# STATA - Create the Sample for Analysis: Part 1

Prof. Tzu-Ting Yang

Institute of Economics, Academia Sinica

March 16, 2022

## Create the Sample for Analysis

• Please see C4\_create\_sample\_part1.do

gen: Creating new variables

Creating new variables

```
gen age_sq = age^2
gen log_incwage = log(incwage)
gen month=10 /* constant value of 10 */
gen id=_n /* id number of observation */
gen total=_N /* total number of observations */
gen bigyear=year if incwage> 100000 /* show the year
if income > 100000 */
```

- gen: creating new variables
- **Line 4:** \_*n* represents the index for a specific observation
- Line 5: \_N represents the total number of observations

id	total
10	cocai
1	60000
2	60000
3	60000
4	60000
5	60000
6	60000
7	60000
8	60000

Creating new variables

```
gen str6 source="CPS" /* string variable */
replace incwage =. if incwage==9999999 /* replace
9999999 to missing */
```

- A couple of things to note:
  - Stata's default data type is float, so if you want to create a variable in some other format (e.g. byte, string), you need to specify this (Line 1)
  - missing numeric observations, denoted by a dot, are interpreted by Stata as a very large positive number.
- You need to pay special attention to such observations when using if statements



source		
CPS		

incwage		
15000		
42000		
29000		
9999999		
9999999		

inc	wage
	15000
	42000
	29000
	-

Creating new variables based on summary measures

```
egen year_inc=total(incwage), by(year)
egen state_inc=mean(incwage), by(statefip)
egen g_state_county=group(statefip county) /*
    generates numeric id variable for state and
    county */
```

- **egen**: this command typically creates new variables based on summary measures, such as sum, mean, min and max
- **Line 1:** Use function **total** to sum the wage income for each year
- Line 2: Use function mean to get average wage income for each state

statefip	state_inc
connect	34858.93
connect	34858.93

Creating new variables based on summary measures

```
egen count_id=count(id)
egen g_state_county=group(statefip county)
egen diff_v = diff(incwage inctot)
```

- Line 1: Use function count to count number of observations
- Line 2: Use function group to generate numeric id variable for state and county
- Line 3: Use function diff to generate a variable indicating whether variables incwage and inctot are different or not

id	count_id	
1	60000	
2	60000	
3	60000	

	incwage	inctot	diff_v
1	0	5760	1
2		99999999	1
3	200000	200251	1
4	55000	55000	0

- The egen command is also useful if your data is in long format
- For example, you want to find the difference in sample *i*'s income and maximum of income within specific state

Creating new variables based on summary measures

The following routine will achieve this:

```
gen temp1=incwage if year==2015
egen temp2=max(temp1), by(statefip)
gen temp3=incwage-temp2 if year==2015
egen diff=max(temp3), by(statefip)
drop temp*
```

	incwage	temp1	temp2	temp3	inctot
1	50000	50000	1099999	-1049999	50144
2	10000	10000	1099999	-1089999	10000
3	0	0	1099999	-1099999	40008

replace: Modifying existing variables

## STATA Command:replace

Modifying existing variables

- **replace**: this command modifies existing variables in exactly the same way as generate creates new variables
- **Line 2:** replace *ln\_inctot* to zero if *ln\_inctot* == .
  - Note that if you apply a transformation to missing data, the result will still be a missing value
  - A transformation that is undefined, e.g. taking the natural log of a negative number creates a missing value

# STATA Command:replace

Modifying existing variables

	yr	year
1	114	2014
2	114	2014
3	114	2014

# STATA Command:replace

Modifying existing variables

	yr	year
1	14	2014
2	14	2014
3	14	2014

label: Put labels on datasets, variables or values

Put labels on datasets, variables or values

```
label data "Data from CPS 2014-2015"
label variable incwage "wage income"
```

- label: This command let you put labels on datasets, variables or values
- This helps to make it clear exactly what the dataset contains
- Line 1: put a label on the current data set
- Line 2: put a label on variable incwage

Put labels on datasets, variables or values

```
tab sex
label define gendercode 1 "Male" 2 "Female"
label values sex gendercode
tab sex
codebook sex
```

