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

**Lecture - 33**
**Fixed Effect Model in Healthcare**

Welcome friends, once again to my NPTEL MOOC module on Exploring Healthcare Data, Healthcare Survey Data. As part of the Minister of Educations Initiative of the NPTEL program where so many unique courses are being floated, I have tried my best to contribute to this gap in different courses given the challenges of physical education.

This time we have considered the model on healthcare, the course outline is on healthcare and understanding healthcare data. These are called healthcare data analytics as well. Here we are dealing with 7th-week lectures. In the 7th week, we have already taken two lectures where we clarified what are called panel data, how to understand or read panel data, and how to know to balance and unbalance panel data.

In the previous lecture, we talked about the common constant model or pool panel model. So, even if it has panel content that is regarded as a cross-section type. In this lecture, we are adding feathers to the existing information on the panel that is a dimension called the effect, fixed effect model.

I am just clarifying the overview of it fixed. The word fixed effect means in the multidimensional set of with time and cross-sectional variation as part of the panel structure of the data. Some variables by structure are defined to be changing with a fixed percentage, with a fixed component, fixed coefficient.

The fixed coefficient may not be captured in the slope coefficient. Slope the change is not exactly captured, but the constant term could be actually identified as the change. Like when you say your consumption pattern every year, it might be hardly changing with different variables.

if your location is going to get changed, I am just giving the various type of examples. Location is changing, your attire is changing, your you are putting on a new dress, for example, your consumption is not going to change much. Even if your income is also

changing the extent of consumption is changed by a certain degree due to a certain dimension of the product,

Some of the yes consumption is a function of income and income change has reflection on consumption for sure. But there is some parallel example in this case though you are a bit confused with this. Some parallel examples like a person's, person height over the time due to a factor of various, as a function of various factors, or consumption as a person's caste for example caste over a period of time.

Caste, religion, etcetera within a family is not going to be changed over time or over a period of time, it is not getting changed much. But when you just try to capture the time effect probably you are not going to identify the exact aspect of the change. But that does not mean it is not changing with a structure, the alpha that is the constant term is somewhere reflecting the standard of living of the person,

No, it is not; though exactly capturing the variation in terms of time, but there might have been a paradigm shift that is captured in terms of the constant term that is an alpha. So, that is where the panel data stands. And the panel data clarify through the fixed effect on the constant changes, on some unit of v top in a change in the constant term, not by the slope coefficient,

That is why we are dealing with this particular lecture and emphasizing on the fixed effect of the panel content. And this is going to be very interesting in your everyday work, especially those who are targeting for the topmost journal, topmost paper, panel paper, and panel data used to be very important to ice icing your model, to give a better direction in your model, to have a better interpretation. So, far as policy recommendation is concerned. That is one direction I said.

Then the second aspect is whoever is doing some kind of an analysis using the large scale data; large scales cross-sectional data. But it has a time component as well. It is a kind of microdata, but the time horizon is very less. It is a kind of sort panel not a long panel where your cross-sectional dimensions are there, but not with by the time.

So, cross-sectional dimensions are obviously, capture certain variations in turn dynamics, So, the within effect is not so significant, but between effect might be there. So, that between

effect and within effect clarification I give it in the previous lecture. So, that aspect you can easily capture through the fixed effect model,

So, the fixed effect model is of course, useful in research and I suggest strongly that you please follow the line, between the space, between the subscript and superscript and understand it very clearly and you will enjoy it like anything, So, let us go ahead and clarify the equation once again. Starting with the very basic fundamental equation of the panel data. The dependent variable is with the subscript i and i that you can easily see.

(Refer Slide Time: 06:43)



I am clarifying bit by bit and giving you the right direction. You might not have received this particular information from my previous year's model and previous model on handling large-scale data. I did not explain much. But now I feel that based on the feedback of the students, I feel that re-emphasis should have been there on each of the coefficients.

