**Marketing Analytics**
**Professor Swagato Chatterjee**
**Vinod Gupta School of Management,**
**Indian Institute of Technology, Kharagpur**
**Lecture 56**
**Text Mining and Sentiment Analytics (Contd.)**

Hello everybody. Welcome to Marketing Analytics course, this is Doctor Swagato Chatterjee from VGSOM, IIT Kharagpur who is taking this course and we are in week 11 and we are still discussing Text Mining and Sentiment Analysis and we will continue the same thing in this particular week also. Till the last week we have done basic sentiment mining, basic text mining and we have also done some bit of spam detection and etcetera using Naive Bayes algorithm here we will actually use that in sentiment mining.

So we will use Naive Bayes for sentiment mining in this session. In the next session we will use certain libraries which has been pre-developed by some researchers for various kinds of application in various places and later point of time we will do a little bit of topic modeling also. So in this particular class we will be using the dataset called hotel_review. So this is the same old dataset which have been used earlier by us. But here I have only included the text part.

(Refer Slide Time: 01:24)



So if you see this dataset carefully that this dataset has 4 columns which is the hotel name, the review title, the overall rating that you have got and the review content the basic text of

the review and these are all content like this place is very nice to stay close to the railway station. The owners are the base and will help you with all your questions.

We would definitely recommend this place blah, blah, blah all of these things, so each of this thing is one review. So the quantitative version of this particular thing, the dataset we have used in week 1 if you remember. In the week 1 we started with a discussion that if the overall rating is there and if you also get let us say the service, location, service, quality you get the ratings of all those stuff.

Then how various service attributes are impacting your ultimate overall satisfaction or overall rating that you have given is sometimes we find out through regression. We also try to find out that which of the aspects are more important in your service context, which aspects are not important, where you are doing good, where you are doing bad and blah, blah, blah. The same dataset we have been using here, but with the text part.

So I want to know that what kind of extra insights we can bring in using the text also till now we are only using quantitative data now I am bringing in qualitative data in this particular model building to find out that how we can get enough information from the qualitative data also to make certain sense out of it. So in this particular class we will do that. So we will be using as usual R programming and the Naive Bayes algorithm to do that.

(Refer Slide Time: 03:21)

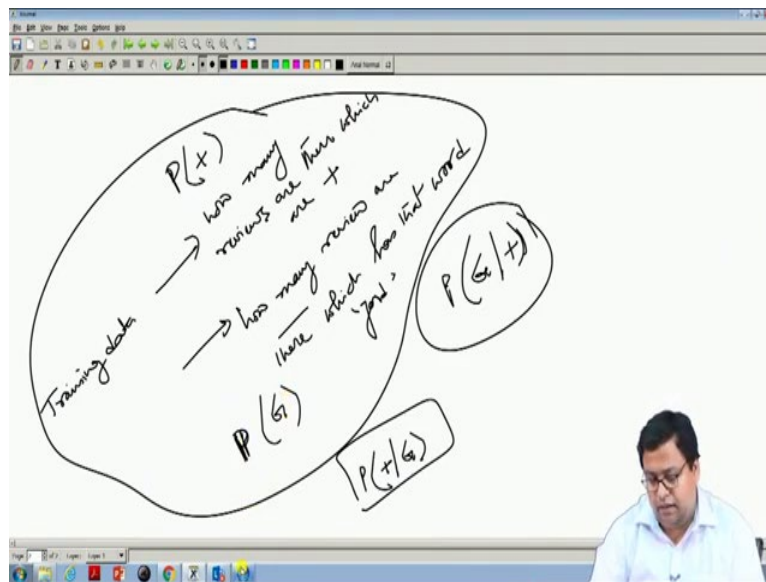So if you remember a small recap of the Naive Bayes algorithm it says that

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A).P(A)}{P(B)}$$

This is something that we will be using.

For example probability of positive sentiment let us say positive sentiment given a word called good.

Let us say if that is so probability of this given G, that is I have write it in this way, this will be $P(G|+)/P(G)$. Now in our old dataset if I have a dataset and if I break into training and testing, there in a training data from the training data I can have an idea that.

(Refer Slide Time: 04:25)



In the training data I can have an idea that how many I would say reviews are there which are positive and I can also get a idea that how many reviews are there which has this word has the word called good. So this will be basically probability of plus and probability of G $P(G|+)$ and then we have to find out by doing another calculation what is this. We can find out this thing then I can calculate probability of positive given $G.P(+|G)$.

So these 3 simple data mining will be done by Naive Bayes algorithm. So the first thing is that so these are P actually I have written it in a different way this is it should be P, okay. So that is something that we will be doing here.

