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

**Lecture - 35**
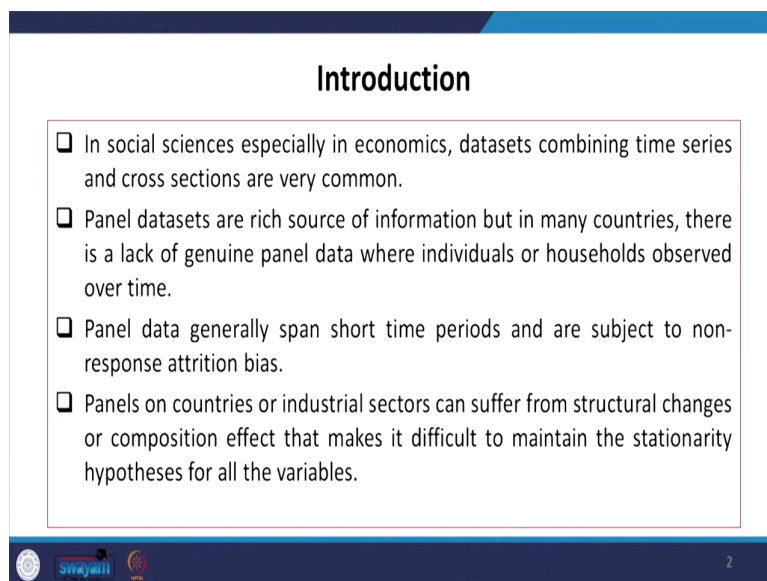**Construction of Pseudo Panel**

Welcome participants, once again to my NPTEL MOOC module on panel survey data on healthcare. We are on the 7th week of understanding panel survey data and the title of the NPTEL MOOC program of my course is exploring healthcare survey data.

The title of this particular week is panel survey data and here we are explaining how to have pseudo panel data. This is quite rare in the existing studies and very less number of studies which could have guided you about the pseudo panel and its construction.

Therefore, we kept this as a specific lecture. In giving you a certain direction about understanding the construction of pseudo panel data. So many large scales cross-sectional data sets could be actually presented with a panel structure.
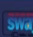
In social sciences especially in economics, datasets combining time-series and cross-sections are very common. Panel data sets are rich source of information, but in many countries, there is a lack of genuine panel data where individuals or households observe over time.
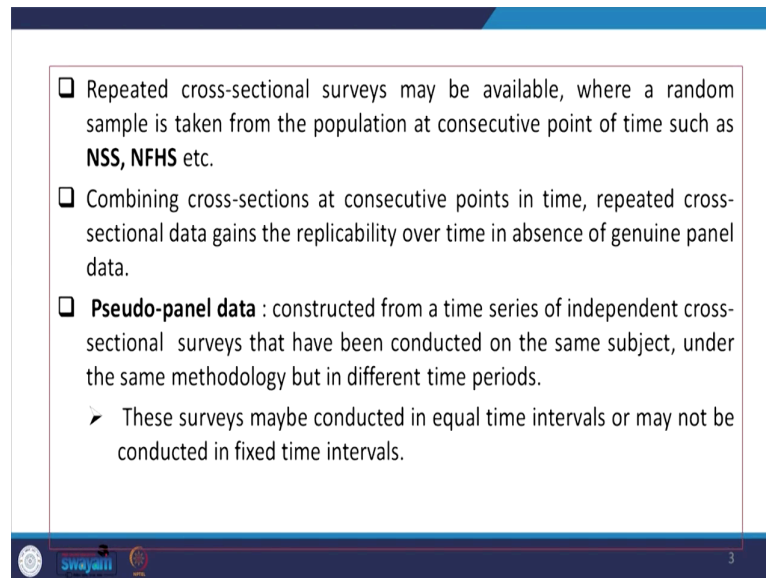
(Refer Slide Time: 01:56)



## Introduction

❑ In social sciences especially in economics, datasets combining time series and cross sections are very common.

❑ Panel datasets are rich source of information but in many countries, there is a lack of genuine panel data where individuals or households observed over time.

❑ Panel data generally span short time periods and are subject to non-response attrition bias.

❑ Panels on countries or industrial sectors can suffer from structural changes or composition effect that makes it difficult to maintain the stationarity hypotheses for all the variables.

Panel data generally span a short period of time and are subject to non-response attrition bias.