Starting with fixed content, then with the slope coefficients etcetera. Now, you can see that this beta term is a is not going to change much, it is constant. So, that means, our by slope by slope dimension it is not identify the changes, but overall the model is going to capture, the certain dimensions in terms of cross-sections are captured in the constant term,

Therefore, it has certain assumptions, So, that the same equation is further elaborated as alpha. So, beta is beta 1 plus beta 2, this is not beta i, each of the beta coefficients is now estimated with its control variables. And these are all our control variables, and this is our

error term. An error term is, of course, attached with each of the dimensions, cross-section dimensions as well as time dimensions.

So, I just said that the slope coefficients are going to be constant, but the intercept varies across individuals not across time. The time component is not bearing the model as such significantly rather it is differentiating the model with its alpha term.
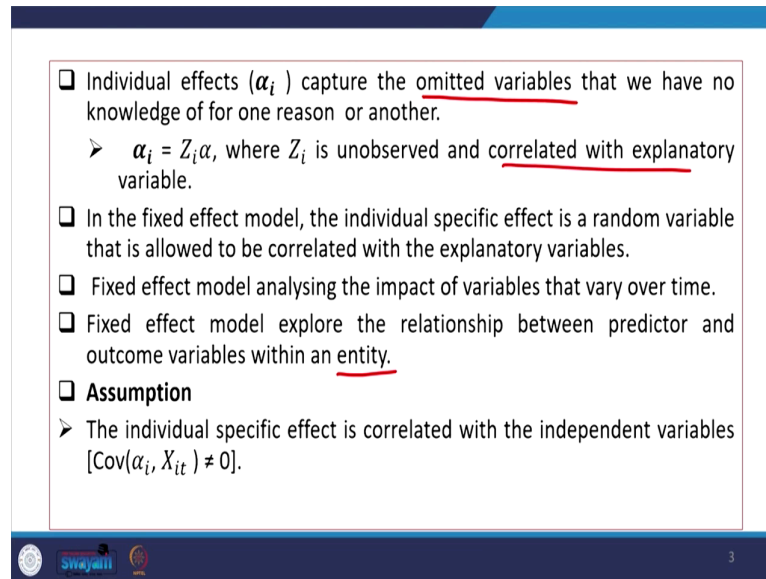
But in the common constant model alpha is actually constant. in the cross-section model or even the ordinary release square model, you could have noticed that alpha is not alpha i, Here alpha i is observed because of the panel content. on the aggregate where you have data, all are cross-sections, there is no question of the same person is indicating or giving certain views over time.

But the same persons when repeated even if you are considering minus the t term, but the same persons have certain cross-sections as well, cross-section variations within the time period as well. So, that is why the alpha term is actually capturing the dynamics of the cross-sections,

So, in the early earlier cases, the alpha term does not consider the dimensionalities of its cross-sectional changes. The above model allows each cross-sectional unit to have a different intercept term through all slopes through all slopes are the same. That is what is called individual effect,

Alpha is in fact, called the individual effect. Each individual's dimensions are actually captured. And this captures the omitted variables. if there are any sort of variables omitted by random procedure or by non-random procedure, in both the case omitted variables have a possibility of dropping or the possibility of deletion that is actually captured through the individual effects. That is why it is called alpha i.

(Refer Slide Time: 10:08)



So, the most important aspect for all of you to note is that omitted variable bias is actually captured in the case of the fixed-effect model. And that we have no knowledge of for one reason or the another about capturing this individual effect and that is captured. Then what is this alpha i?

Alpha i, in fact, is a proportion of the alpha that is through Z i, Z i is actually unobserved. And since it is unobserved consider to be correlated with the explanatory variables, So, that is why there are some problems with the OLS model. So, the normal OLS model is not going to be applied since the Z i is correlated with the explanatory term.

In the fixed-effect model, the individual effect or specific effect is a random variable that is allowed to be correlated with the explanatory variables. Fixed effect model analyzing the impact of variables that vary over time. The fixed-effect model explores the relationship between predictor and outcome variables within an entity.

So, each individual and it is effects are actually captured that is why we say within an entity. And that actually occurs due to the time changes, time dimension, though time-wise there are also changes, in the individual case there are expected to be some variations.

So, there, so the assumption in the specific individual effect is correlated with an independent variable which is why the covariance of alpha i and x is nonzero, and this is a serious error in OLS. So, OLS is not applied.