(Refer Slide Time: 05:42)





So if you check this particular Senti.R this is the R file. So we are using the same old formula of Naive Bayes, but in this context in this particular context which is our review dataset we will be using this thing. So the first job is to read the data. So session, set working directory to source file location and then after setting the working directory you read the data and the data has 942 observations of 4 variables.
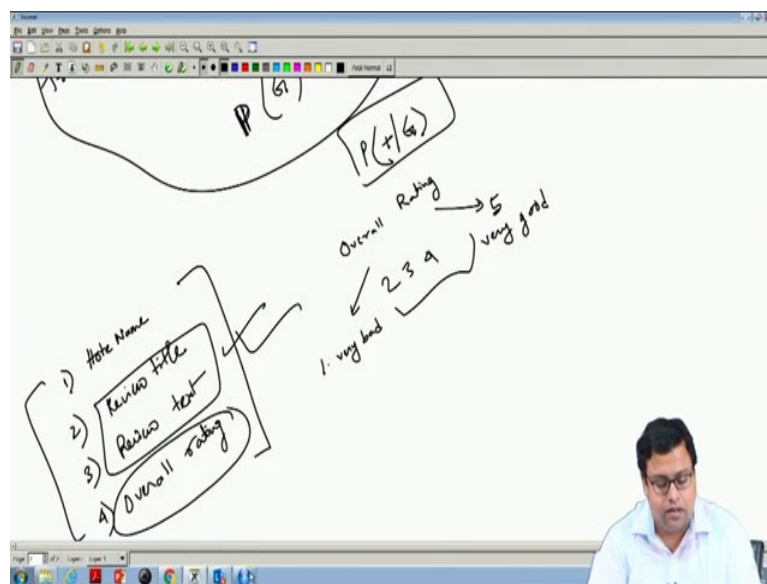
So how the data looks like? The data has the hotel name, the data has the review title, data has the review text and the data also has overall ratings. So these are the 4 things that our dataset has and we will try to see that how this overall rating is related to this two stuff that if you title and then if you text.

(Refer Slide Time: 06:52)





Now if we come to this thing first thing that is we want to do is structure of the data. So the structure of the data says that my dataset has basically if you carefully look at the structure then my dataset has 942 observations of 4 variables there are 4 variables and there are 942 observations means 942 rows and 4 columns are there and the columns names are hotel name city which is a factor variable with 23 levels.
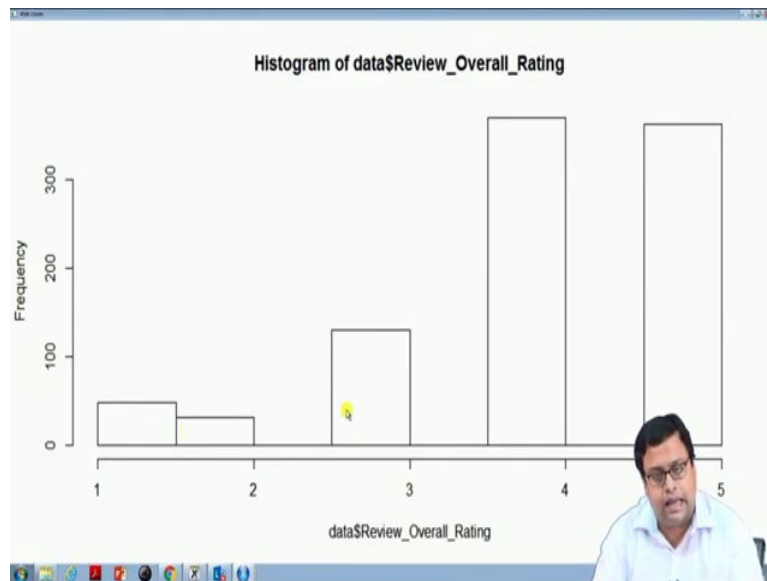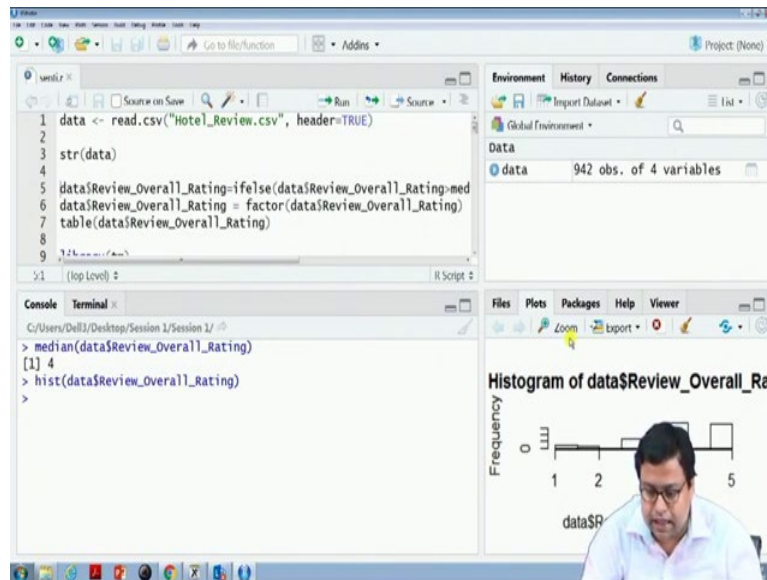
That means there are 23 different hotels data has been captured in this particular dataset and the review title has, review title is a factor of 923 levels that basically there are 942 different kinds of review titles are there it is coming 923 because some review titles are probably same let us good stay has come more than once. So if that is the case then we cannot say that they are different review title.

