Descriptive Analysis

R for Stata Users

Luiza Andrade, Marc-Andrea Fiorina, Rob Marty, Maria Reyes Retana, Rony Rodriguez Ramirez, Luis Eduardo San Martin, Leonardo Viotti
The World Bank | WB Github



Introduction

Initial Setup

If You Attended Session 2

If You Did Not Attend Session 2

1. Go to the dime-r-training-mar2024 folder that you created yesterday, and open the dime-r-training-mar2024 R project that you created there.

Table of contents

- 1. Quick summary statistics
- 2. Descriptive statistics tables
- 3. Exporting descriptive statistics tables
- 4. Formatting tables
- 5. Aggregating observations
- 6. Running regressions
- 7. Exporting regression tables

Workflows for outputs

Not reproducible

Anything that requires

Copy-pasting

Reproducible

R Markdown: dynamic document containing code and text that is exported directly from R into PDF, HTML, Word, Power Point and other formats

LaTeX: typesetting system used for scientific publications that automatically reloads tables and figures every time the document is rendered

Setting the stage

Load the packages that we will use today

```
# Install new packages
install.packages("modelsummary")
install.packages("fixest")
install.packages("janitor")
install.packages("huxtable")
install.packages("openxlsx")
# Load packages
library(here)
library(tidyverse)
library(modelsummary)
library(fixest)
library(janitor)
library(huxtable)
library(openxlsx)
```

Setting the stage

Load the data that we will use today: Stata's census dataset

Tip: Use here, as we saw in the data wrangling session.

```
# Load data
census <-
read_rds(
here(
    "DataWork",
    "DataSets",
    "Final",
    "census.rds"
)</pre>
```

02:00

Taking a peek at the data

```
glimpse(census)
```

```
## Rows: 50
## Columns: 13
## $ state
              <chr> "Alabama", "Alaska", "Arizona", "Arkansas", "California", "Co...
              <chr> "AL", "AK", "AZ", "AR", "CA", "CO", "CT", "DE", "FL", "GA", "...
## $ state2
## $ region
              <fct> South, West, West, South, West, West, NE, South, South, South...
## $ pop
              <int> 3893888, 401851, 2718215, 2286435, 23667902, 2889964, 3107576...
## $ poplt5
              <int> 296412, 38949, 213883, 175592, 1708400, 216495, 185188, 41151...
## $ pop5 17
              <int> 865836, 91796, 577604, 495782, 4680558, 592318, 637731, 12544...
## $ pop18p
              <int> 2731640, 271106, 1926728, 1615061, 17278944, 2081151, 2284657...
## $ pop65p
              <int> 440015, 11547, 307362, 312477, 2414250, 247325, 364864, 59179...
## $ popurban <int> 2337713, 258567, 2278728, 1179556, 21607606, 2329869, 2449774...
## $ medage
             <dbl> 29.3, 26.1, 29.2, 30.6, 29.9, 28.6, 32.0, 29.8, 34.7, 28.7, 2...
## $ death
              <int> 35305, 1604, 21226, 22676, 186428, 18925, 26005, 5123, 104190...
## $ marriage <int> 49018, 5361, 30223, 26513, 210864, 34917, 26048, 4437, 108344...
## $ divorce <int> 26745, 3517, 19908, 15882, 133541, 18571, 13488, 2313, 71579,...
```

Quick summary statistics

Exploring a dataset

```
summary(x, digits)
```

Equivalent to Stata's codebook. Its arguments are:

- x: the object you want to summarize, usually a vector or data frame
- digits: the number of decimal digits to be displayed

Exercise 1

Use the summary() function to describe the census data frame.

