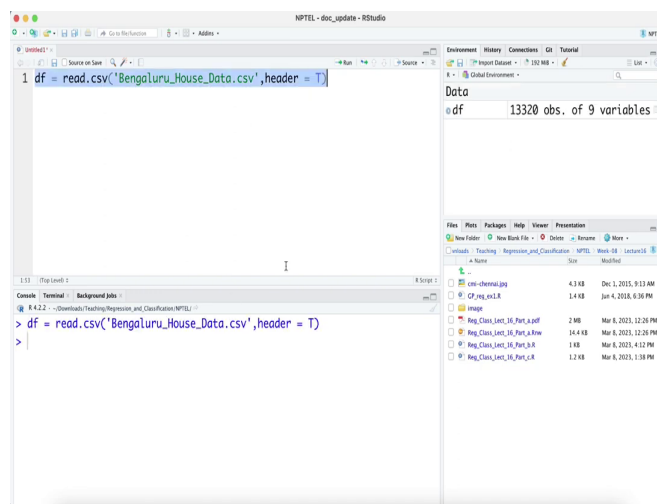**Predictive Analytics - Regression and Classification**
**Prof. Sourish Das**
**Department of Mathematics**
**Chennai Mathematical Institute**

**Lecture - 55**
**Hands on with R : Prediction of Bangalore House Price**

Hello all, in this video we are going to do the Bangalore House Price data analysis using R. In the previous video we did the same analysis with Python. Now, if we do using R, we will get a sense of both language, how both language works and some of the subtle difference in the language.

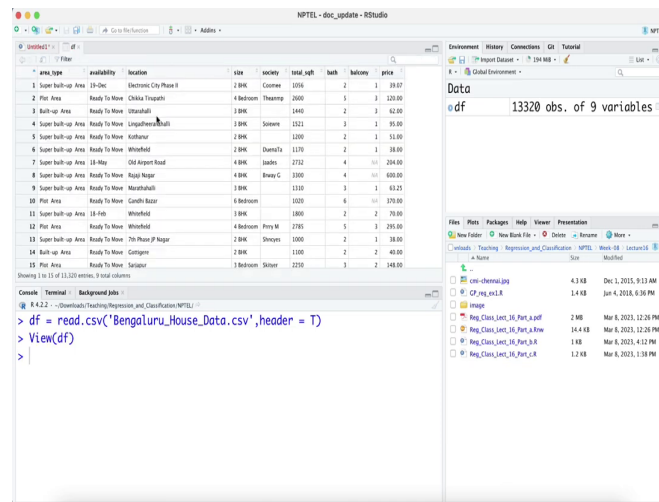(Refer Slide Time: 00:46)



So, first I will call data set read dot csv.

(Refer Slide Time: 01:02)



So, let me just go and take the copy the name of the data set and let me just call header equal to true. So, now you can see the data has come here.
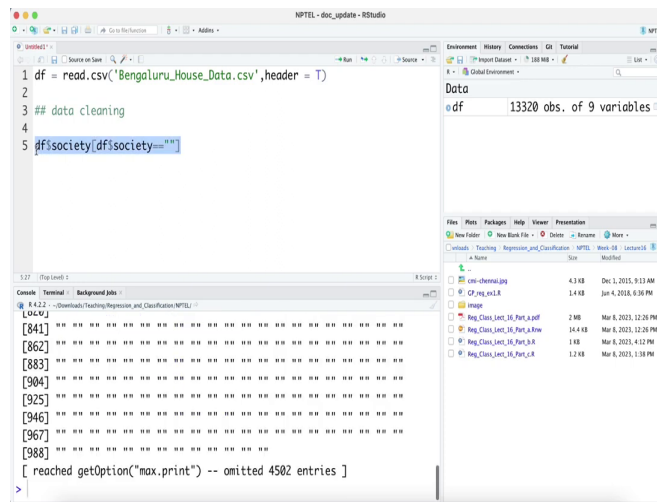
(Refer Slide Time: 01:20)



So, it has 13,320 observation and 9 variables. So, 9 columns you can see all of them here. Now, first thing we will do since we have done quite a view of visualization in Python we will not spend too much time on the visualization. We can do little bit, but first thing probably what we will try to do. We will see we will try to do a some data cleaning.

(Refer Slide Time: 02:05)



First we have to do some data cleaning and if you look into the data you can see in the society there are quite a few you know empty variables are there empty cells are there. So, we can check how many empty cells are there by calling the extracting the society column equated to quote unquote and then take the length of just df society ok.
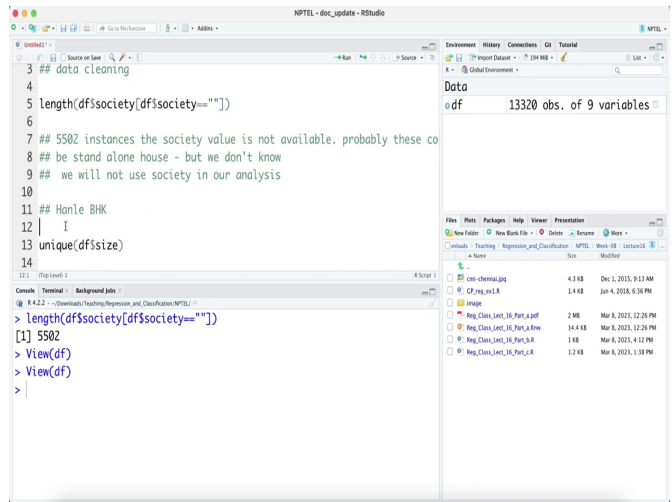
So, this will give you all the cases where you have quote unquote kind of thing and then if you just say length then it will give you. So, 51 5502 cases where society value is not available. So, you can see it in this way we can write it down have a note then write it down with a note in 5502 instances the society value is not available. We probably these are probably these could be; these could be stand alone house, but we do not know ok. So, the next thing is so, probably we will not use society in our analysis.