Panels on countries or industrial sectors can suffer from structural changes or composition effect that makes it very difficult to maintain the stationarity hypotheses for all the variables that are carried forward in the next round.
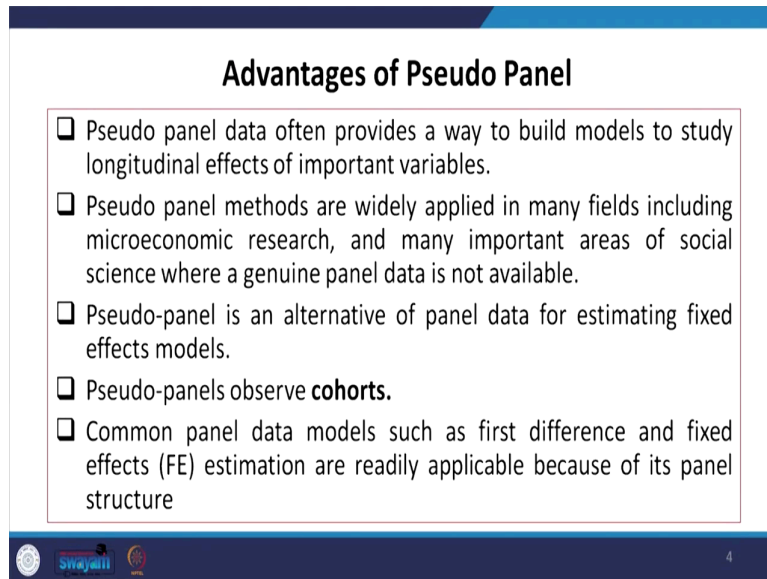
(Refer Slide Time: 02:19)



Repeated cross-sectional surveys may be available, where a random sample is taken from the population at a consecutive point of time such as NSS or NFHS etc.

Combining cross-sections at consecutive points in time repeated, cross-sectional data gains replicability over time in the absence of genuine panel data. Even if there is general panel data, repeated studies are taken in different time periods. Therefore, a pseudo panel can be constructed.

So, the pseudo panel is something where a time series of independent cross-sectional surveys that have been conducted on the same subject under the same methodology, but in different time periods could be conceptualized in a dataset called pseudo panel.

These surveys may be conducted in equal time intervals or may not be conducted in the fixed time interval but a time component is required.
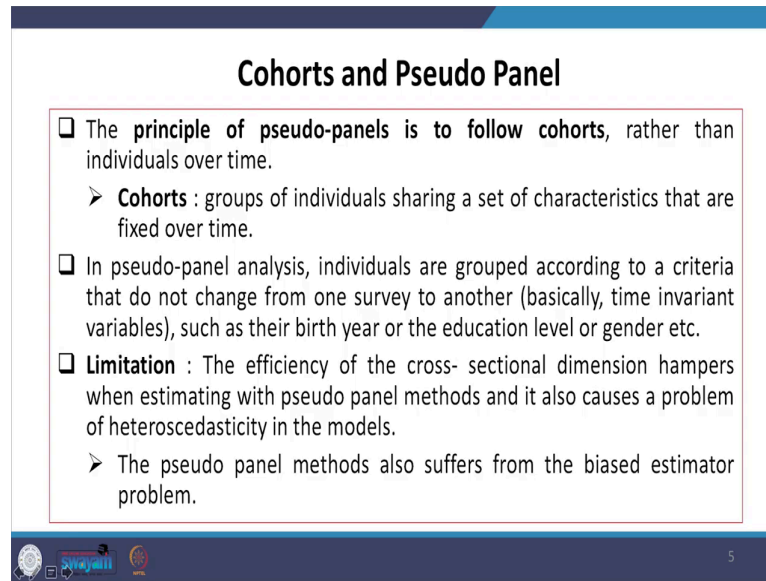
So, the advantage of the pseudo panel is that pseudo panel data often provides a way to build models to study the longitudinal effects of important variables. So, longitudinal components of important variables are continued in this setup.

Pseudo panel methods are widely applied in many fields including microeconomic research and many important areas of social sciences where genuine panel data is not actually available. The pseudo panel is an alternative to panel data for estimating fixed-effect models.

Pseudo panels observe cohorts instead of individual observations. Common panel data models such as first difference and fixed effect estimations are readily applicable because of their panel structure. So, accordingly, methods will be employed.

