Quantitative Data Analysis

Andy Grogan-Kaylor

2021-09-08

Contents

1	Publicly Available Tools for Analysis	1
2	Our Data	1
3	Cleaning Data	2
	3.1 Excel and Google Sheets	2
	3.2 R	3
	3.3 Stata	4
4	Simple Analysis	6

1 Publicly Available Tools for Analysis

Tool	Cost	Ease of Use	Analysis Capabilities	Suitability for Large Data	Keep Track of Complicated Workflows
Excel	Comes installed on many computers	Easy	Limited	Difficult when N > 100	Difficult to Impossible
Google Sheets	Free with a Google account	Easy	Limited	Difficult when N > 100	Difficult to Impossible
R	Free	Challenging	Extensive	Excellent with large datasets	Yes, with script
Stata	Some cost	Learning Curve but Intuitive	Extensive	Excellent with large datasets	Yes, with command file

2 Our Data

We take a look at our simulated data, which has an id number, age, and happiness (on a 5 point scale, with 5 being the happiest.)

id	age	happy	somethingelse
1	73.89	-99	-0.4032
2	200	2	1.091
3	49	-99	-0.662
4	31.48	5	-0.8626
5	60.03	2	0.053
6	56.67	4	1.274

Notice that...

- There are variables in which we may not have interest (e.g. somethingelse).
- None of the variables have informative variable labels. We have to guess at what the variables mean.
- Variables do not seem to have informative *value labels*. While somewhat intuitive, we have to guess at what the values mean.
- Someone appears to 200 years old.
- There appear to be missing values in the variable happy that need to be recoded.

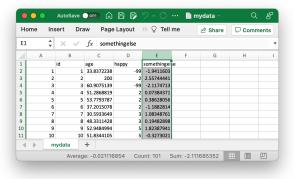
3 Cleaning Data

- 1. Only keep the variables of interest.
- 2. Add variable labels (if we can).
- 3. Add value labels (if we can).
- 4. Recode outliers, values that are errors, or values that should be coded as missing

3.1 Excel and Google Sheets

3.1.1 1. Only keep the variables of interest.

Select the column, or columns, of data that you wish to remove, and right click, or control click, to delete them.



3.1.2 2. Add variable labels (if we can).

We are unable to add informative labels to variables in Excel or Google Sheets.

3.1.3 3. Add value labels (if we can).

We are unable to add informative labels to values in Excel or Google Sheets.

3.1.4 4. Recode outliers, values that are errors, or values that should be coded as missing.

We are likely going to have to use **find and replace** to manually replace problematic values. For example, we will want to replace the 200 in the age column with a . or NA for missing. Similarly, we will want to replace the values of -99 in the happy column with a . or NA for missing.

For small data sets, this will not be difficult, but for larger data sets-especially data with many different kinds of values that need to be recoded-this process will become more difficult and cumbersome.

3.2 R

Much of R's functionality is accomplished through writing *code*, that is saved in a *script*. Notice how—as our tasks get more and more complicated—the saved script provides documentation for the decisions that we have made with the data.

3.2.1 1. Only keep the variables of interest.

We can easily accomplish this with the subset function

id	age	happy
1	73.89	-99
2	200	2
3	49	-99
4	31.48	5
5	60.03	2
6	56.67	4

3.2.2 2. Add variable labels (if we can).

Adding *variable labels* is not well established in R. There are libraries that can add variable labels for some purposes, but not every library in R recognizes value labels.

3.2.3 3. Add value labels (if we can).

In contrast, *value labels* are straightforward in R, and can be accomplished by creating a factor variable. Below we demonstrate how to do this with the happy variable.

id	age	happy	happyFACTOR
1	73.89	-99	NA
2	200	2	Somewhat Unhappy
3	49	-99	NA
4	31.48	5	Very Happy
5	60.03	2	Somewhat Unhappy
6	56.67	4	Somewhat Happy

3.2.4 4. Recode outliers, values that are errors, or values that should be coded as missing.

We can easily accomplish this using Base R's syntax for recoding: data\$variable[rule] <- newvalue.

```
mynewdata$age[mynewdata$age >= 100] <- NA # recode > 100 to NA

mynewdata$happy[mynewdata$happy == -99] <- NA # recode -99 to NA
```

id	age	happy	happyFACTOR
1	73.89	NA	NA
2	NA	2	Somewhat Unhappy
3	49	NA	NA
4	31.48	5	Very Happy
5	60.03	2	Somewhat Unhappy
6	56.67	4	Somewhat Happy

3.3 Stata

3.3.1 1. Only keep the variables of interest.

This is easily accomplished with Stata's drop command. We could also choose to keep our variables of interest.

```
drop somethingelse // drop extraneous variable(s)
```

3.3.2 2. Add variable labels (if we can).

Variable labels can easily be added in Stata.

```
label variable age "Respondent's Age'" // variable label for age

label variable happy "Happiness Score" // variable label for happy

describe // describe the data

Contains data from mydata.dta
```

Observation Variable		100 3		
Variable name	Storage type	Display format	Value label	Variable label
id age happy	long double double	%9.0g %9.0g %9.0g		id Respondent's Age' Happiness Score

3.3.3 Add value labels (if we can)

Value labels are a natural part of Stata.

```
label define happy /// create value label for happy
5 "Very Unhappy" ///
4 "Somewhat Unhappy" ///
3 "Neutral" ///
2 "Somewhat Happy" ///
1 "Very Happy"
label values happy happy // assign value label happy to variable happy
list in 1/10 // list first 10 lines of data
    | id
          age
                              happy |
    |-----|
 1. | 1 45.23996
 2. | 2 200 Somewhat Happy | 3. | 3 58.39718 -99 |
 4. | 4 46.66829 Somewhat Happy | 5. | 5 48.58828 Somewhat Happy |
 6. | 6 41.52565 Somewhat Unhappy |
 7. | 7 45.49311
                            Neutral |
 8. | 8 50.7754
                     Somewhat Happy |
 9. | 9 51.39233 Somewhat Unhappy |
 10. | 10 44.55724 Somewhat Happy |
```

3.3.4 4. Recode outliers, values that are errors, or values that should be coded as missing

```
recode age (100 / max = .) // recode ages > 100

recode happy (-99 = .) // recode -99 to missing
```

4 Simple Analysis