So, we will not use society as we did in Python we will not use society in our analysis as we did in Python. So, I am not ok. Next is BHK and you can see most of these BHKs are you know exactly like at the bedroom sometime BHK. So, we have to handle that. So, we have to handling the handle BHK variable.

(Refer Slide Time: 04:48)



So, how we do that first let us see how many unique BHK columns are there. It is written as size I believe, ok. Yeah size correct.

(Refer Slide Time: 05:06)



So, these are the BHK things.

(Refer Slide Time: 05:14)



So, now what I will do as we did in Python we will create another column called BHK equal to NA first. So, now there are 10th column there will be a 10th column with BHK all values are NA just we created a placeholder, ok.

Now, n equal to nrow of df is the number of row in the data set 13,320, ok. So, first thing we will do is let us take the first case it is to quote unquote BHK right. Now, what I am going to do I am going to called function called string, string split, ok. It split the elements of a character vector, ok.

(Refer Slide Time: 06:20)



So, what it will do is string, (Refer Time: 06:25) string split 1 ok. Let me just first run this.

(Refer Slide Time: 06:32)



I have to give like by what they would split. So, now we can see. So, that this was my original df size to BHK. Now, they splitted into 2 and BHK. So, now, but they have kept it as a list.

(Refer Slide Time: 06:54)



So, we have to put it as a unlist. We have to put it as a unlist, ok.

(Refer Slide Time: 07:10)



Now, if you just put it as a unlist. So, the first element is two next element is BHK. Now, what I am going to do and then if we just put say 1 then it will give me the first element which is 2. So, now what I am going to do I am going to run it through this line basically over the all cases 1 is to n; unlist and this will be like df dollar BHK on the ith cell and this will be ith component. So, let us run this thing. It is done.

Now, if I just go there you see it is filled up. But remember that when it was done it was character.

(Refer Slide Time: 08:08)



So, if I just check the structure of the data set df.

(Refer Slide Time: 08:13)



We will see BHK is kind of used as a is being stored as a character, ok.

(Refer Slide Time: 08:22)



So, now what I am going to do df as a df dollar BHK as dot numeric df dollar BHK. I think no problem. Now, if I just run this.

(Refer Slide Time: 08:51)



Now, you can see BHK is numeric no problem, ok alright. Next we have to if let us look into the data set. So, BHK now size was split and you know put it into two and you know the size of the like number numeric it was converted into numeric. Now, we have to look into the total square feet now because if you can see total square feets are considered as a character variable.

So, because there must there were lot of dash sign or you know very weird way of stored values were done. For example, this one 3010 or 2100 to 2850 this 31 row ok. The 31 row they have done it in a very way characteristic or string way of substituting the data, ok.

(Refer Slide Time: 10:02)



So, if you look into the df dollar total square feet 31 if you look into this it is a very string way of doing data.

(Refer Slide Time: 10:16)



So, now what we can do? We can just say ok string split df dollar total square feet, ok. I have to give up yes we have to give a say dash right.

(Refer Slide Time: 10:34)

(Refer Slide Time: 10:38)



Of course not this is I have to say 31. Sorry about that. So, this is my df 31 and then if I string split. So, by the dash it was splitted. The first element was 2100 and the second element was 2850, but both are character and it included some space and it was kept in screen. So, let us put it in some dummy variable, ok.

(Refer Slide Time: 11:05)



So, let us put it in some dummy variable.

(Refer Slide Time: 11:17)



And then first let us unlist the dummy variable; unlist the dummy variable. Now, it has the two element as a vector, but it both of them have some white space added to it.

(Refer Slide Time: 11:39)



So, what we will do? We will just trim the white space, ok. We will trim white space, ok. So, you can question mark trim white space yeah you see remove leading or trailing white space that is called trimws. So, now if you run it you can see that you know the white space are removed and it just one single element.

(Refer Slide Time: 12:18)



Now, what you can do what we can do is we can convert it to numeric, ok. And then we can take the mean of the 2 2475 is the number. And then what we can do we can just do the same operation over all i because for example, this is the one.

(Refer Slide Time: 12:54)



And then if I just run it through say 30. It is 1025. Now, if I run it through in this is perfectly fine.

(Refer Slide Time: 13:07)



I just convert it to numeric that should be fine, but if I just even if I run it through this process then it is just still 10020 , 1025 square feet. So, I can run this operation through the entire thing ok.