- It can also be helpful to label different values
- Imagine states/countries were coded as numbers (which is the case in many datasets)
- It might be better to label exactly what each value represents
- label define: defining a label (giving it a name and specifying the mapping)
- label values: associating that label with a variable

Put labels on datasets, variables or values

. tab sex

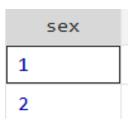
Cum.	Percent	Freq.	sex
48.58	48.58	29,148	1
100.00	51.42	30,852	2
	100.00	60,000	Total

Put labels on datasets, variables or values

. tab sex

sex	Freq.	Percent	Cum.
Male	29,148	48.58	48.58
Female	30,852	51.42	100.00
Total	60,000	100.00	

Put labels on datasets, variables or values



Put labels on datasets, variables or values



rename: Change the names of your variables

### STATA Command:rename

Change the names of your variables

```
rename sex gender rename gender sex
```

- rename: this command can change the names of your variables
- Line 1: rename sex to gender
- Line 2: rename gender to sex

recode: Change the values that certain variables take

Change the values that certain variables take

```
recode sex (1=0) (2=1)
recode incwage inctot (0 = .)
```

- recode: this command can change the values that certain variables take
- This command can also be used to recode missing values to the dot that Stata uses to denote missing
- And you can recode several variables at once

Change the values that certain variables take

```
recode incwage (0 / 50000 = 1) (50001 / 100000 = 2) (100000 / 7000000 = 3)
```

- recode can also make several changes at the same time
- We could, for example, use recode to generate a new variable with categorical income values

Change the values that certain variables take

incwage	
15000	
42000	
29000	

Change the values that certain variables take

incwage	
	1
	1
	1

tostring/destring: Change variable to string/numeric variables

## Numeric and String Variables

- Stata stores or formats data in either of two ways numeric or string
- Numeric will store numbers while string will store text
- Numeric variables are in black/blue color and string variables are in red color
- String variale can also be used to store numbers, but you will not be able to perform numerical analysis on those numbers

# STATA Command:tostring

Change variable to string

- 1 tostring year, replace
  - change variable to string

# STATA Command:tostring

Change variable to string

	year				
1	2014				
2	2014				
3	2014				

# STATA Command:tostring

Change variable to string

	year					
1	2014					
2	2014					
3	2014					

# STATA Command:destring

Change variable to numeric

#### Consolas

```
destring year, gen(year1)
```

• change variable to numeric

# STATA Command:destring

Change variable to numeric

	year	year1		
1	2014	2014		
2	2014	2014		
3	2014	2014		

decode: Create a string variable from a numerical code

## STATA Command: decode

Create a string variable from a numerical code

- decode statefip,gen(state\_name)
  - decode: Create a string variable from a numerical code, as long as the numeric variable has labels attached to each value
  - Line 1: This creates a new variable state\_name, which use label of numerical code as string variable

## STATA Command: decode

Create a string variable from a numerical code

	statefip	state_name
1	maine	maine
2	maine	maine

encode: Converting strings to numerical code

## STATA Command:encode

Converting strings to numerical code

```
encode state_name, gen(state_id)
```

- encode: Converting strings to numerical code
- Line 1: This creates a new variable state\_id, which takes a value of 1 for alabama, 2 for alaska, and so on.

## STATA Command:encode

Converting strings to numerical code

	state_name	state_id
1	maine	maine
2	maine	maine

# substr: Dividing string variables

# Combining String Variables

```
tostring year, gen(yearcode)
gen yearcode1 = string(year)
gen stateyear = state_name + yearcode1
```

- We can create a new string variable whose data is a combination of the data values of other variables
- Line 2: This creates a new variable "stateyear", which combine two string variable "state\_name" and "yearcode1"

# Combining String Variables

	yearcode1	state_name	stateyear
1	2014	maine	maine2014
2	2014	maine	maine2014

### STATA Command:substr

**Dividing String Variables** 

```
gen yr1 = substr(yearcode,3,2)
```

- substr: Divide up a variable or to extract part of a variable to create a new one
  - The first term in parentheses is the string variable that you are extracting from
  - The second term (3) is the position of the first character you want to extract
  - The third term (2) is the number of characters to be extracted