(Refer Slide Time: 11:55)



So, the methods for estimation of the fixed effect model are made through two approaches one is through LSDV model that is famously called LSDV or called Least Square Dummy Variables approach, and the second one is called within effects approach, Now, we are going to experiment with the LSDV approach to take certain variables to be with a dummy content with the dummy variables model.

And we will clarify how those problems can be avoided and then we can apply for an OLS model, The LSDV model incorporates the individual unobserved effects via dummy variables into the model.

(Refer Slide Time: 12:42)



**Least Square Dummy Variable (LSDV)**

❑ The LSDV model incorporates the individual unobserved effects via dummy variables into the model.

❑ The least square dummy variable model allows for heterogeneity among subjects by allowing each entity to have its own intercept value.

❑ Consider the model :

$$Y_{it} = \beta X_{it} + \alpha_i + \varepsilon_{it}$$

➤ Notice α has i subscript suggesting the intercept for entity would be different.

➤ Although the intercept may vary across subjects but does not vary over time i.e. it is time-invariant.

The least-square dummy variable model allows for heterogeneity among subjects by allowing each entity to have its own intercept value, So, through the dummy variables each entity has its own intercept value. Let us consider the model once again. This is equal to beta times X it. So, beta is constant here, but now the other two terms are alpha i and E it is actually having certain variability, but t is not there. So, t if I say it, it is actually not there,

So, t component is not there, but variability is observed through the individual variations and due to what is called if someone observes effects are not captured, so that can be captured through the alpha i. Although the intercept may vary across objects, but does not vary over time. So, that is the reason why it is called time-invariant model. This is also called time invariant. So, there is not term attached with the alpha i, That is why this is called time invariant.

So, the LSDV estimator is actually called a pooled OLS including a set of N dummy variables. This identifies the individuals and enhance an additional N parameter, since we are adding dummies N dummies, So, N parameters are again generated.

Nonetheless, this formulation does not suffer the problem of dummy variable trap since it contains no intercept. We have already read earlier about what is called dummy variable trap. If you take parameters as compared to the number of categories, if they are the same then it is going to have certain problems in the model, that kind of problem of estimation is called a dummy variable trap.

So, the dummy variable trap is avoided depending upon the alpha content we are taking. Like, in this in the first equation Y it is equal to the first beta 1 till X i term.

(Refer Slide Time: 15:02)



❑ The LSDV estimator is pooled OLS including a set of N dummy variables which identify the individuals and hence an additional N parameters. Nevertheless, this formulation does not suffer the problem of dummy variable trap since it contains no intercept.

➤ $Y_{it} = \beta_1 X_{1it} + \beta_2 X_{2it} + ... + \beta_k X_{kit} + \gamma_1 D_{1i} + \gamma_2 D_{2i} + ...... + \gamma_N D_{Ni} + \varepsilon_{it}$

❑ If using separate intercept term then include N-1 dummy variables to avoid dummy variable trap (i.e., the situation of perfect multicollinearity)

➤ $Y_{it} = \alpha_0 + \beta_1 X_{1it} + \beta_2 X_{2it} + .... + \beta_k X_{kit} + \gamma_1 D_{1i} + \gamma_2 D_{2i} + ..... + \gamma_{N-1} D_{(N-1)i} + \varepsilon_{it}$

❑ This type of model is known as **one-way fixed effects** model because only intercept is allowed to differ between individuals.

Now, the alpha unit is expressed into N into a N number of dummies since alpha is not there. So, that one is gamma times 1 D 1 i. First, i is captured, then the second like one X i X 2 i, X 3 i etcetera. Here we are defining each dummy till N dummies to replace that alpha i just to identify it its coefficient so that we can estimate the right coefficient in this model.

If using separate intercept term then in intercept is there and so N minus 1 dummy has to be considered to avoid the dummy variable trap issues, Then, there will be any if you do not consider then dummy variable trap is actually explained through the multi-collinearity problem and multi-collinearity is going to give you 1 as the value, because that there will be a perfect multi-collinearity.

