

Data Preparation Using Stata

Hsueh-Sheng Wu
CFDR Workshop Series
September 18, 2023

BGSU

 Center for Family and Demographic Research

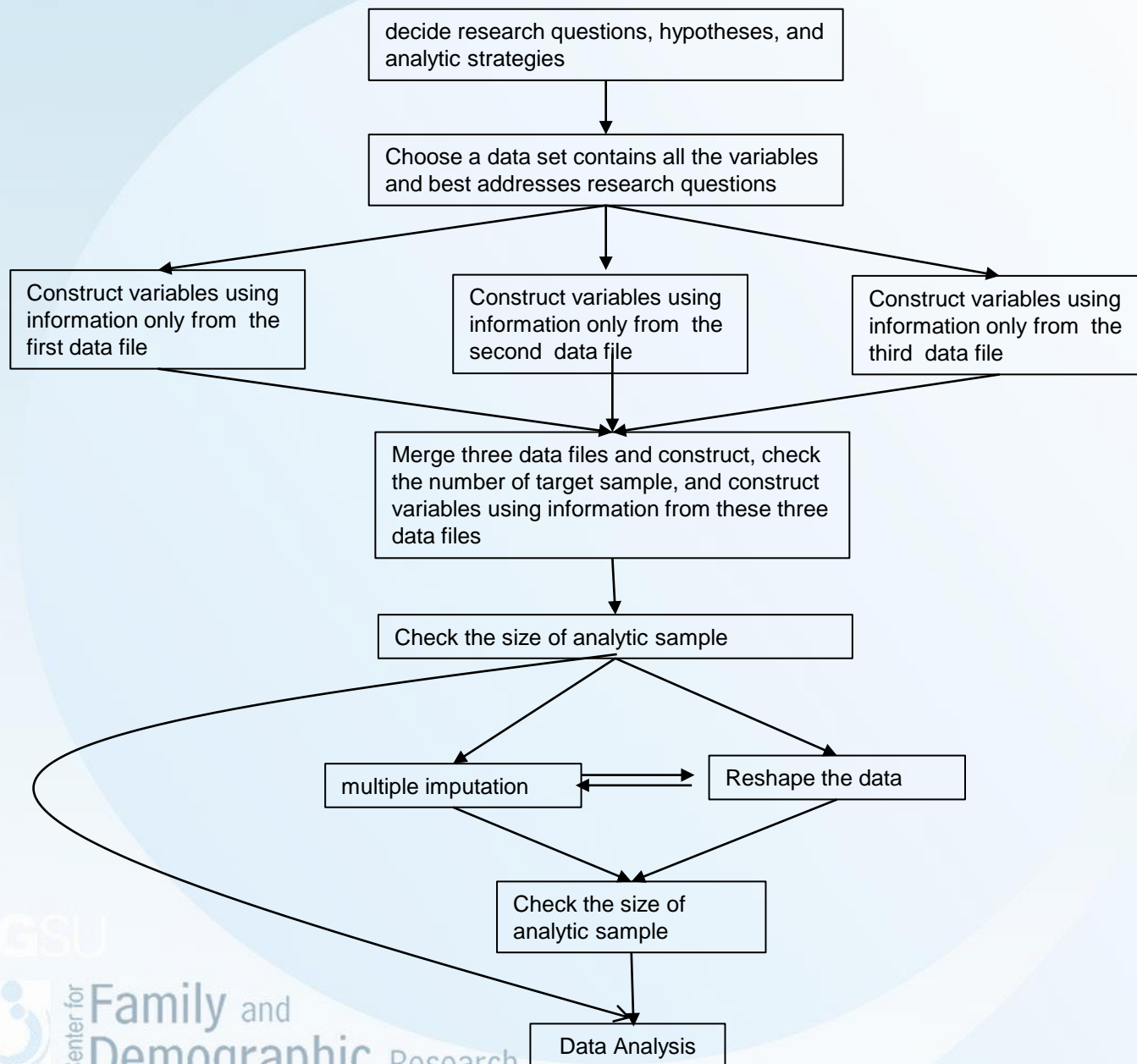
Outline

- Importance of data preparation in research
- A flowchart of data preparation
- Check list at each step of data preparation
 - Preparation works before data preparation
 - Decide on the data file
 - Check the accuracy of each data file
 - Construct variables
 - Merge all the data file
- An example of data preparation
- Conclusion

Importance of Data preparation

- Data preparation plays a critical role in social science research, but received little attention
- The major tasks of Data preparations are
 - Decide Research Questions, Hypotheses, and Analytic Strategies
 - Data Cleaning
 - Variable Transformation
 - Subset Creation
 - Data Documentation and reproduction
- Social Science research often involves lots of decisions on how to analyze imperfect data. Thus, good data preparation skills help social sciences researchers keep track of why and how they code and construct variables and make it easy to modify if necessary, variables.
- Reviewers of manuscripts often raises many methodological questions. Good data preparation skills help researchers justify their decisions on these methodological questions.

