# Introduction to the Stata Language, Part 2

Mark Lunt

Centre for Epidemiology Versus Arthritis
University of Manchester

**CENTRE FOR
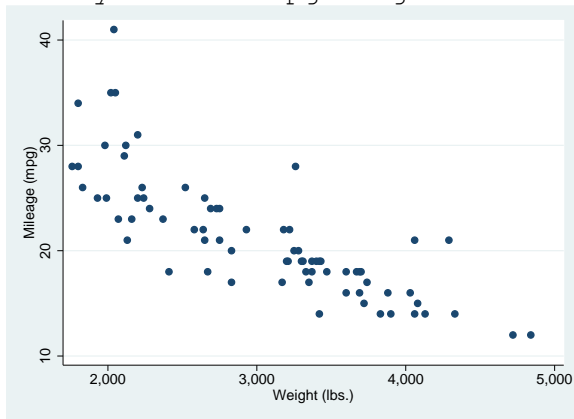EPIDEMIOLOGY
VERSUS
ARTHRITIS**

20/12/2022

Graphics
Summarizing Data
More Stata Syntax
Looping
Reshaping
Other Useful Commands

## Summary

- Graphics
- Summarizing Data
- More Stata Syntax
- Looping
- Reshaping Data
- Other Useful Commands

**CENTRE FOR
EPIDEMIOLOGY
VERSUS
ARTHRITIS**

Graphics
Summarizing Data
More Stata Syntax
Looping
Reshaping
Other Useful Commands

Scatter Plots
Labelling
Overlaying Plots
Schemes
Saving & Exporting
Other Graph Types

## Graphics

- Scatter plots
- Labelling
- Overlaying plots
- Schemes
- Saving & Exporting

**CENTRE FOR**
**EPIDEMIOLOGY**
**VERSUS**
**ARTHRITIS**

Graphics    Scatter Plots
Summarizing Data    Labelling
More Stata Syntax    Overlaying Plots
Looping    Schemes
Reshaping    Saving & Exporting
Other Useful Commands    Other Graph Types

## Scatter Plots

```
twoway scatter mpg weight
```

Graphics
Summarizing Data
More Stata Syntax
Looping
Reshaping
Other Useful Commands

Scatter Plots
Labelling
Overlaying Plots
Schemes
Saving & Exporting
Other Graph Types

## Labelling

Titles `title(), subtitle(), note(), caption()`

Axis names `xtitle, ytitle`

Tick marks `xlabel, ylabel`

CENTRE FOR
EPIDEMIOLOGY
VERSUS
ARTHRITIS

Graphics
Summarizing Data
More Stata Syntax
Looping
Reshaping
Other Useful Commands

Scatter Plots
Labelling
Overlaying Plots
Schemes
Saving & Exporting
Other Graph Types

## Overlaying Graphs

```
twoway (scatter mpg weight) (scatter length
weight, yaxis(2))
```



CENTRE FOR
EPIDEMIOLOGY
VERSUS
ARTHRITIS

| Graphics | Scatter Plots |
| Summarizing Data | Labelling |
| More Stata Syntax | Overlaying Plots |
| Looping | Schemes |
| Reshaping | Saving & Exporting |
| Other Useful Commands | Other Graph Types |

```
twoway lfitci mpg weight || scatter mpg weight
```



CENTRE FOR
EPIDEMIOLOGY
VERSUS
ARTHRITIS

Graphics
Summarizing Data
More Stata Syntax
Looping
Reshaping
Other Useful Commands

Scatter Plots
Labelling
Overlaying Plots
Schemes
Saving & Exporting
Other Graph Types

## Schemes

- Can change appearance of graph:
    - Line thickness
    - Colour or B/W
    - Text size
- Ideal for journal is not ideal for slides
- 11 Schemes provided with stata
- Can write your own by modifying existing ones
- User-written ones also available
- `set scheme` *scheme_name*, `[permanently]`
- Option `scheme(`*scheme_name*`)`

