

## Lab 5: Merging and hypothesis tests

### Objectives

Today we're going to keep working with `cps_2016.dta`, which contains information from the 2016 Current Population Survey.

We're going to merge in county-level unemployment rates from the BLS

By the end of this lab, you should be able to complete the following tasks in Stata:

- Import data from Excel
- Merge data sets
- Test hypotheses for linear combinations of coefficients
- Test the general significance of a regression

### Key commands

command	description
<b>Importing data</b> <code>import excel using "file1.xlsx", firstrow clear</code>	Import the data set file1.xlsx from excel into Stata. The <code>firstrow</code> option tells Stata to use the first row as the variable name. The <code>clear</code> option tells Stata to erase any data already in the set
<b>Identifying duplicates</b> <code>duplicates list var1 var2</code>  <code>duplicates tag var1 var2, gen(d1)</code>	List any observations that are duplicates on the listed variables, <code>var1 var2</code> , etc. Generate a new variable, <code>d1</code> that indicates which variables are duplicates for <code>var1</code> and <code>var2</code>
<b>Merging datasets</b> <code>merge 1:1 var1 var2 using file2</code>  <code>merge m:1 var1 var2 using file2</code>	Perform a one-to-one merge based on <code>var1</code> and <code>var2</code> . There cannot be any duplicates on the variables you are using to merge Perform a many-to-one merge based on <code>var1</code> and <code>var2</code> . There can be duplicate identifiers in the master data set (like if merging state data to individuals), but there should be no duplicates in the using data set

command	description
<b>Converting between string and numeric variables</b>	
<code>decode var1, gen(newvar)</code>	Take a numeric variable with labels and generate a new string variable that is equal to the values of those labels. (You can do the opposite with <code>encode</code> ).
<code>destring var1, replace</code>	Take a string variable, <code>var1</code> and convert it to a numeric variable, replacing the old variable
<code>tostring var2, gen(string_var)</code>	Take a numeric variable, <code>var2</code> and make it a string, but make that into a new variable called <code>string_var</code>
<b>Statistical tests</b>	
<code>test var1 = var2</code>	Run after estimating a regression. Tests the null hypothesis that the coefficient on <code>var1</code> equals the coefficient on <code>var2</code> , against the two-sided alternative.
<code>testparm var1 var2 ...</code>	Run after estimating a regression. Tests the whether all listed variables, <code>var1</code> , <code>var2</code> , etc., are jointly equal to zero, against the two-sided alternative.

### o.o.1 A note on temporary files (optional)

This exercise works by having two data sets stored on your hard drive, then running a `merge` command to unite them. You might notice that the workflow feels clunky and generates extra files - open a data set, save it, open another data set, then merge in the first data set.

You can use temporary files to speed things up! Basically, you can save files in your local memory, and call those files the same way we called local variables. Everything has to be run in the `do-file` for this to work.

A short example (you can paste this in a `do-file` and run it, as it uses Stata files) :

```
tempfile tempauto          // Declare tempfile (needs to run before you try to save)

webuse autosize,clear

save 'tempauto', replace    // save to temp file t1

webuse autoexpense, clear

merge 1:1 make using 'tempauto'  // call tempfile

tab _merge    // check out merge

list
```

```
r htmttools::HTML("{< youtube 8yfXvk8QYy0 >}")
```

## o Lab 5 Worksheet

### o.1.1 What do I submit?

- Your written up answers to the exercise questions. This can be typed or written out then scanned (or photographed), in any reasonable format.
- The do-file you've created that runs this analysis
- A log file that contains the results from this exercise.

### o.1.2 Exercises

1. Visit <https://www.bls.gov/lau/tables.htm> to access 2016 annual **county-level** *annual* unemployment rates.
  1. Download the appropriate table.
  2. Rename variables as needed, and delete any unnecessary cells. If you want your life to be easier, make the first row include your variable names, and then have the data start in second row.<sup>1</sup>
  3. Save your revised file.
2. Open Stata, start a new do-file (or bring in a template). Make sure you add code to start (and end) a log.
3. Import your unemployment excel into Stata and save it as a data file, `unemp.dta`.
4. Open `cps_2016.dta` and restrict the sample to adults (age 18+).
5. Now, merge your unemployment data into `cps_2016.dta` by county. This may not be smooth. A few tips:
  1. The FIPS codes are in different formats between the two data sets. A county code like this "55083" contains a state part (55) and a county part (083).
  2. You can convert a variable to and from a string using the commands `destring var1,replace` and `tostring var2,replace`, respectively.
  3. You can concatenate string variables by adding them like this: `gen newvar = string1 + string2`
  4. Determine whether you need a one-to-one or many-to-one merge.
  5. You may get errors, and you'll need to fix these to have a successful merge.
6. You've done it! Tabulate the new variable `_merge`. What share of observations successfully merge?<sup>2</sup>

<sup>1</sup>You can also sort this out w/ Stata commands if you'd rather work with the raw, unedited file

<sup>2</sup>To get a sense if you've done this right, about 40-45% of observations should match. This is because the CPS will withhold county-level identifiers for very small counties to protect confidentiality.

7. Drop any unmatched observations (you can use `drop if`, as we'll retain this restriction for the rest of the exercise.) What is the average unemployment rate for the entire sample - why would this be different than taking the average of county-level unemployment rates in your excel file?
8. Why can't we use education as a linear variable?
9. Generate three dummy variables. These three variables should be mutually exclusive, and they should not be missing for any people.
  - `lesshs`, a variable equal to one if a person completed *less than* a high school diploma
  - `hsgrad`, a variable equal to one if a person completed at least a high school but less than a Bachelor's degree
  - `colgrad`, a variable equal to one if a person completed a Bachelor's degree or higher

*Note:* Education is coded with **labels**, which means that it is numeric data with a description of what each number means on top. These show up as blue in the Stata browser. To view variables without the labels, add the no-label option: `tab educ, nolabel`.

10. What is the mean of `lesshs`, `hsgrad`, and `colgrad`?
11. Estimate a regression of total personal income on education, using the binary variables you just created. Omit `lesshs`.
12. Set up a hypothesis test for whether both `hsgrad` and `colgrad` are jointly significant. Report the null hypothesis, alternative hypothesis, test statistic, and conclusion.
13. Set up a hypothesis test for whether the returns to being a high-school graduate are the same as the returns to being a college graduate. Report the null hypothesis, alternative hypothesis, test statistic, and conclusion.
14. Is this regression significant overall? Explain how you know.
15. Now add county-level unemployment rate to the previous equation. What is the interpretation of the coefficient on unemployment? Is it statistically significant?
16. Estimate the same equation by regressing total personal income on the education binary variables and county-level unemployment, restricting to those who are currently in the labor force. How does this change the coefficient on unemployment?
17. Identify three *state* or *county-level* variables that are likely to cause omitted variable bias if you want to know whether unemployment affects individual wages.
18. For *one* of the variables you listed above, find the data online, import into Stata, and merge it in.
19. Regress total personal income on the education binary variables, county-level unemployment, and the new variable you found. Restrict your sample to those who are currently in the labor force. How does the inclusion of your new variable affect the coefficient on unemployment?