

# Stata Introduction

Jan Marvin Garbuszus<sup>1</sup>

<https://github.com/JanMarvin/StataFolien>

May 12, 2016

---

<sup>1</sup> I am grateful to Deniz Türker for proofreading and comments on an earlier version of this presentation.

# Literature

- ▶ U. Kohler and F. Kreuter (2012). *Data Analysis Using Stata*. 3rd ed. Lakeway Drive, Texas: Stata Press
- ▶ A. C. Acock (2012). *A Gentle Introduction to Stata*. 3rd ed. Lakeway Drive, Texas: Stata Press

# Introduction

# Introduction Contents

## Datatypes

- Syntaxfiles

- Datafiles

## do-files

- help

- findit

## Working

- Description of data in memory

- Descriptive

# Syntax

We write Syntax.

- ▶ Reproducible
- ▶ Less ...
  - ▶ clicking
  - ▶ repeating
  - ▶ error-prone

# Syntaxfiles

- ▶ Stata syntax is handled in do-files.
- ▶ Starting the editor

`doedit`

- ▶ Fileending of do-files is .do.
- ▶ do-files can be opened and modified with any texteditor.

# secondary data

- ▶ primary data
- ▶ secondary data

We use data. We do not gather data. We will not create or edit data.

# dta-files

- ▶ Data is stored in dta-files.
- ▶ Data filending is .dta.
- ▶ Stata dataobjects are readable only by Stata and some other statistical software e.g. R (R Core Team, 2016) using readstata13 (Garbuszus and Jeworutzki, 2016).



# Importing data

- ▶ Showing the content of a directory

```
dir
```

- ▶ Changing the directory

```
cd "D:/Data/"
```

- ▶ Importing data

```
use <dataname.dta>
```

Avoid directory- or filenames  
containing umlauts or whitespaces.  
Stata does not like those.

# Importing data

To follow the code on this presentation you can use the Stata example dataset `multistage`. Load it on a computer with internet access

```
webuse multistage, clear
```

# do-file 1

- ▶ file header: should contain information like filename, author, creation and/or modification date, a short description of its content and for longer files a short table of contents.

```
/*  
Name: read-soep.do  
Author: Garbuszus  
Date: 2016-02-28  
Content: Importing SOEP-data  
*/
```

# do-file la

Document your code!

- ▶ Teamwork
- ▶ Obliviousness

Nobody likes to study your syntax,  
document it. Document a lot and  
thoroughly, you will forget!

# do-file lb

- ▶ `/*` (begin) und `*/` (end) define a comment. Alternative this can be done at the beginning of a line with `*` or inline with `//`  
  

```
** We write code for Stata version 8 and later  
version 8 // This runs with Stata 8 and might not work with Stata 7
```
- ▶ Commands in comments will be skipped!

# do-file lc

- This is followed by

```
** remove everything from our data memory  
clear all  
** disables more functionality  
set more off
```

clear all erases everything from Stata's memory, securing, that everything in the following do files is run on exactly the data we define and not on some earlier created data artifacts. It's reproducible!

# do-file II

- ▶ Import some testdata

```
use "D:/Data/testdaten.dta"
```

or as shown on frame 10

```
webuse multistage, clear
```

- ▶ Describe

```
describe
```

- ▶ Help for command

```
help describe
```

# help

tabulate **varname1** **varname2** [**if**] [**in**] [**weight**] [, options]

options	Description
<hr/>	
Main	
chi2	report Pearson's chi-squared
exact[(#)]	report Fisher's exact test
...	

regress **depvar** [**indepvars**] [**if**] [**in**] [**weight**] [, options]

options	Description
<hr/>	
Model	
noconstant	suppress constant term
hascons	has user-supplied constant
...	



# findit

And how does one find commands?

```
findit describe
```

findit can be helpfull, but it's  
not guaranteed to find exactly  
what you're looking for.

# do-File III

Different tasks, different do-files.

