







Introductory Stata Training Lesson 2 - Transforming Data

Harry Gimiseve, Rishabh Mukerjee, Emily Schmidt

Port Moresby, Papua New Guinea

March 2025

International Food Policy Research Institute (IFPRI)

1. Review of Lesson 1



1. Review of Lesson 1

2. <u>Transforming existing data</u>



- 1. Review of Lesson 1
- 2. Transforming existing data
- 3. <u>Creating new variables</u>

- 1. <u>Review of Lesson 1</u>
- 2. Transforming existing data
- 3. <u>Creating new variables</u>
- 4. <u>Help window</u>

- 1. <u>Review of Lesson 1</u>
- 2. Transforming existing data
- 3. <u>Creating new variables</u>
- 4. <u>Help window</u>
- 5. Collapse and merge

Today we are going to work mostly with the household roster dataset.



Today we are going to work mostly with the household roster dataset.

First, let me show you the household roster module in the questionnaire (Module 1.1)

MODULE 1 – HOUSEHOLD CHARACTERISTICS Section 1.1: Household demographics, current household members

We would like to record all of the individuals that are members of this household. A household member is "a part of a group of people who live together, may pool some of their money, and eat most meals from the same pot". A member can also be someone who lives most of the time (at least 8 months of the year) and shares meals with the family, but is currently away (e.g., for work, ceremony, travel, etc.) It is possible that there are individuals who are household members that are not related by blood or family ties, but they are still a part of the household because they share in the household resources (such as a servant that lives in the household and eats from the same food stock).

	1	2	3	4a	4b.	4c – 4e			5a.	5b.	5c.	5d.	6.
	Please	What is	What is	What is the age of I understand that you		When was [NAME] born?		For	For	For	For	Do you have a	
	[NAME]	the	the sex	[NAME]?	do not know the age of				children	children	children	children	clinic card / book
	everyone	relation-	of		[NAME], but do you				<5:	<5:	<5:	<5:	for this child?
	in your	ship of	[NAME]	[Mark 0 if [NAME]	think you could give								
	household	[NAME]	?	is less than 1	me an estimate of the				Is the	Who is the	Is the	Who is the	[code:
	starting	to the head?		year.j	age within these			mother of	the child a	father?	yesno_ak_vij		
<u>_</u>	vourself		If ano is <- 5 ->	categories?				the child a [PID]		[DID]	נסוסו		
yourself.				skin to Q5	[code:				current		household	ורוטן	Please ask to see
₽		[code:		omp to Qo	agecategory vll				household		member?		the immunization
son		rel_vl]	[code:	If age is > 5 and	5 5 5 7 1	4c. Day	4d.	4e.	member?				card, you will need
Pe		-	gender	< 12 → skip to			Month	Year					it for future
			_vl]	Q10					lf no		lf no		modules.]
				11		If ago in	[aada:		→skip to		→skip to		
				If age ≥ 12 →		11 age is 'Don't	months		50		6		
				SKIP IO Q7		know'	vl]						
				If age is 'Don't		please							
				know' please write		write							
				888 → answer 4b		888							
1													
2													
3													
4													
-													
5													
6													
7													

- Today we are going to work mostly with the household roster dataset.
- First, let's turn to the household roster module in the questionnaire (Module 1.1 - page 5)
- This module asks questions about each person in the household.
- Earlier, the data was at the household level

Earlier, the data was at the household level, and now it is at the person level

Earlier, the data was at the household level, and now it is at the person level

Do you think we will have more or less observations in our dataset today, compared to previous lesson?

Yesterday, the data was at the household level, and today it is at the person level

Do you think we will have more or less observations in our dataset today, compared to previous lesson?

How can we find out how many observations are in the new dataset?

Yesterday, the data was at the household level, and today it is at the person level

Do you think we will have more or less observations in our dataset today, compared to yesterday?

How can we find out how many observations are in the new dataset? Count

Let's double-click to open s11_long.dta

Let's double-click to open s11_long.dta

How many observations are there?

Let's double-click to open s11_long.dta

How many observations are there? Count

Let's double-click to open s11_long.dta

How many observations are there? Count

count

16,192

16,192 observations (or 16,192 people)

The hh_3_4.dta dataset, contained 2,699 household observations

The hh_3_4.dta dataset, contained 2,699 household observations

Now, the household roster dataset, there are 16,192 person observations (person-level)

The hh_3_4.dta dataset, contained 2,699 household observations

Now, the household roster dataset, there are 16,192 person observations (person-level)

So there are many more observations in this dataset than in the previous data

Let's take a deeper look at some of the variables in this new dataset

- Question: How many men are in this dataset?
- Looking at Window 3 (Variables Window), which variable can tell us this information?

- Question: How many men are in this dataset?
- Looking at Window 3 (Variables Window), which variable can tell us this information? a3 - "3. What is the sex of \${a1}?"

- Question: How many men are in this dataset?
- Looking at Window 3 (Variables Window), which variable can tell us this information? a3 - "3. What is the sex of a1?"
- a1 was the name of the HH member, so a3 is asking their sex

Question: How many men are in this dataset?

What code can we use on a3 to see how many men are in this dataset?

Question: How many men are in this dataset?

What code can we use on a3 to see how many men are in this dataset? tab a3

- What code can we use on a3 to see how many men are in this dataset? tab a3
 - Codebook a3 could also work, but tabulate is the best option because it will also tell us the percentage in addition to the frequency

. tab a3			
3. What is the sex of ?	Freq.	Percent	Cum.
Male Female	8,459 7,733	52.24 47.76	52.24 100.00
Total	16,192	100.00	



	. tab a3			
	3. What is the sex of ?	Freq.	Percent	Cum.
ſ	Male	8,459	52.24	52.24
	Female	7,733	47.76	100.00
-	Total	16,192	100.00	

Question: How many men are in this dataset?

▶ <u>Answer:</u> *8*,459

Question: What percent of this dataset is male?

. tab a3			
3. What is the sex of ?	Freq.	Percent	Cum.
Male	8,459	52.24	52.24
Female	7,733	47.76	100.00
Total	16,192	100.00	

Question: How many men are in this dataset?

Answer: 8,459

Question: What percent of this dataset is male? 52.24%

. tab a3			
3. What is the sex of ?	Freq.	Percent	Cum.
Male	8,459	52.24	52.24
Female	7,733	47.76	100.00
Total	16,192	100.00	

Let's look at another variable in this new dataset.

Question: What is the average age in this dataset?

Let's look at another variable in this new dataset.

Question: What is the average age in this dataset?

Looking at Window 3 (Variables Window), which variable can tell us this information?
Let's look at another variable in this new dataset.

Question: What is the average age in this dataset?

Looking at Window 3 (Variables Window), which variable can tell us this information? a4a - "4. What is the age of \${a1}?"

Question: What is the average age in this dataset?

What code can we use on a4a to see what the average age is?

Question: What is the average age in this dataset?

What code can we use on a4a to see what the average age is? sum a4a

Can also use *sum a4a*, *det* or *codebook a4a*

Question: What is the average age in this dataset?

What code can we use on a4a to see what the average age is? sum a4a

Can also use *sum a4a*, *det* or *codebook a4a*

. sum a4a					
Variable	Obs	Mean	Std. dev.	Min	Max
a4a	16,192	125.095	279.0221	0	888

Question: What is the average age in this dataset?

What code can we use on a4a to see what the average age is? sum a4a

Can also use *sum a4a*, *det* or *codebook a4a*

. sum a4a					
Variable	Obs	Mean	Std. dev.	Min	Max
a4a	16,192	125.095	279.0221	0	888

What's wrong with this output?

Question: What is the average age in this dataset?

What code can we use on a4a to see what the average age is? sum a4a

Can also use *sum a4a, det* or *codebook a4a*

. sum a4a					
Variable	Obs	Mean	Std. dev.	Min	Max
a4a	16,192	125.095	279.0221	0	888

What's wrong with this output? Average age appears to be too high!

