# Getting to Know Your Data Set in Stata 12

## Table of Contents

This is a guideline for Stata 12. Updated versions of Stata may have differing commands.

I. Data Set

For our data analysis, we will use a fake data set entitled **Test Survey**. The variables are the text and numerical values (these will have **num** in front of the variable name) for favorite season (**season** and **numseason**), number of cups of coffee consumed each day (**coffee** and **numcoffee**), favorite restaurant in a local town (**restaurant** and **numrestaurant**), favorite computer's name (**computer** and **numcomputer**), favorite meal in the dining hall (**peircemeals** and **numpeircemeals**), a question pertaining to the fate of the characters on the popular television show *Lost* (**lost** and **numlost**), how much they like Kenyon with respect to another school (**kenyon**  and **numkenyon**), students favorite dessert (**dessert** and **numdessert**), whether they watch a particular television show (**AHS** and **numAHS**), what their favorite season of said television show is (**favoriteAHS** and **numfavoriteAHS**), and students' favorite building on campus (**building** and **numbuilding**). This survey has 50 observations. We are going to use these variables to examine some simple summary statistics and a simple regression. We will use the names of our variables. If you were to replicate this analysis, use your variable names rather than our variable names.

II. Acquiring your Data Set from Qualtrics

Most of your data sets will come from the survey collection website Qualtrics. To get your desired data set, go to kenyon.qualtrics.com.

1. To find your desired data, search the name of the survey in the search bar in the top right corner (it will be next to "Create Project"). Once you have found your survey, click "Data & Analysis then select "Export & Import" and choose "Export Data" in order to acquire the data set.

2. Upon choosing "Export Data", you will then want to select "Download Data Table". Several options will appear as to which format you can download your data as. This is up to you, but we chose "CSV" so that we could keep the choice text. We then opened this data in Excel, and copy and pasted it into Stata's Data Editor.

3. To get the numeric data, we again download the data from Qualtrics, but this time choose "SPSS". Then you will go to the SPSS statistical package and copy the desired numerical values from its "Data View" tab and paste them into Stata's Data Editor.

4. These numerical values will be named "var1", "var2", etc. To help make more sense of the data, use "rename var12 numseason".

```
. rename var12 numseason

. rename var13 numrestaurant

. rename var14 numcomputer

. rename var15 numpeircemeals

. rename var16 numlost

. rename var17 numkenyon

. rename var18 numdessert

. rename var19 numAHS
```

5. Finally, reorder your data set so that choice text is next to the numerical values associated with each variable. To do this, use "order season numseason coffee numcoffee…building numbuilding".

| | season | numseason | coffee | numcoffee |
|---|---|---|---|---|
| 1 | Fall | 3 | 5 or more | 4 |
| 2 | Spring | 4 | 5 or more | 4 |
| 3 | Winter | 1 | 5 or more | 4 |
| 4 | Fall | 3 | 2-3 | 2 |
| 5 | Winter | 1 | 0-1 | 1 |
| 6 | Winter | 1 | 3-4 | 3 |
| 7 | Spring | 4 | 0-1 | 1 |
| 8 | Winter | 1 | 2-3 | 2 |
| 9 | Summer | 2 | 5 or more | 4 |
| 10 | Fall | 3 | 2-3 | 2 |
| 11 | Spring | 4 | 2-3 | 2 |
| 12 | Fall | 3 | 2-3 | 2 |
| 13 | Spring | 4 | 2-3 | 2 |
| 14 | Winter | 1 | 5 or more | 4 |
| 15 | Summer | 2 | 2-3 | 2 |
| 16 | Summer | 2 | 0-1 | 1 |
| 17 | Winter | 1 | 0-1 | 1 |
| 18 | Fall | 3 | 3-4 | 3 |
| 19 | Winter | 1 | 0-1 | 1 |
| 20 | Fall | 3 | 3-4 | 3 |
| 21 | Spring | 4 | 0-1 | 1 |
| 22 | Summer | 2 | 3-4 | 3 |
| 23 | Spring | 4 | 5 or more | 4 |
| 24 | Summer | 2 | 3-4 | 3 |
| 25 | Fall | 3 | 3-4 | 3 |
| 26 | Fall | 3 | 2-3 | 2 |
| 27 | Spring | 4 | 0-1 | 1 |
| 28 | Summer | 2 | 5 or more | 4 |
| 29 | Fall | 3 | 3-4 | 3 |
| 30 | Fall | 3 | 0-1 | 1 |
| 31 | Fall | 3 | 0-1 | 1 |
| 32 | Spring | 4 | 0-1 | 1 |
| 33 | Summer | 2 | 0-1 | 1 |
| 34 | Fall | 3 | 5 or more | 4 |
| 35 | Summer | 2 | 3-4 | 3 |
| 36 | Spring | 4 | 3-4 | 3 |
| 37 | Fall | 3 | 0-1 | 1 |
| 38 | Winter | 1 | 3-4 | 3 |
| 39 | Winter | 1 | 2-3 | 2 |
| 40 | Winter | 1 | 3-4 | 3 |
| 41 | Spring | 4 | 2-3 | 2 |

