**Handling Large-Scale Unit Level Data Using STATA**
**Professor Pratap C. Mohanty**
**Department of Humanities and Social Sciences,**
**Indian Institute of Technology, Roorkee**
**Lecture 14**
**Managing Data in Stata - II**

Welcome friends once again to the NPTEL module on Handling Large-Scale Data Using Stata. We are trying to explain the use of Stata with the core data, the original data. I have already guided you in the last class about some commands, some basic commands, like summarize, codebook, listing, browsing. All those things we have done in the last class. So as a continuous into our previous lecture, we have mentioned Managing Data in Stata II. So here we are going to guide you further. Last class we have shown you codebook and its compact details.

(Refer Slide Time: 01:22)

I am going to tell you the commands for checking unique identifiers. I told you that identifier is very important for analysis. Which are the unique identifiers, how duplicate reports are important for analysis. Like, if you enter duplicate report, let me clarify what you mean by that. If a variable has a repeated value, it will produce a table with number of copies, there is observations, likewise, visible in our table here, these are sample result we developed, copies, observation and surplus. If there are no repeated values, it will simply give you zero surplus.

So zero surplus means all the observations are uniquely identified, likewise this. So duplicate report, I am just going to click that for you to operate duplicate, there is a mistake, let me clarify,

just a minute, so duplicates report space, then the variable we wanted to check, now enter. So this gives information about the copies, the variable we wanted to check has zero surplus. So zero surplus means no duplicates are there in our data.

Let me go through once again. We can also take this with our IHDS dataset. How many duplicates are there? So, IHDS dataset we can check like here. Let me open that in our window. So here is the IHDS one, let me click ok. So here is the IHDS sample data we have just opened. from the IHDS data, again if I go by the same command duplicates report, what you do, you can go by page up command. I have the command here duplicates report here. instead of this, what I do, I will enter the ID for us.

If the second one is ID, HHID is there. So if I just enter it gives me the information like this. Here how many duplicates copies are generated. Look at so many duplicate copies are generated. So it suggests that there are duplicates. So just one variable is not uniquely identified, we have to go by a combination of variables for finding the unique ID.

So, it is not coming so many are there. what I suggest you to go through this interpretation, like if it is there like one copy, 23 observation, and zero surplus that means 23 observation are clearly unique ID. But there are two copies of such, 22 of such observations having two copies are generated. So two copies divided by two, 22 divided by two, we have 11 surplus left.

Similarly, for three copies, type of three copies are of 24, if I divided by three, so it will be eight. Eight already. two copies of 22 observations are there. If I just multiply then with that eight into two, 16 will be the surplus. So these kind of surplus information are visible before us. You can find out and rest if you operate, you will find out differences.

(Refer Slide Time: 06:31)



Let us understand the easiest way of tracking the unique IDs. Let me point out very clearly that it is very important for you to merge the data with another dataset. If you want to merge, unique IDs are very essential. Or even sometimes if you try to make merging for panel information, it is very very essential or some blocks in NSS data there are so many blocks information they give, if you wanted to simply merge those blocks, you need unique ID as well.

Now, check this ISID checks whether the specified variables uniquely identify the observation or not. So we have already told you if surplus is there that means those are not uniquely identified, that number of copies are generated, so those are not uniquely identified. One, let me take a note here that single variable that is variable in the syntax command var command or one variable in the command is used for single identifier.

And if you include varlist, it include the compound identifiers, list of variables if you give that may give that. Combinedly they are the uniquely identifier variables that generate a single variable which is uniquely identified. But individually there may not be, let us test it. In column identifiers, what do you mean by column identifiers? Rows we have already told you they are the observations and in Stata columns are identified by variable names. Variables are generally entered in columns. Variable names are always unique. You cannot have same variable names together in Stata. It always popped up with an error command or to suggest you to change the name since the same name already exist.

So Stata would not allow you to create two variables of the same name. But they often have multiple parts, like the name is there, if you just add another variable with an underscore or with another addition, you have another part of that particular variable to identify whether that is a different variable or not, but connected with the original variable.

We will operate everything steadily in our successive classes. We will keep on using those commands. But at this moment, I am just guiding you which are the simplest way of clarifying some entries.