If you have some difficulties here reading the dummy variable trap, I suggest you to read the simple text bo called by Damodar Gujarati. It will give you the right directions as well to explain to you a dummy variable trap correctly.

So, in this one when we are considering the alpha as the coefficient alpha term, then basically the that since this is there to avoid dummy variable trap we have considered N minus 1 dummy, So, N minus 1 dummies, then each of the estimations is going to give you the correct prediction of the model.

So, this type of model is also known as a one-way fixed-effect model, because only the intercept is allowed to differ between individuals. So, since intercept term is now actually allowed to be estimated through different intercepts that is why I called the one-way fixed effect.

So, obviously, there will be a fixed two-way, fixed effect model, in that case, time effects can be also incorporated with the individual effect, So, like the model includes N plus T dummy variables as well if time effects are also captured.

(Refer Slide Time: 17:27)



So, again without the intercept then there will be N plus T minus 2 dummy because now we have time dimensions as well. So, first of all, when time dimensions are there, this is our basic model. Now, this is without the time, so till N now you can see this is continuing to this till N without the alpha term um.

Now, we have the time dimensions, So, since so dummies are considered to be defined N minus 1 and N minus, in this case, N minus 1 and T minus 1 and so this is basically equal to N plus T minus 2 dummy variables, alright, with the intercept term if you are including the intercept term. If we do not have an intercept term then it will be simply N plus T dummy because the intercept term is not there.

This equation basically is without the intercept term. So, in this equation the dummy variables are N plus T minus N minus 1 and T minus 1 that is N plus T time minus 2

dummies are there. And accordingly, we can have captured each of the differences. So, this is called a two-way fixed-effect model.

There are some advantages and disadvantages. The advantages are it has simplified the model. The disadvantage is that there are too many dummy variables if you have a pool of cross-sections. If so many cross-sections are there then you might have so many dummy variables that lead to a loss of degrees of freedom.

And this may not be able to identify the impact of time in the variant variable. Many dummy variables will lead to the problem of multi-collinearity. So, that is there.

Now, I am also explaining what is called the within-effect estimation. Within effect first, we have explained the concept called LSDV, what is called the least square dummy variable model. We will have practical sessions. Once I explain this within effect estimation I will go for the practical session or add certain results with the original data and we will clarify.

The LSDV model is feasible to implement when the number of cross-section units is small, When the number of cross-section units is large then it is better not to use the LSDV but rather to use a within-group a fixed model that is through the within effect model.

(Refer Slide Time: 20:19)



We eliminate alpha i by expressing the value of the dependent and explanatory variables for each of the individuals as deviations from their respective means, When we know that the individuals are not actually varying, some of the variables are actually not varying much or

not at will varying, and some characteristics of panel data that is some of are not varying over time.

So, if you just deviate their mean value from the individual entries, then that will be defined, which will be considered to be having 0 variations. So, most of your variables are going to be having 0 values, and that way we can avoid many of the problems in the model.

So, the equations as per I just addressed are that this is our first panel data with slope coefficients are constant and then this is the constant term with certain variations, cross-sectional variations. Now, we are trying to take the averages, average overtime, and overtime of each of the explanatory variables. Now, what it is going to give us is that this is your average X bar, 1 i, then 2 i, then k i, alright. So, this is also average over time,

Now, we have taken the difference between the first equation and the second one. So, here it has the next equations to clarify what is this average. Average is basically over the time only T tends to 1 to T in each of the cases and you can find out the average accordingly.

(Refer Slide Time: 22:16)



The resulting values are called demeaned values or time means at each unit i; each by each unit i we have actually derived the mean values. So, the next step is for us to subtract these two equations from 2 to 1, So, 2 is our or 1 to 2 from or from we are subtracting from the equation 1 to 2; 1 is here and then subtracting these two. What is going to give us is that like

we have taken Y it minus Y it minus this Y i bar, so in each case, we have got the difference right.

Now, what is interesting is that if there are no such changes in the constant term especially when we know that it is having a certain fixed effect, fixed within effect; within effect is fixed, in that case, alpha i and alpha bar is going to be 0. When this is going to be completely 0, then the rest of the equation is simply of our OLS type of model, .

