
How to Use Datasets of Georgia RTM¹/MICS Plus 2020-2021 survey

MICS Plus is a longitudinal household survey that collects information from a representative sample of households through interviews on direct phone calls. The same households are interviewed multiple times at regular intervals over a period of one year. Depending on the frequency of calls, up to 12 waves of calls may be made to each sampled household.

Details of the Georgia RTM/MICS Plus 2020-2021 survey methodology will be provided in a separate document. This note summarizes some of the methodological features with respect to the use of micro datasets and accompanies the publicly shared datasets.

Anonymized datasets of the Georgia RTM/MICS Plus 2020-2021 survey are being periodically released for public use after every three waves.

1

Survey Implementation

Households interviewed in Georgia MICS 2018 survey were used as the sample frame of the Georgia RTM/MICS Plus 2020-2021 survey. The response rate for Georgia MICS 2018 was 94 percent. Interviews were completed in 12,270 households, and 95 percent of these households had provided a phone number. Phone numbers were collected from the respondent to the Household Questionnaire, and only one phone number was collected per household.

The sample size for the Georgia RTM/MICS Plus 2020-2021 was determined as 2,118 households at the outset, with the goal that around 2,000 households would be captured with completed interviews. 3 households were selected from each of the 706 clusters of Georgia MICS 2018.

The first wave of calls were completed in 2 stages, as described below:

Stage 1

A total of 2,118 households were selected (3 households from each cluster) from the households interviewed in the Georgia MICS 2018 regardless of the availability of phone number. At the end of this stage, phone number was not available or valid for 6.8 percent of the 2,118 sampled households. Calls were placed to households with available phone number.

Stage 2

At the end of stage 1, there were 734 households with incomplete interviews, including those without phone number. 731 of these households were substituted with other households regardless of availability

¹ Real Time Monitoring

of phone number from the corresponding clusters in the Georgia MICS 2018 data. The substitution was performed by applying a model-based, conditional substitution method that employs Euclidian Distance analysis (available under Nearest Neighbor Analysis in SPSS). For 3 households, it was not possible to identify households that were similar enough in terms of their selected background characteristics. These households were substituted with other households randomly selected from the corresponding cluster.

After the completion of Wave 1, the final effective sample size for the subsequent waves was established at 1,996 households (see Section 7 "Number of Cases in Datasets" for details). The above-mentioned substitution process applies only to the first wave of calls and is not repeated in the subsequent waves.

2

Typology of Questions

MICS Plus interviews are administered to one respondent. Defined as a knowledgeable adult household member, the respondent answers all questions regardless of whether the question is about himself/herself, the household, the dwelling, or a specific household member. While it is common to have the same respondent across waves, note that different respondents may be interviewed across different waves.

One questionnaire is administered to the selected households of MICS Plus surveys. Questions are organized in modules. There are 3 main types of modules, resulting in different units of analysis.

- The module "Call Attempts Panel" captures information related to the calls made to the selected households. A household has as many records in the Call Attempts dataset as the number of call attempts made to that household, regardless of whether the call resulted in a completed interview or not.
- The "List of Household Members" captures key characteristics of each household member in the interviewed households.
- Each wave of questionnaires includes several "Sectoral" modules that capture information on various topics of interest. While the Call Attempts Panel and the List of Household Members are included in all waves, the Sectoral modules may be repeated or changed from one wave to another. Examples of Sectoral modules include "COVID-19", "Education", and "Child wellbeing and Health".

Sectoral modules can have different groups of individuals of interest. For example, the "Education" module targets children age 2-17 years at the beginning of the school year. In the case of this module, one child age 2-17 years at the beginning of the school year is randomly selected from the List of Household Members; questions in this module target the selected child only. Random selection and subsequent application of sample weights (see sub-Section 5 "Sample weight variables") ensure that the data collected is representative of the child population age 2-17 years at the beginning of the school year in Georgia.

In general, a question can be about the household or dwelling, household members, the respondent, or a randomly selected individual as described above. The table below shows the possible typology of questions in MICS Plus surveys.