### STATA Command:substr

Dividing String Variables

#### Consolas

```
gen yr2 = substr(yearcode,-2,2)
```

 Alternatively, you can select your starting character by counting from the end (2 positions from the end instead of 3 positions from the start)

## STATA Command:substr

**Dividing String Variables** 

	yr2	yr1	yearcode	
1	14	14	2014	
2	14	14	2014	

# Create dummy variables

#### Consolas

```
gen largewage1=(incwage_test>=30000)
```

• **Line 1:** create a variable largewage1 if the statement within parentheses is true (*incwage\_test* >= 30000)

	largewage1	incwage
1	0	15000
2	1	42000
3	0	29000

```
tab statefip, gen(state_d)
```

- This creates:
  - a dummy variable "state\_d1" equal to 1 if the state is "alabama" and zero otherwise
  - a dummy variable "state\_d2" if the state is "alaska" and zero otherwise and so on
- We can use "state\_d\*" to represent the all dummy variables

. tab statefip, gen(state\_d)

state (fips code)	Freq.	Percent	Cum.
connecticut	2,799	4.67	4.67
delaware	2,113	3.52	8.19
illinois	4,393	7.32	15.51
indiana	2,137	3.56	19.07
iowa	2,486	4.14	23.21
kansas	1,984	3.31	26.52
maine	1,794	2.99	29.51
maryland	2,319	3.87	33.38

	statefip	state_d1	state_d2	state_d3	state_d4	state_d5	state_d6	state_d7
1	maine	0	0	0	0	0	0	1
2	maine	0	0	0	0	0	0	1
3	maine	0	0	0	0	0	0	1

```
reg incwage i.sex#i.year reg incwage i.sex##i.year
```

- Sometimes we just need to include these dummies in our regression and do not want to create these variable permanently
- Using one # tells Stata to ignore the level for each variable and only report the interaction
- Doubling up the #, i.region##i.year makes the command account for both dummies for each region, each year and each interaction.

. reg incwage i.sex#i.year //i.region#i.year Source SS df MS Number of obs 60,000 F(5, 59994) 18.94 1.6571e+15 3.3141e+14 Prob > F Model 5 0.0000 Residual 1.0495e+18 59,994 1.7494e+13 R-squared 0.0016 0.0015 Adj R-squared Total 1.0512e+18 59,999 1.7520e+13 Root MSE 4.2e+06 Coef. Std. Err. t P>|t| [95% Conf. Interval] incwage sex#year 1#2015 223840.4 60074.95 3.73 0.000 106093.2 341587.5 1#2016 222358.2 60117.63 3.70 0.000 104527.4 340188.9 2#2014 -220957.9 59198.05 -3.73 0.000 -336986.3 -104929.5 0.063 2#2015 -110505.1 59445,27 -1.86 -227018 6007.824 2#2016 109416.8 1.84 0.065 -7016,214 225849.8 59404.48 \_cons 2271781 42675.02 53.23 0.000 2188138 2355424

. reg incwage i.sex##i.year

Source Model Residual	SS 1.6571e+15 1.0495e+18	df 5 59,994	MS 3.3141e+14 1.7494e+13	F(5, <b>4</b> Prob <b>3</b> R-sq	er of obs 59994) > F uared R-squared	= 60,000 = 18.94 = 0.0000 = 0.0016 = 0.0015
Total	1.0512e+18	59,999	1.7520e+1	<b>3</b> Root	MSE	= 4.2e+06
incwage	Coef.	Std. Err.	t	P> t	[95% Con	f. Interval]
sex 2 year	-220957.9	59198.05	-3.73	0.000	-336986.3	-104929.5
2015 2016	223840.4 222358.2	60074.95 60117.63	3.73 3.70	0.000 0.000	106093.2 104527.4	
sex#year 2#2015 2#2016	-113387.5 108016.5	83694.88 83696.57	-1.35 1.29	0.175 0.197	-277429.8 -56029.05	272062.1
_cons	2271781	42675.02	53.23	0.000	2188138	2355424