

# Applied Data Analysis for Public Policy Studies

Michele Fioretti

2020-08-24



# Contents



# Syllabus



Welcome to the “Applied Data Analysis for Public Policy Studies” course at ScPo! On this page we outline the course and present the Syllabus. 2020/2021 is the first time that we teach this course in this format, so comments are appreciated!

## Objective

The objective of this course is to provide students with the statistical tools to understand and analyze public policy critically. The course introduces students to basic concepts of statistical analysis, including multivariate regression, sampling, hypothesis testing, causality, difference-in-differences analysis, and regression discontinuity design. The course will place a strong emphasis on conceptually understanding the application of these tools, rather than their mathematical foundations. So, for instance, when discussing linear regressions, we will focus on what the estimates mean rather than the mathematical proof underpinning the calculations. To this extent, statistical concepts (i) will be presented through multiple examples from academic articles in the sphere of public policy, and (ii) making use of the statistical language R.

The ultimate goal of the course is to: 1. Become a Reader: Be familiar with simple statistical tools, be able to interpret results by critically thinking through the conceptual limitations of the analyses presented in research papers and policy briefs. 2. Become a User: Conduct statistical analysis and data cleaning exercises using R.

## Syllabus and Requirements

We will be using the free statistical computing language R very intensively. Before coming to the first session, please install R and RStudio as explained at the beginning of chapter ??.

## Course Structure

Groups meet once per week for 2 hours. The main purpose of the weekly meetings is to clarify any questions, and to work together through tutorials. The little theory we need will be covered in this book, and **you are expected to read through this in your own time** before coming to class.

## Slides

There are slides for most book chapters at a dedicated github repository.

## This Book and Other Material

What you are looking at is an online textbook. You can therefore look at it in your browser (as you are doing just now), on your mobile phone or tablet, but you can also download it as a **pdf** file or as an **epub** file for your ebook-reader. We don't have any ambition to actually produce and publish a book for now, so you should just see this as a way to disseminate our lecture notes to you. The second part of course material next to the book is an extensive suite of tutorials and interactive demonstrations, which are all contained in the R package that builds this book (and which you installed by issuing the above commands).

## Teaching

Instructor: Michele Fioretti - email: michele dot fioretti at sciencespo dot fr  
Teaching Assistant: Florentin Cognie - email: florentin dot cognie at sciencespo dot fr

## Open Source

The book and all other content for this course are hosted under an open source license on github. You can contribute to the book by just clicking on the appropriate edit symbol in the top bar of this page. The material in this book comes from “Introduction to Econometrics with R” by lorian Oswald and Jean-Marc Robin and Vincent Viers (2018), which can be found at github repository.

## Assessments

100% assessed coursework. 2 problem sets (worth 40%), quizzes on moodle (10%) and take home exam (50%).

## Communication

We will communicate exclusively on our slack group. You will get an invitation email to join from your instructor in due course.





# Chapter 1

## Introduction to R

### 1.1 Getting Started

R is both a programming language and software environment for statistical computing, which is free and open-source. To get started, you will need to install two pieces of software:

1. R, the actual programming language.
  - Chose your operating system, and select the most recent version.
2. RStudio, an excellent IDE for working with R.
  - Note, you must have R installed to use RStudio. RStudio is simply an interface used to interact with R.

The popularity of R is on the rise, and everyday it becomes a better tool for statistical analysis. It even generated this book!

The following few chapters will serve as a whirlwind introduction to R. They are by no means meant to be a complete reference for the R language, but simply an introduction to the basics that we will need along the way. Several of the more important topics will be re-stressed as they are actually needed for analyses.

This introductory R chapter may feel like an overwhelming amount of information. You are not expected to pick up everything the first time through. You should try all of the code from this chapter, then return to it a number of times as you return to the concepts when performing analyses. We only present the most basic aspects of R. If you want to know more, there are countless online tutorials, and you could start with the official CRAN sample session or have a look at the resources at Rstudio or on this github repo.

## 1.2 Starting R and RStudio

A key difference for you to understand is the one between **R**, the actual programming language, and **RStudio**, a popular interface to **R** which allows you to work efficiently and with greater ease with **R**.

The best way to appreciate the value of **RStudio** is to start using **R** without **RStudio**. To do this, double-click on the **R** GUI that you should have downloaded on your computer following the steps above (on windows or Mac), or start **R** in your terminal (on Linux or Mac) by just typing **R** in a terminal, see figure ?? . You've just opened the **R console** which allows you to start typing code right after the **>** sign, called prompt. Try typing **2 + 2** or **print("Your Name")** and hit the return key. And voilà, your first **R** commands!



Figure 1.1: R GUI symbol and R in a MacOS Terminal

```
florian.oswald — /usr/local/bin/R --no-save — R — R --no-save — 80x23
~ R
R version 3.5.1 (2018-07-02) — "Feather Spray"
Copyright (C) 2018 The R Foundation for Statistical Computing
Platform: x86_64-apple-darwin15.6.0 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

> 2 + 2
[1] 4
>
```

Figure 1.2: R GUI symbol and R in a MacOS Terminal

Typing one command after the other into the console is not very convenient as our analysis becomes more involved. Ideally, we would like to collect all command statements in a file and run them one after the other, automatically.

We can do this by writing so-called **script files** or just **scripts**, i.e. simple text files with extension `.R` or `.r` which can be inserted (or sourced) into an R session. RStudio makes this process very easy.

Open **RStudio** by clicking on the **RStudio** application on your computer, and notice how different the whole environment is from the basic R console – in fact, that very same R console is running in your bottom left panel. The upper-left panel is a space for you to write scripts – that is to say many lines of codes which you can run when you choose to. To run a single line of code, simply highlight it and hit **Command + Return**.

We highly recommend that you use **RStudio** for everything related to this course (in particular, to launch our apps and tutorials).

RStudio has a large number of useful keyboard shortcuts. A list of these can be found using a keyboard shortcut – the keyboard shortcut to rule them all:

- On Windows: **Alt + Shift + K**
- On Mac: **Option + Shift + K**

The **RStudio** team has developed a number of “cheatsheets” for working with both R and **RStudio**. This particular cheatseet for Base R will summarize many of the concepts in this document.<sup>1</sup>

### 1.2.1 First Glossary

- **R**: a statistical programming language
- **RStudio**: an integrated development environment (IDE) to work with R
- **command**: user input (text or numbers) that **R** understands.
- **script**: a list of commands collected in a text file, each separated by a new line, to be run one after the other.

## 1.3 Basic Calculations

To get started, we’ll use R like a simple calculator. Run the following code either directly from your RStudio console, or in RStudio by writting them in a script and running them using **Command + Return**.

---

<sup>1</sup>When programming, it is often a good practice to follow a style guide. (Where do spaces go? Tabs or spaces? Underscores or CamelCase when naming variables?) No style guide is “correct” but it helps to be aware of what others do. The more import thing is to be consistent within your own code. Here are two guides: Hadley Wickham Style Guide, and the Google Style Guide. For this course, our main deviation from these two guides is the use of `=` in place of `<-`. For all practical purposes, you should think `=` whenever you see `<-`.

**Addition, Subtraction, Multiplication and Division**

Math	R code	Result
$3 + 2$	<code>3 + 2</code>	5
$3 - 2$	<code>3 - 2</code>	1
$3 \cdot 2$	<code>3 * 2</code>	6
$3/2$	<code>3 / 2</code>	1.5

**Exponents**

Math	R code	Result
$3^2$	<code>3 ^ 2</code>	9
$2^{(-3)}$	<code>2 ^ (-3)</code>	0.125
$100^{1/2}$	<code>100 ^ (1 / 2)</code>	10
$\sqrt{100}$	<code>sqrt(100)</code>	10

**Mathematical Constants**

Math	R code	Result
$\pi$	<code>pi</code>	3.1415927
$e$	<code>exp(1)</code>	2.7182818

**Logarithms**

Note that we will use `ln` and `log` interchangeably to mean the natural logarithm. There is no `ln()` in R, instead it uses `log()` to mean the natural logarithm.

Math	R code	Result
$\log(e)$	<code>log(exp(1))</code>	1
$\log_{10}(1000)$	<code>log10(1000)</code>	3
$\log_2(8)$	<code>log2(8)</code>	3
$\log_4(16)$	<code>log(16, base = 4)</code>	2

**Trigonometry**

Math	R code	Result
$\sin(\pi/2)$	<code>sin(pi / 2)</code>	1
$\cos(0)$	<code>cos(0)</code>	1

## 1.4 Getting Help

In using R as a calculator, we have seen a number of functions: `sqrt()`, `exp()`, `log()` and `sin()`. To get documentation about a function in R, simply put a question mark in front of the function name, or call the function `help(function)` and RStudio will display the documentation, for example:

```
?log
?sin
?paste
?lm
help(lm)    # help() is equivalent
help(ggplot, package="ggplot2") # show help from a certain package
```

Frequently one of the most difficult things to do when learning R is asking for help. First, you need to decide to ask for help, then you need to know how to ask for help. Your very first line of defense should be to Google your error message or a short description of your issue. (The ability to solve problems using this method is quickly becoming an extremely valuable skill.) If that fails, and it eventually will, you should ask for help. There are a number of things you should include when contacting an instructor, or posting to a help website such as Stack Overflow.

- Describe what you expect the code to do.
- State the end goal you are trying to achieve. (Sometimes what you expect the code to do, is not what you want to actually do.)
- Provide the full text of any errors you have received.
- Provide enough code to recreate the error. Often for the purpose of this course, you could simply post your entire `.R` script or `.Rmd` to `slack`.
- Sometimes it is also helpful to include a screenshot of your entire RStudio window when the error occurs.

If you follow these steps, you will get your issue resolved much quicker, and possibly learn more in the process. Do not be discouraged by running into errors and difficulties when learning R. (Or any other technical skill.) It is simply part of the learning process.

## 1.5 Installing Packages

R comes with a number of built-in functions and datasets, but one of the main strengths of R as an open-source project is its package system. Packages add additional functions and data. Frequently if you want to do something in R, and it is not available by default, there is a good chance that there is a package that will fulfill your needs.

To install a package, use the `install.packages()` function. Think of this as buying a recipe book from the store, bringing it home, and putting it on your shelf (i.e. into your library):

```
install.packages("ggplot2")
```

Once a package is installed, it must be loaded into your current R session before being used. Think of this as taking the book off of the shelf and opening it up to read.

```
library(ggplot2)
```

Once you close R, all the packages are closed and put back on the imaginary shelf. The next time you open R, you do not have to install the package again, but you do have to load any packages you intend to use by invoking `library()`.

## 1.6 Code vs Output in this Book

A quick note on styling choices in this book. We had to make a decision how to visually separate R code and resulting output in this book. We decided to prefix all output lines with `#OUT>` to make the distinction. A typical code snippet with output is thus going to look like this:

```
1 + 3
```

```
#OUT> [1] 4
```

```
# everything after a # is a comment, i.e. R disregards it.
```

where you see on the first line the R code, and on the second line the output. As mentioned, that line starts with `#OUT>` to say this is an output, followed by `[1]` (indicating this is a vector of length one - more on this below!), followed by the actual result - `1 + 3 = 4`!

Notice that you can simply copy and paste all the code you see into your R console. In fact, you are strongly encouraged to actually do this and try out **all the code** you see in this book.

Finally, please note that this way of showing output is fully our choice in this textbook, and that you should expect other output formats elsewhere. For

example, in my RStudio console, the above code and output looks like this:

```
> 1 + 3  
[1] 4
```

## 1.7 ScPoEconometrics Package

To fully take advantage of our course, please install the associated R package directly from its online code repository. You can do this by copy and pasting the following three lines into your R console:

```
if (!require("devtools")) install.packages("devtools")  
library(devtools)  
install_github(repo = "ScPoEcon/ScPoEconometrics")
```

In order to check whether everything works fine, you could load the library, and check it's current version:

```
library(ScPoEconometrics)  
packageVersion("ScPoEconometrics")
```

```
#OUT> [1] '0.2.6'
```

## 1.8 Data Types

R has a number of basic data types. While R is not a strongly typed language (i.e. you can be agnostic about types most of the times), it is useful to know what data types are available to you:

- Numeric
  - Also known as Double. The default type when dealing with numbers.
  - Examples: 1, 1.0, 42.5
- Integer
  - Examples: 1L, 2L, 42L
- Complex
  - Example: 4 + 2i
- Logical
  - Two possible values: TRUE and FALSE
  - You can also use T and F, but this is not recommended.
  - NA is also considered logical.
- Character
  - Examples: "a", "Statistics", "1 plus 2."
- Categorical or **factor**
  - A mixture of integer and character. A **factor** variable assigns a label to a numeric value.

- For example `factor(x=c(0,1),labels=c("male","female"))` assigns the string male to the numeric values 0, and the string female to the value 1.

## 1.9 Data Structures

R also has a number of basic data structures. A data structure is either homogeneous (all elements are of the same data type) or heterogeneous (elements can be of more than one data type).

Dimension	Homogeneous	Heterogeneous
1	Vector	List
2	Matrix	Data Frame
3+	Array	nested Lists

### 1.9.1 Vectors

Many operations in R make heavy use of **vectors**. A vector is a container for objects of identical type (see ?? above). Vectors in R are indexed starting at 1. That is what the [1] in the output is indicating, that the first element of the row being displayed is the first element of the vector. Larger vectors will start additional rows with something like [7] where 7 is the index of the first element of that row.

Possibly the most common way to create a vector in R is using the `c()` function, which is short for “combine”. As the name suggests, it combines a list of elements separated by commas. (Are you busy typing all of those examples into your R console? :-) )

```
c(1, 3, 5, 7, 8, 9)
```

```
#OUT> [1] 1 3 5 7 8 9
```

Here R simply outputs this vector. If we would like to store this vector in a **variable** we can do so with the **assignment** operator `=`. In this case the variable `x` now holds the vector we just created, and we can access the vector by typing `x`.

```
x = c(1, 3, 5, 7, 8, 9)
x
```

```
#OUT> [1] 1 3 5 7 8 9
```

As an aside, there is a long history of the assignment operator in R, partially due to the keys available on the keyboards of the creators of the S language.



(Which preceded R.) For simplicity we will use =, but know that often you will see <- as the assignment operator.

Because vectors must contain elements that are all the same type, R will automatically **coerce** (i.e. convert) to a single type when attempting to create a vector that combines multiple types.

```
c(42, "Statistics", TRUE)
```

```
#OUT> [1] "42"          "Statistics" "TRUE"
```

```
c(42, TRUE)
```

```
#OUT> [1] 42 1
```

Frequently you may wish to create a vector based on a sequence of numbers. The quickest and easiest way to do this is with the : operator, which creates a sequence of integers between two specified integers.

```
(y = 1:100)
```

```
#OUT> [1] 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27
#OUT> [33] 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58
#OUT> [65] 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90
#OUT> [97] 97 98 99 100
```

Here we see R labeling the rows after the first since this is a large vector. Also, we see that by putting parentheses around the assignment, R both stores the vector in a variable called y and automatically outputs y to the console.

Note that scalars do not exist in R. They are simply vectors of length 1.

```
2
```

```
#OUT> [1] 2
```

If we want to create a sequence that isn't limited to integers and increasing by 1 at a time, we can use the seq() function.

```
seq(from = 1.5, to = 4.2, by = 0.1)
```

```
#OUT> [1] 1.5 1.6 1.7 1.8 1.9 2.0 2.1 2.2 2.3 2.4 2.5 2.6 2.7 2.8 2.9 3.0 3.1 3.2 3.3 3.4 3.5 3.6 3.7 3.8
```

We will discuss functions in detail later, but note here that the input labels from, to, and by are optional.

```
seq(1.5, 4.2, 0.1)
```

```
#OUT> [1] 1.5 1.6 1.7 1.8 1.9 2.0 2.1 2.2 2.3 2.4 2.5 2.6 2.7 2.8 2.9 3.0 3.1 3.2 3.3 3.4 3.5 3.6 3.7 3.8
```

Another common operation to create a vector is rep(), which can repeat a single value a number of times.

```
rep("A", times = 10)
```

```
#OUT> [1] "A" "A" "A" "A" "A" "A" "A" "A" "A" "A"
```

The `rep()` function can be used to repeat a vector some number of times.

```
rep(x, times = 3)
```

```
#OUT> [1] 1 3 5 7 8 9 1 3 5 7 8 9 1 3 5 7 8 9
```

We have now seen four different ways to create vectors:

- `c()`
- `:`
- `seq()`
- `rep()`

So far we have mostly used them in isolation, but they are often used together.

```
c(x, rep(seq(1, 9, 2), 3), c(1, 2, 3), 42, 2:4)
```

```
#OUT> [1] 1 3 5 7 8 9 1 3 5 7 9 1 3 5 7 9 1 3 5 7 9 1 2 3 42 2 3 4
```

The length of a vector can be obtained with the `length()` function.

```
length(x)
```

```
#OUT> [1] 6
```

```
length(y)
```

```
#OUT> [1] 100
```

Let's try this out! **Your turn:**

### 1.9.1.1 Task 1

1. Create a vector of five ones, i.e. `[1,1,1,1,1]`
2. Notice that the colon operator `a:b` is just short for construct a sequence from a to b. Create a vector the counts down from 10 to 0, i.e. it looks like `[10,9,8,7,6,5,4,3,2,1,0]`!
3. the `rep` function takes additional arguments `times` (as above), and each, which tells you how often each element should be repeated (as opposed to the entire input vector). Use `rep` to create a vector that looks like this: `[1 1 1 2 2 2 3 3 3 1 1 1 2 2 2 3 3 3]`

### 1.9.1.2 Subsetting

To subset a vector, i.e. to choose only some elements of it, we use square brackets, `[]`. Here we see that `x[1]` returns the first element, and `x[3]` returns the third element:

```
x
#OUT> [1] 1 3 5 7 8 9
```

```
x[1]
```

```
#OUT> [1] 1
```

```
x[3]
```

```
#OUT> [1] 5
```

We can also exclude certain indexes, in this case the second element.

```
x[-2]
```

```
#OUT> [1] 1 5 7 8 9
```

Lastly we see that we can subset based on a vector of indices.

```
x[1:3]
```

```
#OUT> [1] 1 3 5
```

```
x[c(1,3,4)]
```

```
#OUT> [1] 1 5 7
```

All of the above are subsetting a vector using a vector of indexes. (Remember a single number is still a vector.) We could instead use a vector of logical values.

```
z = c(TRUE, TRUE, FALSE, TRUE, TRUE, FALSE)
z
```

```
#OUT> [1] TRUE TRUE FALSE TRUE TRUE FALSE
```

```
x[z]
```

```
#OUT> [1] 1 3 7 8
```

R is able to perform many operations on vectors and scalars alike:

```
x = 1:10 # a vector
x + 1    # add a scalar
```

```
#OUT> [1] 2 3 4 5 6 7 8 9 10 11
```

```
2 * x    # multiply all elements by 2
```

```
#OUT> [1] 2 4 6 8 10 12 14 16 18 20
```

```

2 ^ x      # take 2 to the x as exponents

#OUT> [1]      2      4      8     16     32     64    128    256    512 1024

sqrt(x)     # compute the square root of all elements in x

#OUT> [1] 1.000000 1.414214 1.732051 2.000000 2.236068 2.449490 2.645751 2.828427 3.000000

log(x)      # take the natural log of all elements in x

#OUT> [1] 0.0000000 0.6931472 1.0986123 1.3862944 1.6094379 1.7917595 1.9459101 2.0794415

x + 2*x     # add vector x to vector 2x

#OUT> [1]  3  6  9 12 15 18 21 24 27 30

```

We see that when a function like `log()` is called on a vector `x`, a vector is returned which has applied the function to each element of the vector `x`.

## 1.9.2 Logical Operators

	Operator	Summary	Example	Result
<code>x &lt; y</code>		x less than y	<code>3 &lt; 42</code>	TRUE
<code>x &gt; y</code>		x greater than y	<code>3 &gt; 42</code>	FALSE
<code>x &lt;= y</code>		x less than or equal to y	<code>3 &lt;= 42</code>	TRUE
<code>x &gt;= y</code>		x greater than or equal to y	<code>3 &gt;= 42</code>	FALSE
<code>x == y</code>		xequal to y	<code>3 == 42</code>	FALSE
<code>x != y</code>		x not equal to y	<code>3 != 42</code>	TRUE
<code>!x</code>		not x	<code>!(3 &gt; 42)</code>	TRUE
<code>x   y</code>		x or y	<code>(3 &gt; 42)   TRUE</code>	TRUE
<code>x &amp; y</code>		x and y	<code>(3 &lt; 4) &amp; (42 &gt; 13)</code>	TRUE

In R, logical operators also work on vectors:

```

x = c(1, 3, 5, 7, 8, 9)

x > 3

#OUT> [1] FALSE FALSE  TRUE  TRUE  TRUE  TRUE

x < 3

#OUT> [1]  TRUE FALSE FALSE FALSE FALSE FALSE

x == 3

#OUT> [1] FALSE  TRUE FALSE FALSE FALSE FALSE

```

```
x != 3

#OUT> [1] TRUE FALSE TRUE TRUE TRUE TRUE
x == 3 & x != 3
```

```
#OUT> [1] FALSE FALSE FALSE FALSE FALSE FALSE
x == 3 | x != 3
```

```
#OUT> [1] TRUE TRUE TRUE TRUE TRUE TRUE
```

This is quite useful for subsetting.

```
x[x > 3]

#OUT> [1] 5 7 8 9
x[x != 3]
```

```
#OUT> [1] 1 5 7 8 9
sum(x > 3)
```

```
#OUT> [1] 4
as.numeric(x > 3)
```

```
#OUT> [1] 0 0 1 1 1 1
```

Here we saw that using the `sum()` function on a vector of logical `TRUE` and `FALSE` values that is the result of `x > 3` results in a numeric result: you just counted for how many elements of `x`, the condition `> 3` is `TRUE`. During the call to `sum()`, R is first automatically coercing the logical to numeric where `TRUE` is 1 and `FALSE` is 0. This coercion from logical to numeric happens for most mathematical operations.

```
# which(condition of x) returns true/false
# each index of x where condition is true
which(x > 3)
```

```
#OUT> [1] 3 4 5 6
x[which(x > 3)]
```

```
#OUT> [1] 5 7 8 9
max(x)
```

```
#OUT> [1] 9
which(x == max(x))
```

```
#OUT> [1] 6
```

```
which.max(x)
```

```
#OUT> [1] 6
```

### 1.9.2.1 Task 2

1. Create a vector filled with 10 numbers drawn from the uniform distribution (hint: use function `runif`) and store them in `x`.
2. Using logical subsetting as above, get all the elements of `x` which are larger than 0.5, and store them in `y`.
3. using the function `which`, store the indices of all the elements of `x` which are larger than 0.5 in `iy`.
4. Check that `y` and `x[iy]` are identical.

## 1.9.3 Matrices

R can also be used for **matrix** calculations. Matrices have rows and columns containing a single data type. In a matrix, the order of rows and columns is important. (This is not true of data frames, which we will see later.)

Matrices can be created using the `matrix` function.

```
x = 1:9
x
```

```
#OUT> [1] 1 2 3 4 5 6 7 8 9
```

```
X = matrix(x, nrow = 3, ncol = 3)
X
```

```
#OUT>      [,1] [,2] [,3]
#OUT> [1,]    1    4    7
#OUT> [2,]    2    5    8
#OUT> [3,]    3    6    9
```

Notice here that R is case sensitive (`x` vs `X`).

By default the `matrix` function fills your data into the matrix column by column. But we can also tell R to fill rows instead:

```
Y = matrix(x, nrow = 3, ncol = 3, byrow = TRUE)
Y
```

```
#OUT>      [,1] [,2] [,3]
#OUT> [1,]    1    2    3
#OUT> [2,]    4    5    6
#OUT> [3,]    7    8    9
```

We can also create a matrix of a specified dimension where every element is the same, in this case 0.

```
Z = matrix(0, 2, 4)
Z
```

```
#OUT>      [,1] [,2] [,3] [,4]
#OUT> [1,]    0    0    0    0
#OUT> [2,]    0    0    0    0
```

Like vectors, matrices can be subsetted using square brackets, `[]`. However, since matrices are two-dimensional, we need to specify both a row and a column when subsetting.

```
X
```

```
#OUT>      [,1] [,2] [,3]
#OUT> [1,]    1    4    7
#OUT> [2,]    2    5    8
#OUT> [3,]    3    6    9
```

```
X[1, 2]
```

```
#OUT> [1] 4
```

Here we accessed the element in the first row and the second column. We could also subset an entire row or column.

```
X[1, ]
```

```
#OUT> [1] 1 4 7
```

```
X[, 2]
```

```
#OUT> [1] 4 5 6
```

We can also use vectors to subset more than one row or column at a time. Here we subset to the first and third column of the second row:

```
X[2, c(1, 3)]
```

```
#OUT> [1] 2 8
```

Matrices can also be created by combining vectors as columns, using `cbind`, or combining vectors as rows, using `rbind`.

```
x = 1:9
rev(x)
```

```
#OUT> [1] 9 8 7 6 5 4 3 2 1
```

```
rep(1, 9)
```

```
#OUT> [1] 1 1 1 1 1 1 1 1 1
```

```
rbind(x, rev(x), rep(1, 9))
```

```
#OUT>  [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9]
#OUT> x    1    2    3    4    5    6    7    8    9
#OUT>      9    8    7    6    5    4    3    2    1
#OUT>      1    1    1    1    1    1    1    1    1
```

```
cbind(col_1 = x, col_2 = rev(x), col_3 = rep(1, 9))
```

```
#OUT>      col_1 col_2 col_3
#OUT> [1,]      1      9      1
#OUT> [2,]      2      8      1
#OUT> [3,]      3      7      1
#OUT> [4,]      4      6      1
#OUT> [5,]      5      5      1
#OUT> [6,]      6      4      1
#OUT> [7,]      7      3      1
#OUT> [8,]      8      2      1
#OUT> [9,]      9      1      1
```

When using `rbind` and `cbind` you can specify “argument” names that will be used as column names.

R can then be used to perform matrix calculations.

```
x = 1:9
y = 9:1
X = matrix(x, 3, 3)
Y = matrix(y, 3, 3)
X
```

```
#OUT>      [,1] [,2] [,3]
#OUT> [1,]      1      4      7
#OUT> [2,]      2      5      8
#OUT> [3,]      3      6      9
Y
```

```
#OUT>      [,1] [,2] [,3]
#OUT> [1,]      9      6      3
#OUT> [2,]      8      5      2
#OUT> [3,]      7      4      1
X + Y
```

```
#OUT>      [,1] [,2] [,3]
#OUT> [1,]     10     10     10
#OUT> [2,]     10     10     10
#OUT> [3,]     10     10     10
```



```
X - Y
```

```
#OUT>      [,1] [,2] [,3]
#OUT> [1,]   -8  -2   4
#OUT> [2,]   -6   0   6
#OUT> [3,]   -4   2   8
```

```
X * Y
```

```
#OUT>      [,1] [,2] [,3]
#OUT> [1,]     9  24  21
#OUT> [2,]    16  25  16
#OUT> [3,]    21  24   9
```

```
X / Y
```

```
#OUT>      [,1]      [,2]      [,3]
#OUT> [1,] 0.1111111 0.6666667 2.333333
#OUT> [2,] 0.2500000 1.0000000 4.000000
#OUT> [3,] 0.4285714 1.5000000 9.000000
```

Note that `X * Y` is **not** matrix multiplication. It is element by element multiplication. (Same for `X / Y`). Matrix multiplication uses `%*%`. Other matrix functions include `t()` which gives the transpose of a matrix and `solve()` which returns the inverse of a square matrix if it is invertible.

```
X %*% Y
```

```
#OUT>      [,1] [,2] [,3]
#OUT> [1,]    90   54   18
#OUT> [2,]   114   69   24
#OUT> [3,]   138   84   30
```

```
t(X)
```

```
#OUT>      [,1] [,2] [,3]
#OUT> [1,]     1     2     3
#OUT> [2,]     4     5     6
#OUT> [3,]     7     8     9
```

### 1.9.4 Arrays

A vector is a one-dimensional array. A matrix is a two-dimensional array. In R you can create arrays of arbitrary dimensionality N. Here is how:

```
d = 1:16
d3 = array(data = d, dim = c(4,2,2))
d4 = array(data = d, dim = c(4,2,2,3)) # will recycle 1:16
d3
```

```
#OUT> , , 1
#OUT>
#OUT>      [,1] [,2]
#OUT> [1,]    1    5
#OUT> [2,]    2    6
#OUT> [3,]    3    7
#OUT> [4,]    4    8
#OUT>
#OUT> , , 2
#OUT>
#OUT>      [,1] [,2]
#OUT> [1,]    9   13
#OUT> [2,]   10   14
#OUT> [3,]   11   15
#OUT> [4,]   12   16
```

You can see that `d3` are simply two (4,2) matrices laid on top of each other, as if there were two pages. Similarly, `d4` would have two pages, and another 3 registers in a fourth dimension. And so on. You can subset an array like you would a vector or a matrix, taking care to index each dimension:

```
d3[,1,1] # all elements from col 1, page 1
```

```
#OUT> [1] 1 2 3 4
```

```
d3[2:3, , ] # rows 2:3 from all pages
```

```
#OUT> , , 1
#OUT>
#OUT>      [,1] [,2]
#OUT> [1,]    2    6
#OUT> [2,]    3    7
#OUT>
#OUT> , , 2
#OUT>
#OUT>      [,1] [,2]
#OUT> [1,]   10   14
#OUT> [2,]   11   15
```

```
d3[2,2, ] # row 2, col 2 from both pages.
```

```
#OUT> [1] 6 14
```

#### 1.9.4.1 Task 3

1. Create a vector containing 1,2,3,4,5 called `v`.
2. Create a (2,5) matrix `m` containing the data 1,2,3,4,5,6,7,8,9,10. The first row should be 1,2,3,4,5.

3. Perform matrix multiplication of `m` with `v`. Use the command `%*%`. What dimension does the output have?
4. Why does `v %*% m` not work?

### 1.9.5 Lists

A list is a one-dimensional heterogeneous data structure. So it is indexed like a vector with a single integer value (or with a name), but each element can contain an element of any type. Lists are similar to a python or julia `Dict` object. Many R structures and outputs are lists themselves. Lists are extremely useful and versatile objects, so make sure you understand their usage:

```
# creation without fieldnames
list(42, "Hello", TRUE)

#OUT> [[1]]
#OUT> [1] 42
#OUT>
#OUT> [[2]]
#OUT> [1] "Hello"
#OUT>
#OUT> [[3]]
#OUT> [1] TRUE

# creation with fieldnames
ex_list = list(
  a = c(1, 2, 3, 4),
  b = TRUE,
  c = "Hello!",
  d = function(arg = 42) {print("Hello World!")},
  e = diag(5)
)
```

Lists can be subset using two syntaxes, the `$` operator, and square brackets `[]`. The `$` operator returns a named **element** of a list. The `[]` syntax returns a **list**, while the `[[ ]]` returns an **element** of a list.

