

## **Making Life Easier on the SAS® End: Best Practices for Collecting Survey Data**

Julie Plano, Keli Sorrentino, Yale University

### **ABSTRACT**

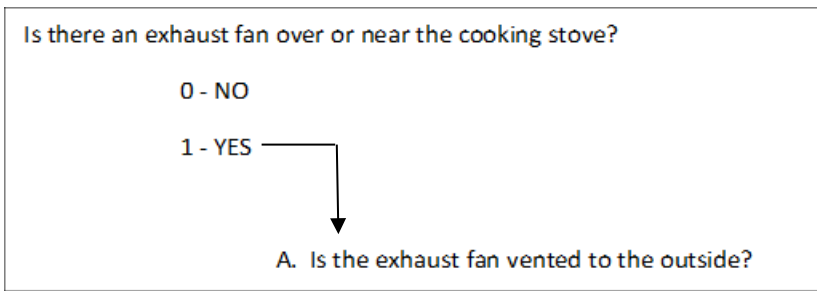
An analysis-ready SAS dataset is the deliverable, the one and only thing the analysts are waiting for. It is the ultimate goal of any research project. Much effort goes into the collecting and cleaning of data before it is ready to be analyzed. Our research team has spent decades learning to design good data collection instruments, resulting in purposeful survey questions with practical response values. Whether the questionnaire is administered in person or collected via online software, it is important to begin with a clear and concise plan. Here we share our expertise and best practices to help research teams design data collection instruments effectively and efficiently from the ground up. Our example will use Qualtrics Research Suite to design and gather data, and SAS 9.4 to create a dataset for further investigation and analysis.

### **INTRODUCTION**

Research begins with a question. An investigation is launched to gather information and find the answer. This process is data collection. When collecting information, the questions asked must be concise and relevant. The variables collected should have purpose to ensure the outcome can be evaluated with accuracy. Building a thoughtful survey is the first step to increasing our knowledge on the subject in question. Considerations should include details about the population being measured, the collection mechanism, and survey layout. Communication is key! This is an optimal time to close the gap between field researchers, programmers, and investigators. Errors should be anticipated; the data will not be perfect but there are several ways to set your project up for success. Utilize tools such as skip patterns, variable constraints, and SAS edit check programming. The data management techniques presented in this paper are intended to take the user from the very beginning of a research study and implement best practices to ensure the datasets imported from the field data entry system are accurate and clean, requiring little or no manipulation before analysis.

### **DATA COLLECTION**

All data collected will need to be linked to a specific person, case or location. Unique identifiers should be assigned at the earliest possible time to enable accurate tracking of collected data for each case. When creating an instrument consider the population being surveyed. Questions should be written at a sixth-grade reading level so that all possible audiences are given the opportunity to comprehend each word or phrase. Don't use abbreviations or slang words as these can cause the meaning of the question to change based on the study subject's perception. Focus on one topic at a time, include only one question, and use a follow-up question if needed for further clarification (Figure 1).



**Figure 1. Example Follow-up Question**

Finally, keep it short; questions should be concise and to the point. Extra verbiage could confuse or even sway the study subject to answer differently. Planning how response values will be assigned and recorded will help streamline the process from field collection to analysis-ready.

### **SUBJECT IDENTIFIER**

The first step in tracking all pieces of your study data is assigning a unique identifier to each case. A study ID of some type is typically used in research to de-identify the data and make it confidential. But this number can be incredibly useful with some thoughtful planning. A numerical study ID should be assigned to every potential study case as soon as a person or location is added to the project. This study ID can be used to track every piece of data associated with the case including personal information, questionnaire data, environmental samples, and results of sample analysis. Projects may require data be stored in multiple, separate datasets and formats. The study ID makes it easy to link every piece of data for each case together in SAS to compile all pieces of information. Having several options for collecting and processing raw data at your disposal enables the research team to adjust to the needs of all aspects of different research projects. Access databases, Excel spreadsheets, and online survey tools such as Qualtrics are some data collection options. Researchers often choose to put all the data in one Excel spreadsheet, adding each set of data to one row as it is collected or developed. While this method can work, input of data is challenging, and errors are common. Field questionnaire data may be collected using paper and pen and then later input into an Access database or collected directly as electronic data using Qualtrics. Figure A1 (see Appendix) illustrates how to use a unique study ID across multiple data collection and storage options. Figure A2 incorporates the unique identifier for a family and includes the individual and samples.

