The corresponding video tutorials are available online:

https://www.youtube.com/playlist?list=PL7BcpOtSe5u_zQctYXz4ee79Zc9r4mfnr



Tutorial #8 Missing values and filter missings

1

pairfam tutorial

8. Missing values and filter missings

Kristin Hajek & Madison Garrett, May 2022

The following tutorial explains the usage of missing values and filter missings.



Missing values

Negative values = missing values

» Exception: kldb2010, isco08, mcs, pcs

» -1, ..., -12, -77

No system missings

» Exception: Identification variables and preload variables

Recoding to system missings

» mvdecode _all , mv(-1=.a \-2=.b \-3=.c)

» **recode** variablex -1=.a -2=.b -3=.c

How does pairfam manage missing values?

Missing values in the pairfam data are typically coded with negative numbers.

Only four generated variables have valid negative values: *kldb2010*, *isco08*, *mcs*, and *pcs*.

Missing values range from -1 to -12. Additionally, *mcs* and *pcs* also contain the missing code -77.

With the exception of identification variables (such as *id* or *pid*) and preload variables, the data include no system missings.

Negative missing values can be easily recoded into system missings with the Stata command *mvdecode* or *recode*. Alternatively, analyses can also be restricted to positive values only.

			pairfam Prond Knalpin of Intimate Millioninips and Earling Of
Miss	sing v	alues	
	-1	Don't know	
	-2	No answer	
	-3	Does not apply	
	-4	Filter error / Incorrect entry	
	-5	Inconsistent value	
	-6	Unreadable answer	
	-7 / -77	Incomplete data	
	-9	Invalid multiple answer	
	-10 / -11	Not in DemoDiff / Not in pairfam	
	-12	Non-response PAPI	

Why are there so many different missing values?

The missing values -1 "Don't know" and -2 "No answer" were assigned if the respondent could not or did not want to answer a question. These codes are the only missing values documented in the questionnaire.

The value -3 "Does not apply" was assigned if a respondent had not been posed the corresponding question, meaning the respondent was "filtered over".

Errors in the CAPI program that erroneously guided respondents to the wrong questions in the interview are indicated by the missing code -4 "Filter error / Incorrect entry", as are incorrect data entries made by interviewers. In waves 12 and 13, the value -4 was also assigned if the respondent incorrectly responded (or incorrectly did not respond) to a question in PAPI mode.

In order to detect inconsistencies between responses, we checked for logically impossible or empirically implausible combinations of values for two or more variables. Inconsistent values were then coded to -5 "Inconsistent value" if it was clear that the value was implausible.

For open answers that were not legible, the value -6 "Unreadable answer" was assigned.

For generated data variables and files, the value -7 "Incomplete data" is used to indicate cases for which the necessary information was lacking to compute a valid value. For the generated variables *mcs* and *pcs*, the value -77 "Incomplete data" is

used.

In the alteri data sets, the value -9 "Invalid multiple answer" was assigned if the respondent checked more than the permitted number of options. For waves 1 and 3, special missing codes indicate differences between the *pairfam* questionnaire and the *DemoDiff* sample. If a question from the *pairfam* questionnaire was not part of *DemoDiff*, the corresponding variable was set to -10 "Not in *DemoDiff*". *DemoDiff* variables not included in the *pairfam* questionnaire are indicated by -11 "Not in *pairfam*".

In waves 12 and 13, anchor respondents interviewed via telephone due to COVID-19 restrictions who did not return the PAPI questionnaire received the missing value -12 "Non-response PAPI" for these questions.

Please be aware that the values -10, -11, and -12 are also part of the classification scheme for *kldb2010* and *isco08*.

Question 59 Variable sin2	Is there anyone you are interested in? Yes	
Responden	No answer	

What exactly are filter missings?

Let's look at an example. Single respondents were asked the question "Is there anyone you are interested in?".

This is a screenshot from the anchor codebook. On the left you can see the question number and the corresponding variable name.

Here you can see the question wording and the available answers as well as their variable values.

If a filter was implemented for a question, there is a box underneath the question stating which respondents were posed the question of interest – in text and numerical form. The variables in parentheses are the ones used for filtering.

In this example, this question was posed only to respondents with no reported relationship in the last three months.

As you can see, auxiliary variables, such as *hp* and *hps2*, can help you to understand question filters and facilitate your analyses.

It is helpful to consult the codebooks and check question filtering to better understand a variable's values.