What's wrong with this output? Average age appears to be too high!

Why is the average age too high?

. sum a4a					
Variable	Obs	Mean	Std. dev.	Min	Max
a4a	16,192	125.095	279.0221	0	888

- What's wrong with this output? Average age appears to be too high!
- Why is the average age too high? Maximum age is an impossibly large number, 888.

. sum a4a					
Variable	Obs	Mean	Std. dev.	Min	Max
a4a	16,192	125.095	279.0221	0	888

- What's wrong with this output? Average age appears to be too high!
- Why is the average age too high? Maximum age is an impossibly large number, 888.

. sum a4a					
Variable	Obs	Mean	Std. dev.	Min	Max
a4a	16,192	125.095	279.0221	0	888

If we look at the questionnaire, we see that data collectors were told to write 888 if age is 'don't know'

2. Transforming Existing Data

Now, we will learn how we can transform/change the data, to find the real average age

2. Transforming Existing Data

Now, we will learn how we can transform/change the data, to find the real average age

But first, let's talk about *missing values*

Which command tells us how many observations have a missing value for a variable?

Which command tells us how many observations have a missing value for a variable? Codebook

So let's see how many missing observations there are in our age variable. What do we type?

So let's see how many missing observations there are in our age variable. What do we type? Codebook a4a

So let's see how many missing observations there are in our age variable. What do we type? Codebook a4a

How many missing observations are there?

. codebook a4a						
 a4a						4a. What is the age of
Type:	Numeric (i	nt)				
Range: Unique values:	[0,888] 97		Miss	Units: 1 ing .: 0/1	6,192	
Mean: Std. dev.:	125.095 279.022					
Percentiles:	10% 5	25% 10	50% 22	75% 45	90% 888	

So let's see how many missing observations there are in our age variable. What do we type? Codebook a4a

How many missing observations are there? 0



• codebook a4a											
a4a						 4a.	What	is	the a	ige	of
Туре:	Numeric	(int)									
Range: Unique values:	[0,888] 97		Miss:	Units: 1 ing .: 0/1	6,192						
Mean: Std. dev.:	125.095 279.022										
Percentiles:	10% 5	25% 10	50% 22	75% 45	90% 888						

But we just said that there are missing observations, and they were input as 888

. codebook a4a												
 a4a							4a.	What	is	the a	age	of
Туре:	Numeric	(int)										
Range: Unique values:	[0,888] 97		Miss	Units: 1 ing .: 0/1	6,192							
Mean: Std. dev.:	125.095 279.022											
Percentiles:	10% 5	25% 10	50% 22	75% 45	90% 888							

- But we just said that there are missing observations, and they were input as 888
- So, we need to change 888 so that they are recognized as missing.

Missing values appear in 2 different ways in Stata, depending on the variable's data type:

- Missing values appear in 2 different ways in Stata, depending on the variable's data type:
 - Numeric/categorical variables
 - Missing values are coded as a period (.)
 - Code: replace [varname]=. if ...

- Missing values appear in 2 different ways in Stata, depending on the variable's data type:
 - Numeric/categorical variables
 - Missing values are coded as a period (.)
 - ► Code: replace [varname]=. if ...

String variables

- Missing values are coded as a blank ("")
- Code: replace [varname]="" if ...

We need to know what value to change the 888 to in the a4a variable, so Stata recognizes it as missing

- We need to know what value to change the -888 to in the a4a variable, so Stata recognizes it as missing
- What type of data is a4a? What code can we type to find out?

- We need to know what value to change the -888 to in the a4 variable, so Stata recognizes it as missing
- What type of data is a4a? What code can we type to find out? Codebook a4a - a4a is numeric

- We need to know what value to change the 888 to in the a4a variable, so Stata recognizes it as missing
- What type of data is a4a? What code can we type to find out? Codebook a4a - a4a is numeric
- So what type of missing value do we want to change the 888 to?

- We need to know what value to change the 888 to in the a4a variable, so Stata recognizes it as missing.
- What type of data is a4a? What code can we type to find out? Codebook a4a - a4a is numeric
- So what type of missing value do we want to change the 888 to? A period (.)

We need to change all 888 to . in variable a4a. There are 2 ways to do this:

We need to change all 888 to . in variable a4a. There are 2 ways to do this:

▶ <u>Code:</u> *recode* a4a 888=.

We need to change all 888 to . in variable a4a. There are 2 ways to do this:

▶ <u>Code:</u> *recode* a4a 888=.

Code: replace a4a=. if a4a==888

We need to change all 888 to . in variable a4a. There are 2 ways to do this:

▶ <u>Code:</u> *recode* a4a 888=.

Code: replace a4a=. if a4==888

Choose one of these codes, and run it

Now, let's run the codebook command again for a4a and see if we succeeded in changing them to missing:

Now, let's run the codebook command again for a4a and see if we succeeded in changing them to missing:

. cod	lebook a4a							
a4a							4a. What is the a	ge of
	Туре:	Numeric (i	nt)					
	Range:	[0,106]			Units: 1			
	Unique values:	96		Miss	ing .: 1,9	04/16,192		
	Mean:	23.4313						
	Std. dev.:	18.0899						
	Percentiles:	10%	25%	50%	75%	90%		
		4	9	18	34	51		

Now, let's run the codebook command again for a4a and see if we succeeded in changing them to missing:

. codebook a4a						
a4a						4a. What is the age of
Type:	Numeric (i	int)				
Range: Unique values:	[0,106] 96		l Miss:	Jnits: 1 ing .: 1,90	04/16,192	
Mean: Std. dev.:	23.4313 18.0899					
Percentiles:	10% 4	25% 9	50% 18	75% 34	90% 51	

How many ages are unknown/missing?

adahaak ata

Now, let's run the codebook command again for a4a and see if we succeeded in changing them to missing:

. coue	DOOK a4a									
a4a								4a.	What is	the age of
	Type:	Numeric (i	nt)							
	Range:	[0,106]			Units: 1					
	Unique values:	96		Miss	ing .: 1,9	04/16,192	J			
	Mean:	23.4313								
	Std. dev.:	18.0899								
	Percentiles:	10%	25%	50%	75%	90%				
		4	9	18	34	51				

How many ages are unknown/missing? 1,904 (REMEMBER THIS, it will be used later)

Another way to see how many missing values there are is to use the "missing" option on the tabulate command
Another way to see how many missing values there are is to use the "missing" option on the tabulate command

▶ tab a4a<u>, missing</u>

Another way to see how many missing values there are is to use the "missing" option on the tabulate command

► tab a4a<u>, missing</u>

The beginning and end of the output is shown here (it's a very long output because there are so many different ages!)

. tab a4a, m			
4a. What is			
the age of	Freq.	Percent	Cum.
0	317	1.96	1.96
1	255	1.57	3.53
2	349	2.16	5.69
3	311	1.92	7.61
4	327	2.02	9.63
5	260	1.61	11.23
91	1	0.01	88.20
92	1	0.01	88.21
93	2	0.01	88.22
96	1	0.01	88.23
105	1	0.01	88.23
106	1	0.01	88.24
	1,904	11.76	100.00
Total	16,192	100.00	

- Another way to see how many missing values there are is to use the "missing" option on the tabulate command
- ► tab a4a<u>, missing</u>
- The beginning and end of the output is shown here (it's a very long output because there are so many different ages!)
- What percent of ages are missing?