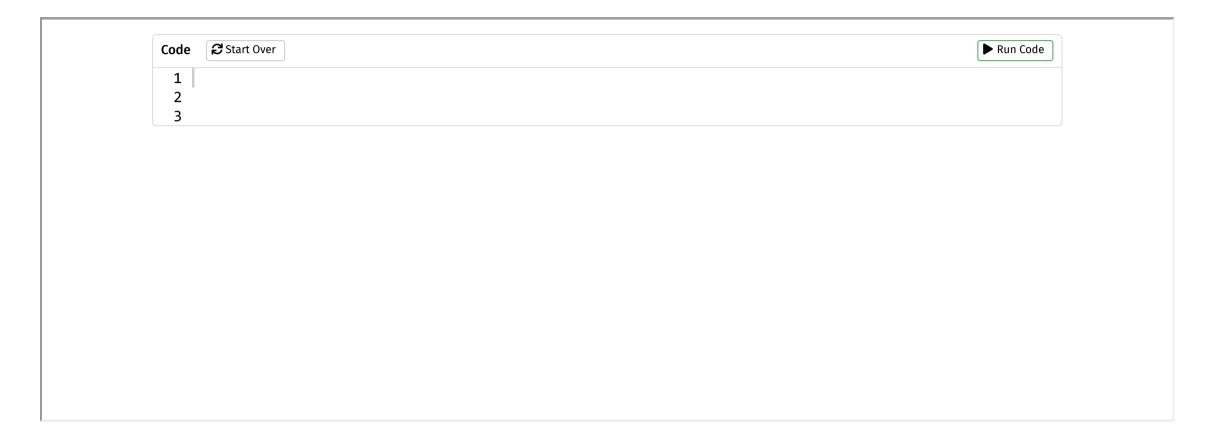
00:45

Exploring a dataset

summary(census)

```
state
                     state2
                                         region
                                                       pop
Length:50
                  Length:50
                                                  Min.
                                                       : 401851
Class : character
                  Class :character
                                     N Cntrl:12
                                                  1st Qu.: 1169218
                                     South :16
Mode :character
                  Mode :character
                                                  Median : 3066433
                                                       : 4518149
                                           :13
                                     West
                                                  Mean
                                                  3rd Ou.: 5434033
                                                        :23667902
                                                  Max.
                                                         pop65p
   poplt5
                    pop5 17
                                       pop18p
     : 35998
                 Min. : 91796
                                   Min. : 271106
                                                     Min. : 11547
1st Qu.: 98831
                 1st Qu.: 257949
                                   1st Qu.: 823702
                                                     1st Qu.: 118660
Median : 227468
                 Median : 629654
                                   Median : 2175130
                                                     Median : 370495
     : 326278
                       : 945952
                                        : 3245920
                                                     Mean : 509503
                 Mean
                                   Mean
3rd Qu.: 361321
                 3rd Qu.:1143292
                                   3rd Qu.: 3858173
                                                     3rd Qu.: 580087
     :1708400
                 Max.
                        :4680558
                                          :17278944
                                   Max.
                                                     Max.
                                                           : 2414250
   popurban
                      medage
                                      death
                                                     marriage
     : 172735
                  Min.
                         :24.20
                                  Min.
                                           1604
                                                  Min. : 4437
1st Qu.: 826651
                  1st Qu.:28.73
                                  1st Qu.: 9087
                                                  1st Qu.: 14840
Median : 2156905
                  Median : 29.75
                                  Median : 26177
                                                  Median : 36279
     : 3328253
                         :29.54
                                       : 39474
                                                  Mean : 47701
                  Mean
                                  Mean
3rd Qu.: 3403450
                  3rd Qu.:30.20
                                  3rd Qu.: 46533
                                                  3rd Qu.: 57338
     :21607606
                        :34.70
                                         :186428
                                                        :210864
                  Max.
                                  Max.
                                                   Max.
   divorce
     : 2142
Min.
1st Qu.: 6898
Median : 17113
     : 23679
3rd Qu.: 27987
     :133541
```

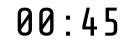
Exploring a dataset



- summary() can also be used with a single variable.
- When used with continuous variables, it works similarly to summarize in Stata.
- When used with categorical variables, it works similarly to tabulate.

Exercise 2

Use the summary() function to display summary statistics for a continuous variable in the census data frame.

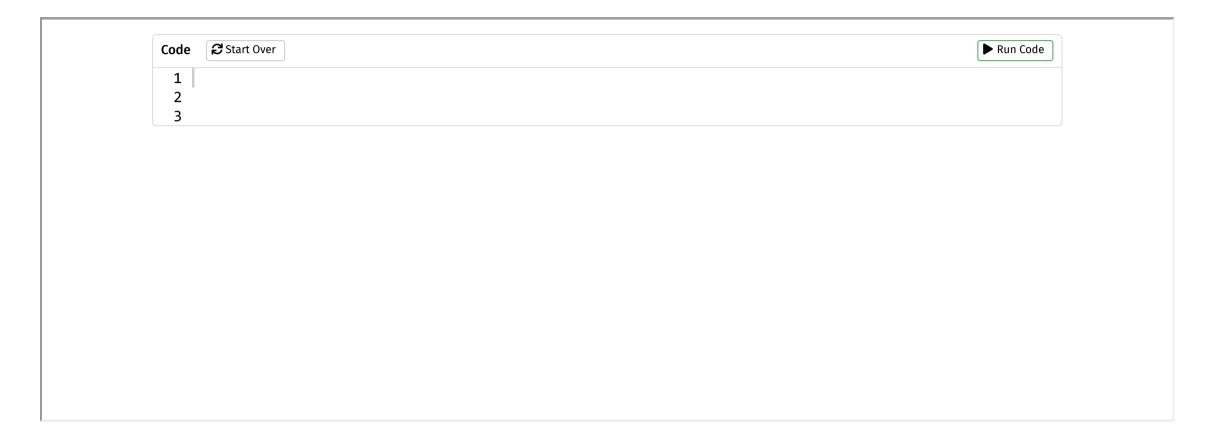


Exercise 2

Use the summary() function to display summary statistics for a continuous variable in the census data frame.