It is simply having beta constant values or beta coefficient having no further changes within it or dynamics within it. So, we will simply estimate the beta values, So, it is not beta i, it is only beta as a beta 1, beta 2, and beta k. So, each of the parameters we can easily estimate since now alpha i is avoided. Now, we will, therefore, use the OLS in the equation to estimate the parameters. One of the disadvantages of the model is the time-invariant variable.

Some of the time-invariant variables are wiped out because of differentiating. Even X i we say X 1 till X i 1 i bar, X 2 i bar etcetera it may be the case that one of these coefficients might be wiped out because of its constant variation, all the variations are constant for the same categories or some coefficient variable. So, that value will be also difficult to estimate because that is time-invariant.

So, another disadvantage is that this method is creating some problems in terms of differentiating. When we differentiate a variable we remove the long-run component from that variable. So, in the long run the time component is actually completely also wiped out. So, the time dimensionality is not going to be interpreted at all in this model. So, now, we have the practical examples and we are trying to show you the results with the help of the data.

The data set and variables are used similarly in all 3 models that are CCM, fixed effect as well as a random effect. The random effect model we are going to go deal with in our next lecture. At this moment we are comparing CCM and FE. FE stands for fixed effect. Now, you might be confused that how to get this data, if I have data how to develop a panel. I will be using WHO OLS IHDS panel data.

We have already shown you how it is defined to be a panel and whether it is strongly a panel or a strongly balanced panel or not, Rest of the details you want to develop your panel data you can refer to my module that was earlier developed which is called handling large-scale data set with stata.

So, the stata command to run fixed effect and random effect is your xtreg, xt regression; xt is often used in the cache, it is used in cache or panel data estimation before using xtreg we need to set stata to handle the panel. So, we need to set the stata we need to set or set the stata with the command xtset to set the stata about to recognize the stata about your data whether it is a panel or not.

So, we want to also understand whether it is a balanced one or not a balanced one, We can still run the model for unbalanced panel data with certain missing observations as well.

(Refer Slide Time: 26:53)



Now, here are all our commands and I am going to show them one by one. We will also compare the fixed effect using the LSDV model, then a comparison between CCM and LSDV fixed effect model as well.

(Refer Slide Time: 27:06)



Now, let us even in the next one it will be within effect also would be also captured. So, let us go for it and then clarify.

(Refer Slide Time: 27:14)



(Refer Slide Time: 27:19)

(Refer Slide Time: 27:20)



(Refer Slide Time: 27:24)

(Refer Slide Time: 27:28)



Here is the data being loaded on your screen and we are taking the Excel data provided by the WHO, example data set. It is balanced data as per the WHO, we will also test that as well. We recognize its variable ID variable as the first row.

(Refer Slide Time: 27:35)

Now, we will also show you the do file that is considered for the analysis, and in the do file, we will be experimenting with the commands and then give you the direction. Now, through the isd through the WHO dataset, we have already explained to you about the CCM model.

CCM model I said that how simple OLS is drawn, but once again we will read this data with stata, xtset should be made first because xtset we will recognize your data as strongly balanced you can see that. And the data is from 93 to 97,

Now, we will be also running the model as per this need regression. We have already made the results earlier for CCM, and now we are clarifying for LSDV model.

(Refer Slide Time: 28:43)



So, LSDV, show the LSDV model is drawn on your screen. You can see how many the dummy variables components are generated, for them to replace the alpha term, so, the number of degrees of freedom is accordingly reduced as I already mentioned, alright.

So, that is one of the wrong disadvantages of this data set. So, dummy variables, the way we explain there are so many dummy variables created, you can easily see this, 33 dummies are there, 32 dummies are therefrom, one is the base category then other dummies are there,

So, the constant term is also there. Show N minus 1 dummy has to be created as I already mentioned since 34 is the N in the model, So, minus 1 is created 30 32 dummies are actually created; the show. And other two so other two are as well there other two variables are also mentioned.

