

# Data Cleaning Guidelines

## for 2<sup>nd</sup> Phase FATE-Surveys: Stata Syntax

*Abridged version*

Author: Patrick Illien<sup>1</sup>

### Contents

Introduction .....	2
Copying Dataset .....	3
Documenting Changes and Getting Started.....	3
Deleting Cases .....	4
Tidying Dataset.....	5
Checking and Correcting Identifiers .....	6
Adding new Identifiers.....	7
Consistency Checks.....	9
Consistency Checks for Part 1.....	9
Consistency Checks for Part 2.....	14
Dealing with Missing Data.....	15
Dealing with Outliers.....	16
Labelling.....	16
Anonymising Data.....	19
Backup.....	20

---

<sup>1</sup> I would like to thank Dr. Christoph Bader and Dr. Maurice Tschopp for valuable inputs.

## Introduction

These guidelines have been developed for the 2<sup>nd</sup> phase FATE surveys to ensure a minimum standard for comparability. The present document has been inspired by the following sources, which contain useful tips and should be consulted for more information:

- ACAPS, 2016. *Data cleaning: technical brief*. Available [here](#).
- Beaver, M., 2012. *Survey Data Cleaning Guidelines: (SPSS and Stata)*. Available [here](#).
- Data Science Primer, 2018. *Chapter 3: Data Cleaning*. Available [here](#).
- IFPRI (International Food Policy Research Institute), 2018. *A guide to data cleaning using Stata*. Available [here](#).

Most of this guide is presented in tables. In each one those, the left columns describe the steps to be taken, whereas the right columns show the corresponding Stata syntax. All Stata codes were created and tested using version 15.1. The guidelines have a logical order and each step should be undertaken after the other because some syntax only works if the previous ones have already been applied successfully.

Of course, the better the survey is designed, the less errors will be made later. Investing time in survey design, coding and enumerator training at the start pays!

Unless otherwise stated, mark problematic observations, don't change values – it is up to the analysis how to deal with them!

## Copying Dataset

- Perform all steps for each dataset (wide and all longs of part 1 & 2).
- **The first thing to do is to make a copy of the original data file (a simple copy without opening the raw data in a programme)! ALWAYS keep the source files in a separate folder called “RAW DATA” and change its attribute to READ-ONLY, to avoid modification of any of the files. Then encrypt the raw-data files and don’t forget to back them up as well.**
- Never touch the original data and only work in the copied file from now on, even if only importing the data into a programme. Discuss with your supervisors if the original data (which are not anonymised) shall be deleted after data analysis is finished.

## Documenting Changes and Getting Started

- Perform all steps in each dataset (wide and all longs of part 1 & 2).
- Be sure to always document any changes you make directly in a R, Stata or SPSS when using any of these software packages. If you use any other software, document any changes in an excel sheet variable by variable.

What?	How? (STATA syntax)
<ul style="list-style-type: none"> <li>• Make one change log per dataset and keep the change logs organised.</li> <li>• Each change log contains all modifications to that dataset and will serve as an audit trail. It will allow a return to the original value if required. Within the change log, store the following information:             <ul style="list-style-type: none"> <li>○ At the top of the file:                 <ul style="list-style-type: none"> <li>▪ Country and survey year (if necessary also survey type)</li> <li>▪ Dataset</li> <li>▪ Date of file status</li> <li>▪ Author of file</li> <li>▪ Purpose of file</li> <li>▪ Comments</li> </ul> </li> <li>○ Throughout:                 <ul style="list-style-type: none"> <li>▪ Variable or observation concerned</li> <li>▪ What the change was (e.g. old and new value)</li> <li>▪ Comments/justification</li> </ul> </li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>• Always document your comments and commands in do-files, the are the change log when working with Stata.</li> <li>• You can organise your do-files in Stata’s project manager. Open in from within the do-file editor on Windows and Unix or from within the main menu on Mac (select File&gt;New&gt;Project....)</li> <li>• Create one master do-file per dataset and name it accordingly: type doedit into the command window to create a do-file where you write your commands. You can add comments after * at the start of a line. Save the do-file into the project manager.</li> <li>• Add all the relevant information at the top as described in the left column.</li> <li>• Apply the structure and cleaning process of this guide and add section titles accordingly. Take over all the relevant commands from each section in this guide. In each section, the datasets in which certain steps have to be performed (i.e. which have to be added to the do-file) are specified in <b>colour</b>.</li> <li>• Describe and justify each change you make to the data.</li> <li>• All the other information will be readable from the commands themselves but it is always good to add some descriptions.</li> <li>• Repeat the entire process for each dataset (all wide and longs of part 1 &amp; 2).</li> </ul>
<ul style="list-style-type: none"> <li>• Increase the number of variables allowed by the software</li> </ul>	<ul style="list-style-type: none"> <li>• <i>clear</i> <i>set maxvar 10000</i></li> </ul>
<ul style="list-style-type: none"> <li>• Increase the size of the results window display</li> </ul>	<ul style="list-style-type: none"> <li>• <i>set scrollbufsize 700000</i></li> </ul>

<ul style="list-style-type: none"> <li>Specify the directory/path to the copy of the raw data</li> </ul>	<ul style="list-style-type: none"> <li>Use the cd command, for example: <i>cd "Z:\Data_Cleaning\Stata"</i></li> </ul>
<ul style="list-style-type: none"> <li>Import the copy of your raw data into your software package</li> </ul>	<ul style="list-style-type: none"> <li>Use the import command with clear in order to start from a clean slate, make sure to specify the bindquotes option with strict as else it might mess up how it delimites double quotes, for example: <i>import delimited Long_format_22.10/FATE_Rwanda_Part_1_221018_final.csv, bindquotes(strict) clear</i></li> </ul>

## Deleting Cases

- Perform these steps only for the wide and long\_final datasets of part 1 & 2.
- As a general rule, you should not delete any data!
- The only exception is if all or almost all entries (including key demographics) are blank for a household. This may be the case if the enumerators wrongly submitted a trial run or an aborted interview. We do this step at the start so that the faulty submissions don't add confusion when doing all subsequent cleaning steps (e.g. adding new identifiers, doing consistency checks, etc.).