- `ex_list[1]` returns a list contain the first element.
- `ex_list[[1]]` returns the first element of the list, in this case, a vector.

```
# subsetting
ex_list$a
```

```
#OUT>      [,1] [,2] [,3] [,4] [,5]
#OUT> [1,]    1    0    0    0    0
#OUT> [2,]    0    1    0    0    0
#OUT> [3,]    0    0    1    0    0
#OUT> [4,]    0    0    0    1    0
```

```
#OUT> [5,] 0 0 0 0 1
```

```
ex_list[1:2]
```

```
#OUT> $a
```

```
#OUT> [1] 1 2 3 4
```

```
#OUT>
```

```
#OUT> $b
```

```
#OUT> [1] TRUE
```

```
ex_list[1]
```

```
#OUT> $a
```

```
#OUT> [1] 1 2 3 4
```

```
ex_list[[1]]
```

```
#OUT> [1] 1 2 3 4
```

```
ex_list[c("e", "a")]
```

```
#OUT> $e
```

```
#OUT>      [,1] [,2] [,3] [,4] [,5]
```

```
#OUT> [1,] 1 0 0 0 0
```

```
#OUT> [2,] 0 1 0 0 0
```

```
#OUT> [3,] 0 0 1 0 0
```

```
#OUT> [4,] 0 0 0 1 0
```

```
#OUT> [5,] 0 0 0 0 1
```

```
#OUT>
```

```
#OUT> $a
```

```
#OUT> [1] 1 2 3 4
```

```
ex_list["e"]
```

```
#OUT> $e
```

```
#OUT>      [,1] [,2] [,3] [,4] [,5]
```

```
#OUT> [1,] 1 0 0 0 0
```

```
#OUT> [2,] 0 1 0 0 0
```

```
#OUT> [3,] 0 0 1 0 0
```

```
#OUT> [4,] 0 0 0 1 0
```

```
#OUT> [5,] 0 0 0 0 1
```

```
ex_list[["e"]]
```

```
#OUT>      [,1] [,2] [,3] [,4] [,5]
```

```
#OUT> [1,] 1 0 0 0 0
```

```
#OUT> [2,] 0 1 0 0 0
```

```
#OUT> [3,] 0 0 1 0 0
```

```
#OUT> [4,] 0 0 0 1 0
```

```
#OUT> [5,] 0 0 0 0 1
```

```
ex_list$d

#OUT> function(arg = 42) {print("Hello World!")}
ex_list$d(arg = 1)

#OUT> [1] "Hello World!"
```

#### 1.9.5.1 Task 4

1. Copy and paste the above code for `ex_list` into your R session. Remember that `list` can hold any kind of R object. Like...another list! So, create a new list `new_list` that has two fields: a first field called “this” with string content “is awesome”, and a second field called “ex\_list” that contains `ex_list`.
2. Accessing members is like in a plain list, just with several layers now. Get the element `c` from `ex_list` in `new_list`!
3. Compose a new string out of the first element in `new_list`, the element under label `this`. Use the function `paste` to print `R is awesome` to your screen.

## 1.10 Data Frames

We have previously seen vectors and matrices for storing data as we introduced R. We will now introduce a **data frame** which will be the most common way that we store and interact with data in this course. A `data.frame` is similar to a python `pandas.dataframe` or a julia `DataFrame`. (But the R version was the first! :-) )

```
example_data = data.frame(x = c(1, 3, 5, 7, 9, 1, 3, 5, 7, 9),
                          y = c(rep("Hello", 9), "Goodbye"),
                          z = rep(c(TRUE, FALSE), 5))
```

Unlike a matrix, which can be thought of as a vector rearranged into rows and columns, a data frame is not required to have the same data type for each element. A data frame is a **list** of vectors, and each vector has a name. So, each vector must contain the same data type, but the different vectors can store different data types. Note, however, that all vectors must have **the same length** (differently from a `list`)!

A **data.frame** is similar to a typical Spreadsheet. There are rows, and there are columns. A row is typically thought of as an observation, and each column is a certain variable, characteristic or feature of that observation.

Let's look at the data frame we just created above:

```
example_data
```

```
#OUT>   x      y      z
#OUT> 1  1   Hello  TRUE
#OUT> 2  3   Hello FALSE
#OUT> 3  5   Hello  TRUE
#OUT> 4  7   Hello FALSE
#OUT> 5  9   Hello  TRUE
#OUT> 6  1   Hello FALSE
#OUT> 7  3   Hello  TRUE
#OUT> 8  5   Hello FALSE
#OUT> 9  7   Hello  TRUE
#OUT> 10 9 Goodbye FALSE
```

Unlike a list, which has more flexibility, the elements of a data frame must all be vectors. Again, we access any given column with the `$` operator:

```
example_data$x
```

```
#OUT> [1] 1 3 5 7 9 1 3 5 7 9
all.equal(length(example_data$x),
           length(example_data$y),
           length(example_data$z))
```

```
#OUT> [1] TRUE
```

```
str(example_data)
```

```
#OUT> 'data.frame': 10 obs. of 3 variables:
#OUT> $ x: num 1 3 5 7 9 1 3 5 7 9
#OUT> $ y: Factor w/ 2 levels "Goodbye","Hello": 2 2 2 2 2 2 2 2 2 1
#OUT> $ z: logi TRUE FALSE TRUE FALSE TRUE FALSE ...
```

```
nrow(example_data)
```

```
#OUT> [1] 10
```

```
ncol(example_data)
```

```
#OUT> [1] 3
```

```
dim(example_data)
```

```
#OUT> [1] 10 3
```

```
names(example_data)
```

```
#OUT> [1] "x" "y" "z"
```

### 1.10.1 Working with data.frames

The `data.frame()` function above is one way to create a data frame. We can also import data from various file types in into R, as well as use data stored in packages.

To read this data back into R, we will use the built-in function `read.csv`:

```
path = system.file(package="ScPoEconometrics","datasets","example-data.csv")
example_data_from_disk = read.csv(path)
```

This particular line of code assumes that you installed the associated R package to this book, hence you have this dataset stored on your computer at `system.file(package = "ScPoEconometrics", "datasets", "example-data.csv")`.

```
example_data_from_disk
```

```
#OUT>      x      y      z
#OUT> 1  1 Hello TRUE
#OUT> 2  3 Hello FALSE
#OUT> 3  5 Hello TRUE
#OUT> 4  7 Hello FALSE
#OUT> 5  9 Hello TRUE
#OUT> 6  1 Hello FALSE
#OUT> 7  3 Hello TRUE
#OUT> 8  5 Hello FALSE
#OUT> 9  7 Hello TRUE
#OUT> 10 9 Goodbye FALSE
```

When using data, there are three things we would generally like to do:

- Look at the raw data.
- Understand the data. (Where did it come from? What are the variables? Etc.)
- Visualize the data.

To look at data in a `data.frame`, we have two useful commands: `head()` and `str()`.

```
# we are working with the built-in mtcars dataset:
mtcars
```

```
#OUT>
#OUT>      mpg cyl  disp  hp drat   wt  qsec vs am gear carb
#OUT> Mazda RX4      21.0   6 160.0 110 3.90 2.620 16.46 0 1   4   4
#OUT> Mazda RX4 Wag   21.0   6 160.0 110 3.90 2.875 17.02 0 1   4   4
#OUT> Datsun 710      22.8   4 108.0  93 3.85 2.320 18.61 1 1   4   1
#OUT> Hornet 4 Drive   21.4   6 258.0 110 3.08 3.215 19.44 1 0   3   1
#OUT> Hornet Sportabout 18.7   8 360.0 175 3.15 3.440 17.02 0 0   3   2
#OUT> Valiant         18.1   6 225.0 105 2.76 3.460 20.22 1 0   3   1
#OUT> Duster 360      14.3   8 360.0 245 3.21 3.570 15.84 0 0   3   4
```

```

#OUT> Merc 240D      24.4  4 146.7 62 3.69 3.190 20.00 1 0  4  2
#OUT> Merc 230      22.8  4 140.8 95 3.92 3.150 22.90 1 0  4  2
#OUT> Merc 280      19.2  6 167.6 123 3.92 3.440 18.30 1 0  4  4
#OUT> Merc 280C     17.8  6 167.6 123 3.92 3.440 18.90 1 0  4  4
#OUT> Merc 450SE    16.4  8 275.8 180 3.07 4.070 17.40 0 0  3  3
#OUT> Merc 450SL    17.3  8 275.8 180 3.07 3.730 17.60 0 0  3  3
#OUT> Merc 450SLC   15.2  8 275.8 180 3.07 3.780 18.00 0 0  3  3
#OUT> Cadillac Fleetwood 10.4  8 472.0 205 2.93 5.250 17.98 0 0  3  4
#OUT> Lincoln Continental 10.4  8 460.0 215 3.00 5.424 17.82 0 0  3  4
#OUT> Chrysler Imperial 14.7  8 440.0 230 3.23 5.345 17.42 0 0  3  4
#OUT> Fiat 128      32.4  4  78.7 66 4.08 2.200 19.47 1 1  4  1
#OUT> Honda Civic   30.4  4  75.7 52 4.93 1.615 18.52 1 1  4  2
#OUT> Toyota Corolla 33.9  4  71.1 65 4.22 1.835 19.90 1 1  4  1
#OUT> Toyota Corona 21.5  4 120.1 97 3.70 2.465 20.01 1 0  3  1
#OUT> Dodge Challenger 15.5  8 318.0 150 2.76 3.520 16.87 0 0  3  2
#OUT> AMC Javelin   15.2  8 304.0 150 3.15 3.435 17.30 0 0  3  2
#OUT> Camaro Z28    13.3  8 350.0 245 3.73 3.840 15.41 0 0  3  4
#OUT> Pontiac Firebird 19.2  8 400.0 175 3.08 3.845 17.05 0 0  3  2
#OUT> Fiat X1-9     27.3  4  79.0 66 4.08 1.935 18.90 1 1  4  1
#OUT> Porsche 914-2 26.0  4 120.3 91 4.43 2.140 16.70 0 1  5  2
#OUT> Lotus Europa  30.4  4  95.1 113 3.77 1.513 16.90 1 1  5  2
#OUT> Ford Pantera L 15.8  8 351.0 264 4.22 3.170 14.50 0 1  5  4
#OUT> Ferrari Dino  19.7  6 145.0 175 3.62 2.770 15.50 0 1  5  6
#OUT> Maserati Bora  15.0  8 301.0 335 3.54 3.570 14.60 0 1  5  8
#OUT> Volvo 142E    21.4  4 121.0 109 4.11 2.780 18.60 1 1  4  2

```

You can see that this prints the entire data.frame to screen. The function `head()` will display the first `n` observations of the data frame.

```
head(mtcars,n=2)
```

```

#OUT>           mpg cyl disp  hp drat   wt  qsec vs am gear carb
#OUT> Mazda RX4      21   6  160 110  3.9 2.620 16.46 0  1   4   4
#OUT> Mazda RX4 Wag  21   6  160 110  3.9 2.875 17.02 0  1   4   4

```

```
head(mtcars) # default
```

```

#OUT>           mpg cyl disp  hp drat   wt  qsec vs am gear carb
#OUT> Mazda RX4      21.0  6  160 110  3.90 2.620 16.46 0  1   4   4
#OUT> Mazda RX4 Wag  21.0  6  160 110  3.90 2.875 17.02 0  1   4   4
#OUT> Datsun 710      22.8  4  108  93  3.85 2.320 18.61 1  1   4   1
#OUT> Hornet 4 Drive  21.4  6  258 110  3.08 3.215 19.44 1  0   3   1
#OUT> Hornet Sportabout 18.7  8  360 175  3.15 3.440 17.02 0  0   3   2
#OUT> Valiant         18.1  6  225 105  2.76 3.460 20.22 1  0   3   1

```

The function `str()` will display the “structure” of the data frame. It will display the number of **observations** and **variables**, list the variables, give the type of each variable, and show some elements of each variable. This information can



also be found in the “Environment” window in RStudio.

```
str(mtcars)
```

```
#OUT> 'data.frame': 32 obs. of 11 variables:
#OUT> $ mpg : num 21 21 22.8 21.4 18.7 18.1 14.3 24.4 22.8 19.2 ...
#OUT> $ cyl : num 6 6 4 6 8 6 8 4 4 6 ...
#OUT> $ disp: num 160 160 108 258 360 ...
#OUT> $ hp : num 110 110 93 110 175 105 245 62 95 123 ...
#OUT> $ drat: num 3.9 3.9 3.85 3.08 3.15 2.76 3.21 3.69 3.92 3.92 ...
#OUT> $ wt : num 2.62 2.88 2.32 3.21 3.44 ...
#OUT> $ qsec: num 16.5 17 18.6 19.4 17 ...
#OUT> $ vs : num 0 0 1 1 0 1 0 1 1 1 ...
#OUT> $ am : num 1 1 1 0 0 0 0 0 0 0 ...
#OUT> $ gear: num 4 4 4 3 3 3 3 4 4 4 ...
#OUT> $ carb: num 4 4 1 1 2 1 4 2 2 4 ...
```

In this dataset an observation is for a particular model of a car, and the variables describe attributes of the car, for example its fuel efficiency, or its weight.

To understand more about the data set, we use the `?` operator to pull up the documentation for the data.

```
?mtcars
```

R has a number of functions for quickly working with and extracting basic information from data frames. To quickly obtain a vector of the variable names, we use the `names()` function.

```
names(mtcars)
```

```
#OUT> [1] "mpg" "cyl" "disp" "hp" "drat" "wt" "qsec" "vs" "am" "gear" "carb"
```

To access one of the variables **as a vector**, we use the `$` operator.

```
mtcars$mpg
```

```
#OUT> [1] 21.0 21.0 22.8 21.4 18.7 18.1 14.3 24.4 22.8 19.2 17.8 16.4 17.3 15.2 10.4 10.4 14.7 32.4 30.4
#OUT> [26] 27.3 26.0 30.4 15.8 19.7 15.0 21.4
```

```
mtcars$wt
```

```
#OUT> [1] 2.620 2.875 2.320 3.215 3.440 3.460 3.570 3.190 3.150 3.440 3.440 4.070 3.730 3.780 5.250 5.250
#OUT> [22] 3.520 3.435 3.840 3.845 1.935 2.140 1.513 3.170 2.770 3.570 2.780
```

We can use the `dim()`, `nrow()` and `ncol()` functions to obtain information about the dimension of the data frame.

```
dim(mtcars)
```

```
#OUT> [1] 32 11
```

```
nrow(mtcars)
```

```
#OUT> [1] 32
```

```
ncol(mtcars)
```

```
#OUT> [1] 11
```

Here `nrow()` is also the number of observations, which in most cases is the sample size.

Subsetting data frames can work much like subsetting matrices using square brackets, `[ , ]`. Here, we find vehicles with `mpg` over 25 miles per gallon and only display columns `cyl`, `disp` and `wt`.

```
# mpg[row condition, col condition]
mtcars[mtcars$mpg > 20, c("cyl", "disp", "wt")]
```

```
#OUT>
#OUT>      cyl  disp  wt
#OUT> Mazda RX4      6 160.0 2.620
#OUT> Mazda RX4 Wag  6 160.0 2.875
#OUT> Datsun 710      4 108.0 2.320
#OUT> Hornet 4 Drive  6 258.0 3.215
#OUT> Merc 240D      4 146.7 3.190
#OUT> Merc 230       4 140.8 3.150
#OUT> Fiat 128       4  78.7 2.200
#OUT> Honda Civic    4  75.7 1.615
#OUT> Toyota Corolla 4  71.1 1.835
#OUT> Toyota Corona  4 120.1 2.465
#OUT> Fiat X1-9      4  79.0 1.935
#OUT> Porsche 914-2  4 120.3 2.140
#OUT> Lotus Europa   4  95.1 1.513
#OUT> Volvo 142E     4 121.0 2.780
```

An alternative would be to use the `subset()` function, which has a much more readable syntax.

```
subset(mtcars, subset = mpg > 25, select = c("cyl", "disp", "wt"))
```

#### 1.10.1.1 Task 5

1. How many observations are there in `mtcars`?
2. How many variables?
3. What is the average value of `mpg`?
4. What is the average value of `mpg` for cars with more than 4 cylinders, i.e. with `cyl>4`?

## 1.11 Programming Basics

In this section we illustrate some general concepts related to programming.

### 1.11.1 Variables

We encountered the term variable already several times, but mainly in the context of a column of a `data.frame`. In programming, a variable denotes an object. Another way to say it is that a variable is a name or a label for something:

```
x = 1
y = "roses"
z = function(x){sqrt(x)}
```

Here `x` refers to the value 1, `y` holds the string “roses”, and `z` is the name of a function that computes  $\sqrt{x}$ . Notice that the argument `x` of the function is different from the `x` we just defined. It is **local** to the function:

```
x
```

```
#OUT> [1] 1
```

```
z(9)
```

```
#OUT> [1] 3
```

### 1.11.2 Control Flow

Control Flow relates to ways in which you can adapt your code to different circumstances. Based on a **condition** being **TRUE**, your program will do one thing, as opposed to another thing. This is most widely known as an **if/else** statement. In R, the if/else syntax is:

```
if (condition = TRUE) {
  some R code
} else {
  some other R code
}
```

For example,

```
x = 1
y = 3
if (x > y) { # test if x > y
  # if TRUE
  z = x * y
  print("x is larger than y")
}
```

```

} else {
  # if FALSE
  z = x + 5 * y
  print("x is less than or equal to y")
}

```

```
#OUT> [1] "x is less than or equal to y"
```

```
z
```

```
#OUT> [1] 16
```

### 1.11.3 Loops

Loops are a very important programming construct. As the name suggests, in a loop, the programming repeatedly loops over a set of instructions, until some condition tells it to stop. A very powerful, yet simple, construction is that the program can count how many steps it has done already - which may be important to know for many algorithms. The syntax of a **for** loop (there are others), is

```

for (ix in 1:10){  # does not have to be 1:10!
  # loop body: gets executed each time
  # the value of ix changes with each iteration
}

```

For example, consider this simple **for** loop, which will simply print the value of the iterator (called *i* in this case) to screen:

```

for (i in 1:5){
  print(i)
}

```

```
#OUT> [1] 1
```

```
#OUT> [1] 2
```

```
#OUT> [1] 3
```

```
#OUT> [1] 4
```

```
#OUT> [1] 5
```

Notice that instead of 1:5, we could have any kind of iterable collection:

```

for (i in c("mangos","bananas","apples")){
  print(paste("I love",i))  # the paste function pastes together strings
}

```

```
#OUT> [1] "I love mangos"
```

```
#OUT> [1] "I love bananas"
```

```
#OUT> [1] "I love apples"
```

We often also see nested loops, which are just what its name suggests:

```
for (i in 2:3){
  # first nest: for each i
  for (j in c("mangos","bananas","apples")){
    # second nest: for each j
    print(paste("Can I get",i,j,"please?"))
  }
}
```

```
#OUT> [1] "Can I get 2 mangos please?"
#OUT> [1] "Can I get 2 bananas please?"
#OUT> [1] "Can I get 2 apples please?"
#OUT> [1] "Can I get 3 mangos please?"
#OUT> [1] "Can I get 3 bananas please?"
#OUT> [1] "Can I get 3 apples please?"
```

The important thing to note here is that you can do calculations with the iterators while inside a loop.

#### 1.11.4 Functions

So far we have been using functions, but haven't actually discussed some of their details. A function is a set of instructions that R executes for us, much like those collected in a script file. The good thing is that functions are much more flexible than scripts, since they can depend on input arguments, which change the way the function behaves. Here is how to define a function:

```
function_name <- function(arg1,arg2=default_value){
  # function body
  # you do stuff with arg1 and arg2
  # you can have any number of arguments, with or without defaults
  # any valid `R` commands can be included here
  # the last line is returned
}
```

And here is a trivial example of a function definition:

```
hello <- function(your_name = "Lord Vader"){
  paste("You R most welcome,",your_name)
  # we could also write:
  # return(paste("You R most welcome,",your_name))
}

# we call the function by typing it's name with round brackets
hello()
```

```
#OUT> [1] "You R most welcome, Lord Vader"
```

You see that by not specifying the argument `your_name`, R reverts to the default value given. Try with your own name now!

Just typing the function name returns the actual definition to us, which is handy sometimes:

```
hello
```

```
#OUT> function(your_name = "Lord Vader"){  
#OUT>   paste("You R most welcome,",your_name)  
#OUT>   # we could also write:  
#OUT>   # return(paste("You R most welcome,",your_name))  
#OUT> }
```

It's instructive to consider that before we defined the function `hello` above, R did not know what to do, had you called `hello()`. The function did not exist! In this sense, we taught R a new trick. This feature to create new capabilities on top of a core language is one of the most powerful characteristics of programming languages. In general, it is good practice to split your code into several smaller functions, rather than one long script file. It makes your code more readable, and it is easier to track down mistakes.

#### 1.11.4.1 Task 6

1. Write a for loop that counts down from 10 to 1, printing the value of the iterator to the screen.
2. Modify that loop to write “i iterations to go” where `i` is the iterator
3. Modify that loop so that each iteration takes roughly one second. You can achieve that by adding the command `Sys.sleep(1)` below the line that prints “i iterations to go”.

## Chapter 2

# Working With Data

In this chapter we will first learn some basic concepts that help summarizing data. Then, we will tackle a real-world task and read, clean, and summarize data from the web.

### 2.1 Summary Statistics

R has built in functions for a large number of summary statistics. For numeric variables, we can summarize data by looking at their center and spread, for example.

```
# for the mpg dataset, we load:  
library(ggplot2)
```

#### Central Tendency

Suppose we want to know the mean and median of all the values stored in the `data.frame` column `mpg$cty`:

	Measure	R	Result
Mean	<code>mean(mpg\$cty)</code>		16.8589744
Median	<code>median(mpg\$cty)</code>		17

#### Spread

How do the values in that column vary? How far spread out are they?

	Measure	R	Result
Variance		<code>var(mpg\$cty)</code>	18.1130736
Standard Deviation		<code>sd(mpg\$cty)</code>	4.2559457
IQR		<code>IQR(mpg\$cty)</code>	5
Minimum		<code>min(mpg\$cty)</code>	9
Maximum		<code>max(mpg\$cty)</code>	35
Range		<code>range(mpg\$cty)</code>	9, 35

## Categorical

For categorical variables, counts and percentages can be used for summary.

```
table(mpg$drv)
```

```
#OUT>
#OUT>  4   f   r
#OUT> 103 106 25
```

```
table(mpg$drv) / nrow(mpg)
```

```
#OUT>
#OUT>          4          f          r
#OUT> 0.4401709 0.4529915 0.1068376
```

## 2.2 Plotting

Now that we have some data to work with, and we have learned about the data at the most basic level, our next tasks will be to visualize it. Often, a proper visualization can illuminate features of the data that can inform further analysis.

We will look at four methods of visualizing data by using the basic `plot` facilities built-in with R:

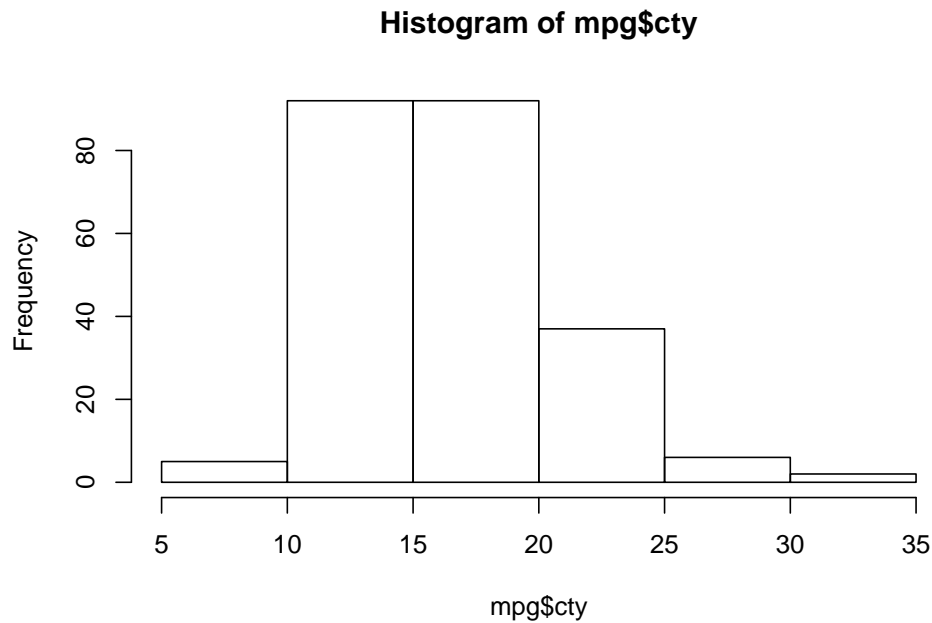
- Histograms
- Barplots
- Boxplots
- Scatterplots

### 2.2.1 Histograms

When visualizing a single numerical variable, a **histogram** is useful. It summarizes the distribution of values in a vector. In R you create one using the `hist()` function:

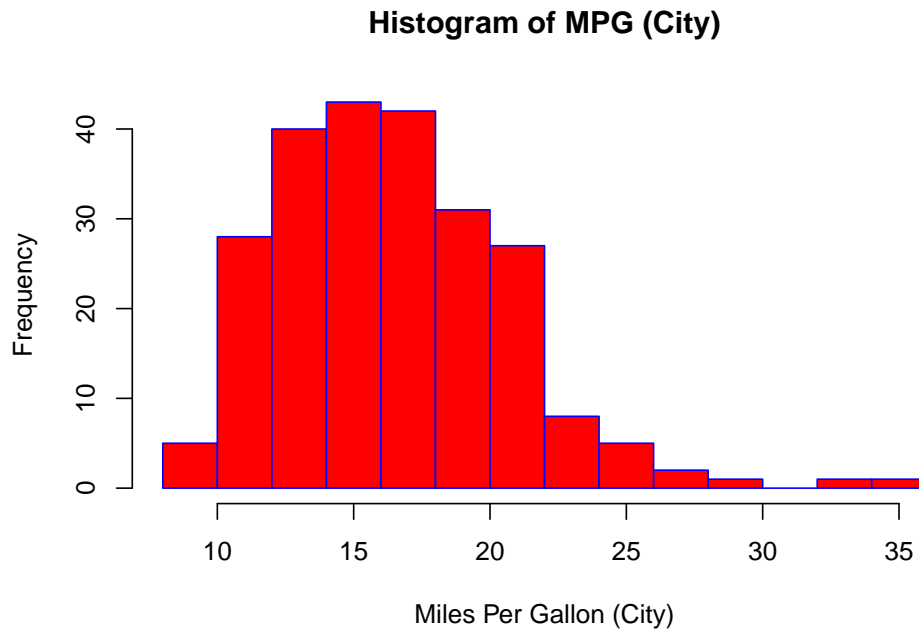


```
hist(mpg$cty)
```



The histogram function has a number of parameters which can be changed to make our plot look much nicer. Use the `?` operator to read the documentation for the `hist()` to see a full list of these parameters.

```
hist(mpg$cty,  
     xlab  = "Miles Per Gallon (City)",  
     main  = "Histogram of MPG (City)", # main title  
     breaks = 12, # how many breaks?  
     col   = "red",  
     border = "blue")
```

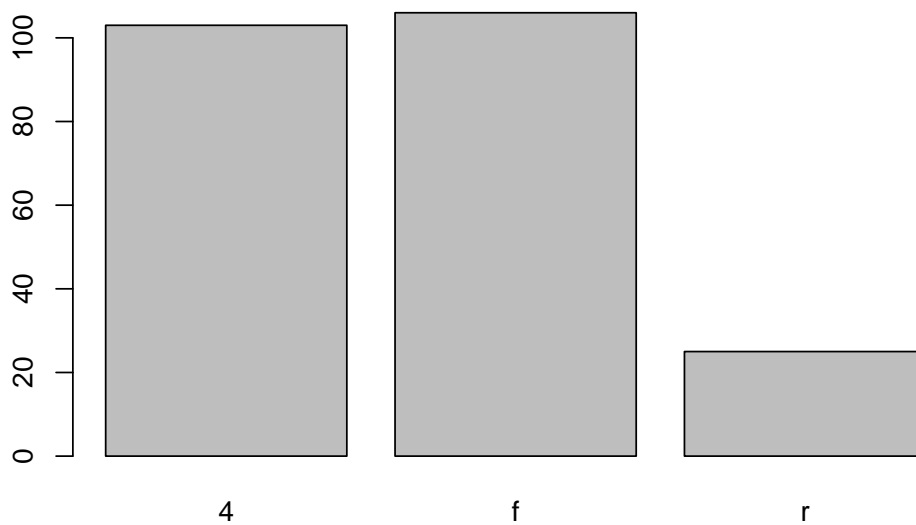


Importantly, you should always be sure to label your axes and give the plot a title. The argument `breaks` is specific to `hist()`. Entering an integer will give a suggestion to R for how many bars to use for the histogram. By default R will attempt to intelligently guess a good number of `breaks`, but as we can see here, it is sometimes useful to modify this yourself.