CENTRE FOR
EPIDEMIOLOGY
VERSUS
ARTHRITIS

| Graphics | Scatter Plots |
| Summarizing Data | Labelling |
| More Stata Syntax | Overlaying Plots |
| Looping | Schemes |
| Reshaping | Saving & Exporting |
| Other Useful Commands | Other Graph Types |

Graphics

Summarizing Data

More Stata Syntax

Looping
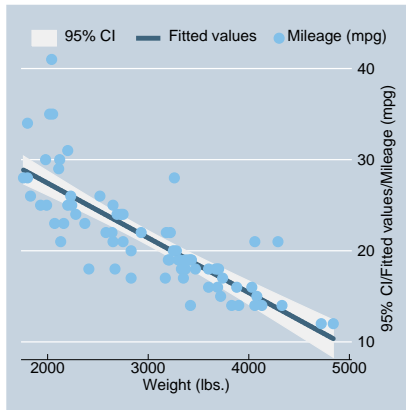
Reshaping

Other Useful Commands

Scatter Plots

Labelling

Overlaying Plots

Schemes

**Saving & Exporting**

Other Graph Types

## Saving Graphs

- Save graphs in stata format with `graph save`
- Save graphs in other formats with `graph export`
- Format used defined by
  - Filename suffix
  - Option `as()` to `graph export`
- Use `help graph export` to find out formats available to you (depends on version and OS).

**CENTRE FOR
EPIDEMIOLOGY
VERSUS
ARTHRITIS**

Graphics
Summarizing Data
More Stata Syntax
Looping
Reshaping
Other Useful Commands

Scatter Plots
Labelling
Overlaying Plots
Schemes
**Saving & Exporting**
Other Graph Types

## Naming graphs

- By default, every graph called "Graph"
- Can store files in memory by renaming:
  - Option `name()` to `graph` commands
  - `graph rename Graph` *newname*
- Recall with `graph display` *name*
- Can display multiple graphs as the same time if they have different names

Graphics

Summarizing Data

More Stata Syntax

Looping

Reshaping

Other Useful Commands

Scatter Plots

Labelling

Overlaying Plots

Schemes

Saving & Exporting

Other Graph Types

## Other Graph Types

graph bar Bar charts

graph box Box and whisker plots

graph matrix Given *n* variables, creates an *n* by *n* matrix of scatterplots, plotting every variable against every other variable.

twoway histogram Histograms

twoway rcap Given two *y*-values for each *x*-value, plots a line between the two *y*-values, with "caps" at each end. Useful for showing confidence intervals if overlaid.

CENTRE FOR
EPIDEMIOLOGY
VERSUS
ARTHRITIS

Graphics
Summarizing Data
More Stata Syntax
Looping
Reshaping
Other Useful Commands

Scatter Plots
Labelling
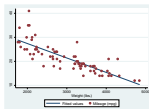Overlaying Plots
Schemes
Saving & Exporting
Other Graph Types

## Other Graph Types
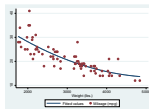
twoway lfit[ci] Linear regression fit to a scatter plot

twoway qfit[ci] Quadratic regression fit to a scatter plot

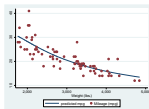twoway fpfit[ci] Fractional polynomial fit to a scatter plot

twoway lowess Nonparametric smoothed fit to a scatter plot



lfit



qfit



fpfit



lowess

CENTRE FOR
EPIDEMIOLOGY
VERSUS
ARTHRITIS

Graphics     Scatter Plots
Summarizing Data     Labelling
More Stata Syntax     Overlaying Plots
Looping     Schemes
Reshaping     Saving & Exporting
Other Useful Commands     Other Graph Types

## Kernel Density

## Summarizing Data

- describe
- codebook
- summarize
- tabulate

Graphics
**Summarizing Data**
More Stata Syntax
Looping
Reshaping
Other Useful Commands

**Describe**
Codebook
Summarize
Tabulate

## describe

- `describe [varlist]`
- Number of observations and variables
- For each variable
  - Name
  - Type
  - Format
  - Labels