What?	How? (STATA syntax)
<ul style="list-style-type: none"> <li>Delete line if there are no observations ("blank row")</li> </ul>	<ul style="list-style-type: none"> <li><i>egen nmiss = rowmiss(_all) display c(k) drop if nmiss==c(k)-1</i></li> <li>Alternatively you can install the "missings" package and type the following: <i>missings dropobs _all, force display r(n_dropped)</i></li> </ul>
<ul style="list-style-type: none"> <li>Delete line if almost all entries are blank</li> </ul>	<ul style="list-style-type: none"> <li><i>ds, return list local varlist=r(varlist) egen count_miss=rowmiss(`varlist') capture describe generate count_miss_prop=count_miss/r(k)*100 list count_miss_prop list submissiondate code_id village today_date a2_nameprim key if count_miss_prop&gt;99 count if count_miss_prop&gt;99 *check the list and double check (if necessary in the dataset) to see if you really want to delete those observations drop if count_miss_prop&gt;99.5 count</i></li> <li>You can also check the percentage of missing observations like this: <i>missings report, observations percent</i></li> </ul>

## Tidying Dataset

- Perform all steps in each dataset (wide and all longs of part 1 & 2).
- Make sure that the dataset is tidied up as follows

What?	How? (STATA syntax)
<ul style="list-style-type: none"> <li>• Fonts have been harmonised</li> </ul>	<ul style="list-style-type: none"> <li>• Automatically done in Stata when importing the csv file.</li> </ul>
<ul style="list-style-type: none"> <li>• Leading and trailing spaces of string variables have been deleted</li> </ul>	<ul style="list-style-type: none"> <li>• <pre>ds, has(type string) foreach var of varlist `r(varlist)' {   replace `var' = strtrim(`var') }</pre></li> </ul>
<ul style="list-style-type: none"> <li>• Put all string variables into lower case</li> </ul>	<ul style="list-style-type: none"> <li>• <pre>ds, has(type string) foreach var of varlist `r(varlist)' {   replace `var' = lower(`var') }</pre></li> </ul>
<ul style="list-style-type: none"> <li>• Put all variable names into lower case and delete spaces</li> </ul>	<ul style="list-style-type: none"> <li>• We don't need to delete spaces as Stata doesn't allow them in variable names anyways, only put in lower case: <pre>rename *, lower</pre></li> </ul>
<ul style="list-style-type: none"> <li>• Replace all the NA and DK and convert only the variables with single numbers to numerics</li> </ul>	<ul style="list-style-type: none"> <li>• <pre>ds, has(type string) foreach var of varlist `r(varlist)' {   replace `var' = ".a" if `var' == "na" &amp; length(`var') == 2   replace `var' = ".b" if `var' == "dk" &amp; length(`var') == 2 }  ds, has(type string) foreach var of varlist `r(varlist)' {   count if strmatch(`var', "*" *) == 1   regexm(`var', "[a-zA-Z]") == 1 &amp; regexm(`var', "\.a") == 0 &amp; regexm(`var', "\.b") == 0   if r(N) == 0 {     destring `var', replace   } }</pre></li> </ul>
<ul style="list-style-type: none"> <li>• Once, all the NA and DK have been changed, check for inconsistencies in alphabetic responses or categories by showing all string responses containing any alphabetical letters</li> </ul>	<ul style="list-style-type: none"> <li>• <pre>ds, has(type string) foreach var of varlist `r(varlist)' {   list `var' if `var' != "" &amp; regexm(`var', "[a-zA-Z]") == 1 }</pre></li> </ul>
<ul style="list-style-type: none"> <li>• Go through the entire list and check for mislabeled categorical string labels, i.e. separate classes that should really be the same, e.g. If there is kgs' and 'kilo' and "kg", you should combine them. Modify this code to correct all mislabeled categorical variables</li> </ul>	<ul style="list-style-type: none"> <li>• <pre>ds, has(type string) foreach var of varlist `r(varlist)' {   replace `var' = "kg" if `var' == "kgs"   `var' == "kilo"     `var' == "kilos" }</pre></li> </ul>
<ul style="list-style-type: none"> <li>• Check for remaining inconsistencies (adjust relevant variable names according to dataset)</li> </ul>	<ul style="list-style-type: none"> <li>• You can exclude variables that you don't want to check as follows (this depends on the dataset): <pre>preserve drop submissiondate a13_codeenumerator village starttime today_date a2_nameprim fam_name* end_time instanceid key</pre></li> </ul>

	<pre> ds, has(type string) foreach var of varlist `r(varlist)' {   list `var' if `var'!="" &amp; `var'!="kg" &amp; regexm(`var',"[a-zA-Z]")==1 } restore </pre>
<ul style="list-style-type: none"> <li>List variables containing multiple answers</li> </ul>	<ul style="list-style-type: none"> <li> <pre> ds, has(type string) foreach var of varlist `r(varlist)' {   count if strmatch(`var',"*")==1 &amp; regexm(`var',"[a-zA-Z]")==0   if r(N)&gt;0 {     list `var' if `var'!=""   } } </pre> </li> </ul>
<ul style="list-style-type: none"> <li>List only the multiple answer observations</li> </ul>	<ul style="list-style-type: none"> <li> <pre> ds, has(type string) foreach var of varlist `r(varlist)' {   list `var' if strmatch(`var',"*")==1 &amp; regexm(`var',"[a-zA-Z]")==0 &amp; `var'!="" } </pre> </li> </ul>
<ul style="list-style-type: none"> <li>Show all string responses that contain any number</li> </ul>	<ul style="list-style-type: none"> <li> <pre> ds, has(type string) foreach var of varlist `r(varlist)' {   list `var' if regexm(`var',"[0/9]")==1 } </pre> </li> </ul>
<ul style="list-style-type: none"> <li>Make sure that there are no variables containing only single numbers that are characterised as strings. The following commands should yield no result</li> </ul>	<ul style="list-style-type: none"> <li> <pre> ds, has(type string) foreach var of varlist `r(varlist)' {   count if strmatch(`var',"*")==1   regexm(`var',"[a-zA-Z]")==1 &amp; regexm(`var',"\.a")==0 &amp; regexm(`var',"\.b")==0   if r(N)==0 {     list `var' if `var'!=""   } } </pre> </li> </ul>
<ul style="list-style-type: none"> <li>Save everything you have done up to now under the name clean_draft</li> </ul>	<ul style="list-style-type: none"> <li>For example: <pre> save Wide_format_22.10/FATE_Rwanda_Part_1_221018_final_WIDE_clean_draft, replace clear </pre> </li> </ul>

## Checking and Correcting Identifiers

- Perform all steps in the wide and long\_final datasets of part 1 & 2).
- Each household will receive its unique identifier. We therefore have to first make sure that the code\_id is attached to the correct household and that it is the same household for a specific code\_id in part 1 and part 2. We will merge certain key characteristics of part 1 and 2 to check if the households match but we will continue working on the separate datasets after.