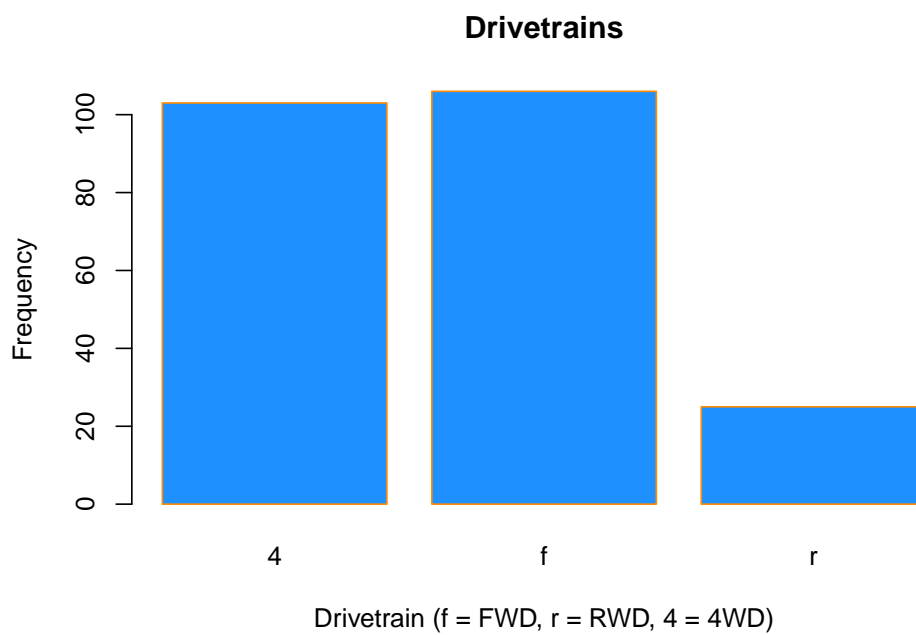
### 2.2.2 Barplots

Somewhat similar to a histogram, a barplot can provide a visual summary of a categorical variable, or a numeric variable with a finite number of values, like a ranking from 1 to 10.

```
barplot(table(mpg$drv))
```



```
barplot(table(mpg$drv),  
  xlab = "Drivetrain (f = FWD, r = RWD, 4 = 4WD)",  
  ylab = "Frequency",  
  main = "Drivetrains",  
  col = "dodgerblue",  
  border = "darkorange")
```



### 2.2.3 Boxplots

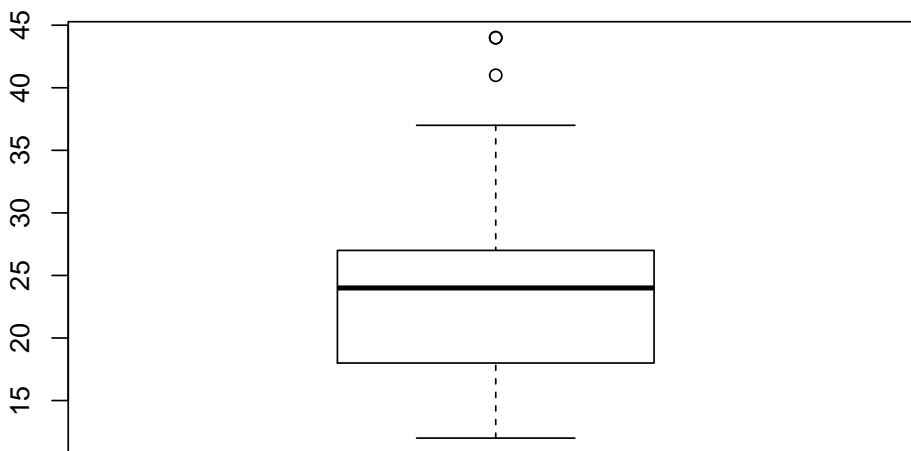
To visualize the relationship between a numerical and categorical variable, one could use a **boxplot**. In the `mpg` dataset, the `drv` variable takes a small, finite number of values. A car can only be front wheel drive, 4 wheel drive, or rear wheel drive.

```
unique(mpg$drv)
```

```
#OUT> [1] "f" "4" "r"
```

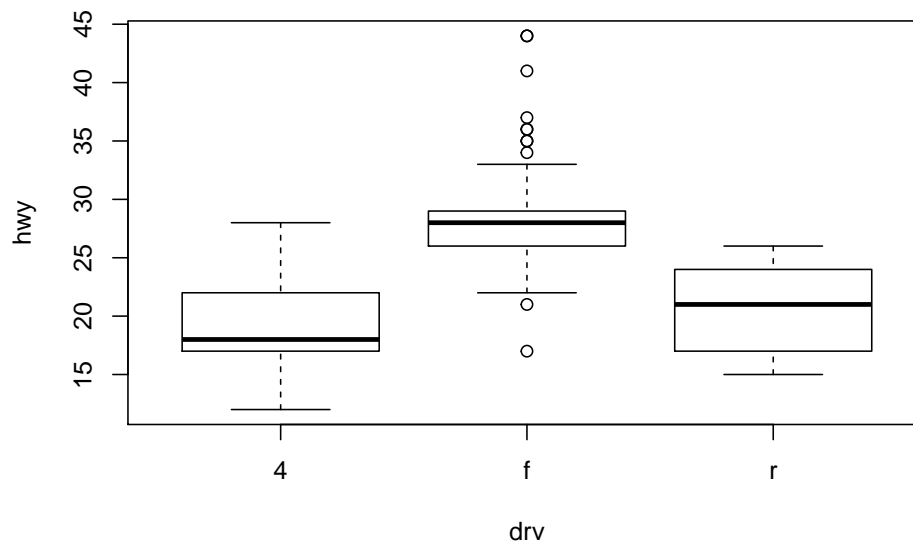
First note that we can use a single boxplot as an alternative to a histogram for visualizing a single numerical variable. To do so in R, we use the `boxplot()` function. The box shows the interquartile range, the solid line in the middle is the value of the median, the whiskers show 1.5 times the interquartile range, and the dots are outliers.

```
boxplot(mpg$hwy)
```



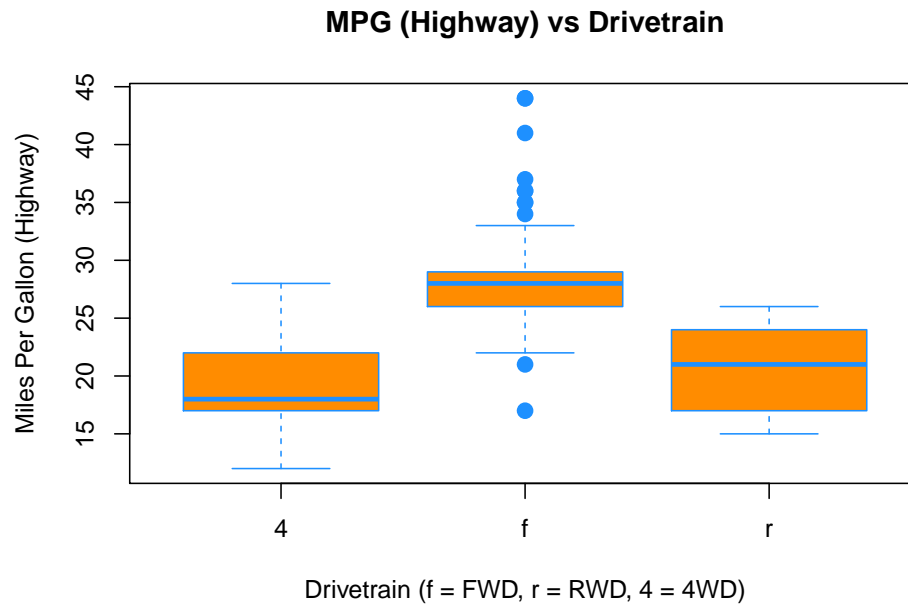
However, more often we will use boxplots to compare a numerical variable for different values of a categorical variable.

```
boxplot(hwy ~ drv, data = mpg)
```



Here used the `boxplot()` command to create side-by-side boxplots. However, since we are now dealing with two variables, the syntax has changed. The R syntax `hwy ~ drv, data = mpg` reads “Plot the `hwy` variable against the `drv` variable using the dataset `mpg`.” We see the use of a `~` (which specifies a formula) and also a `data =` argument. This will be a syntax that is common to many functions we will use in this course.

```
boxplot(hwy ~ drv, data = mpg,
  xlab = "Drivetrain (f = FWD, r = RWD, 4 = 4WD)",
  ylab = "Miles Per Gallon (Highway)",
  main = "MPG (Highway) vs Drivetrain",
  pch = 20,
  cex = 2,
  col = "darkorange",
  border = "dodgerblue")
```

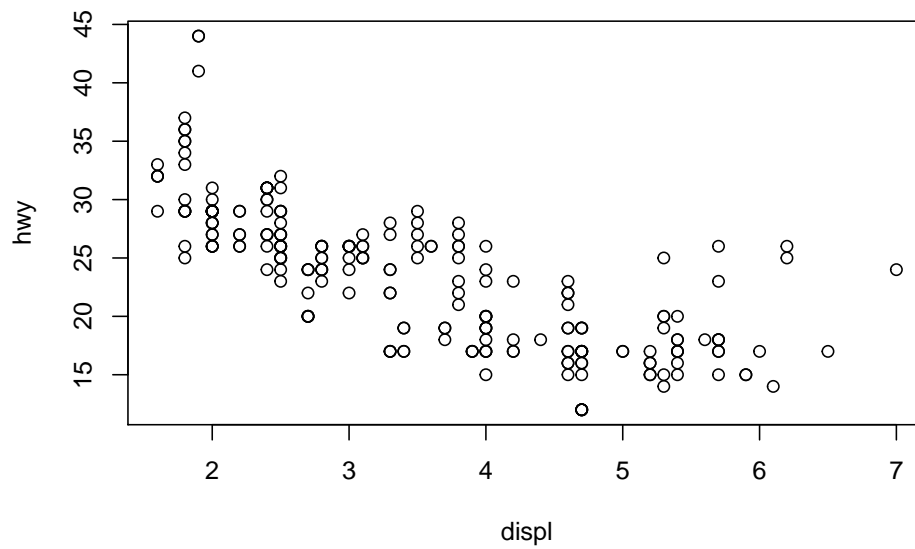


Again, `boxplot()` has a number of additional arguments which have the ability to make our plot more visually appealing.

## 2.2.4 Scatterplots

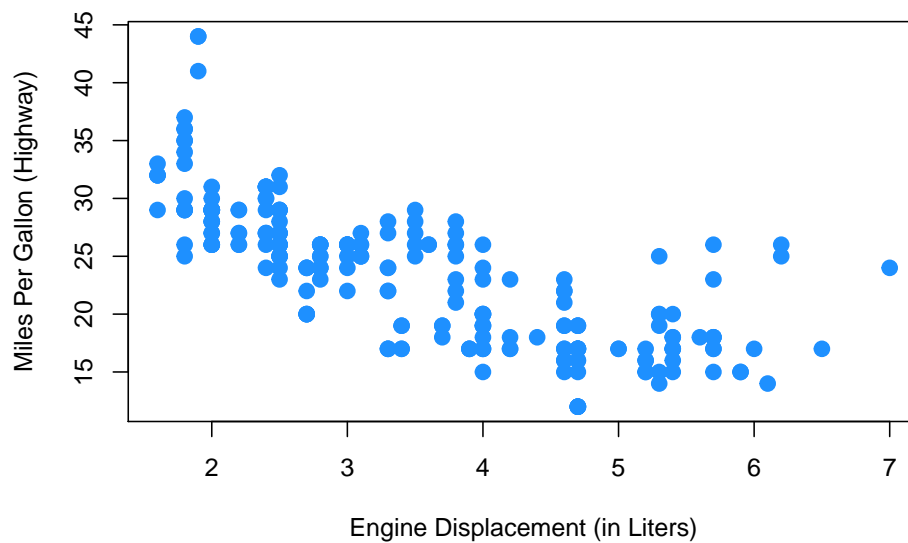
Lastly, to visualize the relationship between two numeric variables we will use a **scatterplot**. This can be done with the `plot()` function and the `~` syntax we just used with a boxplot. (The function `plot()` can also be used more generally; see the documentation for details.)

```
plot(hwy ~ displ, data = mpg)
```



```
plot(hwy ~ displ, data = mpg,  
     xlab = "Engine Displacement (in Liters)",  
     ylab = "Miles Per Gallon (Highway)",  
     main = "MPG (Highway) vs Engine Displacement",  
     pch = 20,  
     cex = 2,  
     col = "dodgerblue")
```

**MPG (Highway) vs Engine Displacement**



### 2.2.5 ggplot

All of the above plots could also have been generated using the `ggplot` function from the already loaded `ggplot2` package. Which function you use is up to you, but sometimes a plot is easier to build in base R (like in the `boxplot` example maybe), sometimes the other way around.

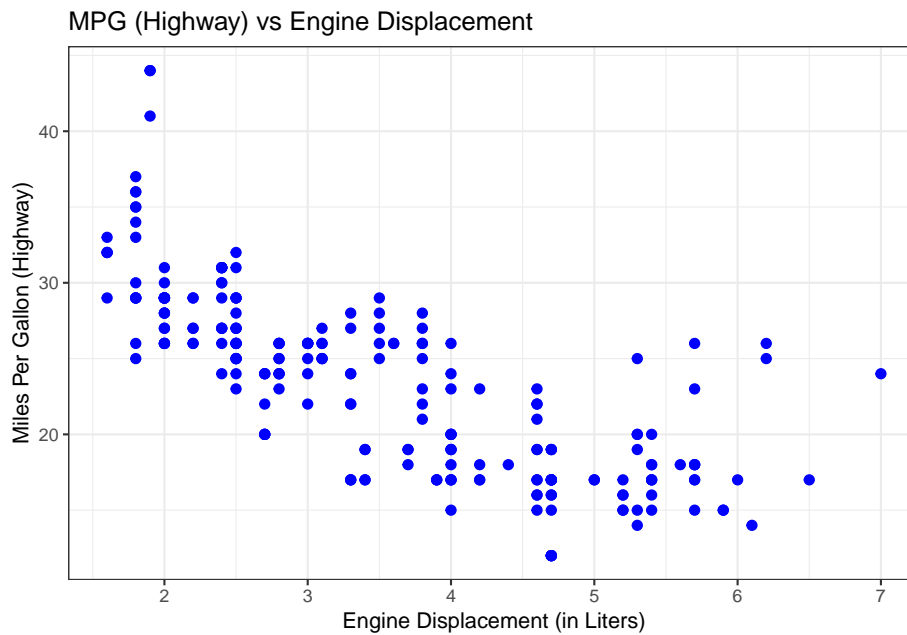
```
ggplot(data = mpg, mapping = aes(x=displ,y=hwy)) + geom_point()
```



`ggplot` is impossible to describe in brief terms, so please look at the package's website which provides excellent guidance. We will from time to time use `ggplot` in this book, so you could familiarize yourself with it. Let's quickly demonstrate how one could further customize that first plot:

```
ggplot(data = mpg, mapping = aes(x=displ,y=hwy)) + # ggplot() makes base plot
  geom_point(color="blue",size=2) + # how to show x and y?
  scale_y_continuous(name="Miles Per Gallon (Highway)") + # name of y axis
  scale_x_continuous(name="Engine Displacement (in Liters)") + # x axis
  theme_bw() + # change the background
  ggtitle("MPG (Highway) vs Engine Displacement") # add a title
```





If you want to see `ggplot` in action, you could start with this and then look at that very nice tutorial? It's fun!

## 2.3 Summarizing Two Variables

We often are interested in how two variables are related to each other. The core concepts here are covariance and correlation. Let's generate some data on  $x$  and  $y$  and plot them against each other:

Taking as example the data in this plot, the concepts covariance and correlation relate to the following type of question:

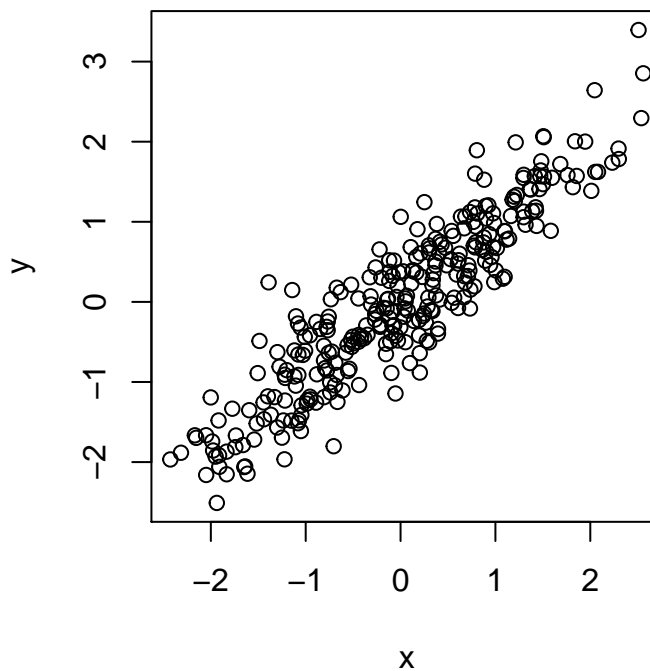
Given we observe value of something like  $x = 2$ , say, can we expect a high or a low value of  $y$ , on average? Something like  $y = 2$  or rather something like  $y = -2$ ?

The answer to this type of question can be addressed by computing the covariance of both variables:

```
cov(x,y)
```

```
#OUT> [1] 1.041195
```

Here, this gives a positive number, 1.04, indicating that as one variable lies above its average, the other one does as well. In other words, it indicates a **positive**

Figure 2.1: How are  $x$  and  $y$  related?

**relationship.** What is less clear, however, how to interpret the magnitude of 1.04. Is that a strong or a weak positive association?

In fact, we cannot tell. This is because the covariance is measured in the same units as the data, and those units often differ between both variables. There is a better measure available to us though, the **correlation**, which is obtained by standardizing each variable. By standardizing a variable  $x$  one means to divide  $x$  by its standard deviation  $\sigma_x$ :

$$z = \frac{x}{\sigma_x}$$

The correlation coefficient between  $x$  and  $y$ , commonly denoted  $r_{x,y}$ , is then defined as

$$r_{x,y} = \frac{\text{cov}(x,y)}{\sigma_x \sigma_y},$$

and we get rid of the units problem. In R, you can call directly

```
cor(x,y)
```

```
#OUT> [1] 0.9142495
```

Now this is better. Given that the correlation has to lie in  $[-1, 1]$ , a value of 0.91 is indicative of a rather strong positive relationship for the data in figure ??

Note that  $x, y$  being drawn from a continuous distribution (they are joint normally distributed) had no implication for covariance and correlation: We can compute those measures also for discrete random variables (like the throws of two dice, as you will see in one of our tutorials).

### 2.3.1 Visually estimating $\sigma$

Sometimes it is useful to estimate the standard deviation of some data without the help of a computer (for example during an exam ;-)). If  $x$  is approximately normally distributed, 95% of its observations will lie within a range of  $\bar{x} \pm$  two standard deviations of  $x$ . That is to say, four standard deviations of  $x$  cover 95% of its observations. Hence, a simple way to estimate the standard deviation for a variable is to look at the range of  $x$ , and simply divide that number by four.

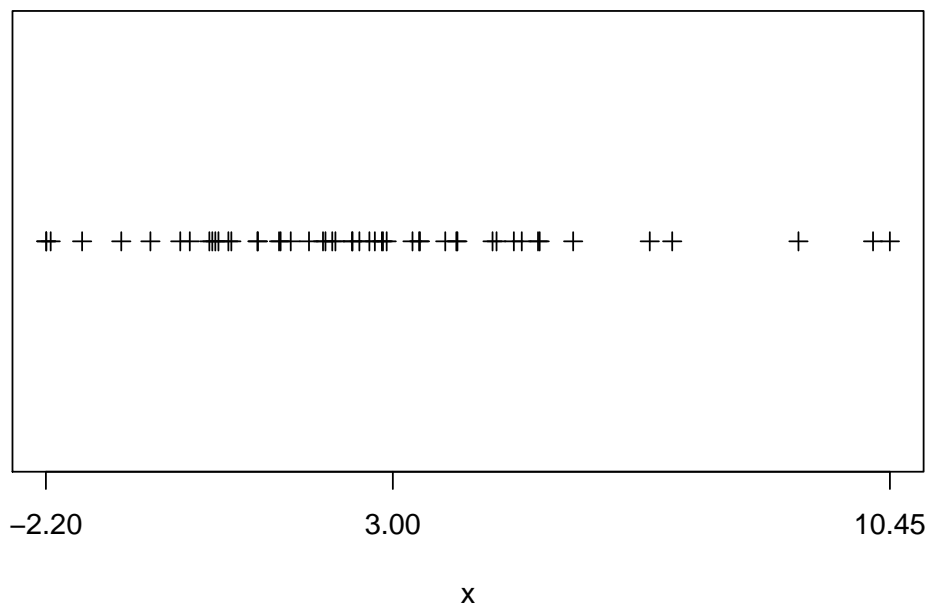


Figure 2.2: visual estimation on  $\sigma$ . The x-axis labels min and max as well as mean of  $x$ .

This is illustrated in figure ?. Here we see that  $\text{range}(x)/4$  gives 3.16 which compares favourably to the actual standard deviation 3.

## 2.4 The tidyverse

Hadley Wickham is the author of R packages `ggplot2` and also of `dplyr` (and also a myriad of others). With `ggplot2` he introduced what is called the grammar of graphics (hence, `gg`) to R. Grammar in the sense that there are **nouns** and **verbs** and a **syntax**, i.e. rules of how nouns and verbs are to be put together to construct an understandable sentence. He has extended the grammar idea into various other packages. The **tidyverse** package is a collection of those packages.

tidy data is data where:

- Each variable is a column
- Each observation is a row
- Each value is a cell

Fair enough, you might say, that is a regular spreadsheet. And you are right! However, data comes to us not tidy most of the times, and we first need to clean, or tidy, it up. Once it's in tidy format, we can use the tools in the **tidyverse** with great efficiency to analyse the data and stop worrying about which tool to use.

### 2.4.1 Reading .csv data in the tidy way

We could have used the `read_csv()` function from the `readr` package to read our example dataset from the previous chapter. The `readr` function `read_csv()` has a number of advantages over the built-in `read.csv`. For example, it is much faster reading larger data. It also uses the `tibble` package to read the data as a tibble. **A tibble is simply a data frame that prints with sanity.** Notice in the output below that we are given additional information such as dimension and variable type.

```
library(readr) # you need `install.packages("readr")` once!
path = system.file(package="ScPoEconometrics", "datasets", "example-data.csv")
example_data_from_disk = read_csv(path)
```

### 2.4.2 Tidy data.frames are tibbles

Let's grab some data from the `ggplot2` package:

```
data(mpg, package = "ggplot2") # load dataset `mpg` from `ggplot2` package
head(mpg, n = 10)
```

```
#OUT> # A tibble: 10 x 11
#OUT>   manufacturer model      displ year   cyl trans      drv    cty   hwy fl    class
#OUT>   <chr>         <chr>    <dbl> <int> <int> <chr>   <chr> <int> <int> <chr> <chr>
```

```
#OUT> 1 audi      a4      1.8 1999  4 auto(l5)  f      18  29 p  compact
#OUT> 2 audi      a4      1.8 1999  4 manual(m5) f      21  29 p  compact
#OUT> 3 audi      a4      2   2008  4 manual(m6) f      20  31 p  compact
#OUT> 4 audi      a4      2   2008  4 auto(av)   f      21  30 p  compact
#OUT> 5 audi      a4      2.8 1999  6 auto(l5)  f      16  26 p  compact
#OUT> 6 audi      a4      2.8 1999  6 manual(m5) f      18  26 p  compact
#OUT> 7 audi      a4      3.1 2008  6 auto(av)   f      18  27 p  compact
#OUT> 8 audi      a4 quattro 1.8 1999  4 manual(m5) 4      18  26 p  compact
#OUT> 9 audi      a4 quattro 1.8 1999  4 auto(l5)  4      16  25 p  compact
#OUT> 10 audi     a4 quattro 2   2008  4 manual(m6) 4      20  28 p  compact
```

The function `head()` will display the first `n` observations of the data frame, as we have seen. The `head()` function was more useful before tibbles. Notice that `mpg` is a tibble already, so the output from `head()` indicates there are only 10 observations. Note that this applies to `head(mpg, n = 10)` and not `mpg` itself. Also note that tibbles print a limited number of rows and columns by default. The last line of the printed output indicates with rows and columns were omitted.

```
mpg
```

```
#OUT> # A tibble: 234 x 11
#OUT>   manufacturer model    displ year   cyl trans      drv   cty   hwy fl   class
#OUT>   <chr>         <chr>    <dbl> <int> <int> <chr>    <chr> <int> <int> <chr> <chr>
#OUT> 1 audi      a4      1.8 1999   4 auto(l5) f      18  29 p  compact
#OUT> 2 audi      a4      1.8 1999   4 manual(m5) f      21  29 p  compact
#OUT> 3 audi      a4      2   2008   4 manual(m6) f      20  31 p  compact
#OUT> 4 audi      a4      2   2008   4 auto(av)   f      21  30 p  compact
#OUT> 5 audi      a4      2.8 1999   6 auto(l5)  f      16  26 p  compact
#OUT> 6 audi      a4      2.8 1999   6 manual(m5) f      18  26 p  compact
#OUT> 7 audi      a4      3.1 2008   6 auto(av)   f      18  27 p  compact
#OUT> 8 audi      a4 quattro 1.8 1999   4 manual(m5) 4      18  26 p  compact
#OUT> 9 audi      a4 quattro 1.8 1999   4 auto(l5)  4      16  25 p  compact
#OUT> 10 audi     a4 quattro 2   2008   4 manual(m6) 4      20  28 p  compact
#OUT> # ... with 224 more rows
```

Let's look at `str` as well to get familiar with the content of the data:

```
str(mpg)
```

```
#OUT> tibble [234 x 11] (S3: tbl_df/tbl/data.frame)
#OUT> $ manufacturer: chr [1:234] "audi" "audi" "audi" "audi" ...
#OUT> $ model        : chr [1:234] "a4" "a4" "a4" "a4" ...
#OUT> $ displ        : num [1:234] 1.8 1.8 2 2 2.8 2.8 3.1 1.8 1.8 2 ...
#OUT> $ year         : int [1:234] 1999 1999 2008 2008 1999 1999 2008 1999 1999 2008 ...
#OUT> $ cyl          : int [1:234] 4 4 4 4 6 6 6 4 4 4 ...
#OUT> $ trans        : chr [1:234] "auto(l5)" "manual(m5)" "manual(m6)" "auto(av)" ...
#OUT> $ drv          : chr [1:234] "f" "f" "f" "f" ...
```

```
#OUT> $ cty      : int [1:234] 18 21 20 21 16 18 18 18 16 20 ...
#OUT> $ hwy      : int [1:234] 29 29 31 30 26 26 27 26 25 28 ...
#OUT> $ fl       : chr [1:234] "p" "p" "p" "p" ...
#OUT> $ class    : chr [1:234] "compact" "compact" "compact" "compact" ...
```

In this dataset an observation is for a particular model-year of a car, and the variables describe attributes of the car, for example its highway fuel efficiency.

To understand more about the data set, we use the `?`  operator to pull up the documentation for the data.

```
?mpg
```

Working with tibbles is mostly the same as working with plain `data.frames`:

```
names(mpg)
```

```
#OUT> [1] "manufacturer" "model"      "displ"      "year"      "cyl"      "trans"      "
#OUT> [9] "hwy"          "fl"          "class"
```

```
mpg$year
```

```
#OUT> [1] 1999 1999 2008 2008 1999 1999 2008 1999 1999 2008 2008 1999 1999 2008 2008 1999 2
#OUT> [26] 2008 2008 2008 2008 2008 1999 1999 1999 2008 1999 2008 2008 1999 1999 1999 1999 2
#OUT> [51] 1999 1999 2008 2008 2008 1999 1999 1999 2008 2008 2008 1999 2008 1999 2008 2008 2
#OUT> [76] 1999 2008 1999 1999 1999 2008 2008 1999 1999 1999 1999 1999 2008 1999 2008 1999 1
#OUT> [101] 1999 1999 1999 1999 2008 2008 2008 2008 1999 1999 2008 2008 1999 1999 2008 1999 1
#OUT> [126] 1999 2008 2008 2008 2008 1999 2008 2008 1999 1999 1999 2008 1999 2008 2008 1999 1
#OUT> [151] 1999 1999 2008 2008 1999 1999 1999 2008 2008 1999 1999 2008 2008 2008 2008 1999 1
#OUT> [176] 1999 1999 2008 2008 1999 1999 2008 2008 1999 1999 2008 1999 1999 2008 2008 1999 1
#OUT> [201] 1999 1999 2008 1999 1999 2008 2008 1999 1999 2008 2008 1999 1999 1999 1999 2008 1
#OUT> [226] 2008 2008 1999 1999 2008 2008 1999 1999 2008
```

```
mpg$hwy
```

```
#OUT> [1] 29 29 31 30 26 26 27 26 25 28 27 25 25 25 25 24 25 23 20 15 20 17 17 26 23 26 25 24 1
#OUT> [43] 24 17 22 21 23 23 19 18 17 17 19 19 12 17 15 17 17 12 17 16 18 15 16 12 17 17 16 12 1
#OUT> [85] 17 16 16 17 15 17 26 25 26 24 21 22 23 22 20 33 32 32 29 32 34 36 36 29 26 27 30 31 2
#OUT> [127] 12 19 18 14 15 18 18 15 17 16 18 17 19 19 17 29 27 31 32 27 26 26 25 25 17 17 20 18 1
#OUT> [169] 26 25 27 25 27 20 20 19 17 20 17 29 27 31 31 26 26 28 27 29 31 31 26 26 27 30 33 35 1
#OUT> [211] 29 24 44 29 26 29 29 29 29 23 24 44 41 29 26 28 29 29 29 28 29 26 26 26
```

Subsetting is also similar to `dataframe`. Here, we find fuel efficient vehicles earning over 35 miles per gallon and only display `manufacturer`, `model` and `year`.

```
# mpg[row condition, col condition]
mpg[mpg$hwy > 35, c("manufacturer", "model", "year")]
```

```
#OUT> # A tibble: 6 x 3
#OUT>   manufacturer model      year
```

```
#OUT>   <chr>      <chr>      <int>
#OUT> 1 honda      civic      2008
#OUT> 2 honda      civic      2008
#OUT> 3 toyota     corolla    2008
#OUT> 4 volkswagen jetta      1999
#OUT> 5 volkswagen new beetle 1999
#OUT> 6 volkswagen new beetle 1999
```

An alternative would be to use the `subset()` function, which has a much more readable syntax.

```
subset(mpg, subset = hwy > 35, select = c("manufacturer", "model", "year"))
```

Lastly, and most *tidy*, we could use the `filter` and `select` functions from the `dplyr` package which introduces the pipe operator `f(x) %>% g(z)` from the `magrittr` package. This operator takes the output of the first command, for example `y = f(x)`, and passes it as the first argument to the next function, i.e. we'd obtain `g(y,z)` here.<sup>1</sup>

```
library(dplyr)
mpg %>%
  filter(hwy > 35) %>%
  select(manufacturer, model, year)
```

```
#OUT> # A tibble: 6 x 3
#OUT>   manufacturer model      year
#OUT>   <chr>      <chr>      <int>
#OUT> 1 honda      civic      2008
#OUT> 2 honda      civic      2008
#OUT> 3 toyota     corolla    2008
#OUT> 4 volkswagen jetta      1999
#OUT> 5 volkswagen new beetle 1999
#OUT> 6 volkswagen new beetle 1999
```

Note that the above syntax is equivalent to the following pipe-free command (which is much harder to read!):

```
library(dplyr)
select(filter(mpg, hwy > 35), manufacturer, model, year)
```

```
#OUT> # A tibble: 6 x 3
#OUT>   manufacturer model      year
#OUT>   <chr>      <chr>      <int>
#OUT> 1 honda      civic      2008
#OUT> 2 honda      civic      2008
```

<sup>1</sup>A pipe is a concept from the Unix world, where it means to take the output of some command, and pass it on to another command. This way, one can construct a pipeline of commands. For additional info on the pipe operator in R, you might be interested in this [tutorial](#).

```
#OUT> 3 toyota      corolla      2008
#OUT> 4 volkswagen jetta        1999
#OUT> 5 volkswagen new beetle   1999
#OUT> 6 volkswagen new beetle   1999
```

All three approaches produce the same results. Which you use will be largely based on a given situation as well as your preference.

#### 2.4.2.1 Task 1

1. Make sure to have the `mpg` dataset loaded by typing `data(mpg)` (and `library(ggplot2)` if you haven't!). Use the `table` function to find out how many cars were built by mercury?
2. What is the average year the audi's were built in this dataset? Use the function `mean` on the subset of column `year` that corresponds to audi. (Be careful: subsetting a `tibble` returns a `tibble` (and not a vector)!. so get the `year` column after you have subset the `tibble`.)
3. Use the `dplyr` piping syntax from above first with `group_by` and then with `summarise(newvar=your_expression)` to find the mean `year` by all manufacturers (i.e. same as previous task, but for all manufacturers. don't write a loop!).

### 2.4.3 Tidy Example: Importing Non-Tidy Excel Data

The data we will look at is from Eurostat on demography and migration. You should download the data yourself (click on previous link, then drill down to database by themes > Population and social conditions > Demograph and migration > Population change - Demographic balance and crude rates at national level (demo\_gind)).

Once downloaded, we can read the data with the function `read_excel` from the package `readxl`, again part of the `tidyverse` suite.

It's important to know how the data is organized in the spreadsheet. Open the file with Excel to see:

- There is a heading which we don't need.
- There are 5 rows with info that we don't need.
- There is one table per variable (total population, males, females, etc)
- Each table has one row for each country, and one column for each year.
- As such, this data is **not tidy**.

Now we will read the first chunk of data, from the first table: total population:

```
library(readxl) # load the library
# Notice that if you installed the R package of this book,
# you have the .xls data file already at
```



```
# `system.file(package="ScPoEconometrics",
#               "datasets", "demo_gind.xls")`
# otherwise:
# * download the file to your computer
# * change the argument `path` to where you downloaded it
# you may want to change your working directory with `setwd("your/directory")
# or in RStudio by clicking Session > Set Working Directory

# total population in raw format
tot_pop_raw = read_excel(
  path = system.file(package="ScPoEconometrics",
                      "datasets", "demo_gind.xls"),
  sheet="Data", # which sheet
  range="A9:K68") # which excel cell range to read
names(tot_pop_raw)[1] <- "Country" # lets rename the first column
tot_pop_raw
```

```
#OUT> # A tibble: 59 x 11
#OUT>   Country                `2008`  `2009`  `2010`  `2011`  `2012`  `2013`  `2014`  `20
#OUT>   <chr>                <chr>   <chr>   <chr>   <chr>   <chr>   <chr>   <chr>
#OUT> 1 European Union (current composition) 500297033 502090235 5031706~ 5029648~ 5040479~ 505163
#OUT> 2 European Union (without United Kingdo~ 438725386 440047892 4406604~ 4399423~ 4405526~ 441257
#OUT> 3 European Union (before the accession ~ 495985066 497780439 4988677~ 4986749~ 4997719~ 500900
#OUT> 4 Euro area (19 countries)            333096775 334470255 3352664~ 3345725~ 3352889~ 3360449~ 33
#OUT> 5 Euro area (18 countries)            329884170 331286399 3321244~ 3315200~ 3322852~ 3330730~ 33
#OUT> 6 Belgium                          10666866 10753080 10839905 11000638 11075889 11137974 1118084
#OUT> 7 Bulgaria                          7518002 7467119 7421766 7369431 7327224 7284552 7245677
#OUT> 8 Czech Republic                    10343422 10425783 10462088 10486731 10505445 10516125 10512
#OUT> 9 Denmark                          5475791 5511451 5534738 5560628 5580516 5602628 5627235
#OUT> 10 Germany (until 1990 former territory ~ 82217837 82002356 81802257 80222065 80327900 8052374
#OUT> # ... with 49 more rows
```

This shows a tibble, which we encountered just above. The column names are Country, 2008, 2009, ..., and the rows are numbered 1, 2, 3, .... Notice, in particular, that all columns seem to be of type `<chr>`, i.e. characters - a string, not a number! We'll have to fix that, as this is clearly numeric data.

### 2.4.3.1 tidyr

In the previous tibble, each year is a column name (like 2008) instead of all years being collected in one column `year`. We really would like to have several rows for each Country, one row per year. We want to `gather()` all years into a new column to tidy this up - and here is how:

1. specify which columns are to be gathered: in our case, all years (note that `paste(2008:2017)` produces a vector like `["2008", "2009",`

"2010",...])

