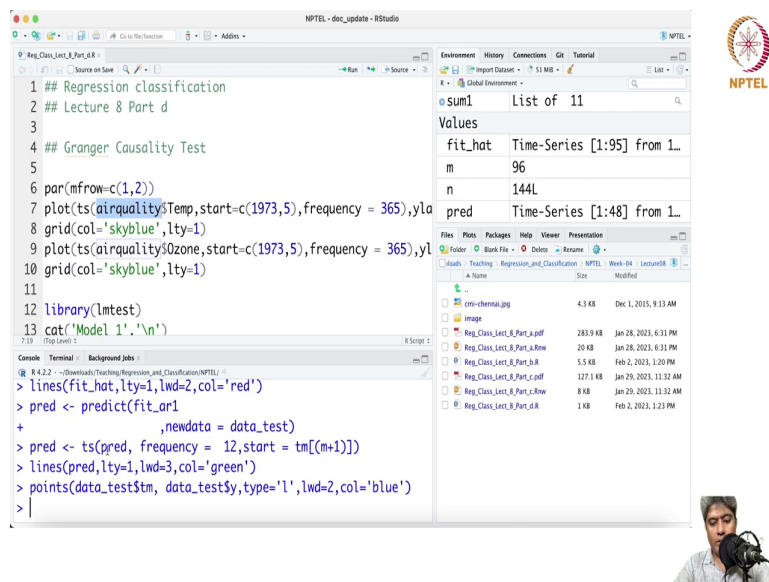**Predictive Analytics - Regression and Classification**
**Prof. Sourish Das**
**Department of Mathematics**
**Chennai Mathematical Institute**

**Lecture - 29**
**Hands on with R Part - 7**

Hello all. Welcome back to the last part of video lecture series of 8. Now, we are going to do the Hands-on for Granger causality test.

(Refer Slide Time: 00:29)



So, first we will consider this airquality dataset.

(Refer Slide Time: 00:45)



Say this test comes with the basic datasets package.

(Refer Slide Time: 00:51)



And in the air, it is if you look into the description of the dataset its New York Air Quality Measurement. So, Daily Air Quality Measurements in New York between May and September of 1973. So, data its a data frame it has 153 observation on 6 variables ozone, solar, wind temperature, month and day.

Now, we are only focusing on temperature and ozone, but you can try with the other variable as well. So, the first thing we are going to plot is temperature dataset.

(Refer Slide Time: 01:42)

(Refer Slide Time: 01:46)



That is how like temperature dataset has.

(Refer Slide Time: 01:52)



This is the ozone dataset as we have seen.

(Refer Slide Time: 01:55)



Looks like there are some missing values are also there in the ozone time series.

(Refer Slide Time: 02:00)



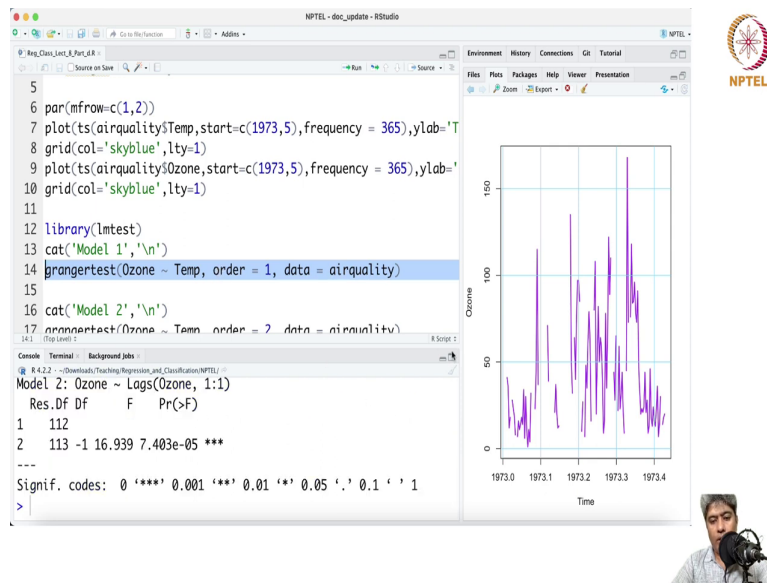So, the first thing we are running here is in from the LM test. The first model that we are trying to fit is the grangertest where Ozone is a function of Temperature and with only lag 1 and if we run this.

And data is over air quality we run this the P value this is the F test run. So, the Model 1 is Ozone with Lag 1 Ozone and Lag 1 Temperature and then the second model is only the null model and then it did a test whether the lag model has effect or not and then it says that ok the null model is not right.

(Refer Slide Time: 02:52)



So, that means, temperature does have a effect on the Ozone. Similarly, we run the second model and in this ozone again as a function of Temperature, but now we are going up to the second order of the granger causal test model.