(Refer Slide Time: 13:24)



For i in 1 is to n. So, let me just copy this and paste it and ok I do not need this guy. I just need to put i here and then into finally, df total square feet i you just put this dummy variable and then you are done essentially. So, let me run this operation there 46 warnings, but you can see that now more or less there being there are no much. So, those 31 is 2475.754. So, now most of the cases have been taken care of no dash business are there. So, there is a big problem solving.
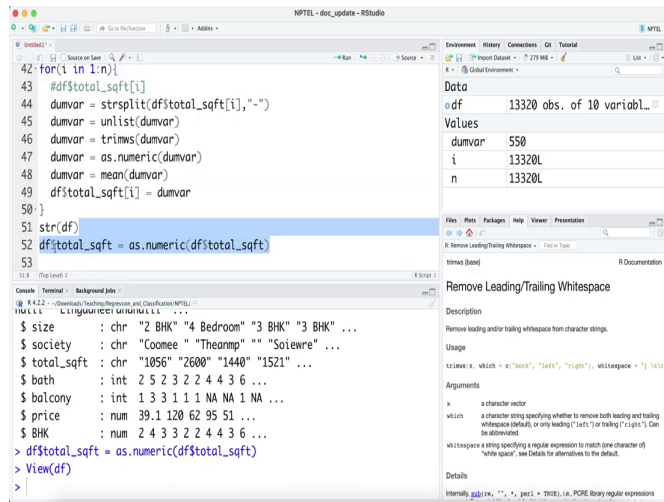
(Refer Slide Time: 14:37)



But still if you look into the structure of the df I think I have to just check whether it is still total character.

(Refer Slide Time: 14:48)



So, what we have to do? We have to say ok boss we need to put as dot numeric alright. And now we can see probably they have indeed gone into as dot numeric ok.

(Refer Slide Time: 15:08)



So, total square feet is now numeric. So, this is a big problem solving. Next we have to handle the location ok. So, location was location is very important and we have to handle the location. Why? Because there are few location few new quite a few location where the number of instances were very few.

So, how to handle it? df dollar location and first what we can do? We can just say table. So, it will give me a table of each location how many instances are there.

(Refer Slide Time: 16:03)



So, if I just say sum it will give me 13,320 you see, but I do not want it. I want the table in each locations are there.

(Refer Slide Time: 16:12)



And then I can just sort them.

(Refer Slide Time: 16:20)



And let me put it in some variable or some object loctn location ok.

(Refer Slide Time: 16:33)

(Refer Slide Time: 16:44)



And so, if you just say head of loctn then it will be like few first few maybe 10. These are all places where you have only one instance. So, these are the cases where you have only single instances and there are one instances where you have no space like location is not known. So, you just have some quote unquote kind of thing.

So, what we can do wherever the location is belongs to this say Anekal or any one Giri Nagar or something the location is less than say 10 or 20 or 15 some number. We consider is or say 20 let us take 20 then we will consider it as that we do not have enough instances in those case we will put it into others ok. So, how can we do that?

(Refer Slide Time: 17:48)



So, let us take one particular location say 3 ok on the third case is Uttarahalli.

(Refer Slide Time: 18:01)



Now, in the location; so, this is the name of the case in the location table if I just put that guy 186; that means, out of 13,320 cases instances there are 186 house, which belongs to the location Uttarahalli ok.

(Refer Slide Time: 18:27)



Now, if what we will try to do? We will try to secure that good boss if this location is say less than 20 then we will say this location as others ok. This is something what we will trying want to do, but we do not want to do it here we will do it in a little bit more fashionable way. We have to also take care of one problem that there is a one quote unquote thing is there. So, and how location table is handling this guy? So, if I just say df location.

(Refer Slide Time: 19:17)



If I just say df location equal to quote unquote and df dollar location, ok. So, there is one case.

Now, if I just put it there if I just put it there then it is a NA; that means, and we have to; that means, handle NA inside the if condition and that could be a bit of a challenge that could be a bit of a challenge.

(Refer Slide Time: 20:09)



So, there will be one NA that we will encounter. So, let us you know try to handle it for i in 1 is to n. First thing what we will do? We will take the location of the ith instances and put it into the location table and we will see what is the count give me the count and that count I will put it as a dum variable. So, either it will give me NA or it will give me an some number ok.

So, if dum variable is dot na then we will say ok boss in this particular case you just put others. Else if dum variable is less than 20; that means, the number of instances in the entire data set for that particular location is less than 20 then also you put it others else you just leave it there you do not worry ok. So, you know if I just run it, it run very fast and you see if there any others have come there may be one or something, ok we will see. So, at least I think the location thing is now kind of ok.

(Refer Slide Time: 22:15)



We as most of the data cleaning thing has been done. Now, we will do the model training. So, let us split the data in train and test. Split the data into train and test ok. Split the data into train and test. Now, so, the first thing I will do m equal to ceiling n is to 0.7 ok. So, that will give me. So, out of 13,320, 3 9324 instances we will take as a test data set sorry training data set. We will take it randomly and rest we will take it as a test data set.
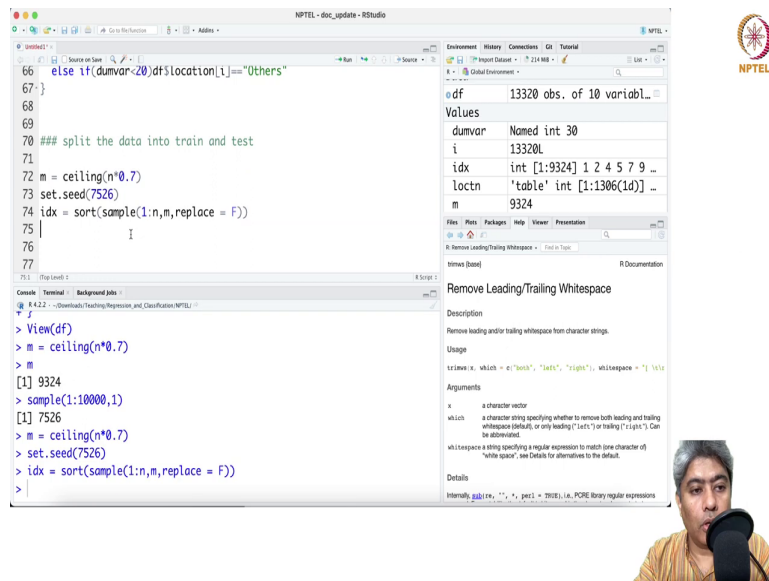
(Refer Slide Time: 23:17)