- ▶ Import data
- ▶ Preparation
- ▶ Descriptive analysis
- ▶ etc.

Not every do-file is required to be rerun all the time. For instance importing and preparation of the data should be a one time task. It's not required to do this for every single crosstable or boxplot. Time is money!

## do-file IV

since do-files may be executed by other do-files nesting some of them in a master file helps with housekeeping

```
** call other do-files  
do "01-Data_Import.do"  
do "02-Data_Preparation.do"  
do "03-Descriptive_Analysis.do"
```

This tightens your syntax and the correct order of do-files is assured.

# Description of data in memory

## Description of data in memory

```
describe  
codebook  
list  
browse
```

Studying describes help page (`help describe`) presents additional options like

```
describe, simple
```

Help files contain information whether or not a command requires a single variable or a list of variables.

# Descriptive Analysis

## Descriptive analysis of variables <sup>2</sup>

```
tabulate var1  
tabulate var2  
tabulate var1 var2  
summarize var1  
summarize var1-var4
```

Studying summarizes help page (`help summarize`) presents additional options like

```
summarize var1, detail
```

---

<sup>2</sup> var1 and var2 are placeholders.

# Variables

# Variables Contents

## Basics

- Generate

- Rename

- Label

## Elaborated

- Generate using Math

- Recode into a new Variable

- Categorize

- Missing Values

- Variable Handling

## Save and Import

- Save

- Import

# Generate

Time and time again generation of variables is required

```
** clone variable
clonevar female = sex

** generate variable
gen gender = .
recode gender . = 1 if sex == 1
recode gender . = 2 if sex == 2

** spot the difference
tab female
tab gender
```



# Rename

Thinking about it we dislike the variable name gender.

```
drop sex // throw this variable away  
** gender is now called sex  
rename gender sex
```

# Moment generate ... = . ?

Missing values in Stata are identified by the dot.

```
clonevar edu = school
```

```
** create a missing
```

```
recode edu 10 = .
```

```
** different missing values (Stata knows .a - .z)
```

```
recode edu 9 = .a
```

```
recode edu 8 = .b
```

```
recode edu 7 = .c
```

Additionally the dot is  $\pm\infty$  and  
is Stata's biggest value.

# A moment please ... if ... == 1 ?

## if-conditions

```
** sex in county 1 and school 1 or 2
tab sex if county == 1 & ///
    (school == 1 | school == 2)
```

The condition is following the if-statement. Different conditions can be linked using & and |. Linked conditions require variable names within each condition!

# and tab ?

tab is short for

```
help tabulate twoway
```

```
tabulate varname1 varname2 [if] [ ...
```

The underlined part shows the minimum letters Stata requires for identification of a command.

# Label

To label a variable, a label has to be defined and afterwards applied

```
** label a variable
label variable sex "Sex"
** define label for values - labsex is the label name
label define labsex 1 "male" 2 "female"
** assign label to variable
label values sex labsex
```

Now labels are assigned

```
tab sex
tab sex, nolabel
```

# Generate II

It is possible to create variables using mathematical expressions

```
// create fantasy weight and height
replace weight = weight / 2.2
replace height = height * 0.004

// bmi for lbs and in
gen bmi = (weight / (height ^ 2))

/* other examples
** householdincome
gen hhinc = inc_female + inc_male
** age in 2016
gen age = 2016 - gebjahr
*/
```

## Generate III

A new variable is created from recoded values. The original variable is not modified.

```
recode school ///
  (10=.) ///
  (1=1 "One") ///
  (2=2 "Two") ///
  (3 5 6 7 8 = 99 "Ninetynine") ///
  (4=3 "Three") ///
  (9=5 "Five"), gen (school_rec)
```

Variable names begin with a character (a-z or A-Z).