### **QUESTION DEVELOPMENT AND LAYOUT**

Good data collection instruments are critical to producing quality research data. Begin with having a clear list of the data you need to collect. Determine any covariates that will be important to the final analysis questions. Once you know what you want to ask, steps should be taken to develop how the questions will be asked and options for collecting responses. This may sound simple, but it can be very tricky to ensure the data you collect will correctly answer your research questions. Bring the entire research team together; each group—principal investigators, data management, and the field team— will have different input which will help make the instrument’s design solid. Having a quality design will help with your data analysis and also ensure consistent data collection. Questions and their answer choices should be clear, so they are not left open to interpretation.

Consider the following when creating your research survey: Am I really asking the question I want to ask? Will the answers I get give me the data I need? Often the original question is too broad in its scope to provide useful data in the end. The question may need to be broken down into a set of questions to get at the answer you really want.

Question A (Vague)	Question B (Clear)
Do you have pets?	Do you keep any of the following pets in your home? 1 = Dog 2 = Cat 3 = Rabbit 4 = Guinea pig, mouse, rat 5 = Non-furbearing animal

**Figure 2. Question Comparison**

In Question A, Figure 2, a yes answer could mean a variety of things. The 'yes' could be a horse they keep elsewhere, or it might be a fish. If you only care about fur-bearing animals in the home, that is the question you need to ask (Question B, Figure 2). Once you determine you are asking the questions you need, you should consider if each question is applicable to everyone or if a lead-in question would reduce confusion (Figure A3).

Once all the survey questions are developed it is important to carefully consider the order of the questions in the instrument. Does the instrument have good flow? Are the questions in a logical order? Will the order of the questions help the subject stay on track when answering? Transitional phrases should be added to help smooth the administration of the instrument (i.e. "Next I would like to ask some questions about appliances in your home").

Test the instrument. Reading through a questionnaire to yourself is very different from administering it to another person. This is an important step and testing should be done several times with different members of the design team. The instrument should also be tested under the same circumstances it will be used in the field, whether that is in person, over the phone, or an online survey completed by the subject. Testing across a wide group of people will help identify any issues with deployment or instrument design and allow time to rectify problems.

At the development phase, communication between the principal investigator (PI), data team (coordinator, manager, programmers, analysts), and field researchers is invaluable. Each have their own goals and perspective. The PI has a question to answer. The data team has the task of drilling down to the core of this issue to guarantee the right questions are asked and the responses provide valid and valuable information. The field researchers have the task of data collection and will consider how the subjects will respond and if the interview can be followed with ease. Working together and communicating early on creates a successful platform to ensure a scientifically sound study.

## **RESPONSE VALUES**

Why create more work? Thoroughly examine variables and their values at the development phase because this is what will make your life easier on the SAS end. Always use numbers as codes instead of multiple yes/no or check boxes. Figure A4, asks "Which of the following best describes the home?" On the left we have a list of choices, the correct value intended to be marked with a check. When viewing this type of data in a table several things could happen. SAS returns it as a text value and if not entered the same each time, more response categories than anticipated will be created. Data entry personnel may assign their own numbering system. The online data entry system captures the correct information but then converts it by autonumbering your response. If all responses are recorded consistently then this may not present a problem. Also consider if the list of responses is not inclusive of all potential options, necessitating the addition of another choice partway through the study. The data entry person could become inconsistent with their numbering if they are not aware of the new response choice. The online entry

mechanism (i.e. Qualtrics) could re-code the values entirely. When discussing instrument design with your team, request one variable for each question and code appropriately with numeric responses. This will reduce error and save time re-coding variables during the SAS programming phase.

For accuracy, it is better not to ask the study subject to do calculations (Figure A5). Instead, get their best answer and then let SAS do any computations. Offer ranges when that will provide accurate enough data (Figure A6). This will help to avoid the “I don’t know” responses.