summary(census\$pop)

```
## Min. 1st Qu. Median Mean 3rd Qu. Max.
## 401851 1169218 3066433 4518149 5434033 23667902
```



Summarizing categorical variables

```
tabyl(x, ...)
```

Equivalent to tabulate in Stata, creates a frequency table. Its main arguments are vectors to be tabulated.

- x: the object you want to summarize, usually a vector or data frame
- ... additional options as show_na, or show_missing_levels.

Exercise 3

Use the tabyl() function to display frequency tables for:

- 1. The variable region in the census data frame
- 2. The variables region and state in the census data frame, simultaneously

01:00

Summarizing categorical variables

One way tabulation

```
census %>%
  tabyl(region)
```

region	n	percent
NE	9	0.18
N Cntrl	12	0.24
South	16	0.32
West	13	0.26

Summarizing categorical variables

Two way tabulation

```
census %>%
  tabyl(state, region)
```

state	NE	N Cntrl	South	West
Alabama	0	0	1	0
Alaska	0	0	0	1
Arizona	0	0	0	1
Arkansas	0	0	1	0
California	0	0	0	1
Colorado	0	0	0	1
Connecticut	1	0	0	0

Descriptives tables

Descriptives tables

What if you want to...

- ...export the summary statistics to another software?
- ...customize which statistics to display?
- ...format the table?

Well, then you will need a few more packages

- There are many packages that can be used both for displaying and exporting summary statistics
- Today we will show you a combination of two packages: modelsummary and huxtable
- We chose this combination because together, they can perform all the tasks we are interested in
- In fact, modelsummary can perform most of them by itself -- with the exception of exporting formatted tables to Excel

The package *modelsummary* contains a family of functions called **datasummary** which can be used to create different types of summary statistics tables. These include:

- datasummary_skim, to create descriptive statistics tables
- datasummary_balance, to create balance tables
- datasummary_correlation, to create a correlation table
- datasummary_crosstab, to create a twoway tabulation
- datasummary, to create customized descriptive statistics tables

datasummary_skim(data, output,)

- data: the data set to be summarized, the only required argument
- **output:** the type of output desired
- ...: additional options allow for formatting customization, such as including notes and titles

```
datasummary_skim(
  data,
  type = "numeric",
  output = "default",
  histogram = TRUE,
  title = NULL,
  notes = NULL,
  ...
)
```

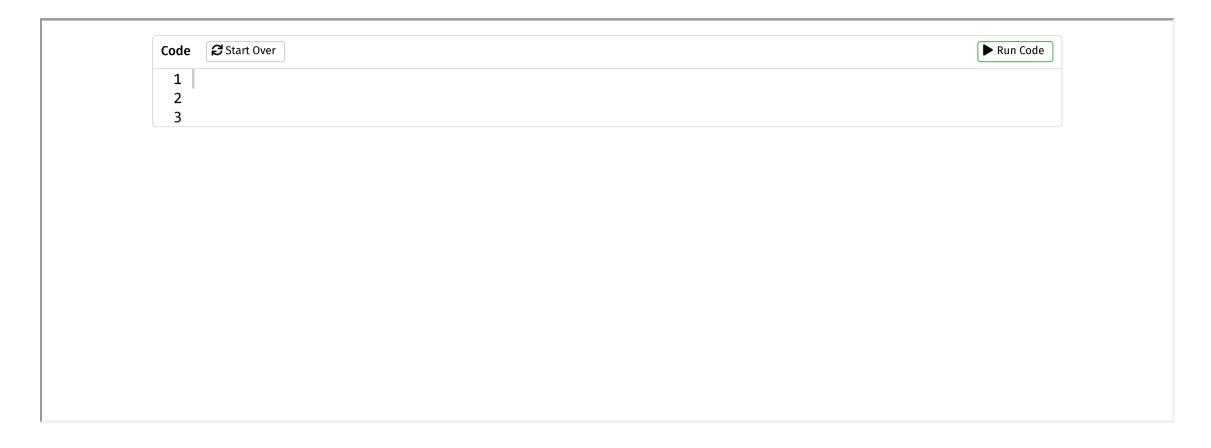
Exercise 4

Use datasummary_skim() to create a descriptive statistics table for the census data.