# A moment please /// ?

Three following *Slashes* indicate to Stata, a command will continue in the following line.

```
** everything in a single line
display 1 + 1
** now with a (unneeded) linebreak
display 1 + ///
1
```

This way two and more lines can be connected. This makes your syntax more readable. Else check help delimit



# Recoding

Sometimes it is helpful to categorize metric variables

```
** create rounded values
gen bmi_r = round(bmi)
recode bmi_r ///
  (0/18.5=1) ///
  (18.5/25.0=2) ///
  (25.0/30.0=3) ///
  (30.0/60=4)

tab bmi_r
```

0/18 means from 0 to 18. 25/30  
from 25 to 30. Be careful with  
the category range!

# Data preparation

Lets have a look at our fantasy weights data. We replace a few values with different values

```
** weight  
replace weight = -1 if weight < 40  
replace weight = -2 if weight > 140  
  
mean weight
```

Now we decode these values as missing values

```
mvdecode weight, mv (-1=.a\ -2=.b)  
mean weight
```

In the SOEP negative values are typically missing values.

# Mistakes

After hours of work you figure that you have made a mistake and want to modify the just created variable `bmi_r`. What now?

- ▶ Luckily you have not touched the original data. Just run your do-file up to the point where the mistake is located.
- ▶ You already spotted the error and want to rerun the variable creation. Simply remove the variable

```
drop bmi_r
```

# Save Data

Save your do-file into a directory on your computer (e.g. *do-files*) and your data into a different folder(e.g. *data*).

```
save "D:/Data/data/testdata_recoded.dta"
```

If a file with the same name exists, Stata does not want to erase anything by mistake and asks you for the `replace` option

```
save "D:/Data/data/testdata_recoded.dta", replace
```

Replaced data is forever gone.  
Never ever replace your original  
data!

# Import Data

With `use ...`, `clear` the currently imported data set will be replaced

```
use "D:/Data/data/testdata.dta", clear
```

Unsaved modifications to the current data set are ignored and are not saved. Using the `clear` option forces Stata to do something it would otherwise warn you about.

# Graphics

# Graphics Contents

Preface

Scatterplot

Boxplots

Histograms

Dot-Charts

Export of Graphics

# Graphics

The next slides will provide an example of different types of graphics. Detailed examples and syntax for modification of graphic title, titles for the axis or combination of multiple graphics see Kohler and Kreuter (2012, S. Chap. 6) and in-depth with lots of images Mitchell (2012).



## Preface ...

### Warning!

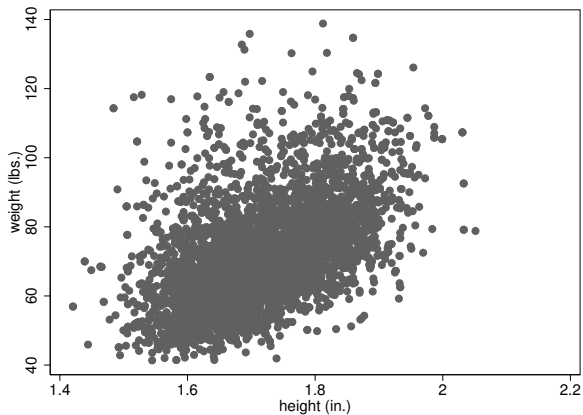
Don't do any pie, pizza, circle or however you name them charts.  
Just don't do it!

# Additionally

```
** otherwise we see colors  
set scheme simono
```

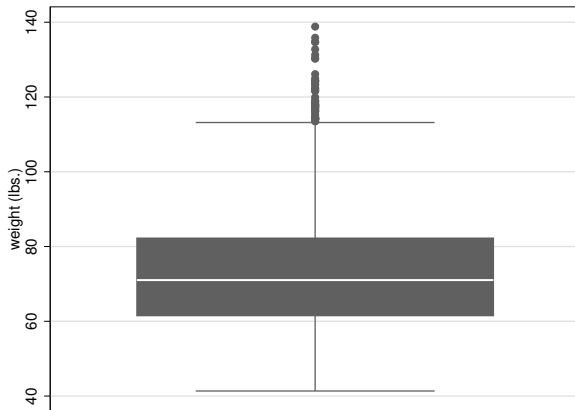
# Scatterplots

```
scatter weight height
```



