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

### Lecture - 22 Data Extraction from ASCII Format

Welcome friends once again to the NPTEL MOOC module on Handling Health Care Data. We are on the fifth week and trying to understand the survey data and their analysis. On this lecture we will be giving you the hand holding of NSS data and their most important aspect is on data extraction especially from the ASCII format database. Myself, Dr. Pratap C. Mohanty, I have been teaching research methodology and these type of courses over 6-7 years.

So, I will be most happy to address all your queries. So, let us move forward and understand what you mean by this data set and how we can set the tuning for understanding it for final use. Starting from the very beginning of understanding data especially in Stata. Stata has it is own format for storing data sets in *.dta* format that has to be noted very clearly. So, *.dta* if that extension you see in the file; that means, this is the Stata data.

(Refer Slide Time: 01:40)

INTRODUCTION
<ul> <li>Stata has its own format for storing data sets, .dta files.</li> <li>These are highly convenient for Stata users because Stata can use them immediately without any need for interpretation.</li> </ul>
But many a times large scale unit level datasets are not available in Stata format. They are available in different formats or in raw formats.
Raw data can be defined as unstructured or unprocessed data.
🔄 Swajan 🔮 2

These are highly convenient for Stata users because Stata can use them immediately without any need for interpretation. But many a times large scale unit level data sets are not available in Stata format. They are available in different formats or in raw formats or in ASCII data format or in simply in notepad format.

Raw data can be defined as unstructured or unprocessed data. So, lots of understanding is required and based on that we can go for processing the data before final analysis. Data extraction is in fact, needed to retrieve the unstructured data for further data processing and analysis.

(Refer Slide Time: 02:27)



Usually, the large-scale survey data sets are available in Excel file, either in .xls format or in .xlsx format or in ASCII text data format. If your data is originally in Excel or some other format, you need to prepare the data before reading it directly into Stata. Because it is labeling and other thing that must be checked first and then you can import it to Stata.

So, then what do you mean by ASCII text data? Some clarification we are trying to give here. These are basically a kind of American Standard Code for Information Interchange that is the full form. It has set of codes usually of 256 characters saved in computers. These includes alphabets, numbers then control characters like enter, shift, space then and some graphic characters as well.

(Refer Slide Time: 03:36)



The ASCII text datasets are raw data, or they are also called delimiters. There are in general three groups of ASCII text depending upon delimiters. Three different types; first one is called free format, then delimited format and fixed format. So, these three we are just going to explain you one by one for your clarity.

(Refer Slide Time: 04:07)



First, in this particular format, that is, free format the data considers spaces in between. There are clear designated spaces involved to separate the variable or the information. So, in the

example dataset you can see the spaces are so clearly identified and these differentiate the entries.

And the format is simple and intuitive enough to be used for small data, usually this is applicable for small data entries. If your data are in free format with variables separated by blanks or comma or tabs you can use one command for extraction that is called *infile*. *Infile* command is useful for this kind of free format data.

So, in that case you can use infile then variable and you can specify the variable names, variable 1 or till 5 or it is byte position if it is given. Then, 'where is your source of data?' should be clearly given. Because if any earlier stored data is available, it may not able to extract it.

At this time since we are not going to show you with *infile* format data or free format data. We are not going to show the exact operation but just for your clarity we have mentioned the command and we will operate the infixed format data and on that basis you can also understand this kind of data set if any, you are having in future.

So, *infile* command always requires a dictionary file through which you can able to extract the data.

(Refer Slide Time: 06:11)



We are going to define dictionary file shortly: What do you mean by dictionary file?, How dictionary file is defined?, How a dictionary should be made that Stata is going to read before

it is extraction? Then we need to note here, that the using modifier in this case requires the data dictionary file that is *.dct* not a raw data file.

So, once you have opened the Stata and your dct file has already specified the source of the data, it is path name has already been discussed in the dct file or in the dictionary file. Then after the using command your dct file is most important, no need to further open your raw data. Both the syntax will give the same result.

On the first page i.e., before this page, we have given the variable name and their byte position. Without the dictionary file on this page we are guiding you through the dictionary file both will give you the same result.

(Refer Slide Time: 07:19)

	<ul> <li>Delimited Format</li> <li>Delimited ASCII text format uses other delimiters than space, where a</li> </ul>	
	special character (the delimiter) separates one value from the next.	
	The tab-delimited and comma-delimited are most common in this format, but any special characters such as @, #, \$, %, ^, &, and * also can be a delimiter.	
	The most common form is CSV (comma separated value) files, where the delimited is a comma, but the delimiter could be a space, a tab, or	
	in theory just about any other character or set of characters.	
	The extension of csv file is .csv.	
	csv is able to deal with complicated and ill-organized data from spreadsheet and database.	
Ó	7 T	

Next format is called delimited format. Delimited format is also important and these are also one form of text ASCII data format. Delimited ASCII text format uses other delimiters than space or comma. There are other delimiters could be used. So, like tab delimited, comma delimited two are most commonly used.

And besides that, some special character can be used as delimiters. Dollar symbol could be also given at the end of the variable or the beginning of the variable to separate the variable name. So, the delimiters could be of these: you may be asked in the exam that, 'can these are considered to be delimiters in some of the data in your objective question?'.

So, be careful about it. You have to be very clear that these are going to be useful when we have a delimited format data. The most common form of delimited data is called CSV that is that Comma Separated Value files, where the delimited is a comma. But the delimiter could be a space, a tab or in theory just about any other characteristic or set of characters used for separating it that is why, it is called comma separated delimiters.

The extension of csv file is called *.csv*. *.csv* is able to deal with complicated and ill-organized data from spread sheet and database.

(Refer Slide Time: 09:12)

	<ul> <li>This format often has a list of variable names at the first line.</li> <li>Example</li> </ul>
	ENTID,State,b2_q204,b2_q210,b2_q211,b2_q212,b2_q213,b2_q216,b2_q218,b2_q219 6950911502,9,1,6,1,2005,1,2,2,2 8383911101,3,1,2,2,2013,1,1,1,1 5116511403,21,2,1,1,2013,1,2,2,2
	<ul> <li>If your data is in delimited format, you can use the insheet/import command.</li> <li>insheet using "filename", clear</li> <li>import delimited "filename", clear</li> </ul>
	Or you can also import it from file menu File -> Import -> text data (delimited, *.csv,) -> browse -> submit/okay.
Ó	swayaan 🚳

The sample is given here like here a list of variables are given and these are separated by comma. And so, delimiter here is our comma but in case of our first example that is in this free format data we discuss about space is the delimiter. If your data is in delimited format you can use the *insheet* command or import command. Both the commands are useful; *insheet* or *import* command will give you the same result. But if you are using import then you have to discuss the delimited import delimited file that has to be mentioned.

If it is *insheet* then *insheet* using "file name" has to be mentioned. Otherwise in Stata, if you go to file then you click on the very first important primary buttons you will find file. And then on the drop down menu you click on import, then text data, then delimited or *.csv* format, then browse and submit, you can get the result.

(Refer Slide Time: 10:39)



Fixed format, the third one which is mostly used and we are also going to use it now for practical purposes as well for showing you the direct extraction result. So, fixed format files are usually of text files, but unlike csv file there is no separator between variables. So, all the data points could be conjoined could be set together without any space or any delimiters.

So, those delimiters even if it is not there, still we can read between the entries because of their space identification in the layout file. We have shown you earlier about layout file. Layout file gives very systematic presentation of your byte positions, byte positions then, even if there is no comma separated files still, you can able to specify the variable names.

Then survey data often come in these formats with one or more records per case and each variable in a fixed position in each record. This recognizes data items by column position i.e., which column position these are entered and also their column ranges. The column positions are fixed for each row. The data for every row files you will find the same position that is why we are saying it is fixed.

The data cannot be read directly, but a layout file is required. So, layout file we have mentioned here. So, we have shown it in our data when I discussed on the very first week of our module.