III. Stata Output

Getting to know your data set:

1. For a quick overview of your data set, use the command "codebook" to get the total number of observations for each variable and the frequency for each possible response.

```
.  codebook

─────────────────────────────────────────────────────────────────────────────
season                                                                   QID1
─────────────────────────────────────────────────────────────────────────────

              type:  string (str34), but longest is str6

      unique values:  4                     missing "":  0/50

         tabulation:  Freq.  Value
                        17   "Fall"
                        13   "Spring"
                         9   "Summer"
                        11   "Winter"

─────────────────────────────────────────────────────────────────────────────
numseason                                                          (unlabeled)
─────────────────────────────────────────────────────────────────────────────

              type:  numeric (byte)

             range:  [1,4]                      units:  1
      unique values:  4                     missing .:  0/50

         tabulation:  Freq.  Value
                        11   1
                         9   2
                        17   3
                        13   4

─────────────────────────────────────────────────────────────────────────────
coffee                                                                   QID2
─────────────────────────────────────────────────────────────────────────────

              type:  string (str58), but longest is str9

      unique values:  4                     missing "":  0/50

         tabulation:  Freq.  Value
                        15   "0-1"
                        11   "2-3"
──more──
```

2. Use the "tab1 season coffee...building" command to get the frequency and percent of each response for each variable. You could also use the "tabulate season" command to just get the frequency and percent for one variable and its responses.

. tab1 season coffee restaurant computer peircemeals lost kenyon dessert AHS favoriteAHS building

-> tabulation of season

| QID1 | Freq. | Percent | Cum. |
|---|---|---|---|
| Fall | 17 | 34.00 | 34.00 |
| Spring | 13 | 26.00 | 60.00 |
| Summer | 9 | 18.00 | 78.00 |
| Winter | 11 | 22.00 | 100.00 |
| Total | 50 | 100.00 | |

-> tabulation of coffee

| QID2 | Freq. | Percent | Cum. |
|---|---|---|---|
| 0-1 | 15 | 30.00 | 30.00 |
| 2-3 | 11 | 22.00 | 52.00 |
| 3-4 | 14 | 28.00 | 80.00 |
| 5 or more | 10 | 20.00 | 100.00 |
| Total | 50 | 100.00 | |

-> tabulation of restaurant

| QID3 | Freq. | Percent | Cum. |
|---|---|---|---|
| Buffalo Wild Wings | 9 | 18.00 | 18.00 |
| Chipotle | 9 | 18.00 | 36.00 |
| Fiesta | 7 | 14.00 | 50.00 |
| Panera | 8 | 16.00 | 66.00 |
| Ruby Tuesday's | 6 | 12.00 | 78.00 |
| Southside Diner | 11 | 22.00 | 100.00 |
| Total | 50 | 100.00 | |

-> tabulation of computer

| QID4 | Freq. | Percent | Cum. |
|---|---|---|---|
| Bunny | 16 | 32.00 | 32.00 |
| Peg | 10 | 20.00 | 52.00 |
| Ruthie | 15 | 30.00 | 82.00 |

. tabulate season

| QID1 | Freq. | Percent | Cum. |
|---|---|---|---|
| Fall | 17 | 34.00 | 34.00 |
| Spring | 13 | 26.00 | 60.00 |
| Summer | 9 | 18.00 | 78.00 |
| Winter | 11 | 22.00 | 100.00 |
| Total | 50 | 100.00 | |

.

3. To compare different variables to one another, use the command "tabulate season coffee, column". This command will give you the frequency and percentage of a data point fulfilling two different criteria (i.e. how many responses answered "Fall" and "0-1" cups of coffee and so forth).

. tabulate season coffee, column

```
Key

  frequency
column percentage
```

|  | | QID2 | | | |  |
|---|---|---|---|---|---|---|
| QID1 | 0-1 | 2-3 | 3-4 | 5 or more | | Total |
| Fall | 5 | 4 | 5 | 3 | | 17 |
|  | 33.33 | 36.36 | 35.71 | 30.00 | | 34.00 |
| Spring | 5 | 4 | 2 | 2 | | 13 |
|  | 33.33 | 36.36 | 14.29 | 20.00 | | 26.00 |
| Summer | 2 | 1 | 4 | 2 | | 9 |
|  | 13.33 | 9.09 | 28.57 | 20.00 | | 18.00 |
| Winter | 3 | 2 | 3 | 3 | | 11 |
|  | 20.00 | 18.18 | 21.43 | 30.00 | | 22.00 |
| Total | 15 | 11 | 14 | 10 | | 50 |
|  | 100.00 | 100.00 | 100.00 | 100.00 | | 100.00 |

.