Despite good question design there is a chance the only valid response is “I don’t know” or it is left missing. Missing values can happen because the question was overlooked in error, the response values did not encompass an accurate choice, or the question was refused. In research it is important to differentiate between different kinds of non-response. If the study subject responds “I don’t know” as their answer this is also important data to capture. Options for collecting missing data could include using a numeric value in the choice list as well as having a code to distinguish the type of missing response. For example, 8 = Don’t know, 9 = Not Applicable. At the SAS level we could look through our response values using edit check programming (see section on Error Reduction) to help the coders determine if the answer was truly missing or an error occurred that can be resolved.

Where possible you should always avoid open-ended questions that lead to text answers. If you can’t use PROC FREQ to list response values easily then your cleaning and/or analysis will require more work. Targeted comment variables are sometimes necessary to collect information that cannot be put into a category. Valid and consistent responses are the goal for any research project. However, there are times when a question may require an “other” to accommodate unforeseen responses. When this happens a ‘comment’ or ‘specify’ variable can be attached to ensure all the information is being recorded. For example, when asking about race (Figure 3), we can numerically categorize what the common responses are but not all individuals feel they fit into one of our choices. The “Other” option can contain a line that will be input at data entry. In SAS we can easily look at “variable\_x” responses and accurately re-code them if applicable or even create a new group. Having the text entered as data makes it easy to look through the responses in a SAS dataset.

1 – American Indian/Alaska Native
2 – Asian
3 – Native Hawaiian or Other Pacific Islander
4 – Black or African American
5 – White
6 – More than one race
7 – Other, Specify: _____
<b>VARIABLE_X</b>

**Figure 3. Example Response Choices for Race**

## DATA ENTRY MECHANISMS

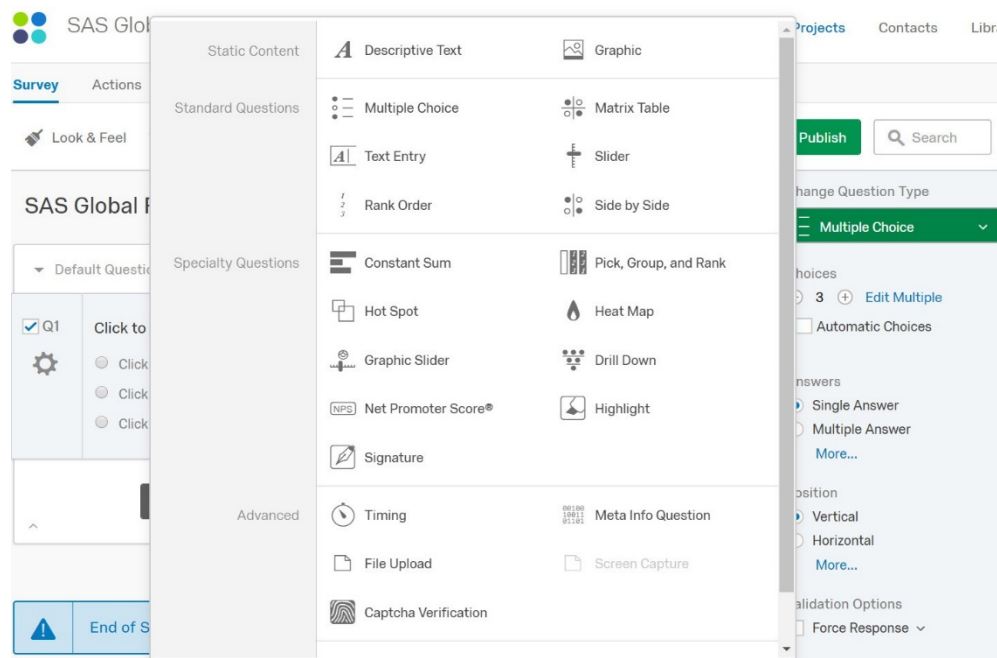
While there are many options for collecting data directly into an electronic format, for most of the research projects done at our Center we have found nothing beats paper and pen to enable the data collector to adequately record all the information provided by the study subject. Once the data are collected, reviewed, and verified for correct coding, it is then entered into an electronic format. Typically, this involves double-entering the data into a form in Access. We have had great success with this method for many large research projects. However, we also find alternatives to collecting data are required in certain circumstances. Qualtrics online survey software is ideal for creating an instrument that can be completed by the study subject or administered by a member of the research team in the field. Qualtrics is simple to use and it is easy to customize a survey to your needs whether small or large. With Qualtrics you can format questions and collect responses in multiple ways. Then, as a first step to reducing errors and making your data analysis ready, you can easily set up skip patterns and place constraints on the responses allowed for a question.

