# Stata Level 1: Fundamentals of Data Analysis

JUAN SEBASTIAN CUERVO

# Goal 1 - Why use Stata?

JUAN SEBASTIAN CUERVO
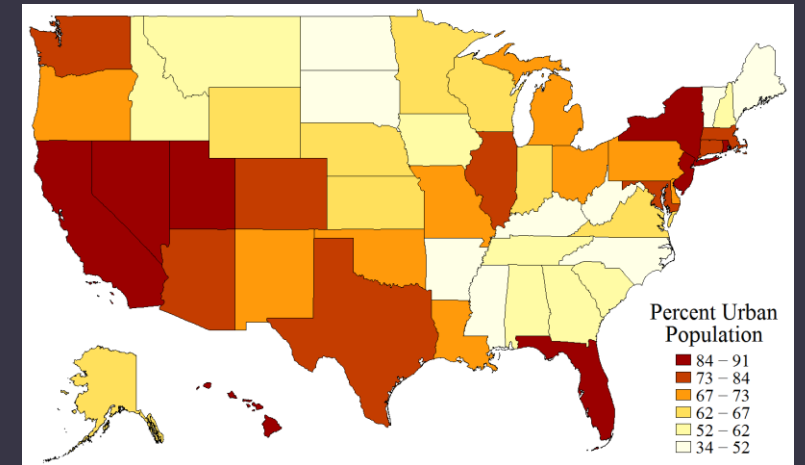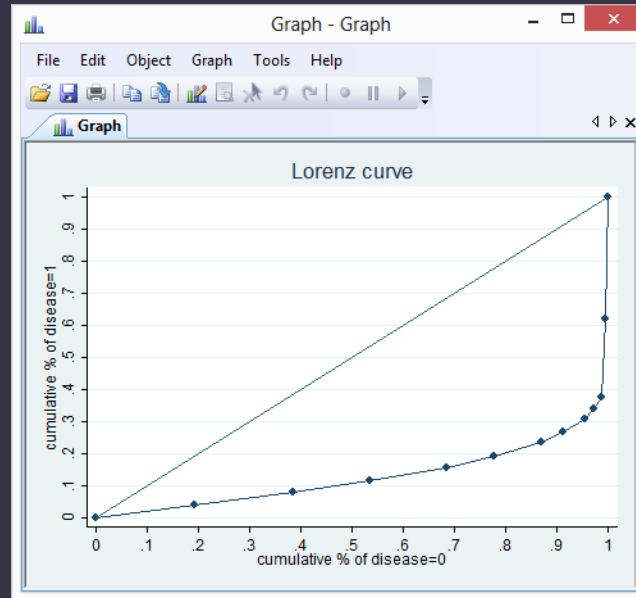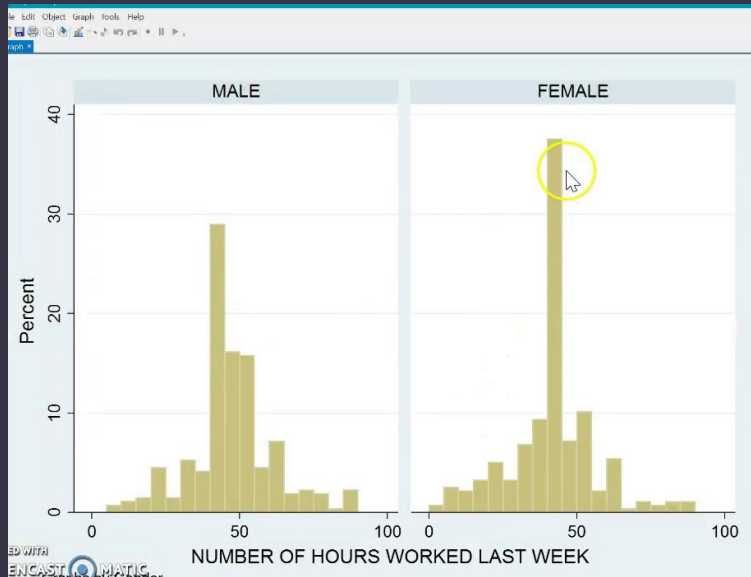
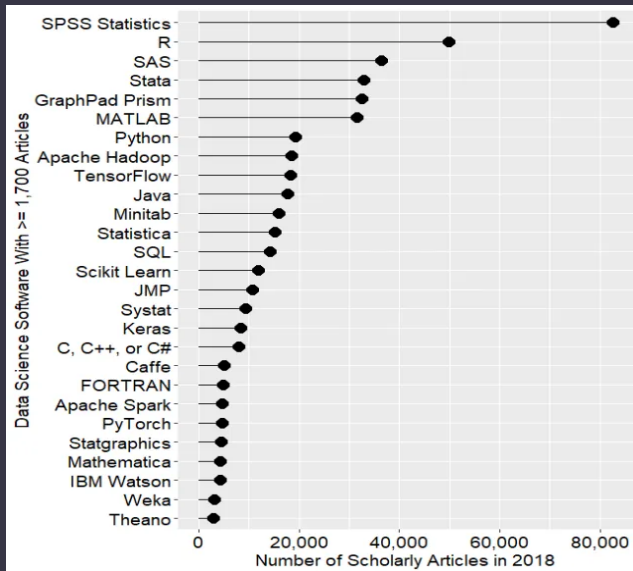# Goal 1: Why use Stata? - Level 1