So, we have to first set the seed. So, how we will do that? So, typically the way we do it I do it is 1 is to say 10,000 I just draw one sample all I need 7526.
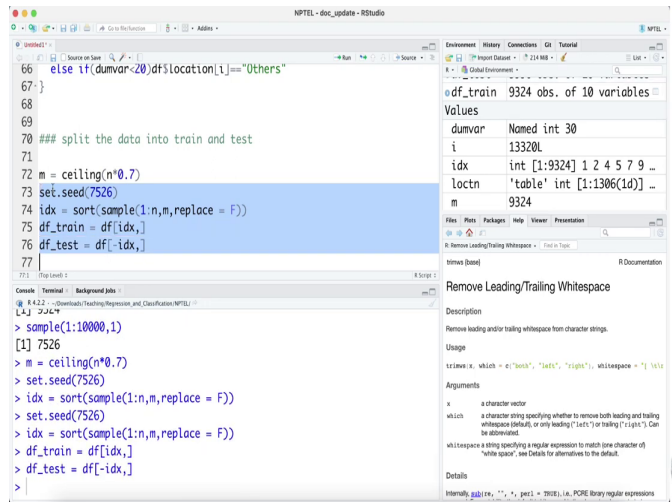
(Refer Slide Time: 23:31)



And that I will set up as a seed. Then what I am going to do? I am going to draw sample from 1 is to n, m, many sample and replace equal to false and then I must do a. It is a one time sort do not worry about that one time sort is fine and I just do idx. So, these are the idx I am going to draw and if you use the same seed you will get the same data set same split exactly my result.

(Refer Slide Time: 24:13)



And your result will be same. So, df train will be df with these idx comma. So, these are the row values and column values all columns will come in and df test will be minus of these idx. So, now this is df train ok and this is df test, ok alright next what I am going to do is we are going to fit our first model mod1 lm price.

(Refer Slide Time: 25:10)



Let me just get the colnames of the df so, price. So, that you know we do not make spelling mistake price total square feet plus area type plus balcony plus and BHK plus your bathroom. I think we have to take bathroom the sometime there are bathroom and then BHK there are yeah.

(Refer Slide Time: 26:25)



So, we can take bathroom instead of bathroom BHK and data equal to df train ok and at the end you just say summary of mod1. Remember note that I did not put location yet ok. I did not put location yet.

(Refer Slide Time: 26:54)



So, if you see either put area total square feet does have a strong correlation, strong effect very strong effect t value is 58, t value is very small area type plot type area has quite significant effect plot area the people the things, which are plot area they have a very little bit premium price people are paying. If bath is a number of bathroom increases then the premium or the price goes up, but if the number of BHK increases it looks like it is a negative this is bit worry the some I do not know why it is negative? And adjusted R squared we got 0.4251 with no location.

(Refer Slide Time: 28:00)



So, now, if we add location if we want to add location so, mod2 equal to what we will do update the model1 you update model1 how you do that dot comma dot plus you just put the location ok. So, what is the col names of df. So, this is the location.

So, you just put location and it should work alright and then let us see summary mod2 of location and you can see that adjusted r square now went up to 0.629 right ok.