There can be some common review title for example very good hotel let us say that was the review title, so we cannot do anything with that. Then comes the review overall rating. Review overall rating is the quantitative rating in a 1 to 5 point scale that somebody has given 1 means in this overall rating 1 means basically very bad and 5 means very good and anything in between that let us say 2, 3, 4 these are mediocre.

So depending on how much is the cutoff we will decide whether it is a positive review or negative review and then comes review content which is the text and there are 941 levels because probably one thing what there is exactly two reviews which are similar to each other otherwise they are dissimilar to each other. So we got basically these 942 observations. Now I just press control L to clean my console control L is to clean my console.

And once I clean the console next what will I do. So what I do right now is I break the dataset the currently given dataset into two halves. One is the positive reviews and another is the negative reviews. So (what do I) how to find out positive review? We call it something called median split. What is median split? Find out the median of a variable and then split it.

(Refer Slide Time: 09:23)





For example let us say if I just find out median of data dollar review overall rating. So that is coming up to be 4. So frankly speaking anything any review rating which is and if I just put a histogram of the same thing this is how it looks like the histogram. So which one will I say positive review and which one will I say negative review.

So can I say that all these 5 are positive, these 4 probably are also positive so majorly biased towards positive review, but 3, 2, 1 are obviously negative reviews these guys are all negative review so I can say that. So if 3, 2, 1 are negative reviews 4 and 5 are positive review.

(Refer Slide Time: 10:11)



So what I will do here is data dollar review overall rating is let us say if they are greater than or let us say greater than equal to if I just write greater than that means it should has to be higher than 4.

If it is not higher than 4 then it will be negative. So data review overall rating is if else data dollar review overall rating greater than the median of that then it is 1 otherwise 0. So change this thing to 1, 0 is what I am trying to do. So I am changing it to 1, 0. So I will later rather say greater than equal to okay because if I just write greater than that means only 5 will be considered as positive review 1, 2, 3, 4 will be considered as negative review which is not the case probably.

So greater than equal to I will say up to 4 it is 4 and 5 it is positive review and otherwise it is negative review. So I am changing it to 1, 0. So now if I just try to see table data dollar review overall rating. I have 210 0s and 732 1s so this is my situation right now. So next what will I do I will make them factor.

This 1, 0 instead of considering as 0 and 1, two numbers we should consider them as two groups, two factors. So I am changing it to two factors. Now if I want to see the table of this thing the same thing whatever I told will be there so 210. If we get an idea 210 / 942 that comes up to be 22 % are negative and around 78 % are positive, but still this is biased but this not heavily biased so that we have to use something else it is not like that it is biased, but we can still use this one. So we will be going ahead.

Now the library that I call is this tm library text mining library. So I will call this library there is no package. So what we will do is we will install this thing, install tm we will take some time. Okay it has install now, so now I will call this library called tm. So I called this library it is saying that tm has been built under 3.6.3, but it is okay as long as the warning message and as long as the other functionalities are there we can go ahead.

So then as usual we are changing this data dollar review content which looks like this. So we are dealing with the review content. Why we are dealing with review content, somebody else will tell you sometimes later I will tell you that review content is something where the text is of significant size. On other than review I should not have actually click review content the whole thing because there are lots of content and it will get printed here, yeah so it got printed.

So this is something which is of huge text size. You will see that this is one content the ambience and the climate over there tells it all the energy blah, blah, blah so all of this the whole content is of significant size.

(Refer Slide Time: 13:51)

On the other hand if I just printed review title review titles were small so there are not enough words to analyze the review title that whether they are impacting the overall rating, but they might be impacting more than other things.

So anyways we will deal with review content so this content we change it to corpus form as we did in the last video. Now once we have change it to corpus form we print the corpus and the corpus has basically a metadata which has 942 documents and if I want to inspect the first 10 documents these are the first 10 documents. So 1, 2 so these are all text that is there. Now you see in this text there are lots of punctuations marks are there.

