**Exploring Survey Data on Health Care**
**Prof. Pratap C. Mohanty**
**Department of Humanities and Social Sciences**
**Indian Institute of Technology, Roorkee**

**Lecture - 32**
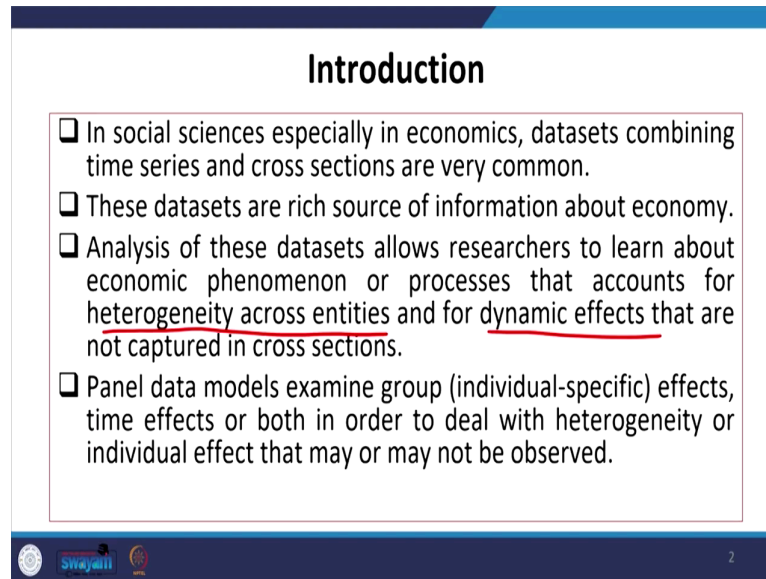**Common Constant Model in Healthcare**

Welcome friends once again to my NPTEL MOOC module on Exploring Healthcare Survey Data. We are on the verge of 7th week of the module; we have already started explaining panel data in healthcare. On this particular lecture we are directed to deal with the different models of panel data and healthcare.

Here we have thought of dealing with one of the very starting points of panel model usually discussed as or called as common constant model. Otherwise called pooled model in panel or in healthcare. So, let us understand to what extent pooled or common constant in a model are different than that of the proper panel model.

This lecture week we are going to give largely on clarifying the concept and with some practical handouts, practical experiments. We will also nurture you with the exact interpretations of these all models and these are highly useful in research specially those are advanced users.

Advanced users in the sense those who are writing their papers and submitting to a good journal usually A category journals. For them I will always suggest to follow this very lecture carefully and then accomplish your paper accordingly. So, here are the background information about this model.

(Refer Slide Time: 02:02)



In social sciences, especially in economics data sets, combining time series and cross sections are usually very prevalent. These data sets are rich source of information regarding the activities of an economy. To analyze these databases, we allow researchers to learn about how economic dimensions or structure accounts for heterogeneity across entities or dynamic effects that are not captured in cross sections.

So, two important points are to be highlight here, one is called heterogeneity across entities. So, it's not just the cross-sectional heterogeneity, there are heterogeneity due to other dimensions as well. In addition to that we are also adding information about dynamic effects which are not captured in cross sections data.