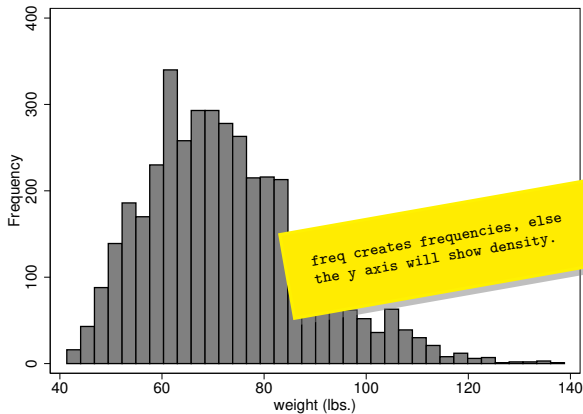
# Boxplots

```
graph box weight
```



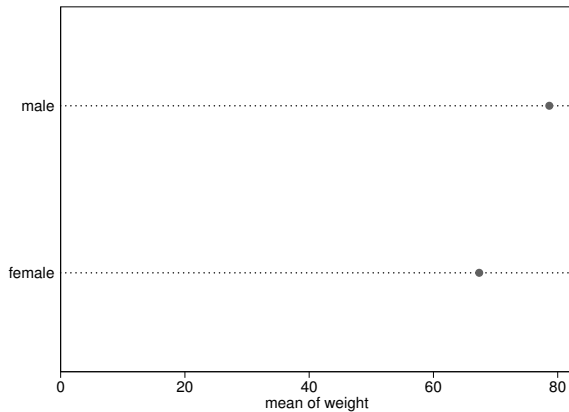
# Histograms

```
hist weight, freq
```



# Dot-Charts

```
graph dot (mean) weight, over(sex)
```



# Export of Graphics

```
graph export "D:/Data/graphics/graph1.pdf", replace
```

- ▶ For export there are many different formats
- ▶ .ps .pes .wmf .emf .pdf .png .tif
- ▶ For usage in MS-Office a good choice is .png
- ▶ For usage in TeX a good choice is either .pdf or .eps

# Improved Datahandling



# Improved Datahandling Contents

by

in

Loops

log-files

Macros

Using

Merge

Operators

by |

## Select subsets

```
summarize weight if sex == 1  
summarize weight if sex == 2
```

## Do it in a single line

```
by edu, sort: summarize hhinc
```

by is not allowed with every command, in doubt contact the help. To use by cases have to be sorted. Alternativ: bysort

by II

It is possible to use more than one variable

```
by sex county, sort: summarize weight
```

in

Select only a few cases

```
** select the first  
list sex in 1  
** select the first ten  
list sex in 1/10
```

Usually single cases are not interesting, but this can be helpful to spot strange cases.

# Loops I

Loops in Stata. Example from Kohler and Kreuter (2012, S. 61ff.)

```
** Generate variables r1 - r10
** Newlist
foreach var of newlist r1-r10 {
    gen `var' = runiform()
}

** Numeric list
foreach num of numlist 1/10 {
    replace r`num' = runiform()
}
```

# Loops II

```
** multiple commands in a loop
foreach var of varlist ybirth income {
    summarize `var', meanonly
    generate `var'_c = `var' - r(mean)
    label variable `var'_c "`var' (centered)"
}

** Forvalues
forvalues num = 1/10 {
    replace r`num' = runiform()
}
```

# log-files

- ▶ Results can be stored in log-files.
- ▶ These contain all commands and their results.
- ▶ Logging begins after

```
log using <dateiname.log>
```

- ▶ log-files can be opened with every texteditor.