. tab a4a, m			
4a. What is			
the age of	Freq.	Percent	Cum.
0	317	1.96	1.96
1	255	1.57	3.53
2	349	2.16	5.69
3	311	1.92	7.61
4	327	2.02	9.63
5	260	1.61	11.23
91	1	0.01	88.20
92	1	0.01	88.21
93	2	0.01	88.22
96	1	0.01	88.23
105	1	0.01	88.23
106	1	0.01	88.24
•	1,904	11.76	100.00
Total	16,192	100.00	

- Another way to see how many missing values there are is to use the "missing" option on the tabulate command
- ► tab a4a<u>, missing</u>
- The beginning and end of the output is shown here (it's a very long output because there are so many different ages!)
- What percent of ages are missing? 11.76%

. tab a4a, m			
4a. What is			
the age of	Freq.	Percent	Cum.
0	317	1.96	1.96
1	255	1.57	3.53
2	349	2.16	5.69
3	311	1.92	7.61
4	327	2.02	9.63
5	260	1.61	11.23
91	1	0.01	88.20
92	1	0.01	88.21
93	2	0.01	88.22
96	1	0.01	88.23
105	1	0.01	88.23
106	1	0.01	88.24
	1,904	11.76	100.00
Total	16,192	100.00	

Now we can find the average age in the dataset.

What code do we use?

Now we can find the average age in the dataset.

► What code do we use? *sum a4a*

. sum a4a					
Variable	Obs	Mean	Std. dev.	Min	Max
a4a	14,288	23.43127	18.08986	0	106

Now we can find the average age in the dataset.

► What code do we use? *sum a4a*

. sum a4a					
Variable	Obs	Mean	Std. dev.	Min	Max
a4a	14,288	23.43127	18.08986	0	106

What is the average age?

▶ Now we can find the average age in the dataset.

► What code do we use? *sum a4a*

. sum a4a					
Variable	Obs	Mean	Std. dev.	Min	Max
a4a	14,288	23.43127	18.08986	0	106

► What is the average age? 23.43 years

Now we can find the average age in the dataset.

► What code do we use? *sum a4a*

. sum a4a					
Variable	Obs	Mean	Std. dev.	Min	Max
a4a	14,288	23.43127	18.08986	0	106

How many observations are being included in this summarize command?

▶ Now we can find the average age in the dataset.

► What code do we use? *sum a4a*

. sum a4a					
Variable	Obs	Mean	Std. dev.	Min	Max
a4a	14,288	23.43127	18.08986	0	106

How many observations are being included in this summarize command? 14,288

Now we can find the average age in the dataset.

► What code do we use? *sum a4a*

. sum a4a					
Variable	Obs	Mean	Std. dev.	Min	Max
a4a	14,288	23.43127	18.08986	0	106

- How many observations are being included in this summarize command? 14,288
 - This is because 1,904 are now missing (they are not included)

What if we want our age variable to have a more intuitive name and label?

- What if we want our age variable to have a more intuitive name and label?
- Renaming a variable
 - <u>Code:</u> rename a4 age

- What if we want our age variable to have a more intuitive name and label?
- Renaming a variable
 - ▶ <u>Code:</u> rename a4 age
- Relabeling a variable
 - <u>Code:</u> label variable age "Age of household member"

- What if we want our age variable to have a more intuitive name and label?
- Renaming a variable
 - Code: rename a4 age
- Relabeling a variable
 - <u>Code:</u> label variable age "Age of household member"
 - Remember, we just changed the name of our variable, so we have to refer to it here as "age" instead of "a4"

Now let's visualize our newly cleaned age variable

Now let's visualize our newly cleaned age variable

Do you remember the code we learned yesterday to show the density of a variable as bar graphs?

Now let's visualize our newly cleaned age variable

Do you remember the code we learned yesterday to show the density of a variable as bar graphs? *Histogram*

Now let's visualize our newly cleaned age variable

Do you remember the code we learned yesterday to show the density of a variable as bar graphs? *Histogram*

Code: hist age





- hist age
- What can we learn from this figure?



- hist age
- What can we learn from this figure?
 - The majority of people in the dataset are 25 years old or younger
 - There are very few people in the dataset older than 80
 - The sample comprise of more young population



Now we are going to learn how to create a new variable

Now we are going to learn how to create a new variable

The type of variable we are going to create is called a dummy variable

A dummy variable (also known as an indicator variable or a binary variable) takes the value 0 or 1 to indicate the absence or presence of some categorical effect

- A dummy variable (also known as an indicator variable or a binary variable) takes the value 0 or 1 to indicate the absence or presence of some categorical effect
- For example, it may be useful to have a variable that easily indicates whether or not someone on the household roster is a child (15 years or younger)

- A dummy variable (also known as an indicator variable or a binary variable) takes the value 0 or 1 to indicate the absence or presence of some categorical effect
- For example, it may be useful to have a variable that easily indicates whether or not someone on the household roster is a child (15 years or younger)
- Examples of what this would let us easily see:
 - How many children are in the sample
 - What percent of children are in school

3. Creating New Variables Generate

Let's make a dummy variable for whether or not someone is a child (15 years or younger)

3. Creating New Variables Generate

- Let's make a dummy variable for whether or not someone is a child (15 years or younger)
- First, we will create the variable and set all values equal to missing (.)

3. Creating New Variables Generate

- Let's make a dummy variable for whether or not someone is a child (15 years or younger)
- First, we will create the variable and set all values equal to missing (.)
- <u>Code:</u> generate child=.
 Shorthand: gen child=.

Let's explore our new variable:

Let's explore our new variable: codebook child

Let's explore our new variable: codebook child

. code	book child				
child					(unlabeled
	type:	numeric (float)			
	range:	[.,.]	units:		
	unique values:	0	missing .:	16,192/16,192	
	tabulation:	Freq. Value 16,192 .			

Next, let's change all observations to 1 if the person is 15 years or younger.

Next, let's change all observations to 1 if the person is 15 years or younger.

What code did we learn to recode or change values?

- Next, let's change all observations to 1 if the person is 15 years or younger.
- What code did we learn to recode or change values? replace...if
- Next, let's change all observations to 1 if the person is 15 years or younger.
- What code did we learn to recode or change values? replace...if
- <u>Code:</u> replace child=1 if age<=15
 Also works: replace child=1 if age<16

Let's explore our new variable again: codebook child

Let's explore our new variable again: codebook child

. codebook child				
child				(unlabeled)
type:	numeric (float)			
range: unique values:	[1,1] 1	units: missing .:	1 10,132/16,192	
tabulation:	Freq. Value 6,060 1 10,132 .			

Next, let's change all observations to 0 if the person is an adult (16 years or older)

Next, let's change all observations to 0 if the person is an adult (16 years or older)





- Next, let's change all observations to 0 if the person is an adult (16 years or older)
- Code: replace child=0 if age>=16
 - Also works: replace child=0 if age>15

Let's explore our new variable again: codebook child

Let's explore our new variable again: codebook child

. codebook child					
child					(unlabeled)
type:	numeri	c (float)			
range: unique values:	[0,1] 2		units: missing .:	1 0/16,192	
tabulation:	Freq. 10,132 6,060	Value 0 1			

Let's explore our new variable again: *codebook child*

. codebook child					
child					(unlabeled)
type:	numeri	c (float)			
range:	[0,1]		units:	1	
unique values:	2		missing .:	0/16,192	
tabulation:	Freq.	Value			
	10,132	0			
	6,060	1			

How many missing values do we have?

Let's explore our new variable again: *codebook child*

. codebook child				
child				(unlabeled)
type:	numeri	c (float)		
range: unique values:	[0,1] 2		units: 1 missing .: 0/16,192	
tabulation:	Freq. 10,132 6,060	Value O 1		

How many missing values do we have? 0

Let's explore our new variable again: codebook child

. codebook child				
child				(unlabeled)
type:	numeric	(float)		
range: unique values:	[0,1] 2		units: 1 missing .: 0/16,192	
tabulation:	Freq. 10,132 6,060	Value 0 1		

How many missing values do we have? 0 - this is a problem! Does anyone know why?