2. say what those columns should be gathered into, i.e. what is the key for those values: we'll call it year.
3. Finally, what is the name of the new resulting column, containing the value from each cell: let's call it counts.

```
library(tidyr)    # for the gather function
tot_pop = gather(tot_pop_raw, paste(2008:2017),key="year", value = "counts")
tot_pop
```

```
#OUT> # A tibble: 590 x 3
#OUT>   Country                                year counts
#OUT>   <chr>                                <chr> <chr>
#OUT> 1 European Union (current composition) 2008 500297033
#OUT> 2 European Union (without United Kingdom) 2008 438725386
#OUT> 3 European Union (before the accession of Croatia) 2008 495985066
#OUT> 4 Euro area (19 countries)                2008 333096775
#OUT> 5 Euro area (18 countries)                2008 329884170
#OUT> 6 Belgium                                2008 10666866
#OUT> 7 Bulgaria                                2008 7518002
#OUT> 8 Czech Republic                        2008 10343422
#OUT> 9 Denmark                                2008 5475791
#OUT> 10 Germany (until 1990 former territory of the FRG) 2008 82217837
#OUT> # ... with 580 more rows
```

That's better! However, counts is still chr! Let's convert it to a number:

```
tot_pop$counts = as.integer(tot_pop$counts)
```

```
#OUT> Warning: NAs introduced by coercion
```

```
tot_pop
```

```
#OUT> # A tibble: 590 x 3
#OUT>   Country                                year counts
#OUT>   <chr>                                <chr> <int>
#OUT> 1 European Union (current composition) 2008 500297033
#OUT> 2 European Union (without United Kingdom) 2008 438725386
#OUT> 3 European Union (before the accession of Croatia) 2008 495985066
#OUT> 4 Euro area (19 countries)                2008 333096775
#OUT> 5 Euro area (18 countries)                2008 329884170
#OUT> 6 Belgium                                2008 10666866
#OUT> 7 Bulgaria                                2008 7518002
#OUT> 8 Czech Republic                        2008 10343422
#OUT> 9 Denmark                                2008 5475791
#OUT> 10 Germany (until 1990 former territory of the FRG) 2008 82217837
#OUT> # ... with 580 more rows
```

Now you can see that column counts is indeed int, i.e. an integer number,

and we are fine. The `Warning: NAs introduced by coercion` means that R converted some values to NA, because it couldn't convert them into `numeric`. More below!

### 2.4.3.2 dplyr

The transform chapter of Hadley Wickham's book is a great place to read up more on using `dplyr`.

With `dplyr` you can do the following operations on `data.frames` and `tibbles`:

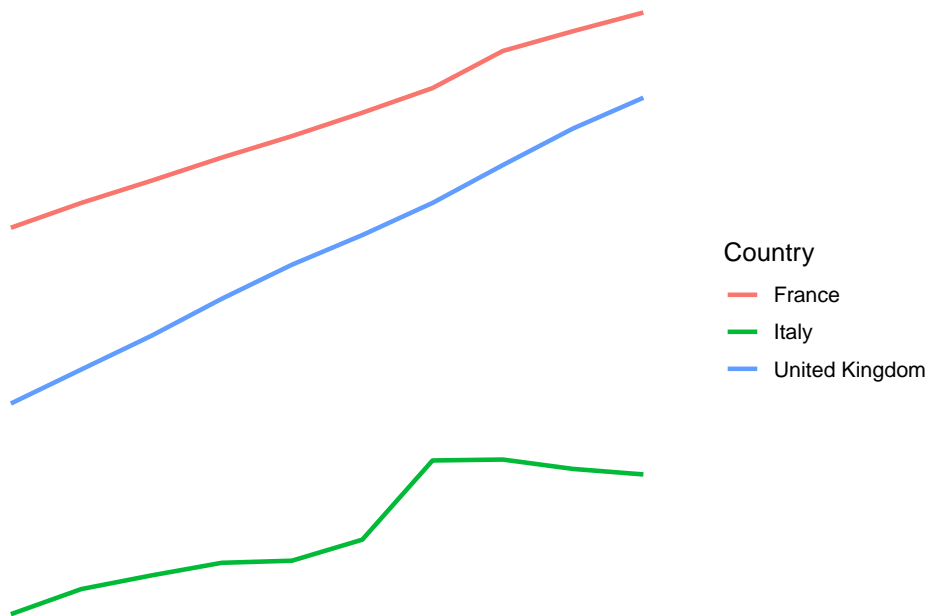
- Choose observations based on a certain value (i.e. subset): `filter()`
- Reorder rows: `arrange()`
- Select variables by name: `select()`
- Create new variables out of existing ones: `mutate()`
- Summarise variables: `summarise()`

All of those verbs can be used with `group_by()`, where we apply the respective operation on a group of the dataframe/tibble. For example, on our `tot_pop` tibble we will now

- filter
- mutate
- and plot the resulting values

Let's get a plot of the populations of France, the UK and Italy over time, in terms of millions of people. We will make use of the piping syntax of `dplyr` which we introduced just above.

```
library(dplyr) # for %>%, filter, mutate, ...
# 1. take the data.frame `tot_pop`
tot_pop %>%
  # 2. pipe it into the filter function
  # filter on Country being one of "France", "United Kingdom" or "Italy"
  filter(Country %in% c("France", "United Kingdom", "Italy")) %>%
  # 3. pipe the result into the mutate function
  # create a new column called millions
  mutate(millions = counts / 1e6) %>%
  # 4. pipe the result into ggplot to make a plot
  ggplot(mapping = aes(x=year, y=millions, color=Country, group=Country)) + geom_line(size=1)
```



### Arrange a tibble

- What are the top/bottom 5 most populated areas?

```
top5 = tot_pop %>%
  arrange(desc(counts)) %>% # arrange in descending order of col `counts`
  top_n(5)

bottom5 = tot_pop %>%
  arrange(desc(counts)) %>%
  top_n(-5)
# let's see top 5
top5
```

```
#OUT> # A tibble: 5 x 3
```

```
#OUT> Country                                year    counts
#OUT> <chr>                                <chr>    <int>
```

```
#OUT> 1 European Economic Area (EU28 - current composition, plus IS, LI, NO)    2017
```

```
#OUT> 2 European Economic Area (EU28 - current composition, plus IS, LI, NO)    2016
```

```
#OUT> 3 European Economic Area (EU28 - current composition, plus IS, LI, NO)    2015
```

```
#OUT> 4 European Economic Area (EU27 - before the accession of Croatia, plus IS, LI, NO) 2014
```

```
#OUT> 5 European Economic Area (EU28 - current composition, plus IS, LI, NO)    2014
```

```
# and bottom 5
```

```
bottom5
```

```
#OUT> # A tibble: 5 x 3
#OUT>   Country    year counts
#OUT>   <chr>      <chr> <int>
#OUT> 1 San Marino 2015  32789
#OUT> 2 San Marino 2014  32520
#OUT> 3 San Marino 2008  32054
#OUT> 4 San Marino 2011  31863
#OUT> 5 San Marino 2009  31269
```

Now this is not exactly what we wanted. It's always the same country in both top and bottom, because there are multiple years per country. Let's compute average population over the last 5 years and rank according to that:

```
topbottom = tot_pop %>%
  group_by(Country) %>%
  filter(year > 2012) %>%
  summarise(mean_count = mean(counts)) %>%
  arrange(desc(mean_count))

top5 = topbottom %>% top_n(5)
bottom5 = topbottom %>% top_n(-5)
top5
```

```
#OUT> # A tibble: 5 x 2
#OUT>   Country                                mean_count
#OUT>   <chr>                                <dbl>
#OUT> 1 European Economic Area (EU28 - current composition, plus IS, LI, NO)      514029320
#OUT> 2 European Economic Area (EU27 - before the accession of Croatia, plus IS, LI, NO) 509813491.
#OUT> 3 European Union (current composition)                                     508502858.
#OUT> 4 European Union (before the accession of Croatia)                         504287028.
#OUT> 5 European Union (without United Kingdom)                                443638309.
bottom5
```

```
#OUT> # A tibble: 5 x 2
#OUT>   Country    mean_count
#OUT>   <chr>      <dbl>
#OUT> 1 Luxembourg  563319.
#OUT> 2 Malta       440467.
#OUT> 3 Iceland     329501.
#OUT> 4 Liechtenstein 37353
#OUT> 5 San Marino   33014.
```

That's better!

### Look for NAs in a tibble

Sometimes data is missing, and R represents it with the special value NA (not available). It is good to know where in our dataset we are going to encounter any missing values, so the task here is: let's produce a table that has three columns:

1. the names of countries with missing data
2. how many years of data are missing for each of those
3. and the actual years that are missing

```
missings = tot_pop %>%
  filter(is.na(counts)) %>% # is.na(x) returns TRUE if x is NA
  group_by(Country) %>%
  summarise(n_missing = n(), years = paste(year, collapse = ", "))
```

```
#OUT> `summarise()` ungrouping output (override with `.groups` argument)
```

```
knitr::kable(missings) # knitr::kable makes a nice table
```

Country	n_missing	years
Albania	2	2010, 2012
Andorra	2	2014, 2015
Armenia	1	2014
France (metropolitan)	4	2014, 2015, 2016, 2017
Georgia	1	2013
Monaco	7	2008, 2009, 2010, 2011, 2012, 2013, 2014
Russia	4	2013, 2015, 2016, 2017
San Marino	1	2010

### Males and Females

Let's look at the numbers by male and female population. They are in the same xls file, but at different cell ranges. Also, I just realised that the special character : indicates missing data. We can feed that to `read_excel` and that will spare us the need to convert data types afterwards. Let's see:

```
females_raw = read_excel(
  path = system.file(package="ScPoEconometrics",
                      "datasets", "demo_gind.xls"),
  sheet="Data", # which sheet
  range="A141:K200", # which excel cell range to read
  na=":" ) # missing data indicator
names(females_raw)[1] <- "Country" # lets rename the first column
females_raw
```

```
#OUT> # A tibble: 59 x 11
```

```

#OUT> Country                `2008`  `2009`  `2010`  `2011`  `2012`  `2013`  `2014`  `2015`
#OUT> <chr>                   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
#OUT> 1 European Union (current composition) 256306485 257178433 2.58e8 2.58e8 2.58e8 2.59e8 2.59e8
#OUT> 2 European Union (without United Kingdo~ 224899187 225552613 2.26e8 2.26e8 2.26e8 2.26e8 2.26e8
#OUT> 3 European Union (before the accession ~ 254071506 254945404 2.55e8 2.56e8 2.56e8 2.56e8 2.57e8
#OUT> 4 Euro area (19 countries)              170579892 171269855 1.72e8 1.72e8 1.72e8 1.72e8 1.72e8
#OUT> 5 Euro area (18 countries)              168854706 169559351 1.70e8 1.70e8 1.70e8 1.71e8 1.71e8
#OUT> 6 Belgium                             5442557  5484429  5.53e6 5.60e6 5.64e6 5.67e6 5.69e6
#OUT> 7 Bulgaria                             3857635  3831264  3.81e6 3.78e6 3.76e6 3.74e6 3.72e6
#OUT> 8 Czech Republic                       5278305  5309105  5.33e6 5.34e6 5.35e6 5.35e6 5.35e6
#OUT> 9 Denmark                             2763125  2779431  2.79e6 2.80e6 2.81e6 2.82e6 2.83e6
#OUT> 10 Germany (until 1990 former territory ~ 41943545 41818073 4.17e7 4.11e7 4.11e7 4.11e7 4.11e7
#OUT> # ... with 49 more rows

```

You can see that R now correctly read the numbers as such, after we told it that the `:` character has the special `missing` meaning: before, it coerced the entire 2008 column (for example) to be of type `chr` after it hit the first `:`. We had to manually convert the column back to `numeric`, in the process automatically coercing the `:`s into `NA`. Now we addressed that issue directly. Let's also get the male data in the same way:

```

males_raw = read_excel(
  path = system.file(package="ScPoEconometrics",
                     "datasets", "demo_gind.xls"),
  sheet="Data", # which sheet
  range="A75:K134", # which excel cell range to read
  na=":" ) # missing data indicator
names(males_raw)[1] <- "Country" # lets rename the first column

```

Next step was to tidy up this data, just as before:

```

females = gather(females_raw, paste(2008:2017), key="year", value = "counts")
males = gather(males_raw, paste(2008:2017), key="year", value = "counts")

```

Let's try to tweak our above plot to show the same data in two separate panels: one for males and one for females. This is easiest to do with `ggplot` if we have all the data in one single `data.frame` (or `tibble`), and marked with a `group identifier`. Let's first add this to both datasets, and then let's just combine both into one:

```

females$sex = "female"
males$sex = "male"
sexes = rbind(males, females) # "row bind" 2 data.frames
sexes

```

```

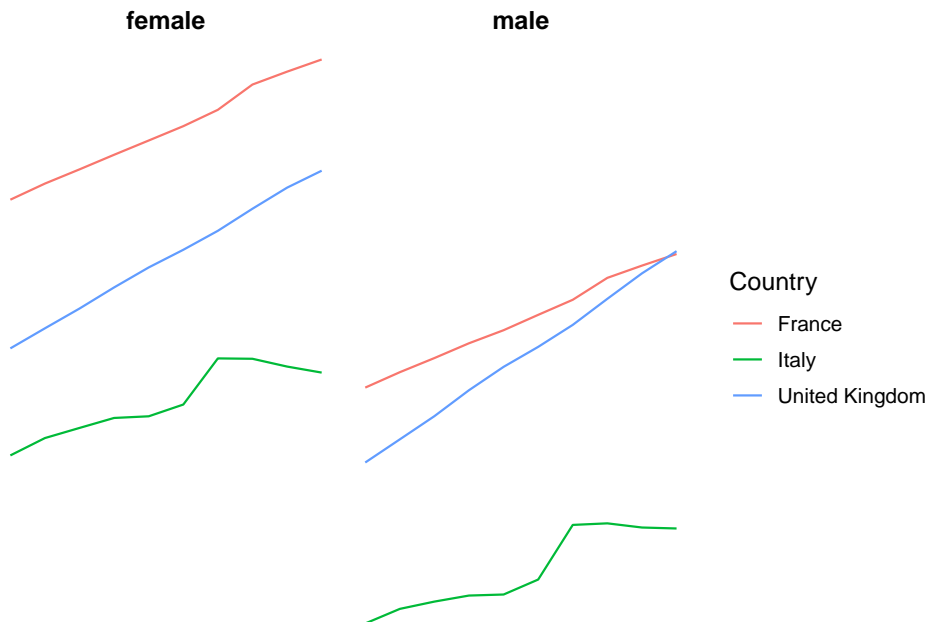
#OUT> # A tibble: 1,180 x 4
#OUT> Country                year    counts sex
#OUT> <chr>                   <chr>   <dbl> <chr>

```

```
#OUT> 1 European Union (current composition)      2008  243990548 male
#OUT> 2 European Union (without United Kingdom)   2008  213826199 male
#OUT> 3 European Union (before the accession of Croatia) 2008  241913560 male
#OUT> 4 Euro area (19 countries)                  2008  162516883 male
#OUT> 5 Euro area (18 countries)                  2008  161029464 male
#OUT> 6 Belgium                                  2008   5224309 male
#OUT> 7 Bulgaria                                 2008   3660367 male
#OUT> 8 Czech Republic                          2008   5065117 male
#OUT> 9 Denmark                                  2008   2712666 male
#OUT> 10 Germany (until 1990 former territory of the FRG) 2008  40274292 male
#OUT> # ... with 1,170 more rows
```

Now that we have all the data nice and tidy in a `data.frame`, this is a very small change to our previous plotting code:

```
sexes %>%
  filter(Country %in% c("France", "United Kingdom", "Italy")) %>%
  mutate(millions = counts / 1e6) %>%
  ggplot(mapping = aes(x=as.Date(year,format="%Y"), # convert to `Date`
                      y=millions,colour=Country,group=Country)) +
    geom_line() +
    scale_x_date(name = "year") + # rename x axis
    facet_wrap(~sex) # make two panels, splitting by groups `sex`
```





### Always Compare to Germany :-)

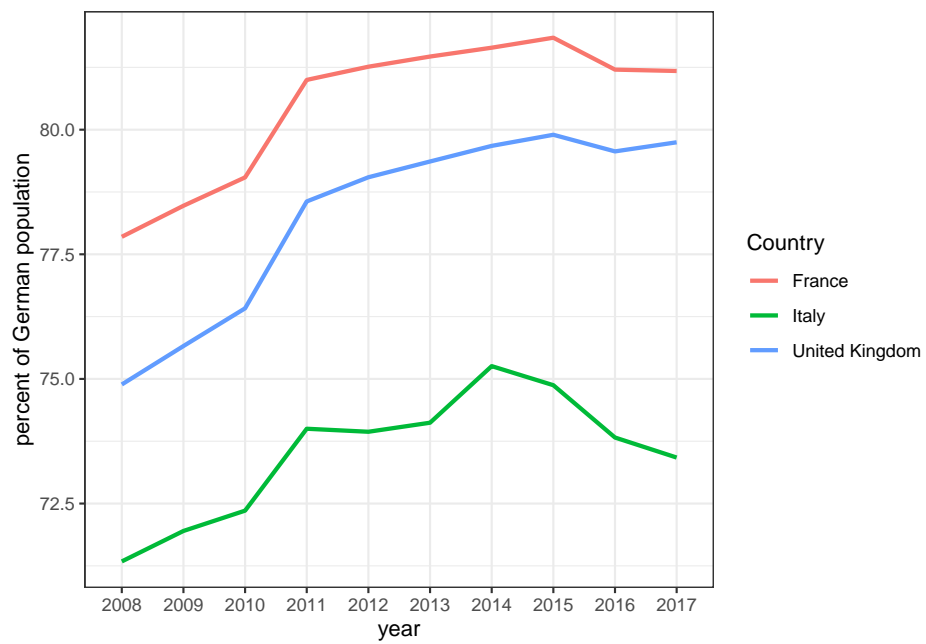
How do our three countries compare with respect to the biggest country in the EU in terms of population? What fraction of Germany does the French population make in any given year, for example?

```
# remember that the pipe operator %>% takes the
# result of the previous operation and passes it
# as the *first* argument to the next function call
merge_GER <- tot_pop %>%
  # 1. subset to countries of interest
  filter(
    Country %in%
      c("France",
        "United Kingdom",
        "Italy")
  ) %>%
  # 2. group data by year
  group_by(year) %>%
  # 3. add GER's count as new column *by year*
  left_join(
    # Germany only
    filter(tot_pop,
      Country %in% "Germany including former GDR"),
    # join back in `by year`
    by="year")
merge_GER

#OUT> # A tibble: 30 x 5
#OUT> # Groups:   year [10]
#OUT>   Country.x   year counts.x Country.y           counts.y
#OUT>   <chr>      <chr>   <int> <chr>              <int>
#OUT> 1 France      2008  64007193 Germany including former GDR 82217837
#OUT> 2 Italy       2008  58652875 Germany including former GDR 82217837
#OUT> 3 United Kingdom 2008  61571647 Germany including former GDR 82217837
#OUT> 4 France      2009  64350226 Germany including former GDR 82002356
#OUT> 5 Italy       2009  59000586 Germany including former GDR 82002356
#OUT> 6 United Kingdom 2009  62042343 Germany including former GDR 82002356
#OUT> 7 France      2010  64658856 Germany including former GDR 81802257
#OUT> 8 Italy       2010  59190143 Germany including former GDR 81802257
#OUT> 9 United Kingdom 2010  62510197 Germany including former GDR 81802257
#OUT> 10 France     2011  64978721 Germany including former GDR 80222065
#OUT> # ... with 20 more rows
```

Here you see that the merge (or join) operation labelled `col.x` and `col.y` if both datasets contained a column called `col`. Now let's continue to compute what proportion of german population each country amounts to:

```
names(merge_GER)[1] <- "Country"
merge_GER %>%
  mutate(prop_GER = 100 * counts.x / counts.y) %>%
  # 5. plot
  ggplot(mapping =
    aes(x = year,
        y = prop_GER,
        color = Country,
        group = Country)) +
  geom_line(size=1) +
  scale_y_continuous("percent of German population") +
  theme_bw() # new theme for a change?
```



## Chapter 3

# Linear Regression

In this chapter we will learn an additional way how one can represent the relationship between outcome, or dependent variable  $y$  and an explanatory or independent variable  $x$ . We will refer throughout to the graphical representation of a collection of independent observations on  $x$  and  $y$ , i.e., a dataset.

### 3.1 How are $x$ and $y$ related?

#### 3.1.1 Data on Cars

We will look at the built-in `cars` dataset. Let's get a view of this by just typing `View(cars)` in Rstudio. You can see something like this:

```
#OUT>   speed dist
#OUT> 1     4    2
#OUT> 2     4   10
#OUT> 3     7    4
#OUT> 4     7   22
#OUT> 5     8   16
#OUT> 6     9   10
```

We have a `data.frame` with two columns: `speed` and `dist`. Type `help(cars)` to find out more about the dataset. There you could read that

The data give the speed of cars (mph) and the distances taken to stop (ft).

It's good practice to know the extent of a dataset. You could just type

```
dim(cars)
```

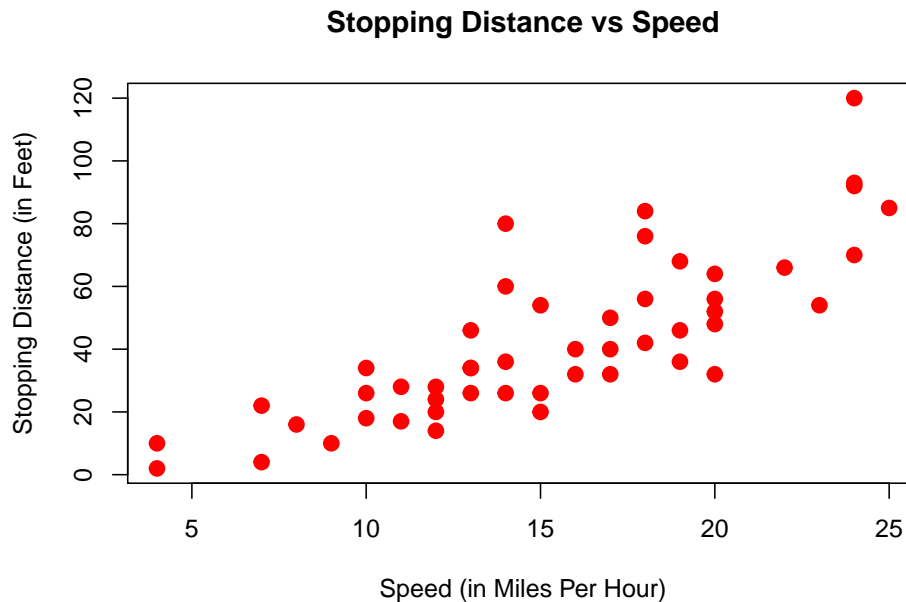
```
#OUT> [1] 50 2
```

to find out that we have 50 rows and 2 columns. A central question that we want to ask now is the following:

### 3.1.2 How are speed and dist related?

The simplest way to start is to plot the data. Remembering that we view each row of a data.frame as an observation, we could just label one axis of a graph **speed**, and the other one **dist**, and go through our table above row by row. We just have to read off the x/y coordinates and mark them in the graph. In R:

```
plot(dist ~ speed, data = cars,
      xlab = "Speed (in Miles Per Hour)",
      ylab = "Stopping Distance (in Feet)",
      main = "Stopping Distance vs Speed",
      pch = 20,
      cex = 2,
      col = "red")
```

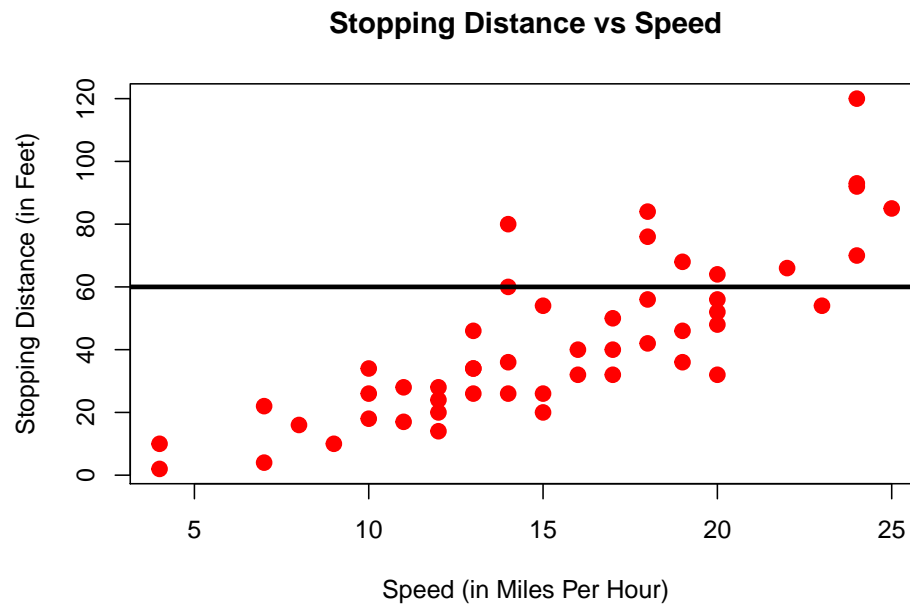


Here, each dot represents one observation. In this case, one particular measurement **speed** and **dist** for a car. Now, again:

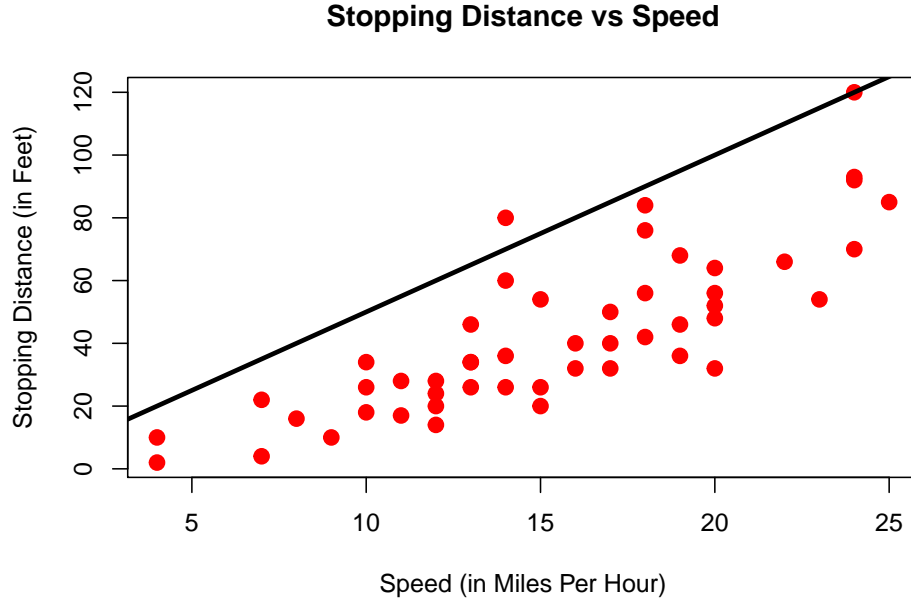
How are **speed** and **dist** related? How could one best summarize this relationship?

One thing we could do, is draw a straight line through this scatterplot, like so:

```
plot(dist ~ speed, data = cars,  
     xlab = "Speed (in Miles Per Hour)",  
     ylab = "Stopping Distance (in Feet)",  
     main = "Stopping Distance vs Speed",  
     pch = 20,  
     cex = 2,  
     col = "red")  
abline(a = 60, b = 0, lw=3)
```



Now that doesn't seem a particularly good way to summarize the relationship. Clearly, a better line would be not be flat, but have a slope, i.e. go upwards:



That is slightly better. However, the line seems at too high a level - the point at which it crosses the y-axis is called the intercept; and it's too high. We just learned how to represent a line, i.e. with two numbers called intercept and slope. Let's write down a simple formula which represents a line where some outcome  $z$  is related to a variable  $x$ :

$$z = b_0 + b_1x \quad (3.1)$$

Here  $b_0$  represents the value of the intercept (i.e.  $z$  when  $x = 0$ ), and  $b_1$  is the value of the slope. The question for us is now: How to choose the number  $b_0$  and  $b_1$  such that the result is the **good** line?

### 3.1.3 Choosing the Best Line

In order to be able to reason about good or bad line, we need to denote the output of equation (??). We call the value  $\hat{y}_i$  the predicted value for observation  $i$ , after having chosen some particular values  $b_0$  and  $b_1$ :

$$\hat{y}_i = b_0 + b_1x_i \quad (3.2)$$

In general it is likely that we won't be able to choose  $b_0$  and  $b_1$  in such a way as to provide a perfect prediction, i.e. one where  $\hat{y}_i = y_i$  for all  $i$ . That is, we expect to make an error in our prediction  $\hat{y}_i$ , so let's denote this value  $e_i$ . If we acknowledge that we will make errors, let's at least make them as small as possible! Exactly this is going to be our task now.

Suppose we have the following set of 9 observations on  $x$  and  $y$ , and we put the best straight line into it, that we can think of. It would look like this:

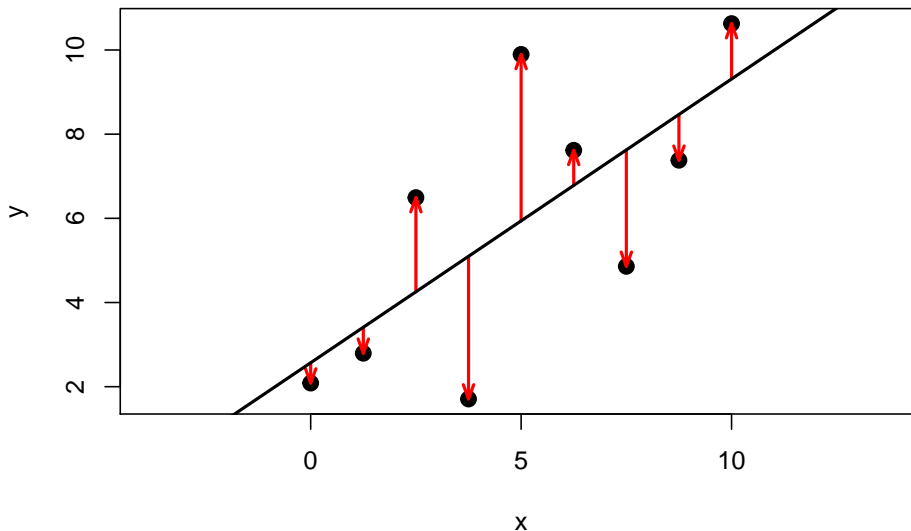


Figure 3.1: The best line and its errors

Here, the red arrows indicate the **distance** between the prediction (i.e. the black line) to each data point, in other words, each arrow is a particular  $e_i$ . An upward pointing arrow indicates a positive value of a particular  $e_i$ , and vice versa for downward pointing arrows. The errors are also called residuals, which comes from the way can write the equation for this relationship between two particular values  $(y_i, x_i)$  belonging to observation  $i$ :

$$y_i = b_0 + b_1 x_i + e_i \quad (3.3)$$

You realize of course that  $\hat{y}_i = y_i - e_i$ , which just means that our prediction is the observed value  $y_i$  minus any error  $e_i$  we make. In other words,  $e_i$  is what is left to be explained on top of the line  $b_0 + b_1 x_i$ , hence, it's a residual to explain  $y_i$ . Here are  $y$ ,  $\hat{y}$  and the resulting  $e$  which are plotted in figure ??:

x	y	y_hat	error
0.00	2.09	2.57	-0.48
1.25	2.79	3.41	-0.62
2.50	6.49	4.25	2.24
3.75	1.71	5.10	-3.39
5.00	9.89	5.94	3.95
6.25	7.62	6.78	0.83
7.50	4.86	7.63	-2.77
8.75	7.38	8.47	-1.09
10.00	10.63	9.31	1.32

If our line was a **perfect fit** to the data, all  $e_i = 0$ , and the column **error** would display 0 for each row - there would be no errors at all. (All points in figure ?? would perfectly line up on a straight line).