## log-files II

- ▶ Define a name and path where the log-file is saved. Usually the log-file name should reflect the do-files name.

```
log using read-soep.log, replace
```

- ▶ At the end of the do-file

```
** commands will not be logged after this  
log close
```

- ▶ Thrown errors stop the do-files execution, therefore a manual closing of the log-file might be required because Stata cannot handle two open log-files

```
log close
```

capture log close atop log using  
ensures that all open log-files  
will be closed.



# Macros

With longer pathes and bigger projects it is helpful to use so called Macros. Macros may be global (it is known everywhere inside your code) or local (it is known only in a specific loop). We use macros for paths

```
** global Macroname Pathname  
global DATA "D:/Data/data/original/"  
global OUT "D:/Data/data/bearbeitet/"
```

Now the macros can be called like this

```
cd "${DATA}"
```

# Macros II

It is especially helpfull handling long or different directorys for input, output, images or log files. Some additional examples

```
use "${DATA}testdata.dta", clear  
save "${OUT}testdata_recoded.dta", replace
```

The syntax readability benefits a lot from this

You will save a lot of time.  
Long paths are written only once,  
which minimizes typos. A correct  
macro will be correct for the  
entire document. Still macros  
won't help you thinking.

## using

Stata can handle only a limited range of variables, so a little housekeeping is advised.<sup>3</sup> Handling the SOEP this limit is not as far away as you might believe.

```
** Import a minimal ppfad
** kepp only hhnr persnr sex gebjahr
use hhnr persnr sex gebjahr using "${DATA}ppfad.dta"
** Sort cases by persnr and gebjahr
sort persnr gebjahr
```

---

<sup>3</sup> Different Stata versions have different limits: IC (the version you're working with) can handle 2,047, SE can handle 32,767 variables. Handling more variables is rather expensive. Prices for 2013: \$189 and \$395 for students. Non academic single user licenses \$1545 and \$2090.

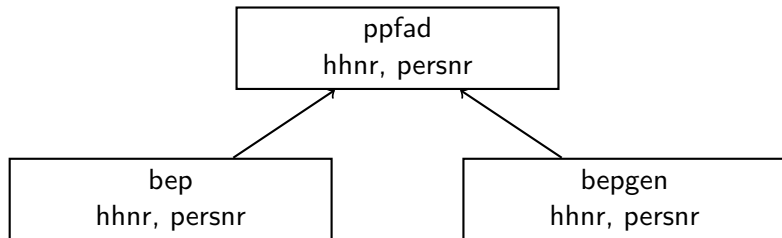
## merge

In the Sozio-oekonomischen Panel (SOEP) (s. Wagner et al., 2007) each years data is split into different smaller data files. This has historical and maybe even practical reasons.

Starting from a file different of these data files share, we can combine the smaller files to generate a complete data set. This combination is called *merging*.

## merge II

As example: We merge **ppfad** with **bep** and **bepgen**. **ppfad** is the key file for the person or p-context. **ppfad** is the data file every p-file shares. To make sure that we combine our data correctly for the p-context, we seek for matches in household-id (hhnr) and person-id (persnr).



# merge III

## Running merge

```
** merge ppfad with 2014 bep  
merge 1:1 persnr using "${DATA}bep.dta"
```

Result	# of obs.	
not matched	25,834	
from master	25,834	( <code>_merge==1</code> )
from using	0	( <code>_merge==2</code> )
matched	10,471	( <code>_merge==3</code> )

This means

- ▶ 25,834 cases of ppfad were not found in bep.
- ▶ 0 cases of bep were not found in ppfad
- ▶ 10,471 cases of bep were found in ppfad

## merge IV

A new variable `_merge` is created containing the `_merge==` values. For now we remove all cases where we could not find a match from using which is `_merge==2`. These are potentially problematic as we have no information about these persons sex or year of birth. Luckily at our current stage these are zero cases. Afterwards we remove `_merge`.

```
** Drop if using does not match master
drop if _merge==2
** drop _merge
drop _merge
```

The coice which merge is dropped  
is depending on the analysis.