Being ready for the unforeseen is critical. We provide our field staff with paper versions of all electronic data collection tools in case technical difficulties arise. Data can be entered to Qualtrics when the staff member returns to the study office.

## ERROR REDUCTION

### PHASE 1

Reduce the number of errors during data collection by making use of the design options in Qualtrics. The first step is to choose the question type and begin creating the survey (Display 1). Use the “Text Entry” question type to collect data. This format provides the best control over the variable options available for data entry. If the standard “Multiple Choice” question format is used, then Qualtrics will assign the variable values which can cause issues if additions or changes are required later.



Display 1. Qualtrics Research Suite – Question Type

Display 2 shows how to add several constraints to variable responses. Choose if the question must be answered by forcing a response. Use Custom Validation to limit the types of data that can be input or Content Validation to format the input of the data into a date, phone number or other option. Use skip patterns and display logic to collect only the data required for each situation. Errors will be reduced when questions do not display if the response is not applicable. This also simplifies the instrument completion for the research staff and study subject.

The screenshot displays the Qualtrics Research Suite interface for configuring a question. The main area shows a question titled "Have you enjoyed the SAS 2019 Global Forum?" with a response type of "Text Entry". Below the question, a condition is defined: "Condition: Have you enjoyed the SAS 20... Is Equal to 1. Skip To: Will you be planning to attend the SA...". The right-hand sidebar contains configuration options for the question, including "Text Type" (Single Line, Multi Line, Essay Text Box, Form, Password), "Validation Options" (Force Response), "Validation Type" (None, Minimum Length, Maximum Length, Character Range, Content Validation, Custom Validation), and "Actions" (Add Page Break, Add Display Logic, Add Skip Logic, Copy Question). The "Force Response" and "Custom Validation" options are circled in red.

## Display 2. Qualtrics Research Suite – Reducing Errors

## PHASE 2

Creating SAS error check programs are an efficient way to automate the process of scanning data for errors. These can be written as soon as the instruments are field ready. Programs are written for each separate instrument and can be adjusted to accommodate any changes as a study progresses. These types of programs can be as complicated or as simplistic as needed. Example Code A looks at one variable, Race. The overall method can be applied to an entire interview.

### **Example Code A**

The following code returns errors found for the variable Race (Figure 3). First the macro "Invalid" is created to customize the output for any errors found. Study ID, the variable name and the variable value are printed. There are seven valid responses for the variable Race1. An IF/THEN statement locates any values not within range. If a question is asked several times with the same ranges—for instance, here we ask the race of up to six people—an array can be included to run through each variable. Below, Output 1 shows an example of the output if an error is found within the data:

```
%MACRO INVALID(VARNAME,VALUE); PUT
    @8 Study_id
    @20 "INVALID " &VARNAME
    @65 &VARNAME " = " &VALUE " _____" /;
%MEND INVALID;
QUIT;

IF NOT (1<= Race1 <=7) THEN DO; %INVALID(' Race1 ', Race1); END;

array R6 (6) Race1 Race2 Race3 Race4 Race5 Race6; IF
COUNT = 6 THEN DO;
do i = 1 to 6;
IF R6[I] NOT in (1,2,3,4,5,6,7) THEN DO; %INVALID('R6[I] ',R6[I]); END;
END;
END;
END;
```

8972

INVALID RACE1

Q3\_C1 = 0

### **Output 1. The output from error macro %INVALID**

Example Code B checks the variable ranges as well as if skip patterns are being applied accurately. The question (HQ32) asked if the study subject had an experience. If yes (1), then the following sub-questions are asked and a response is required for each one. If any of the sub-questions are missing, an error is printed. Looking at a question as a whole makes it easier for the coders to locate exactly where the error is occurring. The macro here is written to print the three variables associated with the question.