Type of question	Example	Representative of
About household	WS4. In the last 30 days, has there been any time when your household did not have sufficient quantities of drinking water?	Households in Georgia
About household members	CVH26. Have you or any one of the members of this household lost his/her job since December 2020?	Household population of Georgia
About respondent	CV10. Since the last (day of the week), have you been able to keep distance from people when in public places? Would you say: always, very often, sometimes, rarely or never?	Knowledgeable adult household members who were interviewed. Since respondents are not selected randomly, such questions do not yield representative data on Georgia's adult population's characteristics, opinions, and behavioural patterns.
About a randomly selected individual	ED9. Since the last (day of the week), did (name ²) watch any TV lessons?	Population age 2-17 years at the beginning of the school year in Georgia

3

Datasets

Due to the different units of analysis, MICS Plus data have a hierarchical structure. The Georgia RTM/MICS Plus digital data collection system was developed using [CSPPro](#), a software programme that well handles hierarchical datasets. For analysis, CSPPro datasets are exported to the following three SPSS datasets:

HH - Household dataset

Unit of analysis: Households

Includes: Information on household characteristics, such as ownership of consumer goods, energy use, water and sanitation, as well as sectoral modules such as COVID-19, Distance learning/Education, Healthcare services, and Child wellbeing and Health. Variables corresponding to questions on the household and dwelling, as well as those sectoral modules that target the respondent or randomly selected individuals, are included in this dataset.

HL - Household members dataset

Unit of analysis: Household members

Includes: Characteristics of individual members of the households such as date of birth, age, sex and relationship to the household head.

CA - Call attempts dataset

Unit of analysis: Call attempts

Includes: Information on all calls made to the households, regardless of the interview result, such as date and time of call, call outcome, and consent for interview.

² Name of the randomly selected individual. In this example, name of the selected child age 2-17 years at the beginning of the school year.

All datasets are anonymized, e.g., unique identifiers such as location, phone numbers, and names collected during the interviews are removed from all datasets to ensure confidentiality and privacy.

4

Dataset Naming Conventions

MICS Plus datasets in SPSS format are distributed in a compressed WINZIP folder and are uniquely named using the following naming conventions:

- [CCC] 3-digit Standard ISO Country Code
- [DD] Dataset type: HH (Household), HL (Household members), CA (Call attempts)
- [W###] Wave number (W01 – first wave, W02 – second wave...)

For example, SPSS datasets of the first wave are stored in the WINZIP file named "GEOMICSPUSW01.ZIP" and are named as GEOHHW01.SAV, GEOHLW01.SAV, and GEOCAW01.SAV. In addition to the SPSS datasets, the WINZIP files also include a technical note with contact details.

CSPRO datasets are not shared publicly.

5

Variable Naming Conventions, Variable Construction

Contents of datasets correspond to questionnaire contents for each wave.

The most common correspondence is the "one question-to-one variable" form. For example, a variable named "EU1" in the dataset represents the question numbered "EU1" in the questionnaire.

Other specific types of variables are explained below.

Identification (ID) variables

Each dataset has case-identifiers that uniquely identify each unit of analysis (household, household member, or call attempt) and allow merging of datasets when the relationship is logically possible. Details on merging datasets can be found in the Section 9 "Merging Datasets".

ID variable in the Household dataset - HHID, which is created by combining the Wave Number (HH0), Cluster number (HH1), Household Number (HH2), and Wave 1 Stage Number³ (HH0A), as shown below:

³ Only the Wave 1 household dataset has up to 2 stages due to the specific implementation approach of this wave, as explained in the Section 1 "Survey Implementation". For other waves, the variable "HH0A" (Wave 1 Stage Number) is set to "0 – Not applicable".

HHID	HH0	HH1	HH2	HH0A
1001071	1	1	7	1
1095021	1	95	2	1
1266192	1	266	19	2

ID variables in the Household Members dataset – HHID and HHMEMID, which are combinations of Wave Number (HH0), Stratum (HH1), Household Number (HH2), Wave 1 Stage Number³ (HH0A), and Line Number of Household Member (HL1)

HHID	HHMEMID	HH0	HH1	HH2	HH0A	HL1
2003020	200302001	2	3	2	0	1
2003020	200302002	2	3	2	0	2
2003020	200302003	2	3	2	0	3

ID variables in the Call Attempts dataset – HHID and CAID, which are combinations of Wave Number (HH0), Stratum (HH1), Household Number (HH2), Wave 1 Stage Number³ (HH0A), and Call Attempt Number (CA1).