Organize large quantities of information

# Goal 1: Why use Stata? -Level 2
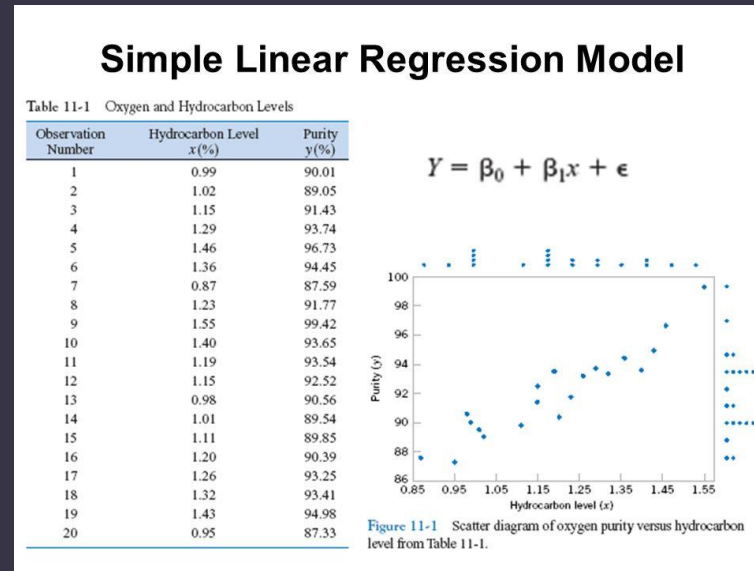


Visualize data more easily

# Goal 1: Why use Stata? - Level 3
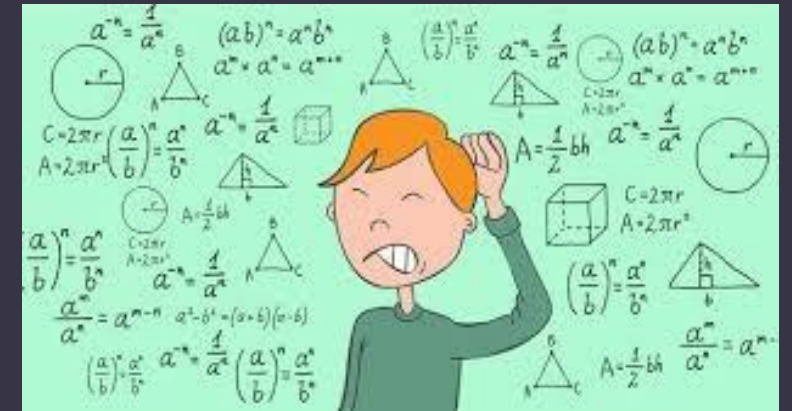
Stata in Data Science          Statistical inference          Less fear to quantitative

Generate insights in a data-oriented world - without overwhelm

# Scope of course - Level 1



Organize large quantities of information

# Course Methodology

Why? How? Sure?