How many missing values do we have? 0 - this is a problem! Does anyone know why?

How many missing values were changed when we recoded the 888 to missing in the "age" variable?

How many missing values do we have? 0 - this is a problem! Does anyone know why?

How many missing values were changed when we recoded the 888 to missing in the "age" variable? 1,904

How many missing values do we have? 0 - this is a problem! Does anyone know why?

How many missing values were changed when we recoded the -888 to missing in the "age" variable? 1,908

So, there are 1904 people who we don't know the ages of, but now they are coded as either "yes" a child or "no", not a child in our new "child" variable...

How many missing values do we have? 0 - this is a problem! Does anyone know why?

Stata treats missing values in numeric variables as "infinite", and so when we use > or >=, all missing values are included.

How many missing values do we have? 0 - this is a problem! Does anyone know why?

- Stata treats missing values in numeric variables as "infinite", and so when we use > or >=, all missing values are included.
- How can we change these values in the "child" variable back to missing, if they were missing in the "age" variable?

How many missing values do we have? 0 - this is a problem! Does anyone know why?

- Stata treats missing values in numeric variables as "infinite", and so when we use > or >=, all missing values are included.
- How can we change these values in the "child" variable back to missing, if they were missing in the "age" variable? *replace child=. if age==*.

How many missing values do we have? 0 - this is a problem! Does anyone know why?

- Stata treats missing values in numeric variables as "infinite", and so when we use > or >=, all missing values are included.
- How can we change these values in the "child" variable back to missing, if they were missing in the "age" variable? *replace child=. if age==*.

Note the use of the single = compared to the double ==

Let's explore our new variable again: codebook child

Let's explore our new variable again: codebook child

. codebook child				
child				(unlabeled)
type:	numeric (float)			
range:	[0,1]	units:	1	
unique values:	2	missing .:	1,904/16,192	
tabulation:	Freq. Value			
	8,228 0			
	6,060 1			
	1,904 .			

Let's explore our new variable again: codebook child

. codebook child				
child				(unlabeled)
type:	numeric (float)			
range: unique values:	[0,1] 2	units: missing .:	1 1,904/16,192	
tabulation:	Freq. Value 8,228 0 6,060 1 1,904 .			

Now what do we still need to add?

Let's explore our new variable again: codebook child

. codebook child				
child				(unlabeled)
type:	numeric (float)			
range: unique values:	[0,1] 2	units: missing .:	1 1,904/16,192	
tabulation:	Freq. Value 8,228 0 6,060 1 1,904 .			

Now what do we still need to add? Labels!

Let's first add a variable label. What is the code for this?

Let's first add a variable label. What is the code for this?

Label variable child "Is the household member 15 years or younger?"

Shorthand: lab var child "Is the household member..."

Now, let's label the 0 values as "No" and the 1 values as "Yes"

Now, let's label the 0 values as "No" and the 1 values as "Yes"

First, we have to create a value label definition

Now, let's label the 0 values as "No" and the 1 values as "Yes"

First, we have to create a value label definition

<u>Code:</u> label define no_yes 0 "No" 1 "Yes"

- Now, let's label the 0 values as "No" and the 1 values as "Yes"
- First, we have to create a value label definition
- Code: label define no_yes 0 "No" 1 "Yes"
- Right now this label definition is only saved in Stata's memory - it has not been applied to any variables yet

Now we will write the code to apply this value label to our new "child" variable

Now we will write the code to apply this value label to our new "child" variable

Code: label values child no_yes

Now we will write the code to apply this value label to our new "child" variable

Code: label values child no_yes

This code applies our new label definition "no_yes" to the variable "child"

- Now we will write the code to apply this value label to our new "child" variable
- Code: label values child no_yes
- This code applies our new label definition "no_yes" to the variable "child"
- Let's look one more time at our new variable: codebook child



Let's explore our new variable: codebook child

. codebook child												
child					Is th	he ho	usehold	member	15	years	or	younger?
type:	numeric	(float)										
label:	noyes											
range:	[0,1]			units:	1							
unique values:	2			missing .:	1,904/16,3	192						
tabulation:	Freq.	Numeric	Label									
	8,228	0	No									
	6,060	1	Yes									
	1,904											

et's explore ou	r nev	w va							
codebook child			riab	le: cod	ebook ch	ild usehold member	15 years or you	nger?	
type: r label: 1	numeric noyes	(float)							
range:	[0 , 1]			units:	1				
unique values: 2	2			missing .:	1,904/16,192				
tabulation: I	Freq. 1	Numeric	Label						
٤	8,228	0	No						
e	6,060	1	Yes						
1	1,904	•							






Question: What percent of children have 'student' listed as their primary activity?

Question: What percent of children have 'student' listed as their primary activity?

What code might help us find this answer?

Question: What percent of children have 'student' listed as their primary activity?

What code might help us find this answer? Tabulate

- Question: What percent of children have 'student' listed as their primary activity?
- ► What code might help us find this answer? *Tabulate*
- Code: tab a8 if child==1

Question: What percent of children have 'student' listed as their primary activity?

. tab a8 if child == 1			
<pre>b8 : What is the current PRIMARY activity / occupation for \${b_name2}?</pre>			
	Freq.	Percent	Cum.
Any agricultural (farmer, gardener, sub	125	10.37	10.37
Forestry or fishery labor	1	0.08	10.46
Transport and storage laborers	1	0.08	10.54
Services and sales worker	1	0.08	10.62
Security guard	1	0.08	10.71
Student	1,041	86.39	97.10
Domestic work/housewife	7	0.58	97.68
Unemployed and NOT looking for work	6	0.50	98.17
Unemployed and ACTIVELY looking for wor	1	0.08	98.26
Community leader (counselor, peace medi	1	0.08	98.34
Not applicable	18	1.49	99.83
Other	2	0.17	100.00
Total	1,205	100.00	

Question: What percent of children have 'student' listed as their primary activity?

. tab a8 if child == 1			
<pre>b8 : What is the current PRIMARY activity / occupation for \${b_name2}?</pre>			
	Freq.	Percent	Cum.
Any agricultural (farmer, gardener, sub	125	10.37	10.37
Forestry or fishery labor	1	0.08	10.46
Transport and storage laborers	1	0.08	10.54
Services and sales worker	1	0.08	10.62
Security guard	1	0.08	10.71
Student	1,041	86.39	97.10
Domestic work/housewife	7	0.58	97.68
Unemployed and NOT looking for work	6	0.50	98.17
Unemployed and ACTIVELY looking for wor	1	0.08	98.26
Community leader (counselor, peace medi	1	0.08	98.34
Not applicable	18	1.49	99.83
Other	2	0.17	100.00
Total	1,205	100.00	



Question: What percent of children have 'student' listed as their primary activity?

. tab a8 if child == 1			
<pre>b8 : What is the current PRIMARY activity / occupation for \${b_name2}?</pre>	Freq.	Percent	Cum.
Any agricultural (farmer, gardener, sub	125	10.37	10.37
Forestry or fishery labor	1	0.08	10.46
Transport and storage laborers	1	0.08	10.54
Services and sales worker	1	0.08	10.62
Security guard	1	0.08	10.71
Student	1,041	86.39	97.10
Domestic work/housewife	7	0.58	97.68
Unemployed and NOT looking for work	6	0.50	98.17
Unemployed and ACTIVELY looking for wor	1	0.08	98.26
Community leader (counselor, peace medi	1	0.08	98.34
Not applicable	18	1.49	99.83
Other	2	0.17	100.00
Total	1,205	100.00	

Answer: 86.39% of children aged 15 and younger are primarily students

Challenge question: Which province has the highest percentage of children?

Challenge question: Which province has the highest percentage of children?

Let's use a figure to find the answer!

Challenge question: Which province has the highest percentage of children?

Let's use a figure to find the answer!