00:45

datasummary_skim(census)

	Unique (#)	Missing (%)	Mean	SD	Min	Median	Мах	
pop	50	0	4518149.4	4715037.8	401851.0	3066433.0	23667902.0	
poplt5	50	0	326277.8	331585.1	35998.0	227467.5	1708400.0	L .
pop5_17	50	0	945951.6	959372.8	91796.0	629654.0	4680558.0	L .
рор18р	50	0	3245920.1	3430531.3	271106.0	2175130.0	17278944.0	L
pop65p	50	0	509502.8	538932.4	11547.0	370495.0	2414250.0	
popurban	50	0	3328253.2	4090177.9	172735.0	2156905.0	21607606.0	_
medage	37	0	29.5	1.7	24.2	29.8	34.7	
death	50	0	39474.3	41742.3	1604.0	26176.5	186428.0	_
marriage	50	0	47701.4	45130.4	4437.0	36279.0	210864.0	
divorce	50	0	23679.4	25094.0	2142.0	17112.5	133541.0	_



- modelsummary summarizes only numeric variables by default.
- To summarize categorical variables, use the argument type

```
datasummary_skim(census, type = "categorical")
```

data	N	%
NE	9	18.0
N Cntrl	12	24.0
South	16	32.0
West	13	26.0

You can also customize the variables and statistics to include using a **formula** with the datasummary() function.

```
datasummary(formula, data, output, ...)
```

- formula: a two-sided formula to describe the table: rows ~ columns
- data: the data set to be summarized
- output: the type of output desired
- ...: additional options allow for formatting customization

```
datasummary(
  var1 + var2 + var3 ~ stat1 + stat2 + stat3 + stat4,
  data = data
)
```

Exercise 5

Create a table showing the number of observations, mean, standard deviation, minimum, maximum and median value for all the population, number of deaths, number of marriage and number of divorces in the census data.

```
datasummary(
  var1 + var2 + var3 ~ stat1 + stat2 + stat3 + stat4,
  data = data
)
```

Tip: some of the allowed statistics are N, Mean, SD, Min, Max, Median, P0, P25, P50, P75, P100, Histogram

```
datasummary(
   pop + death + marriage + divorce ~ N + Mean + SD + Median + Min + Max,
   data = census
)
```

	N	Mean	SD	Median	Min	Max
рор	50	4518149.44	4715037.75	3066433.00	401851.00	23667902.00
death	50	39474.26	41742.35	26176.50	1604.00	186428.00
marriage	50	47701.40	45130.42	36279.00	4437.00	210864.00
divorce	50	23679.44	25094.01	17112.50	2142.00	133541.00

```
datasummary(
  All(census) ~ N + Mean + SD + Median + Min + Max,
  data = census
)
```

	N	Mean	SD	Median	Min	Мах
рор	50	4518149.44	4715037.75	3066433.00	401851.00	23667902.00
poplt5	50	326277.78	331585.14	227467.50	35998.00	1708400.00
pop5_17	50	945951.60	959372.83	629654.00	91796.00	4680558.00
pop18p	50	3245920.06	3430531.31	2175130.00	271106.00	17278944.00
pop65p	50	509502.80	538932.38	370495.00	11547.00	2414250.00
popurban	50	3328253.18	4090177.93	2156905.00	172735.00	21607606.00
medage	50	29.54	1.69	29.75	24.20	34.70
death	50	39474.26	41742.35	26176.50	1604.00	186428.00
marriage	50	47701.40	45130.42	36279.00	4437.00	210864.00
divorce	50	23679.44	25094.01	17112.50	2142.00	133541.00

Balance tables with modelsummary

```
census_rct <-
  census %>%
 mutate(
   treatment = as.numeric(runif(n()) > 0.5)
  ) 응>응
  select(
    -starts_with("state")
datasummary_balance(
  ~ treatment,
  data = census_rct
```

Balance tables with modelsummary

		0		1			
		Mean	Std. Dev.	Mean	Std. Dev.	Diff. in Means	Std. Error
pop		6127525.7	6824185.9	3828416.7	3351348.8	-2299109.0	1850820.2
poplt5		421233.2	474619.5	285582.6	244984.0	-135650.6	129353.6
pop5_17		1241312.6	1354292.9	819368.3	718505.4	-421944.3	370167.5
pop18p		4464979.9	5000055.2	2723465.8	2393654.8	-1741514.1	1352924.9
pop65p		710303.4	793981.0	423445.4	365507.6	-286858.0	214112.3
popurban		5063526.5	6138886.0	2584564.6	2587414.0	-2478961.8	1644284.7
medage		30.0	2.0	29.3	1.5	-0.7	0.6
death		53856.7	60274.9	33310.3	29745.6	-20546.4	16354.9
marriage		65459.0	59965.6	40091.0	35439.1	-25368.0	16601.5
divorce		32523.5	35141.1	19889.1	18701.3	-12634.4	9608.3
		N	Pct.	N	Pct.		
region	NE	4	26.7	5	14.3		
	N Cntrl	2	13.3	10	28.6		
	South	3	20.0	13	37.1		
	West	6	40.0	7	20.0		