(Refer Slide Time: 29:00)
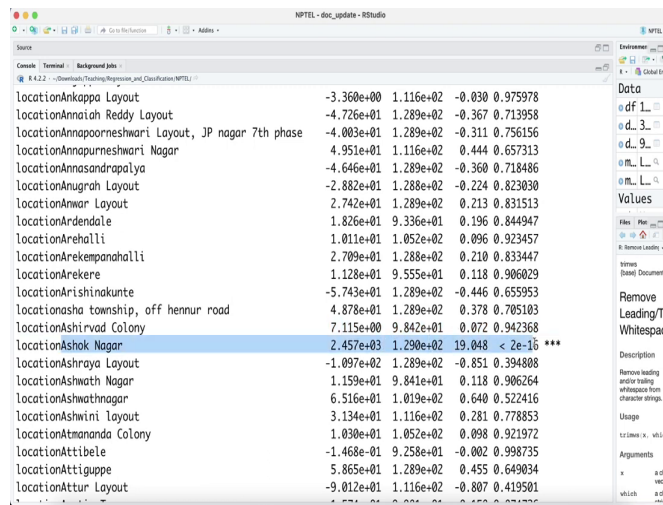


```
locationBannerghatta Road            2.210e+01  9.156e+01   0.241 0.809291
locationBapuji Layout                1.158e+02  1.290e+02   0.898 0.369216
locationBapuji Nagar                -8.589e+00  1.289e+02  -0.067 0.946876
locationBasapura                     6.042e+00  9.663e+01   0.063 0.950143
locationBasava Nagar                 1.611e+01  1.052e+02   0.153 0.878275
locationBasavanagara                 9.783e-01  1.052e+02   0.009 0.992583
locationBasavanapura                -1.452e+01  1.019e+02  -0.143 0.886672
locationBasavangudi                  1.007e+02  9.299e+01   1.083 0.278668
locationBasavanna Nagar              3.909e+00  1.116e+02   0.035 0.972056
locationBasaveshwara Nagar           6.086e+01  9.363e+01   0.650 0.515663
locationBasaveshwara Nagar Yelahanka 2.757e+01  9.982e+01   0.276 0.782391
locationbasaveshwarnagar             2.131e+01  1.289e+02   0.165 0.868651
locationBasaveswarnagar              2.051e+02  1.289e+02   1.591 0.111687
 [ reached getOption("max.print") -- omitted 946 rows ]
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 91.1 on 8092 degrees of freedom
  (86 observations deleted due to missingness)
Multiple R-squared:  0.675,     Adjusted R-squared:  0.629
F-statistic: 14.68 on 1145 and 8092 DF,  p-value: < 2.2e-16

>
```
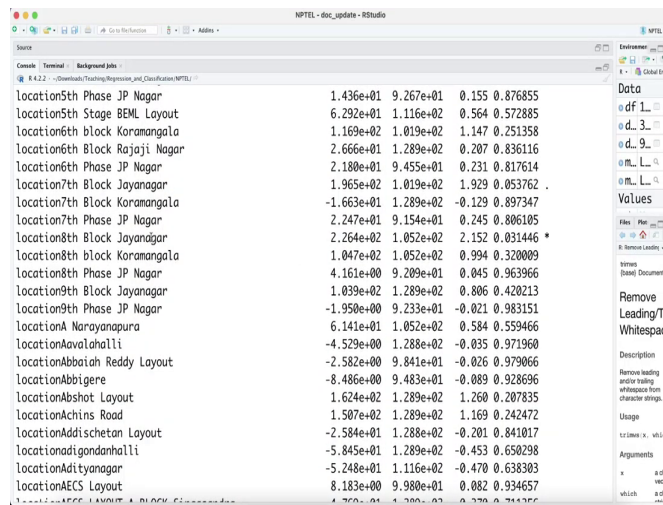
(Refer Slide Time: 29:06)



See Ashok Nagar has very strong positive premium looks like Ashok Nagar ok.

(Refer Slide Time: 29:30)

(Refer Slide Time: 29:35)



8th Block Jayanagar 5th Block Jayanagar 4 Bedroom house in Bengaluru ok then these are the in some of the location still required you know I think still the bedroom things are 7th block, 6th block this 4 Bedroom Farm House in Bagalur ok they have put the bedroom information in the location which is not right second stage are in the ok.

So, alright I think there are lot of things are being missed actually lot of areas are being missed, but adjusted R squared as now 0.629, which is similar to that of similar to that of Python.

(Refer Slide Time: 30:47)



Now, what I am going to do is very simple that I am going to like check some price like test data ok test data I am going to create some df test 1. So, this is df test data now this test data is a test point that we are going to try out the location we are going to use say first Rajaji Nagar.

I remember in Rajaji Nagar case we got a value of near about 200 lakhs or 2.4 crores something like that let us see in this time what happens what is the situation. So, now, we have changed it was Uttarahalli, but now in the test data we have given a location Rajaji Nagar.

(Refer Slide Time: 31:53)



And the total square feet will be total. So, total square feet will be say 1000 square feet ok. And then bathroom will be there will be two bathroom and it will be 2 BHK it will be 2 BHK.

(Refer Slide Time: 32:33)



So, let us run this through and now if I just predict; predict mod2 with newdata equal to test data 220.45. So, it is gives me a price of 220 lakh rupees so; that means, 2.2 crore in Python it was giving me 2.4 crore. So, it is somewhat in the same range remember that we also had made a model 1 initial model base model what happens if we use model 1. And model1 had a very low R square 45 percent in model1 it is giving me 52 lakh. So, model1 is looks like giving me a not so good also model.

(Refer Slide Time: 33:30)



Accuracy for model1 is adjusted R square is very low compared to may at least 20 percent lower than model2. So, location is very important location if you based on the location the house price varied extremely ok. Now one thing we can do is we can look into the train dataset ok in the train dataset among the df train location equal to Rajaji Nagar and we can look into the price. So, these are the instances which is the houses which belongs to the Rajaji Nagar and their prices.

(Refer Slide Time: 34:28)



Now, if I take the mean price or median price. So, the median price is near about 250.

And mean price average or average price is 320. So, given this I think predict model predicted model is giving 220 I think is a decently the correct value probably. And we can try few other models for example, we have not tried log linear model.

(Refer Slide Time: 35:16)



So, we can try log linear model. So, let us try that. So, what we can do? We can just take the mod1 and this will be the third model and we are trying log of price right log of price and log of total square feet. So, this will be log of bathroom I am giving 1 because in case there is 0 bathroom it should not happen [FL], let me not take if it there is 0 then it will throw error, but let us take this for just for fun let us see how this turns ok. I have to put location also.

(Refer Slide Time: 36:13)



Plus location ok, so it will take little time because it will create all the dummy variables and then it will fit the model, alright.

(Refer Slide Time: 36:31)



Summary mod3 so, now, after log written it has gone up to 85 percent ok so, this is a huge jump another 20 percent.