Now, back to our claim that this particular line is the best line. What exactly characterizes this best line? We now come back to what we said above - how to make the errors as small as possible? Keeping in mind that each residual  $e_i$  is  $y_i - \hat{y}_i$ , we have the following minization problem to solve:

$$e_i = y_i - \hat{y}_i = y_i - \underbrace{(b_0 + b_1 x_i)}_{\text{prediction}} \quad (3.4)$$

$$e_1^2 + \dots + e_N^2 = \sum_{i=1}^N e_i^2 \equiv \text{SSR}(b_0, b_1) \quad (3.5)$$

$$(b_0, b_1) = \arg \min_{\text{int, slope}} \sum_{i=1}^N [y_i - (\text{int} + \text{slope } x_i)]^2 \quad (3.6)$$

The best line chooses  $b_0$  and  $b_1$  so as to minimize the sum of **squared residuals** (SSR).

Wait a moment, why squared residuals? This is easy to understand: suppose that instead, we wanted to just make the sum of the arrows in figure ?? as small as possible (that is, no squares). Choosing our line to make this number small would not give a particularly good representation of the data – given that errors of opposite sign and equal magnitude offset, we could have very long arrows (but of opposite signs), and a poor resulting line. Squaring each error avoids this (because now negative errors get positive values!)

We illustrate this in figure ??. This is the same data as in figure ??, but instead of arrows of length  $e_i$  for each observation  $i$ , now we draw a square with side  $e_i$ , i.e. an area of  $e_i^2$ . We have two apps for you at this point, one where you have to try and find the best line by choosing  $b_0$  and  $b_1$ , only focusing on the sum of errors (and not their square), and a second one focusing on squared errors:

```
library(ScPoEconometrics)
launchApp("reg_simple_arrows")
launchApp("reg_simple") # with squared errors
launchApp("SSR_cone") # visualize the minimization problem from above!
```

Most of our apps have an associated **about** document, which gives extra information and explanations. After you have looked at all three apps, we invite you thus to have a look at the associated explainers by typing



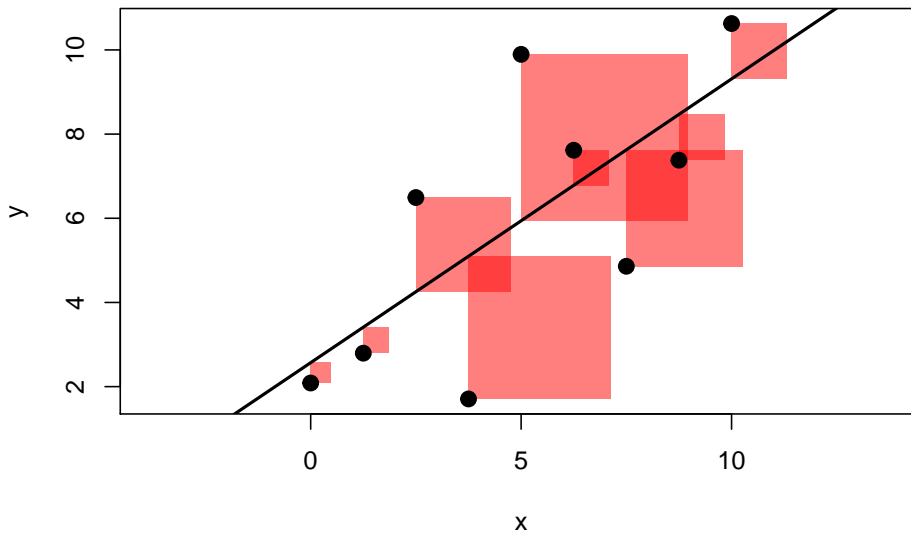


Figure 3.2: The best line and its SQUARED errors

```
aboutApp("reg_simple_arrows")
aboutApp("reg_simple")
aboutApp("SSR_cone")
```

## 3.2 Ordinary Least Squares (OLS) Estimator

The method to compute (or estimate)  $b_0$  and  $b_1$  we illustrated above is called Ordinary Least Squares, or OLS.  $b_0$  and  $b_1$  are therefore also often called the OLS coefficients. By solving problem (??) one can derive an explicit formula for them:

$$b_1 = \frac{\text{cov}(x, y)}{\text{var}(x)}, \quad (3.7)$$

i.e. the estimate of the slope coefficient is the covariance between  $x$  and  $y$  divided by the variance of  $x$ , both computed from our sample of data. With  $b_1$  in hand, we can get the estimate for the intercept as

$$b_0 = \bar{y} - b_1 \bar{x}. \quad (3.8)$$

where  $\bar{z}$  denotes the sample mean of variable  $z$ . The interpretation of the OLS slope coefficient  $b_1$  is as follows. Given a line as in  $y = b_0 + b_1 x$ ,

- $b_1 = \frac{dy}{dx}$  measures the change in  $y$  resulting from a one unit change in  $x$
- For example, if  $y$  is wage and  $x$  is years of education,  $b_1$  would measure the effect of an additional year of education on wages.

There is an alternative representation for the OLS slope coefficient which relates to the correlation coefficient  $r$ . Remember from section ?? that  $r = \frac{\text{cov}(x,y)}{s_x s_y}$ , where  $s_z$  is the standard deviation of variable  $z$ . With this in hand, we can derive the OLS slope coefficient as

$$b_1 = \frac{\text{cov}(x, y)}{\text{var}(x)} \quad (3.9)$$

$$= \frac{\text{cov}(x, y)}{s_x s_x} \quad (3.10)$$

$$= r \frac{s_y}{s_x} \quad (3.11)$$

In other words, the slope coefficient is equal to the correlation coefficient  $r$  times the ratio of standard deviations of  $y$  and  $x$ .

### 3.2.1 Linear Regression without Regressor

There are several important special cases for the linear regression introduced above. Let's start with the most obvious one: What is the meaning of running a regression without any regressor, i.e. without a  $x$ ? Our line becomes very simple. Instead of (??), we get

$$y = b_0. \quad (3.12)$$

This means that our minimization problem in (??) also becomes very simple: We only have to choose  $b_0$ ! We have

$$b_0 = \arg \min_{\text{int}} \sum_{i=1}^N [y_i - \text{int}]^2,$$

which is a quadratic equation with a unique optimum such that

$$b_0 = \frac{1}{N} \sum_{i=1}^N y_i = \bar{y}.$$

Least Squares **without regressor**  $x$  estimates the sample mean of the outcome variable  $y$ , i.e. it produces  $\bar{y}$ .

### 3.2.2 Regression without an Intercept

We follow the same logic here, just that we miss another bit from our initial equation and the minimisation problem in (??) now becomes:

$$b_1 = \arg \min_{\text{slope}} \sum_{i=1}^N [y_i - \text{slope } x_i]^2 \quad (3.13)$$

$$\mapsto b_1 = \frac{\frac{1}{N} \sum_{i=1}^N x_i y_i}{\frac{1}{N} \sum_{i=1}^N x_i^2} = \frac{\bar{x}\bar{y}}{\bar{x}^2} \quad (3.14)$$

Least Squares **without intercept** (i.e. with  $b_0 = 0$ ) is a line that passes through the origin.

In this case we only get to choose the slope  $b_1$  of this anchored line.<sup>1</sup> You should now try out both of those restrictions on our linear model by spending some time with

```
launchApp("reg_constrained")
```

### 3.2.3 Centering A Regression

By centering or demeaning a regression, we mean to subtract from both  $y$  and  $x$  their respective averages to obtain  $\tilde{y}_i = y_i - \bar{y}$  and  $\tilde{x}_i = x_i - \bar{x}$ . We then run a regression without intercept as above. That is, we use  $\tilde{x}_i, \tilde{y}_i$  instead of  $x_i, y_i$  in (??) to obtain our slope estimate  $b_1$ :

$$b_1 = \frac{\frac{1}{N} \sum_{i=1}^N \tilde{x}_i \tilde{y}_i}{\frac{1}{N} \sum_{i=1}^N \tilde{x}_i^2} \quad (3.15)$$

$$= \frac{\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2} \quad (3.16)$$

$$= \frac{\text{cov}(x, y)}{\text{var}(x)} \quad (3.17)$$

This last expression is identical to the one in (??)! It's the standard OLS estimate for the slope coefficient. We note the following:

<sup>1</sup>This slope is related to the angle between vectors  $\mathbf{a} = (\bar{x}, \bar{y})$ , and  $\mathbf{b} = (\bar{x}, 0)$ . Hence, it's related to the scalar projection of  $\mathbf{a}$  on  $\mathbf{b}$ .

Adding a constant to a regression produces the same result as centering all variables and estimating without intercept. So, unless all variables are centered, **always** include an intercept in the regression.

To get a better feel for what is going on here, you can try this out now by yourself by typing:

```
launchApp("demeaned_reg")
```

### 3.2.4 Standardizing A Regression

Standardizing a variable  $z$  means to demean as above, but in addition to divide the demeaned value by its own standard deviation. Similarly to what we did above for centering, we define transformed variables  $\tilde{y}_i = \frac{y_i - \bar{y}}{\sigma_y}$  and  $\tilde{x}_i = \frac{x_i - \bar{x}}{\sigma_x}$  where  $\sigma_z$  is the standard deviation of variable  $z$ . From here on, you should by now be used to what comes next! As above, we use  $\tilde{x}_i, \tilde{y}_i$  instead of  $x_i, y_i$  in (??) to this time obtain:

$$b_1 = \frac{\frac{1}{N} \sum_{i=1}^N \tilde{x}_i \tilde{y}_i}{\frac{1}{N} \sum_{i=1}^N \tilde{x}_i^2} \quad (3.18)$$

$$= \frac{\frac{1}{N} \sum_{i=1}^N \frac{x_i - \bar{x}}{\sigma_x} \frac{y_i - \bar{y}}{\sigma_y}}{\frac{1}{N} \sum_{i=1}^N \left( \frac{x_i - \bar{x}}{\sigma_x} \right)^2} \quad (3.19)$$

$$= \frac{Cov(x, y)}{\sigma_x \sigma_y} \quad (3.20)$$

$$= Corr(x, y) \quad (3.21)$$

After we standardize both  $y$  and  $x$ , the slope coefficient  $b_1$  in the regression without intercept is equal to the **correlation coefficient**.

And also for this case we have a practical application for you. Just type this and play around with the app for a little while!

```
launchApp("reg_standardized")
```

## 3.3 Predictions and Residuals

Now we want to ask how our residuals  $e_i$  relate to the prediction  $\hat{y}_i$ . Let us first think about the average of all predictions  $\hat{y}_i$ , i.e. the number  $\frac{1}{N} \sum_{i=1}^N \hat{y}_i$ . Let's

just take (??) and plug this into this average, so that we get

$$\frac{1}{N} \sum_{i=1}^N \hat{y}_i = \frac{1}{N} \sum_{i=1}^N b_0 + b_1 x_i \quad (3.22)$$

$$= b_0 + b_1 \frac{1}{N} \sum_{i=1}^N x_i \quad (3.23)$$

$$= b_0 + b_1 \bar{x} \quad (3.24)$$

$$(3.25)$$

But that last line is just equal to the formula for the OLS intercept (??),  $b_0 = \bar{y} - b_1 \bar{x}$ ! That means of course that

$$\frac{1}{N} \sum_{i=1}^N \hat{y}_i = b_0 + b_1 \bar{x} = \bar{y}$$

in other words:

The average of our predictions  $\hat{y}_i$  is identically equal to the mean of the outcome  $y$ . This implies that the average of the residuals is equal to zero.

Related to this result, we can show that the prediction  $\hat{y}$  and the residuals are uncorrelated, something that is often called **orthogonality** between  $\hat{y}_i$  and  $e_i$ . We would write this as

$$Cov(\hat{y}, e) = \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - \bar{y})(e_i - \bar{e}) = \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - \bar{y})e_i \quad (3.26)$$

$$= \frac{1}{N} \sum_{i=1}^N \hat{y}_i e_i - \bar{y} \frac{1}{N} \sum_{i=1}^N e_i = 0 \quad (3.27)$$

It's useful to bring back the sample data which generate figure ?? at this point in order to verify these claims:

```
#OUT>      y y_hat error
#OUT> 1  2.09  2.57 -0.48
#OUT> 2  2.79  3.41 -0.62
#OUT> 3  6.49  4.25  2.24
#OUT> 4  1.71  5.10 -3.39
#OUT> 5  9.89  5.94  3.95
#OUT> 6  7.62  6.78  0.83
```

```
#OUT> 7  4.86  7.63 -2.77
#OUT> 8  7.38  8.47 -1.09
#OUT> 9 10.63  9.31  1.32
```

Let's check that these claims are true in this sample of data. We want that

1. The average of  $\hat{y}_i$  to be the same as the mean of  $y$
2. The average of the errors should be zero.
3. Prediction and errors should be uncorrelated.

```
# 1.
all.equal(mean(ss$error), 0)
```

```
#OUT> [1] TRUE
```

```
# 2.
all.equal(mean(ss$y_hat), mean(ss$y))
```

```
#OUT> [1] TRUE
```

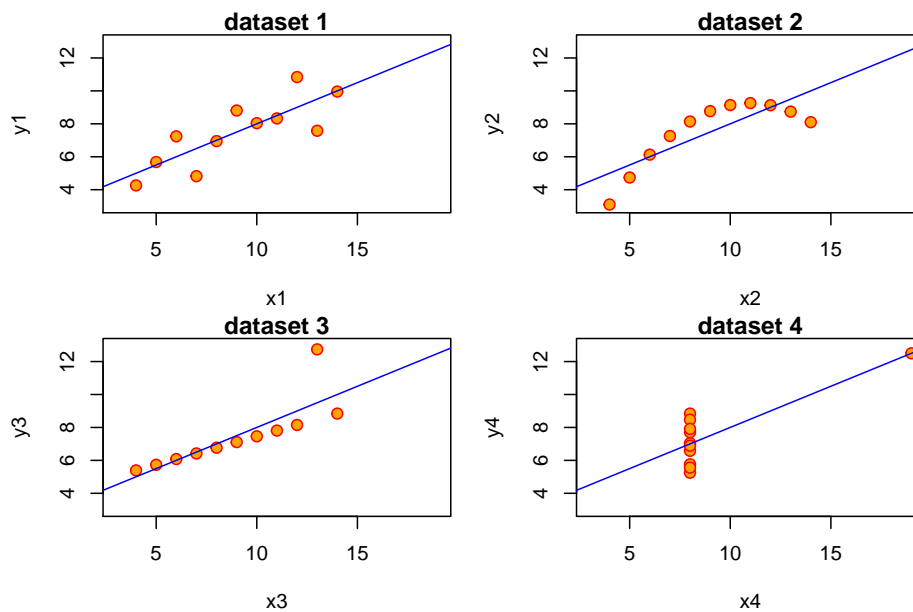
```
# 3.
all.equal(cov(ss$error, ss$y_hat), 0)
```

```
#OUT> [1] TRUE
```

So indeed we can confirm this result with our test dataset. Great!

### 3.4 Correlation, Covariance and Linearity

It is important to keep in mind that Correlation and Covariance relate to a linear relationship between  $x$  and  $y$ . Given how the regression line is estimated by OLS (see just above), you can see that the regression line inherits this property from the Covariance. A famous exercise by Francis Anscombe (1973) illustrates this by constructing 4 different datasets which all have identical **linear** statistics: mean, variance, correlation and regression line are identical. However, the usefulness of the statistics to describe the relationship in the data is not clear.



The important lesson from this example is the following:

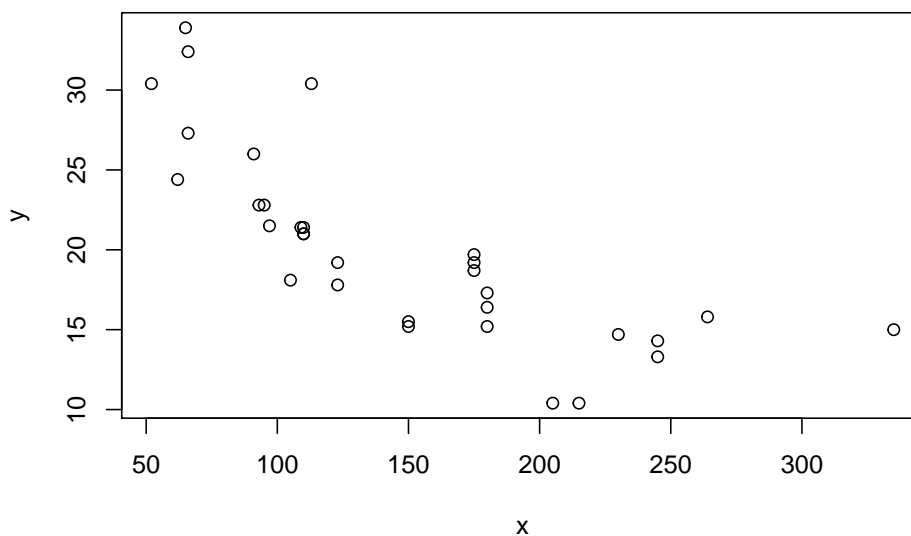
Always **visually inspect** your data, and don't rely exclusively on summary statistics like mean, variance, correlation and regression line. All of those assume a **linear** relationship between the variables in your data.

The mission of Anscombe has been continued recently. As a result of this we can have a look at the `datasauRus` package, which pursues Anscombe's idea through a multitude of funny data sets, all with the same linear statistics. Don't just compute the covariance, or you might actually end up looking at a Dinosaur! What? Type this to find out:

```
launchApp("datasaurus")
aboutApp("datasaurus")
```

### 3.4.1 Non-Linear Relationships in Data

Suppose our data now looks like this:



Putting our previous best line defined in equation (??) as  $y = b_0 + b_1x + e$ , we get something like this:

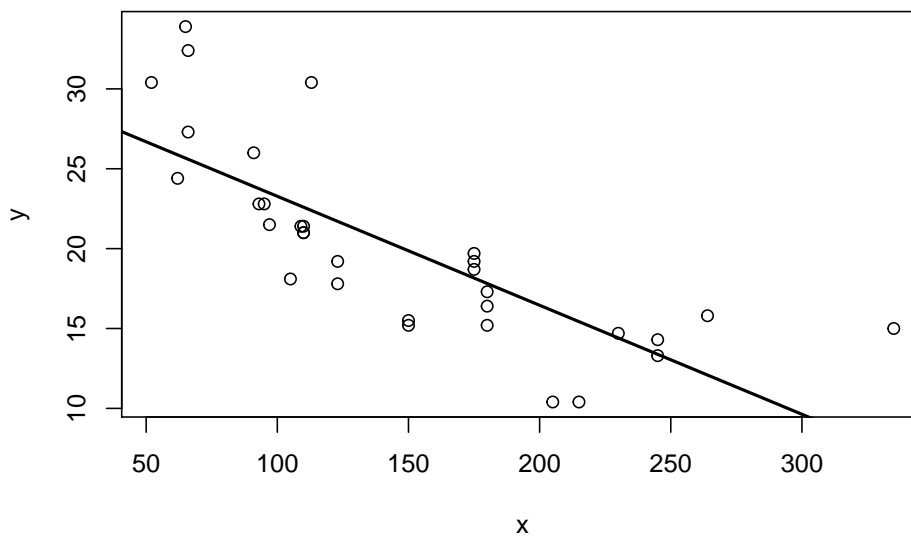


Figure 3.3: Best line with non-linear data?

Somehow when looking at ?? one is not totally convinced that the straight line is a good summary of this relationship. For values  $x \in [50, 120]$  the line seems to low, then again too high, and it completely misses the right boundary. It's easy to address this shortcoming by including higher order terms of an explanatory variable. We would modify (??) to read now



$$y_i = b_0 + b_1x_i + b_2x_i^2 + e_i \quad (3.28)$$

This is a special case of multiple regression, which we will talk about in chapter ???. You can see that there are multiple slope coefficients. For now, let's just see how this performs:

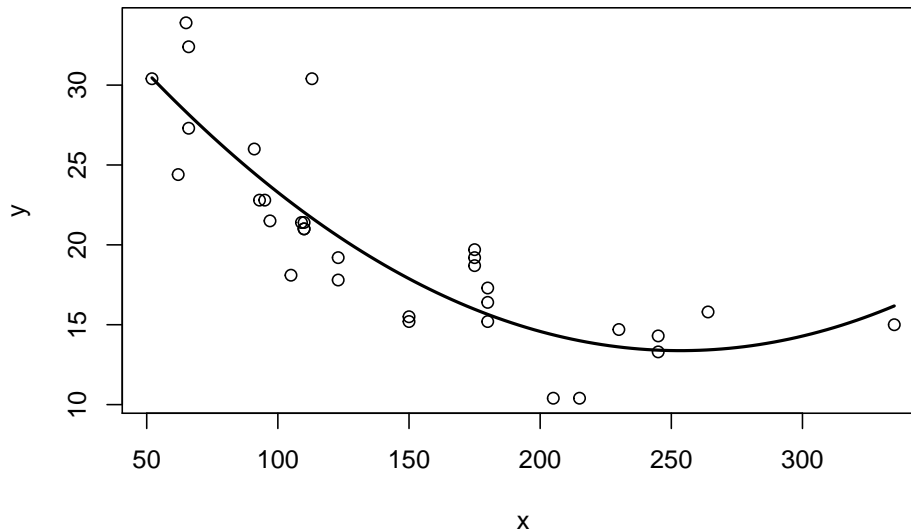


Figure 3.4: Better line with non-linear data!

### 3.5 Analysing $Var(y)$

Analysis of Variance (ANOVA) refers to a method to decompose variation in one variable as a function of several others. We can use this idea on our outcome  $y$ . Suppose we wanted to know the variance of  $y$ , keeping in mind that, by definition,  $y_i = \hat{y}_i + e_i$ . We would write

$$Var(y) = Var(\hat{y} + e) \quad (3.29)$$

$$= Var(\hat{y}) + Var(e) + 2Cov(\hat{y}, e) \quad (3.30)$$

$$= Var(\hat{y}) + Var(e) \quad (3.31)$$

We have seen above in ?? that the covariance between prediction  $\hat{y}$  and error  $e$  is zero, that's why we have  $Cov(\hat{y}, e) = 0$  in (??). What this tells us in words is that we can decompose the variance in the observed outcome  $y$  into a part that relates to variance as explained by the model and a part that comes from

unexplained variation. Finally, we know the definition of variance, and can thus write down the respective formulae for each part:

- $Var(y) = \frac{1}{N} \sum_{i=1}^N (y_i - \bar{y})^2$
- $Var(\hat{y}) = \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - \bar{y})^2$ , because the mean of  $\hat{y}$  is  $\bar{y}$  as we know. Finally,
- $Var(e) = \frac{1}{N} \sum_{i=1}^N e_i^2$ , because the mean of  $e$  is zero.

We can thus formulate how the total variation in outcome  $y$  is apportioned between model and unexplained variation:

The total variation in outcome  $y$  (often called SST, or total sum of squares) is equal to the sum of explained squares (SSE) plus the sum of residuals (SSR). We have thus **SST = SSE + SSR**.

### 3.6 Assessing the Goodness of Fit

In our setup, there exists a convenient measure for how good a particular statistical model fits the data. It is called  $R^2$  (R squared), also called the coefficient of determination. We make use of the just introduced decomposition of variance, and write the formula as

$$R^2 = \frac{\text{variance explained}}{\text{total variance}} = \frac{SSE}{SST} = 1 - \frac{SSR}{SST} \in [0, 1] \quad (3.32)$$

It is easy to see that a good fit is one where the sum of explained squares (SSE) is large relative to the total variation (SST). In such a case, we observe an  $R^2$  close to one. In the opposite case, we will see an  $R^2$  close to zero. Notice that a small  $R^2$  does not imply that the model is useless, just that it explains a small fraction of the observed variation.

### 3.7 An Example: A Log Wage Equation

Let's consider the following example concerning wage data collected in the 1976 Current Population Survey in the USA.<sup>2</sup> We want to investigate the relationship between average hourly earnings, and years of education. Let's start with a plot:

```
data("wage1", package = "wooldridge")    # load data

# a function that returns a plot
plotfun <- function(wage1, log=FALSE, rug = TRUE){
```

<sup>2</sup>This example is close to the vignette of the wooldridge package, whose author I hereby thank for the excellent work.

```

y = wage1$wage
if (log){
  y = log(wage1$wage)
}
plot(y = y,
     x = wage1$educ,
     col = "red", pch = 21, bg = "grey",
     cex=1.25, xaxt="n", frame = FALSE,      # set default x-axis to none
     main = ifelse(log,"log(Wages) vs. Education, 1976","Wages vs. Education, 1976"),
     xlab = "years of education",
     ylab = ifelse(log,"Log Hourly wages","Hourly wages"))
axis(side = 1, at = c(0,6,12,18))           # add custom ticks to x axis
if (rug) rug(wage1$wage, side=2, col="red")  # add `rug` to y axis
}

par(mfcol = c(2,1)) # set up a plot with 2 panels
# plot 1: standard scatter plot
plotfun(wage1)

# plot 2: add a panel with histogram+density
hist(wage1$wage,prob = TRUE, col = "grey", border = "red",
     main = "Histogram of wages and Density",xlab = "hourly wage")
lines(density(wage1$wage), col = "black", lw = 2)

```

Looking at the top panel of figure ??, you notice two things: From the red ticks on the y axis, you see that wages are very concentrated at around 5 USD per hour, with fewer and fewer observations at higher rates; and second, that it seems that the hourly wage seems to increase with higher education levels. The bottom panel reinforces the first point, showing that the estimated pdf (probability density function) shown as a black line has a very long right tail: there are always fewer and fewer, but always larger and larger values of hourly wage in the data.

You have seen this shape of a distribution in the tutorial for chapter 2 already! Do you remember the name of this particular shape of a distribution? (why not type `ScPoEconometrics::runTutorial('chapter2')` to check?)

Let's run a first regression on this data to generate some intuition:

$$\text{wage}_i = b_0 + b_1 \text{educ}_i + e_i \quad (3.33)$$

We use the `lm` function for this purpose as follows:

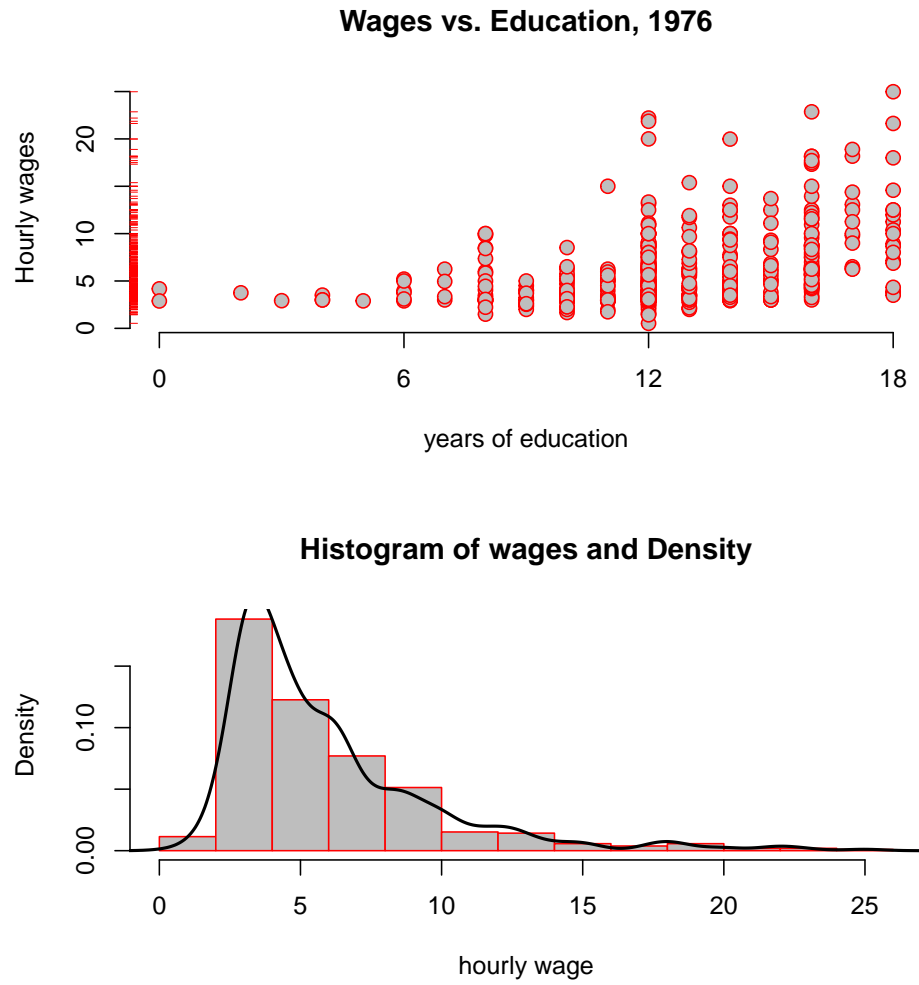


Figure 3.5: Wages vs Education from the wooldridge dataset wage1.

```
hourly_wage <- lm(formula = wage ~ educ, data = wage1)
```

and we can add the resulting regression line to our above plot:

```
plotfun(wage1)
abline(hourly_wage, col = 'black', lw = 2) # add regression line
```

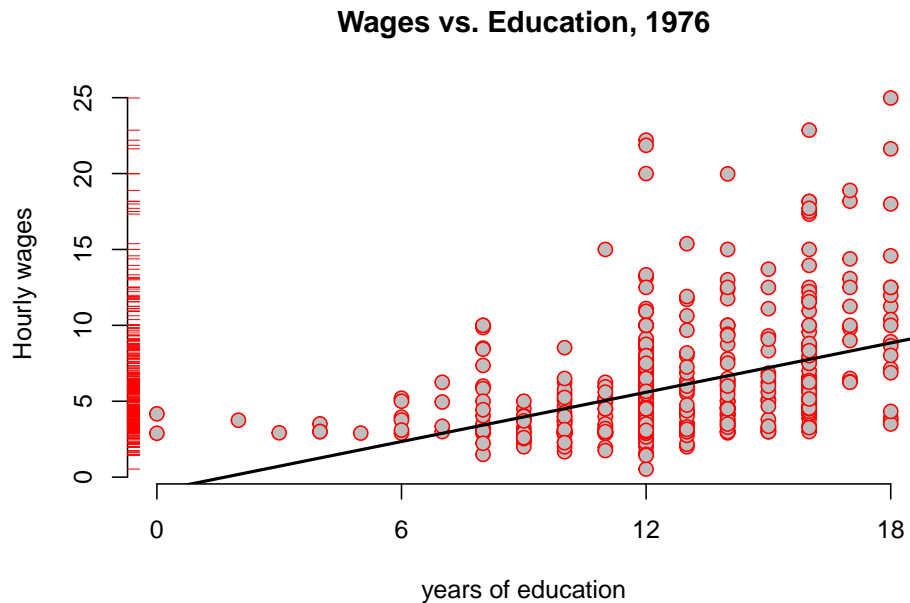


Figure 3.6: Wages vs Education from the wooldridge dataset wage1, with regression

The `hourly_wage` object contains the results of this estimation. We can get a summary of those results with the `summary` method:

```
summary(hourly_wage)
```

```
#OUT>
#OUT> Call:
#OUT> lm(formula = wage ~ educ, data = wage1)
#OUT>
#OUT> Residuals:
#OUT>      Min       1Q   Median       3Q      Max
#OUT> -5.3396 -2.1501 -0.9674  1.1921 16.6085
#OUT>
#OUT> Coefficients:
#OUT>              Estimate Std. Error t value Pr(>|t|)
#OUT> (Intercept) -0.90485    0.68497  -1.321    0.187
#OUT> educ         0.54136    0.05325  10.167 <2e-16 ***
```

```
#OUT> ---
#OUT> Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#OUT>
#OUT> Residual standard error: 3.378 on 524 degrees of freedom
#OUT> Multiple R-squared:  0.1648, Adjusted R-squared:  0.1632
#OUT> F-statistic: 103.4 on 1 and 524 DF,  p-value: < 2.2e-16
```

The main interpretation of this table can be read off the column labelled Estimate, reporting estimated coefficients  $b_0, b_1$ :

1. With zero year of education, the hourly wage is about -0.9 dollars per hour (row named (**Intercept**))
2. Each additional year of education increase hourly wage by 54 cents. (row named **educ**)
3. For example, for 15 years of education, we predict roughly  $-0.9 + 0.541 \times 15 = 7.215$  dollars/h.