Exporting tables

Exporting modelsummary table to LaTeX

To export the tables we created, we can simply use the option output:

```
descriptives <-
 All(census) ~ N + Mean + SD + Median + Min + Max
datasummary(
  descriptives,
  data = census,
 output = here( # file path to output file
   "DataWork",
   "Output",
    "Raw",
    "summary-stats.tex"
```

Exporting modelsummary table

Other valid output formats include:

- .docx
- .pptx
- .html
- .md

Exporting modelsummary table

Other valid output formats include:

- docx
- .pptx
- .html
- .md
- ... but not .xls

Exporting modelsummary table to Excel

- To export the table to Excel, we will first convert it into an object of type huxtable
- huxtable is another R package, one that allows not only for exporting tables, but also for extensive customization
- Before getting to the customization part, however, let's export this table:

```
# Create the huxtable object
summary stats table <-
  datasummary(
    descriptives,
    data = census,
    output = "huxtable"
# Export it to Excel
quick xlsx(
  summary stats table, # object to be exported
 file = here( # file path to output file
    "DataWork".
    "Output",
    "Raw".
    "summary-stats.xlsx"
```

Exporting tables

A similar code can also export the same table to a self-standing LaTeX document

```
# Export to LaTeX
quick_latex(
    summary_stats_table,
    file = here(
        "DataWork",
        "Output",
        "Raw",
        "summary-stats.tex"
)
```

Exporting tables to different Excel tabs

```
# Start a new workbook
wb <- createWorkbook()
# Add one sheet to it
wh <-
 as Workbook(
    summary_stats_table,
    Workbook = wb,
    sheet = "Summary stats"
# Add another sheet to it
wh <-
  as_Workbook(
   hux("Mock", "table"),
    Workbook = wb,
    sheet = "Other sheet"
# Save the workbook
saveWorkbook(
  wb, # object to be saved
  file = here( # file path to output file
    "DataWork",
    "Output",
    "Raw",
    "summary-stats.xlsx"
  overwrite = TRUE # replace existing file
```

Exporting tables to LaTeX fragment

```
summary_stats_table %>%
  print_latex() %>% # See LaTeX code

# Save LaTeX code

capture.output(
  file = here(
    "DataWork",
    "Output",
    "Raw",
    "summary-stats.tex"
  )
)
```

You will also need to load the required LaTeX packages. To copy the code that creates a preamble with all of them, use this code:

```
report_latex_dependencies()
```

Formatting tables

Beautifying tables

- huxtable also allows you to customize table formatting so it can be exported with the same layout to multiple software
- Before we do that, however, we will create a version of the data where the variable names are the Stata labels

```
# Fxtract variable labels from data frame
labels <- names(census)
names(labels) <- attributes(census)$var.labels</pre>
# Rename the variables
census labelled <-
  census %>%
  rename(
   labels
# Create a labelled summary table
summary stats table <-
  datasummary(
    All(census_labelled) ~ N + Mean + SD + Median + Min + Max,
    data = census_labelled,
    output = "huxtable"
```

Beautifying tables

The code below shows the table summary_stats_table can be formatted

```
# Format table
summary_stats_table %>%
# Use first row as table header
set header rows(1, TRUE) %>%
 # Use first column as row header
set_header_cols(1, TRUE) %>%
 # Don't round large numbers
set number format(everywhere, 2:ncol(.), "%9.0f") %>%
# Centralize cells in first row
set_align(1, everywhere, "center") %>%
# Set a theme for quick formatting
theme_basic()
```

	N	Mean	SD	Median	Min	Мах
Population	50	4518149	4715038	3066433	401851	23667902
Pop, < 5 year	50	326278	331585	227468	35998	1708400
Pop, 5 to 17 years	50	945952	959373	629654	91796	4680558
Pop, 18 and older	50	3245920	3430531	2175130	271106	17278944
Pop, 65 and older	50	509503	538932	370495	11547	2414250
Urban population	50	3328253	4090178	2156905	172735	21607606
Median age	50	30	2	30	24	35
Number of deaths	50	39474	41742	26177	1604	186428
Number of marriages	50	47701	45130	36279	4437	210864
Number of divorces	50	23679	25094	17113	2142	133541