Now, we are clarifying the cohorts and pseudo panel. The principle of pseudo panels is to follow the cohorts instead of the individual as the individual persons over time. So, the cohorts which I have just emphasized are going to be defined. These are something defined in groups of individuals sharing a set of characteristics that are actually fixed over time.

In pseudo panel analysis, individuals are grouped according to a criterion that do not change from one survey to another basically time-invariant variable are considered to be the cohorts and these are like birth year or education level or gender etc. Usually, there is not much change over time.

So, the limitations of this data are that the efficiency of cross-sectional dimension hampers when estimating with pseudo panel methods and it also occurs a problem of heteroscedasticity in the models. So, there might be some heteroscedasticity problems since and accordingly the efficiency of the model is compromised to some extent. The pseudo panel models also suffer from a biased estimator problem because of these issues.

In the pseudo panel instead of the individual component, we have a cohort component. Cohort we have taken c stands for cohort and t is the time dimension for each of the variables.

So, we are basically looking at the R and Rct. So, c stands for cohort and t stands for time variation and accordingly we will estimate its R-value based on these two components.
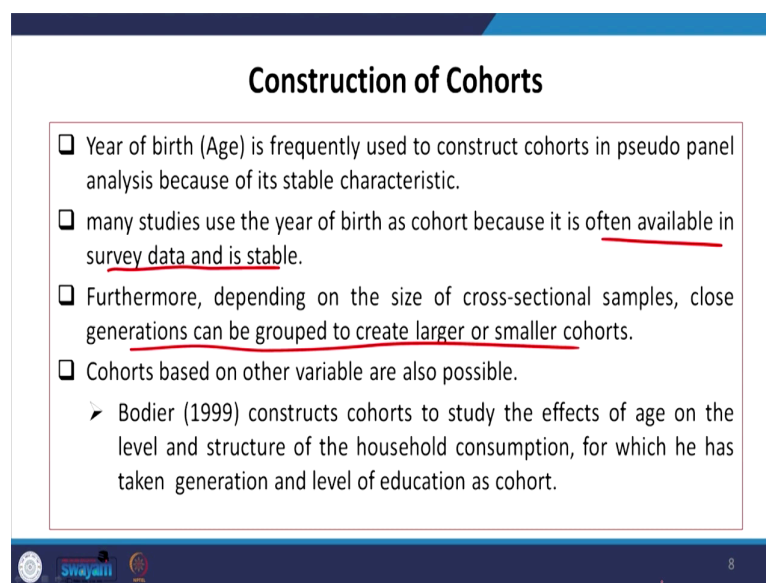
We are now guiding the construction of cohorts. The first step is to construct a cohort-based on the selection of the variables. Those variables' values usually do not have many changes over time and that is and cohorts are constructed on a variable which is actually time-invariant.

Each individual in the sample must be placed with exactly one cohort, without that it is difficult to define. It must follow the assumption that the cohort term $\alpha$ is fixed over time. So, $\alpha$ is fixed when the true cohorts contain the same individuals at each point time period.
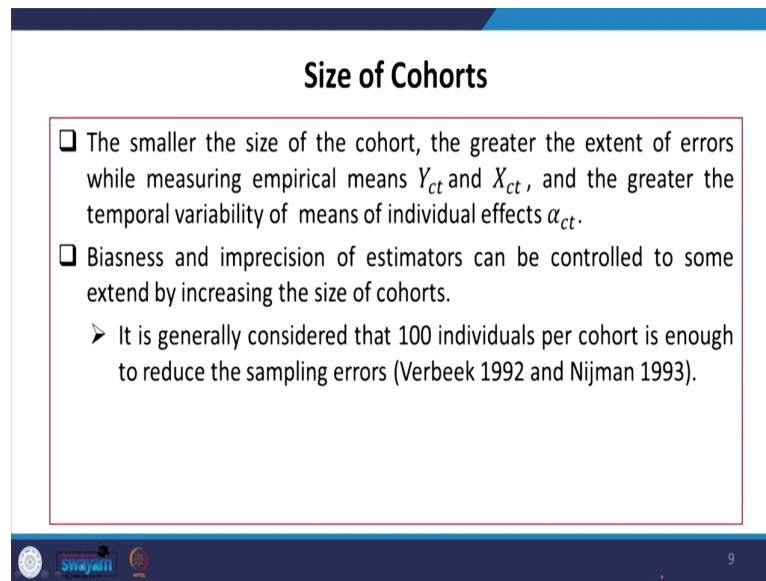
(Refer Slide Time: 07:48)



In the construction of the cohort year of birth is frequently used to construct cohorts in pseudo panel analysis because of its stable characteristics. Many studies use the year of birth as the cohort because of it is often available in survey data and is usually stable.