Lot of numbers are there, lots of various other things are there. So these things may not have any meaning. So the first job as we did any dataset is to probably do a little of preprocessing of the data.

So this is very important mode so in the text context that you have to preprocess your data. So what we will do to preprocess? I first content transform to lower that means I change the whole dataset to the lower case.

So this is equal to sign, this is basically equal to sign. We have discussed in the first class that this is also how equal to signs are written. So I run that and in corpus clean the all everything that is there in corpus clean is basically a in the lower form. So if I just write corpus clean you see that here this i has been changed to its lower case. Let us say visited in end April vary so everything after full stop has become lower case.
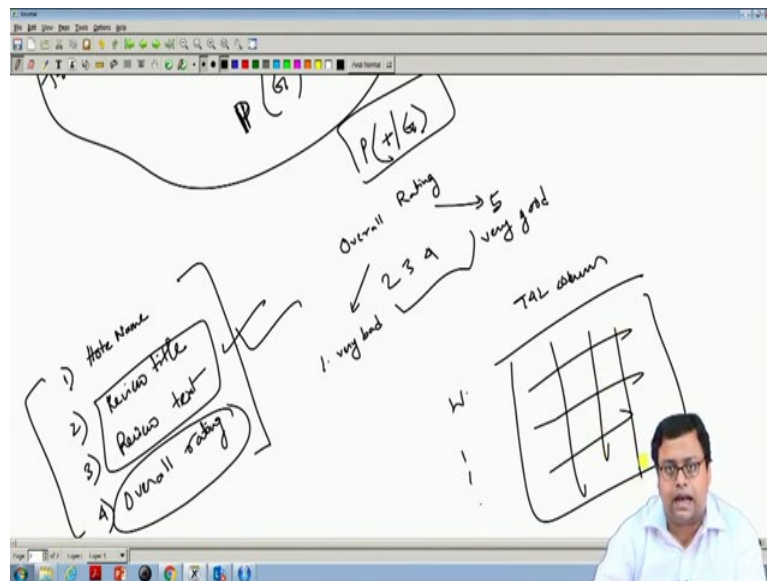
So the whole dataset is right now in lower case then I clean, I remove the numbers I have remove the numbers then I remove the stop words so what are the stop words. Stop words are generic words which do not have any contribution in the