

Data Visualization With Stata

Andy Grogan-Kaylor

10 Jul 2020

Introduction

- Stata is a powerful and intuitive data analysis program.
- Learning how to graph in Stata is an important part of learning how to use Stata. Yet, the default graphs in Stata can sometimes be less than optimal.
- This document is an introduction to (a) basic graphing ideas in Stata; and (b) a quick note on the use of schemes to make your Stata graphs look more professional.

What are Variables?

- By variables, I simply mean the columns of data that you have.
- For our purposes, you may think of variables as synonymous with questionnaire items, or columns of data.

Variable Types

- *categorical variables* represent unordered categories like *race*, *ethnicity*, *neighborhood*, *religious affiliation*, or *place of residence*.
- *continuous variables* represent a continuous scale like *income*, a *mental health scale*, or a *measure of life expectancy*.

A Data Visualization Strategy

Once we have discerned the type of variable that have, there are two followup questions we may ask before deciding upon a graphing strategy:

- Is our graph about **one thing at a time**?
 - How much of x is there?
 - What is the distribution of x ?
- Is our graph about **two things at a time**?
 - What is the relationship of x and y ?
 - How are x and y associated?



Figure 1: Norway Spruce and Larch Forest in Austrian Alps

Data Source

Image Source: <https://ec.europa.eu/jrc/en/research-topic/forestry/qr-tree-project/norway-spruce>

The data used in this example are derived from the R package *Functions and Datasets for “Forest Analytics with R”*.

According to the documentation, the source of these data are: “von Guttenberg’s Norway spruce (*Picea abies* [L.] Karst) tree measurement data.”



Figure 2: Old Tjikko, a 9,550 Year Old Norway Spruce in Sweden

The documentation goes on to further note that:

“The data are measures from 107 trees. The trees were selected as being of average size from healthy and well stocked stands in the Alps.”

```
. use gutten.dta, clear
```

Variables

site Growth *quality* class of the tree’s habitat. 5 levels.

location Distinguishes tree *location*. 7 levels.

tree An identifier for the tree within location.

age_base The tree age taken at ground level.

For some purposes, it might be best to use a centered age variable, centered at the grand mean of tree age:

```

. egen ageMEAN = mean(age_base)

. generate ageCENTERED = age_base - ageMEAN

```

height Tree height, m.

dbh_cm Tree diameter, cm.

volume Tree volume.

age_bh Tree age taken at 1.3 m.

tree.ID A factor uniquely identifying the tree.

One Continuous Thing At A Time

```

. histogram height, title("Tree Height")
(bin=30, start=1.5, width=1.4)

. graph export myhistogram.png, width(500) replace
(file myhistogram.png written in PNG format)

```

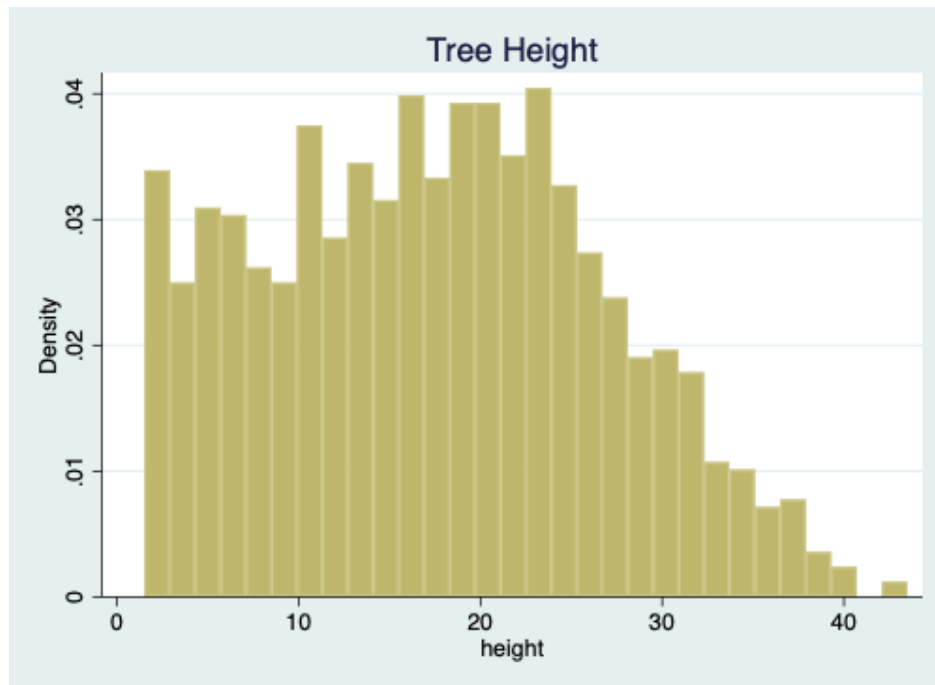


Figure 3: histogram of tree height

One Categorical Thing At A Time

```

. graph bar, over(location) title("Tree Location")

. graph export mybargraph.png, width(500) replace
(file mybargraph.png written in PNG format)