What?	How? (STATA syntax)
<ul style="list-style-type: none"> <li>First correct any wrong IDs of part 1, e.g. if they exist twice, manually in accordance with your sample frame. Part 1 code_id must be unique so that we can match part 1 and 2 of the survey correctly for each household. If you have discovered some manifest errors in the identity of the respondent, such as wrong code_ids or gender, correct them directly in the wide and long datasets of part 1 and 2 via the use command in your</li> </ul>	<ul style="list-style-type: none"> <li>Run do-files provided by the project</li> </ul>

do-file. Inspect households that were not available and have most entries missing (avail!=1) and mark them with “delete”.	
<ul style="list-style-type: none"> <li>Then prepare part 2 for the merger. Inspect households that were not available and have most entries missing (avail!=1) and mark them with “delete”. You don’t need to check anything here, we will do that in the merged file.</li> </ul>	
<ul style="list-style-type: none"> <li>Merge the key identifiers of part 1 and 2 in a new dataset to check if each household is correct</li> <li>Now undertake as many cross-checks as possible in order to ensure that each household of part 1 corresponds to the correct household of part 2. Check any inconsistencies manually with the sampling frame and describe the problem. This is the most tedious cleaning step and might take some time.</li> <li>If you have discovered some manifest errors in the identity of the respondent, such as wrong code_ids or gender, correct them directly in the relevant wide and long datasets of part 1 and 2 via the use command in your do-file. Other inconsistencies (such as different household types or agricultural production), do not have to be changed as it is up to the analyst how to deal with them (they might exclude them or make certain assumptions such as that every household has a vegetable garden). You should nevertheless describe these inconsistencies in a new variable called “remarks”. They might in part be based on different interpretations of the question.</li> <li>Make sure to save all the keys if the entire code_id will be deleted in a new dataset that you can later use for the long do-files of part 1 (other than long_final), see below.</li> </ul>	
<ul style="list-style-type: none"> <li>Now merge the corrected datasets and run all checks again to see if they now yield the correct results.</li> </ul>	

## Adding new Identifiers

- Perform all steps in each dataset (wide and all longs of part 1 & 2).
- When you are sure that all your households are identified correctly and the data cleaning has been done, create new unique household identifiers in the following way (we will not create unique identifiers for each individual). This is to create simple, totally anonymised and unique identifiers across all datasets. We do this step before the consistency checks so that we have already corrected any possibly wrong household entries.

What?	How? (STATA syntax)
<ul style="list-style-type: none"> <li>Create new unique identifiers: <ul style="list-style-type: none"> <li>Country code / last two digits of year in which survey took place/random consecutive three-digit household number (preceded by zeros as</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>Run do-files provided by the project</li> <li>This is the key step in it creating new unique identifiers: <pre> set seed 2345 gen random=runiform() sort code_id, stable by code_id: gen group_random=random[1] egen sequence=group(group_random), missing gen survey="418" tostring sequence, replace format(%03.0f) gen new_identifier=survey+sequence destring new_identifier, replace sort new_identifier </pre> </li> </ul>

<i>placeholders if necessary)</i>	
<ul style="list-style-type: none"> <li>• Create an identifier template that you can use after to add the unique identifiers into each dataset.</li> </ul>	<ul style="list-style-type: none"> <li>• This step is included in the do-files provided by the project:  <pre>preserve keep code_id key_part1 dup_key_part1 new_identifier order new_identifier, first duplicates drop if dup_key_part1&gt;0 drop dup_key_part1 rename key_part1 key save Identifiers/New_identifiers_part1_key, replace restore</pre> </li> </ul>
<ul style="list-style-type: none"> <li>• Make sure to use the correct file</li> </ul>	<ul style="list-style-type: none"> <li>• For example:  <pre>use Wide_format_22.10/FATE_Rwanda_Part_1_221018_final_WID E_idcheck, clear</pre> </li> </ul>
<ul style="list-style-type: none"> <li>• Merge unique identifier variables into all datasets using the key variable</li> </ul>	<ul style="list-style-type: none"> <li>• <b>Part 1, wide and long final formats:</b>  <pre>merge m:1 key using "Identifiers/New_identifiers_part1_key" order new_identifier, first list if new_identifier==. /*this list should be empty*/ list if _merge!=3 /*this list should be empty*/ save Wide_format_22.10/FATE_Rwanda_Part_1_221018_final_WID E_idcheck, replace</pre> </li> <li>• <b>Part 1, all other long formats:</b>  <pre>merge m:1 parent_key using "Identifiers/New_identifiers_part1_parent_key" order new_identifier, first *check if _merge 1 observations are households that have been deleted preserve keep if _merge==1 keep new_identifier parent_key _merge rename _merge _merge_master merge m:1 parent_key using "Identifiers/Deleted_households" /*now go manually through list and identify all anomalies*/ restore drop if _merge==1 list if new_identifier==. /*this list should be empty*/ list if parent_key=="" /*this list should be empty*/ save Long_format_22.10/FATE_Rwanda_Part_1_221018_final- consent_given-hired_labour-repeat_cc5_idcheck, replace</pre> </li> <li>• <b>Part 2, wide and long formats:</b>  <pre>merge m:1 key using "Identifiers/New_identifiers_part2_key" order new_identifier, first list if new_identifier==. /*this list should be empty*/ list if _merge!=3 /*this list should be empty*/ save Wide_format_22.10/FATE_Rwanda_Part_2_221018_final_WID E_idcheck, replace</pre> </li> </ul>



## Consistency Checks

- We undertake some key consistency checks at this stage. However, we don't check for all possible inconsistencies. First, many rules were already inserted in the ODK file in order to avoid logical problems from the start and we don't have to re-check for most of them (some were also inserted in order to facilitate the question presentation but are not relevant for analysis and can be ignored). Second, it is up to the analyst to detect and deal with outliers. We therefore went through all the survey questions and prioritised the following consistency checks that you should undertake now. The rest is up to the analyst.
- Mark all the inconsistencies you find but don't change any data. It is up to the analyst to effect those changes and to justify and document them.

What?	How? (STATA syntax)
<ul style="list-style-type: none"> <li>• Undertake one check after the other</li> </ul>	<ul style="list-style-type: none"> <li>• Within a do-file, each check is a separate activity. The checks should be run, one at a time, where the list is checked and the problems identified and documented before the next check is run. DO NOT RUN the complete do-file at once! You will get garbage and you will not be able to figure out which listings to work on before working on others.</li> </ul>
<ul style="list-style-type: none"> <li>• <b>Skipping:</b> Certain variables should only have values if the answer to a previous filter question appropriate. However, as the skipping functions have been inserted in the original xls files, we will only double check module filter questions in part 1 as they would have the largest impact (see priority skipping checks below). Make sure that the original xls files contains all the necessary skipping functions, if that is not the case or if you want to add skipping checks, here is an example:</li> </ul>	<ul style="list-style-type: none"> <li>• Let's assume that if the filter question has a value of 1, there should be data in the subordinate questions and that if the filter question does not have a value of 1, there should be no data in the subordinate questions. We can then check if the skipping worked correctly with the following command (which can be adapted to match other skipping cases):  <pre>foreach var of varlist consent {   gen s`var'=.   replace s`var'=1 if `var'==1   replace s`var'=0 if s`var'!=1   list `var' if s`var'==0 }</pre></li> </ul>