Further, depending on the size of the cross-sectional samples, close generations can be grouped to create larger or smaller cohorts. Based on the size we can be able to define bigger cohorts or smaller cohorts.

Cohorts based on other variables are also possible like Bodier 1999 constructs cohorts to study the effects of age on the level and structure of the household consumption, for which he has taken generation and level of education as a cohort.

(Refer Slide Time: 08:59)



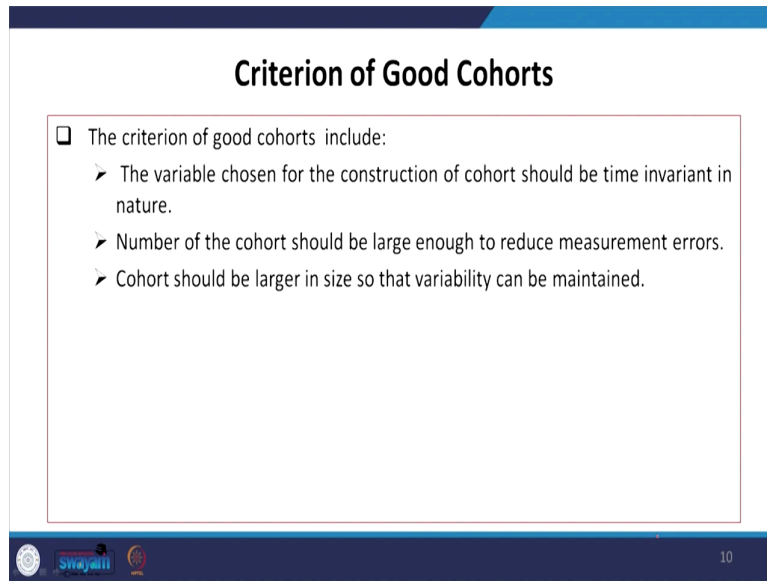So, what about the size of cohorts, what should be the ideal size some suggestions will count here. The smaller the size of the cohort the greater the extent of errors while measuring empirical means such as Yct and Xct and the greater the temporal variability of means of the individual variability of the means of the individual effects.

So, there will be a greater temporal variability if the size of the cohort is actually smaller. Biasness and imprecision of estimators can be controlled to some extent by increasing the size of the cohort.

It is generally considered that 100 individuals per cohort are enough to reduce the sampling errors as per the suggestions given by these 2 authors.
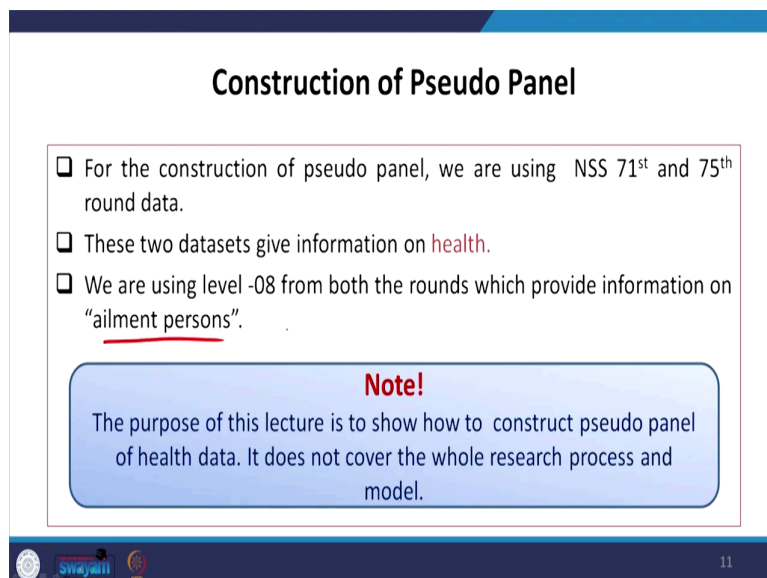
(Refer Slide Time: 10:13)



