# Basics of STATA Software

Introduction to the STATA interface and its functions. Entering and importing data into STATA from economic databases and various data files. Basic data management: opening, viewing, and editing variables. Generating and labeling variables, creating data subsets

Tomáš Oleš

Department of Economic Policy
Faculty of Economics and Finance

January 29, 2025

- Get accustomed to the STATA interface.
- Enter and import data into STATA.
- Get comfortable with using the STATA command language.
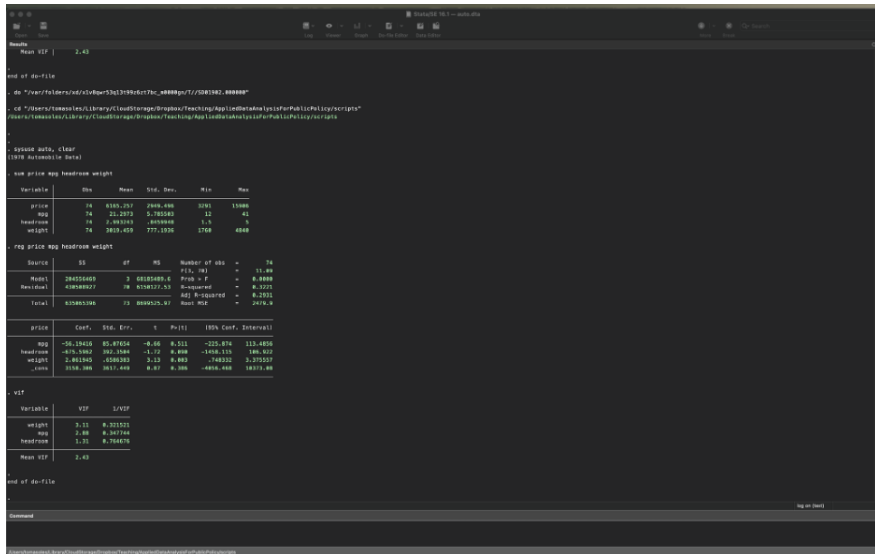- Learn common data management commands in STATA.

# Why Research Needs to be Reproducible?

- "The natural (social) scientist is concerned with a particular kind of phenomena ... he has to confine himself to that which is reproducible ... I do not claim that the reproducible by itself is more important than the unique. But I do claim that the unique exceeds the treatment by scientific method. Indeed it is the aim of this method to find and test natural laws." — Wolfgang Pauli

- "Kohn's Second Law: An experiment is reproducible until another laboratory tries to repeat it." — Alexander Kohn

# A Key to Reproducibility is to Learn to Code < 3

# Create Your First Code and Run It

```
display("Hello Word!")
```

# Using log-file

```
*Load data and create a log-file
log using "~/AppliedDataAnalysisForPublicPolicy/log/my_first_log.log", ///
replace
sum price mpg headroom weight turn displacement
log close
```

## Entering data

```
/*inputting new data for the replace command*/
clear
input str10 exammark mark
"" 93
"" 92
"" 83
"" 76
end
```

## Importing Data

**Default Data Format:** Stata primarily uses `.dta` files but supports multiple formats, including:

- `.xls`, `.xlsx` (Excel)
- `.csv` (Comma-separated values)
- `.shp`, `.dbf` (Shapefiles & dBase)
- ...

**What if my data is in an unsupported format?** Convert it to a Stata-compatible format (e.g., `.csv`) using Python or another tool.

**Limitation:** Stata can handle only one dataset (table) at a time.

**Tasks:**

- Import `datasets/highest-points-by-state.csv` into Stata.
- Import `datasets/workout1.dta`.

