

# ***Using Stata on Long-Format Data***

Hsueh-Sheng Wu  
CFDR Workshop Series  
May 21, 2018

BGSU



Center for  
**Family and  
Demographic** Research

# Outline

- Long format versus wide format data
- Why do we need data in long format?
- Variable description
- Common Stata commands for working with data in long format
- Examples:
  - Cross-sectional family relation data
    - What is the size of the family?
    - Do householders have kids living at home?
    - What is the gender of the spouse of the household head?
  - Finding change in an variable and create discrete-time event history data
    - Using the information from the previous time point to substitute the missing value
    - Determine whether the respondents entered marriage
    - Determine how many times respondents transited into marriage
    - Determine the time when respondents entered the first marriage
    - Remove data that occurred after respondents entered the first marriage
- Conclusions

# Long Format versus Wide Format Data

Sample data in Long Format			Sample Data in Wide Format								
Family ID	pernum	sex	Family ID	sex1	sex2	sex3	sex4	sex5	sex6	sex7	sex8
9	1	2	9	2	2	.	.	.	.	.	.
9	2	2	8212	1	2	1	1	2	2	1	2
8212	1	1									
8212	2	2									
8212	3	1									
8212	4	1									
8212	5	2									
8212	6	2									
8212	7	1									
8212	8	2									

- Long format has 10 rows of data; wide format has two rows.
- Both formats have the case ID variable: Family ID
- Long format has an index variable: pernum
- Long format has one sex variable, but wide format has eight because the information of the index variable is incorporated into the sex variables
- Both formats provide same information for individuals
- The index variable can indicate the ID within family or for age/time
- The ID variable indicates a family or individual.

# Why Do We Need Data in Long Format?

- Data in long format preserve the original data structure, so it is easier to check the data and construct variables.
- Certain analyses (e.g., discrete-time survival analysis) require data in long format.
- Because Stata reads in the whole data set at the same time, Stata has special commands to construct variables for data in long format.
- For types of variable and data constructions:
  - Reshape data between wide format and long format
  - Generate summary statistics
  - Copy information from one row to another (e.g., the gender of the spouse for the household head)
  - Create a time interval capturing the change of a status (e.g., the duration for individuals to change from being single to being married)

# Variable Description

<b>Variable Descriptio of the Sample Data</b>			
<b>Definition</b>	<b>Name</b>	<b>Value</b>	<b>Value label</b>
<b>Survey Year</b>	<b>year</b>	<b>2016</b>	
<b>Survey month</b>	<b>month</b>	<b>1-12</b>	<b>January - December</b>
<b>Household serial number</b>	<b>fam_id</b>	<b>1- 94097</b>	
<b>Person number in sample unit</b>	<b>pernum</b>	<b>1-16</b>	
<b>Person ID</b>	<b>id</b>	<b>20141000000601 - 2016120746002</b>	
<b>Relationship to household head</b>	<b>relate</b>		
		<b>101</b>	<b>Head/headholder</b>
		<b>201</b>	<b>Spouse</b>
		<b>301</b>	<b>Child</b>
		<b>501</b>	<b>Parent</b>
		<b>701</b>	<b>Sibling</b>
		<b>901</b>	<b>Grandchild</b>
		<b>1001</b>	<b>Other relatives, n.s.</b>
		<b>1114</b>	<b>Unmarried partner</b>
		<b>1115</b>	<b>Housemate/roomate</b>
		<b>1241</b>	<b>Roomer/boarder/lodger</b>
		<b>1242</b>	<b>Foster children</b>
		<b>1260</b>	<b>Other nonrelatives</b>
<b>Age</b>	<b>age</b>		
		<b>0-85</b>	<b>0-85 years old</b>
<b>Gender</b>	<b>sex</b>		
		<b>1</b>	<b>Male</b>
		<b>2</b>	<b>Female</b>
<b>Marital Status</b>	<b>marst</b>		
		<b>1</b>	<b>Married, spouse present</b>
		<b>2</b>	<b>Married, spouse absent</b>
		<b>3</b>	<b>Separated</b>
		<b>4</b>	<b>Divorced</b>
		<b>5</b>	<b>Widowed</b>
		<b>6</b>	<b>Never married/single</b>
		<b>9</b>	<b>Not in a union</b>

# Common Stata Commands for Working with Data in Long Format

Command commands for Working with Data in Long Format	
<code>reshape wide year month marst sex age t_time, i(id) j(time)</code> <code>reshape long year month marst sex age t_time, i(id) j(time)</code>	Change data between the wide format and the long format
<code>duplicates report fam_id pernum</code>	Check whether each record is a unique one
<code>sort id time</code>	Arrange data in a special order
<code>by id: gen n=_n</code>	Create an indicator variable for each record of the person or family
<code>by id: gen N=_N</code>	Calculate the total number of records of the person or family
<code>gen marst_r = marst</code> <code>replace marst_r = marst[_n-1] if marst_r ==. &amp; marst_r[_n-1] ~=.</code>	Handling the missing value
<code>by id: replace c_mar = 1 if marst_r[_n-1] &gt;=4 &amp; (marst_r ==1   marst_r ==2)</code>	Code the transition into marriage
<code>by id: gen i_c_mar = sum(c_mar)</code>	Create an indicator for the time of enering marriage
<code>by id: egen s_c_mar = sum(c_mar)</code>	Create an indicator for number of times of entering marriage
<code>by id: gen time_mar1 = time if c_mar ==1 &amp; i_c_mar ==1</code>	Extract the time when the first marriage take place
<code>by id: egen m_time_mar1 = max(time_mar1)</code>	Expand the time of the first transition in marital status to all records of the individual
<code>by id: drop if time &gt; m_time_mar1</code>	Remove records that occurred after the first transition in marital status
<code>by id: replace s_sex =. if pernum ~=1</code>	Replace the values in irrelevant records

# Examples

- See accompanying Stata commands and log files



# Conclusions

- With data in long format, researchers usually create summary statistics, copy information from one row to another, or create a time interval capturing the change of a status.
- Stata commands greatly reduce the difficulty of accomplishing these tasks.
- When working with data in long format, be sure that there are no duplicate records per family or individual and that data are sorted correctly.
- You can use the `–reshape-` command to jump between long format and wide format to speed up the variable construction process.
- If you have any questions about long format data, please come see me at 5D, Williams Hall or send me an email ([wuh@bgsu.edu](mailto:wuh@bgsu.edu)).