### 3.8 Scaling Regressions

Regression estimates  $(b_0, b_1)$  are in the scale of the data. The actual value of the estimates will vary, if we change the scale of the data. The overall fit of the model to the data would not change, however, so that the  $R^2$  statistic would be constant.

Suppose we wanted to use the above estimates to report the effect of years of education on annual wages instead of hourly ones. Let's assume we have full-time workers, 7h per day, 5 days per week, 45 weeks per year. Calling this factor  $\delta = 7 \times 5 \times 45 = 1575$ , we have that  $x$  dollars per hour imply  $x \times \delta = x \times 1575$  dollars per year.

What would be the effect of using  $\tilde{y} = wage \times 1575$  instead of  $y = wage$  as outcome variable on our regression coefficients  $b_0$  and  $b_1$ ? Well, let's try!

% Table created by stargazer v.5.2.2 by Marek Hlavac, Harvard University. E-mail: hlavac at fas.harvard.edu % Date and time: Mon, Aug 24, 2020 - 17:43:30

Let's call the coefficients in the column labelled (1) as  $b_0$  and  $b_1$ , and let's call the ones in column (2)  $b_0^*$  and  $b_1^*$ . In column (1) we see that another year increases hourly wage by 0.54 dollars, as before. In column (2), the corresponding number is 852.64, i.e. another year of education will increase annual wages by 852.64 dollars, on average. Notice however, that  $b_0 \times \delta = -0.9 \times 1575 = -1425.14 = b_0^*$  and that  $b_1 \times \delta = 0.54 \times 1575 = 852.64 = b_1^*$ , that is we just had to multiply both coefficients by the scaling factor applied to original outcome  $y$  to obtain our new coefficients  $b_0^*$  and  $b_1^*$ ! Also, observe that the  $R^2$ s of both regressions are identical! So, really, we did not have to run the regression in column (2) at

Table 3.1: Effect of Scaling on Coefficients

	<i>Dependent variable:</i>	
	wage (1)	annual_wage (2)
educ	0.541*** (0.053)	852.641*** (83.866)
Constant	-0.905 (0.685)	-1,425.141 (1,078.824)
Observations	526	526
R <sup>2</sup>	0.165	0.165
Adjusted R <sup>2</sup>	0.163	0.163
Residual Std. Error (df = 524)	3.378	5,320.963
F Statistic (df = 1; 524)	103.363***	103.363***
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01	

all to make this change: multiplying all coefficients through by  $\delta$  is enough in this case. We keep the identically same fit to the data.

Rescaling the regressors  $x$  is slightly different, but it's easy to work out how different, given the linear nature of the covariance operator, which is part of the OLS estimator. Suppose we rescale  $x$  by the number  $c$ . Then, using the OLS formula in (??), we see that we get new slope coefficient  $b_1^*$  via

$$b_1^* = \frac{Cov(cx, y)}{Var(cx)} \quad (3.34)$$

$$= \frac{cCov(x, y)}{c^2Var(x)} \quad (3.35)$$

$$= \frac{1}{c}b_1. \quad (3.36)$$

As for the intercept, and by using (??)

$$b_0^* = \bar{y} - b_1^* \frac{1}{N} \sum_{i=1}^N c \cdot x_i \quad (3.37)$$

$$= \bar{y} - b_1^* \frac{c}{N} \sum_{i=1}^N x_i \quad (3.38)$$

$$= \bar{y} - \frac{1}{c} b_1^* c * \bar{x} \quad (3.39)$$

$$= \bar{y} - b_1 * \bar{x} \quad (3.40)$$

$$= b_0 \quad (3.41)$$

That is, we change the slope by the inverse of the scaling factor applied to regressor  $x$ , but the intercept is unaffected from this. You should play around for a while with our rescaling app to get a feeling for this:

```
library(ScPoEconometrics)
launchApp('Rescale')
```

### 3.9 A Particular Rescaling: The log Transform

The natural logarithm is a particularly important transformation that we often encounter in economics. Why would we transform a variable with the log function to start with?

1. Several important economic variables (like wages, city size, firm size, etc) are approximately log-normally distributed. By transforming them with the log, we obtain an approximately normally distributed variable, which has desirable properties for our regression.
2. Applying the log reduces the impact of outliers.
3. The transformation allows for a convenient interpretation in terms of percentage changes of the outcome variable.

Let's investigate this issue in our running example by transforming the wage data above. Look back at the bottom panel of figure ?? : Of course you saw immediately that this looked a lot like a log-normal distribution, so point 1. above applies. We modify the left hand side of equation (??):

$$\log(\text{wage}_i) = b_0 + b_1 \text{educ}_i + e_i \quad (3.42)$$

Let's use the `update` function to modify our previous regression model:

```
log_hourly_wage = update(hourly_wage, log(wage) ~ ., data = wage1)
```

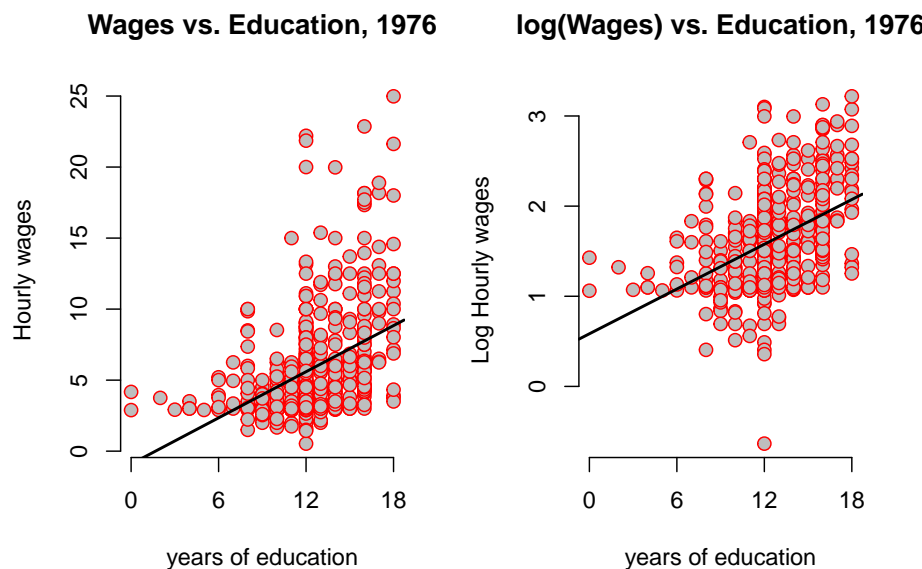


The `update` function takes an existing `lm` object, like `hourly_wage` here, and updates the `formula`. Here the `.` on the right hand side means leave unchanged (so the RHS stays unchanged). How do our pictures change?

```
par(mfrow = c(1,2))

plotfun(wage1,rug = FALSE)
abline(hourly_wage, col = 'black', lw = 2) # add regression line

plotfun(wage1,log = TRUE, rug = FALSE)
abline(log_hourly_wage, col = 'black', lw = 2) # add regression line
```



```
par(mfrow = c(1,1))
```

It looks as if the regression line has the same slope, but beware of the different scales of the y-axis! You can clearly see that all y-values have been compressed by the log transformation. The log case behaves differently from our scaling by a constant number case above because it is a nonlinear function. Let's compare the output between both models:

```
% Table created by stargazer v.5.2.2 by Marek Hlavac, Harvard University. E-
mail: hlavac at fas.harvard.edu % Date and time: Mon, Aug 24, 2020 - 17:43:31
```

The interpretation of the transformed model in column (2) is now the following:

Table 3.2: Log Transformed Equation

	<i>Dependent variable:</i>	
	wage (1)	log(wage) (2)
educ	0.541*** (0.053)	0.083*** (0.008)
Constant	-0.905 (0.685)	0.584*** (0.097)
Observations	526	526
R <sup>2</sup>	0.165	0.186
Adjusted R <sup>2</sup>	0.163	0.184
Residual Std. Error (df = 524)	3.378	0.480
F Statistic (df = 1; 524)	103.363***	119.582***
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01	

We call a regression of the form  $\log(y) = b_0 + b_1x + u$  a log-level specification, because we regressed the log of a variable on the level (i.e not the log!) of another variable. Here, the impact of increasing  $x$  by one unit is to increase  $y$  by  $100 \times b_1$  **percent**. In our example: an additional year of education will increase hourly wages by 8.3%. Notice that this is very different from saying ...increases log hourly wages by 8.3%, which is wrong.

Notice that the  $R^2$  slightly improved, so have a better fit to the data. This is due the fact that the log compressed large outlier values. Whether we apply the *log* to left or right-hand side variables makes a difference, as outlined in this important table:

#### Common Regression Specifications

Specification	Outcome Var	Regressor	Interpretation of $b_1$	Comment
Level-level	$y$	$x$	$\Delta y = b_1 \Delta x$	Standard
Level-log	$y$	$\log(x)$	$\Delta y = \frac{b_1}{100} \Delta x$	less frequent
Log-level	$\log(y)$	$x$	$\% \Delta y = (100b_1) \Delta x$	Semi-elasticity
Log-Log	$\log(y)$	$\log(x)$	$\% \Delta y = \% \Delta b_1 x$	Elasticity

You may remember from your introductory micro course what the definition of the elasticity of  $y$  with respect to  $x$  is: This number tells us by how many percent  $y$  will change, if we change  $x$  by one percent. Let's look at another example from the `wooldridge` package of datasets, this time concerning CEO salaries and their relationship with company sales.

```
data("ceosal1", package = "wooldridge")
par(mfrow = c(1,2))
plot(salary ~ sales, data = ceosal1, main = "Sales vs Salaries", xaxt = "n", frame = FALSE)
axis(1, at = c(0,40000, 80000))
rug(ceosal1$salary, side = 2)
rug(ceosal1$sales, side = 1)
plot(log(salary) ~ log(sales), data = ceosal1, main = "Log(Log) Sales vs Salaries", xaxt = "n", frame = FALSE)
```



Figure 3.7: The effect of log-transforming highly skewed data.

Referring back at table ??, here we have a log-log specification. Therefore we interpret this regression as follows:

In a log-log equation, the slope coefficient  $b_1$  is the elasticity of  $y$  with respect to changes in  $x$ . Here: A 1% increase in sales is associated to a 0.26% increase in CEO salaries. Note, again, that there is no log in this statement.



## Chapter 4

# Multiple Regression

We can extend the discussion from chapter ?? to more than one explanatory variable. For example, suppose that instead of only  $x$  we now had  $x_1$  and  $x_2$  in order to explain  $y$ . Everything we've learned for the single variable case applies here as well. Instead of a regression line, we now get a regression plane, i.e. an object representable in 3 dimensions:  $(x_1, x_2, y)$ . As an example, suppose we wanted to explain how many miles per gallon (**mpg**) a car can travel as a function of its horse power (**hp**) and its weight (**wt**). In other words we want to estimate the equation

$$mpg_i = b_0 + b_1 hp_i + b_2 wt_i + e_i \quad (4.1)$$

on our built-in dataset of cars (**mtcars**):

```
head(subset(mtcars, select = c(mpg, hp, wt)))
```

```
#OUT>           mpg  hp   wt
#OUT> Mazda RX4    21.0 110 2.620
#OUT> Mazda RX4 Wag 21.0 110 2.875
#OUT> Datsun 710    22.8  93 2.320
#OUT> Hornet 4 Drive 21.4 110 3.215
#OUT> Hornet Sportabout 18.7 175 3.440
#OUT> Valiant      18.1 105 3.460
```

How do you think **hp** and **wt** will influence how many miles per gallon of gasoline each of those cars can travel? In other words, what do you expect the signs of  $b_1$  and  $b_2$  to be?

With two explanatory variables as here, it is still possible to visualize the regression plane, so let's start with this as an answer. The OLS regression plane through this dataset looks like in figure ??:



Figure 4.1: Multiple Regression - a plane in 3D. The red lines indicate the residual for each observation.

This visualization shows a couple of things: the data are shown with red points and the grey plane is the one resulting from OLS estimation of equation (??). You should realize that this is exactly the same story as told in figure ?? - just in three dimensions!

Furthermore, multiple regression refers the fact that there could be more than two regressors. In fact, you could in principle have  $K$  regressors, and our theory developed so far would still be valid:

$$\hat{y}_i = b_0 + b_1x_{1i} + b_2x_{2i} + \cdots + b_Kx_{Ki} \quad (4.2)$$

$$e_i = y_i - \hat{y}_i \quad (4.3)$$

Just as before, the least squares method chooses numbers  $(b_0, b_1, \dots, b_K)$  to as to minimize SSR, exactly as in the minimization problem for the one regressor case seen in (??).

## 4.1 All Else Equal

We can see from the above plot that cars with more horse power and greater weight, in general travel fewer miles per gallon of combustible. Hence, we observe

a plane that is downward sloping in both the weight and horse power directions. Suppose now we wanted to know impact of hp on mpg in isolation, so as if we could ask

Keeping the value of *wt* fixed for a certain car, what would be the impact on *mpg* be if we were to increase **only** its *hp*? Put differently, keeping **all else equal**, what's the impact of changing *hp* on *mpg*?

We ask this kind of question all the time in econometrics. In figure ?? you clearly see that both explanatory variables have a negative impact on the outcome of interest: as one increases either the horse power or the weight of a car, one finds that miles per gallon decreases. What is kind of hard to read off is how negative an impact each variable has in isolation.

As a matter of fact, the kind of question asked here is so common that it has got its own name: we'd say "ceteris paribus, what is the impact of hp on mpg". ceteris paribus is latin and means the others equal, i.e. all other variables fixed. In terms of our model in (??), we want to know the following quantity:

$$\frac{\partial mpg_i}{\partial hp_i} = b_1 \quad (4.4)$$

The  $\partial$  sign denotes a partial derivative of the function describing mpg with respect to the variable hp. It measures how the value of mpg changes, as we change the value of hp ever so slightly. In our context, this means: keeping all other variables fixed, what is the effect of hp on mpg?. We call the value of coefficient  $b_1$  therefore also the partial effect of hp on mpg. In terms of our dataset, we use R to run the following **multiple regression**:

```
#OUT>
#OUT> Call:
#OUT> lm(formula = mpg ~ wt + hp, data = mtcars)
#OUT>
#OUT> Residuals:
#OUT>      Min       1Q   Median       3Q      Max
#OUT> -3.941 -1.600 -0.182  1.050  5.854
#OUT>
#OUT> Coefficients:
#OUT>              Estimate Std. Error t value Pr(>|t|)
#OUT> (Intercept)  37.22727     1.59879   23.285 < 2e-16 ***
#OUT> wt         -3.87783     0.63273   -6.129 1.12e-06 ***
#OUT> hp          -0.03177     0.00903   -3.519  0.00145 **
#OUT> ---
#OUT> Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#OUT>
#OUT> Residual standard error: 2.593 on 29 degrees of freedom
```

```
#OUT> Multiple R-squared:  0.8268, Adjusted R-squared:  0.8148
#OUT> F-statistic: 69.21 on 2 and 29 DF,  p-value: 9.109e-12
```

From this table you see that the coefficient on `wt` has value -3.87783. You can interpret this as follows:

Holding all other variables fixed at their observed values - or ceteris paribus - a one unit increase in *wt* implies a -3.87783 units change in *mpg*. In other words, increasing the weight of a car by 1000 pounds (lbs), will lead to 3.88 miles less travelled per gallon. Similarly, a car with one additional horse power means that we will travel 0.03177 fewer miles per gallon of gasoline, all else (i.e. *wt*) equal.

## 4.2 Multicollinearity

One important requirement for multiple regression is that the data be **not linearly dependent**: Each variable should provide at least some new information for the outcome, and it cannot be replicated as a linear combination of other variables. Suppose that in the example above, we had a variable `wtplus` defined as `wt + 1`, and we included this new variable together with `wt` in our regression. In this case, `wtplus` provides no new information. It's enough to know *wt*, and add 1 to it. In this sense, `wt_plus` is a redundant variable and should not be included in the model. Notice that this holds only for linearly dependent variables - nonlinear transformations (like for example  $wt^2$ ) are exempt from this rule. Here is why:

$$y = b_0 + b_1 wt + b_2 wtplus + e \quad (4.5)$$

$$= b_0 + b_1 wt + b_2 (wt + 1) + e \quad (4.6)$$

$$= (b_0 + b_2) + wt(b_1 + b_2) + e \quad (4.7)$$

This shows that we cannot identify the regression coefficients in case of linearly dependent data. Variation in the variable `wt` identifies a different coefficient, say  $\gamma = b_1 + b_2$ , from what we actually wanted: separate estimates for  $b_1, b_2$ .

We cannot have variables which are linearly dependent, or perfectly colinear. This is known as the **rank condition**. In particular, the condition dictates that we need at least  $N \geq K + 1$ , i.e. more observations than coefficients. The greater the degree of linear dependence amongst our explanatory variables, the less information we can extract from them, and our estimates becomes less precise.



### 4.3 Log Wage Equation

Let's go back to our previous example of the relationship between log wages and education. How does this relationship change if we also think that experience in the labor market has an impact, next to years of education? Here is a picture:

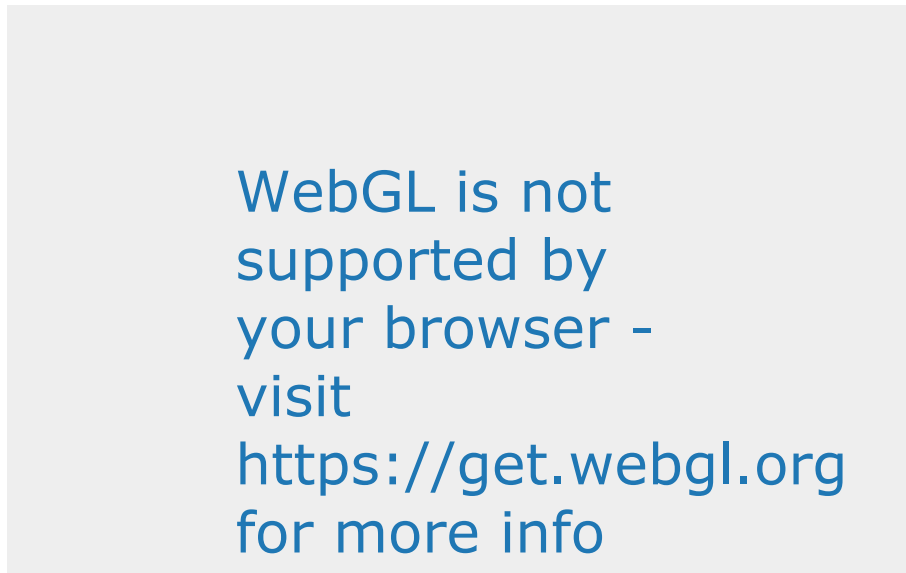


Figure 4.2: Log wages vs education and experience in 3D.

Let's add even more variables! For instance, what's the impact of experience in the labor market, and time spent with the current employer? Let's first look at how those variables co-vary with each other:

```
cmat = round(cor(subset(wage1, select = c(lwage, educ, exper, tenure))), 2) # correlation matrix
corrplot::corrplot(cmat, type = "upper", method = "ellipse")
```

The way to read the so-called correlation plot in figure ?? is straightforward: each row illustrates the correlation of a certain variable with the other variables. In this example both the shape of the ellipse in each cell as well as their color coding tell us how strongly two variables correlate. Let us put this into a regression model now:

```
educ_only <- lm(lwage ~ educ, data = wage1)
educ_exper <- lm(lwage ~ educ + exper, data = wage1)
log_wages <- lm(lwage ~ educ + exper + tenure, data = wage1)
stargazer::stargazer(educ_only, educ_exper, log_wages, type = if (knitr::is_latex_output()) "latex")
```

% Table created by stargazer v.5.2.2 by Marek Hlavac, Harvard University. E-mail: hlavac at fas.harvard.edu % Date and time: Mon, Aug 24, 2020 - 17:43:46

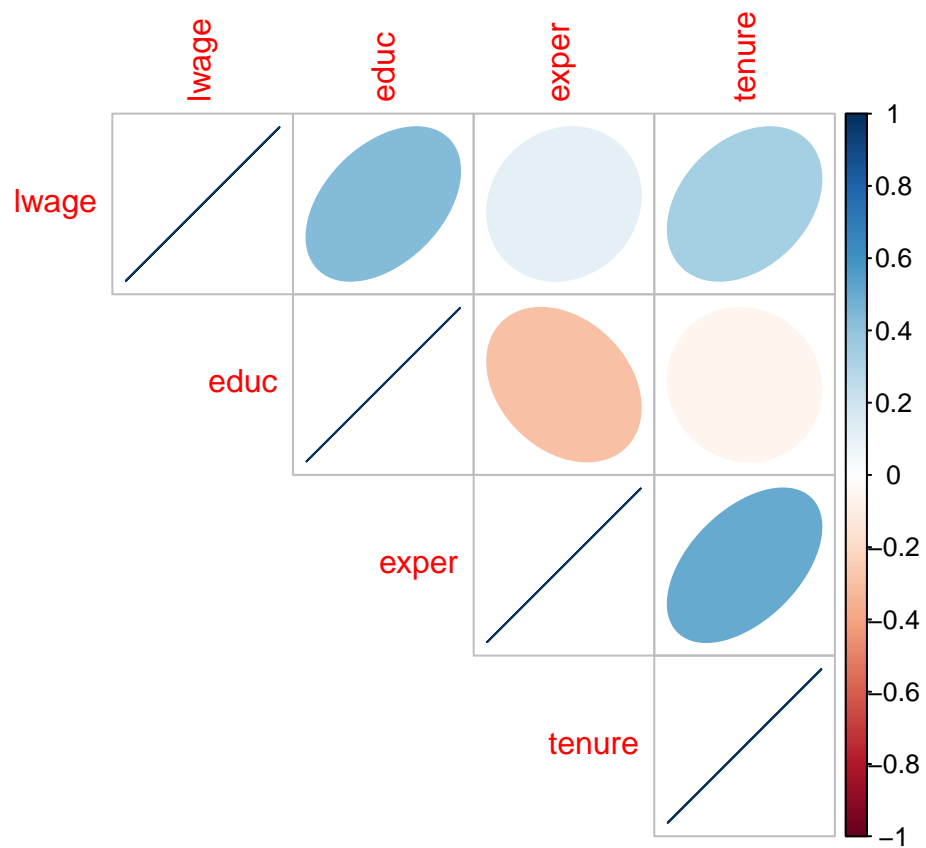


Figure 4.3: correlation plot

Table 4.1

	<i>Dependent variable:</i>		
	lwage		
	(1)	(2)	(3)
educ	0.083*** (0.008)	0.098*** (0.008)	0.092*** (0.007)
exper		0.010*** (0.002)	0.004** (0.002)
tenure			0.022*** (0.003)
Constant	0.584*** (0.097)	0.217** (0.109)	0.284*** (0.104)
Observations	526	526	526
R <sup>2</sup>	0.186	0.249	0.316
Adjusted R <sup>2</sup>	0.184	0.246	0.312
Residual Std. Error	0.480 (df = 524)	0.461 (df = 523)	0.441 (df = 522)
F Statistic	119.582*** (df = 1; 524)	86.862*** (df = 2; 523)	80.391*** (df = 3; 522)
<i>Note:</i>		*p<0.1; **p<0.05; ***p<0.01	

Column (1) refers to model (??) from the previous chapter, where we only had `educ` as a regressor: we obtain an  $R^2$  of 0.186. Column (2) is the model that generated the plane in figure ?? above. (3) is the model with three regressors. You can see that by adding more regressors, the quality of our fit increases, as more of the variation in  $y$  is now accounted for by our model. You can also see that the values of our estimated coefficients keeps changing as we move from left to right across the columns. Given the correlation structure shown in figure ??, it is only natural that this is happening: We see that `educ` and `exper` are negatively correlated, for example. So, if we omit `exper` from the model in column (1), `educ` will reflect part of this correlation with `exper` by a lower estimated value. By directly controlling for `exper` in column (2) we get an estimate of the effect of `educ` net of whatever effect `exper` has in isolation on the outcome variable. We will come back to this point later on.

## 4.4 How To Make Predictions

So suppose we have a model like

$$\text{lwage} = b_0 + b_1(\text{educ}) + b_2(\text{exper}) + b_3(\text{tenure}) + \epsilon$$

How could we use this to make a prediction of log wages, given some new data? Remember that the OLS procedure gives us estimates for the values  $b_0, b_1, b_2, b_3$ . With those in hand, it is straightforward to make a prediction about the conditional mean of the outcome - just plug in the desired numbers for `educ`, `exper` and `tenure`. Suppose you want to know what the mean of `lwage` is conditional on `educ = 10`, `exper=4` and `tenure = 2`. You'd do

$$E[\text{lwage} | \text{educ} = 10, \text{exper} = 4, \text{tenure} = 2] = b_0 + b_1 10 + b_2 4 + b_3 2 \quad (4.8)$$

$$= 1.27. \quad (4.9)$$

I computed the last line directly with

```
x = c(1,10,4,2) # 1 for intercept
pred = coef(log_wages) %*% x
```

but R has a more complete prediction interface, using the function `predict`. For starters, you can predict the model on all data points which were contained in the dataset we used for estimation, i.e. `wage1` in our case:

```
head(predict(log_wages)) # first 6 observations of wage1 as predicted by our model
```

```
#OUT>      1      2      3      4      5      6
#OUT> 1.304921 1.523506 1.304921 1.819802 1.461690 1.970451
```

Often you want to add that prediction to the original dataset:

```
wage_prediction = cbind(wage1, prediction = predict(log_wages))
head(wage_prediction[, c("lwage", "educ", "exper", "tenure", "prediction")])
```

```
#OUT>      lwage educ exper tenure prediction
#OUT> 1 1.131402  11     2      0   1.304921
#OUT> 2 1.175573  12    22      2   1.523506
#OUT> 3 1.098612  11     2      0   1.304921
#OUT> 4 1.791759   8    44     28   1.819802
#OUT> 5 1.667707  12     7      2   1.461690
#OUT> 6 2.169054  16     9      8   1.970451
```

You'll remember that we called the distance in prediction and observed outcome our residual  $e$ . Well here this is just `lwage - prediction`. Indeed,  $e$  is such an important quantity that R has a convenient method to compute  $y - \hat{y}$  from an `lm` object directly - the method `resid`. Let's add another column to `wage_prediction`:

```
wage_prediction = cbind(wage_prediction, residual = resid(log_wages))
head(wage_prediction[, c("lwage", "educ", "exper", "tenure", "prediction", "residual")])
```

```
#OUT>      lwage educ exper tenure prediction  residual
#OUT> 1 1.131402  11     2      0   1.304921 -0.17351850
#OUT> 2 1.175573  12    22      2   1.523506 -0.34793289
#OUT> 3 1.098612  11     2      0   1.304921 -0.20630832
#OUT> 4 1.791759   8    44     28   1.819802 -0.02804286
#OUT> 5 1.667707  12     7      2   1.461690  0.20601725
#OUT> 6 2.169054  16     9      8   1.970451  0.19860271
```

Using the data in `wage_prediction`, you should now check for yourself what we already know about  $\hat{y}$  and  $e$  from section ??:

1. What is the average of the vector `residual`?
2. What is the average of `prediction`?
3. How does this compare to the average of the outcome `lwage`?
4. What is the correlation between `prediction` and `residual`?



## Chapter 5

# Categorical Variables

Up until now, we have encountered only examples with continuous variables  $x$  and  $y$ , that is,  $x, y \in \mathbb{R}$ , so that a typical observation could have been  $(y_i, x_i) = (1.5, 5.62)$ . There are many situations where it makes sense to think about the data in terms of categories, rather than continuous numbers. For example, whether an observation  $i$  is male or female, whether a pixel on a screen is black or white, and whether a good was produced in France, Germany, Italy, China or Spain are all categorical classifications of data.

Probably the simplest type of categorical variable is the binary, boolean, or just dummy variable. As the name suggests, it can take on only two values, 0 and 1, or TRUE and FALSE.

### 5.1 The Binary Regressor Case

Even though this is an extremely parsimonious way of encoding that, it is a very powerful tool that allows us to represent that a certain observation  $i$  **is a member** of a certain category  $j$ . For example, let's imagine we have income data on males and females, and we would create a variable called `is.male` that is `TRUE` whenever  $i$  is male, `FALSE` otherwise, and similarly for women. For example, to encode whether subject  $i$  is male, one could do this:

$$\text{is.male}_i = \begin{cases} 1 & \text{if } i \text{ is male} \\ 0 & \text{if } i \text{ is not male.} \end{cases} ,$$

and similarly for females, we'd have

$$\text{is.female}_i = \begin{cases} 1 & \text{if } i \text{ is female} \\ 0 & \text{if } i \text{ is not female.} \end{cases}$$

By definition, we have just introduced a linear dependence into our dataset. It will always be true that  $\text{is.male}_i + \text{is.female}_i = 1$ . This is because dummy variables are based on data being mutually exclusively categorized - here, you are either male or female.<sup>1</sup> This should immediately remind you of section ?? where we introduced multicollinearity. A regression of income on both of our variables like this

$$y_i = b_0 + b_1 \text{is.female}_i + b_2 \text{is.male}_i + e_i$$

would be invalid because of perfect colinearity between  $\text{is.female}_i$  and  $\text{is.male}_i$ . The solution to this is pragmatic and simple:

In dummy variable regressions, we remove one category from the regression (for example here: `is.male`) and call it the reference category. The effect of being male is absorbed in the intercept. The coefficient on the remaining categories measures the difference in mean outcome with respect to the reference category.

Now let's try this out. We start by creating the female indicator as above,

$$\text{is.female}_i = \begin{cases} 1 & \text{if } i \text{ is female} \\ 0 & \text{if } i \text{ is not female.} \end{cases}$$

and let's suppose that  $y_i$  is a measure of  $i$ 's annual labor income. Our model is

$$y_i = b_0 + b_1 \text{is.female}_i + e_i \tag{5.1}$$

and here is how we estimate this in R:

```
# x = sample(x = c(0, 1), size = n, replace = T)
dta$is.female = factor(x) # convert x to factor
dummy_reg = lm(y~is.female,dta)
summary(dummy_reg)
```

```
#OUT>
```

```
#OUT> Call:
```

```
#OUT> lm(formula = y ~ is.female, data = dta)
```

---

<sup>1</sup>There are transgender individuals where this example will not apply.



```

#OUT>
#OUT> Residuals:
#OUT>      Min       1Q   Median       3Q      Max
#OUT> -2.4253 -0.6551  0.1321  0.7217  2.8129
#OUT>
#OUT> Coefficients:
#OUT>              Estimate Std. Error t value Pr(>|t|)
#OUT> (Intercept)    2.0568     0.2125   9.680 7.27e-13 ***
#OUT> is.female1     -3.0169     0.3203  -9.418 1.74e-12 ***
#OUT> ---
#OUT> Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#OUT>
#OUT> Residual standard error: 1.124 on 48 degrees of freedom
#OUT> Multiple R-squared:  0.6489, Adjusted R-squared:  0.6416
#OUT> F-statistic: 88.7 on 1 and 48 DF,  p-value: 1.737e-12

```

Notice that R displays the level of the factor to which coefficient  $b_1$  belongs here, i.e. `is.female1` means this coefficient is on level `is.female = 1` - the reference level is `is.female = 0`, and it has no separate coefficient. Also interesting is that  $b_1$  is equal to the difference in conditional means between male and female

$$b_1 = E[y|\text{is.female} = 1] - E[y|\text{is.female} = 0] = -3.0169.$$

A dummy variable measures the difference or the offset in the mean of the response variable,  $E[y]$ , **conditional** on  $x$  belonging to some category - relative to a baseline category. In our artificial example, the coefficient  $b_1$  informs us that women earn on average 3.756 units less than men.