Export beautified tables

```
# Format table
summary stats table <-
  summary stats table %>%
  set header rows(1, TRUE) %>% # Use first row as table header
  set header cols(1, TRUE) %>% # Use first column as row header
  set number format(everywhere, 2:ncol(.), "%9.0f") %>% # Don't round large numbers
  set align(1, everywhere, "center") %>% # Centralize cells in first row
  theme basic() # Set a theme for quick formatting
quick xlsx(
  summary stats table,
 file = here(
    "DataWork"
   "Output",
    "Raw".
    "summary-stats-basic.xlsx"
```

Export beautified tables

Before

4	Α	В	С	D	Е	F
1	skim_varia	Mean	Median	SD	Min	Max
2	рор	4520000	3070000	4720000	402000	23700000
3	poplt5	326000	227000	332000	36000	1710000
4	pop5_17	946000	630000	959000	91800	4680000
5	рор18р	3250000	2180000	3430000	271000	17300000
6	рор65р	510000	370000	539000	11500	2410000
7	popurban	3330000	2160000	4090000	173000	21600000
8	medage	29.5	29.8	1.69	24.2	34.7
9	death	39500	26200	41700	1600	186000
10	marriage	47700	36300	45100	4440	211000
11	divorce	23700	17100	25100	2140	134000
9	medage death marriage	29.5 39500 47700	29.8 26200 36300	1.69 41700 45100	24.2 1600 4440	34. 18600 21100

After

4	Α	В	С	D	Е	F	G
1		N	Mean	SD	Median	Min	Max
2	Population	50	4518149	4715038	3066433	401851	23667902
3	Pop, < 5 year	50	326278	331585	227468	35998	1708400
4	Pop, 5 to 17 years	50	945952	959373	629654	91796	4680558
5	Pop, 18 and older	50	3245920	3430531	2175130	271106	17278944
6	Pop, 65 and older	50	509503	538932	370495	11547	2414250
7	Urban population	50	3328253	4090178	2156905	172735	21607606
8	Median age	50	30	2	30	24	35
9	Number of deaths	50	39474	41742	26177	1604	186428
10	Number of marriages	50	47701	45130	36279	4437	210864
11	Number of divorces	50	23679	25094	17113	2142	133541

Other themes to play with

jams

Туре	Price	Sugar content
Strawberry	1.90	40.00%
Raspberry	2.10	35.00%
Plum	1.80	50.00%

theme_article

Туре	Price	Sugar content
Strawberry	1.90	40.00%
Raspberry	2.10	35.00%
Plum	1.80	50.00%

theme_green

Туре	Price	Sugar content
Strawberry	1.90	40.00%
Raspberry	2.10	35.00%
Plum	1.80	50.00%

theme_plain

Туре	Price	Sugar content
Strawberry	1.90	40.00%
Raspberry	2.10	35.00%
Plum	1.80	50.00%

theme_bright

Туре	Price	Sugar content
Strawberry	1.90	40.00%
Raspberry	2.10	35.00%
Plum	1.80	50.00%

theme_mondrian

	Price	Sugar content
Strawberry	1.90	
Raspberry	2.10	35.00%
Plum	1.80	50.00%

theme_basic

Туре	Price	Sugar content
Strawberry	1.90	40.00%
Raspberry	2.10	35.00%
Plum	1.80	50.00%

theme_grey

Туре	Price	Sugar content
Strawberry	1.90	40.00%
Raspberry	2.10	35.00%
Plum	1.80	50.00%

theme_orange

Туре	Price	Sugar content
Strawberry	1.90	40.00%
Raspberry	2.10	35.00%
Plum	1.80	50.00%

theme_compact

Type	Price Su	ıgar content
Strawberry	1.90	40.00%
Raspberry	2.10	35.00%
Plum	1.80	50.00%

theme_blue

Туре	Price	Sugar content
Strawberry	1.90	40.00%
Raspberry	2.10	35.00%
Plum	1.80	50.00%

theme_striped

Туре	Price	Sugar content
Strawberry	1.90	40.00%
Raspberry	2.10	35.00%
Plum	1.80	50.00%

Ok, can we run some regressions now?!

The base R command for linear regressions is called lm

lm(formula, data, subset, weights, ...)