The advantage of fixed format file is that they are smaller and no space is wasted or no separator is required. Therefore, the byte space it consumes because of that file is usually smaller. The disadvantage of this is that they are completely useless without the corresponding data dictionary.

If data dictionary is not defined correctly then it is actually very difficult. Since we have large number of variable and there are entries requires, a layout file specifies the byte position by its column distance or column position.

So, every time it is very difficult to remember and operate. So, a data dictionary is required. If corresponding data dictionary is not supported then this is going to be very cumbersome. If there are many variables to be read it is efficient to write a separate file to provide necessary information, the file is actually called data dictionary.

So, this is what I was saying. If so many variables are mentioned by their byte position it is better to prepare a separate file with all those details. We are going to guide it shortly. Those two type of file are called data dictionary. This defines variable type, name, position, format, level, value level etc.

So, data dictionary includes first, the variable type, whether these are in string or numeric. Then name of the variable we wanted to specify in the dictionary file. Then its position, byte position and its duration then we can also specify their format or also we can specify their value level etc. can also be given on the dictionary file as well.

A data dictionary is defined as the file that explains how the variables are read and write. In India, most of the large scale unit level raw data sets are published in this format some examples are NSS data sets, then economic census data sets etc.

(Refer Slide Time: 15:19)



Understanding extraction of unit level data; we have used NSS 75th round for extraction which we have said from the beginning of our module that we will be exercising and this module on NSS 75th round.

Similarly, you can also extend it to any other round. This round is on health care and as on NSS 75th round there are in fact, two sub-round that is, one is on health care another is on education.

### (Refer Slide Time: 16:12)



So, micro data you need to get for the raw data. Once you click on this, the micro data catalogue gives you all available micro data of NSS, but once you click on this download tables data you will be redirected to the concerned page. But as we have already said it requires a minimum i.e., a very basic registration which will be further approved based on your registration.

Usually it is not paid, it is freely available for the public. For other health round data, go to the micro data catalogue and register yourself to access these data. With the page which we have already guided earlier will look like this. At this moment we have also other data sets available like PLFS.

(Refer Slide Time: 17:12)

>	Unit Level Data of Periodic Labour Force Survey (PLFS), July 2019-June 2020 Click on this
>	Unit Level Data of Time Use Survey, January 2019-December 2019
>	Unit Level Data of Periodic Labour Force Survey (PLFS), July 2018-June 2019
>	Unit Level data & Report on NSS 75th Round for Schedule- 25.0, July 2017 - June 2018, (Social Consumption: Health)
>	Unit Level data & Report on NSS 75th Round for Schedule- 25.2, July 2017 - June 2018, (Social Consumption: Education)
>	Unit Level data & Report on NSS 76th Round for Schedule 26.0 (Survey of Persons with Disabilities)
>	Unit Level data & Report on NSS 76th Round for Schedule 1.2 (Drinking Water, Sanitation, Hygiene and Housing Conditic
>	Periodic Labour Force Survey(Micro Level Data)
>	Document related to national account
>	National Accounts Data
>	Statistical Year Book

We are here trying to download this one. So, you just click on this link, and you will be downloading the required data.

(Refer Slide Time: 17:22)

Download all the files	J.
Unit Level data & Report on NSS 75th Round (July 2017 -June 2018) for 25.0(Social Consumption: Health)	Schedule-
Appendix-4	<b>a</b>
Appendix-II	
Chapter 1	<b>—</b>
Chapter 2	<b>—</b>
Chapter-4	<b>—</b>
Chapter-5	<b>a</b>
Consolidated Corrigendum	<b>a</b>
README 75 Round Schedule 25.0 (Download Data Layout & Unit Level Data)	≡ E
Estimation Procedure NSS 75	<b>a</b>
Nic Amendment 2008	<b></b>
Schedule 0.0	<b>e</b>
Schedule 25.0	<b>a</b>
Key Indicators Report of Household Social Consumption in India Health, 75th Round-(July 2017 - June 2018)/p>	<b>—</b>
💿 swym 🧕	14

Now, next one is to open, once you click on that you will be redirected to this particular page. You will open this page and within this you have a link called layout and README file. So, you just click on this, a PDF will be opened in another page. On this PDF you can see a different guidance. (Refer Slide Time: 17:50)



(Refer Slide Time: 17:51)



So, once you click on this you will be shown a complete page. I will also show it here. So, you simply click on each of the file, on this NSS 75th data there are 13 levels information. Why levels are important? I have already told you that the ministry has categorized the different subsets of information.

Like household is separated from individual, then within individual they will also categorize for those persons who accessed health care or having some forms of element over last 15 days or in last 365 days or there is certain information about deaths.

So, these information are important. Just for your clarity once again; once you see these data having number of cases is 555352. And this is in fact, the highest, that does indicate that this must be an individual file where it includes all the individuals.

Then before that usually NSS provides on the second this one is your household information. Before that the very basic information about the recoding of the different record its FSU units, Second stage stratum units etc.

Then now you might be confused that why it is so less? It is because very less number of responses must have been there since this is on deaths. Information about deaths i.e., those who have died in last 365 days and if the family member reported it. So, accordingly the data reported is very less.

So, a quick check for you that once you open this README file you should download these as well. Keep your cursor over here on data layout, you will be redirected to download. Similarly, on each of these 13-txt file you will have to download all those details.

So, once again I am guiding those who have already done it earlier like to convert the unit level data ASCII format to useable form, for that we need these supporting files as well. So, supporting files are very-very essential. So, like schedule file is required. Schedule means what? I have already guided in one of the lectures is on schedule versus questionnaire. Schedule is a systematic and complete document of information.

It is not just the questionnaire; it is also including a structured format of questionnaire. The Readme file which I have just shown you and a data layout file which you can able to download from the link which gives the byte position of all the information or the all the variables. Then the Readme file basically gives raw data.

So, we need to download the raw data very clearly and you need to mark the common primary key. Common primary keys are essential because these holds in merging different blocks. So, different blocks of information why it is required because if you have simply here for example, 113823 are the household.

If you do not merge with the individual file, you cannot get the how complete information. Like from the household you will get household size, household religion and household caste. Then you have to mix it with or joint or merge the data of individual information like individual gender. Then there are so many individual identification like individual education.

From the household you can get the standard of living or the income, NSS usually give a expenditure information. But in different sub blocks also or blocks of information you will also get information about their expenditure on different heads as well.

(Refer Slide Time: 22:22)

			round he	alth d:	ata have
	unic_75 inc.		Touria ne	until ut	
sep	arate data file fo	or each level.			
	Name	Date modified	Туре	Size	
	R75250L01	25-09-2021 12:14	Text Document	16,007 KB	
	R75250L02	25-09-2021 12:15	Text Document	16,007 KB	
	R75250L03	25-09-2021 12:15	Text Document	78,097 KB	
	R75250L04	25-09-2021 12:16	Text Document	357 KB	
	R75250L05	25-09-2021 12:16	Text Document	13,209 KB	
	R75250L06	25-09-2021 12:16	Text Document	13,209 KB	
	R75250L07	25-09-2021 12:17	Text Document	13,209 KB	
	R75250L08	25-09-2021 12:17	Text Document	6,081 KB	
	R75250L09	25-09-2021 12:18	Text Document	6,081 KB	
	R75250L10	25-09-2021 12:18	Text Document	6,081 KB	
	R75250L11	25-09-2021 12:18	Text Document	6,014 KB	
	R75250L12	25-09-2021 12:19	Text Document	9,881 KB	
			Tout Discussion	4 527 VP	

After saying so, I think we need to move it to the respective files. Once you have downloaded the raw data which from the previous page I have already shown you, once you keep your cursor on each of the txt file you can able to download. These data files are in ASCII format and downloaded from Readme file of 75th round. And the NSS 75th round data have separate data files for each level. So, for each 13 level we have separate data files.

### (Refer Slide Time: 23:56)



So, how it looks like? I will show everything once again with the original data. But at this moment I am just showing you the sample of it snapshot of that raw data. While each row in the data set represents a single observation with varying information about each unit without extraction it is impossible to do any meaningful analysis on these data sets. So, without reading between these things it is very difficult to understand in fact.