**CENTRE FOR**
**EPIDEMIOLOGY**
**VERSUS**
**ARTHRITIS**

Graphics
**Summarizing Data**
More Stata Syntax
Looping
Reshaping
Other Useful Commands

Describe
**Codebook**
Summarize
Tabulate

## codebook

- More detail on each variable:
  - All variables: type, range, unique values, missing values, units
  - Continuous vars: mean, SD, percentiles
  - Categorical vars: frequency table / sample values

Graphics
**Summarizing Data**
More Stata Syntax
Looping
Reshaping
Other Useful Commands

Describe
Codebook
**Summarize**
Tabulate

## summarize

```
summarize [varlist]
```

- Gives mean, SD, min, max, non-missing values
- Option `detail` gives fuller summary

```
summarize  price mpg headroom trunk

    Variable |     Obs       Mean    Std. Dev.      Min       Max
-------------+-------------------------------------------------------
       price |      74    6165.257    2949.496      3291      15906
         mpg |      74     21.2973    5.785503        12         41
    headroom |      74    2.993243    .8459948       1.5          5
       trunk |      74    13.75676    4.277404         5         23
```

CENTRE FOR
EPIDEMIOLOGY
VERSUS
ARTHRITIS

Graphics
**Summarizing Data**
More Stata Syntax
Looping
Reshaping
Other Useful Commands

Describe
Codebook
Summarize
**Tabulate**

## tabulate

- `tabulate variable` gives a frequency table
- `tabulate var1 var2` give a cross-tabulation
- Option `ro` and `co` give row and column percentages respectively
- Option `chi2` gives $\chi^2$-test.

CENTRE FOR
EPIDEMIOLOGY
VERSUS
ARTHRITIS

Graphics
Summarizing Data
More Stata Syntax
Looping
Reshaping
Other Useful Commands

## More Stata Syntax

```
[by varlist]:  command varlist [if
expression][, options]
```

- `by` repeats an analysis for each subgroup
- `if` selects a single subgroup to analyse.

Graphics
Summarizing Data
More Stata Syntax
Looping
Reshaping
Other Useful Commands

## Logical Operators

| Operator | Meaning |
|----------|---------|
| & | and |
| \| | or |
| == | equal |
| ~=, != | not equal |
| < | less than |
| <= | less than or equal |
| > | greater than |
| >= | greater than or equal |

Graphics
Summarizing Data
More Stata Syntax
Looping
Reshaping
Other Useful Commands

## Missing Values

- Missing values are bigger than any "real" value
- Using variables in logical expressions is dangerous if missing values exist
- E.g. (price > 15000) is true if price is missing.
- `gen hi_price = price > 15000 if price < .`
- Be very careful when categorising continuous variables.

Graphics
Summarizing Data
More Stata Syntax
Looping
Reshaping
Other Useful Commands

## The `by` *varlist* clause

- Produces results for each subgroup defined by *varlist* separately
- Data needs to be sorted for `by` to work
- Command `bysort` will do it for you
- Can replace a lot of `if` clauses
- Complex expression can only be used with `if`
- Does not work with every command

**CENTRE FOR
EPIDEMIOLOGY
VERSUS
ARTHRITIS**

Graphics
Summarizing Data
More Stata Syntax
Looping
Reshaping
Other Useful Commands

## Subscripting

- Square brackets ([]) after a variable name used pick out an observation by its number
- `weight[7]` means the weight of the seventh observation
- `_n` means the number of the current observation
- `_N` means the number of observations in the data (or `by` group)

CENTRE FOR
EPIDEMIOLOGY
VERSUS
ARTHRITIS

Graphics
Summarizing Data
**More Stata Syntax**
Looping
Reshaping
Other Useful Commands

## Lagged Variables

- `varname[_n - 1]` means the value of the variable `varname` in the previous observation
- `bysort idno (fupno): replace haq = haq[_n - 1] if haq == .`
- `bysort idno (fupno): gen diff = haq - haq[_n-1]`