(Refer Slide Time: 36:50)



Just because I am getting, but that may not be a correct thing because remember that.

(Refer Slide Time: 36:52)



These are in the log scale log scale it is linear it is becoming linear and that is why and in the log bath and log BHK they are all now positive and everything is now pretty cool actually. So, maybe it will give us a much better how accuracy. So, what I will do? I will just take this thing.

(Refer Slide Time: 37:24)



And we can just use this model 3 to predict the same data set. So, test data is still that same Rajaji Nagar with 100 square feet and 2 bathroom and 2 BHK, but now it will give me log scale the price will be log scale.

So, I have to put it in e to the power remember that I have to put it in e to the power. So, it is now giving me 107s ok. So, I do not know how good or how bad it is. So, this is looks like this could model looks good, but it is looks slighlty slightly giving me under value now is it a right model. So, best way to check it out what will be the test cases in the test data frame what would be the predicted price.
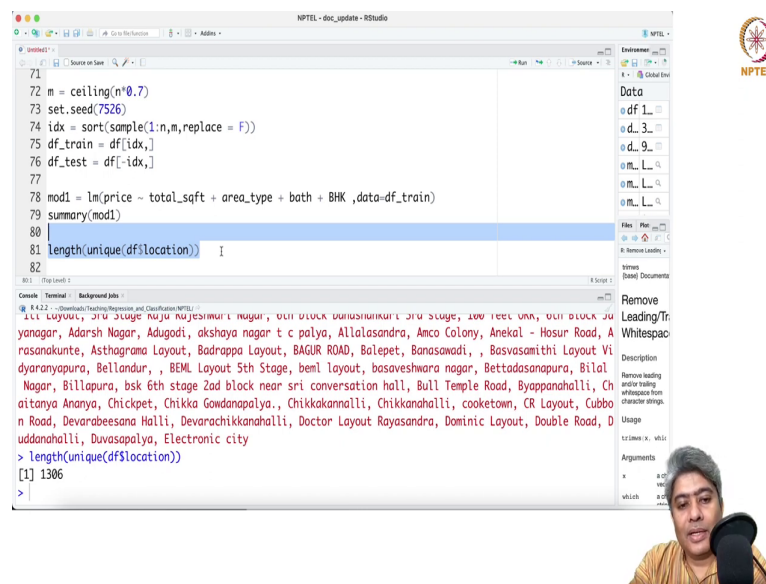
So, test df test dollar price hat what I will do? I will just put the test hat1 ok predict 1 df test here instead of that I will just put df test and then df test 2 that will be model 2 ok.

There is some error its throwing because in the test second model its same location has new level ok that is this is bit of a problem always I have found when we have too many you have too many. So, basically what is happening the issue is I will tell you if there are cases where the test dataset has some values which are not there in the not there in the train data set.

So, when the split happened the split was not correctly I mean has not considered every everybody in from the thing. So, one way of doing it is from each test area each area we can split it in train and test. So, that is a one way of definitely one way of doing it that will be a yeah then that is a good idea actually let me try that yeah.

So, this splitting what is happening let me repeat the problem ok please pause your video you may take a break come back, but let me give you tell you this problem I face this problem of 10 time, but let me tell you this problem. In the location there are too many locations.
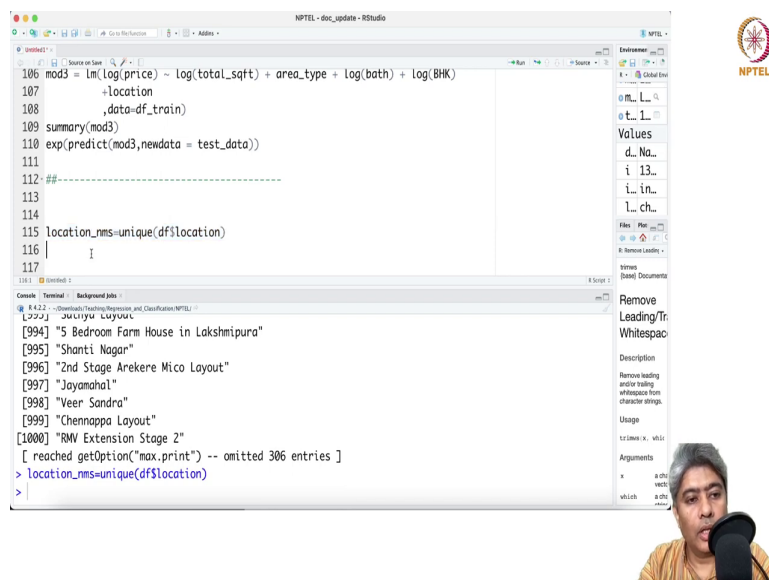
(Refer Slide Time: 40:42)



In fact, we can look into the location there are df dollar location and if I just say unique length ok length there are 1306 unique location. Now, when I am splitting the data I am randomly splitting and quite a few of them few locations all samples of locations have gone into the test so; that means, those locations are completely missing in the training dataset.

So, for those location the training dataset did not able to train and as a result now it is throwing that in the test dataset there are because if you look into what it is saying it is saying

that model data frame default error in model factors location has new level the test dataset has new levels, which is not available in the training datasets.

So, that is what it is talking about. So, we have to split the data in so, way that all 1306 location will split into 7030 way then each location will have a representative of 7 representative of samples in their training dataset ok you understand. So, let us try to make this thing it is going to be a bit of a let me try this let us try this ok.