(Refer Slide Time: 23:36)



### (Refer Slide Time: 23:38)



So, let me just guide you some other important aspects. The most important file that is key to your work or to your extraction is the Readme file and the layout file. The layout file guides how the data is arranged and in which column you get the variable and their information. This is how the layout file looks like. This is how it is visible.

So, like for example, you asked about individual file about the gender, their age, their marital status, their relation to head, their byte positions are important. You have to look at these bytes position like gender as we know that it would be either 1, 2 or 3. So, at maximum 1 byte space it will consume. So, byte space is 43 to 43; that means, on the data they have entered they have only given 1 digit space; that means, at maximum it can take value from 0 till 9.

So, then similarly others other entries you can able to read. We are now just going to guide you about common ID that is quite important for merging and we will explain you in detail.

(Refer Slide Time: 24:57)



Maybe in the next last we will clarify very clearly about common ID. At this moment we are trying to extract it. One of the tips is given here that we need to always read the layout file very thoroughly. Another important file is called schedule file that is called the questionnaire file. It describes the sequence of question upon which information is collected. It also helps in labelling the variables and the variable values.

(Refer Slide Time: 25:29)



Then we will guide you once you have done those basic check or the basic files next aspect for you to check with the dictionary file. You have to create a dictionary file if you are handling a bigger database, without dictionary file you can also able to extract data but what we are suggesting you that it will minimize your time and your result will be always be more correct if you follow an appropriate dictionary file. Then how to go for it for creating a dictionary file? There are different methods of extraction, I am going to discuss about it.

But, we can now open the data set and I will show you both together. I will guide you how to extract and which method we can apply. And I will then do simultaneous attempts to clarify all those details.

(Refer Slide Time: 26:33)



(Refer Slide Time: 26:39)

	K		<ul> <li>Variables</li> </ul>
iter commands have			🔨 Eitter vanablen beer
Command re		Copyright 1905-2017 StateCorp LLC	Name Label
There are no dems to show.	Statistics/vata Analysis	4905 Lakeway Drive	There are no densi to show.
	Special Edition	College Station, Texas 77045 USA 200-STATA-PC http://www.stata.com	
		979-696-4600 stataBstata.com	
		979-496-4601 (EAH)	
	In some time is supported to be		
	transfer in the former		
	Notes:		
	<ol> <li>Unicode is supported; see</li> <li>Maximum number of variable</li> </ol>	help unicode_advice. # 1# set to 5000; see help met maximut.	
			Properties
			Properties B + -
			Properties B = + + + Motore
			Properties B = + Voldate I torre Later
			Popular =
			Popperan A totale Total Line
			Populari di - Vadan San San San San San San San S
			Brayer for ■ Function Gauss
			Pagenia B     
			Bragentine ■ Future Count Co
			Reparting Sources Former Example Ex
			Reporter
			Pupperture Control
			Purport on ■
			Pupperture Pupperture torus
			Purport on The second
	Concurd		Pupping 

### (Refer Slide Time: 26:46)

ttwork							
Stata15							
SAGE_India_wave_2_dataset							
GAGE_India_wave_1_dataset							
NSS 75th Health							
(8 Drive (D.)							
Local Disk (F.)							
Local Disk (E)							
USB Drive (D.)							
Local Disk (C-)							
fideas							
Notares							
Masic							
Develoads							
recoments							
PERMIT P		] rough.do	25-09-2521 (09.4)	00 i de	1.03		
no vagens		eptel feal	28-09-2021 09-41	Test Document	6.1.8		
is ru		ngtel, do.do	28-09-2021 10:19	DO Har	4.83		
		estd	28-09-2021 09:26	Tort Decomment	14.02		
neDrive		Level of Block State	28-05-2221 09-14	Dia hie	11/188		
neDrive		Direct 12,2444	21-04-2021 11-04	COLUMN THE	13 52 1 68		
			28-29-2021 10.19	DOM FOR	10,291 10		
Lecture 10		Block1,2,Lifts	28-09-2021 10:01	251.1 m	11,751.08		
Foundation Day PP%		block 7(a,b) dta	28-09-2021 10-14	004 File	15,618 KB		
Dr. Y. C. Silvastav		Bick12,0,05,67.4%	28-89-282110.18	255. File	1,62,142.68		
12 Prof. GN Pillai		block 1,2,3,4,5 dta	28-09-2021 10:17	\$54 File	79,545 KB		
DOAA	1	block 1,2,3,4.dts	28-09-2021 10:05	DBA For	71,786.KB		
Pictures	1	Supplementry_Documents	27-10-2021 11:45	File folder			
Documents	1	Reports	27-10-2021 11-65	Für feider			
Develoads	1	extraction	27 10 2021 11 45	FileFalder			
Desktop	1	detfile	27-10-2021 11-45	File folder			
wick access		anci data	27-10-2021 11-45	File Folder			
		Nove	Date modified	lype	Size		

So, first of all, here is 75th round. We are also opening a Stata window here and we are also opening NSS 75th on health. And the PPT you can also open simultaneously will guide you all those thing side by side.

So, let me go quickly to the respective page which I have already said. So, first what I do? I will open and show you what is this data file, the important files which we have been discussing.

(Refer Slide Time: 27:09)



First of all, the supplementary document which I have said you have to open and read it very carefully.

(Refer Slide Time: 27:24)

Home Tools Data_layout_File RL. ×			🕐 🥼 Sign		
🖺 🕁 🕀 🖶 🔍	🕆 🕘 🔔 🕨 🔿 🕒	) · · · · · · · · · · · · · · · · · · ·	<i>C</i> _ 🖂		
0	Search 'Add Link'				
			Export PDF		
	Gover	nment of India	Edit PDF		
	Data Qualit	y Assurance Division	Create PDF		
	Nationa 164 Conal Lal Ti	l Statistics Office holun Bood, Kollioto 108	📮 Comment		
	164, Gopal Lal Thakur Koad, Kolkata-108. Phone No. 2577-1128				
	NSS 75 <sup>th</sup> Round				
	Final Multiplier-posted unit-level data for Schedule- 25.0 of NSS 75 <sup>th</sup> round				
			Adobe Sign		
A) Data for Sch.	25.0 (Social Consumpti	ion: Health).	🔏 Fill & Sign		
There are 12	data filas halanain	a to 12 different levels as nor levent	🔓 Send for Commer		
(detelow75, 250 V	tata mes belongin	g to 15 different levels as per layout	hore Tools		
(uatalay / 5_250.A	1.5).				
	File names	No. of Records			
	R75250L01.TXT	113823			
	R75250L02.TXT	113823			
			v		

(Refer Slide Time: 27:29)

☆ ⊕ ⊕ Q	🕆 🕘 🗆 🕨 🕙 🕤 🕑	114 · 🔂 🗇 🖉 🖉 🖉 🖉 🖓	<i>d</i> _ 🖂
☆ � ⊕ Q	©         ©         ©         ©         ©         ©         ©         R         7	m     k     T     D     k     k     b     0       113823     113823       555352       2537       93925       93925       93925       43240       43240       43240       42762       70258	Control File     C
	R75250L13.TXT	32257	
	Total	1342307	

README file I have just opened on the screen. This looks like the one which we have shown it on the PPT. They are the different text, the layout file is available here, once I take my cursor it gives me the thumb to click on it and you can download it. So, this is the one you can download. And these are all file, if you keep your cursor over here, it asks for downloading.