1. Section goals

2. Short video lessons

3. Comprehension Check

# Goal 2 - Where do I Start?

JUAN SEBASTIAN CUERVO

# Goal 2 - Where do I Start ?



Stata interface and windows

Data editor and browse

# Goal 2 - Where do I Start ?



Help and search in Stata

Establish a working directory

# Goal 3 - How to import data in Stata?

JUAN SEBASTIAN CUERVO

# Goal 3: How to import data in Stata?



Stata (.dta)

Text (.txt)

Excel (.dta)

Save files

# Goal 4 - Why keep track of your work in Stata?

JUAN SEBASTIAN CUERVO

# Goal 4: Why keep track of your work in Stata?



Data reproducibility - Do file



Log files

# Goal 5 - Why do I want to alter my dataset?

JUAN SEBASTIAN CUERVO

# Goal 5: Why do I want to alter my dataset?



What are these variables?

What does these values mean?

# Goal 6 - How to analyze my data ?

JUAN SEBASTIAN CUERVO

# Goal 6: How to analyze my data?

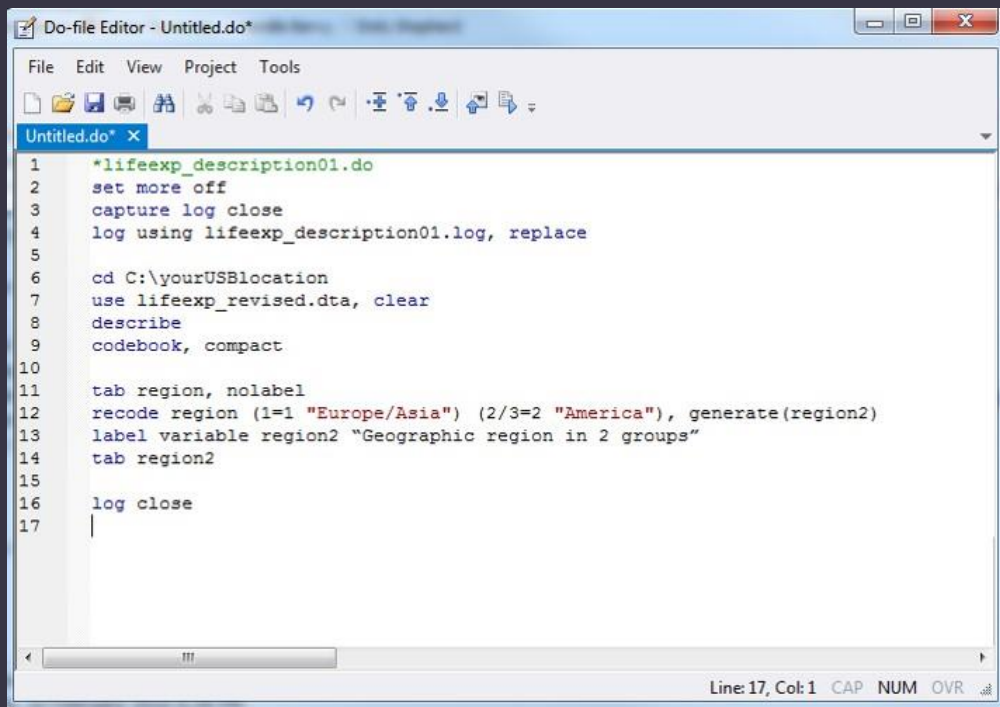| Variable | Obs | Mean | Std. Dev. | Min | Max |
|---|---|---|---|---|---|
| make | 0 | | | | |
| price | 74 | 6165.257 | 2949.496 | 3291 | 15906 |
| mpg | 74 | 21.2973 | 5.785503 | 12 | 41 |
| rep78 | 69 | 3.405797 | .9899323 | 1 | 5 |
| headroom | 74 | 2.993243 | .8459948 | 1.5 | 5 |
| trunk | 74 | 13.75676 | 4.277404 | 5 | 23 |
| weight | 74 | 3019.459 | 777.1936 | 1760 | 4840 |
| length | 74 | 187.9324 | 22.26634 | 142 | 233 |
| turn | 74 | 39.64865 | 4.399354 | 31 | 51 |
| displacement | 74 | 197.2973 | 91.83722 | 79 | 425 |
| gear_ratio | 74 | 3.014865 | .4562871 | 2.19 | 3.89 |
| foreign | 74 | .2972973 | .4601885 | 0 | 1 |