### Consistency Checks for Part 1

What?	How? (STATA syntax)
<b>Checks in wide</b>	
<ul style="list-style-type: none"> <li>• We repeat several of the commands (e.g. forval) below on string variables with almost the same name (e.g. b1_name_1 to b1_name_11). We include as many numbers of the variable as there are string variables, e.g. if b1_name_12 is a numeric variable because there are only missing values (i.e. no household had 12 members), we cannot include it in the command as it would lead to a type mismatch.</li> <li>• The checks in this section are only possible in the wide format of part 1, however, <b>you should also add all the checks of the long formats of part 1 (see below) to the do-file for the wide format</b> so that you can check independently. The syntax will have to be adapted accordingly.</li> </ul>	
<ul style="list-style-type: none"> <li>• Priority skipping check module A: Check that there</li> </ul>	<ul style="list-style-type: none"> <li>• <code>list new_identifier if hh_people_nb==0</code></li> </ul>

are household members	
• Priority skipping check module B	• <i>list new_identifier if wage_employment!=1 &amp; wage_employ!=</i> "
• Priority skipping check module C	• <i>list new_identifier if cd2_labour_exchange!=1 &amp; howmany_exchange!=</i> "
• Priority skipping check module Da	• <i>list new_identifier if cd1_agr_production!=1 &amp; cd1_how_many!=</i> "
• Priority skipping check module Db	• <i>list new_identifier if cd1_agr_production!=1 &amp; cb1_crop!=</i> .
• No priority skipping check for module Dc necessary	
• No priority skipping check for module De necessary	
• Priority skipping check module E	• <i>list new_identifier if cc5!=1 &amp; activities_hired_work!=</i> "
• Priority skipping check module F	• <i>list new_identifier if da1_ownland!=1 &amp; daa_plots!=</i> . • <i>list new_identifier if da1_ownland!=1 &amp; da9_landcoop!=</i> . • <i>list new_identifier if da4_rent_to!=1 &amp; da4_rent_plot!=</i> . • <i>list new_identifier if da5_otherland!=1 &amp; da5_owner_other!=</i> .
• No priority skipping check for module G necessary	
• No priority skipping check for module Z necessary	
• Check household type with gender of head	• <i>forval i=1/45 {</i> • <i>list new_identifier if (b02_hh_member_gender_`i`==2 &amp; b03_hhrelationship_`i`==1 &amp; a7_hhtype==3)  </i> • <i>(b02_hh_member_gender_`i`==1 &amp; b03_hhrelationship_`i`==1 &amp; a7_hhtype==2)</i> • <i>}</i>
• Check that each paid working household member is older than 14 years old	• <i>forval i=1/11 {</i> • <i>forval j=1/4 {</i> • <i>list b1_name_`i` if b1_name_`i`==name_display_`j` &amp; name_display_`j`!=</i> " & b04_hh_member_age_`i`<14 • <i>}</i> • <i>}</i>
• Check if any numeric variables are negative • If any other variables other than GPS location have negative values, mark them	• <i>preserve</i> • <i>drop gpslatitude</i> • <i>ds, has(type numeric)</i> • <i>foreach var of varlist `r(varlist)` {</i> • <i>list `var` if `var`&lt;0</i> • <i>}</i> • <i>restore</i>
• Check that each unpaid working household member is older than 14 years old	• <i>forval i=1/11 {</i> • <i>forval j=1/2 {</i> • <i>list b1_name_`i` if b1_name_`i`==name_display2_`j` &amp; name_display2_`j`!=</i> " & b04_hh_member_age_`i`<14 • <i>}</i> • <i>}</i>
• List all the production units other than "kg": first check that everything that should be "kg" is	• <i>ds, has(type string)</i> • <i>foreach var of varlist `r(varlist)` {</i> • <i>list `var` if strmatch("`var`","*prod_unit*")==1 &amp; `var`!=</i> " & `var`!="kg" • <i>}</i>

<p>spelled correctly, you then might have to convert the rest using a locally appropriate conversion factor or deal with them in some other way in the analysis.</p>	
<ul style="list-style-type: none"> <li>Show all the area units that are not pre-coded. You might have to convert them using a locally appropriate conversion factor or deal with them in some other way in the analysis.</li> </ul>	<ul style="list-style-type: none"> <li><pre>ds, has(type string) foreach var of varlist `r(varlist)' { list `var' if strmatch("`var'", "*unit_other*")==1 &amp; `var'!= "" }</pre></li> </ul>
<ul style="list-style-type: none"> <li>Check that if others were selected, that they are the same activities for across module E</li> </ul>	<ul style="list-style-type: none"> <li><pre>preserve forval i=1/3 { forval j=1/5 { tostring cc2_2_activity_hired `i' `j', replace replace cc2_2_activity_hired `i' `j'="" if cc2_2_activity_hired `i' `j'="" gen test_`i' `j'=1 if activities_hired_work `i'==activities_most_work `i' &amp; activities_hired_work `i'==cc2_2_activity_hired `i' `j' list activities_hired_work `i' activities_most_work `i' cc2_2_activity_hired `i' `j' if test_`i' `j'=1 } } restore</pre></li> </ul>
<ul style="list-style-type: none"> <li>If there are child-headed households verify the age</li> </ul>	<pre>forval i=1/45 { list b1_name_`i' if b03_hhrelationship_`i'=1 &amp; a7_hhtype==4 &amp; b04_hh_member_age_`i'&gt;18   b03_hhrelationship_`i'=1 &amp; a7_hhtype==4 &amp; b04_hh_member_age_`i'&lt;10   b03_hhrelationship_`i'=1 &amp; a7_hhtype==5 &amp; b04_hh_member_age_`i'&gt;18   b03_hhrelationship_`i'=1 &amp; a7_hhtype==5 &amp; b04_hh_member_age_`i'&lt;10 }</pre>
<ul style="list-style-type: none"> <li>Check that the household size is equal to the maximum of household members</li> </ul>	<ul style="list-style-type: none"> <li><pre>preserve forval i=1/45 { tostring b1_name_`i', replace gen count_`i'=1 if b1_name_`i'!= "" &amp; b1_name_`i'!="." } egen total=rowtotal(count_1-count_45) list new_identifier if hh_people_nb!=total &amp; hh_people_nb!=. restore</pre></li> </ul>
<ul style="list-style-type: none"> <li>Check gender of primary respondent with gender in household list</li> </ul>	<ul style="list-style-type: none"> <li><pre>forval i=1/11 { list new_identifier if a2_nameprim==b1_name_`i' &amp; a2_nameprim!=" &amp; a8_gender!=b02_hh_member_gender_`i' }</pre></li> </ul>
<ul style="list-style-type: none"> <li>Check that primary</li> </ul>	<ul style="list-style-type: none"> <li><pre>forval i=1/11 {</pre></li> </ul>