The criteria of good cohorts include the variable chosen for the construction of the cohort should be time-invariant in nature. The number of cohorts should be large enough to reduce the measurement errors which we have just guided and the cohort should be larger in size so that the variability can be also maintained in the data.
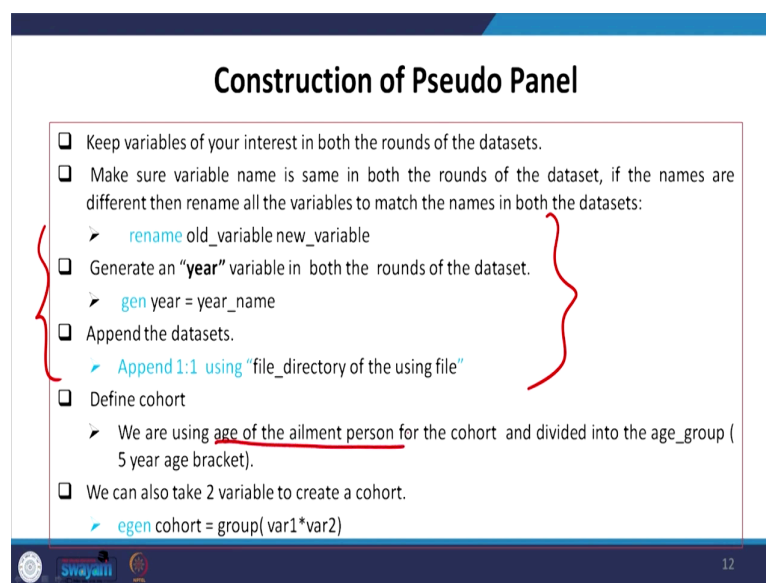
(Refer Slide Time: 10:51)



Now we are guiding about the construction of the pseudo panel. For the construction of the pseudo panel, we are using two rounds of data that is NSS 71st round and the NSS 75th round on healthcare.

These two data sets give information on health. We are using level 8 from both rounds, which provides information about the ailment of the person. The purpose of this lecture is to show how to construct a pseudo panel of healthcare data. It does not cover the whole research process and model.

We are not interpreting everything or not to guide about to complete analysis. The construction of this panel will now guide you, some of the processes we are not operating these are very common.

(Refer Slide Time: 11:44)



We need to keep variables of your interest that are going to be used in the panel data of both rounds. We need to make sure that the variable name is same in both rounds of the data sets. So, that panel can be made otherwise some errors might be generated.

If the names are different, then rename all the variables to match the names in both datasets which is essential. So, we can rename by the command "rename old variable new variable" and we need to generate the year variable in both the rounds of the dataset. So, like year 1 or year 2 or year that particular year name we can write it down and the year variable is also generated.

Then we need to append these two datasets with 1 is to 1 like "append using file_directory". So, one is to one append without director name you can give it.

Now, we are guiding very particularly on cohort generation. So, to define cohort we are using the age of the ailment persons for cohort and dividing it into the age group. Once we define that then we are actually dividing them or classifying them into age groups with the 5-year age bracket.

So, all of them with certain different 5 years' age brackets we can do that, then we can also take 2 variables as well. Instead of just 1 variable as a cohort, you can also generate 2 variables simultaneously. In that case, the command is "egen".
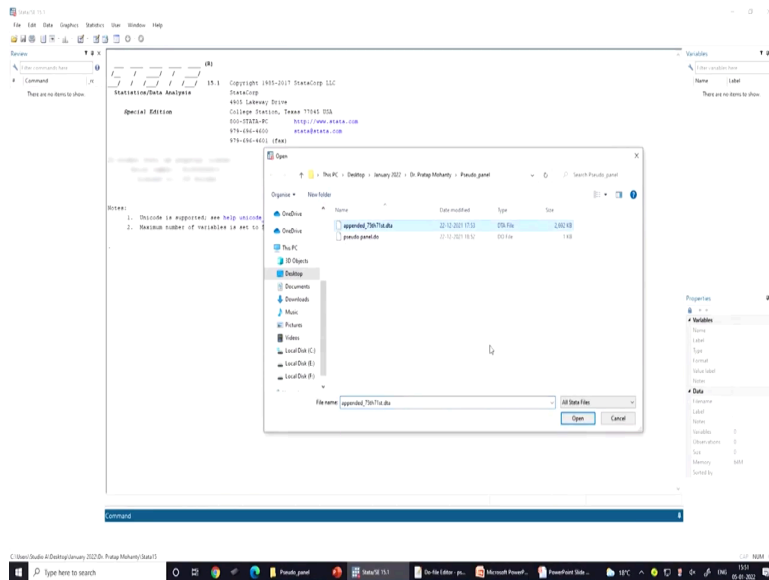
(Refer Slide Time: 13:40)



So, the example dataset we are using as I already mentioned is NSS 71 and 75th round. Now, I am just giving you the direction at this moment then we will use it directly from there. First, we drop the unnecessary variable then we will sort the age variable.