<u>Code:</u> graph bar child, over(prov)

graph bar child, over(prov)



graph bar child, over(prov)



- Notice that the bar labels on x-axis are on top of each other. To produce a neat graph, we can use the following code to avoid the overlap:
- graph bar child, over(prov, label(angle(45)))

graph bar child, over(prov, label(angle(45)))



graph bar child, over(prov, label(angle(45)))



 Because values are 0 and 1, a bar graph (ranging from 0 to 1) shows the prevalence of the dummy variable (1=100%)



Morobe has the highest prevalence of children

 Because values are 0 and 1, a bar graph (ranging from 0 to 1) shows the prevalence of the dummy variable (1=100%)

Now, let's try to graph the prevalence of women by province.



Now, let's try to graph the prevalence of women by province.

► How might we write the code?



Now, let's try to graph the prevalence of women by province.

How might we write the code? graph bar a3, over(prov)

Graph bar a3, over(prov)



► Graph bar a3, over(prov)



Does this show the prevalence of women in each province?

► Graph bar a3, over(prov)



Does this figure show the prevalence of women in each province? No - all of the bars are over 1!

Does this figure show the prevalence of women in each province? No - all of the bars are over 1!

Why is this figure different from the other? Let's look again at our gender variable.

Does this figure show the prevalence of women in each province? No - all of the bars are over 1!

Why is this figure different from the other? Let's look again at our gender variable. Codebook a3

Does this figure show the prevalence of women in each province? No - all of the bars are over 1!

Why is this figure different from the other? Let's look again at our gender variable. Codebook a3

· · · · -							
. codebook a3							
23					2 Wh	at is th	o say of
as					3. WI		ie sex of
Type:	Numeric	(byte)					
Label:	b3						
	F4 01						
Range:	[1,2]			Units: 1			
Unique values:	2			Missing .: 0/16,192			
-				-			
Tabulation	Fred	Numeric	Label				
		Numer 10	Laber				
	8,459	1	Maie				
	7,733	2	Female				
1	-						

Does this figure show the prevalence of women in each province? No - all of the bars are over 1!

Why is this figure different from the other? Let's look again at our gender variable. Codebook a3

• codebook a3		
a3	3. What is the sex of ?	The value
Type: Numeric (byte) Label: b3		labels are 1/2
Range: [1,2] Units: 1 Unique values: 2 Missing .: 0/16,192		instead of 0/1
Tabulation: Freq. Numeric Label 8,459 1 Male 7,733 2 Female		

Let's make a new dummy variable for whether someone is a woman.

Let's make a new dummy variable for whether someone is a woman.

► What's first?

Let's make a new dummy variable for whether someone is a woman.

► What's first? *generate woman*=.

Let's make a new dummy variable for whether someone is a woman.

▶ What's first? *generate woman*=.

How do we set the values to 1 for woman and 0 for man?

Let's make a new dummy variable for whether someone is a woman.

► What's first? *generate woman*=.

How do we set the values to 1 for woman and 0 for man?

- ► Replace woman=1 if a3==2
- ► Replace woman=0 if a3==1

Let's create variable label:

Let's create variable label: lab var woman "Is the household member a woman?"

Let's create variable label: lab var woman "Is the household member a woman?"

Remember, we already have the value label definition "no_yes" saved in Stata. Will this value definition work for this variable too?

Let's create variable label: lab var woman "Is the household member a woman?"

Remember, we already have the value label definition "no_yes" saved in Stata. Will this value definition work for this variable too? Yes, it would work for all dummy variables

- Let's create variable label: lab var woman "Is the household member a woman?"
- Remember, we already have the value label definition "no_yes" saved in Stata. Will this value definition work for this variable too? Yes, it would work for all dummy variables
- How can we apply this value label definition to our new "woman" variable?
Let's create variable label: lab var woman "Is the household member a woman?"

Remember, we already have the value label definition "no_yes" saved in Stata. Will this value definition work for this variable too? Yes, it would work for all dummy variables

How can we apply this value label definition to our new "woman" variable? *label values woman no_yes*

Let's look at our new variable:

Let's look at our new variable: codebook woman

Let's look at our new variable: codebook woman

woman						Is	the	household	member	а	woman?
	Type: Label:	Numeric no_yes	(float)								
Ur	Range: nique values:	[0,1] 2			Units: Missing .:	1 0/16,19	2				
	Tabulation:	Freq. 8,459 7,733	Numeric 0 1	Label No Yes							

Let's try the bar graph again with our new variable:

Let's try the bar graph again with our new variable: graph bar woman, over(prov, label(angle(45)))

Let's try the bar graph again with our new variable: graph bar woman, over(prov, label(angle(45)))



Let's try the bar graph again with our new variable: graph bar woman, over(prov, label(angle(45)))



Does it show prevalence now?

Let's try the bar graph again with our new variable: graph bar woman, over(prov, label(angle(45)))



Does it show
 prevalence now? Yes!

Let's try the bar graph again with our new variable: graph bar woman, over(prov, label(angle(45)))



Does it show prevalence now? Yes!

Which province has the highest prevalence of women in the sample?

Let's try the bar graph again with our new variable: graph bar woman, over(prov, label(angle(45)))



Does it show prevalence now? <u>Yes!</u>

Which province has the highest prevalence of women in the sample? Milne Bay province

Help windows explain how to use each code, provide different options, examples, etc.

Help windows explain how to use each code, provide different options, examples, etc.

<u>Code:</u> help [command]



Help windows explain how to use each code, provide different options, examples, etc.

Code: help [command]

Example: help tabulate twoway

Example: help tabulate twoway

► Syntax

Syntax

Two-way table

tabulate varname1 varname2 [if] [in] [weight] [, options]

- Example: help tabulate twoway
 - ► Syntax

► Options

	options	Description
- 1	Main	
	<u>ch</u> i2	report Pearson's chi-squared
	<u>e</u> xact[(#)]	report Fisher's exact test
	gamma	report Goodman and Kruskal's gamma
	<u>lr</u> chi2	report likelihood-ratio chi-squared
	<u>t</u> aub	report Kendall's tau-b
	V	report Cramér's V
	<u>cchi</u> 2	report Pearson's chi-squared in each cell
	<u>co</u> lumn	report relative frequency within its column of each cell
	row	report relative frequency within its row of each cell
	clrchi2	report likelihood ratio shi squared in each cell
	cell	report the relative frequency of each cell

Example: help tabulate twoway

► Syntax

► Options

► Examples

Two-way table of frequencies . tabulate region agecat

Include row percentages
. tabulate region agecat, row

Include column percentages
. tabulate region agecat, column

Include cell percentages
 . tabulate region agecat, cell

Now let's look at the help window for summarize





Now let's look at the help window for summarize

Code: help summarize



Now let's look at the help window for summarize

Code: help summarize

Syntax

► Options

[R]	summarize — Summa	ary statistics
	(<u>Vie</u> v	v complete PDF manual entry)
Syn	tax	
	<u>su</u> mmarize [<u>va</u>	rlist] [<u>if</u>] [<u>in</u>] [<u>weight</u>] [, options]
	options	Description
	Main	
	<u>d</u> etail	display additional statistics
	<u>mean</u> only	suppress the display; calculate only the mean; programmer's option
	format	use variable's display format
	<pre>separator(#)</pre>	<pre>draw separator line after every # variables; default is separator(5)</pre>
	display options	control spacing, line width, and base and empty cells

What if we want to know some things at the household level that need to be constructed from this person-level data?

- What if we want to know some things at the household level that need to be constructed from this person-level data?
 - Household size?

- What if we want to know some things at the household level that need to be constructed from this person-level data?
 - Household size?
 - Age of youngest household member?

- What if we want to know some things at the household level that need to be constructed from this person-level data?
 - Household size?
 - Age of youngest household member?
 - Percent of children and percent of women in the household?