The panel data models examine group effects, time effects or both in order to deal with heterogeneity or individual effect that may or may not be observed. So, like it captures individual specific effects as well as group effects, both the cases are going to be discussed through this lecture.
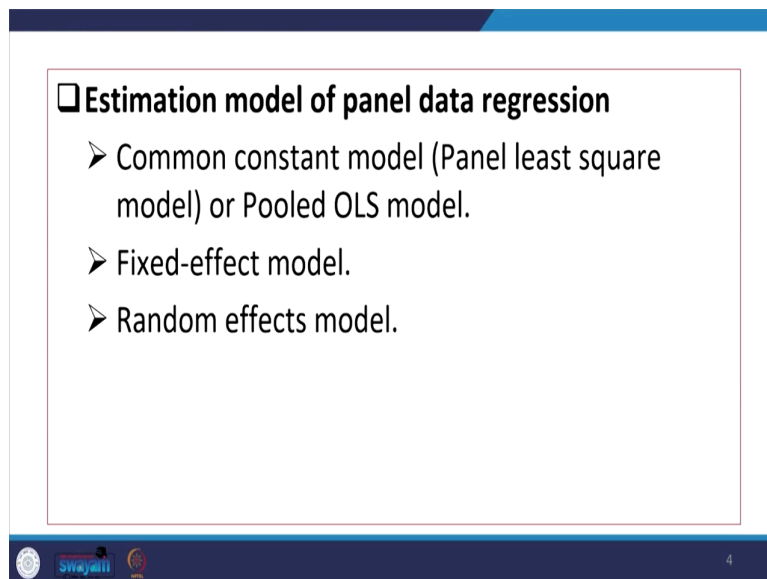
The panel data methods can be divided into two broad categories one is called homogeneous panel data; second one is called heterogeneous panel data. Homogeneous, as the word clarifies this assumes that the model parameters are actually homogeneous. In the sense they are common across individuals like we have seen in ordinarily square method or ordinary regression methods or in a general cross-sectional data.

These are also called pooled models. So, homogeneous panel data are discussed in the context of pooled models as well. Whereas, in the case of heterogeneous panel data models, this allows the model to have different parameters across individuals.

The individuals are actually varying with their predicted values when there are some forms of dynamism, some forms of differences across individuals may be due to the differences across time or across the cross sections.

So, when both are understood they are actually captured in the heterogeneous panel model. Another aspect is called fixed effect and random effect models within the heterogeneous panel model data. These both models actually are example of panel data models, heterogeneous panel data models.
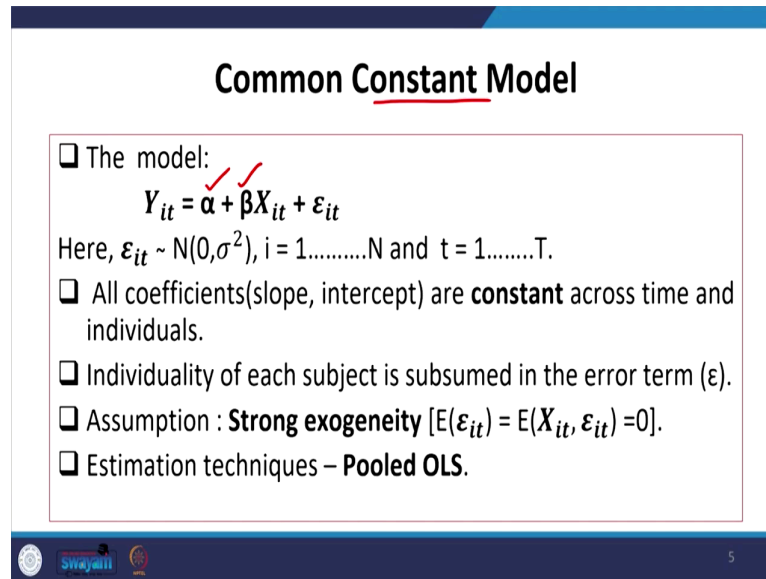
(Refer Slide Time: 05:11)



To estimate the panel data regression, we consider three important models. One is called Common constant model also called Panel least square model. These are also called Pooled OLS regression model. Then another one is called Fixed effect model and the third one is called Random effect model.

## Common Constant Model

❑ The model:
$$Y_{it} = \alpha + \beta X_{it} + \varepsilon_{it}$$
Here, $\varepsilon_{it} \sim N(0,\sigma^2)$, i = 1..........N and  t = 1........T.

❑ All coefficients(slope, intercept) are **constant** across time and individuals.

❑ Individuality of each subject is subsumed in the error term (ε).

❑ Assumption : **Strong exogeneity** [$E(\varepsilon_{it})$ = $E(X_{it}, \varepsilon_{it})$ =0].

❑ Estimation techniques – **Pooled OLS**.