### (Refer Slide Time: 27:58)



### (Refer Slide Time: 28:08)

Home Tools	Data_layout_File RE ×			🕐 🥼 Sign
B ☆ ♠ €	<del>)</del> Q			de 🖂
	<ol> <li>Weights The weights</li> </ol>	(or multipliers) are given at the end of each record from 133th byte onw hts (multipliers) are Sub-sample-wise, details of which are as given belo	/ards. w:	Search 1Add Link
	(For descrip field staff)	tion of Sub-sample, please see Instructions Manual, NSS $75^{\rm th}$ Round	l, for	Edit PDF
	NSS, NS	C and Sub-sample-wise weights (all sub-round multipliers):		📮 Comment
	NSS = B	ytes 127-129 (3 bytes)		Combine Files
	NSC = B	ytes 130-132 (3 bytes)		Delete inset: extract and institle pages.
	MLT = E	ytes 133-142 (10 bytes, assumed two places of decimal)		Try now
	All recor	ds of a second stage stratum will have same weight figure.	E	Compress PDF
	8. Use of S	ub-sample-wise weights (all sub-round multipliers)	D,	🔏 Redact
	For generati	ng Sub-sample-wise estimates based on data of all sub-rounds taken tog	ether,	Protect
	either Sub-s	unple-1 FSU's or Sub-sample-2 FSU's are to be considered at one time.	Sub-	Adobe Sign
	sample cod	e is available in the data file at 26th byte (Please see layout of data	1 í.e.,	🔏 Fil & Sign
	datalay75_2	50.XLS).		🔓 Send for Comme.
	Apply final	weight for Sub-sample wise estimates as follows:		
	Final We	ight ⇒ MLT/100		Convert, edit and e-sign P forms & agreements
	I mai we	ight WE1/100		Free 7-Day Trial

### (Refer Slide Time: 28:18)

me Tools Data_lay	out_Fie RE ×		🕑 🥼 Sign
☆ ⊕ <b>⊖</b> Q	(a) (b) (c) (c) (c) (c) (c) (c) (c) (c) (c) (c		<i>8</i> 🖬
F	or generating Sub-sample-combin- ogether, both Sub-sample-1 FSU's a	ed estimates based on data of all sub-rounds taken ind Sub-sample-2 FSU's are to be considered.	Search Add Link
A	apply final weight for Sub-sample c	ombined estimates as follows:	Edit PDF
	Final weight = MLT/100, if NSS	S=NSC	🔀 Create PDF 🗸 🗸
	= MLT/200 otherw	vise	📮 Comment
0		•	Combine Files
9	. Common Primary Key for identif	fication of a record for any schedule is:	📶 Organize Pages 🖌
	FSU Serial Number	= $4(5)$ (i.e., offset = 4th byte, length = 5 bytes)	Delete, insert, extract and rotate pages.
	Segment Number	= 31(1)	Try now
	Second Stage Stratum Number	= 32(1)	<ul> <li>Compress PDF</li> </ul>
	household Number	= 33(2)	🔏 Redact
	Level Number	= 35(2)	O Protect
	Person sho/hosnital sho/ailmen	st slng = 27(5)	Adobe Sign
	r erson sino/nospital sino/armien	-57(5)	🔏 Fill & Sign
1	0.List of Documents		Send for Comme.
	a) General Information	README75_250.doc	·
	b) Text Data Layout	datalay75_250.XLS	Convert, edit and e-ugo l forms & agreements
	c) Blank schedule 0.0	sch0.0.pdf	Free 7-Gay Stal

Then other important aspect you must read that is, they have guided about multiplier. A multiplier in this case is considered to be the weight for the analysis; how the weight is defined what do you mean by final weight. And they have guided you about the final weight here they multiplied divided by 100 if NSS and NSC are equal otherwise it is multiplied divided by 200.

And what is this multiplier file is already defined on the data I will show it on the layout file. The next important information here is to know about the common primary ID. I will also emphasize this several times. The common primary ID as per the suggestions in the data are consisting of these 6 important indicators

FSU, then segment number then second stage stratum number, household number, level number then person ID or element number serial number. So, if you are working of on the personal ID or on the individual file then you may be concerned for this particular position.

### (Refer Slide Time: 29:17)

e Tools Data_layout_Fil	R_ ×		🕐 🌲 Sign
🕆 🕀 🖯 🍳	🔁 👶 3 / 3 🖡 🕙 🕞 💓 me		<i>C</i>
10. <u>I</u>	ist of Documents		Search 'Add Link'
a	) General Information REA	DME75_250.doc	💫 Export PDF 🗸
t	) Text Data Layout datal	ay75 250.XLS	Edit PDF
c	) Blank schedule 0.0 sch	D.0.pdf	Create PDF
	) Blank schedule 25.0	5.0 pdf	Comment
· · · ·	) Estimation many lum acts for 75th and	Estimation Drassland NSS 75 - 46	Combine Files
c	) Estimation procedure note for 75 <sup>th</sup> round	Estimation_Procedure_N88_/5.pdf	C) Organze rages 7
f	) Instructions Manual for field staff, NSS 7	5 <sup>th</sup> Round.	Delete, insert, extract and rotate pages.
	*******	***	Try now
			Compress PDF
			🔏 Redact
			Protect
			Adobe Sign
			🔏 Fill & Sign
			Send for Comme.
			Convert, edit and e-sign P forms & agreements
			Free T-Day Inal

Then other important documents they have also suggested us to read. So, we are going to read all those things like your schedule, your layout file etc.

(Refer Slide Time: 29:48)

hataslave 💓 😸 🕬 - 🖓 - 😵 - P		deniey 15,2	50 (I) - Cen	pathility Mo	<b>4</b> •	P Search										Ramesh Arbo	unten 🖡		-	0
le Home Inset Draw Page Layou	formulas	Data	Review	View III	fφ													d She	e 🖸	Comme
A Cut	- A' A'		p. 11	Was list	General		Nor	nal She	iomal	Rad	600	đ		-	2 🖽	∑ Autabar	- 60	0	(0)	
Copy -		- =	·			in a factor	Ligo									🚺 fill -	2.4	200	1225	
- 🗳 Format Painter 🛛 🖁 🦉 🗉 🐇 - 🔮	2 · 4 ·	5 5 3	2 2 B	Merge & Cer	ter - 188 - % 🔊	Farmating	* Table *	0.51	accuration	uzatta	υp					🖗 Clear -	Filter -	Select -		
Cipboard 5 Fort	5		Alignment		G Number	5			Shies					Cel			Látina		Analysis	
	D																			
4 × 1 × 4																				
A B	с	D E	F	GH	1 J.		(   L		м	N	0	P	Q	R	8	T	U	٧	W	X
Text I	Data Layou	ut																		
N88 75th Reuni	(July 2017-Je	une 2018)																		
		Total	no.of level	s = 13																
Sch. 25.0 : LEVEL - 01(Blocks 1 and	2)	Recor	nd Length -	142+1																
srt. Item	Schedule	relecence	Leagth	Byte pos	tion Remarks															
1 Center Rand	DIRKS	treas (Col.	-	1	1 Centruit															
2 FUI Smid No	- i	1			2 Generates															
1 Road	- i				10 '75' General	4														
4 Schedule	·	÷.		. n	11 *250* Genera	ed														
5 Sample				14 -	14 Generated from Sc	h00														
6 Sector				15 -	15 Generated from Sci	h00														
7 NSS-Region				16 -	18 Generated from Sci	h00														
8 District				19 -	20 Generated from Sci	h00														
9 Stratum				21 -	22 Generated from Sc	h00														
10 Sub-stratum				23 -	24 Generated from Sc	h0.0														
11 Sub-Round				- 25 -	25 Generated from Sc	h00														
12 Sub-sample				26,-	26 Generated from Sc	h0.0														
13 FOD-Sub-Region				27 -	30 Generated from Sc	h0.0														
14 Hamlet group/ Sub-block no	1	4		31 -	31															
15 Second-stage-stratum no.	- h	2		32 -	32															
10 Sample Mail No	· · · · · · · · ·	· •		33 -	34 "Common-It															
12 Eða				37	4 20000 centra	ted.														
19 Sine of information in cel 1 block (i)	i i	7		42.	43															
30 Response Code	- ÷			44 .	44															
21 Survey Code	·	9		45 -	45															
22 Substitution Code/ Casualty code	i i	10		46 -	45															
23 Employee code	2.1	(a).(a)	3 4	47 -	50															
24 Employee code	2.1	(a)(b)	4 4	51 -	54															
25 Employee code	2.1	(t).(d)	3	55 -	58															
26 Date of Sarvey	2	2(6)	3	59 -	64 "DD MM 17															
27 Date of Despatch	. 31	2(iv)	4 4	65 -	TO MAL DO															
28 Time to canvass (minutes)	2	÷		71 -	13															
77 Ne. at givestigators (FT ASO) in the team		3		/4 -	74															
24 Kemarks in block 12/13	2	0(i)	,	13 -	13															
isyout75_25.0 (+)																				
1											_					5		Ð	+	· ·
																				138

I am going to open the layout file then just to show you the exact file and for your understanding. So, this is in Excel format and this gives information about byte positions and the column wise position is given. I have just opened it for your reference.