- What if we want to know some things at the household level that need to be constructed from this person-level data?
 - Household size?
 - Age of youngest household member?
 - Percent of children and percent of women in the household?

We can use the *collapse* command to find these things!

Collapse creates a dataset of the summary statistics that you specify, at the level that you specify

Collapse creates a dataset of the summary statistics that you specify, at the level that you specify

Examples of summary statistics that can be used are: mean, median, sum, count, max, and min

- Collapse creates a dataset of the summary statistics that you specify, at the level that you specify
- Examples of summary statistics that can be used are: mean, median, sum, count, max, and min
- ▶ We will collapse this dataset to the household level.

- Collapse creates a dataset of the summary statistics that you specify, at the level that you specify
- Examples of summary statistics that can be used are: mean, median, sum, count, max, and min
- ▶ We will collapse this dataset to the household level.
- Which summary statistic and variable would we use to find the number of household members?

- Collapse creates a dataset of the summary statistics that you specify, at the level that you specify
- Examples of summary statistics that can be used are: mean, median, sum, count, max, and min
- ▶ We will collapse this dataset to the household level.
- Which summary statistic and variable would we use to find the number of household members? (count) a_pid

- Collapse creates a dataset of the summary statistics that you specify, at the level that you specify
- Examples of summary statistics that can be used are: mean, median, sum, count, max, and min
- ▶ We will collapse this dataset to the household level.
- Which summary statistic and variable would we use to find the age of the youngest HH member?

- Collapse creates a dataset of the summary statistics that you specify, at the level that you specify
- Examples of summary statistics that can be used are: mean, median, sum, count, max, and min
- ▶ We will collapse this dataset to the household level.
- Which summary statistic and variable would we use to find the age of the youngest HH member? (min) age

- Collapse creates a dataset of the summary statistics that you specify, at the level that you specify
- Examples of summary statistics that can be used are: mean, median, sum, count, max, and min
- ▶ We will collapse this dataset to the household level.
- Which summary statistic and variables would we use to find the percent of women & children in each HH?

- Collapse creates a dataset of the summary statistics that you specify, at the level that you specify
- Examples of summary statistics that can be used are: mean, median, sum, count, max, and min
- ▶ We will collapse this dataset to the household level.
- Which summary statistic and variables would we use to find the percent of women & children in each HH? (mean) child woman

<u>Code:</u> collapse (count) a_pid (min) age (mean) child woman, by(hhid)



- Code: collapse (count) a_pid (min) age (mean) child woman, by(hhid)
- How many observations do we have now in the dataset?
- Code: collapse (count) a_pid (min) age (mean) child woman, by(hhid)
- How many observations do we have now in the dataset? 2,699

Now, let's look at what happened in our Window 3 (Variable Window) after the collapse

Now, let's look at what happened in our Window 3 (Variable Window) after the collapse

Variables		тд×
K Filter variables h	ere	
Name	Label	
hhid	Household identifier	
pid	(count) pid	
age	(min) age	
child	(mean) child	
woman	(mean) woman	

Now, let's look at what happened in our Window 3 (Variable Window) after the collapse

Variables		▼ ₽×
Y Filter variable	es here	
Name	Label	
hhid	Household identifier	
pid	(count) pid	
age	(min) age	
child	(mean) child	
woman	(mean) woman	

Variable names remained the same

Now, let's look at what happened in our Window 3 (Variable Window) after the collapse



- Variable names remained the same
- Variable labels tell us the summary statistic used in the collapse (instead of keeping the old labels)

Let's create more meaningful variable names and labels:

Let's create more meaningful variable names and labels:

What does the variable a_pid tell us now?

Let's create more meaningful variable names and labels:

What does the variable a_pid tell us now? The count of people in the household = household size!

Let's create more meaningful variable names and labels:

What does the variable a_pid tell us now? The count of people in the household = household size!

rename a_pid hhsize

Let's create more meaningful variable names and labels:

What does the variable a_pid tell us now? The count of people in the household = household size!

rename a_pid hhsize

lab var hhsize "Household size"

Let's create more meaningful variable names and labels:

What does the variable age tell us now?

Let's create more meaningful variable names and labels:

What does the variable age tell us now? The minimum age out of the ages of each person in the HH = the age of the youngest person in the HH

Let's create more meaningful variable names and labels:

What does the variable age tell us now? The minimum age out of the ages of each person in the HH = the age of the youngest person in the HH

rename age age_youngest

Let's create more meaningful variable names and labels:

What does the variable age tell us now? The minimum age out of the ages of each person in the HH = the age of the youngest person in the HH

rename age age_youngest

lab var age_youngest "Age of the youngest household member"

Let's create more meaningful variable names and labels:

What does the variable child tell us now?

Let's create more meaningful variable names and labels:

What does the variable child tell us now? The average/mean of the dummy variable "child" = the percent of children in each household

Let's create more meaningful variable names and labels:

What does the variable child tell us now? The average/mean of the dummy variable "child" = the percent of children in each household

rename child perc_children

Let's create more meaningful variable names and labels:

What does the variable child tell us now? The average/mean of the dummy variable "child" = the percent of children in each household

rename child perc_children

lab var perc_children "Percent of children 15 years and younger in the household"

Let's create more meaningful variable names and labels:

What does the variable woman tell us now?

Let's create more meaningful variable names and labels:

What does the variable woman tell us now? The average/mean of the dummy variable "woman" = the percent of women in each household

Let's create more meaningful variable names and labels:

What does the variable woman tell us now? The average/mean of the dummy variable "woman" = the percent of women in each household

rename woman perc_women

Let's create more meaningful variable names and labels:

What does the variable woman tell us now? The average/mean of the dummy variable "woman" = the percent of women in each household

rename woman perc_women

lab var perc_women "Percent of women in the household"

Question: What is the average household size?

Question: What is the average household size?

sum hhsize

Question: What is the average household size?

► sum hhsize

. sum hhsize					
Variable	Obs	Mean	Std. dev.	Min	Max
hhsize	2,699	5.999259	2.603812	1	22



Question: What is the average household size?

► sum hhsize

. sum hhsize					
Variable	Obs	Mean	Std. dev.	Min	Max
hhsize	2,699	5.999259	2.603812	1	22



Question: What is the average age of the youngest household member?

Question: What is the average age of the youngest household member?

Sum age_youngest

Question: What is the average age of the youngest household member?

Question: What is the average age of the youngest household member?

Sum age_youngest - 8.7 years old

Question: What is the average prevalence of <u>children</u> in a household?

Question: What is the average age of the youngest household member?

Sum age_youngest - 8.7 years old

Question: What is the average prevalence of <u>children</u> in a household?

Sum perc_children

Question: What is the average age of the youngest household member?

- Question: What is the average prevalence of <u>children</u> in a household?
 - Sum perc_children 40.8% children

Question: What is the average age of the youngest household member?

- Question: What is the average prevalence of <u>children</u> in a household?
 - Sum perc_children 40.8% children
- Question: What is the average prevalence of women in a household?

Question: What is the average age of the youngest household member?

- Question: What is the average prevalence of <u>children</u> in a household?
 - Sum perc_children 40.8% children
- Question: What is the average prevalence of women in a household?
 - Sum perc_women

Question: What is the average age of the youngest household member?

- Question: What is the average prevalence of <u>children</u> in a household?
 - Sum perc_children 40.8% children
- Question: What is the average prevalence of women in a household?