(Refer Slide Time: 09:18)

So let us convert string variable to numeric variable. As I already told you repeatedly that string variable, mathematical operations are just not possible directly. We need to convert first then can operate. Let me remind you once again that, there are some categorical variables also entered as string. So that can be further converted to label numeric values. So convert number variables stored as string to numeric. Basically those are numeric variables or number variables, but stored as string. We need to convert them to numeric to do our operation.

So it is very common for numeric variables to be imported into Stata as string. Before you can do much work with them, they need to be converted into numeric variable as I just mentioned. This

can happen when one of the numbers was mistakenly entered as a later or a non-numeric code is used for missing. There are various ways of getting a string variable in our data. Besides the manual entry as string, we have also other possibilities like mentioned here. So mistakenly we can enter letter also in our data. So that converts the data to string. Similarly, a non-numeric code also uses for missing also read as string.

So let us understand how to convert it. The most important command here is destring. Destring command is highlighted in blue color, is the easiest way to convert string to the numbers. Variable is string and by various ways you can understand whether it is string, through codebook you can understand whether it is string or not, through summary statistics also you can understand whether it is string or not, simply browsing the data in the browser window, you can get whether it is string or not by color. if it is red, you will get string.

Once you have observed that variable is string, you need to destring it. You need to convert it to a numeric value. So the command is destring, then the variable name or variable list if you have so many variables and you simply replace, that will replace the original value and it will give you the new that is the destring value. But when you want both the variables together, you need to generate a new variable, let me operate here.

The simplest way to check whether it is string or not, I am just going to the data. This is the IHDS data. It is a bigger one. Let me start with, string variables are highlighted in red color. the blue color is in labeled numeric and the black color are in numeric values or continuous variables.

 I can just change, suppose I know that state ID, I know that IDHH and ID person, these are in string data. But let me operate from the NSS at this moment. We will clarify later on from the IHDS data. At this moment, let me operate from the National Sample Survey, because we have already sampled and filtered the data for quick understanding.

So what I do, I will simply destring. I told you already destring, then variable name, then replace. So then you just go through destring, like this destring, then variable name. The variable are string or not, once again you can check. This is the variable. Let me come back once again I will open that, open once again the variable I will open the NSS then I will operate.

Let me check once that whether that is a string or not. And I know that there are so many variables in string. State, for example, state is in string. Let me convert it to a destring one. So the command is here destring, then state you just click on state. if I just add comma to it, comma then replace, it will replace.

It is always suggested that you keep both the variables. If you have already filtered and shortlisted your data, your dataset is very minimal, so try to keep both. Otherwise, if you are doubly sure your string data is not going to be useful at all and you go by a new variable then you may continue with that, continue with this command.

Let me first start with a destring then variable then generate. Then I will show you, like for example here, let me put a comma here, then generate, I am trying to generate a new variable. Then within bracket a variable name state, so statenew. Bracket is not there. Let me give the bracket, shift bracket, then bracket close. If I just enter here, destring state generate a new variable enter. look at another statenew variable is defined here. I am just putting my cursor. It might be visible to you, state then next statenew.

I will just show you from the browse window or for simplicity, for your simple understanding, br then statenew enter. It has shown you in numeric. this is no longer in red color. Whereas you just compare these two together instead of one just compare two together, like here you just enter state once again and enter, we will get both together. Initially it was state with red color, now statenew is in black color that means it has confirmed us that your data is converted to a numeric number. Now you can apply your mathematical operations for results.

Let me move on further. for your clarity, if I just do the same command destring, it is already defined. I am just going to do it, like if I just say replace instead of generating a new one, it will convert. Let us experiment it. So, destring state replace, what it gives, it has replaced the original string to the new state now is in numeric format. Look at this state is now here state you compare state and newstate. State since we have already replaced with a numeric variable, this is highlighted in black color.

Let me proceed further we have already clarified. So we will suggest that you please go by this then later on you can drop that particular variable if you want or you can replace later. But at this moment a new variable is suggested to generate and operate likewise we did.

It is also interesting to note that the string the way we did is not enough. In case of some categorical values are entered in the string format or string variable, then in that case encoding is important. Encoding string into numeric variable is very very important just string is not enough. String command is not going to work.