The first level is going to give us about the server related information its records. Like records on sampling related information, digit related information, round related information, district related information. So, this one is actually help us to understand the primary key.

From second onwards it is very important. Second level is for the household. Here, household type, then religion then social group then household size etc. and their byte positions are given. Byte positions like it is here, suppose you are saying household type, it has the byte position from 53 to 53, then religion it has the position 54 to 54.

So, these are important. I will also guide you about schedule. So, as I already told you during the understanding about data schedule 25 is meant for health care in all the rounds they have defined schedule 25 for health care.

# (Refer Slide Time: 31:02)



# (Refer Slide Time: 31:11)

ome Tools	Data_layout_File RE	Sch_25.0.pdf x						🕐 🥼 Sig
0 ☆ <b>@</b> €	Ð Q	(†) (i) 1 / 16	k 🕙 🖸 🟵	1945. +	<b>∀</b> . 1	₽ L & O O O		da 🖂
	3. sub-	district/tehsil/town:"		7. nam	e of head	of household:		<u>î</u>
	4. villa	ge name:		8. nam	e of infor	mant:		2
	[1] ide	ntification of sample household						
	item no.	item	co	de	item no.	item	code	
	1.	srl. no. of sample FSU			6.	sample household number		
	2.	round number	7	5	7.	serial number of informant (as in column 1 of block 4)		
	3.	schedule number	2	5 0	8.	response code		ľ
	4.	sample hg/sb number			9.	survey code		
	5.	second-stage stratum number			10.	reason for substitution of original household (code)		
	CODE	S FOR BLOCK 1						
	item 8:	response code: informant: co-op others -9.	perative and cap	able -1, c	o-operat	ive but not capable -2, busy -3, re	eluctant -4,	
	item 9:	survey code: original -1, substitu	ite -2, casualty -	-3.				
	item 10	reason for substitution of origi	nal household:	informan	t busy - L	members away from home -2. in	nformant	

# (Refer Slide Time: 31:13)

	Data_layout_File RE							0 #	Sigr
☆ ଚ 🖶	Q	(†) (k) 1 / 16	• 🕘 G	) ()	1945. +	<b>∀</b> . 🕁	<i>₽ ℓ &amp; </i>	P.	
	2.	round number	7		5	7.	(as in column 1 of block 4)		^
	3.	schedule number	2	5	0	8.	response code		2
	4.	sample hg/sb number				9.	survey code		
	5	second-stage stratum number				10	reason for substitution of		
							original household (code)		
		others -9.					,,		
	item 9: item 16	others -9. survey code: original -1, substitute ): reason for substitution of original non-cooperative -3, others -9.	2, casua househ	ilty -3. old: in	forman	• t busy -1	l, members away from home -2, informant		
	item 9: item 16 * tick n	others -9. survey code: original -1, substitute -, : reason for substitution of original non-cooperative -3, others -9. nark ( $$ ) may be put in the appropriat	2, casua <b>househ</b> te place	ilty -3. old: in	forman	t busy -1	l, members away from home -2, informant		4
	item 9: item 16 * tick n	others -9. survey code: original -1, substitute - : reason for substitution of original non-cooperative -3, others -9. mark ( $\vec{v}$ ) may be put in the appropriat	2, casua <b>househ</b> te place	ilty -3.	forman	t busy -1	l, members away from home -2, informant		
	item 9: item 16 * tick n	others -9. survey code: original -1, substitute - ): reason for substitution of original non-cooperative -3, others -9. nark ( $$ ) may be put in the appropriat	2, casua <b>househ</b> te place	lty -3.	forman	t busy -1	l, members away from home -2, informant		
	item 9: item 16 * tick n	others -9. survey code: original -1, substitute - Preason for substitution of original non-cooperative -3, others -9. mark ( $$ ) may be put in the appropriat	2, casua <b>househ</b> te place	ilty -3.	forman	t busy -1	l, members away from home -2, informant		

Then this is the questionnaire or this is the schedule where you exactly going to get it. On the very first identification related information available in our first layout or the and the layout on the first block of information.

(Refer Slide Time: 31:22)

Home Tools	Data_layout_File	RE Sch_25.0.pdf ×									0	8	Sig	n In
🖺 🖒 🖗	0 Q	(b) 2 / 16 k 🕙 🖸 G	) 1945 ·	<u>ы</u> . (	Ţ	1 6 1	0 0 (					$\partial_{\mathbf{a}}$		0
	[2] p	articulars of field operations											î	(
	sl. no.	item	Fie Juni	ld In ior S	vestigat atistical (JSO)	or (FI) / Officer	Fiel Sta	d Offic itistical	er (FO Office	)/ Senior er (SSO)			2	8
	(1)	(2)			(3)				(4)					ľ
	1.(a)	(i) name (block letters)												(
		(ii) code	-	Γ	Τ									
		(iii) signature												Ľ
	1.(b)	(i) name (block letters)												(
		(ii) code											•	1
		(iii) signature												4
	2.	date(s) of:	DD		MM	YY	DD	N	ſM	YY				C
		(i) survey/ inspection												1
		(ii) receipt												
		(iii) scrutiny					П							
		(iv) despatch												
	3.	number of additional sheet(s) attached												
	4.	total time taken to canvass the schedule by the						_					v	ŀ

### (Refer Slide Time: 31:24)

A 🖶 Q		1	2 / 16   🖡 🕘 🖸 🕖 1945 - 🚼 🖓 🗮	1/2000 20	3
Ĺ					
	4.	total time taken to canvas	the schedule by the		
		team of investigators (FI/	50)		
		(in minutes) [no decimal ]	oint]		ł
	5.	number of investigators (I	/JSO) in the team who		
		canvassed the schedule			
	6.	whether any remark has	(i) in block 12:13		
		been entered by	(ii) elsewhere in the		
		officer (ver. 1, no. 2)	(ii) eisennere in nie		
		onicer (jes-1: 10-2)	Schedule		
_					
1	12] r	emarks by investigator (l	I/JSO)		
L.	131.2	ommonte by supervisory	(ficor(s)		
1	1914	ounieurs by supervisory	uncer(s)		

From the second one it is all about the field operation from the third it is level third it is your household information.

(Refer Slide Time: 31:30)

D & 9	) 🖶 Q	(*) (*) (*)	k 👌 🖸	) () 1945		2.		6
							٨	6
	[3] household c	haracteristics						
	1. household siz	e	8. t	type of la	trine usually used (code)			Ц
	<ol> <li>whether the share for child non-household f last 365 days? (y)</li> </ol>	household paid major birth expenses for any emale member(s) during res-1, no-2)	if o ite	code in tem 8 is 01-09	9. access to latrine: exclusive use-1, common use of households in the building-2, public/community latrine-3, others-9		4	
	3. principal	description:			10. how many members use the latrine?			ľ
	(NIC-2008)		11.	. major sc	urce of drinking water (code)			a
	(1110-2000)	code (5-digit)	12	arranger	nent of garbage disposal (code)			0
	<ol> <li>principal occupation (NCO-2004)</li> </ol>	description:	13.	. primary 30 days	source of energy for cooking during the last (code)		4	0
		code (3-digit)	14.	was then (see list one hous (yes-1, )	e a sudden outbreak of communicable disease * below) in the community afflicting at least sehold member during last 365 days? 10-2)			
	5. household ty	pe (code)	15.	amount	of medical insurance premium paid for ld members during last 365 days (Rs.)			.0
	6. religion (code	2)	16.	househo	old's usual monthly consumer expenditure			
	7. social group	code)		(Rs.)				