Sum perc_women - 47.2% women

5. Collapse and Merge Merge (same number of observations)

What if we want to know if household size is associated with the type of roof that a household has?
- What if we want to know if household size is associated with the type of roof that a household has?
- We have household size in this dataset, but roof type is in the dataset that we worked with in Lesson 1

- What if we want to know if household size is associated with the type of roof that a household has?
- We have household size in this dataset, but roof type is in the dataset that we worked with in Lesson 1
- We can merge the two datasets together and now they have the same number and type of observations (2,699 households)

s11_long (collapsed; n=2,699)			
hhid	age_youngest		
0100RSPN	11		

s11_long (collapsed; n=2,699)			
hhid age_youngest			
0100RSPN	11		

hh_3_4 (n=2,699)				
hhid	i1			
0100RSPN	metal roof			

s11_long (collapsed; n=2,699)		hh_3_4 (n=2,	699)
hhid	age_youngest	hhid	i1
0100RSPN	11	0100RSPN	metal roof

merged data (n=2,699)					
hhid age_youngest i1					
0100RSPN	11	metal			
		roof			



Help merge

► Help merge

[D] merge — Merge datasets

(View complete PDF manual entry)

Syntax

One-to-one merge on specified key variables

merge 1:1 varlist using filename [, options]

Help merge

[D] merge — Merge datasets (<u>View complete PDF manual entry</u>) Syntax One-to-one merge on specified key variables merge 1:1 varlist using filename [, options]

We are merging <u>one</u> household to <u>one</u> household (we have exactly one observation for each household in both datasets) - so we will use "1:1" syntax

Which variable will we be using to match observations from one dataset to the other?

Which variable will we be using to match observations from one dataset to the other? *hhid* - "*Household Identifier*"

- Which variable will we be using to match observations from one dataset to the other? *hhid* - "*Household Identifier*"
- <u>Code:</u> merge 1:1 hhid using "C:\Users\RMukerjee\IFPRI Dropbox\Rishabh Mukerjee\ STATA_training_PNG2025\ Stata_training_png_March2025\datasets\hh_3_4.dta"

merge 1:1 hhid using "C:\Users\RMukerjee\IFPRI Dropbox\Rishabh Mukerjee\ STATA_training_PNG2025\ Stata_training_png_March2025\datasets\hh_3_4.dta"

<pre>. merge : > raining_PNG2025\Stata_ > the s33 dataset from [(label community already</pre>	1:1 hhid using "C:\Usen training_png_March2025 Lesson 1, matching base defined)	rs\RMukerjee\IFPRI Dropbox\R \Datasets\hh_3_4.dta" ed on hhid (household ID)	ishabh //Merge	Mukerjee\Stata_1 the dataset with
Result	Number of obs			
Not matched	0	(mongo2)		
	2,699	(_merge==>)		

merge 1:1 hhid using "C:\Users\RMukerjee\IFPRI Dropbox\Rishabh Mukerjee\ STATA_training_PNG2025\ Stata_training_png_March2025\datasets\hh_3_4.dta"

> r > (la	merg raining_PNG2025\Stat the s33 dataset fro abel community alrea	ge 1:1 hhid using "C:\Use ta_training_png_March2025 om Lesson 1, matching bas ady defined)	rs\RMukerjee\IFPRI Dropbox\Risha \Datasets\hh_3_4.dta" //Me ed on hhid (household ID)	bh Mukerjee\Stata_ rge the dataset wit
	Result	Number of obs		
	Not matched	0		
	Matched	2,699	(_merge==3)	

Stata output tells us how many observations were and were not matched. How many were matched?

merge 1:1 hhid using "C:\Users\RMukerjee\IFPRI Dropbox\Rishabh Mukerjee\ STATA_training_PNG2025\ Stata_training_png_March2025\datasets\hh_3_4.dta"

 merge raining_PNG2025\Stata the s33 dataset from (label community alread) 	e 1:1 hhid using "C:\User a_training_png_March2025\ m Lesson 1, matching base dy defined)	rs\RMukerjee\IFPRI Dropbox \Datasets\hh_3_4.dta" ed on hhid (household ID)	\Rishabh //Merge	Muke the	erjee∖St dataset	ata_1 with
Result	Number of obs					
Not matched Matched	0 2,699	(_merge==3)				

Stata output tells us how many observations were and were not matched. How many were matched? 2,699 (all of them!)

merge 1:1 hhid using "C:\Users\RMukerjee\IFPRI Dropbox\Rishabh Mukerjee\ STATA_training_PNG2025\ Stata_training_png_March2025\datasets\hh_3_4.dta"

<pre>. merg > raining_PNG2025\Stat > the s33 dataset fro (label community alrea</pre>	ge 1:1 hhid using "C:\Use a_training_png_March2025 om Lesson 1, matching bas ady defined)	rs\RMukerjee\IFPRI Dropbo \Datasets\hh_3_4.dta" ed on hhid (household ID)	x\Rishabh //Merge	Mukerjee\Stata_ the dataset wit
Result	Number of obs			
Not matched Matched	0 2,699	(_merge==3)		

With every merge, Stata creates a new variable called _merge

merge 1:1 hhid using "C:\Users\RMukerjee\IFPRI Dropbox\Rishabh Mukerjee\ STATA_training_PNG2025\ Stata_training_png_March2025\datasets\hh_3_4.dta"

 merge raining_PNG2025\Stata the s33 dataset from (label community alread) 	e 1:1 hhid using "C:\Use a_training_png_March2025 n Lesson 1, matching bas dy defined)	rs\RMukerjee\IFPRI Dropbox\Rishabh \Datasets\hh_3_4.dta" //Merge ed on hhid (household ID)	Mukerjee\Stata_ the dataset wit
Result	Number of obs		
Not matched	0		
Matched	2,699	(_merge==3)	

- With every merge, Stata creates a new variable called _merge
 _merge equals 1 or 2 for un-matched observations
 - _merge equals 3 for matched observations

merge 1:1 hhid using "C:\Users\RMukerjee\IFPRI Dropbox\Rishabh Mukerjee\ STATA_training_PNG2025\ Stata_training_png_March2025\datasets\hh_3_4.dta"

. merge 1:1 hhid using "C:\Users\RMukerjee\IFPRI Dropbox\Rishabh Mukerjee\Stata_1
> raining_PNG2025\Stata_training_png_March2025\Datasets\hh_3_4.dta" //Merge the dataset with
> the s33 dataset from Lesson 1, matching based on hhid (household ID)
(label community already defined)

Result Number of obs
Not matched 0
Matched 2,699 (_merge==3)

With every merge, Stata creates a new variable called _merge

_merge equals 1 or 2 for un-matched observations

_merge equals 3 for matched observations

If all matched, it's a good practice to drop _merge

Question: On average, do households with a thatched roof or a corrugated roof have more children?

Question: On average, do households with a thatched roof or a corrugated roof have more children?

Option 1: two "sum...if" codes

Question: On average, do households with a thatched roof or a corrugated roof have more children?

Option 1: two "sum...if" codes

sum perc_children if i1==1

sum perc_children if i1==2

Question: On average, do households with a thatched roof or a corrugated roof have more children?

Option 1: two "sum...if" codes

sum perc_children if i1==1 - on average in households with thatched roofs, <u>43%</u> of the household members are aged 15 years or younger

sum perc_children if i1==2 - on average in households with corrugated metal roofs, <u>38%</u> of the household members are aged 15 years or younger

Question: On average, do households with a thatched roof or a corrugated roof have more children?

Option 2: bysort

Bysort repeats a Stata command on a subset of data

Question: On average, do households with a thatched roof or a corrugated roof have more children?

Option 2: bysort

- Bysort repeats a Stata command on a subset of data
- We can repeat the "summarize perc_children" command on all of the different values of i1

Question: On average, do households with a thatched roof or a corrugated roof have more children?