# merge V

Example: merge with the generated variable highest educational achievement `bepsbil` from data-file **bepgen**.

```
merge 1:1 hhnr persnr using "${DATA}bepgen.dta", keepusing(bepsbil)
drop if _merge==2
drop _merge
```

bepgen contains 65 variables, but  
all we need is the educational  
achievement.



# Operators

Stata lists the following operators, which can be used with variables.  
Using `var1^2` every value of `var1` will be squared.

```
== // equals
!= // not equal ~=
& // and
| // or
! // not
< // smaller
> // greater
<= // smaller equal
>= // greater equal
+ // plus
- // minus
/ // divide
* // multiply
^ // power
```

# Functions

Additionally Stata contains a range of built in functions (help functions). A selection:

```
abs() // absolute value/modulus |-2| == 2
max() // greatest value
min() // smallest value
exp() // e-function
ln() // logarithm
round() // round
sin() // Sinus
sqrt() // square root
runiform() // random number
```

# Statistics

# Statistics Contents

Statistic

Tables

Cross tabs

Chi, V and Phi

Inference

# Statistic I

## Arithmetic mean and some additional statistics

```
mean height  
summarize height  
summarize height, detail
```

mean returns the standard  
error:  $s/\sqrt{n}$ . summarize  
returns the standard deviation:  
 $\sqrt{\sum (x - \bar{x})^2 / n}$ .

# Statistic II

Minimum, maximum, arithmetic mean, median, number of observations and quartils.

```
tabstat height, statistic(min max mean median p50)  
tabstat height, statistic(min max range mean count q) by(county)
```

$$\text{range} = \text{max} - \text{min}$$

# Statistics III

Standard deviation, standard error, variance und interquartil range.

```
tabstat height, statistic(sd sem var q iqr)
```

Skewness and kurtosis

```
tabstat height, statistics(skewness kurtosis)
```

$$\text{iqr} = q3 - q1$$

# Tables

We already had `summarize` and `tab`, let us combine them

```
tab county, summarize(height)
```



# Cross tabs

Create a cross tab with

```
tab school county
```

And some more tables

```
** tabs of each variable  
tab1 sex county school  
** cross tabs of ab bc ac  
tab2 sex county school
```

# Cross tabs II

Let us recapitulate Statistics I, we will do some tests <sup>4</sup>

```
** chi^2
tab sex county, chi
** cramers v
tab sex county, V

** manual
tab sex county, exp col row
** phi and for 2x2 tabs V
di (1013*1002 - 925*1131) / (2144*1927*1938*2133)^(1/2)
** chi^2
di 4071 *(-.00753735)^2
** usual way to compute V
di (.23128021 / 4071 * (2 - 1) )^(1/2) // V
```

---

<sup>4</sup> or see e.g. Agresti and Finlay (2009).

# Numlabel

You can add numeric values to your variable labels

```
** numeric labels on for all variables  
numlabel _all, add  
** and off  
numlabel _all, remove
```

# T-Test

And finally let us have a look at some test-statistics

```
** compare  
tab sex, sum(bmi)  
ttest bmi, by(sex)
```

And much more to come ...

## Regression

`regress`  
`logit`  
`mlog`

To be continued.

# Index and Literature

browse, 20  
 by, 50, 51  
     bysort, 50  
  
 Categorize, 33  
 clear all, 14  
 comment, 11, 12  
     commenting, 13  
  
 Describe, 20  
     codebook, 20  
     describe, 15, 20  
     detail, 69  
     summarize, 21, 69  
 directory  
     cd, 9  
     dir, 9  
 display, 32  
 do, 19  
 do-file, 6, 11, 18, 19  
     doedit, 6  
 drop, 35  
 dta-file, 8  
  
 file header, 11  
 findit, 17  
 Functions, 66  
  
 Generate, 24, 30  
     clonevar, 24, 33  
     gen, 24, 30, 31  
     generate, 24