```

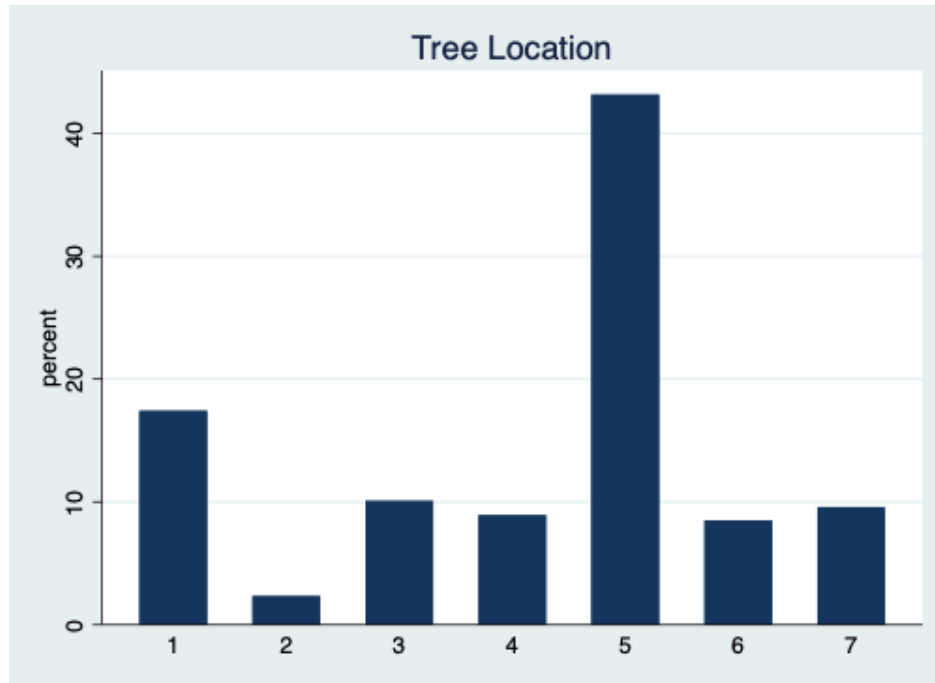


Figure 4: bar graph of tree location

Continuous by Continuous

```
. twoway scatter height age_base, title("Tree Height by Age")
. graph export myscatter.png, width(500) replace
(file myscatter.png written in PNG format)
```

Categorical by Categorical

```
. graph bar, over(site) over(location) title("Tree Site Growth Quality by Location")
. graph export mybargraph2.png, width(500) replace
(file mybargraph2.png written in PNG format)
```

Continuous by Categorical

```
. graph bar height, over(location) title("Tree Height by Location")
. graph export mybargraph3.png, width(500) replace
(file mybargraph3.png written in PNG format)
```

Schemes

Stata *graph schemes* can substantially improve the look of a graph. Built in graph schemes include `sj` and `economist`. `lean2` is a user written scheme that is helpful when preparing graphics for publication. I have written a Stata Michigan graph scheme that can be installed.