It is instructive to reconsider this example graphically:

In figure ?? we see that this regression simplifies to the straight line connecting the mean, or the expected value of  $y$  when `is.femalei = 0`, i.e.  $E[y|\text{is.female}_i = 0]$ , to the mean when `is.femalei = 1`, i.e.  $E[y|\text{is.female}_i = 1]$ . It is useful to remember that the unconditional mean of  $y$ , i.e.  $E[y]$ , is going to be the result of regressing  $y$  only on an intercept, illustrated by the blue line. This line will always lie in between both conditional means. As indicated by the red arrow, the estimate of the coefficient on the dummy,  $b_1$ , is equal to the difference in conditional means for both groups. You should look at our app now to deepen your understanding of what's going on here:

```

library(ScPoEconometrics)
launchApp("reg_dummy")

```

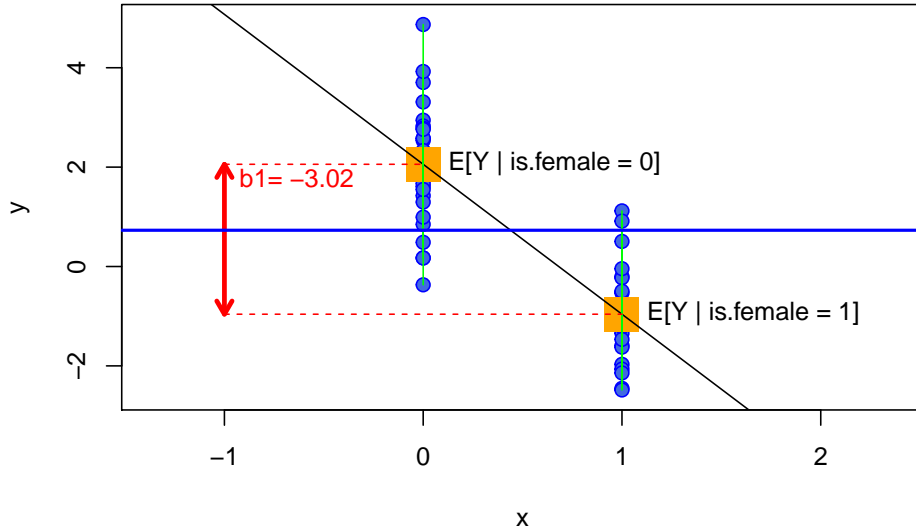


Figure 5.1: regressing  $y \in \mathbb{R}$  on  $\text{is.female}_i \in \{0, 1\}$ . The blue line is  $E[y]$ , the red arrow is the size of  $b_1$ . Which is the same as the slope of the regression line in this case and the difference in conditional means!

## 5.2 Dummy and Continuous Variables

What happens if there are more predictors than just the dummy variable in a regression? For example, what if instead we had

$$y_i = b_0 + b_1 \text{is.female}_i + b_2 \text{exper}_i + e_i \quad (5.2)$$

where  $\text{exper}_i$  would measure years of experience in the labor market? As above, the dummy variable acts as an intercept shifter. We have

$$y_i = \begin{cases} b_0 + b_1 + b_2 \times \text{exper}_i + e_i & \text{if is.female}=1 \\ b_0 + \quad + b_2 \times \text{exper}_i + e_i & \text{if is.female}=0 \end{cases} \quad (5.3)$$

so that the intercept is  $b_0 + b_1$  for women but  $b_0$  for men. We will see this in the real-world example below, but for now let's see the effect of switching the dummy on and off in this app:

```
library(ScPoEconometrics)
launchApp("reg_dummy_example")
```

## 5.3 Categorical Variables in R: factor

R has extensive support for categorical variables built-in. The relevant data type representing a categorical variable is called **factor**. We encountered them as basic data types in section ?? already, but it is worth repeating this here. We have seen that a factor categorizes a usually small number of numeric values by labels, as in this example which is similar to what I used to create regressor `is.female` for the above regression:

```
is.female = factor(x = c(0,1,1,0), labels = c(FALSE,TRUE))
is.female
```

```
#OUT> [1] FALSE TRUE  TRUE  FALSE
#OUT> Levels: FALSE TRUE
```

You can see the result is a vector object of type **factor** with 4 entries, whereby 0 is represented as **FALSE** and 1 as **TRUE**. An other example could be if we wanted to record a variable sex instead, and we could do

```
sex = factor(x = c(0,1,1,0), labels = c("male","female"))
sex
```

```
#OUT> [1] male  female female male
#OUT> Levels: male female
```

You can see that this is almost identical, just the labels are different.

### 5.3.1 More Levels

We can go beyond binary categorical variables such as **TRUE** vs **FALSE**. For example, suppose that  $x$  measures educational attainment, i.e. it is now something like  $x_i \in \{\text{high school, some college, BA, MSc}\}$ . In R parlance, high school, some college, BA, MSc are the **levels of factor**  $x$ . A straightforward extension of the above would dictate to create one dummy variable for each category (or level), like

$$\begin{aligned}\text{has.HS}_i &= \mathbf{1}[x_i == \text{high school}] \\ \text{has.someCol}_i &= \mathbf{1}[x_i == \text{some college}] \\ \text{has.BA}_i &= \mathbf{1}[x_i == \text{BA}] \\ \text{has.MSc}_i &= \mathbf{1}[x_i == \text{MSc}]\end{aligned}$$

but you can see that this is cumbersome. There is a better solution for us available:

```
factor(x = c(1,1,2,4,3,4), labels = c("HS", "someCol", "BA", "MSc"))
```

```
#OUT> [1] HS      HS      someCol MSc      BA      MSc
#OUT> Levels: HS someCol BA MSc
```

Notice here that R will apply the labels in increasing order the way you supplied it (i.e. a numerical value 4 will correspond to “MSc”, no matter the ordering in `x`.)

### 5.3.2 Log Wages and Dummies

The above developed `factor` terminology fits neatly into R’s linear model fitting framework. Let us illustrate the simplest use by way of example.

Going back to our wage example, let’s say that a worker’s wage depends on their education as well as their sex:

$$\ln w_i = b_0 + b_1 educ_i + b_2 female_i + e_i \quad (5.4)$$

```
data("wage1", package = "wooldridge")
wage1$female = as.factor(wage1$female) # convert 0-1 to factor
lm_w = lm(lwage ~ educ, data = wage1)
lm_w_sex = lm(lwage ~ educ + female, data = wage1)
stargazer::stargazer(lm_w, lm_w_sex, type = if (knitr::is_latex_output()) "latex" else "
```

% Table created by stargazer v.5.2.2 by Marek Hlavac, Harvard University. E-mail: hlavac at fas.harvard.edu % Date and time: Mon, Aug 24, 2020 - 17:43:47

We know the results from column (1) very well by now. How does the relationship change if we include the `female` indicator? Remember from above that `female` is a `factor` with two levels, 0 and 1, where 1 means that’s a female. We see in the above output that R included a regressor called `female1`. This is a combination of the variable name `female` and the level which was included in the regression. In other words, R chooses a reference category (by default the first of all levels by order of appearance), which is excluded - here this is `female==0`. The interpretation is that  $b_2$  measures the effect of being female relative to being male. R automatically creates a dummy variable for each potential level, excluding the first category.

Figure ?? illustrates this. The left panel is our previous model. The right panel adds the `female` dummy. You can see that both male and female have the same upward sloping regression line. But you can also see that there is a parallel downward shift from male to female line. The estimate of  $b_2 = -0.36$  is the size of the downward shift.

Table 5.1

	<i>Dependent variable:</i>	
	lwage	
	(1)	(2)
educ	0.083*** (0.008)	0.077*** (0.007)
female1		-0.361*** (0.039)
Constant	0.584*** (0.097)	0.826*** (0.094)
Observations	526	526
R <sup>2</sup>	0.186	0.300
Adjusted R <sup>2</sup>	0.184	0.298
Residual Std. Error	0.480 (df = 524)	0.445 (df = 523)
F Statistic	119.582*** (df = 1; 524)	112.189*** (df = 2; 523)

Note:

\*p&lt;0.1; \*\*p&lt;0.05; \*\*\*p&lt;0.01

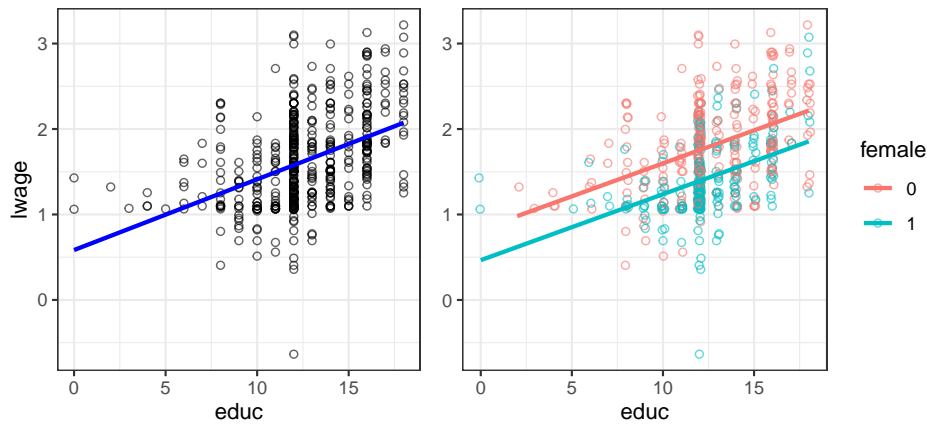


Figure 5.2: log wage vs educ. Right panel with female dummy.

## 5.4 Interactions

Sometimes it is useful to let the slope of a certain variable to be dependent on the value of another regressor. For example consider a model for the sales prices of houses, where **area** is the livable surface of the property, and **age** is its age:

$$\log(\text{price}) = b_0 + b_1 \text{area} + b_2 \text{age} + b_3 (\text{area} \times \text{age}) + e \quad (5.5)$$

In that model, the partial effect of **area** on  $\log(\text{price})$ , keeping all other variables fixed, is

$$\frac{\partial \log(\text{price})}{\partial \text{area}} = b_1 + b_3 (\text{age}) \quad (5.6)$$

If we find that  $b_3 > 0$  in a regression, we conclude that the size of a house values more in older houses. We call  $b_3$  the **interaction effect** between area and age. Let's look at that regression model now.

```
data(hprice3, package = "wooldridge")
summary(lm(lprice ~ area*age, data = hprice3))
```

```
#OUT>
#OUT> Call:
#OUT> lm(formula = lprice ~ area * age, data = hprice3)
#OUT>
#OUT> Residuals:
#OUT>      Min       1Q   Median       3Q      Max
#OUT> -1.27226 -0.16538 -0.00298  0.20673  0.83985
#OUT>
#OUT> Coefficients:
#OUT>              Estimate Std. Error t value Pr(>|t|)
#OUT> (Intercept)  1.071e+01  6.633e-02 161.448 < 2e-16 ***
#OUT> area         3.647e-04  2.875e-05  12.686 < 2e-16 ***
#OUT> age        -7.377e-03  1.358e-03  -5.434 1.1e-07 ***
#OUT> area:age     9.168e-07  4.898e-07   1.872  0.0622 .
#OUT> ---
#OUT> Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#OUT>
#OUT> Residual standard error: 0.2925 on 317 degrees of freedom
#OUT> Multiple R-squared:  0.5586, Adjusted R-squared:  0.5545
#OUT> F-statistic: 133.7 on 3 and 317 DF,  p-value: < 2.2e-16
```

In this instance, we see that indeed there is a small positive interaction between **area** and **age** on the sales price: even though **age** in isolation decreases the sales value, bigger houses command a small premium if they are older.

### 5.4.1 Interactions with Dummies: Differential Slopes

It is straightforward to extend the interactions logic to allow not only for different intercepts, but also different slopes for each subgroup in a dataset. Let's go back to our dataset of wages from section ?? above. Now that we know how to create an interaction between two variables, we can easily modify equation (??) like this:

$$\ln w = b_0 + b_1 \text{female} + b_2 \text{educ} + b_3 (\text{female} \times \text{educ}) + e \quad (5.7)$$

The only peculiarity here is that `female` is a factor with levels 0 and 1: i.e. the interaction term  $b_3$  will be zero for all men. Similarly to above, we can test whether there are indeed different returns to education for men and women by looking at the estimated value  $b_3$ :

```
lm_w_interact <- lm(lwage ~ educ * female , data = wage1) # R expands to full interactions model
summary(lm_w_interact)
```

```
#OUT>
#OUT> Call:
#OUT> lm(formula = lwage ~ educ * female, data = wage1)
#OUT>
#OUT> Residuals:
#OUT>      Min       1Q   Median       3Q      Max
#OUT> -2.02673 -0.27468 -0.03721  0.26221  1.34740
#OUT>
#OUT> Coefficients:
#OUT>              Estimate Std. Error t value Pr(>|t|)
#OUT> (Intercept)   8.260e-01  1.181e-01   6.997 8.08e-12 ***
#OUT> educ          7.723e-02  8.988e-03   8.593  < 2e-16 ***
#OUT> female1     -3.601e-01  1.854e-01  -1.942  0.0527 .
#OUT> educ:female1 -6.408e-05  1.450e-02  -0.004  0.9965
#OUT> ---
#OUT> Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#OUT>
#OUT> Residual standard error: 0.4459 on 522 degrees of freedom
#OUT> Multiple R-squared:  0.3002, Adjusted R-squared:  0.2962
#OUT> F-statistic: 74.65 on 3 and 522 DF,  p-value: < 2.2e-16
```

We will in the next chapter learn that the estimate for  $b_3$  on the interaction `educ:female1` is difficult for us to distinguish from zero in a statistical sense; Hence for now we conclude that there are no significantly different returns in education for men and women in this data. This is easy to verify visually in this plot, where we are unable to detect a difference in slopes in the right panel.

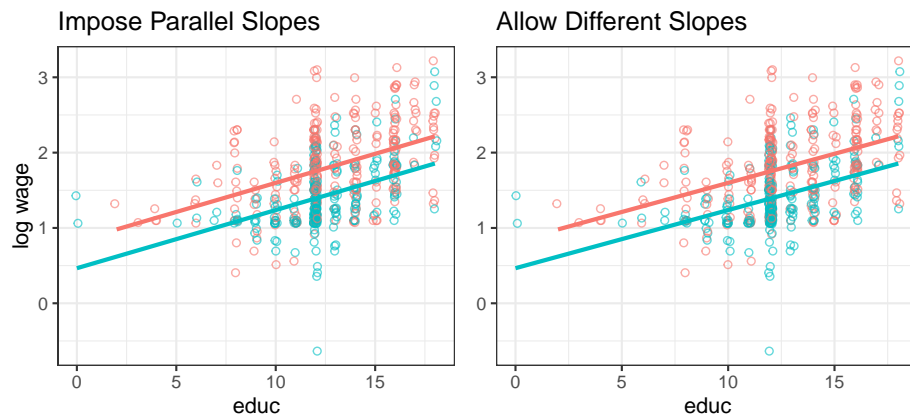


Figure 5.3: log wage vs educ. Right panel allows slopes to be different - turns out they are not!

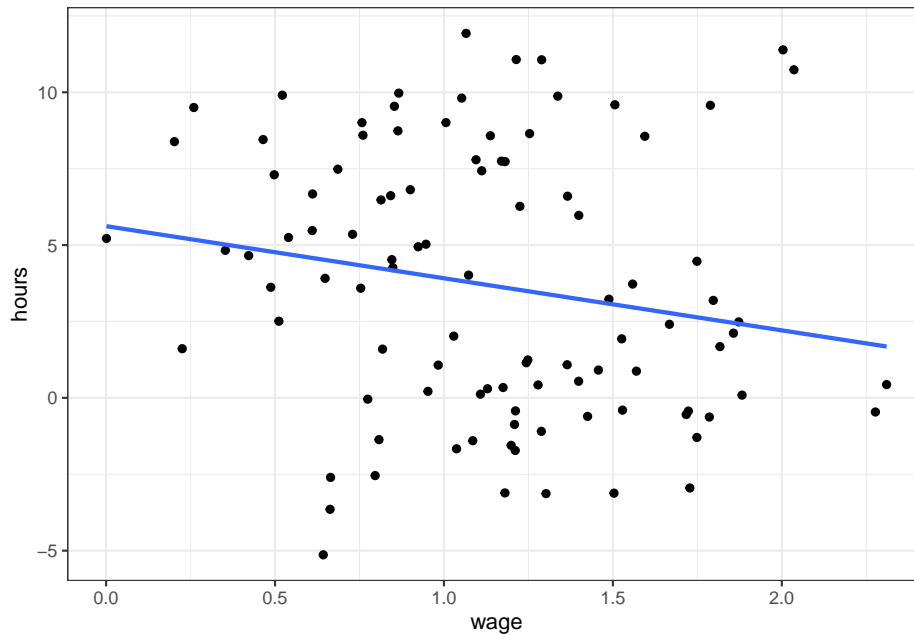
## 5.5 (Unobserved) Individual Heterogeneity

Finally, dummy variables are sometimes very important to account for spurious relationships in that data. Consider the following (artificial example):

1. Suppose we collected data on hourly wage data together with a the number of hours worked for a set of individuals.
2. We plot want to investigate labour supply behaviour of those individuals, hence we run regression `hours_worked ~ wage`.
3. We expect to get a positive coefficient on `wage`: the higher the wage, the more hours worked.
4. You know that individuals are members of either group `g=0` or `g=1`.

```
#OUT> `geom_smooth()` using formula 'y ~ x'
#OUT> `geom_smooth()` using formula 'y ~ x'
#OUT> `geom_smooth()` using formula 'y ~ x'
```





Here we observe a slightly negative relationship: higher wages are associated with fewer hours worked? Maybe. But what is this, there is a group identifier in this data! Let's use this and include  $g$  as a dummy in the regression - suppose  $g$  encodes male and female.

This is an artificial example; yet it shows that you can be severely misled if you don't account for group-specific effects in your data. The problem is particularly acute if we don't know group membership - we can then resort to advanced methods that are beyond the scope of this course to estimate which group each individual belongs to. If we do know group membership, however, it is good practice to include a group dummy so as to control for group effects.

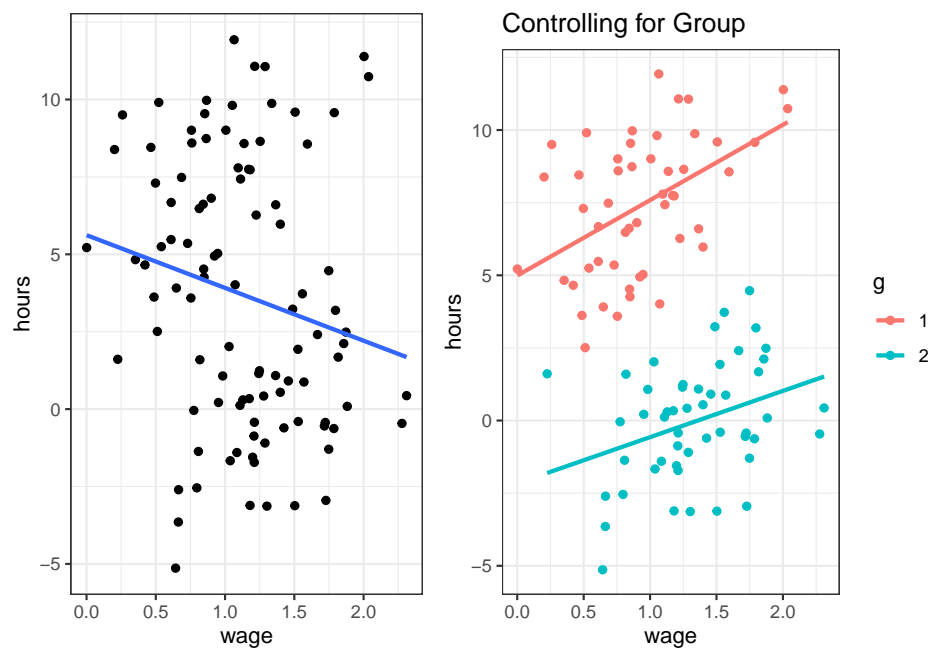


Figure 5.4: Left and right panel exhibit the same data. The right panel controls for group composition.

## Chapter 6

# Regression Inference

In this chapter we want to investigate uncertainty in regression estimates. We want to understand what the precise meaning of the **Std. Error** column in a typical regression table is telling us. In terms of a picture, we want to understand better the meaning of the shaded area as in this one here:

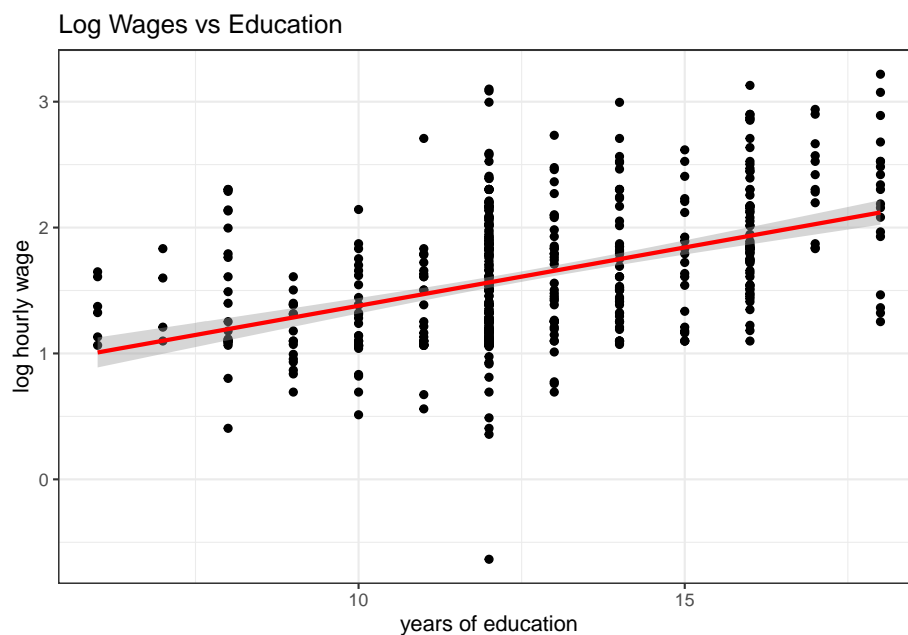


Figure 6.1: Confidence bands around a regression line.

In order to fully understand this, we need to go back and make sure we have a good grasp of sampling. Let's do this first.

## 6.1 Sampling

In class we were confronted with a jar of Tricolore Fusilli pasta as picture in figure ??.<sup>1</sup> We asked ourselves a question which, secretly, many of you had asked themselves at one point in their lives, namely:

What is the proportion of **green** Fusilli in a pack of Tricolore Fusilli?

Well, it's time to find out.



Figure 6.2: A glass jar filled with Fusilli pasta in three different colors.

Let's call the fusilli in this jar our study population, i.e. the set of units about which we want to learn something. There are several approaches to address the question of how big a proportion in the population the green Fusilli make up. One obvious solution is to enumerate all Fusilli according to their color, and compute their proportion in the entire population. It works perfectly well as a solution, but is a long and arduous process, see figures ?? and ??.

Additionally, you may draw worried looks from the people around you, while you are doing it. Maybe this is not the right way to approach this task?<sup>2</sup>

<sup>1</sup>This part is largely based on moderndive, to which I am giving full credit hereby. Thanks for this great idea.

<sup>2</sup>Regardless of the worried onlookers, I did what I had to do and I carried on to count the green pile. I know exactly how many greens are in there now! I then computed the weight of



Figure 6.3: Manually separating Fusilli by their color is very costly in terms of effort and cost.



Figure 6.4: Heaps of Fusilli pasta ready to be counted.

### 6.1.1 Taking One Sample From the Population

We started by randomly grabbing a handful of Fusilli from the jar and by letting drop exactly  $N = 20$  into a paper coffee cup, pictured in ???. We call  $N$  the sample size. The count and corresponding proportions of each color in this first sample are shown in the following table:

Color	Count	Proportion
Red	7	0.35
Green	5	0.25
White	8	0.4

So far, so good. We have our first estimate of the population proportion of green Fusilli in the overall population: 0.25. Notice that taking a sample of  $N = 20$  was much quicker and much less painful than performing the full count (i.e. the census) of Fusilli performed above.



Figure 6.5: Taking one sample of 20 Fusilli from the jar.

Then, we put my sample back into the jar, and we reshuffled the Fusilli. Had we taken another sample, again of  $N = 20$ , would we again have gotten 7

---

20 Fusilli (5g), and backed out the number of Fusilli in the other piles. I will declare those numbers as the true numbers. (Sceptics are free to recount.)

Red, 5 Green, and 8 White, just as in the first sample? Maybe, but maybe not. Suppose we had carried on for several times drawing samples of 20 and counting the colors: Would we also have observed 5 green Fusilli? Definitely not. We would have noted some degree of variability in the proportions computed from our samples. The sample proportions in this case are an example of a sample statistic.

**Sampling Variation** refers to the fact that if we randomly take samples from a wider population, the random composition of each sample will imply that we obtain statistics that vary - they take on potentially different values in each sample.

Let's see how this story evolved as we started taking more samples at a time.

## 6.2 Taking Eleven Samples From The Population

We formed teams of two students in class who would each in turn take samples from the jar (the population) of size  $N = 20$ , as before. Each team computed the proportion of green Fusilli they had in their sample, and we wrote this data down in a table on the board. Then, we drew a histogram which showed how many samples had fallen into which bins.

We looked at the histogram in figure ?? and we noted several things:

1. The largest proportions were 0.3 green
2. The smallest proportion was 0.15 green.
3. Most samples found a proportion of 0.25 green fusilli.
4. We did think that this looked suspiciously like a **normal distribution**.

We collected the sample data into a data.frame:

```
pasta_samples <- data.frame(group = 1:11, replicate = 1:11, prop_green = c(0.3,0.25,0.25,0.3,0.15,0.3,0.25,0.25,0.3,0.2,0.3))
pasta_samples
```

```
#OUT>      group replicate prop_green
#OUT> 1         1         1      0.30
#OUT> 2         2         2      0.25
#OUT> 3         3         3      0.25
#OUT> 4         4         4      0.30
#OUT> 5         5         5      0.15
#OUT> 6         6         6      0.30
#OUT> 7         7         7      0.25
#OUT> 8         8         8      0.25
#OUT> 9         9         9      0.20
```



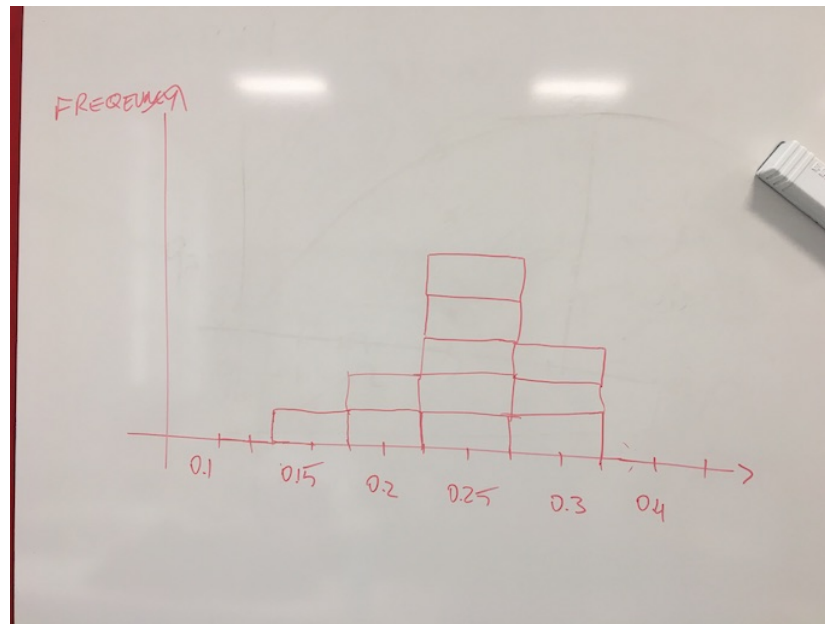
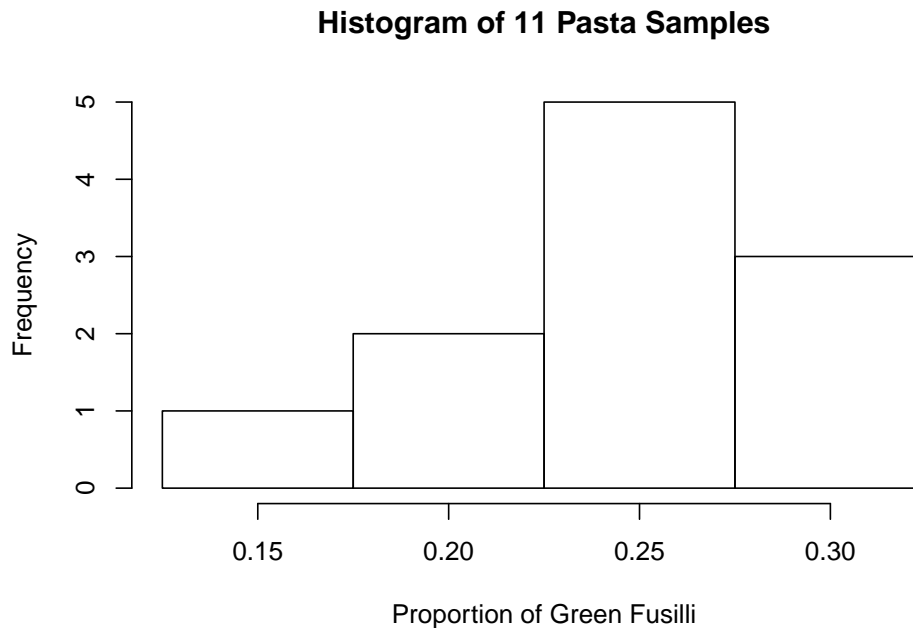


Figure 6.6: Taking eleven samples of 20 Fusilli each from the jar, and plotting the histogram of obtained sample proportions of Green Fusilli.

```
#OUT> 10    10      10      0.25
#OUT> 11    11      11      0.20
```

This produces an associated histogram which looks very much like the one we draws onto the board:





### 6.2.1 Recap

Let's recapitulate what we just did. We wanted to know what proportion of Fusilli in the glass jar in figure ?? are green. We acknowledged that an exclusive count, or a census, is a costly and cumbersome exercise, which in most circumstances we will try to avoid. In order to make some progress nonetheless, we took a random sample from the full population in the jar: we randomly selected 20 Fusilli, and looked at the proportion of green ones in there. We found a proportion of 0.25.

After replacing the Fusilli from the first sample in the jar, we asked ourselves if, upon drawing a new sample of 20 Fusilli, we should expect to see the same outcome - and we concluded: maybe, but maybe not. In short, we discovered some random variation from sample to sample. We called this **sampling variation**.

The purpose of this little activity was three-fold:

1. To understand that random samples differ and that there is sampling variation.
2. To understand that bigger samples will yield smaller sampling variation.
3. To illustrate that the sampling distribution of any statistic (i.e. not only the sample proportion as in our case) computed from a random sample converges to a normal distribution as the sample size increases.

The value of this exercise consisted in making **you** perform the sampling activity yourself. We will now hand over to the brilliant **moderndive** package, which will further develop this chapter.

### 6.3 Handover to Moderndive

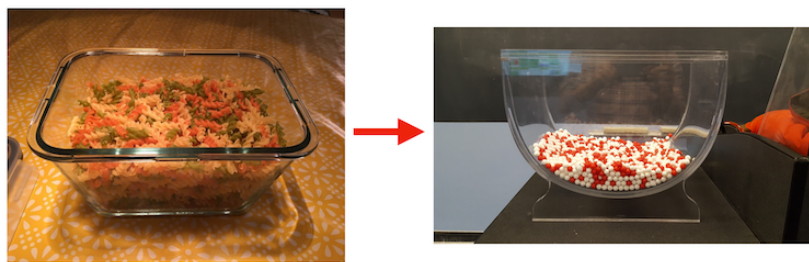


Figure 6.7: The Moderndive package used red and white balls instead of fusilli pasta.

The sampling activity in **moderndive** was performed with red and white balls instead of green fusilli pasta. The rest is identical. We will now read sections 7.2 and 7.3 in their book, as well as chapter 8 on confidence intervals and bootstrapping, and chapter 9 on hypothesis testing.

### 6.4 Uncertainty in Regression Estimates

In the previous chapters we have seen how the OLS method can produce estimates about intercept and slope coefficients from data. You have seen this method at work in R by using the `lm` function as well. It is now time to introduce the notion that given that  $b_0$ ,  $b_1$  and  $b_2$  are estimates of some unknown population parameters, there is some degree of **uncertainty** about their values. Another way to say this is that we want some indication about the precision of those estimates. The underlying issue is that the data we have at hand are usually samples from a larger population.