HHID	CAID	HH0	HH1	HH2	HH0A	CA1
3175070	317507001	3	175	7	0	1
3175070	317507002	3	175	7	0	2
3175070	317507003	3	175	7	0	3

Variables Specific to Wave 1

Variable "SampleType" (Whether a household was initially selected or a substitute household) has been created to indicate whether a household was initially selected (1) or was a substitute household (2).

When analyzing Wave 1 datasets, one needs to select only the relevant records, depending on whether the analysis is intended to be based on households before or after substitution. For this purpose, variables "aftersub" (For analysis – household after substitution) and "beforesub" (For analysis – household before substitution) have been constructed, with values of "0 – No" and "1 - Yes". Therefore, if analysis will be based on households after the substitution, one must select only households with the value "1" on the variable "aftersub"– these are the households that the published results of the survey are based on. For analysis based on households before substitution, households with the value "1" on the variable "beforesub" should be selected.

There are two variables that have the letter “B” in their names and the text “before substitution” in their labels. These are the variables that need to be used for analyses based on households before substitution. For example, variable “hhweight” is the sample weight variable to be used for analyses after substitution (the final survey sample), while variable “hhweightB” is the sample weight variable to be used for analyses before substitution. The same applies to variables “chschange217weight” and “chschange217weightB” – sample weight variables for randomly selected children age 2-17 years at the beginning of school year (see sub-Section 5 “Sample weight variables”).

The following is an example of SPSS code that selects households for 1) analysis after substitution and 2) analysis before substitution and applies sample weights.

```
* 1) Example of selecting households for analysis after substitution and applying sample weight.

* open the household dataset.
get file = "GEOHHW01.sav".

* select only the households for analysis after substitution.
select if (aftersub = 1).

* select only the interviewed households.
select if (HH17 = 1).

* applying sample weight.
weight by hhweight.

* 2) Example of selecting households for analysis before substitution and applying sample weight.

* open the household dataset.
get file = "GEOHHW01.sav".

* select only the households for analysis before substitution.
select if (beforesub = 1).

* select only the interviewed households.
select if (HH17 = 1).

* applying sample weight.
weight by hhweightB.
```

Multiple Response Questions

Multiple response questions are those questions that allow the coding of multiple answers to a single question. For such questions, response categories are alphabetical. When a multiple response question is asked, the interviewer does not read the response categories; the respondent provides answers that fit one or more response categories. The interviewer then records the most appropriate response code and probes until the respondent has no more responses. For multiple response questions, each response category has a designated string variable in the relevant dataset. These variables are named ending with the letter that corresponds to the response category.

For example, variables "CH5A", "CH5B", ..., "CH5X" represent, respectively, response categories "A", "B", ..., "X" of the question "CH5" below. Possible values for each of these variables in the datasets are the letters corresponding to the response category (e.g., "A" for the variable "CH5A") and a blank space (when that response category is not selected). In addition to these variables, a variable with "NR" in the name (e.g., CH5NR) and with the value "?" is created to code cases when there is no response to such questions.

<p>CH5. What equipment do the members of your household use to access internet?</p> <p><i>Probe: Anything else?</i> <i>Multiple responses are allowed.</i> <i>Do not read out the response categories.</i></p>	DESKTOP COMPUTER	A
	LAPTOP.....	B
	TABLET.....	C
	SMART PHONE.....	D
	SMART TV	E
	OTHER (<i>specify</i>)	X

For a household where the respondent has responded that household members use desktop computers, tablets, and smartphones to access the internet, the following values would apply for the corresponding variables:

CH5A	A
CH5B	[blank]
CH5C	C
CH5D	D
CH5E	[blank]
CH5X	[blank]
CH5NR	[blank]

Multipart Questions

Multipart questions are questions that contain two or more sub-questions, grouped together as items, with a leading question. In such cases, the response categories are the same for all sub-questions, which are usually items in relation to the leading question. Response categories can be numeric or alphabetical. For such questions, each sub-question has a designated numeric/string variable in the relevant dataset, and those variables are named ending with the letter that corresponds to the relevant sub-question.