Graphics  
     graph dot over(), 46  
     blackandwhite, 42  
     Boxplots, 44  
     Dot-Charts, 46  
     export, 47  
     graph box, 44  
     graph dot, 46  
     hist, 45  
     hist freq, 45  
     Histograms, 45  
     scatter, 43  
     Scatterplot, 43  
     set scheme, 42  
  
 Help, 16  
     help, 15, 16, 20, 21, 28  
  
 if, 27  
     and, 27  
     or, 27  
 Import, 9, 15, 37  
     clear, 37  
     use, 15, 37  
     use using, 59  
 Interquartil range, 71  
     iqr, 71  
  
 Kurtosis, 71  
     kurtosis, 71  
  
 Label

label, 29  
 label define, 29  
 label values, 29  
 label variable, 29  
 numlabel, 75  
  
 list, 20, 52  
 log, 55, 56  
     close, 56  
     log-files, 55  
     using, 56  
 Loops, 53, 54  
     foreach, 53, 54  
     forvalues, 54  
  
 Macros, 57, 58  
     global, 57  
 Maximum, 70  
     max, 70  
 Mean, 69, 70  
     arithmetic mean, 70  
     mean, 69, 70  
     Median, 70  
 Merge  
     \_\_merge, 63  
     merge, 60–64  
     merge, keepusing, 64  
 Minimum, 70  
     min, 70  
 Missing Values, 34  
     ., 26  
     mv(), 34  
     mvdecode, 34

Mistakes, 35  
  
 newline, 32  
     ///, 32  
     delimit, 32  
  
 Operators, 65  
  
 Quantile, 70  
     q, 70, 71  
  
 Range, 70  
     range, 70  
 Recode  
     recode, 24, 26, 31, 33  
 Regression  
     logit, 77  
     mlog, 77  
     regress, 77  
 Rename, 25  
     rename, 25  
  
 Save, 36  
     replace, 36  
     save, 36  
 Search  
     findit, 17  
     set more off, 14  
 Skewness, 71  
     skewness, 71  
 Sort  
     sort, 59

Standard deviation, 71  
     sd, 71  
 Standard error, 71  
     sem, 71  
 Statistic, 69  
  
 T-Test, 76  
     ttest, 76  
 Table, 72  
      $\chi^2$ , 74  
     chi, 74  
     chi-square, 74  
     Cramers V, 74  
     cross tab, 73, 74  
     nolabel, 29  
     tab, 24, 27, 28, 72–74  
     tab col, 74  
     tab exp, 74  
     tab row, 74  
     tab, summarize(), 72  
     tab1, 73  
     tab2, 73  
     tabulate, 21  
 tabstat, 70, 71  
  
 Umlauts, 9  
     use, 9  
  
 Variance, 71  
     var, 71  
  
 Whitespaces, 9



# Literature I



Acock, A. C. (2012). *A Gentle Introduction to Stata*. 3rd ed. Lakeway Drive, Texas: Stata Press.



Agresti, A. and B. Finlay (2009). *Statistical Methods for the Social Science*. 4. Aufl. London: Pearson.



Garbuszus, J. M. and S. Jeworutzki (2016). *readstata13: Import Stata Data Files*. R package version 0.8.2.



Kohler, U. and F. Kreuter (2012). *Data Analysis Using Stata*. 3rd ed. Lakeway Drive, Texas: Stata Press.



Mitchell, M. N. (2012). *A Visual Guide to Stata Graphics*. 3. Aufl. Lakeway Drive, Texas: Stata Press.



R Core Team (2016). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing. Vienna, Austria.

# Literature II



Wagner, G. G., J. R. Frick, and J. Schupp (2007). *The German Socio-Economic Panel Study (SOEP): Scope, Evolution and Enhancements*. SOEPpapers on Multidisciplinary Panel Data Research 1. DIW Berlin, The German Socio-Economic Panel (SOEP).