respondent is over 18 years old	<i>list new_identifier if a2_nameprim==b1_name_`i` &amp; a2_nameprim!="" &amp; b04_hh_member_age_`i`&lt;18</i> }
<ul style="list-style-type: none"> <li>Check that the primary respondent is a household member (mistakes could be due to spelling error!)</li> </ul>	<ul style="list-style-type: none"> <li><i>list new_identifier if a2_nameprim!=b1_name_1 &amp; a2_nameprim!=b1_name_2 &amp; a2_nameprim!=b1_name_3 &amp; a2_nameprim!=b1_name_4 &amp; a2_nameprim!=b1_name_5 &amp; a2_nameprim!=b1_name_6 &amp; a2_nameprim!=b1_name_7 &amp; a2_nameprim!=b1_name_8 &amp; a2_nameprim!=b1_name_9 &amp; a2_nameprim!=b1_name_10 &amp; a2_nameprim!=b1_name_11</i></li> </ul>
<b>Checks in long_final</b>	
<ul style="list-style-type: none"> <li>Households should only have been replaced if nobody was available in the household</li> </ul>	<ul style="list-style-type: none"> <li><i>list new_identifier if avail!=1 &amp; replacement_id==.</i> <i>list new_identifier if avail==1 &amp; replacement_id!=.</i></li> </ul>
<ul style="list-style-type: none"> <li>Check that only one id was given per household</li> </ul>	<ul style="list-style-type: none"> <li><i>list new_identifier if code_id!=. &amp; replacement_id!=.</i></li> <li>Alternative ways of checking that could be employed throughout: <i>assert code_id!=. if replacement_id!=.</i></li> <li>Or: <i>count if code_id!=. &amp; replacement_id!=.</i></li> </ul>
<ul style="list-style-type: none"> <li>Check if consent was given</li> <li>If the consent is missing but questions were answered, keep the household. If the consent was no or if almost no questions have been answered, delete the household</li> </ul>	<ul style="list-style-type: none"> <li><i>list consent if consent!=1</i></li> </ul>
<ul style="list-style-type: none"> <li>Check that there are no duplicates of primary respondents</li> <li>If there are duplicates you might check if they are from different households after unique identifiers have been developed (see below)</li> </ul>	<ul style="list-style-type: none"> <li><i>duplicates list a2_nameprim if a2_nameprim != ""</i></li> </ul>
<ul style="list-style-type: none"> <li>Check gender of primary respondent with status</li> </ul>	<ul style="list-style-type: none"> <li><i>list a2_nameprim if (a8_gender==2 &amp; status==5)   (a8_gender==1 &amp; status==4)</i></li> </ul>
<ul style="list-style-type: none"> <li>Check household type with primary respondent</li> </ul>	<ul style="list-style-type: none"> <li><i>list a2_nameprim if (a8_gender==2 &amp; status==1 &amp; a7_hhtype==3)   (a8_gender==2 &amp; status==1 &amp; a7_hhtype==4)   (a8_gender==1 &amp;</i></li> </ul>

	<i>status==1 &amp; a7_hhtype==2)   (a8_gender==1 &amp; status==1 &amp; a7_hhtype==4)</i>
<ul style="list-style-type: none"> <li>Check that all households have at least 1 member</li> </ul>	<ul style="list-style-type: none"> <li><i>list new_identifier if hh_people_nb&lt;=0</i></li> </ul>
<ul style="list-style-type: none"> <li>Check that years are formatted correctly</li> </ul>	<ul style="list-style-type: none"> <li><i>foreach var of varlist da2g_surfaceprod da5_aquires1{ list `var' if (`var'!=. &amp; `var'&lt;1900)   (`var'!=. &amp; `var'&gt;2018) }</i></li> </ul>
<ul style="list-style-type: none"> <li>Check that percentages are correct</li> </ul>	<ul style="list-style-type: none"> <li><i>list df9_interest if df9_interest!=. &amp; df9_interest&gt;100 &amp; df9_interest&lt;0</i></li> </ul>
<b>Checks in long_final-consent_given-repeat_hh</b>	
<ul style="list-style-type: none"> <li>Check for duplicate household members</li> </ul>	<ul style="list-style-type: none"> <li><i>duplicates list b1_name if b1_name!=</i>"</li> </ul>
<ul style="list-style-type: none"> <li>Check age difference between oldest child and head</li> </ul>	<ul style="list-style-type: none"> <li><i>by new_identifier, sort: egen head_age=min(b04_hh_member_age) if b03_hhrelationship==1 by new_identifier, sort: egen oldest_child=max(b04_hh_member_age) if b03_hhrelationship==3 by new_identifier, sort: egen min_head_age=min(head_age) by new_identifier, sort: egen max_oldest_child=max(oldest_child) by new_identifier, sort: egen difference=min_head_age-max_oldest_child list difference b1_name new_identifier if difference &lt; 15</i></li> </ul>
<ul style="list-style-type: none"> <li>Check age of head</li> </ul>	<ul style="list-style-type: none"> <li><i>list b1_name if b04_hh_member_age&lt;18 &amp; b03_hhrelationship==1</i></li> </ul>
<ul style="list-style-type: none"> <li>Check age of spouse</li> </ul>	<ul style="list-style-type: none"> <li><i>list b1_name if b04_hh_member_age&lt;18 &amp; b03_hhrelationship==2</i></li> </ul>
<ul style="list-style-type: none"> <li>Check age and marital status</li> </ul>	<ul style="list-style-type: none"> <li><i>list b1_name if b6_marstat!=1 &amp; b6_marstat!=.a &amp; b6_marstat!=.b &amp; b6_marstat!=. &amp; b04_hh_member_age&lt;18</i></li> </ul>
<ul style="list-style-type: none"> <li>Check age and education status</li> </ul>	<ul style="list-style-type: none"> <li><i>list b1_name if b04_hh_member_age&lt;10 &amp; b8_education&gt;16 &amp; b8_education!=96 &amp; b8_education!=.</i></li> </ul>
<b>Checks in long_final-consent_given-repeat_wage_employment</b>	
<ul style="list-style-type: none"> <li>No consistency checks specific to this dataset</li> </ul>	
<b>Checks in long_final-repeat_multiplicity</b>	
<ul style="list-style-type: none"> <li>Check that there are no more than 31 working days per month</li> </ul>	<ul style="list-style-type: none"> <li><i>list bb7_work_day_permonth if bb7_work_day_permonth&gt;31 &amp; bb7_work_day_permonth!=.</i></li> </ul>
<ul style="list-style-type: none"> <li>Check that there are no more than 24 working hours per day</li> </ul>	<ul style="list-style-type: none"> <li><i>list bb8_work_hours_perday if bb8_work_hours_perday&gt;24 &amp; bb8_work_hours_perday!=.</i></li> </ul>
<b>Checks in long_final-consent_given-repeat_cd</b>	
<ul style="list-style-type: none"> <li>Check that there are no more than 365 working days per year</li> </ul>	<ul style="list-style-type: none"> <li><i>list cd6_work_days if cd6_work_days&gt;365 &amp; cd6_work_days!=.</i></li> </ul>
<ul style="list-style-type: none"> <li>Check that there are no more than 24 working hours per day</li> </ul>	<ul style="list-style-type: none"> <li><i>list cd7_hours if cd7_hours&gt;24 &amp; cd7_hours!=.</i></li> </ul>
<b>Checks in long_final-consent_given-repeat_agrprod</b>	