### **Example Code B**

The following code returns errors found for a question with an option for a skip pattern. The initial response must be a '1' for the subsequent response to need values. If the follow-up

questions are missing an error is printed. The macro "Conflict3" prints several variables required for investigation:

```
%MACRO CONFLICT3
  (VARNAME_1,VALUE_1,VARNAME_2,VALUE_2,VARNAME_3,VALUE_3); PUT
  @8 STUDY_ID @13 "CONFLICTING DATA FOR --> " &VARNAME_1
  @63 STUDY_ID &VARNAME_1 " = " &VALUE_1 @85 " _____" /
  @63 STUDY_ID &VARNAME_2 " = " &VALUE_2 @85 " _____" /
  @63 STUDY_ID &VARNAME_3 " = " &VALUE_3 @85 " _____" /;
%MEND CONFLICT3;

IF HQ32 NOT IN (0,1) THEN DO; %INVALID ('HQ32', HQ32); END; IF HQ32 = 1
THEN DO;
IF HQ32A = '' THEN DO;
  %CONFLICT3 ('HQ32',HQ32,'HQ32A',HQ32A,'HQ32A_1',HQ32A_1);END;
```

## OUTPUT DATASETS

Qualtrics offers multiple file extensions for the downloading of data collected online (Display 3). Moving raw data out of Qualtrics to a tab-separated value file extension creates a simple transfer of data. Another option for downloading survey data is to use SAS to query the API and generate an XML map. In this example, the .tsv file is saved locally and we use SAS to import the survey data:

```
PROC IMPORT OUT= WORK.Example
  DATAFILE= "C:\QDATA\Int_111318.tsv" DBMS=TAB REPLACE;
  GETNAMES=YES;
  DATAROW=2
  ; RUN;
```

### Download Data Table

[Use Legacy Exporter](#)

CSV	<b>TSV</b>	XML	SPSS	User Submitted Files	Tableau
-----	------------	-----	------	----------------------	---------



#### Tab separated values

This is a .tsv file that can be imported into other programs. Each value in the response is separated by a tab and each response is separated by a newline character. If your responses contain special characters and you will open this export in Microsoft Excel we recommend using this TSV export because Qualtrics TSV exports use UTF-16 encoding.

[Learn More](#)

- Download all fields
- Use numeric values
- Use choice text

[More Options](#)

Close

[Download](#)

## Display 3. Qualtrics Research Suite – Download Data Table



The SAS import program is an opportune place to not only import data but to format variables, change variable types, account for missing information, assign labels, and create new variables required for analysis. The finished product can be completed at this step.

## **CONCLUSION**

Good planning is the first step to having a clean SAS dataset which is ready for analysis in the shortest time frame. Having members of the field team and the data management team participate in the planning will help streamline the data flow from the field. Considering question layout, response options, data collection methods, and data processing are all critical. Implementing SAS edit check programming early on will ensure data is being entered correctly and consistently. The SAS programmer can have vital input when developing data collection instruments. Attention to all these details will expedite getting that final SAS dataset to the investigators.

## **ACKNOWLEDGMENTS**

We would like to acknowledge our colleagues and friends at the Center for Perinatal, Pediatric & Environmental Epidemiology.

## **CONTACT INFORMATION**

Your comments and questions are valued and encouraged. Contact the author at:

Julie Plano  
Yale School of Public Health  
[Julie.plano@yale.edu](mailto:Julie.plano@yale.edu)

Keli Sorrentino  
Yale School of Public Health  
[Keli.sorrentino@yale.edu](mailto:Keli.sorrentino@yale.edu)

SAS and all other SAS Institute Inc. product or service names are registered trademarks or trademarks of SAS Institute Inc. in the USA and other countries. ® indicates USA registration.