. browse

. summarize mpg price

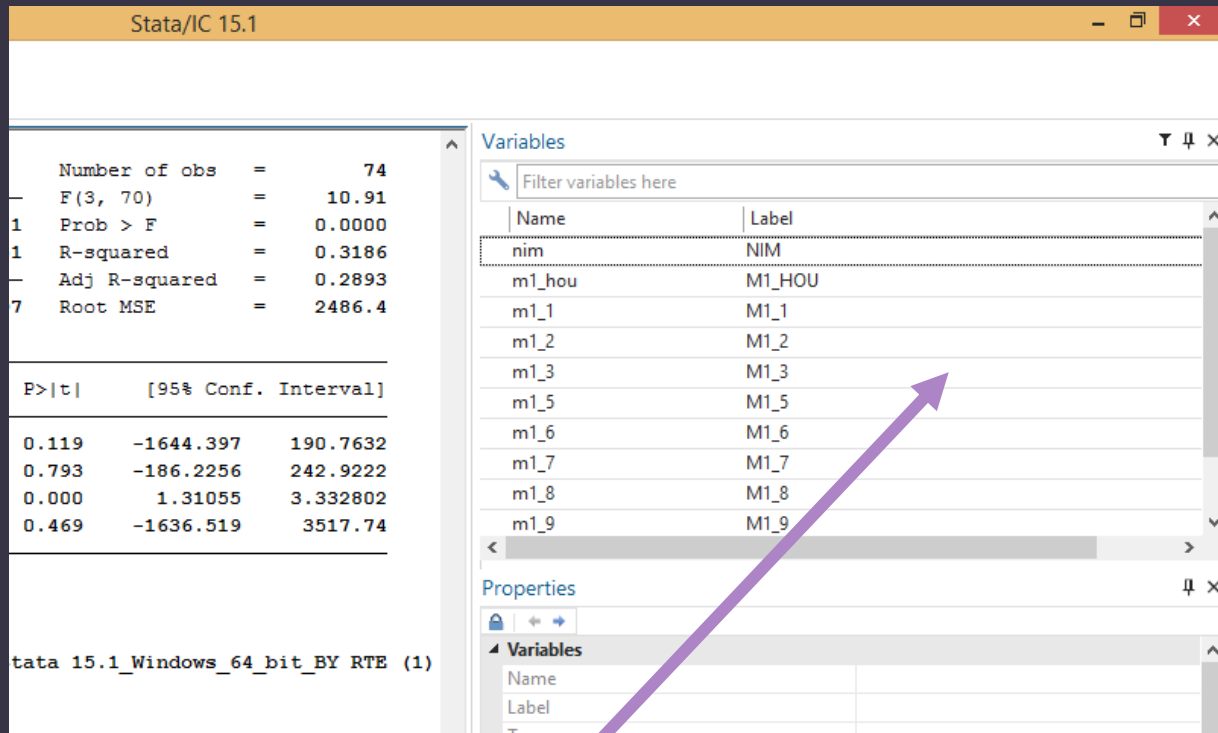| Variable | Obs | Mean | Std. Dev. | Min | Max |
|---|---|---|---|---|---|
| mpg | 74 | 21.2973 | 5.785503 | 12 | 41 |
| price | 74 | 6165.257 | 2949.496 | 3291 | 15906 |

.

. tab V083098x

| J1x. SUMMARY: R Party Identification | Freq. | Percent | Cum. |
|---|---|---|---|
| -1. INAP, -9 in J1; -8,-9 in J1a; -8,-9 | 40 | 1.72 | 1.72 |
| 0. Strong Democrat (1;1;-1) | 580 | 24.98 | 26.70 |
| 1. Weak Democrat (1;5;-1) | 393 | 16.93 | 43.63 |
| 2. Independent-Democrat (3,4,5,-8;-1;5) | 392 | 16.88 | 60.51 |
| 3. Independent-Independent (3,4,5,-8;-1 | 264 | 11.37 | 71.88 |
| 4. Independent-Republican (3,4,5,-8;-1; | 223 | 9.60 | 81.48 |
| 5. Weak Republican (2;5;-1) | 200 | 8.61 | 90.09 |
| 6. Strong Republican (2;1;-1) | 230 | 9.91 | 100.00 |
| Total | 2,322 | 100.00 | |

.