(Refer Slide Time: 42:30)



Now, so, first is location names equal to what I will do? I will just take a yeah this guy and put it there. So, here I have all the names now I will just put it in there ok.

(Refer Slide Time: 43:17)



And then what I am going to do is for it is going to be a bit of a challenge df train df train equal to df 1 and kind of df test equal to df 1 we will just later put it as a drop these first row just to I just need a data structure at a framework ok. So, if I now I have just the framework the structure.

(Refer Slide Time: 43:59)



So, if you just see I have some structure and later I will keep appending on that. So, that is that is my idea and     later I will just drop the first row from both frame.

(Refer Slide Time: 44:15)



Now what I am going to do I am going to say that ok df sub df equal to df dollar location equal to location names comma. So, if I just do that. So, these are the cases where you have bunch of names.

(Refer Slide Time: 44:52)



Say electronic phase let me just df call it df sub ok. So, let me just put it in df sub.

So, what is that 132 observations. So, these are the instances there are 132 observations, which all belongs to Electronic City Phase 2 this location. Now, I have to split this guy into 2.

(Refer Slide Time: 45:25)



So, what I am going to do n1 equal to nrow of df sub oopsie df sub. And then what I am going to do is m1 is equal to ceiling m n1 times 0.7 and then idx equal to 1 sort sample 1 is to n n1 comma m1 replace equals to false alright.

So, this is what I am going to do and then df train is equal to df sub rbind dot data dot frame and df train comma df sub with idx that we have here ok. So, if I just now do that.

(Refer Slide Time: 47:09)



So, now, you have df train as this so, the first thing and then ok.

(Refer Slide Time: 47:31)

(Refer Slide Time: 47:42)



And then similarly df test equal to df test and minus of that is going to give me the all test data set that half 40, 40 out of 94 of them goes to or 93 of them goes to train data set and 39 of them goes to test data set ok, alright.

(Refer Slide Time: 48:00)



So, m m1 is 93 and n1 minus m1 n1 minus m1 is 39. So, 39 goes to so, that is why it is 39 plus above on 40 and it is 93 plus 140 because we have one row added first we will drop that later.

So, now, whole thing I am going to repeat in a row for loop for i in 1 is to length of n g t h length of location names ok. And this is going to be the I alright and you have df train and test here and I just need a set dot seed here remember that this is this method is ok.

(Refer Slide Time: 49:15)

(Refer Slide Time: 49:28)



Now, you have train is this many and test is this many and what I will do? I will just say in df train equal to df train minus 1.

So, I have 9921. Now, it is 9920 because remember that the first row I want to drop right I want to drop the first row.

(Refer Slide Time: 50:06)



And similarly df test I want to drop the first row. Now my things are ready.

(Refer Slide Time: 50:15)



So, let me now run the model let me just copy this model and run it here ok.

(Refer Slide Time: 50:25)

```
Residuals:
    Min      1Q  Median      3Q     Max
-2462.68  -25.82   -9.02   11.04 2814.30

Coefficients:
                            Estimate Std. Error t value Pr(>|t|)
(Intercept)               -58.158524   3.905626 -14.891  < 2e-16 ***
total_sqft                  0.074339   0.001089  68.248  < 2e-16 ***
area_typeCarpet  Area       4.990933  14.832313   0.336    0.737
area_typePlot  Area        73.173201   4.191334  17.458  < 2e-16 ***
area_typeSuper built-up  Area  2.889033   3.077621   0.939    0.348
bath                       25.297041   1.984348  12.748  < 2e-16 ***
BHK                        -8.969711   2.043348  -4.390 1.15e-05 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 114.8 on 9820 degrees of freedom
  (93 observations deleted due to missingness)
Multiple R-squared:  0.4679,     Adjusted R-squared:  0.4675
F-statistic:  1439 on 6 and 9820 DF,  p-value: < 2.2e-16

>
```
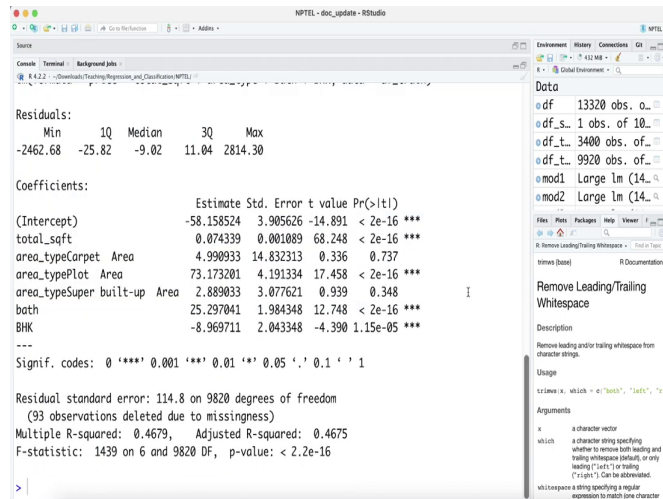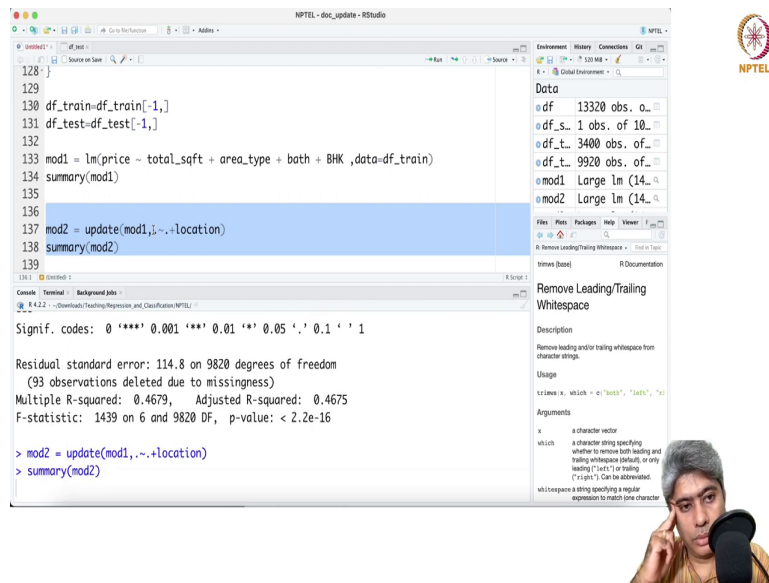
So, this is 46 percent about 46 percent accuracy total square feet BHK is negative which is bit of a I do not know this model has some problem and then we go and try to fit model2 ok.

