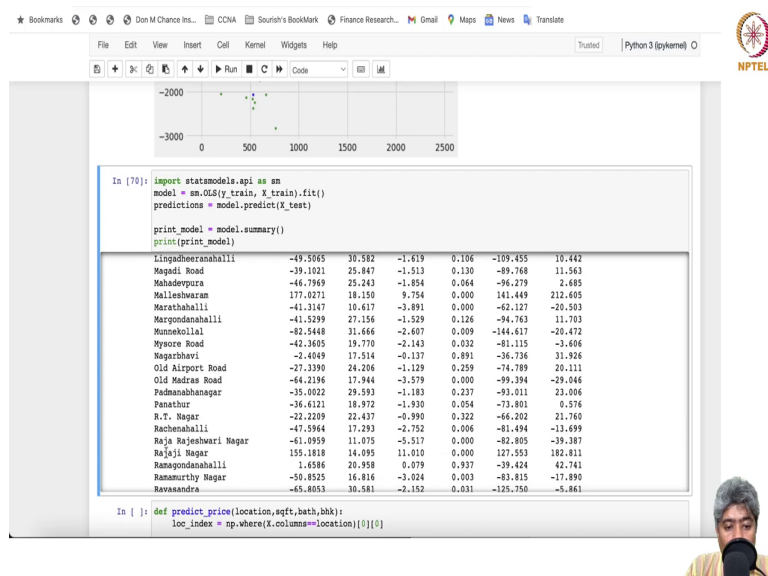terval of the coefficient is somewhere between 127 and 182. Guys, this place looks like really expensive. Tell me I do not live in Bangalore, I live in I am from Chennai. I teach at Chennai Math Institute.

(Refer Slide Time: 29:26)



So, people from Bangalore can tell me whether this place is really expensive place in terms of house price. Ok. So, ok.

(Refer Slide Time: 29:47)



So, I think you guys can figure out whether you live in a not so expensive place or not so or really expensive place.

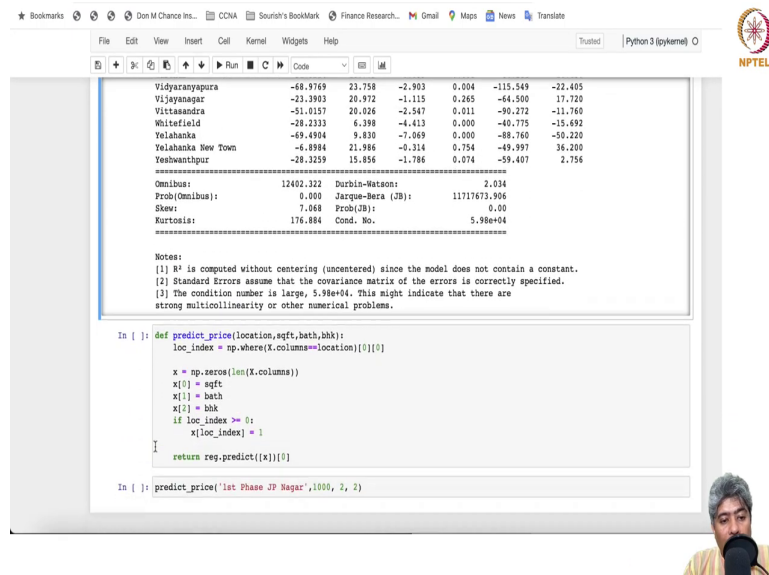Alright. Now, suppose we want to place predict price for a place where I will go give location, square feet, bath and BHK and what could be the price. So, this piece of code, this function will write give you the that price.

(Refer Slide Time: 30:30)



And if I just say predict, it will give you the 107 lakhs. So, if it is first phase of JP Nagar, 1000 square feet apartment with 2 BHK, 2 bathroom and 2 BHK. Then, it will cost you 107 lakh almost 1 point which is 1 crore 7 lakh rupees expected price. So, we can play with similarly you can take this say if you go to Rajaji Nagar, I think this was the place which was we found bit expensive.

Let us take Rajaji Nagar ok. For the same apartment for 1000 square feet and 2 bathroom, 2 BHK, it will cost you 2.4 crore or 241 lakh. It is very expensive place to be whereas, if you in some other place maybe Raja Rajeshwari Nagar. Let us see how it displays. You need only 33 lakh rupees to buy apartment in Raja Rajeshwari Nagar.

So, we can see that you know based on the location, the price same apartment can cost very different. Guys, I have no clue whether these values what this model is predicting, it does

make sense or not. So, tell me if this makes sense you guys and let me know if it makes sense, then we will know that ok these models these predictive models do work and that will be really fun. So, I will stop here.

Thank you very much. See you in the next video.