meaning of the dataset so we remove the stop word, we remove punctuations, we remove white space. So one by one I am removing lots of warning messages are coming up right now we will not bother later we will see.
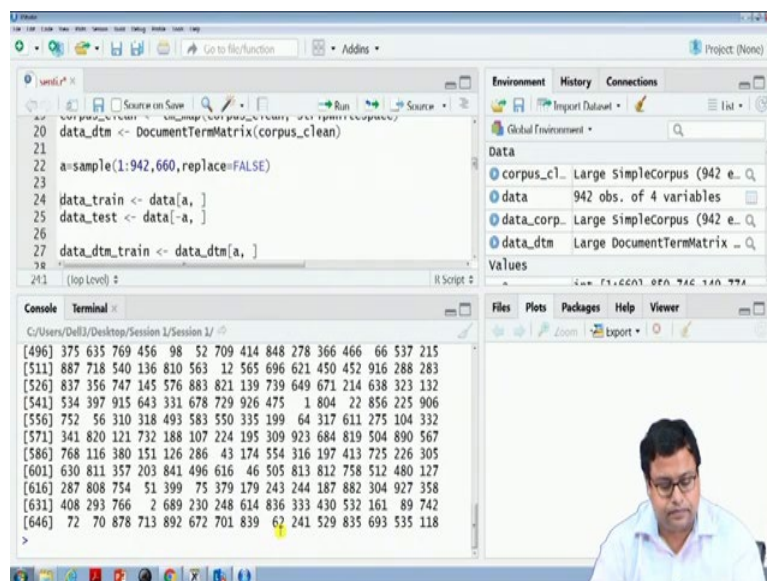
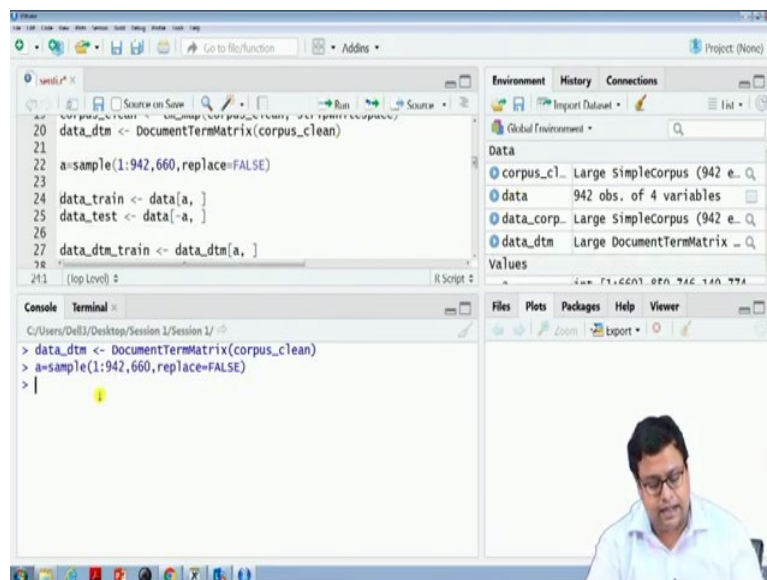And then what I do is with this dataset I create a document term matrix. So in my document term matrix it is also a large document and basically it is document of matrix looks like this that there are lots of words, word 1, 2 and there are 942 columns and this is basically a matrix where the counts are there, all the counters how many times one particular word is occurring in which particular review that is there.
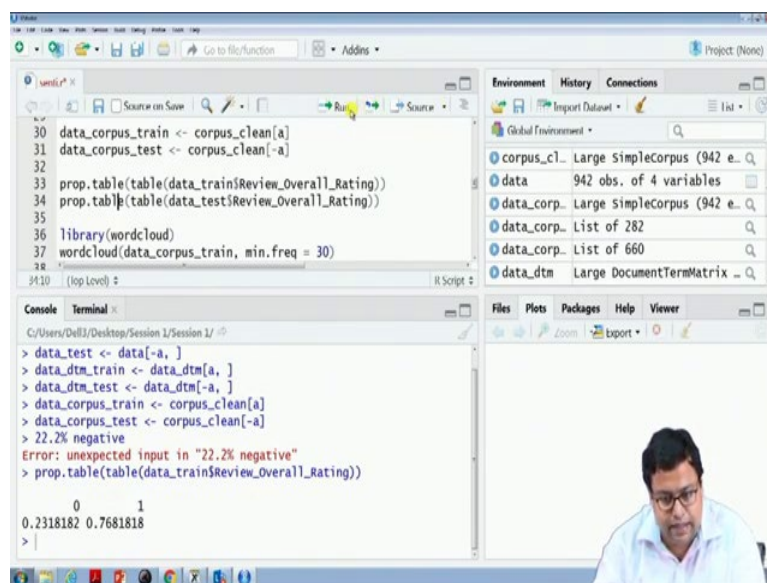
Now obviously I have to break it in training and testing. So if you remember in the last class we have done training testing by taking the first 70 % or 80 % as training and the last 20 % as testing that we cannot do here because there are various people. So your past review is actually sometimes impact your future review. So how the person who has reviewed last, the last review that I see.