(Refer Slide Time: 19:25)



In Stata, variables used as either independent or dependent variable typically cannot be strings and should be made numeric. Typically, the string operation is not going to make string. But can be made numeric. These are generally categorical variables and are often represented by string variables, which we have already clarified in the last lecture with the names of categorical, categories as the values. The categories as the values are entered, but not numeric one.

So the encode command is in place of destring we did in the last slide. Encode, rest of the commands are same, encode is going to help you. Encode command converts string variables to numeric, by assigning a numeric value to each string that is categories and then applying value labels of the original string values. Once it has converted to a numeric value, we can then apply the label, on our own convenience based on the given variable description in the original data. So once the label is defined, we can also clarify in our converted data also.

The numeric version of the variable can either be a new variable using option, like the way we did generate the new variable name within brackets or a replacement or replace we did in the last slide likewise this you can do it for the string variable using an option replace.

(Refer Slide Time: 21:12)



We did so far one interesting part that is called string. Many researchers usually commit that mistake and get frustrated with the data, that the data is becoming difficult to convert. Some errors are coming. They repeatedly search here and there I think the last couple of slides are going to help you a lot in understanding the data and converting to a numeric. The standard rule

is always to convert to a numeric number for quicker operation. There are some advanced commands that can also continue with the string value as well, but we are not dealing with that part.

Organizing dataset is another important aspect of understanding Stata and the large dataset. The sort variable, sort command of the variable is going to help you a lot. Sorting of, likewise in Stata, likewise in Microsoft Excel, we used to sort the variables. If there are so many variables just need to analyze together, we need to sort one and accordingly match whether others are following or not.

So sort observations in ascending order in a dataset is generally essential. It requires a variable or list of variables for sorting. If there are so many variables you are sorting together, it by default consider first as your sorting variable. And on the basis of the first variable, it sorts other variables also. You might have operated in excel also.

So the command for sort is sort then variable name. Let the variable is sector. Then how it is like this. Now look at what I do browse then sector, it will show you it has already sorted look at this I am going down, one then two is there. So since it is in ascending order it has already sorted.

So let me move on. If you give more than one variable with this command, it will by default sort the first variable and then accordingly saved the data with second and third variable as per the first variable. So you must sort each dataset by the linking variable prior to that of merging.

I think we will also clarify while merging and appending and we have some dedicated lectures also for it. We will clarify that it is an essential step before merging that you need to sort the data. Without sort, the variable, the uniquely identified variables cannot be sorted, it cannot match the corresponding unique ID of another dataset.

So once you have sorted it that means by each row all are same in both the datasets. So then only linking variables as well as the original variables both should be sorted before merging and then only merging is going to be successful.

(Refer Slide Time: 24:48)

Let us move on understanding some other operations like renaming the variable. A good variable name tells you very clearly that what the variable contains. If you have named very appropriately, for example, you are working with GDP, but you have named a state, but, probably you will be confused later on when you open and you have already operated with many other information.

State maybe, state name is there, but if you have not mentioned GDP, you will be confused. But, like agriculture, output is there, but you have given gender. Generally, this mistake we do not do but some close name we generally give, but still make sure that this is very unique in your understanding that most appropriate value must be given for better and quicker results.

For our syntax like rename is the command you are supposed to give. Rename old, basically if an old name is there, you need to rename that old name. The dataset has given certain by default names as per their convenience, but you wanted to continue with a new name that name maybe very suitable as per your explanation, your model so you add a new name. So, old name will be renamed with a new one, I will tell you right now.

The variable name cannot contain spaces as I told you, there are two ways of renaming variables. One is renaming is camel case, another is called snake case. Likewise camel has a continuity, so each name you need not give spaces, you just capitalize the second starting letter. It is not a continual flow, there is a capital letter. And in case of snake, there is a clear space. This is

different than that of another one. So you need to give a underscore. Underscore identifies for separating two different words in the name, we will clarify right now.

So let us operate. Suppose I just wanted to rename here, statenew is there. I already defined statenew. Let me rename it to another one like stateone or maybe stateIndia if you are working with other countries together, your data containing other countries, state of India, like state, two approaches as I told you, camel view, so India will be capital like this. So if you enter it will be renamed like this. Now it has already been renamed, you can just have a look from your variable view. I have pointed out, I have put my cursor here. It is getting visible to you.