CENTRE FOR
EPIDEMIOLOGY
VERSUS
ARTHRITIS

Graphics
Summarizing Data
More Stata Syntax
**Looping**
Reshaping
Other Useful Commands

## Looping

```
foreach macname in list {
    list of stata commands
}
```

- Opening { must be on first line
- Command(s) must start on next line
- Final } must have its own line

CENTRE FOR
EPIDEMIOLOGY
VERSUS
ARTHRITIS

Graphics
Summarizing Data
More Stata Syntax
**Looping**
Reshaping
Other Useful Commands

## Other forms of `foreach`

- `foreach var of varlist ...`
- `foreach var of newlist ...`
- `foreach num of numlist ...`

Graphics
Summarizing Data
More Stata Syntax
**Looping**
Reshaping
Other Useful Commands

## Examples of `foreach`

```
foreach visit in 1 2 {
    summarize bp if when == `visit'
} label define yesno 0 "No" 1 "Yes"


foreach x of varlist *_pain {
    label values `x' yesno
}
```

Graphics
Summarizing Data
More Stata Syntax
Looping
**Reshaping**
Other Useful Commands

## Reshaping Data

- Long to wide: very easy
- Wide to long: slightly trickier
- Long form more efficient for storage: only need space for followups that exist
- Long form also normally best for analysis

Graphics
Summarizing Data
More Stata Syntax
Looping
**Reshaping**
Other Useful Commands

## Long Form

| ID | Gender | Anniversary | Score |
|--------|--------|-------------|-------|
| 900108 | 1 | 1 | 7 |
| 900108 | 1 | 2 | 15 |
| 900108 | 1 | 5 | 19 |
| 900113 | 2 | 1 | 0 |
| 900113 | 2 | 2 | 18 |
| 900114 | 1 | 1 | 0 |
| 900114 | 1 | 2 | 0 |

**CENTRE FOR
EPIDEMIOLOGY
VERSUS
ARTHRITIS**

Graphics
Summarizing Data
More Stata Syntax
Looping
**Reshaping**
Other Useful Commands

## Long to wide

Need to specify:

- Unique identifier which shows which observations belong together: `id`
- Which "repeat" a given observation corresponds to: `anniversary`
- Which variables change between visits: `score`

```
reshape wide score, i(id) j(anniversary)
```

Graphics
Summarizing Data
More Stata Syntax
Looping
**Reshaping**
Other Useful Commands

## Wide Form

| ID | Gender | Score1 | Score2 | Score5 |
|--------|--------|--------|--------|--------|
| 900108 | 1 | 7 | 15 | 19 |
| 900113 | 2 | 0 | 18 | . |
| 900114 | 1 | 0 | 0 | . |

Graphics
Summarizing Data
More Stata Syntax
Looping
**Reshaping**
Other Useful Commands

## Wide to long

Need to specify:

- Unique identifier which shows which observations belong together: id
- The name of a new variable to contain "repeat" info: anniversary
- Which variables are in wide form: score
- If suffixes are strings, need to use the string option.

```
reshape long score, i(id) j(anniversary)
```

CENTRE FOR
EPIDEMIOLOGY
**VERSUS
ARTHRITIS**

Graphics
Summarizing Data
More Stata Syntax
Looping
Reshaping
Other Useful Commands

## Other Useful Commands

| | |
|---|---|
| display | Make things appear in the results window. |
| | Can be used as a calculator |
| expand | Produce multiple copies of each observation |
| cmdlog | Make a do-file of all the commands you are entering. |

CENTRE FOR
EPIDEMIOLOGY
VERSUS
ARTHRITIS

Graphics
Summarizing Data
More Stata Syntax
Looping
Reshaping
Other Useful Commands

## Expand

| Exposed | Cases | Controls |
|---------|-------|----------|
| No      | 20    | 40       |
| Yes     | 30    | 10       |

| exposed | case | frequency |
|---------|------|-----------|
| 0       | 0    | 40        |
| 0       | 1    | 20        |
| 1       | 0    | 10        |
| 1       | 1    | 30        |