And that is HC3 and DALE. So, all like IID 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14 all till 30, yes 33 is mentioned, And in the next one, we will be operating with the next command that is we are we could have drawn the CCM as well. CCM is here,

So, just to regress the model we are trying to compare the equation and we will store this so that that is going to be helpful for us to compare, Now, this is being stored in our memory. Now, we will give the conditional command with i dot command the way we do since we have considered the dummy variables model.

Now, the result is displayed. So, now, we are also storing the estimates of this dummy with the i command, i dot command that we have given i dot id is given. So, it is basically comparing with the first categories of the dummy. Now, this is going to give us what we have already stored we can easily compare. Now, estimates will compare each of the tables that is OLS, OLS dummy, and with it a star and at what level this varies we can easily do it and the result is on your screen,

Now, OLS and the dummy variables model, so how this is important and how each of these is different, you can see that in the normal OLS or with the CCM model. It has only two coefficients and estimating, it is similar to that of the OLS model. But in the case of dummy one, it estimates so many coefficients.

Though we have run the OLS model, we can easily identify each of the effects differently as compared to the base category. And disadvantages of this we have already discussed therefore, we are going to explain the fixed effect model. So, the fixed-effect model we are we are now explaining.

In that case, we need to identify the stata per the correct specification with fe. We need to understand whether fe is specified or not. But there are some tests to understand whether it is a fixed effect or random effect that we will do in our next class,

So, the exact Horseman Test how that is significantly identifying or not that aspect we will discuss. Let us assume that it has certain fixed effects or within effects and time-variant effects. So, that is what we have assumed, and based on that we have run the regression.
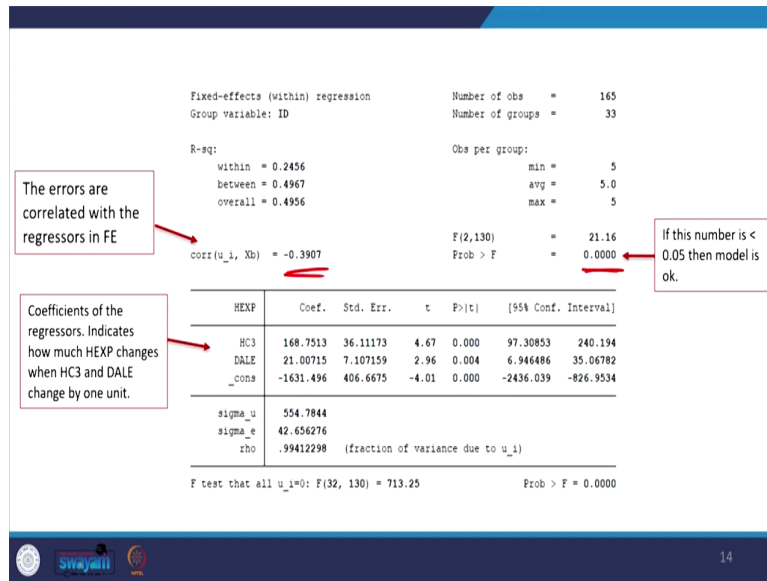
(Refer Slide Time: 33:25)



So, in the fixed-effect model, the correlation is more important. The correlation between the error term and the explanatory variable is there, it is not 0, So, the correlation is identifying the fact that it is, in fact, a model where the alpha i content, alpha actually content is correlated with the explanatory variables,

So, let me just explain to you through the model ones again, through the PPT ones again, then I will come back to the discussion. We have covered all those aspects. These are the command you have to follow the way we did it, we have kept capturing everything in the slide and that will help you to read between the results.

Now, in the within the model, within effect model we have already clarified what is called the LSDV model and what is called within effect model. And xtreg, r e g, xtreg command with finally, at ending with command fe will capture the fixed effect. And adding the robust option is going to control the heteroscedasticity.

So, usually, there are heteroscedasticity in the model, as I already told you. there are it is not avoiding the problem of heteroscedasticity. In order to control the heteroscedasticity, we usually take the command called robust fe, So, xtreg in at the end robust fe is going to control certain heteroscedasticity to run the fixed effect model. This is what it will lo like. And I will come back to it. We have already done this.