CH11. Does your household have :	YES	NO
[A] A fixed telephone line?	FIXED TELEPHONE LINE.....1	2
[B] A radio?	RADIO.....1	2
[C] A wardrobe?	WARDROBE	1 2
[D] A cupboard?	CUPBOARD	1 2
[E] A table?	TABLE.....1	2
[F] A chair?	CHAIR.....1	2
[G] A bed?	BED.....1	2

For example, variables "CH11A", "CH11B", ..., "CH11G" represent, respectively, sub-questions "A", "B", ..., "G" of the question "CH11" above. Each of these variables have values of "1 - Yes" or "2 - No".

Recoded variables

In addition to variables that correspond to the questions in the questionnaires, each dataset also has a number of recoded or computed variables that are necessary for analysis. Names of such variables are as self-explanatory as possible. For example, the variable "area" represents "Area"; variable "headage" stands for "Age of household

head"; and the variable "windex5" is for "Wealth index quintiles". Most of these recoded/computed variables are used as background characteristics in the tabulations. More details on the construction of these variables can be found in the Section 8 "Background characteristics".

Sample weight variables

MICS Plus survey samples are not self-weighting; therefore, datasets also contain sample weight variables, which need to be used while carrying out analysis (unless the unweighted analysis is needed for a specific purpose). Separate sample weights are calculated for households and randomly selected individuals (see Section 2 "Typology of Questions"). The table below lists the sample weight variables for each wave, given that different waves have different groups of individuals selected randomly.

Variable name	Variable description	Wave 1	Wave 2	Wave 3	Wave 4	Wave 5	Wave 6
hhweight	Sample weights of households	x	x	x	x	x	x
chschage217weight	Sample weights of randomly selected children age 2-17 years at the beginning of school year	x				x	x
ch118weight	Sample weights of randomly selected children age 1-18 years		x				
ch017weight	Sample weights of randomly selected children age 0-17 years			x			
ch519weight	Sample weights of randomly selected children age 5-19 years				x		

Sample weights of households for each wave are calculated by adjusting the basic sample weights (inverse of selection probabilities) that are calculated for each cluster (HH1) by the non-response for that cluster in the wave. Sample weights of selected individuals are calculated by multiplying the sample weights of households by the number of individuals of interest. For example, $ch118weight = hhweight * \text{Number of children age 1-18 years in the household}$, ensuring that the data obtained from randomly selected individuals are expanded to the survey population in the selected group of interest.

6

Special Values

Various codes are used to describe special values that apply to a large set of variables, as a matter of convention.

"Not applicable" and "Missing" values

A "Not applicable" value occurs when a question is not supposed to be asked or should be skipped according to the flow of the questionnaire. On the other hand, a "Missing" value for a question applies when the question is supposed to be asked, but either the respondent has refused to provide an answer, or the question was not asked due to a technical error. "Not applicable" values are treated as system-missing in SPSS datasets, while "Missing"

values are coded with a 9, 99, 999, or 9999 depending on the field length of the variable (or with question marks for a string variable) and are treated as user-missing values (different statistical software may handle the "not applicable" and "missing" values differently). It is important to note that the "Not applicable" and "Missing" values can be included or excluded in analyses depending on the indicator of interest. Hence, one needs to pay careful attention to the selection of the denominator and the treatment of the "Missing" values for matching the results of MICS Plus surveys.

Other special values

In addition to the "Not applicable" and "Missing" values, there are often other special values that are usually pre-coded in the questionnaire with the following conventions:

	Numeric variable	String variable
Other	6, 96, 996, 9996, etc.	"X"
Don't know	8, 98, 998, 9998, etc.	"Z"

7

Numbers of Cases in Datasets

As explained in Section 1 "Survey Implementation", Wave 1 was implemented in 2 stages. Each record in the Wave 1 household dataset (HH) represents a single stage for a specific household. There were 2,118 initially selected households (Stage 1), of which 734 were substituted (Stage 2). Out of the 2,118 initially selected households, 1,384 were successfully interviewed in Stage 1 (before substitution) and 452 in Stage 2, yielding a total number of 1,836 interviewed households after the substitution. As a result, there are a total of 2,852 cases in the Wave 1 household dataset (HH).