Let us start with the common constant model which we could have written as pooled. You might have further raised certain confusion to do that we use the term common constant model as per the book and as per the standard definition or explanation.

So, the word common constant means your coefficient of estimation is going to be constant throughout the model for the particular estimation. Now, you can see which are the common estimate that is the constant estimator and this is the coefficient or slope estimator. Both are actually going to be constant slope as well as intercept are constant across time as well as individuals in such cases it is called common constant model.

Now assumption similar to that of the OLS model is that the error term though that varies across cross section and time. But in such type of data, we simply assume that it is normally distributed and then with standard normal distribution with it is mean as 0 and standard variance is sigma square.

So, here i varies from 1 to N across cross-sections and t varies from 1 to capital T time points. So, the individuality of each subject is; of each subject is subsumed in the error term. So, any sort of individual differences those are actually captured in error term. Therefore, the assumption in this case is that there must be strong heterogeneity. Heterogeneity between the explanatory variable and the error term. So, expected value of the explanatory variable and the error term is equal to 0.
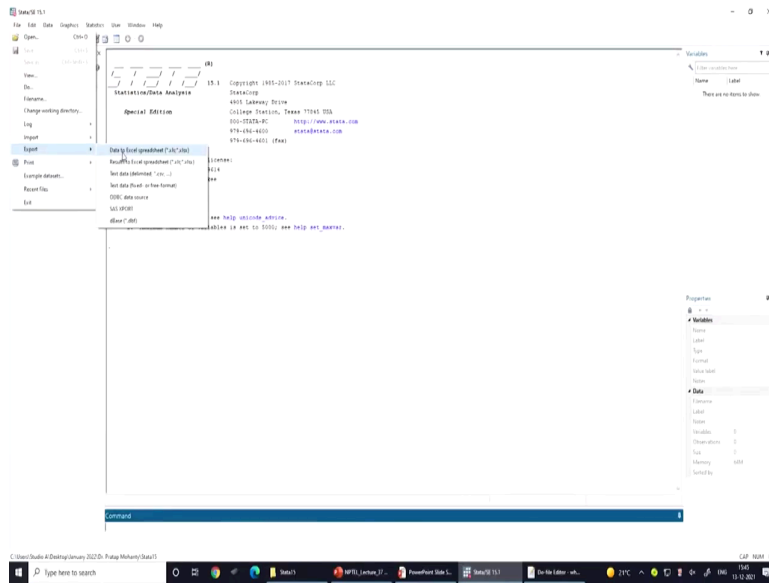
So, the estimation technique we consider is called pooled OLS model because simply each item by varying by time or by varying by cross section. These are simply pooled together and defined with a common model and that is why they are called pooled.

Now, you might have questions on this that how come they are even if pooled there are differences. Our assumption is forcefully making it with no endogeneity, it has strong exogeneity. Therefore, we go for OLS. Even if we do it, we cannot just convince ourselves that it is going to have 0 exogeneity or 0 endogeneity. It has strong exogeneity there are some practical examples we wanted to show it to you and we are going to experiment with it.

Like here we are giving a pool kind of data. We have also cited the data from WHO panel data sets those are available. And we will also make you available with this practical dataset in your portal and all of you could able to find the data for your own experimentation. Along with that we are going to also show UN and run regression results with IHDS panel data. I just want to mention you just give you a note again that if you have still some confusions about forming a panel data through IHDS you need to refer to my earlier module i.e., handling large scale dataset using STATA that is still running at this moment. And you may also follow from YouTube that how panel data could be generated. IHDS was initially not providing the panel data to the public. But nowadays the structured panel data is available for us to work with.