<ul style="list-style-type: none"> <li>No consistency checks specific to this dataset</li> </ul>	
<b>Checks in long_final-consent_given-hired_labour-repeat_cc5</b>	
<ul style="list-style-type: none"> <li>Check that there are no more than 365 working days per year</li> </ul>	<ul style="list-style-type: none"> <li><code>foreach var of varlist hire_onetime m1_b m1_d m1_f m1_h m1_j m1_l { list `var' if `var'&gt;365 &amp; `var'!=. }</code></li> </ul>
<b>Checks in long_final-consent_given-repeat_cc2</b>	
<ul style="list-style-type: none"> <li>Check that there are no more than 365 working days per year</li> </ul>	<ul style="list-style-type: none"> <li><code>foreach var of varlist cc4_work_day cc14_number { list `var' if `var'&gt;365 &amp; `var'!=. }</code></li> </ul>

## Consistency Checks for Part 2

What?	How? (STATA syntax)
<b>Checks in wide &amp; long_final</b>	
<ul style="list-style-type: none"> <li>Check if any numeric variables are negative</li> <li>If any other variables have negative values, mark them</li> </ul>	<ul style="list-style-type: none"> <li><code>ds, has(type numeric) local varlist `r(varlist)' local toexclude r1_1 r1_2 r1_3 r1_4 r1_5 r1_6 gpslatitude local varlist: list varlist - toexclude foreach var of local varlist { list `var' if `var'&lt;0 }</code></li> </ul>
<ul style="list-style-type: none"> <li>Check if the person was available</li> </ul>	<ul style="list-style-type: none"> <li><code>list avail if avail!=1</code></li> </ul>
<ul style="list-style-type: none"> <li>Check if consent was given</li> <li>If the consent is missing but questions were answered, keep the household. If the consent was no or if almost no questions have been answered, delete the household</li> </ul>	<ul style="list-style-type: none"> <li><code>list consent if consent!=1</code></li> </ul>
<ul style="list-style-type: none"> <li>Check that there are no duplicates of respondents</li> <li>If there are duplicates you might check if they are from different households after unique identifiers have been developed (see below)</li> </ul>	<ul style="list-style-type: none"> <li><code>duplicates list a3_primary if a3_primary!=""</code></li> </ul>
<ul style="list-style-type: none"> <li>Check household type with respondent</li> </ul>	<ul style="list-style-type: none"> <li><code>list a3_primary if (a3b_gender==2 &amp; a3c_relationship==1 &amp; a4_hhtype==3)   (a3b_gender==2 &amp; a3c_relationship==1 &amp; a4_hhtype==4)   (a3b_gender==1 &amp; a3c_relationship==1 &amp; a4_hhtype==2)   (a3b_gender==1 &amp; a3c_relationship==1 &amp; a4_hhtype==5)</code></li> </ul>

<ul style="list-style-type: none"> <li>If there are child-headed households verify the age</li> </ul>	<ul style="list-style-type: none"> <li><i>list a3_primary if a4_hhtype==4 &amp; a3b_age&gt;18   a4_hhtype==4 &amp; a3b_age&lt;10   a4_hhtype==5 &amp; a3b_age&gt;18   a4_hhtype==5 &amp; a3b_age&lt;10</i></li> </ul>
<ul style="list-style-type: none"> <li>Check age of head</li> </ul>	<ul style="list-style-type: none"> <li><i>list a3_primary if a3b_age&lt;18 &amp; a3c_relationship==1</i></li> </ul>
<ul style="list-style-type: none"> <li>Check age of spouse</li> </ul>	<ul style="list-style-type: none"> <li><i>list a3_primary if a3b_age&lt;18 &amp; a3c_relationship==2</i></li> </ul>
<ul style="list-style-type: none"> <li>Check if any household head or spouse was replaced</li> </ul>	<ul style="list-style-type: none"> <li><i>list a3_primary a3c_relationship a3b_age if a3c_relationship!=1 &amp; a3c_relationship!=2 &amp; a3c_relationship!=.</i></li> </ul>
<ul style="list-style-type: none"> <li>Check crop and coffee questions</li> </ul>	<ul style="list-style-type: none"> <li><i>foreach var of varlist i03_inputplot i04_tecno i05_fert i06_labor i07_harvest i08_selling i08_seed i08_income { list `var' if i01_crop1_1!=1 &amp; i01_crop1_2!=1 &amp; i01_crop1_3!=1 &amp; i01_crop1_4!=1 &amp; i01_crop1_5!=1 &amp; `var'!=. }</i></li> </ul>
<ul style="list-style-type: none"> <li>Check finances and household type</li> </ul>	<ul style="list-style-type: none"> <li><i>list new_identifier if a4_hhtype!=1 &amp; i11_finance1!=.</i></li> </ul>
<ul style="list-style-type: none"> <li>Check ownership and household type</li> </ul>	<ul style="list-style-type: none"> <li><i>foreach var of varlist i1_houseownership i2_hhgoods i3_productiveassets i4_livestock i6_whobuys_agr i8_whobuysmajor i9_whobuyspurchases { list `var' if (a4_hhtype!=1 &amp; `var'==3)   (a4_hhtype!=1 &amp; `var'==5) }</i></li> </ul>
<ul style="list-style-type: none"> <li>There is a large number of possible time-use inconsistencies depending on the context and they cannot be checked here. Random checks have been effected in the field which must suffice. The rest is up to the analyst.</li> </ul>	
<ul style="list-style-type: none"> <li>Check that there are no more than 365 days per year</li> </ul>	<ul style="list-style-type: none"> <li><i>list excessive_hours if excessive_hours&gt;365 &amp; excessive_hours!=.</i></li> </ul>

The data verification will continue into the actual analysis. Some problems cannot be identified until analysis has begun.