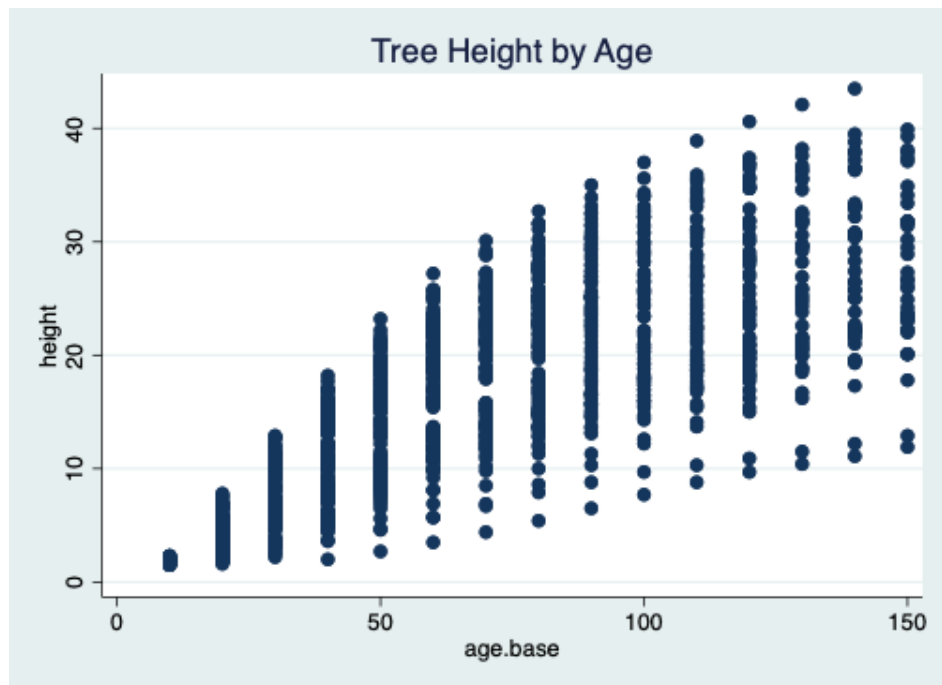


Figure 5: scatterplot of tree height by age

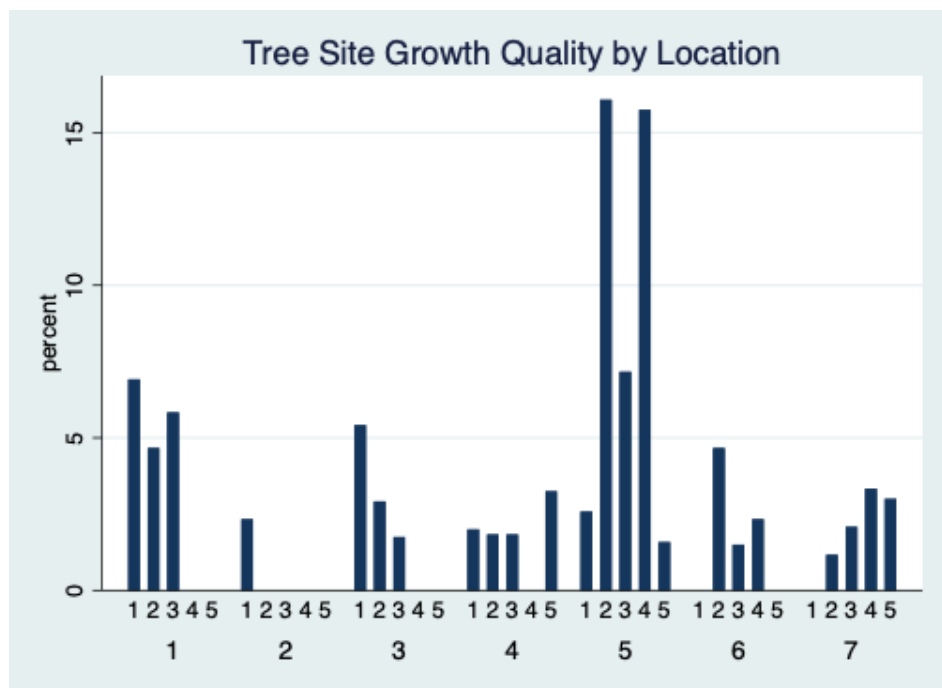


Figure 6: bar graph of tree site by location

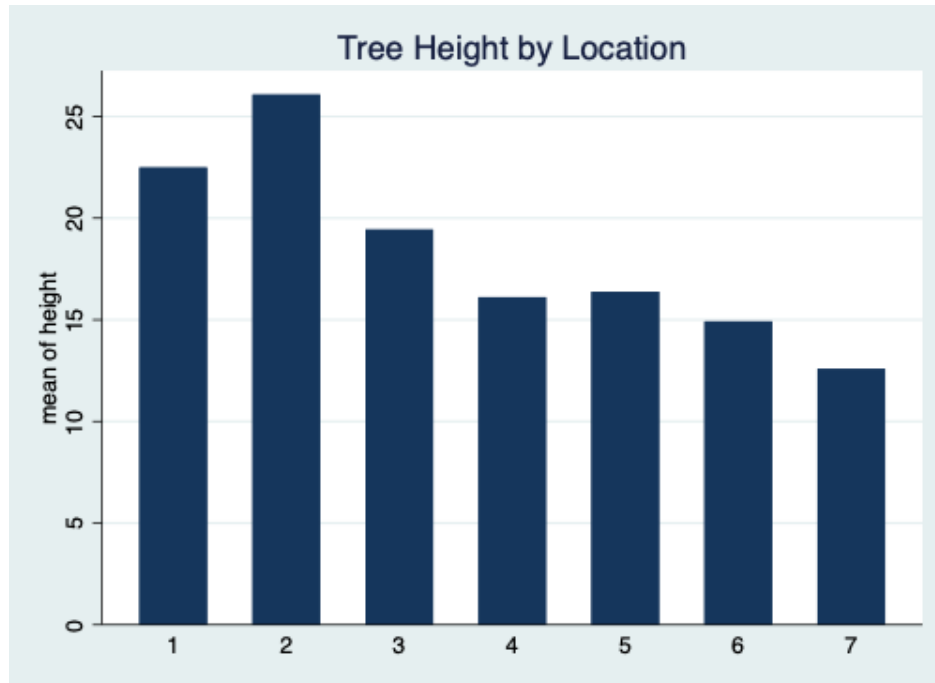


Figure 7: bar graph of mean tree height by location

Continuous by Continuous

```
. twoway scatter height age_base, title("Tree Height by Age") scheme(michigan)

. graph export myscatterM.png, width(500) replace
(file myscatterM.png written in PNG format)
```

Continuous by Categorical