4. To acquire summary statistics regarding the variables, you will have to use the numerical values. The command needed would be "summarize numseason numcoffee numrestaurant numcomputer numpeircemeals numlost numkenyon numdessert numAHS numfavoriteAHS numbuilding". This will give you summary statistics such as mean, number of observations, standard deviation, minimums, and maximums.

. summarize numseason numcoffee numrestaurant numcomputer numpeircemeals numlost numkenyon numdess
> ert numAHS numfavoriteAHS numbuilding

| Variable | Obs | Mean | Std. Dev. | Min | Max |
|---|---|---|---|---|---|
| numseason | 50 | 2.64 | 1.102132 | 1 | 4 |
| numcoffee | 50 | 2.38 | 1.122861 | 1 | 4 |
| numrestaur~t | 50 | 3.48 | 1.668924 | 1 | 6 |
| numcomputer | 50 | 2.28 | 1.107304 | 1 | 4 |
| numpeircem~s | 50 | 3.46 | 1.528104 | 1 | 6 |
| numlost | 50 | 2 | .6998542 | 1 | 3 |
| numkenyon | 50 | 2.14 | .7561989 | 1 | 3 |
| numdessert | 50 | 3.18 | 1.662184 | 1 | 6 |
| numAHS | 50 | 1.5 | .5050763 | 1 | 2 |
| numfavorit~S | 25 | 3.6 | 1.755942 | 1 | 6 |
| numbuilding | 50 | 4.22 | 1.87671 | 1 | 7 |

.

5. Use the "tabulate numseason, summarize(numcoffee)" to get summary statistics for a variable with respect to another variable.

. tabulate numseason , summarize(numcoffee)

| numseason | Summary of numcoffee Mean | Std. Dev. | Freq. |
|---|---|---|---|
| 1 | 2.5454545 | 1.2135598 | 11 |
| 2 | 2.6666667 | 1.118034 | 9 |
| 3 | 2.3529412 | 1.1147408 | 17 |
| 4 | 2.0769231 | 1.1151636 | 13 |
| Total | 2.38 | 1.1228608 | 50 |

.

6. In order to compare three variables to each other, use the commands "contract numseason numcoffee numrestaurant, percent(percent)" and then "tabdisp numseason numcoffee numrestaurant, c(percent)". The first command will contract the other variables so that only the ones of interest will be displayed. The second command will then create a display table that shows the number of responses that were in each combination of the three different variable's responses.

```
. contract numseason numcoffee numrestaurant , percent(percent)


. tabdisp numseason numcoffee numrestaurant , c(percent)
```

| | | | | | | | | | numrestaurant and numcoffee | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 1 | | | | 2 | | | | 3 | | | | 4 | | |
| numseason | 1 | 2 | 3 | 4 | 1 | 2 | 3 | 4 | 1 | 2 | 3 | 4 | 1 | 2 | 3 | 4 |
| 1 | | | | | | 2.00 | | | 6.00 | | 2.00 | | | 2.00 | | |
| 2 | | 2.00 | | | | 2.00 | 2.00 | | | | | | 2.00 | | 2.00 | 2.00 |
| 3 | | 2.00 | 2.00 | | 4.00 | | | 2.00 | | 2.00 | 2.00 | | 2.00 | | 2.00 | 2.00 |
| 4 | 2.00 | 4.00 | 4.00 | | 4.00 | 2.00 | | | 2.00 | | | | | 2.00 | | 2.00 |

| | | | numrestaurant and numcoffee | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | 5 | | | | 6 | | |
| numseason | 1 | 2 | 3 | 4 | 1 | 2 | 3 | 4 |
| 1 | | 4.00 | 2.00 | 2.00 | | | 2.00 | |
| 2 | 2.00 | | 2.00 | | | 2.00 | | |
| 3 | 2.00 | 2.00 | 2.00 | | 2.00 | 2.00 | 2.00 | 2.00 |
| 4 | 2.00 | | | 2.00 | | | | |

7. Finally, to run a regression to examine the causal relationships between two or more variables, use the command "regress numcoffee numseason", taking care to put the dependent variable first and the explanatory variable(s) after that. In other words, put the variable that you are trying to predict first and the variable(s) that help you predict this variable after.

```
. regress numcoffee numseason
```

| Source | SS | df | MS | | Number of obs = | 50 |
|---|---|---|---|---|---|---|
| | | | | | F( 1, 48) = | 1.39 |
| Model | 1.73430108 | 1 | 1.73430108 | | Prob > F = | 0.2448 |
| Residual | 60.0456989 | 48 | 1.25095206 | | R-squared = | 0.0281 |
| | | | | | Adj R-squared = | 0.0078 |
| Total | 61.78 | 49 | 1.26081633 | | Root MSE = | 1.1185 |

| numcoffee | Coef. | Std. Err. | t | P>|t| | [95% Conf. Interval] | |
|---|---|---|---|---|---|---|
| numseason | -.1706989 | .1449736 | -1.18 | 0.245 | -.4621878 | .12079 |
| _cons | 2.830645 | .4141274 | 6.84 | 0.000 | 1.997986 | 3.663304 |