## Dealing with Missing Data

- Don't delete missing data, however, you cannot simply ignore missing values in your dataset either.
- In general, do not assume missing observations to mean 0 (for example, most "how many" questions only allowed an integer as answer so no answer might mean 0 but it is also possible that the person didn't know the answer but there was no choice for DK) unless this is heavily suggested by answers to previous questions or other variable within the household.
- Analysts should report how many observations are missing and inform the reader how missing data was handled.

What?	How? (STATA syntax)
<ul style="list-style-type: none"> <li>Do not drop observations that have missing values or impute the missing values based on other observations! You should always tell your algorithm that a value was missing because missingness is informative.</li> </ul>	<ul style="list-style-type: none"> <li>Stata automatically assigns a "." for missingness to each missing numeric value upon import and a "" (blank) for missingness for each string value.</li> <li>In addition to the default ".", which is called the "system missing value" or sysmiss, Stata has 26 other numeric missing values: ".a", ".b", ".c", ..., ".z", which are called the "extended missing</li> </ul>

	values". This helps us to assign different reasons for missingness, such as NA (.a) or DK (.b), see below.
<ul style="list-style-type: none"> <li>Skipped questions: If certain questions did not appear in the question path, mark them as missing values.</li> </ul>	<ul style="list-style-type: none"> <li>Stata automatically assigns a "." for answers in relations to questions that have been skipped.</li> </ul>
<ul style="list-style-type: none"> <li>We do not distinguish here if a missing answer has been skipped or was left out despite the question having appeared. However, if needed, the analyst can always go back and test for each case if a variable was skipped or is missing despite the question having appeared. This might be the case if there is a very large number of missing values. The analyst can use descriptive statistics to see how many are missing and should look for meaning in non-random missing values. Maybe the respondents are indicating something important by not answering one of the questions or maybe an enumerator consistently left out questions despite them having appeared. There is a variety of statistical methods available for handling missing data. It is up to the analyst to use, document and justify the most appropriate one for his/her case.</li> </ul>	

## Dealing with Outliers

- It always helps to use descriptive statistics to get to know your data better and to check for problems and potential outliers. Here are some useful commands in that regard: codebook, describe, summarize, list, tabstat, tabulate (helpful explanations are provided in IFPRI, 2018). Plotting numeric variables might also be particularly helpful in identifying outliers.
- However, don't change any answers! Leave all outliers in the dataset. It is up to each researcher to deal with outliers in their variables. This is because different analysis goals require different outlier treatments and it is up to each analyst to choose the method most appropriate for his/her case.
- Exception: if you are very sure that an entry is wrong (e.g. because it is logically incoherent as for example shown by a consistency check), you can change it by the correct value or make it a missing variable, however, you MUST document this in the change log including the original value so that it is always possible to go back if necessary.

## Labelling

In order to help you with the labelling process, we first prepare a new and separate dataset (i.e. not in your do-file) based on the choices sheet of the xls forms used in the survey. Proceed as follows:

- Copy the original xls form of part 1 used in the survey. Rename the copy to end in \_replaced. Only work in the copy from now on. Open it with Excel.
- In the choices sheet, select the entire name column (column B), go to Home/Find & Select/Replace... under "Find what:" type "NA", under "Replace with:" type ".a". Click "Options >>" and make sure to tick the box that says "Match case" and "Match entire cell contents" and then click "Replace All".
- In the choices sheet, select the entire name column (column B), go to Home/Find & Select/Replace... under "Find what:" type "DK", under "Replace with:" type ".b". Click "Options >>" and make sure to tick the box that says "Match case" and "Match entire cell contents" and then click "Replace All".
- Save the xls file.
- Open a new dataset in your software package. Now copy the entire 4 first columns of the choices sheet of the xls form \_replaced of part 1 used in the survey. Special paste them (selecting "Tab" as delimiter) into the data-editor (in the edit mode) of your software package. Save that dataset under the name "Value\_labels\_part1".
- Repeat the same procedure for part 2.



Perform all steps in the wide formats of part 1 & 2 and then copy the relevant labels into all corresponding long datasets of part 1 & 2.

### How? (STATA syntax)

- Run the following commands in the “Value\_labels\_part1” dataset in order to create an answer list with the appropriate syntax to copy and paste the value labels rather than having to type them each time:

```
drop var3
replace var4= `"" + var4 + `""
forvalues i=10/100 {
  replace var4=regexpr(var4,"`i.'",")
}
forvalues i=1/9 {
  replace var4=regexpr(var4,"`i.'",")
}
foreach var of varlist var4 {
  replace `var'=strtrim(`var')
}
egen value_label = concat(var2 var4), p(" ")
save Z:\Data_Cleaning\Stata\Labels\Value_labels_part1, replace
```

- Export the dataset under the name “Value\_labels\_part1” to excel (export data to excel spreadsheet) selecting “Save variable names to first row in Excel file” and open the excel file. You can copy the value label list in the “value\_label” column into the do-file as needed.
- Repeat the same procedure for part 2.

- The #delimiter ; command is useful in do-files if you want to paste long lists (e.g. of variable names or value labels) from excel and you only have them in column. Stata will read everything before a semi-colon as one line and you don’t have to type all the names in one row.
- However, once we change the line delimiter to semicolon, all lines, even short ones, must end in semicolons. Stata treats carriage returns as no different from blanks. We can change the delimiter back to carriage return by typing #delimiter cr.
- Attach value labels to your variables in the following way, using var1 (list name) of “Value\_labels\_part1/2” as the name for the value labels :

```
#delimiter;
label define gender
  1 "Male"
  2 "Female";
foreach var of varlist a8_gender a14_gender enumerator {;
  label val `var' gender;
};
```

- The loop attaches the same label to multiple variables and can be expanded as needed.
- Use the survey question in the “label” column of the “survey” sheet of the original xls file to label the variables and label variable. You can also use a loop for multiple-choice dummies or variables in repeat sections of the wide format.
- SurveyCTO also created a Stata\_do\_template. While we cannot use that directly for our data cleaning here, you can copy and paste parts of the labelling section so that you don’t have to copy each survey question from the original xls file.
- You might have to include syntax so that only numeric variables are labelled (otherwise, if string variables are included in the loop, e.g. because another variable with \_other attached is included, you will encounter an error message).

```
ds, has(type string)
local strings `r(varlist)'
#delimiter;
```