The following table presents the frequency distribution of households by the result of the interview after substitution (HH17) at the end of Wave 1.

Total	2,118
Interviewed	1,836
Refused	29
No eligible respondent	1
Phone number does not belong to sampled household	33
Phone number inactive	65
Respondent busy / postponed	0
No response after repeated call attempts or phone(s) turned off	93
No phone number available for sampled household ⁴	60
Other	1

⁴ These are the households with no phone number available in the Georgia MICS 2018 data and no calls were made for them.

Upon completion of Wave 1, the final effective sample size for the consecutive waves was established at 1,996 households by deciding to exclude households with the following interview result codes: "Refused", "Phone number does not belong to sampled household", and "No phone number available for sampled household". Therefore, there are a total of 1,996 cases in all household datasets starting from Wave 2.

The call attempts datasets (CA) contain one record for every call attempt made for each selected household with an available phone number.

The household members datasets (HL) have one record for each household member listed in the interviewed households. The same line number of each household member is kept across different waves. These datasets include information on household members who are currently living in the household, as well as those who are reported as migrated or deceased during a particular wave. Any analysis pertaining to the current household composition (e.g., estimating the number of household members) should only include household members who are reported as currently living in the household.

8

Background characteristics

Results of the Georgia RTM/MICS Plus 2020-2021 survey are presented in the form of tabulations that include disaggregation of indicators by background characteristics. Variables for these background characteristics are created by recoding other existing variables. Below are the most common background characteristics.

Area

Variable name: area (Area)

Relevant waves: all

This variable indicates whether the household's address at the time of the interview was from an urban or rural settlement.

Given the possibility that households can change their addresses between waves, in each wave, the interviewed households are asked whether they are still living at the same address as they were living during the previous wave (variable "CA9C"). If not, the new address is recorded. This variable refers to the area of residence of the households at the time of interview of a particular wave, while the variable "HH6" is the area of residence of the household at the time of Georgia MICS 2018 data collection.

Sex and Age of Household Head

Variable name: headsex (Sex of household head), headage (Age of household head)

Relevant waves: all

In Wave 1, demographic information on all household members (date of birth, age, sex, relationship to household head) have been updated based on the information available from the Georgia MICS 2018. In

all subsequent waves, the respondent is asked whether there has been a change in the composition of the household (variable "HL0"). If there has been a change, the household list is updated, and a question is asked to ascertain whether the household head recorded in the previous wave is still the household head (variable "HL6B"), and if so, who the household head is at the time of the interview of that particular wave (variable "HL6C").

Sex and Age of Respondent

Variable name: respsex (Sex of respondent), respage (Age of respondent)

Relevant waves: all

MICS Plus surveys aim to have a knowledgeable adult member living in the household as the respondent to the questionnaire. While it is common to have the same respondent across waves, note that different respondents may be interviewed across different waves.

Wealth Index Quintile

Variable name: windex5 (Wealth index quintile)

Relevant waves: Wave 2 and consecutive waves

The wealth index is a composite indicator of wealth. Starting from Wave 2, questionnaires have included questions on ownership of consumer goods (variables "CH11A" to "CH14I"), energy use (variables "EU1" to "EU6"), and water and sanitation (variables "WS1" to "WS6"). Based on these and a few more variables (persons per room, access to the internet, etc.) that are related to household wealth, the wealth index is constructed⁵.

Note that the questions for constructing the wealth index were asked for all households only during the Wave 2. In the subsequent waves, these questions are administered only to households for which data on household characteristics were not collected in one of the previous waves that included these questions. Those can be 1) households who have changed their address of living between the waves and 2) households who are interviewed for the first time in that particular wave. Wealth index construction is repeated for each wave.