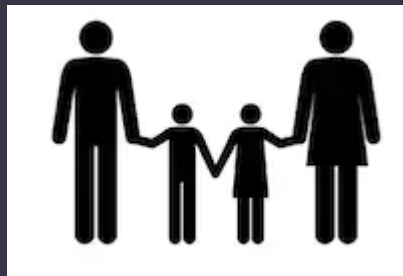Obtain summary statistics          Obtain frequencies, percents

# Goal 7 - Why do I need to join databases?

JUAN SEBASTIAN CUERVO

# Goal 7: Why do I need to join databases?

Household database

People database

Cross information

| Household | People |
|-----------|--------|
| 1 | 1 |
| 1 | 2 |
| 1 | 3 |
| 1 | 4 |
| 2 | 1 |
| 2 | 2 |

# Goal 7: Why do I need to join databases?

## Database 1

| Key | Var 1 | Var 2 |
|-----|-------|-------|
| 1   |       |       |
| 2   |       |       |
| 3   |       |       |
| 4   |       |       |
| 5   |       |       |

## Database 2

| Key | Var 3 | Var 4 |
|-----|-------|-------|
| 1   |       |       |
| 2   |       |       |
| 3   |       |       |
| 4   |       |       |
| 5   |       |       |

Merge

Key variable

## Resulting database

| Key | Var 1 | Var 2 | Var 3 | Var 4 |
|-----|-------|-------|-------|-------|
| 1   |       |       |       |       |
| 2   |       |       |       |       |
| 3   |       |       |       |       |
| 4   |       |       |       |       |
| 5   |       |       |       |       |

# Goal 7: Why do I need to join databases?

## Database 1

| Key | Var 1 | Var 2 |
|-----|-------|-------|
| 1   |       |       |
| 2   |       |       |
| 3   |       |       |
| 4   |       |       |
| 5   |       |       |

## Database 2

| Key | Var 1 | Var 2 |
|-----|-------|-------|
| 6   |       |       |
| 7   |       |       |
| 8   |       |       |

Append →

## Resulting database

| Key | Var 1 | Var 2 |
|-----|-------|-------|
| 1   |       |       |
| 2   |       |       |
| 3   |       |       |
| 4   |       |       |
| 5   |       |       |
| 6   |       |       |
| 7   |       |       |
| 8   |       |       |

# Goal 7: Why do I need to join databases?

## Database 1

| Year | Production |
|------|-----------|
| 1992 | 10.000 |
| 1992 | 16.000 |
| 1993 | 5.000 |
| 1993 | 7.000 |
| 1993 | 4.000 |
| 1993 | 4.000 |
| 1994 | 3.000 |
| 1994 | 3.000 |

Collapse Statistic (Mean)

## Results database

| Year | Mean Production |
|------|-----------------|
| 1992 | 13.000 |
| 1993 | 5.000 |
| 1994 | 3.000 |

# Goal 7: Why do I need to join databases?

Reshape

## Database 1

| Var 1 | Var 2 | Income |
|-------|-------|--------|
| 1 | 1 | 10.000 |
| 1 | 2 | 30.000 |
| 1 | 3 | 40.000 |
| 2 | 1 | 15.000 |
| 2 | 2 | 8.000 |

Reshape →

## Results database

| Var 1 | Income 1 | Income 2 | Income 3 |
|-------|----------|----------|----------|
| 1 | 10.000 | 30.000 | 40.000 |
| 2 | 15.000 | 8.000 | - |

# Recap Course

JUAN SEBASTIAN CUERVO

# Goal 1: Why use Stata? - Level 1



Stata interface and windows



Data editor and browse

# Goal 2 - Where do I Start ?



Stata interface and windows



Data editor and browse

# Goal 4: Why keep track of your work in Stata?



Data reproducibility - Do file

Log files

# Goal 5: Why do I want to alter my dataset?



What are these variables?

What does these values mean?

# Goal 6: How to analyze my data?



Obtain summary statistics



Obtain frequencies, percents

# Goal 7: Why do I need to join databases?

## Database 1

| Key | Var 1 | Var 2 |
|-----|-------|-------|
| 1   |       |       |
| 2   |       |       |
| 3   |       |       |
| 4   |       |       |
| 5   |       |       |

## Database 2

| Key | Var 3 | Var 4 |
|-----|-------|-------|
| 1   |       |       |
| 2   |       |       |
| 3   |       |       |
| 4   |       |       |
| 5   |       |       |

Merge

Key variable

## Resulting database

| Key | Var 1 | Var 2 | Var 3 | Var 4 |
|-----|-------|-------|-------|-------|
| 1   |       |       |       |       |
| 2   |       |       |       |       |
| 3   |       |       |       |       |
| 4   |       |       |       |       |
| 5   |       |       |       |       |