- formula: an object of class "formula" containing a symbolic description of the model
- data: a data frame containing the variables indicated in the formula
- subset: an optional vector specifying a subset of observations to be used in the regression
- weights: an optional vector of weights to be used in the regression

Formulas can take three specifications:

- $y \sim x1 + x2$ regresses variable y on covariates x1 and x2
- y ~ x1:x2 regresses variable y on the interaction of covariates x1 and x2
- $y \sim x1*x2$ is equivalent to $y \sim x1 + x2 + x1:x2$

Exercise 6

Using the **census** data, run a regression of the number of divorces on population, urban population and number of marriages.

```
lm(y \sim x1 + x2, data)
```

01:00

Exercise 6

Using the **census** data, run a regression of the number of divorces on population, urban population and number of marriages.

```
reg1 <-
lm(
   divorce ~ pop + popurban + marriage,
   census
)</pre>
```

- The output of regression commands is a list of relevant information.
- By default, it prints only a small portion of this information.
- The best way to visualize results is to store this list in an object and then access its contents using the function summary

Residual standard error: 7466 on 46 degrees of freedom
Multiple R-squared: 0.9169, Adjusted R-squared: 0.9115
F-statistic: 169.2 on 3 and 46 DF, p-value: < 2.2e-16</pre>

```
reg1 <-
  lm(
    divorce ~ pop + popurban + marriage,
    census
summary(reg1)
## Call:
## lm(formula = divorce ~ pop + popurban + marriage, data = census)
## Residuals:
                1Q Median
       Min
## -22892.3 -1665.1 796.5 4138.0 17212.2
## Coefficients:
              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 1.207e+02 1.838e+03 0.066
                                         0.948
       1.044e-03 1.633e-03 0.639
## pop
                                          0.526
## popurban 1.954e-03 1.796e-03 1.088
                                           0.282
## marriage 2.587e-01 5.958e-02 4.342 7.7e-05 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

The feols command from package fixest allows for more flexibility in model specification

feols(formula, data, subset, weights, ...)

- formula: an object of class "formula" containing a symbolic description of the model
- data: a data frame containing the variables indicated in the formula
- *vcov:* one of "iid", "hetero" (or "HC1"), "cluster", "twoway", "NW" (or "newey_west"), "DK" (or "driscoll_kraay"), or "conley"
- subset: an optional vector specifying a subset of observations to be used in the regression
- weights: an optional vector of weights to be used in the regression
- *cluster*: a list of vectors, a character vector of variable names, a formula or an integer vector specifying how to cluster standard errors
- ...

Formulas for feols are more complex, and take the following format: y ~ x1 + x2 | fe1 + fe2 | x3 ~ iv3

- $y \sim x1 + x2$ takes all the same formulas as lm
- fe1 + fe2 list the variables to be included as fixed effects
- x3 ~ iv3 uses instrument iv3 for variable x3

Exercise 7

Using the census data, run a regression of the number of divorces on population, urban population and number of marriages controlling for region fixed effects.

```
feols(
  y ~ x1 + x2 | fe1 + fe2,
  data
)
```

01:00

Exercise 7

Using the census data, run a regression of divorce on population, urban population and number of marriages controlling for region fixed effects.

```
reg2 <-
feols(
    divorce ~ pop + popurban + marriage | region,
    census,
    se = "iid"
)
summary(reg2)</pre>
```

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Within R2: 0.935434

RMSE: 6,257.6 Adj. R2: 0.927695

```
reg2 <-
 feols(
  divorce ~ pop + popurban + marriage | region,
  census,
  se = "iid"
summary(reg2)
## OLS estimation, Dep. Var.: divorce
## Observations: 50
## Fixed-effects: region: 4
## Standard-errors: IID
       Estimate Std. Error t value Pr(>|t|)
       ## pop
```

Some notes on regressions

- Whenever a factor is included in the list of covariates, it is treated as a categorical variable, i.e., as if you had written i.x in Stata.
- Whenever a boolean is included in the list of covariates, it is treated as a dummy variable, where TRUE is 1 and FALSE is 0.

huxtable also has a quick wrapper for regression tables

huxreg(...)

- ...: Models, or a single list of models. Names will be used as column headings.
- number_format: Format for numbering. See number_format() for details.
- stars: Levels for p value stars.
- bold_signif: Where p values are below this number, cells will be displayed in bold.
- note: Footnote for bottom cell, which spans all columns.
- *statistics:* A vector of summary statistics to display.
- coefs: A vector of coefficients to display. To change display names, name the coef vector: c("Displayed title" = "coefficient name", ...)

huxreg(reg1, reg2)

	(1)	(2)	
(Intercept)	120.730		
	(1838.216)		
рор	0.001	0.000	
	(0.002)	(0.002)	
popurban	0.002	2 0.004	
	(0.002)	(0.002)	
marriage	0.259 ***	0.184 **	
	(0.060)	(0.058)	
N	50	50	
R2	0.917	0.937	
logLik	-514.766	-508.024	
AIC	1039.531	1030.048	
*** p < 0.001; ** p < 0.01; * p < 0.05.			

Formatting regression tables