Option 2: bysort

Bysort repeats a Stata command on a subset of data

- We can repeat the "summarize perc_children" command on all of the different values of i1
- Code: bysort i1: sum perc_children

Question: On average, do households with a thatched roof or a corrugated roof have more children?

bysort i1: sum perc_children

bysort i1: sum perc_children						
-> i1 = Thatched	d roof (gras	s, palm fron	ds, sago leave	s, etc.)		
Variable	Obs	Mean	Std. dev.	Min	Max	
perc_child~n	1,387	.43039	.2671725	0	1	
-> i1 = Corrugat	ted metal ro	of				
Variable	Obs	Mean	Std. dev.	Min	Max	
perc_child~n	1,209	.3842051	.2353648	0	1	

Question: On average, do households with a thatched roof or a corrugated roof have more children?



- Question: On average, do households with a thatched roof or a corrugated roof have more children?
- Answer: Descriptive statistics suggest that households with thatched roofs have a higher percentage of children than households with corrugated metal roofs.

What if we want to know how many people in our dataset live in a house with a thatched roof?

- What if we want to know how many people in our dataset live in a house with a thatched roof?
- We would want to look at the roof type data at the person-level

- What if we want to know how many people in our dataset live in a house with a thatched roof?
- We would want to look at the roof type data at the person-level
- We can merge the housing characteristics dataset to the household roster dataset, maintaining the dataset at the person-level (16,192 people)

Hh_3_4 (n=2,699)			
hhid	i1		
0100RSPN	metal roof		



s11_long (n=16,192)							
hhid		a_pid	a4				
0100RSPN	$ \rangle$	1	Μ				
0100RSPN		2	F				
0100RSPN		3	F				
0100RSPN		4	F				
0100RSPN		5	Μ				
0100RSPN		6	Μ				
0100RSPN		7	Μ				

				<u>s11_long (n=16,192)</u>				
				hhid	a_pid	a4		
	2 1 1			0100RSPN	1	Μ		
Hn_3_4 (n=2,699)				0100RSPN	2	F		
hhid		i1		0100RSPN	3	F		
		motal rea	£	0100RSPN	4	F		
UTUUKSPN		metal 100		0100RSPN	5	М		
				0100RSPN	6	М		
				0100RSPN	7	М		
	merged dat	a (n=6,057)						
hhid	a_pid	a4	i1					
0100RSPN	1	Μ	metal roof					
0100RSPN	2	F	metal roof					
0100RSPN	3	F	metal roof					
0100RSPN	4	F	metal roof					
0100RSPN	5	Μ	metal roof					
0100RSPN	6	Μ	metal roof					
0100RSPN	7	Μ	metal roof			/		



Let's start with the household characteristics dataset this time
- Let's start with the household characteristics dataset this time
 - Close out of Stata

Let's start with the household characteristics dataset this time

Close out of Stata

Double click on hh_3_4.dta to open this file

- Let's start with the household characteristics dataset this time
 - Close out of Stata
 - ▶ Double click on *hh_3_4.dta* to open this file
 - What level are the observations in this dataset?

- Let's start with the household characteristics dataset this time
 - Close out of Stata
 - ▶ Double click on *hh_3_4.dta* to open this file
 - What level are the observations in this dataset? Householdlevel

- Let's start with the household characteristics dataset this time
 - Close out of Stata
 - ▶ Double click on *hh_3_4.dta* to open this file
 - What level are the observations in this dataset? Householdlevel
- We are merging it with s11_long.dta (the household roster file we have been working with today). What level are the observations in this dataset?

- Let's start with the household characteristics dataset this time
 - Close out of Stata
 - ▶ Double click on *hh_3_4.dta* to open this file
 - What level are the observations in this dataset? Householdlevel
- We are merging it with s11_long.dta (the household roster file we have been working with today). What level are the observations in this dataset? *Person-level*

Before, our merge command was: merge 1:1 hhid using [datafile]

- Before, our merge command was: merge 1:1 hhid using [datafile]
- What part of the command do you think we will change for this merge?

- Before, our merge command was: merge 1:1 hhid using [datafile]
- What part of the command do you think we will change for this merge? Change 1:1 to 1:m

- Before, our merge command was: merge 1:1 hhid using [datafile]
- What part of the command do you think we will change for this merge? Change 1:1 to 1:m
- It's now "1:m" because we merging ONE household observation (in the HH characteristics file) to MANY household observations (one per person in the household roster file)



<u>Code:</u> merge 1:1 hhid using "C:\Users\RMukerjee\IFPRI Dropbox\Rishabh Mukerjee\ STATA_training_PNG2025\ Stata_training_png_March2025\datasets\hh_3_4.dta"

- Code: merge 1:1 hhid using "C:\Users\RMukerjee\IFPRI Dropbox\Rishabh Mukerjee\ STATA_training_PNG2025\ Stata_training_png_March2025\datasets\hh_3_4.dta"
- We are still using *hhid* as our matching identifier across the two datasets

merge 1:1 hhid using "C:\Users\RMukerjee\IFPRI Dropbox\Rishabh Mukerjee\ STATA_training_PNG2025\ Stata_training_png_March2025\datasets\hh_3_4.dta"

Result	Number of obs	
Not matched	0	
Matched	16,192	(_merge==3)

merge 1:1 hhid using "C:\Users\RMukerjee\IFPRI Dropbox\Rishabh Mukerjee\ STATA_training_PNG2025\ Stata_training_png_March2025\datasets\hh_3_4.dta"

Result	Number of obs	
Not matched	0	
Matched	16,192	(_merge==3)

• <u>Question:</u> Did all of the observations match?

merge 1:1 hhid using "C:\Users\RMukerjee\IFPRI Dropbox\Rishabh Mukerjee\ STATA_training_PNG2025\ Stata_training_png_March2025\datasets\hh_3_4.dta"

Result	Number of obs	
Not matched	0	
Matched	16,192	(_merge==3)

Question: Did all of the observations match? Yes (all 16,192)

merge 1:1 hhid using "C:\Users\RMukerjee\IFPRI Dropbox\Rishabh Mukerjee\ STATA_training_PNG2025\ Stata_training_png_March2025\datasets\hh_3_4.dta"

Result	Number of obs	
Not matched	0	
Matched	16,192	(_merge==3)

- Question: Did all of the observations match? Yes (all 16,192)
- Question: What are the values of the new _merge variable?

merge 1:1 hhid using "C:\Users\RMukerjee\IFPRI Dropbox\Rishabh Mukerjee\ STATA_training_PNG2025\ Stata_training_png_March2025\datasets\hh_3_4.dta"

Result	Number of obs
Not matched	0
Matched	16,192 (merge==3

Question: Did all of the observations match? Yes (all 16,192)

Question: What are the values of the new _merge variable? All 3s

Question: How many people in the dataset live in houses with a thatched roof?

Question: How many people in the dataset live in houses with a thatched roof?



Question: How many people in the dataset live in houses with a thatched roof?



- Question: How many people in the dataset live in houses with a thatched roof?
- Code: tab i1



- Question: How many people in the dataset live in houses with a thatched roof?
- ► <u>Code:</u> tab i1
- Answer: 8,328 people (51.43% of people in the sample live in houses with a thatched roof)

Question: How many women in the dataset live in houses with a thatched roof?

Question: How many women in the dataset live in houses with a thatched roof?



- Question: How many women in the dataset live in houses with a thatched roof?
- Code:
 - ▶ tab i1 if a3==2
 - Can use codebook a3 to remember how "woman" is coded
 - ▶ tab i1 a3, column
 - ▶ tab a3 i1, row

Question: How many women in the dataset live in houses with a thatched roof?

► <u>Code:</u>

▶ tab i1 if a3==2

Can use codebook a3 to remember how "woman" is coded

▶ tab i1 a3, column

▶ tab a3 i1, row

Answer:

Question: How many women in the dataset live in houses with a thatched roof?

► <u>Code:</u>

- ▶ tab i1 if a3==2
 - Can use codebook a3 to remember how "woman" is coded
- ▶ tab i1 a3, column
- ▶ tab a3 i1, row

Answer: 3,961 women live in a house with a thatched roof (51.22% of women in the sample live in a house with a thatched roof)