A Flowchart of Data preparation



Check list at each step of data preparation

Deciding Research Questions, Hypotheses, and Analytic Strategies

- Decide research questions
 - What research questions are?
 - Who the target population of these research questions is?
- Decide hypotheses
 - What research hypotheses are for each research question?
 - What variables are needed?
 - How these variables have been conceptualize and measured in the past?
- Analytic Strategies
 - What analyses are need for respective hypotheses?
 - What data format(s) are required for testing the hypotheses?

Check list at each step of data preparation (continued)

Decide on the data file

- Locate data files contain variables you need. For constructing complex variables, you need to make sure you have all the variables needed.
- Check out wordings and response categories of questions in codebooks to determine if these questions adequately measure each variable of interest.
- Examine the sampling frame and decide if the data file covers the target population
- Look at the sampling methods and decide if the analysis result should be weighted.
- Depending on research questions, other data attributes like data collection years, the number of waves may also need to be considered.
- Decide the data file or files to use
- Document what file each variable is from and identify possible inconsistencies across variables

Check list at each step of data preparation (continued)

Check the accuracy of each data file

- Make sure that the statistical software can read in data files without any error message
- Check if the number of observations in the data file match that of the codebook
- locate the ID variable or variables and check how many observations are associated with each ID value

Check list at each step of data preparation (continued)

Construct variables

- Check the attributes and/or frequencies of the ID variable, independent, dependent, control, and weight variables. (i.e., string or numerical; value labels, variable labels)
- Generate new variables from these variables
- Add value and variable labels to generated variables
- Recode or replace the values of new variables if needed
- Check if new variables are created correctly in each data file
- Construct variables and note the number of valid observation in the data file
- Create a new data set containing just variables needed for answering your research questions.

Check list at each step of data preparation (continued)

Merge all the data files into a new data file

- Merge data files together
- Identify the sample of target population
- Check whether variables that should not vary across data files actually do not vary across data files.
- Construct variables using different sources of data if necessary
- Identify the size of analytic sample
- If necessary, use multiple imputation to reduce the problem of missing data
- If necessary, reshape the data for analysis using data in long format
- If multiple imputation or reshape command are used, recheck the size of analytic sample.

An example of data preparation

- Research questions:

Whether race and gender have different impact on the trajectory of depressive mood between ages 18 and 45.

- Data: Add Health public data at Waves 3,4, and 5.
- Variables needed: age, gender, race, and depressed mood
- Statistical technique: hierarchical linear model
- Final data format: the long format.

An example of data preparation (continued)

Table 1. The variables and Data files to be used							
		W3		W4		W5	
		<u>original data</u>	<u>new data</u>	<u>original data</u>	<u>new data</u>	<u>original data</u>	<u>new data</u>
	data files	<u>21600-0008-Data.dta</u>	<u>w3.dta</u>	<u>21600-0022-Data</u>	<u>w4.dta</u>	<u>21600-0032-Data</u>	<u>w5.dta</u>
		original variable	new variable	original variable	new variable	original variable	new variable
	ID variable	aid		aid		aid	
variables to be constructed	available variables						
Age	Interview Year	iyear3	Interview_y3	iyear4	Interview_y4	iyear5	Interview_y5
	Interview Month	imonth3	Interview_m3	imonth4	Interview_m4	imonth5	Interview_m5
	Birth Year	h3od1y	birth_y3	h4od1y	birth_y4	h5od1y	birth_y5
	Birth Month	h3od1m	birth_m3	h4od1m	birth_m4	h5od1m	birth_m5
Gender	Gender	bio_sex3	female_w3	bio_sex4	female_w4	h5od2a	female_w5
Race	White	h3od4a	white_w3			h5od4a	white_w5
	Black	h3od4b	black_w3			h5od4b	black_w5
	Indians	h3od4c	indian_w3			h5od4f	indian_w5
	Asians/islander	h3od4d	asian_islander_w3			h5od4d	asian_w5
	asian only					h5od4e	islander_w5
	islander						
	Best Category for multiple racial background	h3od6	best_category_w3			h5od8	best_category_w5
	Hispanic	h3od2	hispanic_w3			h5od4c	hispanic_w5
	Others					h5od4g	other_w5
Depressed Mood	Depressed	h3sp9	depressed_w3	h4mh22	depressed_w4	h5ss0b	depressed_w5
	Sad	h3sp12	sad_w3	h4mh26	sad_w4	h5ss0d	sad_w5
	Cannot shake of the blues	h3sp6	blues_w3	h4mh19	blues_w4	h5ss0a	blues_w5

An example of data preparation (continued)

- See the accompanying command and log files

BGSU



Center for **Family** and
Demographic Research

Conclusions

- Data preparation skills are extremely important to social sciences researchers
- Research questions and hypotheses provides the roadmap for how variables and data sets should be constructed. Thus, it is critical to think them through before working on data.
- Data preparation skills focus on Data preparation, Data Cleaning., Variable Transformation, Subset Creation, and Data Documentation and reproduction. Each of these skills requires different but simple Stata commands
- The use of command and log files is very important in achieving these data preparation tasks
- When working on data, always keep track of which respondents are in the sample and which respondents have missing values on what variables. Such information identifies the size of final analytic sample.
- If you have any questions about data preparation, please send an email to me (wuh@bgsu.edu) or drop by my office. I am glad to help.