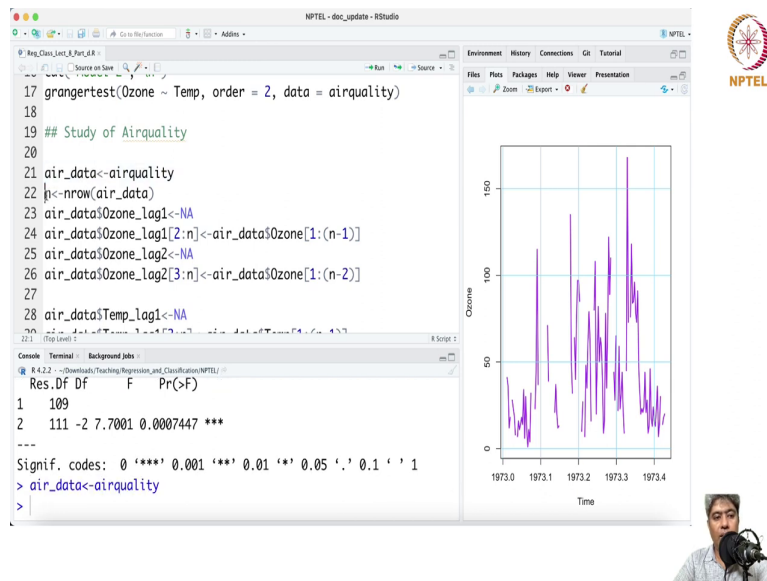(Refer Slide Time: 03:15)



Its essentially autoregressive model essentially its autoregressive model and P value is still small. So, we can say that temperature does have effect on ozone.

(Refer Slide Time: 03:31)



So, in order to understand a little bit more how this whole thing is working. So, we take the air quality data.

(Refer Slide Time: 03:43)



n is the number of samples. First what we and now if you look into the let us look into the air_data.

(Refer Slide Time: 03:53)

(Refer Slide Time: 03:56)



So, that is how the data looks like, ok.

(Refer Slide Time: 04:01)



That is how the Ozone Solar Solar.R radiation Wind Solar.R stands for radiation Wind Temperature Month and Day. Now, there are some NA observations are available.
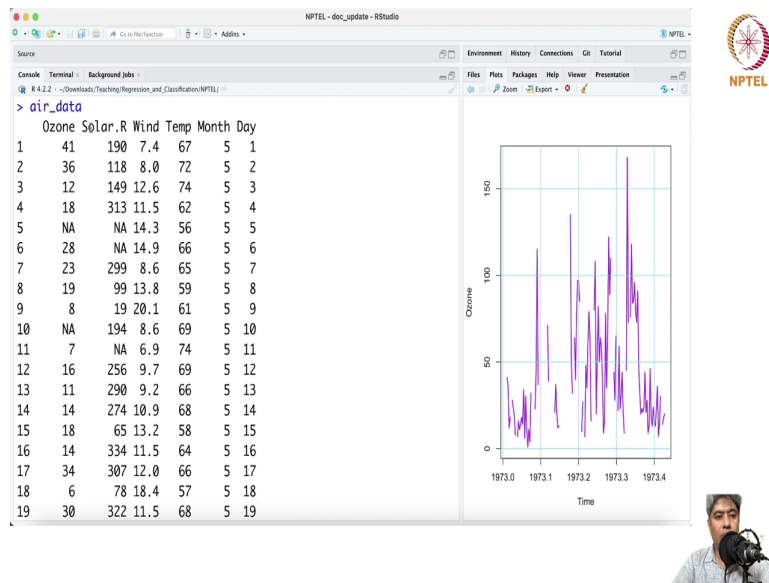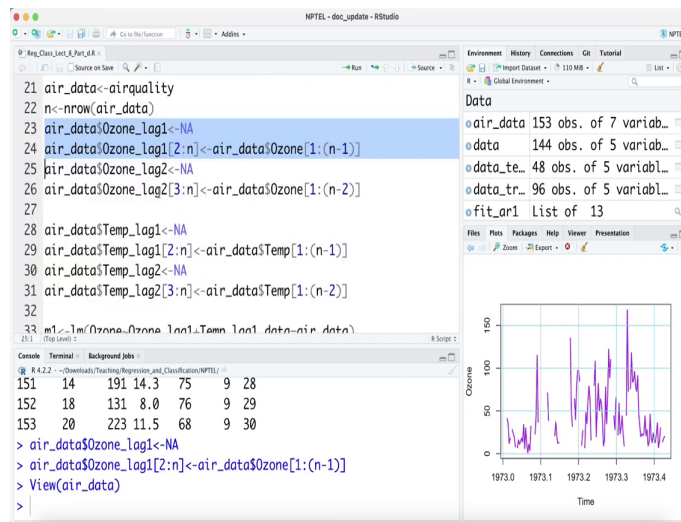
So, it is always bit difficult how you do this you know imputation, but for the time being we are not handling the missing data. We are only going to use the data which is fully available to us, but for the time being. So, we suppose we want to fit the first model we create the lag data set ok.