Household size, household type, religion etc. these are mentioned. Along with that this also gives information about their principal status, working in industries or their occupation, principal occupation information as well.

## (Refer Slide Time: 31:53)



Next to this is your individual file like demographic particulars of household members. So, household member and their marital status, their age etc. is given.

(Refer Slide Time: 32:01)

φ φ			and the state of t					0	÷	sigi
	⊖ ପ୍		•	3 /	8 k 🕙 🕞 🥹 1305 - 🕍 - 🕎	¢ l	á 🔁 🗊 O		da.	^
	CODES	FOR BLOCK	14							
	col. 3:	relation to lawimothe	head: self - 1, sponse of er-in-law - 7, brother/sister/	head • brothe	2, married child - 3, spouse of married ch r-in-law/sister-in-law/other relatives - 8, serve	ild = 4, um int/employe	narried child - 5, grandchild - 6, father/mother/father es/other non-relatives - 9	-in-		
	col 4	gender: ma	le-1, female-2, transgender-	3						
	col. 6:	marital stat	us: never married - 1, curre	ently m	arried - 2, widowed - 3, divorced/separated-	1				
	col. 7:	general edi	ication: not literate -01, literate without any s	choolii	ng -02.		as it at			6
			literate without formal so literate with formal so	al scho choolii	oling: firongh NFEC -03, literate firongh . ig: below primary -06, primary -07, upper p diploma /certificate course (upto secondary diploma/certificate course(graduation & ab	rimary/mid vimary/mid )+12, diplom ove) +14, gr	-04, others -02; dle -08, secondary -10, higher secondary -11, sa'certificate coursefhigher secondary)-13, raduate -15, post graduate and above -16			
	col. 8: u	sual principal	activity status:							
	worked i own acco	n h.h. enterpri onnt worker	ise (self-employed):	-11	worked as casual wage labour: in public wor	ks -41	attended domestic duties and was also engaged in free collection of goods (vegetables, roots,	-93		4
	worked employed	in h.h. e	nterprise (self-employed)	-12	worked as casual wage labour: in other types of work	-51	jirewood, cattle jeed, etc.), sewing, tailoring, weaving, etc. for household use			
	worked a (unpaid)	as helper in h.) family workers	h. enterprise )	-21	did not work but was seeking and/or availab work	ble for -81	rentiers, pensioners, remittance recipients, etc.	-94		
	worked a	s regular sala	wied/wage employee	-31	attended educational institution	-91	not able to work due to disability	-95		
					attended domestic duties only	-92	others (including begging, prostitution, etc.)	.97		
	col 12:	whether hou	sehold paid major share for	childl	birth expenses: yes-1, no-2, pregnancy contin	ming-3				
	col. 13:	whether suffe suffered from syn	ered from any communicab Malaria-1, Viral Hepatitis drome-7, others -9 (Typhoic 8	<b>He dise</b> Jaund 1, Hooi	<b>ase :</b> ice-2, Acute Diarrhoeal Diseases/Dysentery-3 kworm Infection, Filariasis, Tuberculosis etc.)	Dengue fer	er – 4 Chikungunya-5, Measles-6, Acute Encephalitis			
	col. 17:	whether cove	o red by any scheme for heal	th exp	enditure support: government sponsored (e.g.	RSBY. Aro	ovasri, etc.)+1. oovernment/PSU as an employer (e.o.			

### (Refer Slide Time: 32:15)



And their different leveling, these codes are given, how these codes are entered in the text file is understood very clearly. Then your important aspect is to read, if you check with the 5th block here as per the schedule or in 4th block in the text file you have the information about the person who died during the last 365 days.

(Refer Slide Time: 32:36)

ome Tools Data_layou	ie RE Sch_25.0.pdf ×							0	8	Sign In
j ☆ � 🖶 Q	•	8 / 16 🖡 🕙 🖸 🕢 130% -	<b>K</b> - Ţ	₽ 2 6	200	G			la 1	3
						Schedula	25.0: 8			° (
	[7] expenses incurred durin	g like last 365 days for treatment of	members as	in-patient of	medical inst	ifution	4			
	1. sil ab of the inspital	and a case ins in neight offer of			,	•	,			
	2. 10.30 00 denotes 100	products are in item 2, block o)								
	3. JANN CVIENTION (DVI, INT) INVED	3. block (t)								
	<ol> <li>A. Duffner any medical A item main facility -1</li> </ol>	arvice provided free (fally partly) n/a third - Chicandile/MGO Timer								
	we happital + 2 both	1. 10 -1)								
	expenditure for torational	during stay at hospital (in whole ou	mber of Rs.)							
	5. packagi sothpohent i	Re.)								
	non juritage component ()	k.)	E Arreni kingi		in the second second	konstand dari	ingé			0
	6. docioir y purgeou's fe	e (hospital staff" other specialists)								٩
	2. menticines									
	8. dasgmöblic inns-	_								
	<ol> <li>liệd draryst</li> </ol>									
	10 other medical intering	ne (mtehatani charges, physio-								
	thempy, permual man	tical ppliances blood organi,								
	etc. i									
	11 medical expenditure	(Rs.): total (Items 5-10)								
	12 thumper La patanné	its.)								
	13. office mat-medical wa	penais mentual by the licensition					-			
	(rapistration fee, fee	d transport for others expenditure								
	on escart lostging cl	arges 0 ouv_atc.) (Rs.)								
	14. expenditure (Rs.): to	4a) (deuts (1-13)								¥.

### (Refer Slide Time: 32:38)

ie RE	0 #	Sign In
⑦ ◎ I == ▶ ⑦ ○ ○ == - K ♥ ■ ℓ ℓ ₺ □ ∩	0a	⊠ 0
15. total amount reimbursed by medical insurance company or employer (Rs.)		Î Q
16. major source of finance for expenses (code)		D.
<ol> <li>2<sup>nd</sup> most important source of finance for expenses (code)</li> </ol>		8
18. place of hospitalisation (code)		
19. If code is 5 in item 18, then state code (page 15)		
20. loss of household income, if any, due to hospitalisation (Rs.)		
CODES FOR BLOCK 7		1
item 16 &1 <sup>+</sup> : source of finance for expenses: household income/savings -1 sale of physical assets -3		۵
borrowings -2 contributions from friends and relatives -4 other sources -9		4
Item 18: place of hospitalisation:		, 0
same district (rural area) -1 volum state afferent district (rural area) -5 some district (when area) -2 within state different district (urban area) -4		Ð
other state -5		a
		0
		~
	With the state of longet for expenses (cole)         10       If code is Statication (code)         11       2 <sup>10</sup> is marked in momente company         12       2 <sup>10</sup> is marked in momente company         13       If code is Statication (code)         14       2 <sup>10</sup> is more of finance for expenses (code)         15       16         16       17         17       2 <sup>10</sup> is more of finance for expenses (code)         18       place of longetalisation (code)         19       16         10       17         10       18         11       2 <sup>10</sup> is so of homebold momene, if my, due to hospitalisation         10       16         10       17         11       2 <sup>10</sup> is so of homebold momene, if my, due to hospitalisation         11       2 <sup>10</sup> is so of homebold momene, if my, due to hospitalisation         10       16         11       17< is more of home for expense:	Image:

Similarly, next to this all information are related to individual file. I think we have done all those basic checks then we need to understand about ASCII data. ASCII data which I have already clarified, these are entered in the text file in the notepad file.

(Refer Slide Time: 33:01)



### (Refer Slide Time: 33:18)



So, ASCII data like if I just open one file over here let it be these two. The sample which we have shown it is also visible here. How can you read it? Suppose we say that till the end, suppose I just keep my cursor over here, it will read your position. Your line number and it is column number can be shown at the bottom of this particular page.