This is the one we derived and I am interpreting it through your PPT once again. And in the PPT, the first aspect is that the errors are correlated with the regressions in the fixed-effect model, And that is due to the fixed effect component, captured in what is called the individual effects or the constant term.

And that is why it has certain correlation with the error term as well. So, it is nonzero, therefore, there are some problems with the model estimation. We cannot just run the run a normal ordinarily square regression. Now, the second one is if this number is the probability the significance value overall significance of the model, of this, is less than 0.05 then the model is fine,

Then, the coefficients have been interpreted the way we usually do it. This coefficient indicates like the household the health expenditure per capita changes when the disability-adjusted life expectations and the education etcetera is changed by one unit, So, whether they are significant or not etcetera can also be interpreted accordingly.
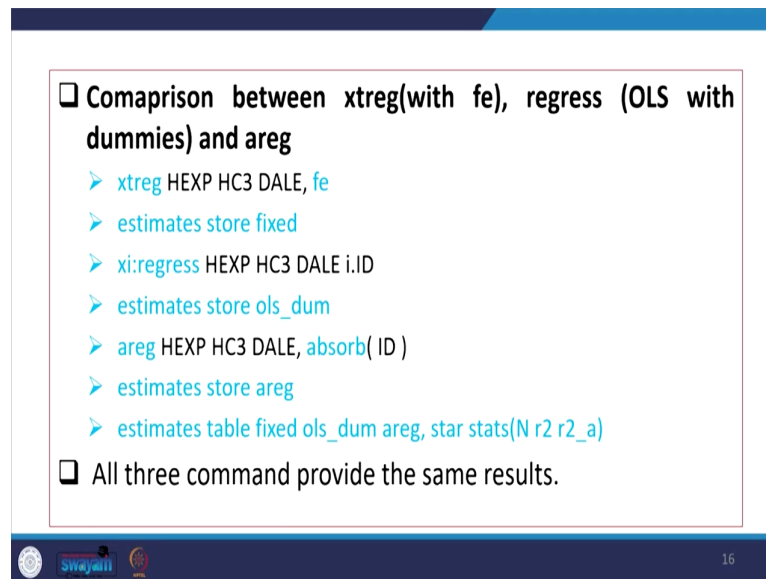
(Refer Slide Time: 36:50)



Another way to estimate the fixed effect using areg command. areg command actually this is going to observe the problems and it is going to hide the binary variables for each entity, So, the since binary variables are there that is going to observe it correctly. Although its output is less informative than regression with explicit dummy variables since explicit dummy variables are there areg command is the base fit and this does have two advantages.

It speeds the speeds of explanatory work providing quick feedback about whether a dummy variable approach is worthwhile. The second one is called when the second one is on when the variable of interests has many values creating dummies for each of them could lead to too many variables or too large a model in that case as per the suggestions Hamilton that you please go through the areg command, that observes the problems with the too many indicator variables,

So, that is also going to give you the result. The commands are already on your screen. You can apply on the same data, I am sure you will get it,

(Refer Slide Time: 38:15)



So, that we can have a comparison. Let me also experiment with it and going to interpret it accordingly, . Now, we are giving you the comparison, we will draw on each of them. Then, so the fixed effect result is being derived at this moment and we will also store this, . We will also store this. So, with the name fixed, and now once again we will draw the LSDV approach with this conditional command with i dot id. And again, we will store with an OLS dummy that is LSDV regression.

Then, we will also compare with the try to observe the identification variable or the dummy variables. Then a areg is the command and we will also store with areg. And now we can after doing so, we can estimate a table with the comparison of each of them. So, each of them now is being stored and now you can each of the they is estimated. Now, you can just see the differences.

So, we have all derived. The results are derived on your screen. So, let us see once,

(Refer Slide Time: 39:49)



So, all 3 are in fact, presented on your screen. This is I think you might have understood the way we proceeded. All 3 are stored. The first one is fixed-effect model the fe g, fe command is given, then the LSDV model is here. Now, LSDV is also reduced or is interpreted in a better way through areg, that has been made.