(Refer Slide Time: 04:47)



So, here we have created the lag. So, 41 was here we just created lag 36 just brought it down by 1.

(Refer Slide Time: 04:58)



Then we created the lag2 variable.

(Refer Slide Time: 05:00)



So, now you have see the lag2 variables have been created ok.

(Refer Slide Time: 05:07)



Then similarly we create the Temperature_lag1 and Temperature_lag2 variables.

(Refer Slide Time: 05:12)



So, Temperature_lag1, Temperature_lag2 variable have been created.

(Refer Slide Time: 05:19)



Now, if you fit the model with Ozone_lag1 and Temperature _ag1 ok.

(Refer Slide Time: 05:36)

(Refer Slide Time: 05:12)



Ozone_lag1 and Temperature_lag1 and similarly m1 say for null I will use 0, but just copy this entire thing and but instead of Temperature_lag1 I am dropping the Temperature_lag1. So, I am just saying the Ozone is only function of its own ok. Now, you see anova if you just run anova between m10 versus m1 ok.

You can see that it is the F test reject the null hypothesis and it says that the lag1 temperature does have effect on Ozone. So, you can do the Granger causal test in this way as well.

Similarly, this is m2 this test with the 2 lag and then I can define the model with null model here, but only with the lag Ozone not with the temperature I am dropping the temperature and then we run the anova m test or F test 2 to comma m20 ok.

(Refer Slide Time: 07:11)



And you can see you can reject the F test. So, lag2 have effect on the Ozone.

(Refer Slide Time: 07:21)



Now, if I compute the AIC of Model 1 and Model 2 clearly the Model 1 has a more AIC than the Model 2 the Model 2 has a lower mic AIC. So, we can we would prefer Model 2 over the Model 1.

(Refer Slide Time: 07:51)



Similarly, you can go for the third variable third lag we have to create a lag here. So, we can always create these lines.

(Refer Slide Time: 08:00)



First, we have to create these lines. So, maybe lag3, lag3, 2, 4, 3.

(Refer Slide Time: 08:26)

So, we have to create these lags for ozone and then now we have to create these lags for temperature third lag plus 3. And then what we need is model for third lag and Ozone plus 3 lag plus 3 and ok.

(Refer Slide Time: 09:10)



And the 3 naught null we just Ozone it will be only function of Ozone and then we will see if it is still effective.

(Refer Slide Time: 09:34)



And now when we are doing third lag we are adding third lag its not effective anymore.

(Refer Slide Time: 09:38)



But can we choose m3 what is the AIC of model third model.

(Refer Slide Time: 09:56)



Now, if you interestingly if you see what is happening. You see what is happening here the AIC is constantly dropping whereas, model; that means, third model is probably better than the first and second model.

(Refer Slide Time: 10:21)



However, what we are seeing that third model the lag does not have any effect. So, that means, we cannot really use we cannot say that third lag has any effect on the 3 lag model of a granger causal model does any effect any of temperature has no effect on the Ozone.

So, that means, as we are putting as we are you know put more and more lags. So, naturally what is happening the model complexity increases and its going to higher and higher dimension as model complexity increases it is doing some overfitting because my AIC is constantly going down. But we know, but most likely it is effectively losing its interpretability and that is why probably its not running its not giving any effect.

(Refer Slide Time: 11:21)

But we though we know that you know lag1 and lag2 does have a effect. So, we can see that you know you know very high M3 in case of third model except the lag1 model does not have any effect.

(Refer Slide Time: 11:49)

(Refer Slide Time: 11:53)



Whereas, m2 if you see m2 lag1 temperature does have a effect we see it does have a effect whereas, the lag2 does not have a effect.

(Refer Slide Time: 12:00)



Also, if you look into the standard error standard error in third model is 0.577 whereas, in the second model for lag 1 is 0.50. So, that means, standard error is increasing. So, there is a high possibility that a multicollinearity also creeping in because of the multicollinearity these lags are not any more effective because their standard error is going up because of the multicollinearity.

So, you have to be very careful about how you do the interpretation of these statistical inference and these you have to be very careful about these statistical inferences when particularly when you are adding more and more features just because every time you are increasing lag.

That means, your model is getting complex your complexity of your model is increasing it will have a higher tendency to overfit because end of the day your training datasets is finite

most of the time your data training dataset is not increasing. So, as a result your model will overfit and as it overfits you do not want your model to overfit because your bias will when it happens your bias will be small.

But you will have a very high variance in out of the sample it will not do very well. So, you have to be very careful about model fitting when for adding more complex lag variable and it is better to obtain a parsimonious small model and just stop there, ok. So, far this week we will this is that is how this much we will be discussing. Next week we will see you with a new video with a new chapter.

Thank you very much. Take care. Bye.