Similarly, all its positions can be well read, like this is your first position, this is your second position, this is your third position etc. So, there are some spaces in given. Spaces you need not worry about it. In the data itself they have specified about how many positions the spaces have taken. So, we need not worry.

We need to read between the position byte position specified in the layout file. So, then the next aspect after guiding all those things I will go to your PPT once again. I have shown all those details. We need to create a dictionary file. So, now, we are handling the fixed format data, not free format, not delimited format.

We are dealing with the fixed form data where the fixed position is actually defined. So, its byte position is well defined in the layout file. Now, how we can create a dictionary file? There are different approaches to do it. But is it really very essential? You may not even require it, there are three methods where we can extract the data. One of the approaches is through the dct file.

(Refer Slide Time: 34:52)

I Save		×
← → × ↑ 🦲 « Desktop → NSS 75th Health v 💍	,O Search NSS 75th	Health
Organize 👻 New folder	8	:: • 🕜
Dropbox	Date modified	Туре
articles on medi     Level_wise_data_in_ASCII_FORMAT     detailed content     Reports	25-09-2021 12:22 25-09-2021 11:27	File folder File folder
Screenshots Supplementry_Documents OneDrive	25-09-2021 11:27	File folder
Inis PC		
Desktop v K		>
File name: LEVEL - 03 (Block 4) Save as type: Octionary (*.dct)		~

(Refer Slide Time: 34:58)

Three Methods of Extracting the Fixed Format Data	
□ First Method Launch stata -> File -> Import -> Text data in fixed format -> Browse dictionary file (.dct) if you have created / else click on specifications -> type your chosen variable name along with storage type and column positions -> browse your dataset file name as we have not specified it in the dictionary file -> check on replace data in memory -> submit/okay	
õ swayati 🔮 24	

Let us start with a very basic of extraction. From the click-based approach we can able to extract the data as well. This is basically called the first method. I have also shown you some previous pages those who are already discussed. I am going to come back to the dictionary file once again.

Let me just stick to these first. Launching a Stata then we will go to File then Import then we will specify the dct file if it is defined. Otherwise, you can also specify the fixed format

information then we can download. Like I can show it over here in your Stata window. Here we will go to File.

(Refer Slide Time: 35:56)

10:         10: <th><b>T U</b> X</th>	<b>T U</b> X
(a) (a) (b) (a) (a) (b) (a) (b) (a) (b) (b) (b) (b) (b) (b) (b) (b) (b) (b	T U X
Bit	<b>T 2</b> X
International System         International System         Operating	Dev.
Name         Image: State         State         Segments         Segments <t< td=""><td>244</td></t<>	244
Dec.	Der.
Team         Operating Work         Operating Work <td>DW.</td>	DW.
Operanding dentry         Special Edition         Collapse frames, Trans 1918 Edit           log         +         600 The First Not 1918 First	
Ling         00-15221-87         MASS/1/WW-MASS-COM           Model         10-15221-87         MASS/1/WW-MASS-COM           Model         10-15221-87         MASS/1/WW-MASS-COM	
Import         Finity predicted ("Ait" also)         933-644-6400         #Example field (Finity also)           Finity predicted ("Ait" also)         933-6464-6400         (Finity also)         (Finity also)	
979-696-4602 (24A)	
Lyport · bet data idelemented "cm/ )	
Even     A     A     A     A     A     A	
Francisco detrouter Nort data or found tormat with a decisionary 814	
to represent states to the states of the sta	
Report Nas * 553.07001	
List Indext Reserve Concernic Data (FED)	
Hour Analytic database fee help unicode_sdrice.	
OBC data source bles is set to 5000 see help set maximum.	
diaxe (".dd)	
Projekt           B           B           C           C           C           C           C           C           C           C           C           C           C           C           C<	•*
Solar P	
Commond I	
Cliber/Sube Albeitsplanway 2020: Watay Menery/Statis	JUM OIR
11 / lypehere to starch O 😫 🎯 🐔 🕐 L MS 700. 🖗 MH U.L. 🔮 Province. 🙀 Stark Ko. 🔮 Societies. 🔮 Societies. 🔮 Societies. 🔮 Societies. 🔮 Societies. 🔮 Societies.	5

This is the first method then we will Import, but import from where? Import from the data in fixed format. I have said that our data is loaded with fixed format information. If it is on csv format or delimited format, then you could have been clicked here. If it is in an Excel format, but some specification needs to be very careful about it. Since it is in fixed format you have to simply click here.

(Refer Slide Time: 36:21)



So, with this click it asks for a dictionary file. So, if your dictionary file has already been installed or you have already defined then you can open it. I will guide you how to open it and then you can operate it. Otherwise, you have to give the specifications. What do you mean by specification? You have to give the command which we are going to operate in our next page.

So, once we submit and end the source of the data. Here browsing means you are going to specify where the data we wanted to import, and your data has to be a micro data file.

(Refer Slide Time: 36:58)



So, micro data file as I already told you it is ASCII data. So, here all files then all text files is visible and you can open it the respective page. So, let me just go back to our PPT once again and I will guide you rest of the details. This is one method.

(Refer Slide Time: 37:20)

Main if/in	
Use dictionary file:     Browse	
○ Specifications:	
Examples: rate 1-4 speed 6-7 acc 9-11 2 lines 1: id 1-6 str name 7-36 2: age 1-2 sex 4	
Text dataset filename: (required if not specified in dictionary) Browse	
Replace data in memory	
🖉 🚯 🗋 OK Cancel Submit	

This is how it looks like.

(Refer Slide Time: 37:21)

Second Method
Use infix command without data dictionary
If an ASCII file has a few variables in a simple format, you just need to list the variables' names, types, and column ranges.
infix str FSU 4-8 str Segment 31-31 str Gender 43-43 using"C:\Users\admin\Desktop\NSS75thHealth\ascii_data\R75250L03.tx t", clear
Stata reads string variable FSU from column 4 through 8; string variable Segment from column 3 through 4 and so on.
💿 swajan 🔮 26

Then second method where we can go with the manual command. So, manual command on the command window or command space. So, we will have to give the infix command. infix command and its position. So, I am just going to show it to you here.

So, what we will do? We will first do couple of check like we can type the path name the way we have defined here. For the time being we are just going to open the command we have

saved it for your quick reference. You can easily type infix then using command, but we are going to show it how this looks like and how you can operate.

(Refer Slide Time: 38:21)

(Refer Slide Time: 38:33)



### (Refer Slide Time: 38:40)



So, basically, we are guiding you on the path names we have already saved it for our reference. So, I am just going to show it over here. So, here it is and I will guide you what exactly you are supposed to write down. So, over here it you have to write down infix. Then using then on the using first party you are creating a dictionary file.

So, if you have already created a dictionary file then you could have returned the with the path, I will tell you what all about the dictionary file is. Then after comma it is your using file and them you have to browse raw data. So, the raw data should be ended with a dot txt.

And another sensitive aspect is that it should start with a inverted comma bracket and that has to be closed at the end. And another one is the backward slash has also to be given very correctly and the file path name has to be specified. Once we do it you can able to import the data.

### (Refer Slide Time: 39:46)



We can just click on this and data will be extracted. Now, this is what is one of the methods where we have used the dct file. But if you do not use the dct file, you can simply type the command with space. Command with a byte position like infix you type infix another method which is manual method like infix we have to clear these first.

So, then clear the stored data then we type one by one; one by one you can just see these byte position in the layout file. Then you right in the command box like what it is written is infix. Then before using do remember that before using you have to mention the type of the data we are deliberately writing *str* stands for string.

We are trying to make the data extracted in string format because that helps in not distributing the data that is compressing the data in string format and that further we can get the numeric form for further analysis. So, str for string we are making all the variables in string. At this moment we just wanted to show for two variables FSU and Segment and its position as I already told you FSU is from 4 to 8 that you can just check.

FSU is here 4 to 8. Position number 4 to 8 even it was guided on the serial number 2 or on the row number 10 of Excel sheet 4 to 8 is your FSU, then second one is your Segment is 31 to 31. So, this is what is given, 31 till 31 is your segment and that is also guided in the schedule file as well in the README file as well.