⁵ To construct the wealth index, principal components analysis is performed by using information on the ownership of consumer goods, dwelling characteristics, water and sanitation, and other characteristics that are thought to relate to the household's wealth, to generate weights (factor scores) for each of the items used. First, initial factor scores are calculated for the total sample. Then, separate factor scores are calculated for households in urban and rural areas. Finally, the urban and rural factor scores are regressed on the initial factor scores to obtain the combined, final factor scores for the total sample. This is carried out to address the urban bias in the wealth index values. Each household in the total sample is then assigned a wealth score based on the assets owned by that household and on the final factor scores obtained as described above. The survey household population is then ranked according to the wealth score of the household they are living in and is finally divided into 5 equal parts (quintiles) from lowest (poorest) to highest (richest). The wealth index is assumed to capture the underlying long-term wealth through information on the household assets and is intended to produce a ranking of households by wealth, from poorest to richest. It does not provide information on absolute poverty, current income, consumption, or expenditure levels. The wealth scores calculated are only applicable for the particular dataset they are based on.

9

Merging Datasets

For analysis purposes, it is possible to combine two or more MICS Plus datasets from different waves, if needed. It is also possible to match different types of datasets within a specific wave. This is only necessary when variables required for the analysis are not present in one file but are present in another. It should be noted that care has been taken to add the number of variables that are considered important for the analysis from one dataset to another. For example, variables on household and sample characteristics from the Household dataset (HH) are already included in the Household members dataset (HL). Nonetheless, there are occasions when data users have to merge different datasets to obtain the variables they need for a particular analysis. This section provides more details and examples on how to accomplish that task.

When merging datasets, the correct use of ID variables and identification of key variables are critical (see Section 5 "Variable Naming Conventions, Variable Construction" for more details on identification variables). Key variables are common variables between all source datasets, which link the observations of one data file to those of the other. The key variables must have the same names in all data files that are being merged. If names are not unique, renaming key variables in one or more datasets is required.

Another important step when merging MICS Plus datasets is to determine the type of relationship between two files, as well as to define the desired unit of analysis. For example, a relationship between households and household members is such that one entity (household) relates to several others (members of the household). There may be one or more household members for each household. This is a "one to many" relationship. On the other hand, in a "one to one" relationship, one dataset's entity is associated with one and only one entity in another dataset. For example, the household dataset from Wave 2 and the household dataset from Wave 3 have a "one to one" relationship.

Merging Datasets from Different Waves

There are two ways of merging datasets from different waves. One is to combine files containing the same variables but different cases. This is particularly useful when analyzing the same questions included in the questionnaires from different waves. With this type of merging, all cases from different datasets are concatenated, and in the resulting dataset all cases from one file are added to the end of all cases from another file.

The following example of SPSS code concatenates cases from Wave 2 and Wave 3 "HH" datasets. The resulting file contains all cases from the Wave 2 "HH" dataset and all cases from the Wave 3 "HH" dataset.

```

* Example of combining Wave 2 and Wave 3 variables on water and sanitation from household datasets.

* open the Wave 2 household dataset.
get file = "GEOHHW02.sav".

* sort data by ID variables.
sort cases by HHID.

* save working dataset and keep only the variables of interest; Household ID (HHID), wave number (HH0), result of interview (HH17), main source of drinking water (WS1), type of toilet facility (WS5), and sample weight variable (hhweight).
save outfile = "tmpHHW2.sav"
  /keep HHID HH0 HH17 WS1 WS5 hhweight.

* repeat above steps for the Wave 3 household dataset.
get file = "GEOHHW03.sav".
sort cases by HHID.
save outfile = "tmpHHW3.sav"
  /keep HHID HH0 HH17 WS1 WS5 hhweight.

* open the working Wave 2 dataset.
get file = "tmpHHW2.sav".

* combine files and note that it is not necessary to identify the ID variable.
add files
  file = *
  /file = "tmpHHW3.sav".

* save the combined dataset that now has cases from both waves.
save outfile = "tmpHHW2and3.sav".

* erase temporarily created files.
erase file = "tmpHHW2.sav".
erase file = "tmpHHW3.sav".