So, both the way we will be emphasising and clarifying your doubt. So, IHDS is the only longitudinal survey large scale survey available in India. There are in fact, two waves of this survey IHDS 1 and IHDS 2. So, let us go for it and explain a bit about it.
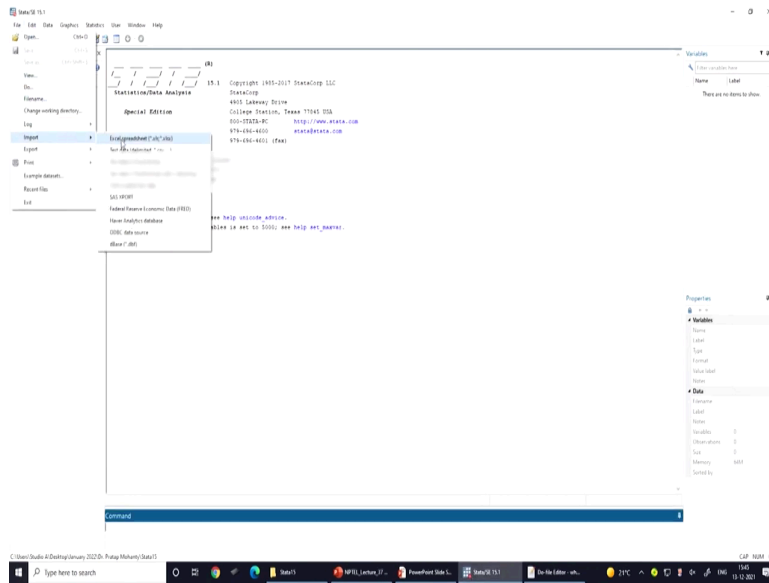
(Refer Slide Time: 09:56)


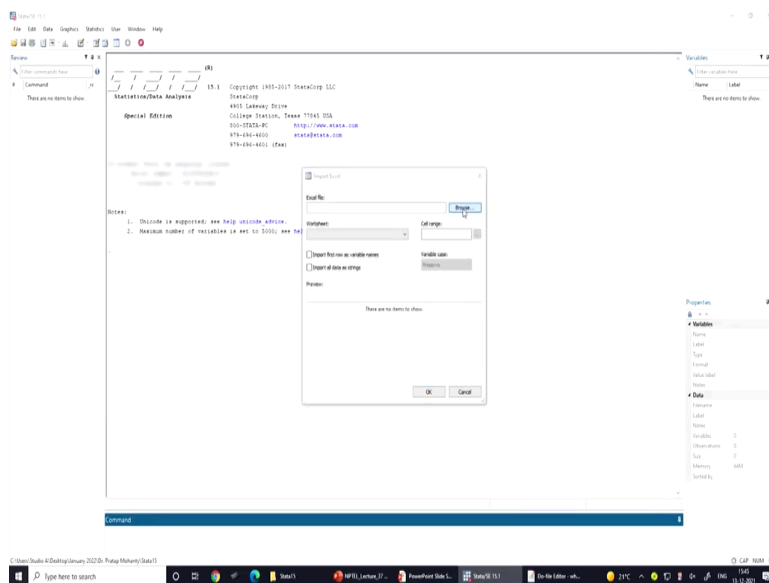
(Refer Slide Time: 09:58)