When I go and write my review is also see the last review and if I see that review is positive sometimes my mind also gets diverted towards positivity or negativity. So all I am trying to say is that other people's review might impact you. So that is why we cannot take the top 400 as testing and lower 200 as testing or training something like it we have to do it randomly. So for that I am using a sample function.

So sample function samples 660 observations from 1 to 942 in this column I create vector called 1 to 942 out of that 660 observations will be randomly picked up and replacement is equal to false means once it is picked up it will never be picked up again. It is not a sampling with replacing, it is sampling without replacing means if 66 has been selected once 66 this particular number cannot be selected once more.
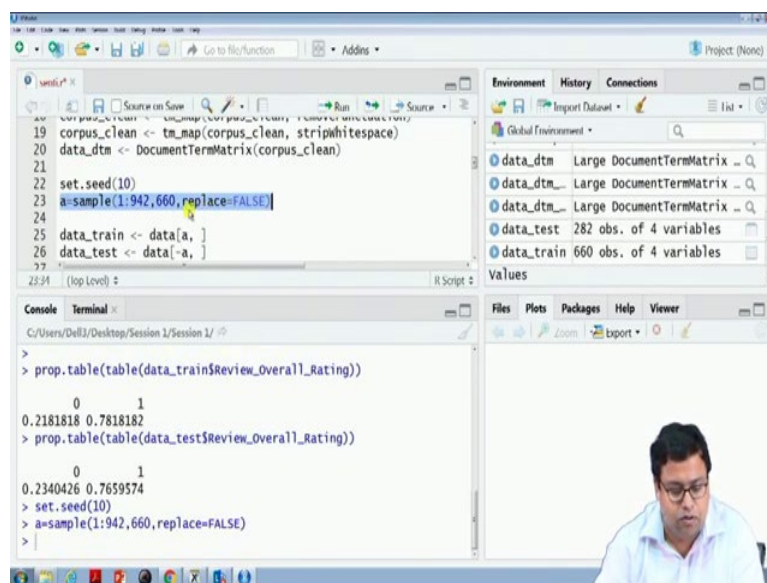
So that kind of sampling we are doing and putting it in a. So my a is right now some random numbers that got generated from 1 to 942 and the length is 660.
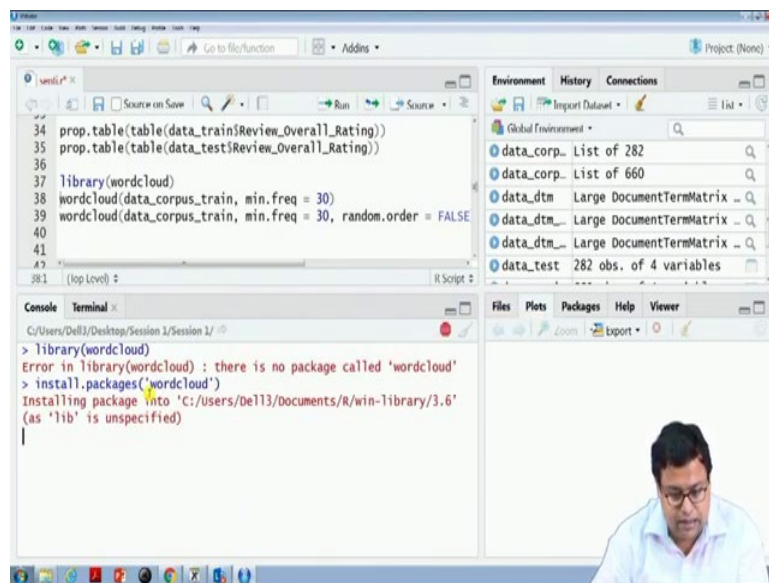
(Refer Slide Time: 18:51)

Now with this row names I am creating the training data and -a means the rest of the row names goes to the testing data so data dollar train. Similarly, if you remember in the last class also we created training, testing for the DTM and training testing for the clean corpus.

Because we will create I would say certain kind of word clouds using the corpus and to see that whether word is a way to find out whether that it is a positive review or negative review and then we will be doing our Naive Bayes using the document term matrix. So next what I do is as I told that there are if you remember 22.2 % positive, sorry negative this was there 22.2 % negative.

So if I just see currently what is the case you will see that 23.1 % are negative in your training data and your testing data 20.2 %. So not very different it is little bit different not very different. So I can probably sample once more and then I can do it to check that whether the prop table has improved a little bit or not. So let me just try out that so 21 and 23 it is okay I think that is fine.

So to have a replaceable or something which is obvious whatever I get and you get will be same you write set seed set .seed is let us say some number any number 10. So now my results and your results will come same. So now set seed and this once more and then the training data, testing data once more then the prop tables and etcetera. So the document term matrix and proportion table. So 22.77 very good and here also 21 and 0.79 so we will go ahead with this.

Now what I do is I call a library called word cloud. So let us see okay so this package not there, so install . packages word cloud. So let us install this one, so it got installed. Now with this package I run this thing.

(Refer Slide Time: 21:22)

Okay, I have to call the package first and with this package I run this thing. So the most common word is basically rooms fair enough. And then valley, friendly, helpful, some night some other words are also there which is coming up. Now rooms has to be the most common word so I should have probably remove that as my stop words, but it is okay. Now if I do the same thing for the training data without this randomness then hotel, room etcetera comes in the middle. So hotel, room, good food, room and rooms are same. So I should have done something called stemming to remove this room and rooms thing and then the rest.
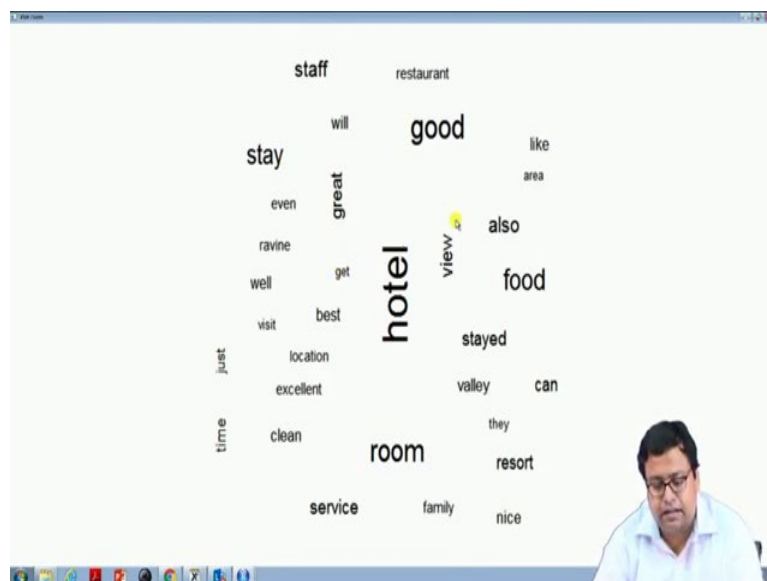
Now I break the dataset based on high and low. So this is high rating while the overall rating is 1 and the low rating of training data while the overall rating is 0. So I am breaking the training dataset into further two halves high and low. Then I am creating a word cloud for high and I will create a word cloud for low.