```
huxreg(
 reg1, reg2,
  # Show variable labels instead of names
  coefs = c(
     "Population" = "pop",
     "Urban population" = "popurban",
     "Number of marriages" = "marriage"
   statistics = c("N. obs." = "nobs")) %>%
   add_rows(
    c("Region FE", "No", "Yes"),
    after = 7
```

	(1)	(2)	
Population	0.001	0.000	
	(0.002)	(0.002)	
Urban population	0.002	0.004	
	(0.002)	(0.002)	
Number of marriages	0.259 ***	0.184 **	
	(0.060)	(0.058)	
Region FE	No	Yes	
N. obs.	50	50	
*** p < 0.001; ** p < 0.01; * p < 0.05.			

Exercise 8

Export a regression table with the results of your estimations using lm and feols:

- Use huxreg to combine reg1 and reg2.
- Use quick xlsx or quick latex to export the output of huxreg to your preferred format.

```
# Combine regression results with huxreg
reg table <- huxreg(reg1, reg2)</pre>
# Export to Excel
quick_xlsx(reg_table, file = here(
                                "DataWork".
                                "Output",
                                "Raw".
                                "regression_table.xlsx"
```

```
# Or, export to LaTeX
quick latex(reg table, file = here(
                                "DataWork",
                                "Output".
                                "Raw".
                                "regression table.tex")
```

References and recommendations

- Econometrics with R https://www.econometrics-with-r.org/index.html
- modelsummary documentation: https://vincentarelbundock.github.io/modelsummary/index.html
- Introduction to huxtable: https://cran.r-project.org/web/packages/huxtable/vignettes/huxtable.html
- Using huxtable for regression tables: https://cran.r-project.org/web/packages/huxtable/vignettes/huxreg.html
- Sample code for tables in R: https://github.com/RRMaximiliano/r-latex-tables-sum-stats
- More sample code for tables in R: https://evalsp20.classes.andrewheiss.com/reference/regtables/
- Johns Hopkins Exploratory Data Analysis at Coursera: https://www.coursera.org/learn/exploratory-data-analysis
- Udacity's Data Analysis with R: https://www.udacity.com/course/data-analysis-with-r--ud651

Since we talked about LaTeX so much...

- DIME LaTeX templates and trainings: https://github.com/worldbank/DIME-LaTeX-Templates
- All you need to know about LaTeX: https://en.wikibooks.org/wiki/LaTeX

Thank you!

Appendix

- If you want to show aggregated statistics, the function **summarise** is a powerful tool.
- It is similar to datasummary in that it calculates a series of statistics for a data frame.
- However, it does not have pre-defined statistics, so it requires more manual input.
- On the other hand, its output is a regular data frame, so it is also useful to create constructed data sets.
- Its Stata equivalent would be collapse

```
summarise(.data, ...,)
```

- data: the data frame to be summarized
- ...: Name-value pairs of summary functions. The name will be the name of the variable in the result.

The "name-value" pairs mentioned under ... look like this: new_variable = function(existing_variable), where possible
functions include:

- Center: mean(), median()
- Spread: sd(), IQR(), mad()
- Range: min(), max(), quantile()
- Count: n(), n_distinct()

```
region_stats <-
  census %>%
  group_by(region) %>%
  summarise(
    `Number of States` = n_distinct(state),
    `Total Population` = sum(pop)
)
```

region	Number of States	Total Population
NE	9	49135283
N Cntrl	12	58865670
South	16	74734029
West	13	43172490

Exercise 9

Recreate the region_stats data set, now including the average and the standard deviation of the population.

01:30

```
region_stats <-
  census %>%
  group_by(region) %>%
  summarise(
      `Number of States` = n_distinct(state),
      `Total Population` = sum(pop),
      `Average Population` = mean(pop),
      `SD of Population` = sd(pop)
)
```

region	Number of States	Total Population	Average Population	SD of Population
NE	9	49135283	5459476	5925235
N Cntrl	12	58865670	4905473	3750094
South	16	74734029	4670877	3277853
West	13	43172490	3320961	6217177

Exercise 9

Use huxtable to format and export the object region_stats.

02:00

```
region_stats_table <-
 region_stats %>%
 rename(Region = region) %>%
  as_hux %>%
  set_header_cols("Region", TRUE) %>%
  theme_bright()
quick_xlsx(
 region_stats_table,
   file = here(
     "DataWork",
     "Output",
     "Raw",
     "region-stats.xlsx"
```