On the start of window, we are deriving results. I am going to open the first on the WHO dataset.
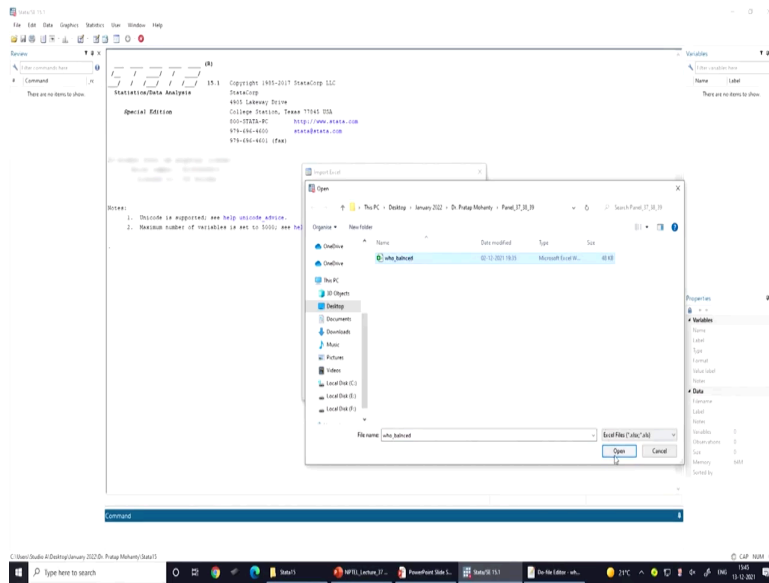
(Refer Slide Time: 10:02)
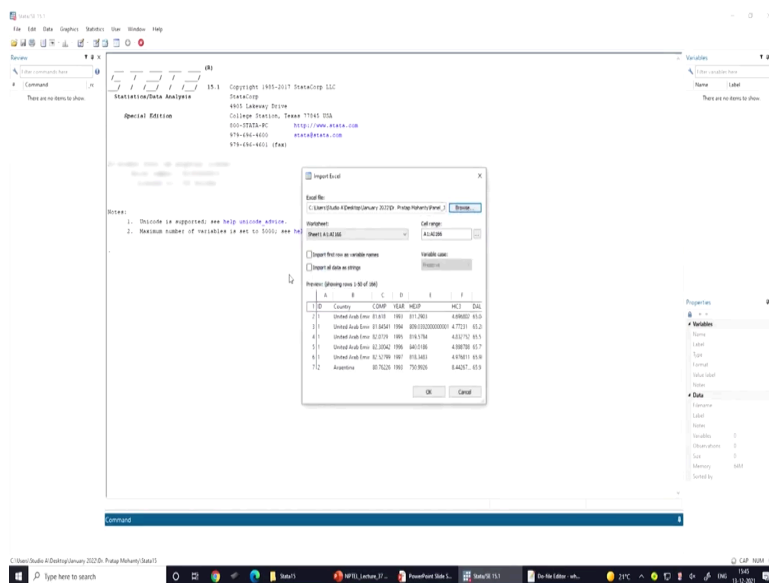


(Refer Slide Time: 10:05)



Now so I will open the file and import the data on your screen.
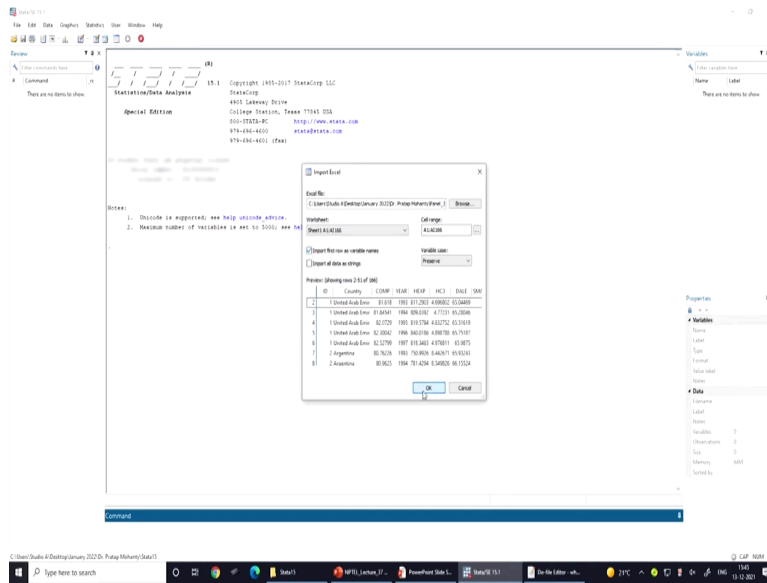
(Refer Slide Time: 10:07)



(Refer Slide Time: 10:11)