<pre> label define read 1 "Cannot read and cannot write" 2 "Can write only" 3 "Can read only" 4 "Can read and write " 96 "Other, specify" .a "NA" .b "DK"; #delimit cr foreach rgvar of varlist b9_read_* {     label variable `rgvar' "Can \\${b1_name} read and write?" }  unab want: b9_read_* local numerics : list want - strings label values `numerics' read </pre>	
<ul style="list-style-type: none"> <li>Multiple-choice variables are read as strings by Stata, however, string variables cannot be labelled in Stata. Fortunately, SurveyCTO automatically split all the select_multiple answers into additional dummy variables (both in the long and wide csv formats). Dummy variables have been created for each choice in the choice list, regardless if that choice was ever chosen or not, and whether that question was ever asked or not. You don't have to do anything else with multiple-choice variables other than labelling them. It is up to the analyst to decide how to deal with multiple-choice questions.</li> <li>Since multiple-choice variables have been split, you have to define a value label for each split variable and attach that value label to the correct split variable. Make sure that the value label always has the number 1 in front of it as a dummy variable by definition is only 1 or 0 (or missing). You might have to include syntax so that only numeric variables are labelled (otherwise, if string variables are included in the loop, e.g. the original multiple-choice variable, you will encounter an error message).</li> </ul> <pre> foreach rgvar of varlist cd4_return_* {     label variable `rgvar' "Does \\${name_display2} receive anything in return from this employer?" } label define return_1           1 "No return" label define return_2           1 "Cash" label define return_3           1 "Goods" label define return_4           1 "Labour" label define return_5           1 "Food or drink" label define return_96          1 "Other, specify"     forval i=1/5 {         foreach var of varlist cd4_return_`i'_* {             capture confirm numeric variable `var'             if! _rc {                 label values `var' return_`i'             }         }     }     foreach var of varlist cd4_return_96_* {         capture confirm numeric variable `var'         if! _rc {             label values `var' return_96         }     } } </pre>	<ul style="list-style-type: none"> <li>For example: <i>save Wide_format_22.10/FATE_Rwanda_Part</i></li> </ul>
<ul style="list-style-type: none"> <li>Save everything you have done up to now under the name <i>clean_final</i></li> </ul>	

	_1_221018_final_WIDE_clean_final, replace
--	---

## Anonymising Data

- Steps to be performed in each dataset indicated below.
- Once you have finished the data cleaning process, a unique identifier has been created for each household and you have made sure that it has been correctly added to each dataset, you can delete all the individual names and other identifying information from the dataset. Be aware that they are spread all over the datasets due to repeat functions. Proceed in the following way in order to find them:

What?	How? (STATA syntax)
<ul style="list-style-type: none"> <li>• Randomise the village and enumerator codes</li> </ul>	<ul style="list-style-type: none"> <li>• <b>Part 1 &amp; 2, wide and long_final:</b> Run do-files provided by the project</li> </ul>
<ul style="list-style-type: none"> <li>• Make sure to use the correct file</li> </ul>	<ul style="list-style-type: none"> <li>• <b>Part 1 &amp; 2, wide and long_final:</b>  <pre>use Wide_format_22.10/ FATE_Rwanda_Part_1_221018_final_WIDE_clean_final_village_enu, clear</pre> </li> </ul>
<ul style="list-style-type: none"> <li>• Find all the variables containing individual names by identifying all the variables with the word "name" in variable name (those are the ones containing people's names)</li> </ul>	<ul style="list-style-type: none"> <li>• <b>Part 1 &amp; 2, all wide and long formats:</b> <ul style="list-style-type: none"> <li>• <code>lookfor name</code></li> <li>• Or:  <pre>ds, has(type string) foreach var of varlist `r(varlist)' { list `var' if strmatch("`var'", "*name*")==1 &amp; `var'!=""} </pre> </li> </ul> </li> </ul>
<ul style="list-style-type: none"> <li>• Now drop all of those variables</li> </ul>	<ul style="list-style-type: none"> <li>• <b>Part 1 &amp; 2, all wide and long formats:</b>  <pre>drop *name*</pre> </li> </ul>
<ul style="list-style-type: none"> <li>• Drop other variables containing identifying information and drop old code_ids as they add confusion (adjust relevant variable names according to dataset). Also drop variable <code>avail</code> as unavailable households have been deleted and the variable does not have meaning anymore.</li> </ul>	<ul style="list-style-type: none"> <li>• <b>Part 1, wide and long_final:</b>  <pre>drop avail submissiondate code_id today_date starttime end_time *gps* instanceid key setofrepeat*</pre> </li> <li>• <b>Part 2, wide and long_final:</b>  <pre>drop avail submissiondate code_id a12_starttime a13_datesecondvisit a3_primary end_time *gps* instanceid key</pre> </li> </ul>
<ul style="list-style-type: none"> <li>• Drop automatically created variables with label <code>"reserved_name_for_field_list_labels"</code></li> </ul>	<ul style="list-style-type: none"> <li>• <b>Part 2, wide and long_final:</b>  <pre>lookfor reserved_name_for_field_list_labels foreach var of varlist `r(varlist)' { drop `var' }</pre> </li> </ul>
<ul style="list-style-type: none"> <li>• List all the string variables in order to check that no sensitive information remains (especially in the "other" answers), else you might have to anonymise manually.</li> </ul>	<ul style="list-style-type: none"> <li>• <b>Part 1 &amp; 2, all wide and long formats:</b>  <pre>ds, has(type string) foreach var of varlist `r(varlist)' { count if strmatch(`var', "* *")==1 &amp; regexm(`var', "[a-zA-Z]")==0   strmatch(`var', "* *")==1 &amp; regexm(`var', "\.a")==1   strmatch(`var', "* *")==1 &amp; regexm(`var', "\.b")==1 if r(N)==0 { list `var' if `var'!=" &amp; `var'!="kg"</pre> </li> </ul>

	<pre>} }</pre>
<ul style="list-style-type: none"> <li>List all the value labels to check that they contain no sensitive information (such as the village name which has already been replaced above), else you might have to anonymise manually.</li> </ul>	<ul style="list-style-type: none"> <li><b>Part 1 &amp; 2, all wide and long formats:</b> <i>label list</i></li> </ul>
<ul style="list-style-type: none"> <li>Sort the table by the new identifiers and make sure that there are no duplicates in part 1 and no more than 1 duplicate in part 2</li> </ul>	<ul style="list-style-type: none"> <li><b>Part 1 &amp; 2, all wide and long_finals:</b> <i>sort new_identifier</i> <i>duplicates report new_identifier</i></li> </ul>
<ul style="list-style-type: none"> <li>Save everything you have done up to now under the name <i>clean_final_anonymised</i></li> </ul>	<ul style="list-style-type: none"> <li>For example: <i>save</i> <i>Wide_format_22.10/FATE_Rwanda_Part_1_221018_final_WI</i> <i>DE_clean_final_anonymised, replace</i></li> </ul>

## Backup

- Make sure to backup your cleaned datasets properly.
- Related information files (do-files for data cleaning, coding books, etc.) should be included together with the backups. Copies in cloud systems without adequate security (i.e. dropbox, google drive) should be avoided.
- Make sure you still have the encrypted raw-data files and discuss with your supervisor if they should be deleted after the analysis is finished.