(Refer Slide Time: 50:51)



It will take little bit time because it will create all the dummy variables.

(Refer Slide Time: 51:04)



Now, it is 63.8 percent accuracy 63.8 percent accuracy.

(Refer Slide Time: 51:13)



And Ashok Nagar has a very high p value.

(Refer Slide Time: 51:21)



```
location4th Block Jayanagar                                   2.052 0.042176 *
location4th Block Koramangala                                 1.661 0.096766 .
location4th Phase JP Nagar                                    0.394 0.693736
location4th T block Jayanagar                                 0.963 0.335777
location5th Block Hbr Layout                                  0.443 0.657522
location5th Block Jayanagar                                   5.665 1.52e-08 ***
location5th block Koramangala                                 2.027 0.042733 *
location5th Phase JP Nagar                                    0.090 0.928481
location5th Stage BEML Layout                                 0.303 0.762000
location6th block banashankari 3rd stage, 100 feet ORR       2.602 0.009288 **
location6th Block Jayanagar                                   0.874 0.382343
location6th block Koramangala                                 1.102 0.270623
location6th Block Rajaji Nagar                                0.054 0.957078
location6th Phase JP Nagar                                    0.411 0.681251
location7th Block Jayanagar                                   2.159 0.030869 *
location7th Block Koramangala                                 0.011 0.991312
location7th Phase JP Nagar                                    0.299 0.764853
location8th Block Jayanagar                                   1.970 0.048833 *
location8th block Koramangala                                 1.002 0.316550
location8th Phase JP Nagar                                    0.046 0.962978
location9th Block Jayanagar                                   0.575 0.565213
location9th Phase JP Nagar                                    0.005 0.996238
locationA Narayanapura                                        0.210 0.833737
locationAavalahalli                                          -0.005 0.996035
```

(Refer Slide Time: 51:26)



Jayanagar, Koramangala has a very high p value alright. So, we will see what is it is happening here ok and then finally, we will take the third model and we will see what is happening here.

(Refer Slide Time: 51:58)

(Refer Slide Time: 52:13)



So, in the third model it is giving us accuracy up to 82 percent, but these are in sample accuracy remember that. So, you have to be bit careful about these kind of accuracies. So, if you run. So, first we have to check out of the sample accuracy. So, let me just copy these two lines and paste it here ok.

And similarly if I see now there is no problem now there is no problem and it is it gives us the exact accuracy thing we wish it.

(Refer Slide Time: 53:00)



And then this is the third one third model and if we just look into the test.

So, we have price hat from first model second model and the third model ok. So, the third model looks like it is little bit maybe it is not that great model I do not know some issue is there it is all same places maybe, we will see what is happening there not test data.

(Refer Slide Time: 53:53)



We want a df test we want df test to be here df test there was a mistake now it is doing well alright.

Now, what I am going to do now if you look into this particular case the price is 150 price hat to the first model predicted 251, second model predicted 220 and third model predicted 138. So, this from this one case looks like the first third model has done good job, but we have to first compute the root mean square error. So, what we will do?

We will take the df test price 1 hat minus df test dollar price square it mean it take a mean of that take square root of that. So, ok I have to take na dot rm equal true some places probably it did not able to. So, this is my RMSE1.

(Refer Slide Time: 55:22)



Out of the sample RMSE RMSE1 this is out sample remember that this is my out sample RMSE and second is from the second model. This is little less and the third one is from the third model ok and this is very less. So, third model is not only doing well it is doing pretty well to great extent.

So, log; that means, basically the total square feet bathroom BHK and price they all have a log linear relationship and that is why even in the out of the sample they are doing far superior than the just sample linear regression ok. There and let us we can also compute the correlation between the two.

(Refer Slide Time: 56:35)



We can just take cor complete yeah and if you just take square. So, this is really less and then if you have just take 2 42 percent and then 3 this is approximate r square is 62 percent out of the sample approximate R square these are like out of the sample out sample approximate R square s q u a r e.

So, it shows that third model is has pretty good predictive power in even out of the sample it can predict the price with 62 percent variability can be explained, which is very good some out of over fitting is happening, but that is can be handled. So, with that I will stop here I hope you enjoyed this hands-on session. So, see you in the next video.

Thank you.