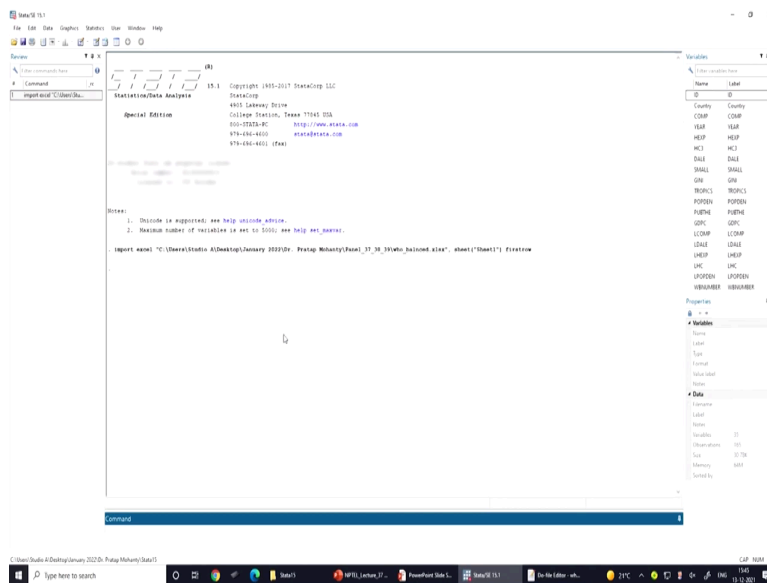And this is the data, we are going to show it to you or going to provide you.

(Refer Slide Time: 10:13)



And now we can define that which are your row variable.

(Refer Slide Time: 10:17)



And the row variable is basically that identifies your variable names correctly. Now we will open and share with you the constant common model analysis.

(Refer Slide Time: 10:35)



And in this WHO dataset we will be running some of the commands we have given here for you. And in order to cut short our time and to explain the best in the lecture we have kept a do file for your reference as well. Here the first one we have already imported and now we are doing some pre-estimation. And first pre-estimation we are doing with this data is describing the data with three important variable household expenditure is the first one that is household expenditure on health.

(Refer Slide Time: 11:09)

(Refer Slide Time: 11:17)



So, H stands for health expenditure. So, that is here and describe is going to give you the nature of the data on your screen. Then next will be clarifying you about the extricate we need to set whether which kind of panel it is what is the panel Id's.

So, Id and the year are the two variables defined as the panel units of panel variables for us, here t is the time variable, Id is the cross sectional variable. So, now, we have defined "xt set" is mandatory for you to run at first to recognize the STATA to go with the panel data estimation.

(Refer Slide Time: 11:52)

Now, this once you set it. This suggest that yes this is a strongly balanced data. So, again some of those clarification you can get it from our previous lecture as well as you can refer to my previous module. Now the data is from 93, 1993 to 1997 and the change by time is of one unit.

So, one unit time effect change is given. Now we are going to run the third important pre-estimation technique that is the summary of the data of these three important variables that is x t sum. We are going to operate xt sum here.

(Refer Slide Time: 12:42)



Now, the result is on your screen. The xt sum with these three variables that summary is giving information. So, we are getting in this panel one of the important aspect is that we are going to get the overall model within effect and between effects as well. So, since this is a pool data, it is not capturing the between and within. It is simply considering the overall estimation. Whereas, in case of fixed effect and random effect you can easily differentiate these categories.

Now, but it has given its standard deviation in each category. By mean, it has considered the overall mean, but by standard deviation there would be certainly some differences in the minimum value, maximum value, and the number of time periods are 5 in each case since it is a strongly balanced panel.

So, these are all the description for you. Now we are going to run with the proper CCM model on your screen. So, CC common constant model is a kind of pool data as I told you. So, the simple regression command is going to give you the result.

