

# Preparing Secondary Datasets for Analysis

## for class projects using STATA

#	Description of Step	STATA Syntax
1	Create a <b>project folder</b> and set it as your <b>Working Directory</b> . Gather and examine the documentation to identify what each case represents and how data was collected.	<code>cd 'path to your project folder'</code> ** cd = change directory ** ex. <code>cd 'E:/my_class/my_project/'</code>
2	Identify <b>variables of interest</b> . Select up to twice the number required as some may be problematic Create a new data file with only your selected variables, plus the unique ID variable.	<code>keep id_var var1 var2 var3 var4 var5</code> <code>compress</code> <code>save my_data</code>
3	Use descriptive statistics, value labels, and the codebook to classify each <b>variable</b> as categorical or numeric and find the meaning of each <b>value</b> . If most values are unlabeled, it is probably numeric.	<code>tab1 var1 var2 var3 var4 var5</code> <code>tabulate categorical, plot</code> <code>histogram numeric</code> <code>summarize numeric</code>
4	Determine if any values mean <b>non-answers</b> like “Not Applicable” or “Refused”. Check the <b>smallest</b> (negative) or <b>largest</b> (9. 99, or 999) values. Make sure these are <b>treated as missing</b> by the software.	** <i>In Stata, missing values are (or start with) periods and are hidden unless requested.</i> <code>mvdecode var1 var2, mv( 8 9 )</code> <code>tab1 var1 var2, missing</code>
5	Re-run descriptive statistics. <b>Drop variables</b> if they: a. have far fewer Obs than others (i.e., more missing values) b. are categorical with over 90% of cases in one group c. have nonsensical values or response patterns	<code>codebook, compact</code> ** <i>Obs = valid values</i> <code>drop problem_vars</code> <code>save my_data, replace</code>
6	See how many cases have <b>no missing values</b> : preferably 20 for each variable (5 vars→100 cases) and representative of the original cases. Drop variables with missing values and/or keep only cases with none, so analyses have the same n.	<code>misstable pattern</code> <code>egen miss = rowmiss(*)</code> <code>keep if miss == 0</code>
7	Use frequencies from Step 3 to identify <b>ordinal</b> vars and numeric vars with 7 or less values. Carefully consider whether to treat as ordinal (if appropriate analyses were taught), numeric (if mean and sd make sense), or else categorical.	** <i>If an ordinal variable is not linearly related to other values, try treating as nominal.</i>
8	Look at descriptive statistics and histograms for <b>numeric variables</b> and consider whether to <b>group values</b> if the distribution is neither normal nor flat. Numeric vars store more information, but may not best represent the responses.	<code>tabstat var, stat(mean sd med min max n)</code> <code>recode var (1/10=1 Lo) (11/20=2 Hi), g(in2)</code>
9	Look at frequencies for <b>categorical variables</b> and consider whether to <b>combine groups</b> . Have 2 to 5 groups, none with less than 10% of cases. If appropriate, compare one group to “all others” (recode to 1 vs 0).	<code>recode var (1 2=1 'Lo') (3 4 5=2 'Hi'), g(in2)</code> <code>generate is3 = (var == 3)</code> <code>label define is3 1 'Group 3' 0 'Others'</code> <code>label values is3 is3</code>
10	Before running your final analysis, examine the relationship between each pair of vars (bivariate). The choice for analysis depends on whether your X and Y is categorical or numeric and which is the predictor vs response.	<code>tab categorical categorical, row chi2</code> <code>tab categorical, sum( numeric )</code> <code>scatter numeric numeric</code>
		<code>help fvvarlist // info to use cat vars in regress</code>