And I will compare that whether they are coming similar or not so let us create a word cloud for high. So high in this case hotel is common, but good room, great these are the words that are coming up just check remember this. So good food, staff, stay, service, great these are the words that are coming up.

Now if I just try out the same thing with a different I will just clean this and then I will run with a different which is low then what comes up. So hotel room, rooms are common in both the cases, but see this good, great has gone away fair enough staff, food, room, rooms hotel these are all, but great is there in the much smaller size and practically good is not there if you see carefully.

So then this adjectives not the nouns, but probably the adjectives are something which can contribute towards your positiveness or negative and that makes sense that when we do a review of anything.

If you say that okay I am a good teacher, teacher is a noun, good is an adjective. So by this good or bad or lazy or very energetic these are the words that we use to define how this teacher is. So teacher is not creating any kind of positive or negative sentiment emotion it is not creating this adjectives are creating so that can be one learning from basic explanatory analysis of putting a word to.

(Refer Slide Time: 24:36)





Now what I will do I will find out the frequent terms up to 5 frequent terms. So there are 1607 entries which has occurred at least 5 times make a dictionary with this words I am making a dictionary with these words and then what I will do is I will create a document term matrix only with this word for the training data also and for the testing data also.  So two
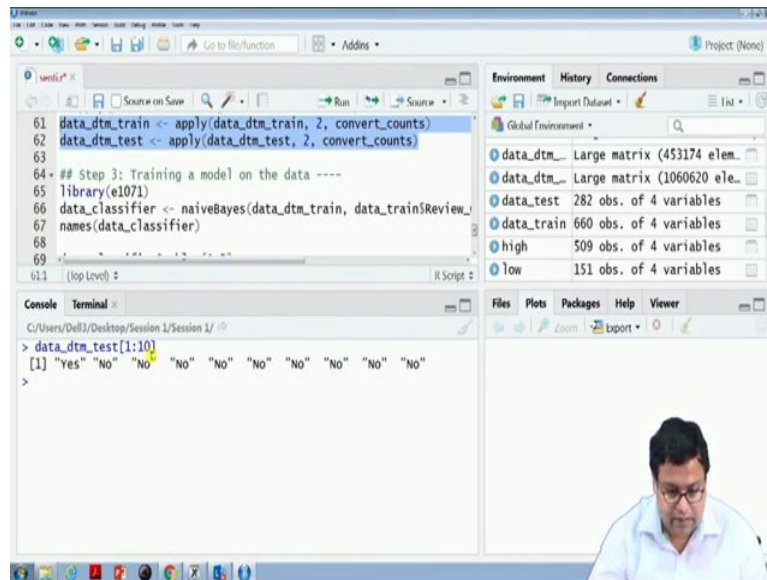
document term matrix once more for the training data and for the testing data, but only with this 1607 words.
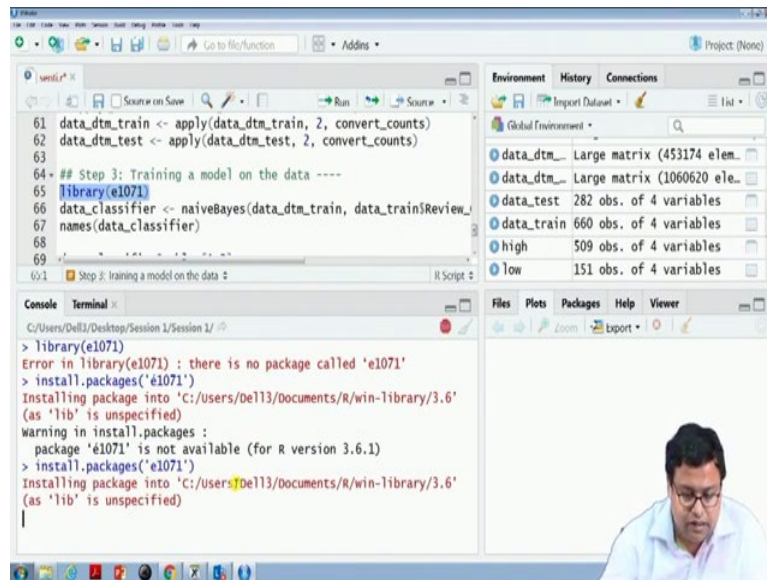
(Refer Slide Time: 25:10)