And now you can see you will get the same result, through the areg. areg and fixed effect model are going to give you when you have already captured the dummy variables problems, So, any of the models you can do it, run it and the any of these two models. You can do it and estimate it, fine. These are all for this lecture. We probably, we have if any things are there we can also explain you otherwise we will; yes.

We can have some comparison as well with the constant common model results, it is not just the 3 within effects of LSDV. In fact, we can also compare the simple regression estimation as well as the fixed effect coefficient as differences. To test whether fe is better than CCM or not we can apply restra icted F test, .

So, the restricted F test is going to give you the direction. The null hypothesis, in this case, is whether alpha i is equal to alpha because alpha is the one in case of CCM, and alpha is the one in the case of FE. So, whether these two are actually equal or not. S, so the null hypothesis is going to be considered as equal. So, our hypothesis is in fact, alpha is equal to alpha.

If the null hypothesis is not rejected, then CCM is appropriate, If both are the same then what is the point of going for a panel regression; if they are different then, of course, you should go for fe test or fixed-effect model. The output of the within effect model with xtreg command a conducted this F test by default,

If the calculated F is greater than the tabulated F, then we will reject the null hypothesis. So, this is the one that is very important; when your calculated F is is different, or greater than that of your tabulated one; the tabulated F value then, of course, that is going to reject the null hypothesis. And fe is going to be considered to be better.

In this case, you can see that here is the calculated value which is 21.16 that is the F value.

(Refer Slide Time: 42:54)



With the xtreg command with fe command, we have given, . If this is this is the one then we will compare with the F value with its degrees of freedom, degrees of freedom with the numerator and denominator values, k minus 1 in the numerator and N minus K in the denominator the calculated value we can compare. If the tabulated value calculated value is different than that of the tabulated value, then we will go for rejecting the null hypothesis,

(Refer Slide Time: 43:25)



So, that next one is called testing of heteroscedasticity which is also equally important in case of the fixed-effect model. We have already said that it does not avoid the problems of

heteroscedasticity. So, a test is quite essential. A test of heteroscedasticity is available for the fixed effect model because of its assumptions, So, the commands we usually give it is xttest3, then we need to actually install it and then run the xttest3 because it is not come up with the version by default.

Then we can run with the xtreg and its fe and that and at the end after the regression, we need to just test xttest3, I will do that here, but I am just interpreting it if that is the case that the assumption is that should be homogeneous, that should not be any heterogeneous.

But if it is significant if the p values are significant, the null hypothesis is that the standard deviation the sigma is actually equal for all is, but if that is significant; that means, they are not equal, so heteroscedasticity exists. So, this is what I am going to show on the screen.

Once, then we can go for it and then explain it, So, we will install that first SSE install xttest3. So, this is being installed, Now, this is and verifying not installed. We can do it then.

(Refer Slide Time: 45:18)



Then, we will now run it and now the installation is completed. Now, we are running the fixed effect model in our case. So, we have already run the fixed-effect model and now we are going to test the xttest3 that is for the heteroscedasticity. Now, it is confirmed in the case that our significance level, the sigma values are actually differing that is differing significantly therefore, there exists heteroscedasticity in the model.

So, that is another way of proving that yes we are running the fixed-effect model and it has heteroscedasticity in the model and it is not going by the common constant model, The common constant model is not the right fit when this kind of test are identifying the heteroscedasticity, alright. So, the fixed-effect model is best chosen.

And in the next class, we will discuss choosing between choice between fixed-effect model and the random effect model, and we will also give you various other dimensions of clarifying random effect in the healthcare database. We will be also using IHDS data. So, we have also kept the IHDS data. Your own reference you can also test on your own.

We have only operated through the WHO data set. If you still have some confusion you can also refer to my earlier model, we have used IHDS data. This time we have we are also keeping the same IHDS data for your reference. We have operated with one database, you can operate on your own and find out the difference. The do-file of both the cases has been uploaded and that will be going to be useful for you to run it.

I hope this is how you can get the best understanding out of it. And if still you have difficulties do not hesitate and come back to us in our live session or in our chat session. Our team is quite charged to deal with all your queries.

Thank you.