Filter missings:	Examp	le 1		
<pre>. use anchor12_capi, clea . tab sin2 hp, m</pre>	ır			
· cab sinz np, m				
Interested in	Respondent has			
potential partner	partn	er		
(Qu. 59)	0 No	1 Yes	Total	
-4 Filter error / In	4	0 1	4	
-3 Does not apply	70	3,909	3,979	
-2 No answer	16	0	16	
-1 Don't know	47	0 [47	
1 Yes	617	0 [617	
2 No	1,131	0 1	1,131	
Total	1,885	3,909	5 , 794	

To further this example, let's look at the anchor12 capi data set in Stata.

Here I have tabulated the variables sin2 and hp.

Respondents who have a partner receive the missing code -3 "Does not apply" for *sin2* as they were not posed this question.

Only a small portion of participants, the ones who were single at the time of the interview for wave 12, answered this question.

Filter missings: Example 2 . tab sdp17 hp, m Partner: Highest level Respondent has of school education partner (since previous wave) 0 No 1 Yes Total -4 Filter error / Inc 0 2 2	
Partner: Highest level Respondent has of school education partner (since previous wave) 0 No 1 Yes Total	
of school education partner (since previous wave) 0 No 1 Yes Total	
of school education partner (since previous wave) 0 No 1 Yes Total	
(since previous wave) 0 No 1 Yes Total	
-4 Filter error / Inc 0 2 2	
-4 Filter error / Inc 0 2 2	
0.5	
-3 Does not apply 1,885 3,493 5,378	
-2 No answer 0 1 1	
-1 Don't know 0 13 13	
4 Secondary general 0 41 41	
5 Intermediate schoo 0 146 146	
7 Certificate from m 0 11 11	
8 Entrance qualifica 0 31 31	
9 General or subject 0 168 168	
10 Other school leav 0 3 3	

The next example illustrates why consulting the codebooks to understand the filtering of certain questions can be important, and how generated variables can be useful for your analyses.

Here, I tabulated the partner's highest level of school education – variable sdp17.

Naturally, anchor respondents without a partner all have the missing code -3 "Does not apply".

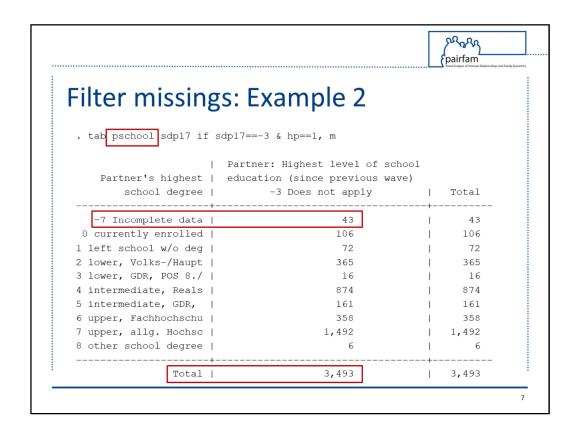
However, there are also almost 3,500 respondents who do have a partner with this missing value as well.

If you read the question closely, you can see that the variable *sdp17* asks for the highest level of school education since the previous wave.

The filter is defined as follows: "Respondents with partner from previous wave who received a school leaving certificate since the previous wave, or with new partner who received a school leaving certificate".

As the goal is not to overburden respondents, they are asked to only report their partner's highest school education *if* there has been a change since the previous wave

Does this mean that you will have to check each wave of the anchor data and copy the highest level of school education to all the other waves? No!

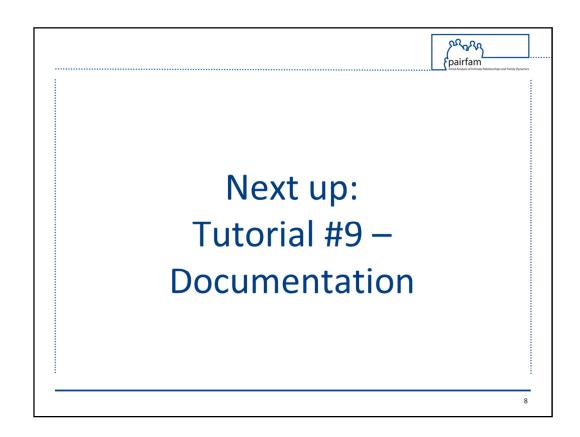


Because we want to make data analysis as easy as possible for our users, the generated variable *pschool* sums up all information on the partner's highest school degree across all waves.

Here, I tabulated *pschool* for the almost 3,500 respondents who had the missing value -3 for *sdp17*.

As you can see, only 43 respondents still have missing values while the rest have valid codes.

We highly recommend consulting the generated variables delivered in the data, as they really can facilitate data analysis!



We've reached the end of the eighth tutorial.

Tutorial number 9 will present the available documentation included in the Scientific Use File.