Now remember document term matrix gives you counts I need 10 in the last week we have done that so I convert it to basically counts rather than from counts to 1, 0 yes no I am converting it to that. So it has done that. Now if I just check what is data DTM _ test if I just try to see that okay so this I will not be able to see this here.
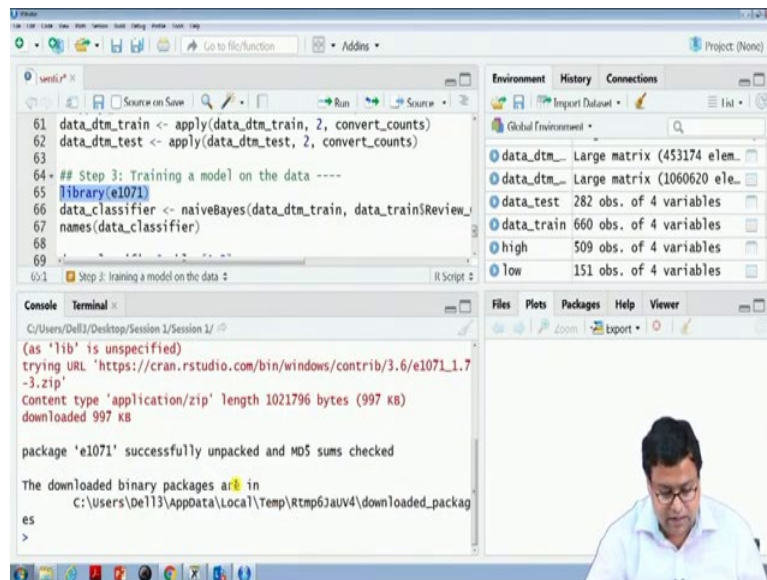
(Refer Slide Time: 25:41)



So let us say data 10 1 to 10. Now this is not how I will see basically the data this is my training data. So these are one to the word characters and their corresponding these are the terms and these are basically the corresponding values 1, 2, 3, 4, 5.

So I will now use this document term matrix training to my model building. So the library that we will use for Naive Bayes algorithm is E1071 that is not there. So install packages E1071 oh sorry I should have written the e properly. Okay so it is installing E1071 right now it has installed it.

(Refer Slide Time: 26:48)





Now I will call this library it is okay warning message, but we do not care. Now _classifier the model name is Naive Bayes this is the algorithm name while this is my x axis which is the term matrix and the y variable is basically the overall rating that you got 1, 0 rating. So if I now run a simple this time I got the classifier.

And if I want to see the probabilities of the classifier I can see for at least first two you see bathrooms this is 0.91, 0.95 and this is 0.8 I do not think bathroom is something they are based on which things are changing a lot. Big is not at all changing, but why do not I find out group let us find out group what this group is.

(Refer Slide Time: 27:38)

So if I am not wrong I will find out group here good here and good is 44. So if I just see data .classifier 44 instead of 1 and 2. If I just try to see the 44<sup>th</sup> one that is the good things and you can see carefully that here when good is actually occurring yes then it is 8 % times this is a positive rating and 4 only 4.7 % case.

So there is a huge odd ratios difference 92 is to 8 divided 95 is to 4, so I can say that good is something which leads to classification. So we are doing it with the only the review content you can do the same stuff exactly same stuff with the review title.

So now I use this to predict my data here DTM dm testing data, testing document term matrix wee se for the prediction purpose and then I will call the library again the installed library installed packages. So see I am installing because today I am using a different system actually. If you have once installed in your system you do not have to install again. So library g-models and then I will do the cross tabulation.

So the cross tabulation is saying that I have my overall accuracy level is 213+33 which is around 246/282. So this many times I am overall correct 87 % of time based on this and if I try to focus on which is basically 75.5 + 117 this one. Now if I try to focus on that this is my actual this is my predicted. Now if I am trying to see that how out of actual ones how many times you have correct the predicted one which is basically 213/223 comes up to be 0.955.

So once positive rating I am properly identifying negative ratings if I just focus on then negative rating identification is not so good because there 59 observations which are negative basically, but I am currently predicting only 33 of them so which might not be good. So we

have to try to improve this classification better using further models. So now this is an example how we can do text mining.

Now I can use this particular text mining technique which is this Naive Bayes to create another data term matrix and can predict that whether that is a positive rating or negative rating. So this is how basically we create a prediction easy prediction mechanism to predict the sentiment of a particular text. In the next session, we will discuss about various kinds of libraries that we can use for testing purpose. Thank you very much I will see you in the next video.