So, in the README file you have all such information. So, 31 to 31 is your segment number. Now after specifying this you can able to extract your data and make sure that you have used in after that using. If you have more variables then you can use all those variables over here.

So, this method is less prepared it this does not require a dictionary file. So, you have to do it manually. At the end you need to give the path name of this text file that is your raw data or the micro data dot txt has to be mentioned and inverted comma has to be closed. So, now, we can enter it and the data will be extracted.

(Refer Slide Time: 42:50)



(Refer Slide Time: 42:56)



So, you can just see whether data is extracted or not. Just browse it and you can get the information. So, now, in the browse since we extracted 2 variables that is FSU and Segment, so data with two variables have been extracted now. Now question arises here how to define a dictionary file, how to develop a dictionary file.

So, for that you do one thing. You simply write down all those things. In fixed dictionary on a do file you create a do file simply you click here create new and simply clear.

(Refer Slide Time: 43:42)



Let us open, since do file has already been opened, it is here you open another do file. And you just write down this way; infix like you simply on this screen infix dictionary and Stata is very case sensitive you have to write very correctly. Then on the next row the bracket has to be started.

Then on the next row all the variables and the another bracket has to be started. Then your variables names though variable name which is displayed on your screen like str you have to write down str then FSU. Where you are going to get this? You are going to get this from the layout file.

The layout file has given all those information like suppose you are extracting the individual file these things has to be written. Now first of all you should get the minimum information like these in the README file I already said you that your common primary information has

to be given FSU, Segment number, SSS number, household number, level number everything you should keep it first.

Keep it first, its byte positions are given 4 till 8 that is 4, 5, 6, 7, 8. Till 8 position you need to write down 4 to 8 on the first screen it is here. So, like str we have defined str FSU 4 to 8 then str Segment 31 to 31 and etc. Your all those variables will follow at the end you should remember that in a multiplier file should also be written very correctly.

Multiplier files has already been given in each entry serial number 24 here, on level number 3 multiplier 5 position number 132, to 142. So, that has to be also defined very clearly. So, this is 132, 142 is defined and at the end the bracket must be closed. If we can just copy it and paste it over here and, that is all about our dictionary file.

And you can manually do initially then later on that is what is our dictionary file. So, bracket has to be similar in both the case. The starting bracket and the end bracket has to be same. So, it has read that you have given the entries correctly.

(Refer Slide Time: 46:34)





How to save this? You need to go to Save As, but on the Save As, by default it takes do file you need to change it to dot dct and that if you are extracting the information about household then you write down household dct or level 3, level 2 dct. You have to write down the appropriate name here. Then you save it and at the time of the appropriate command you can give the command and do the extraction. So, rest it is easy, I have already guided you.

(Refer Slide Time: 47:05)



So, these are all I think I have said Stata will extract. In this case suppose another aspect like if you do not want entire information, you want to get very specific information from 1 till

100 observations then it has to be mentioned. Like infix string FSU it is everything you have defined.

And the string position number you have already defined. At the end if you want that you need to extract 1 upon 100 or starting till 100 observations if you have to specify this then your data will be extract till 1 to 100.

Generally, we do not suggest this because we always suggest that you, please download extract all the information from the beginning. When you are doubly sure about your data and your data is confined to 100 observations only then you can go for this approach.

(Refer Slide Time: 48:09)

Third Method	
Use infix command with data dictionary	
<ul> <li>If there are many variables in a complicated format, you can be from writing a data dictionary.</li> </ul>	enefit
This approach is highly recommended for the purpose of data management.	
swayani 🧐	28

Third method is highly recommended that is the dictionary method. And though initially it seems to be complicated, but actually it helps us in getting a better direction and bigger databases can be extracted in a every less time.

(Refer Slide Time: 48:27)



(Refer Slide Time: 48:43)



(Refer Slide Time: 48:48)

🖥 Save			×	
$\vdash$ $\rightarrow$ $\checkmark$ $\Uparrow$ es	ktop > NSS 75th Health v Ö		h Health	
Organize 👻 New folder		[	iii • 🕜	
Dropbox	Name	Date modified	Туре	
🧟 articles on medi-	ascii_data	25-09-2021 12:22	File folder	
👧 detailed content	Reports	25-09-2021 11:27	File folder	
Screenshots	Supplementry_Documents	25-09-2021 11:27	File folder	
OneDrive				
💻 This PC				
3D Objects				
Desktop 🗸	<		>	
File name: level 0.	3 (block_04)		~	
Save file as a Save as type: Stata D	ata (*.dta)		~	
Stata				
Data(*.dta)				
		Save	Cancel	
i) swayan 💮				31

So, dictionary file the way we have already defined and saved it I think it will be very useful to you. So, infix using a dictionary file is required once you have created your dictionary file and the path name of the data should be also given in the command. This is how it looks like, and we have shown it just a couple of minutes back and then you need to save it accordingly and it will be saved.

(Refer Slide Time: 48:56)

2 🖬 🛞	Dia 📝	i lools												
	FSU[1]		67776											
_	¥SU	Segment	\$55	Household	PersonID	Relation_t-d	Gender	λge	Marital_st-s	General_ed-n	UPA_status	hospitaliz-s n^	Variables	
1	67776	1	1	01	01	1	1	55	2	01	11		+ Filter variables h	
2	67776	1	1	01	02	2	2	50	2	01	92	2	R Name I	ahel
3	67776	1	1	01	03	3	1	24	2	13	61	2	IFI FSU	
4	67776	1	1	01	04	4	3	2.2	2	10	92	1	Segment	
5	67776	1	1	01	05	6	1		1	01	99	2	₽ \$55	
¢	67776	1	1	01	06	\$	2	16	4	10	91	2	Household	
7	67776	1	1	01	07	5	2	14	1	08	91	2	PersonID	
	67776	1	1	01	08	5	2		1	06	91	2	Relation_to_h	
9	67776	1	2	01	01	1	1	6.9	2	01	51	2	Gender	
10	67776	1	2	01	02	2	2	67	2	01	92	1	Mage	
11	67776	1	2	01	03	3	1	36	3	08	\$1	2	General educ	
1.2	67776	1	2	01	04	4	2	27	2	06	92	2	VPA status	
13	67776	1	2	01	05	6	1		1	06	91	2	Variables Snapshot	ä
14	67776	1	2	01	06	6	2	2	1	06	91	2	Descention	
15	67776	1	2	01	07	5	1	27	1	10	61	2	eroperties	
16	67776	1	2	01	08		1	2.6	1	10	61	2	Name	FSU
17	67776	1	2	01	09	5	1	23	1	10	61	2	Label	
10	67776	4	2	01	10	5	2	16	1	06	91	2	Type	strő
19	67776	1	2	02	01	1	1	40	2	01	51	2	Format	%95
20	67776	1	2	02	02	2	2	45	2	01	92	1	Value label	
21	67776	1	2	02	03	5	1	10	1	06	91	2	Notes	
22	67776	1	2	03	04	5	2	14	1	06	91	1	- Data	Jevel (31/hts
23	67776	1	2	02	05	5	2	14	1	06	91	2	Label	
24	67776	4	3	01	01	1	1	37	2	01	51	2	Notes	
25	67776	1	3	01	02	2	2	35	2	01	92	2	Variables	24
:												,	Observations	555,352
												,		

So, this is again repetition of my guidance. I have shown how these variables look like and how you can go for extraction alright. So, if you have any difficulties you may redirect to us,

or you may give all your queries in the chat box of your NPTEL path. And we will be very happy to address it and with this format you can able to extract all the database of NSS and there will be no difficulties.

But I suggest that you should do enough practice the way we have guided and rest all those supporting files with it is PPT with commands layouts or everything we have guided and you should make everything ready before extraction. So, these are all for today and on the next class we will be coming up with content for systematic guidance for you to merge the data sets. And then the next class onwards we will be ready for analysis of the data.

Thank you.