## Solution: Importing Data into Stata

**Importing a CSV File:**

```
import delimited "datasets/highest-points-by-state.csv", clear
```

**Importing a Stata (.dta) File:**

```
use "datasets/workout1.dta", clear
```

**Converting Unsupported Formats:** Use Python (e.g., Pandas) to convert files:

```
import pandas as pd
df = pd.read_json("data.json")
df.to_csv("data.csv", index=False)
```

## Examining the Data

```
use datasets/workout1,clear
describe v03 v04
describe using workout1
codebook v03
browse v01 v02 v03 v04, nolabel
edit v01 v02 v03 v04, nolabel
list v01 v02 v03 v04, nolabel
misstable sum v01 v02 v03 v04
```

## Making Changes to Variables

```
/*recode examples, not based on a dataset*/
recode var1 (-999=.) or recode var1 -999=.
recode _all (-999=.) or recode * (-999=.)
mvdecode _all, mv(-999) or mvdecode *, mv(-999)
mvencode _all, mv(-999) or mvencode *, mv(-999)

/*back to using the workout1 dataset*/
use datasets/workout1,clear
recode v04 (1/3=1) (4/6=2) (7/9=3)
```

## Replacing the Data

```
/*inputting new data for the replace command*/
clear
input str10 exammark mark
"" 93
"" 92
"" 83
"" 76
end
replace exammark="very good" if mark>90
replace exammark="good" if mark<90
/*back to using the workout1 dataset*/
use datasets/workout1,clear
rename v03 Education
rename Education, lower
rename *, upper
```

# Mathematical Operators in Stata

| Arithmetic | | Logical | | Relational | |
|---|---|---|---|---|---|
| $+$ | addition | & | and | $>$ | greater than |
| $-$ | subtraction | $\mid$ | or | $<$ | less than |
| $*$ | multiplication | $\neg$ | not | $\geq$ | greater or equal |
| $/$ | division | $\sim$ | not | $\leq$ | less or equal |
| $\hat{\ }$ | power | | | $==$ | equal |
| $-$ | negation | | | $\neq$ | not equal |
| $+$ | string concatenation | | | $\sim=$ | not equal |

Table: The three types of mathematical operators used in Stata

# Generating variables: workout1.dta

```
/*hypothetical examples on gen command*/

gen age2=age^2  //Age squared
gen id=_n //numbers observations
gen loghours=log(hours) //Log of hours
gen pdollar=price/6  //Price (in Norwegian currency) to dollars
gen agecar=2015-year //The age of car in 2015
```

```
*back to using the workout1 dataset*/
use datasets/workout1,clear
/*generate a new variable using gen and recode*/
recode v04 (1/3=1) (4/6=2) (7/9=3), gen(inccat)
tab inccat

/*generate a new variable using gen and replace*/
gen inccat2=.
replace inccat2=1 if (v04<=3)
replace inccat2=2 if (v04>=4) & (v04<=6)
replace inccat2=3 if (v04>=7) & (v04<.)
tab inccat2

/*hypothetical example showing labelling values of several variables */
label define lablikert 1"disagree" 6"agree"
label values var1-var5 lablikert
```

## Appending Data

```
/*inputting new data manually*/
clear
input id data var1 var2
1 1 3 2
2 1 4 3
3 1 5 1
end
save dataset1,replace
clear
input id data var1 var2
4 2 3 1
5 2 5 3
6 2 5 4
end
save dataset2,replace
clear
```

```
/*appending data*/
append using dataset1 dataset2,gen(dataset3)
save dataset3
list,sep(0)
```

# Merging Data

```
/*inputting new data manually*/
clear
input id v1_14 v2_14
1 3 5
2 4 5
3 2 3
4 1 2
5 1 2
end
save data14,replace
```

```
clear
input id v1_15 v2_15
1 4 5
2 5 5
3 3 4
4 2 3
5 2 3
end
save data15,replace
clear
```

```
/*merging data*/
use data14,clear
merge 1:1 id using data15
save data1415,replace
drop _merge // drops this variable
list, sep(0)
```

# Reshaping Data

```
/*reshape from wide to long*/
use data1415,clear //we use the data from above since it is in
                    a wide format
drop _merge
list
reshape long v1_ v2_ , i(id) j(year)
list
```