```

For certain analytical purposes, datasets from different waves can be merged by combining files that contain the same cases but different variables. For example, one might want to merge information on "Method of protection against COVID-19" (variables "CV31A" to "CV31NR") from Wave 1 household "HH" dataset to Wave 2 and Wave 3 household "HH" datasets. With this type of merging, key variables between datasets must be identified and used to match observations between them. For example, Cluster Number (HH1) and Household Number (HH2) are key variables that indicate which case from one household data file corresponds to which case from another. Similarly, Cluster Number (HH1), Household Number (HH2), and Line Number of Household Member (HL1) are key variables that match cases between household members (HL) datasets from different waves. Furthermore, the relationship between respective datasets from different waves is "one-to-one". This means that the new, merged dataset contains added variables of interest and has the same number of cases as the original dataset.

The SPSS syntax below demonstrates how to merge data from the Wave 1 Household "HH" dataset onto the Wave 2 Household "HH" dataset.

* Example of merging information on "Method of protection against COVID-19" from Wave 1 household dataset onto Wave 2 household dataset.

* open the Wave 1 household dataset.

```
get file = "GEOHHW01.sav".
```

* select only households after substitution, as they are present in consecutive waves.

```
select if (aftersub = 1).
```

* sort data by key variables.

```
sort cases by HH1 HH2.
```

* save working dataset and keep only the variables of interest; the key variables (HH1, HH2), methods of protection (CV31A to CV31NR).

```
save outfile = "tmpHHW1.sav"
```

```
  /keep HH1 HH2 CV31A CV31B CV31C CV31D CV31E CV31F CV31G CV31H CV31I  
    CV31J CV31K CV31X CV31Y CV31Z CV31NR.
```

* open the Wave 2 household dataset.

```
get file = "GEOHHW02.sav".
```

* sort data by key variables.

```
sort cases by HH1 HH2.
```

* merge files and note that it is critical to identify the key variables.

```
match files
```

```
  /file = *
```

```
  /table = " tmpHHW1.sav"
```

```
  /by HH1 HH2.
```

* save updated Wave 2 dataset with added information on methods of protection against COVID-19.

```
save outfile = "tmpHHW2.sav".
```

* erase temporarily created files.

```
erase file = "tmpHHW1.sav".
```

Merging Datasets of Same Wave

As described above, there will be instances when different types of datasets will have to be merged to obtain the variables that meet analysis needs. In the MICS Plus context, that would usually relate to adding household level information from the Household "HH" dataset onto the Household Members "HL" dataset or to the Call Attempts "CA" dataset. It is also possible to merge aggregated information from one of the datasets onto another one. For example, information on the number of household members per household can be added from the Household Members "HL" dataset onto the Household "HH" dataset, first by aggregating information on the number of household members for each household from the "HL" dataset, and second by adding that information onto the "HH" dataset.

When merging household members and call attempts datasets with their households, one needs to use Household ID (HHID) as a key variable. Since there is a "one-to-many" relationship between households and household

members, as well as between households and call attempts, it is possible to merge the "HH" dataset onto "HL" or "CA" datasets, but not the other way around.

The SPSS code that merges Wave 2 "HH" dataset onto the Wave 2 "HL" dataset is provided in the following example.

```

* Example on how to merge variable "Main source of drinking water" (WS1) from the household dataset onto the household
members dataset.

* open the household dataset.
get file = "GEOHHW02.sav".

* select only the interviewed households.
select if (HH17 = 1).

* sort according to the ID variables.
sort cases by HHID.

* save the working dataset temporarily by keeping only the variables of interest.
save outfile = "tmp hh.sav"
  /keep HHID WS1.

* open the household members dataset.
get file = "GEOHLW02.sav".

* sort according to the ID variables.
sort cases by HHID HMEMID.

* perform the matching by the ID variable that is common to both (tmp hh and HL) datasets.
match files
  /file = *
  /table = "tmp hh.sav"
  /by HHID.

* save the household members dataset that now has the variable WS1 added.
save outfile = "GEOHLW02.sav".

* erase temporarily created files.
erase file = "tmp hh.sav".

```