(Refer Slide Time: 13:50)



So, on your screen we have derived the reg where the health expenditure data is the dependent variable. Now you can easily compare this result with the ordinarily square data. This is in fact, the ordinarily square model and the coefficient we can interpret easily. The P values can be interpret by the way we usually do it. So, there is no major difference as compared to the OLS model.
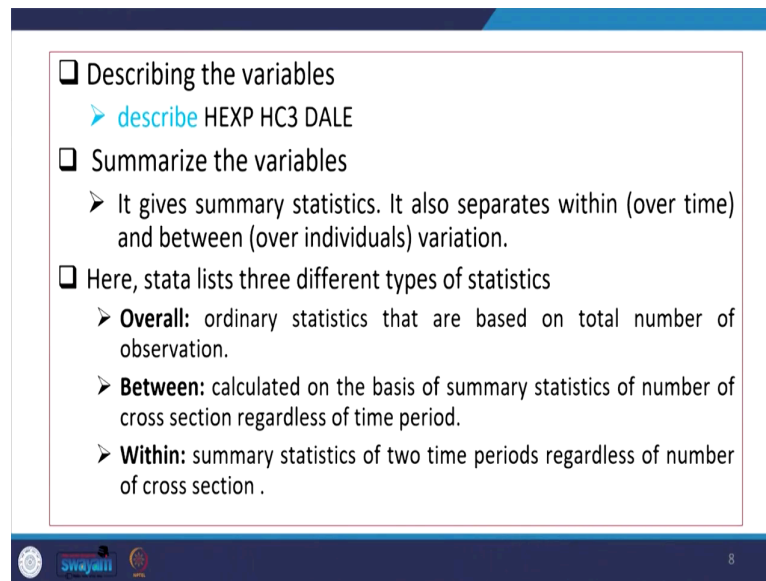
(Refer Slide Time: 14:28)



Now, let us move to our PPT once again. I will also show it with explanation in writing. So, we are explaining the WHO panel data just to observe the effect of education attainment and disability adjusted life expectancy on per capita health expenditure. So, the health expenditure variable was taken as the dependent variable. So, dependent variable is per capita health expenditure.

Whereas, the independent variable as are like educational attainment and dale disability adjusted life expectancy. The purpose of this lecture is to show how to do panel analysis with health data, it does not cover the whole research process and model.

So, we will steadily capture other directions of panel data, but at this moment we are just trying to give you the background information about a common constant model which is essentially required for discussion in order to start with fixed effect and random effect model.

Now, just to clarify I have already discussed about describe command and sum/summarize command that effects in summarizing the statistics. So, basically in summary it also separates within or over time and between variations.

(Refer Slide Time: 15:49)



So, the standard deviation value within and within basically over time. Because within the cross sections, how the same cross section unit is varying over time that is why it is called within. Whereas between basically between the individuals. So, how individuals are actually and their responses are varying across individuals has nothing to do with the time factor.

So, between is not taking comparison of the time whereas, within is actually emphasising the time aspects. So, these are all three types of statistics- overall, then between, and then within. Between as I told you calculated on the basis of some statistics of number of cross section regardless of the time period. And within is basically considering the time factor and casing the time changes.

(Refer Slide Time: 16:44)



So, the xtsum command we have clarified already derived on your screen and you can check that in between and within effect, the variations are observed. So, some basic ideas you can take like across the individual units since it is there are so many observations 165 observations.

The between variations are more than that of the within variations. So, over the time changes are not much. So, far as per capita expenditure is concerned. Whereas, education and disabilities aspect is constant you can also check the differences accordingly.

(Refer Slide Time: 17:28)

Now, the CCM model we have already derived this treats a dataset like any other cross sectional data. This ignores the data that has a time and individual dimensions that is why the assumption that are similar to that of the ordinary least square model.

(Refer Slide Time: 17:46)



So, these are the results we explained and we have clarified how the coefficient could be interpreted and the P value could be compared. So, further to this you may also refer to my previous lecture as well previous lecture on handling large scale data using STATA. So, this is a similar to that of the explanation made in OLS regression that is why we carry with the command reg alright, reg dependent variable and in control variables. If any are there and accordingly you can run the regression and interpret it. There are no more things to add. Rest of the things we will explore and find out.

If you have any difficulties, we will be happy to address it. These are all for this lecture. In the next lecture we will be clarifying on fixed effect model that is going to be very interesting. And I will explain you the equation very carefully and I will also discuss about i term, t term and it will be more fascinating for you to attend.

Thank you, let me stop here.