Qualtrics and all other Qualtrics product or service names are registered trademarks or trademarks of Qualtrics, Provo, UT, USA. <https://www.qualtrics.com>

## APPENDIX

Access DB #1	Excel SS #1	Access DB #2	Access DB #3	Access DB #4	Access DB #4	Excel SS #2	Excel SS #3
Participant Tracking	GPS Data	Instrument Data	Field Biological Sample Collection	Sample Receipt at Lab	Aliquot Processing & Storage	Lab #1 Results	Lab #2 Results
<b>Stud_Num</b>	Stud_Num	Stud_Num	<b>Sample_ID</b>	Sample_ID	Sample_ID	Sample_ID	Sample_ID
50421	50421	50421	50421-101	50421-101	50421-202 50421-203 50421-305	50421-202 50421-203	50421-305

**Figure A1** shows how a unique study number (Study\_Num) is used to identify data collected from one study subject. The same number is used across all tables, forms and databases. The unique sample identifier is followed by the study number and can be found in the last 3 digits of the variable (Sample\_ID)

Access DB #1	Excel SS #1	Access DB #2	Access DB #3	Access DB #4	Access DB #4	Excel SS #2	Excel SS #3
Participant Tracking	GPS Data	Instrument Data	Field Biological Sample Collection	Sample Receipt at Lab	Aliquot Processing & Storage	Lab #1 Results	Lab #2 Results
<b>Family_ID</b>	Family_ID	<b>Individual_ID</b>	Sample_ID	<b>Sample_ID</b>	Sample_ID	Sample_ID	Sample_ID
50421	50421	50421-01	50421-01-101	50421-01-101	50421-01-202 50421-01-203 50421-01-305	50421-01-202 50421-01-203	50421-01-305
Individual_ID							
50421-01							
Next Individual							
Access DB #1	Excel SS #1	Access DB #2	Access DB #3	Access DB #4	Access DB #4	Excel SS #2	Excel SS #3
Participant Tracking	GPS Data	Instrument Data	Field Biological Sample Collection	Sample Receipt at Lab	Aliquot Processing & Storage	Lab #1 Results	Lab #2 Results
Family_ID	Family_ID	Individual_ID	Sample_ID	Sample_ID	Sample_ID	Sample_ID	Sample_ID
50421	50421	50421-02	50421-02-101	50421-02-101	50421-02-202 50421-02-203 50421-02-305	50421-02-202 50421-02-203	50421-02-305
Individual_ID							
50421-02							

**Figure A2** includes a unique identifier by family (Family\_ID). Unique individuals within a family are coded with two digits following the family identifier (Individual\_ID). The sample identifier is followed by the individual id (Sample\_ID)

## Problem

## Solution

<p>If you live in a shared building, to your knowledge, did anyone smoke inside the building?</p> <p><input type="checkbox"/> Yes <input type="checkbox"/> No <input type="checkbox"/> Don't know</p>	<p>Do you live in a shared building?</p> <p>0 – No 1 – Yes</p> <p>↓</p> <p>Did anyone smoke in the common areas of the building? This would include hallways, stairwells, lobby, etc.</p> <p>0 – No 1 - Yes</p>
---	---

**Figure A3 Reduce confusion by asking the question in multiple parts**

<p>Which of the following best describes the home? (Check correct answer)</p> <p><input type="checkbox"/> A one-family house detached from any other house <input type="checkbox"/> A one-family house attached to one or more houses (for example Row house/Town house) <input type="checkbox"/> A building with 10 or fewer apartments <input type="checkbox"/> A building with more than 10 apartments <input type="checkbox"/> Mobile home, Boat, RV, van, etc.</p>	<p>Which of the following best describes the home?</p> <p>1 – A one-family house detached from any other house 2 – A one-family house attached to one or more houses (Row house/Town house) 3 – A building with 10 or fewer apartments 4 – A building with more than 10 apartments 5 – Mobile home, Boat, RV, van, or other</p>
---	---

**Figure A4 Use a single variable to reduce data errors**

<p>How long have you lived in your home? _____</p>	<p>What year did you move into your home?</p> <table border="1"><tr><td> </td><td> </td><td> </td><td> </td></tr></table>				

**Figure A5 Let the computer do the calculations**

<p>When was your home first built? _____</p>	<p>About when was your home first built?</p> <p>1 – 2014 or later 2 – 2000 - 2013 3 – 1980 - 1999 4 – 1960 - 1979 5 – 1940 - 1959 6 – 1920 - 1939 7 – 1900 - 1919 8 – 1899 or earlier</p>
--	---

**Figure A6 Giving ranges for answers can reduce missing data circumstances**