Therefore, the bysort command is used. Then bysort year with age group year and also we have taken the mean value of the variable that is nature of ailment and then we have defined these steps in front of you level of care then level of nature of the treatment, then nature of ailment three important aspects we have also expressed in by their mean values.
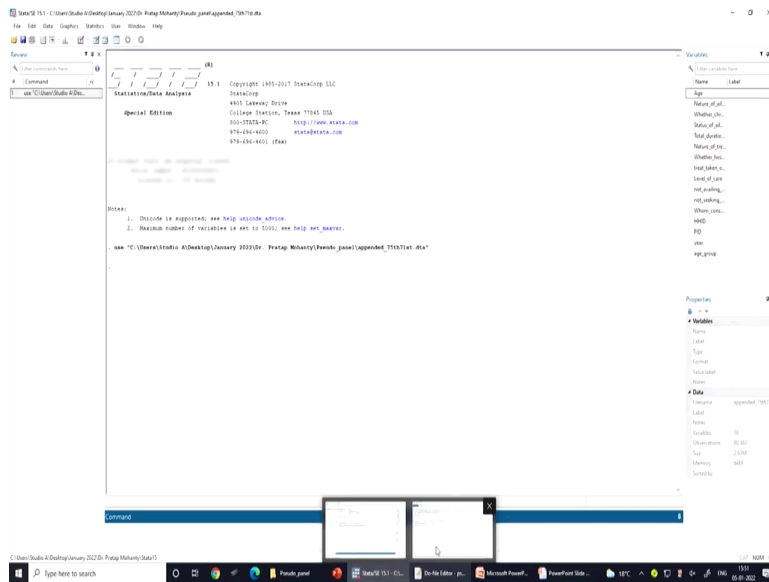
We have sort the data by the group year of that age variable, then we will summarize that and we will generate the cohort by the age group and then we will move to understand or define the data as a cohort as the id variable.

Then accordingly we will run the regression with a fixed-effect model because it is time invariance in nature.

(Refer Slide Time: 15:23)



(Refer Slide Time: 15:31)

(Refer Slide Time: 15:32)



Now, these are the steps which I am just going to show on the screen. We will open the do file on your screen and we can show that these are all important and how you can use them for your work directly.

So, we can drop the variables which are not relevant for our analysis, you can also drop on your own way. So, we have dropped it, and now the number of variables have been dropped. Now we will sort it by age group, years and mean value generated for each of the variables.

(Refer Slide Time: 16:02)

Now, three you can see the nature of the treatment, nature of ailment and level of care alright. Now next is our summarizing our data.

(Refer Slide Time: 16:18)



In the summary, you can see each of the observations and their mean value and their standard deviation and everything is presented on your screen.

The variation is not much higher here, because the mean value has been taken in the case. The next one is to group the variable that is the cohort. We are generating an age group and it has it is confined to the age variable in the grouping of the variable and that is considered as the panel component.

We will define the command here with the xtset cohort as the panel variable and this is how the stata has read that this is the variable panel variable and this is an unbalance data. Now we can run the regression with its fixed effect model because we know that we have taken a time-invariant model.

(Refer Slide Time: 17:35)



 The fixed-effect model has been drawn and you can now see that the p-value is significant and the rest of the interpretation we usually do is also applied here as well. This is how it works and as we already said that we are not going to interpret everything as like a paper.

We have generated or constructed a pseudo panel with the same technique you can apply with so many other databases and once your panel data is generated you can also work for the number of methodologies like policy evaluation techniques in your paper and those are going to be highly appreciated by the review because of its policy implications.

So, these are all for today and we have covered the pseudo panel for you. If there are any sort of queries in this regard, please do not hesitate and try to come to us. We will try to deal with our team and will be most happy to address it.

Thank you.