Otherwise, you can also do underscore operation like the snake view. You just give underscore India. The name has already been changed. So statenew is not there, it is now stateIndia. So I am just going to change that stateIndia to state_india. Now this has changed to state_india. So what is important now the variable has been changed with a new name as per your own convenience. And this is very essential and usually very useful.

(Refer Slide Time: 29:26)



❑ **recode and replace**
- ❑ Variables are often not coded the way we want, often with too many categories or with values out of order.
- ❑ This command can also be used to recode missing values to the dot that Stata uses to denote missings.
- ❑ recode secondaryEducation (.a= .)
- ❑ With string variables, however, you need to use the replace command .e.g.
  replace city="prayagraj" if city == "allahabaad"

Then recode and replace. Recoding and replace is equally important for work and for operation in our module. Like, for example, if you have a data like gender, for example, gender is there in your variable, 1 is coded for male and 2 is for female, but you want it to, when your result, your

interpretation will be higher, positively linked towards, since higher weight is, as per the numeric number higher weightage with female. But it might be becoming very inconvenient for you.

Like, if I go by how many males are impacting the society in any particular policy framework, your data is related to any particular policy analysis, you wanted to understand or interpret through male, not female. But where female is important, you wanted to answer through female, your weight is higher in female. there is no harm in dealing with that, but for some convenience, if you recode it to, like female you make 1, recode to 1 and male at 2.

Suppose you wanted to do that, now your regression coefficient is coming out to be positive. your likelihood analysis will be when the society is more of female, the likelihood of getting advantage or disadvantage is positively linked because of higher male members. Since you have already given a positive or higher value to male, so your relationship is expected to be positive there. But depending upon, it is not a standard rule at all.

There are many sort of example I will do it in between let me go by that. How we are recoding it, I will explain it later. Recoding variables are often not coded the way we want, as I just said. Like in some dataset education attainments they give it in reverse order. Educational attainment like class 10th they have given 1, 12 they have given 2, for higher education 3. If that is there, higher weightage is there already. If the reverse order is there already in the data like 3 weight is for 10th, then 2 is for secondary and 1 is for higher secondary, then that might be inconvenient for the researcher to analyze so you need to recode it.

So, this command, the recode command can also be used to recode missing values to the dot that Stata uses to denote missings. Like as I already pointed out earlier that Stata usually reads dot as the missing values. But if there are some other entries, by mistake also there might be some entries, some alphabet are entered in the numeric data and you have checked that some alphabet has been entered and converted to a string value and you have detected that.

So it is very difficult that which number has replaced with that alphabet. So in that case the easiest way to convert that number to when you do not know the exact replacement, you simply convert it to a missing value as dot. So recode that particular entry as dot. So recoding command is going to be very useful.

I will tell you some of the sample entries for our clarification. So recode secondary education, if secondary education is our variable name, your command will be recode then the variable name, whatever the variable name you have. If you have a .a entry, as I just pointed out a couple of minutes back that you just wanted to read this as a missing value, you simply read as .a equal to simply dot. Recode, in the command recode .a or if 98 is entered in the data or usually 8888 are considered to be the missing value in different dataset, 88, 89, even these are the standard missing values in our different datasets. I think I told you already while guiding the datasets. We have already guided you, reading the missing values in the data. So you can recode it to dot.

So its string variables, however, you need to use the replace command, not the recode. If it is numeric, so we can convert it. But if it is a string variable, the standard recode command is not going to work. In that case replace that. If city is a string variable and with some entry Allahabad, you want to replace it by Prayagraj you have to enter as replace, not recode. Recode is going to be erroneous. It is not going to give you the right result.

So if you want to change it to another one, just change it by that name. But again it is a text. It is in string variable still now. Even if you have recoded, it still converted to a string variable. What we will do for your analysis, we have to convert the other way for analysis to destring, in that case you need to go by the encode command, or sometimes wherever it is only in string but not in label, not in categorical labeling or categorical values in that case if it is only in string format with characters, your conversion will be through destring I have already pointed out.

So here for recode and replace, recode is useful for numeric variables, replace is useful for the variables which are string in nature. We will continue rest of this from keep and drop command and there are other entries in our analysis in our next class. So thank you very much for patient listening.