Note that in the graph below, I have used the `asyvars` option to give different colors to the different bars.

```
. graph bar height, over(location) asyvars title("Tree Height by Location") scheme(michig
> an)

. graph export mybarM.png, width(500) replace
(file mybarM.png written in PNG format)
```

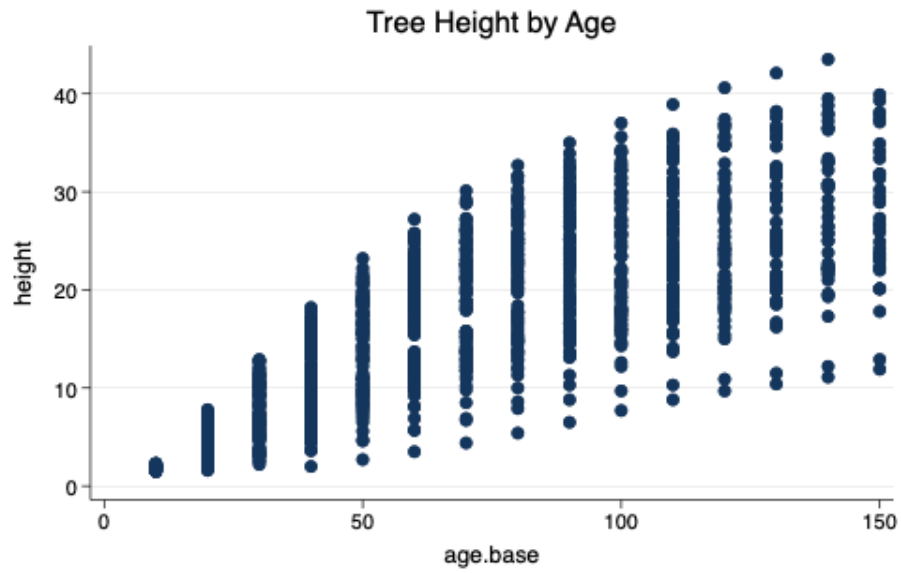


Figure 8: scatterplot of tree height by age

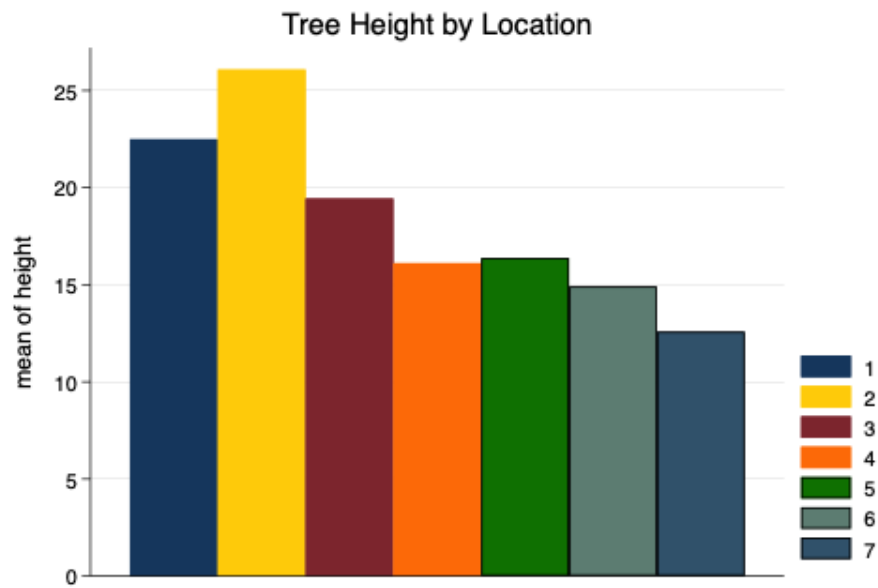


Figure 9: bar graph of mean tree height by location