How confident should we be about the estimated values  $b$ ?

## 6.5 What is true? What are Statistical Models?

A **statistical model** is simply a set of assumptions about how some data have been generated. As such, it models the data-generating process (DGP), as we have it in mind. Once we define a DGP, we could simulate data from it and see how this compares to the data we observe in the real world. Or, we could change the parameters of the DGP so as to understand how the real world data would change, could we (or some policy) change the corresponding parameters in reality. Let us now consider one particular statistical model, which in fact we have seen so many times already.

## 6.6 The Classical Regression Model (CRM)

Let's bring back our simple model (??) to explain this concept.

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \quad (6.1)$$

The smallest set of assumptions used to define the classical regression model as in (??) are the following:

1. The data are **not linearly dependent**: Each variable provides new information for the outcome, and it cannot be replicated as a linear combination of other variables. We have seen this in section ???. In the particular case of one regressor, as here, we require that  $x$  exhibit some variation in the data, i.e.  $Var(x) \neq 0$ .
2. The mean of the residuals conditional on  $x$  should be zero,  $E[\varepsilon|x] = 0$ . Notice that this also means that  $Cov(\varepsilon, x) = 0$ , i.e. that the errors and our explanatory variable(s) should be uncorrelated. It is said that  $x$  should be **strictly exogenous** to the model.

These assumptions are necessary to successfully (and correctly!) run an OLS regression. They are often supplemented with an additional set of assumptions, which help with certain aspects of the exposition, but are not strictly necessary:

3. The data are drawn from a **random sample** of size  $n$ : observation  $(x_i, y_i)$  comes from the exact same distribution, and is independent of observation  $(x_j, y_j)$ , for all  $i \neq j$ .
4. The variance of the error term  $\varepsilon$  is the same for each value of  $x$ :  $Var(\varepsilon|x) = \sigma^2$ . This property is called **homoskedasticity**.
5. The error is normally distributed, i.e.  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$

Invoking assumption 5. in particular defines what is commonly called the normal linear regression model.

### 6.6.1 $b$ is not $\beta$ !

Let's talk about the small but important modifications we applied to model (??) to end up at (??) above:

- $\beta_0$  and  $\beta_1$  are intercept and slope parameters
- $\varepsilon$  is the error term.

First, we assumed that (??) is the correct representation of the DGP. With that assumption in place, the values  $\beta_0$  and  $\beta_1$  are the true parameter values which generated the data. Notice that  $\beta_0$  and  $\beta_1$  are potentially different from  $b_0$  and  $b_1$  in (??) for a given sample of data - they could in practice be very close to each other, but  $b_0$  and  $b_1$  are estimates of  $\beta_0$  and  $\beta_1$ . And, crucially, those estimates are generated from a sample of data. Now, the fact that our data  $\{y_i, x_i\}_{i=1}^N$  are a sample from a larger population, means that there will be sampling variation in our estimates - exactly like in the case of the sample mean estimating the population average as mentioned above. One particular sample of data will generate one particular set of estimates  $b_0$  and  $b_1$ , whereas another sample of data will generate estimates which will in general be different - by how much those estimates differ across samples is the question in this chapter. In general, the more observations we have the greater the precision of our estimates, hence, the closer the estimates from different samples will lie together.

### 6.6.2 Violating the Assumptions of the CRM

It's interesting to consider in which circumstances we might violate those assumptions. Let's give an example for each of them:

1. No Perfect Collinearity. We have seen that a perfect collinearity makes it impossible to compute to OLS coefficients. Remember the example about adding `wtplus = wt + 1` to the `mtcars` dataset? Here it is:

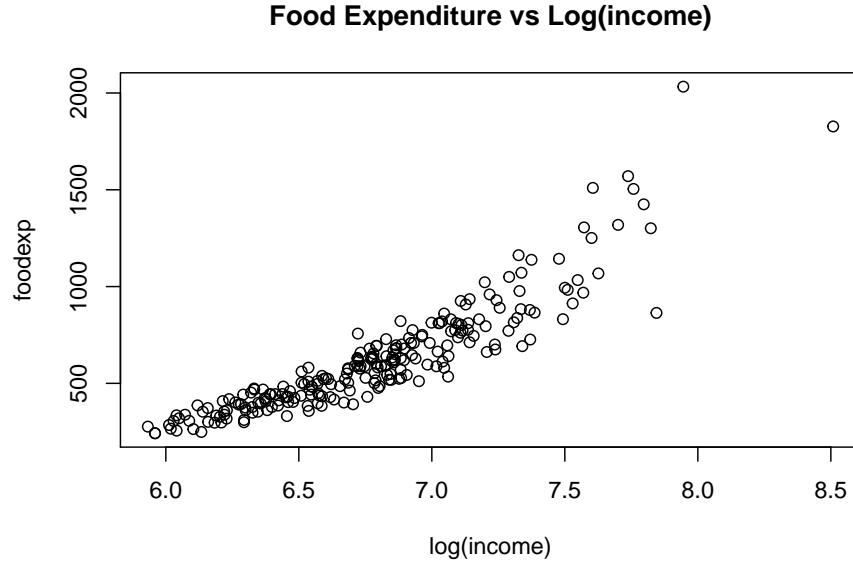
```
library(dplyr)
mtcars %>%
  mutate(wtplus = wt + 1) %>%
  lm(mpg ~ wt + wtplus, data = .)
```

```
#OUT>
#OUT> Call:
#OUT> lm(formula = mpg ~ wt + wtplus, data = .)
#OUT>
#OUT> Coefficients:
#OUT> (Intercept)          wt          wtplus
#OUT>      37.285      -5.344           NA
```

That the coefficient on `wtplus` is NA is the result of the direct linear dependence. (Notice that creating `wtplus2 = (wt + 1)^2` would work, since

that is not linear!)

2. Conditional Mean of errors is zero,  $E[\varepsilon|x] = 0$ . Going back to our running example in figure ?? about wages and education: Suppose that each individual  $i$  in our data has something like innate ability, something we might wish to measure with an IQ-test, however imperfectly. Let's call it  $a_i$ . It seems reasonable to think that high  $a_i$  will go together with high wages. At the same time, people with high  $a_i$  will find studying for exams and school work much less burdensome than others, hence they might select into obtaining more years of schooling. The problem? Well, there is no  $a_i$  in our regression equation - most of time we don't have a good measure of it to start with. So it's an unobserved variable, and as such, it is part of the error term  $\varepsilon$  in our model. We will attribute to **educ** part of the effect on wages that is actually caused by ability  $a_i$ ! Sometimes we may be able to reason about whether our estimate on **educ** is too high or too low, but we will never know its true value. We don't get the ceteris paribus effect (the true partial derivative of **educ** on **lwage**). Technically, the assumption  $E[\varepsilon|x] = 0$  implies that  $Cov(\varepsilon, x) = 0$ , so that's the part that is violated.
  
3. Data from Random Sample. One common concern here is that the observations in the data could have been selected in a particular fashion, which would make it less representative of the underlying population. Suppose we had ended up with individuals only from the richest neighborhood of town; Our interpretation the impact of education on wages might not be valid for other areas.
  
4. Homoskedasticity. For correct inference (below!), we want to know whether the variance of  $\varepsilon$  varies with our explanatory variable  $x$ , or not. Here is a typical example where it does:



As income increases, not all people increase their food consumption in an equal way. So  $Var(\varepsilon|x)$  will vary with the value of  $x$ , hence it won't be equal to the constant  $\sigma^2$ .

5. If the distribution of  $\varepsilon$  is not normal, it is more cumbersome to derive theoretical results about inference.

## 6.7 Standard Errors in Theory

The standard deviation of the OLS parameters is generally called standard error. As such, it is just the square root of the parameter's variance. Under assumptions 1. through 4. above we can define the formula for the variance of our slope coefficient in the context of our single regressor model (??) as follows:

$$Var(b_1|x_i) = \frac{\sigma^2}{\sum_i (x_i - \bar{x})^2} \quad (6.2)$$

In practice, we don't know the theoretical variance of  $\varepsilon$ , i.e.  $\sigma^2$ , but we form an estimate about it from our sample of data. A widely used estimate uses the already encountered SSR (sum of squared residuals), and is denoted  $s^2$ :

$$s^2 = \frac{SSR}{n-p} = \frac{\sum_{i=1}^n (y_i - b_0 - b_1 x_i)^2}{n-p} = \frac{\sum_{i=1}^n e_i^2}{n-p}$$

where  $n-p$  are the degrees of freedom available in this estimation.  $p$  is the number of parameters we wish to estimate (here: 1). So, the variance formula

would become

$$Var(b_1|x_i) = \frac{SSR}{(n-p) \sum_i^N (x_i - \bar{x})^2} \quad (6.3)$$

We most of the time work directly with the standard error of a coefficient, hence we define

$$SE(b_1) = \sqrt{Var(b_1|x_i)} = \sqrt{\frac{SSR}{(n-p) \sum_i^N (x_i - \bar{x})^2}} \quad (6.4)$$

You can clearly see that, as  $n$  increases, the denominator increases, and therefore variance and standard error of the estimate will decrease.





## Chapter 7

# Causality

In this chapter we take on a challenging part of our course. Remember that in the first set of slides we introduced Econometrics as the economist's toolkit to answer questions like does  $x$  **cause**  $y$ ? Let's illustrate the issues at stake with a question from epidemiologie and public health:

Does smoking **cause** lung cancer?

Just in case you were wondering: Yes it does! However, for a very long time the causal impact of smoking on lung cancer was hotly debated, and it's instructive for us to look at this history.<sup>1</sup>

Let's go back to the 1950's. We are at the start of a big increase in deaths from lung cancer. At the same time cigarette consumption was growing very fast. With the benefit of hindsight, we can now draw this graph:

However, time series graphs are poor tools to make causal statements. Many other things had changed from 1900 to 1950, all of which could equally be responsible for the rise in cancer rates:

1. Tarring of roads
2. Inhalation of motor exhausts (leaded gasoline fumes)
3. General greater air pollution.

We call those other factors **confounders** of the relationship between smoking and lung cancer.

So, there were a series of sceptics around who at the time were contesting the existing evidence. That evidence consisted in general of the following:

---

<sup>1</sup>This chapter is drawn from chapter 5 of The Book of Why by Judea Pearl.

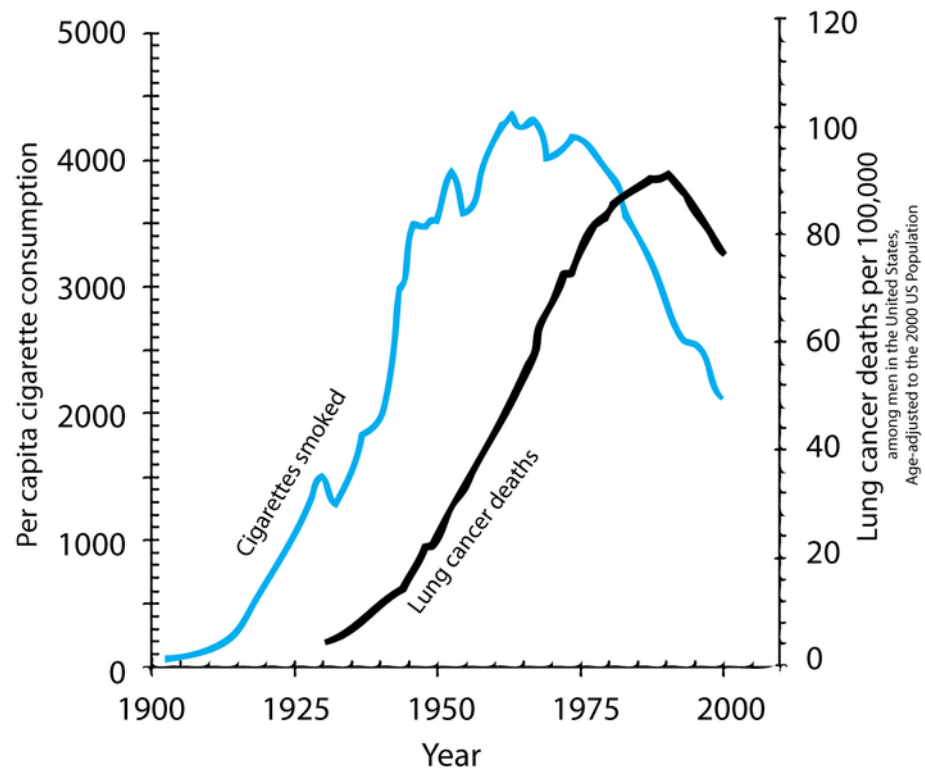


Figure 7.1: Two time series showing cigarette consumption per capita and incidence of lung cancer in the USA.

1. **Case-Control studies:** British Epidemiologists Richard Doll and Austin Bradford Hill started to compare people already diagnosed with cancer to those without, recording their history, and observable characteristics (like age and health behaviours). In one study, out of 649 lung cancer patients interviewed, all but 2 had been smokers! In that study, a cancer patient was 1.5 million times more likely to be have been a smoker than a non-smoker. Still, critics said, there are several sources of bias:
  - Hospital patients could be a selected sample of the general (smoking) population.
  - Patients could suffer from recall bias, affecting their recollection of facts.
  - So, while comparing cancer patients to non-patients and controlling for several important confounders (like age, income and other observable characteristics), there was still scope for bias.
  - Moreover, replicating those studies, as Doll and Hill attempted, would not have solved this issue.
2. Next they attempted what doctors call a **Dose-Response Effect** study. In 1951 they sent out 60,000 questionnaires to British physicians asking about their smoking habits. Then they followed them over time:
  - Only 5 years on, heavy smokers had a death rate from lung cancer that was 24 times higher than for nonsmokers.
  - People who had smoked and then stopped reduced their risk by a factor of 2.
  - Still, notorious sceptics like R.A. Fisher were unconvinced. The studies still failed to compare **otherwise identical** smokers to non-smokers. There were still important unobserved confounders out there which could invalidate the conclusion that we observed indeed a **causal** relationship.

Let's put a some structure on this problem now, so we can make progress.

## 7.1 Directed Acyclical Graphs (DAG)

A DAG is a tool to visualize a causal relationship. It is a graph where nodes are connected via arrows, where an arrow can run in one direction only (hence, directed graph). If an arrow starts at node  $x$  and ends at node  $y$ , we say that  $x$  causes  $y$ . Here is a simple example of such a DAG:

Now consider this setting, where there is a third variable,  $z$ . It could be possible that also  $z$  has a direct influence on  $y$ :

Now let's change this and create a path from  $z$  to both  $x$  and  $y$  instead. We call  $z$  a confounder in the relationship between  $x$  and  $y$ :  $z$  confounds the direct causal impact of  $x$  on  $y$ , by affecting them both at the same time. What is more, there is no arrow from  $x$  to  $y$  at all, so the only real explanatory variable

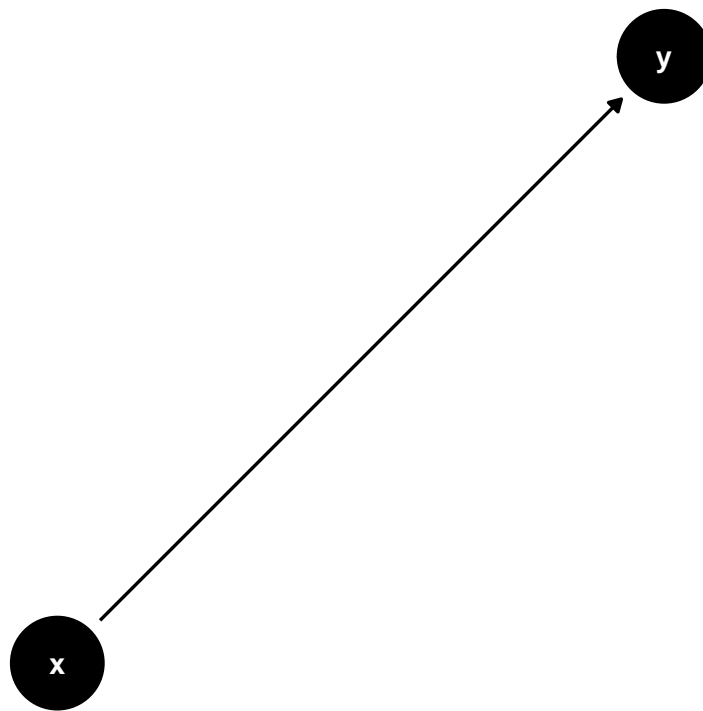


Figure 7.2: A simple DAG showing the causal impact of  $x$  on  $y$ .

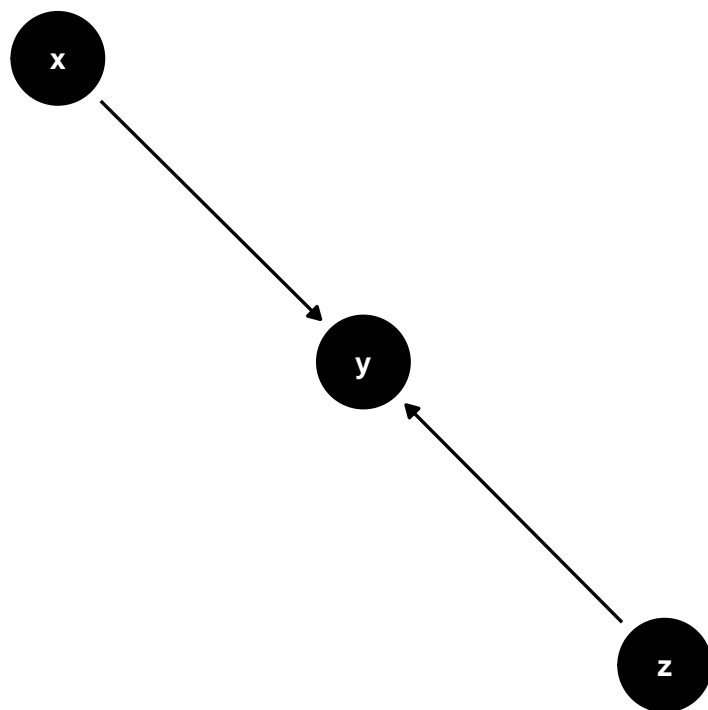


Figure 7.3: A simple DAG with with 2 causal paths: Both  $x$  and  $z$  have a direct impact on  $y$ .

here is in fact  $z$ . Attributing any explanatory power to  $x$  would be wrong in this setting.

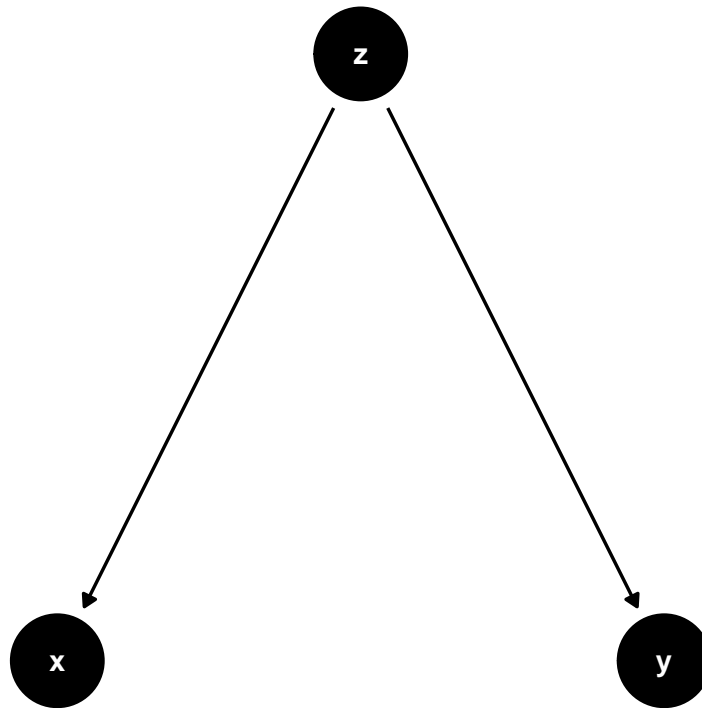


Figure 7.4: A simple DAG where  $z$  is a confounder. There is no causal path from  $x$  to  $y$ , and any correlation we observe between those variables is completely induced by  $z$ . We call this spurious correlation.

Here is a second example where  $z$  is a confounder, but slightly different.

In ?? there is an arrow from  $x$  to  $y$ . In this setting, if we are able to observe  $z$ , we can adjust the correlation we observe between  $x$  to  $y$  for the variation induced by  $z$ . In practice, this is precisely what multiple regression will do: holding  $z$  fixed at some value, what is the partial effect of  $x$  on  $y$ . Notice that  $z$  cedes to be a confounder in this situation, and interpreting our regression coefficient on  $x$  as causal is correct.

## 7.2 Smoking in a DAG

Let's use this and cast our problem as a DAG now. What the scientists in the 1950s faced where two competing models of the relationship between smoking and lung cancer:

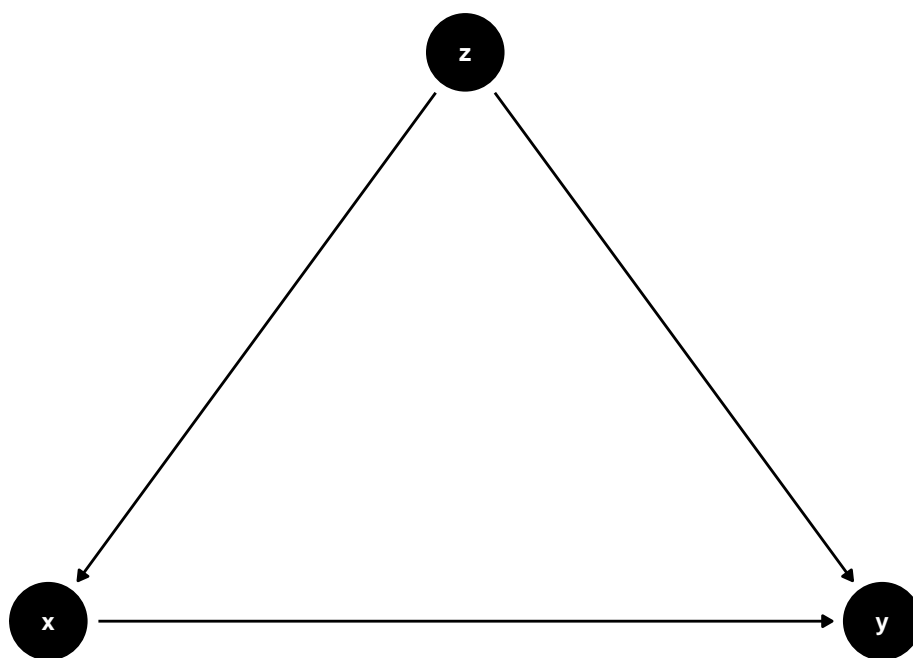


Figure 7.5:  $z$  is still a confounder here, but there is a causal link from  $x$  to  $y$  now. If we observed  $z$ , we can control for it.

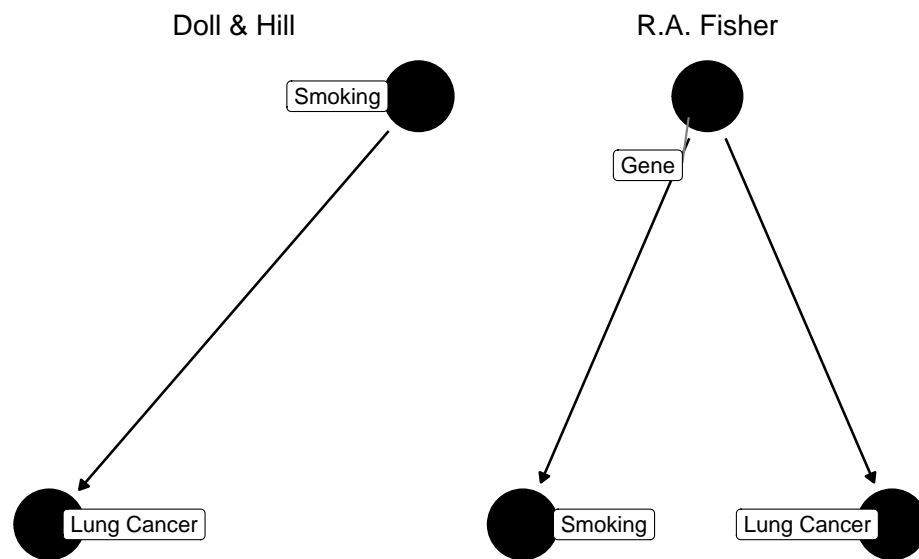


Figure 7.6: Two competing causal graphs for the relationship between smoking and lung cancer. In the right panel Lung Cancer is directly impacted by a genetic factor, which at the same time also influences smoking. This is a stark representation of Fisher's view. Another version would have an additional arrow from Smoking to Lung Cancer in the right panel.



Basically, what critics like Fisher were claiming was that the existing studies did not compare like for like. In other words, our *ceteris paribus* assumption was not satisfied. They were worried that smoking was not the only relevant difference between a population of smokers and one of non-smokers. In particular, they worried that people **self-selected** into smoking, and that the choice to become a smoker may be influenced by other, unobserved, underlying forces - like genetic predisposition, for example. That could mean that smokers were also more likely to take risks, or more likely to be heavy drinkers, or engage in other behaviours that might be conducive to develop lung cancer. They did not formulate it in terms of genetics at the time, because they could not know until the 2000's, when the human genome was sufficiently mapped to establish this fact (and indeed there **is** a smoking gene! But that's beside the point), but they worried about this factor.

The argument was settled in the eyes of most physicians, when Jerome Cornfield in 1959 wrote a rebuttal of Fisher's points. Cornfield's strategy was to allow Fisher to have his unobserved factor, but to show that there was an upper bound to how important it could be in determining the outcome. Here goes:

1. Suppose there is indeed a confounding factor "smoking gene", and that it completely determines the risk of cancer in smokers.
2. Suppose smokers are observed to have 9 times the risk of non-smokers to develop lung cancer.
3. The smoking gene needs to be at least 9 times more prevalent in smokers than in non-smokers to explain this difference in risk.

But now consider what this implies. Let's suppose that around 11% of all non-smokers have the smoking gene. That means that  $9 \times 11 = 99\%$  of smokers need to have it! What's even more worrying, if only even 12% of non smokers have the gene, then the argument breaks down because it would require  $9 \times 12 = 108\%$  of smokers to have it, which is of course impossible.

This argument was so important that it got a name: **Cornfield's inequality**. It left of Fisher's argument nothing but a pile of rubble. It's impossible to think that genetic variation alone could be so important in determining a complex choice of becoming a smoker or not. Looking back at the right panel of figure ??, the link from smoking to lung cancer was much too strong to be explained by the genetic hypothesis alone.

## 7.3 Randomized Control Trials (RCT) Primer

We now present a quick introduction to Randomized Control Trials (RCTs). The history of randomization is fascinating and goes back a long time, again involving R.A. Fisher from above.<sup>2</sup> Suffice it to say that RCTs have become

<sup>2</sup>I refer the interested student to the introduction of the potential outcomes model of Scott Cunningham's mixtape, which heavily influences this section.

so important in Economics that the Nobel Price in Economics 2019 has been awarded to three exponents of the RCT literature, Duflo, Banerje and Kremer. RCTs are widely used in Medicine, where they originate from (in some sense). But, what are RCTs?

A randomized controlled trial is a type of scientific experiment that aims to reduce certain sources of bias when testing the effectiveness of some intervention (treatment or policy); this is accomplished by randomly allocating subjects to two or more groups, treating them differently, and then comparing them with respect to a measured response.

That sounds really intuitive. If we randomly allocate people to receive treatment, there can be no concern of unobserved confounders, as we have relieved the subjects of making the choice to get treated. Remember the cigarette smokers above: The concern was that an unobserved genetic predisposition correlated with both choosing to become a smoker but also with other potentially cancer-inducing behaviours like drinking or risk taking. Imagine for a moment that we could randomly select people at some young age to be selected for treatment (smoking for 30 years, say). The genetic predisposition will be equally prevalent in both treatment and control group. However, only the treatment group is allowed (and indeed forced) to smoke. Observing higher cancer rates in the treatment group would provide causal evidence for the effect of smoking on lung cancer.

Thankfully, such an experiment is impossible to run on ethical grounds. We could never subject individuals to such severe and prolonged health risks for the sake of a research study. That's why the question took so long to be settled!

Let's introduce a formal framework now to think more about RCTs.

## 7.4 The Potential Outcomes Model

The Potential Outcomes Model, often named after one of its inventors the Rubin Causal Model, posits that there are two states of the world - the potential outcomes. A first state, where a certain intervention is administered to an individual, and a second state, where this is not the case. Formally, this idea is expressed with superscripts 0 and 1, like this:

- $Y_i^1$ : individual  $i$  has been treated
- $Y_i^0$ : individual  $i$  has **not** been treated

Denoting with  $D_i \in \{0, 1\}$  the treatment indicator which is one if  $i$  is indeed treated, the observed outcome  $Y_i$  is then

$$Y_i = D_i Y_i^1 + (1 - D_i) Y_i^0 \quad (7.1)$$

This simple equation is able to formalize a rather deep question. We only ever observe one outcome of events for a given individual  $i$ , say  $Y_i = Y_i^1$  in case treatment was given. The deep question is: what would have happened to  $i$ , had they **not** received treatment? You will realize that this a very natural question for us humans to put to ourselves, and to subsequently answer:

- How long would the trip have taken, had I chosen another metro line?
- What would have happened, had I chosen to study a different subject?
- What would have happened, had Neo taken the blue pill instead?

Our ability to make those considerations distinguishes us from animals. It's one of the biggest challenges for machines when trying to be intelligent.

What makes this question so hard to answer for machines and animals alike is the fact that one has to imagine a parallel universe where the actions taken were different, **without** having observed that precise situation before. Neo did not take the blue pill, and whatever happened after that originated from this decision - so how are we to tell what would have happened? It's easy for us and still hard for machines.

Potential outcome  $Y_i^0$  above is what is known as the counterfactual outcome. What would have happened to subject  $i$ , had they **not** received treatment  $D$ ?

Following Rubin, let us define the **treatment effect** for individual  $i$  as follows:

$$\delta_i = Y_i^1 - Y_i^0 \quad (7.2)$$

Notice our insistence about talking about a single individual  $i$  throughout here. Keeping the potential outcome model (??) in mind, i.e. the fact that we only observe one of both outcomes, we face the **fundamental identification problem of program evaluation**:

Given we only observe one potential outcome, we cannot compute the treatment effect  $\delta_i$  for any individual  $i$ .

That's pretty dire news. Let's see if we can do better with an average effect instead. Let's define three average effects of interest:

1. the Average Treatment Effect (ATE):

$$\delta^{ATE} = E[\delta_i] = E[Y_i^1] - E[Y_i^0]$$

2. the Average Treatment on the Treated (ATT):

$$\delta^{ATT} = E[\delta_i | D_i = 1] = E[Y_i^1 | D_i = 1] - E[Y_i^0 | D_i = 1]$$

3. the Average Treatment on the Untreated (ATU):

$$\delta^{ATU} = E[\delta_i | D_i = 0] = E[Y_i^1 | D_i = 0] - E[Y_i^0 | D_i = 0]$$

Notice that none of those can be computed from data either, because all of them require data on individual  $i$  from both scenarios. Let's focus on the ATE for now. Fundamentally we face a **missing data problem**: either  $Y_i^1$  or  $Y_i^0$  are missing from our dataset. Nevertheless, let's setup the following naive simple difference in means estimator  $\hat{\delta}$ :

$$\hat{\delta} = E[Y_i^1 | D_i = 1] - E[Y_i^0 | D_i = 0] \quad (7.3)$$

$$= \frac{1}{N_T} \sum_{i \in T} T_i - \frac{1}{N_C} \sum_{j \in C} Y_j \quad (7.4)$$

in other words, we just difference the mean outcomes in both treatment (T) and control (C) groups. That is,  $N_C$  is the number of people in the control group,  $N_T$  is the same for treatment group.

Now let's consider what randomly choosing people for treatment does. The key consideration here is that the true  $\delta_i$  is potentially different for each person. That is, some people will have a high effect of treatment, while others may have a small (or even negative!) effect. To learn about the true  $\delta^{ATE}$  from our naive  $\hat{\delta}$ , it matters who ends up being treated!

Imagine that individuals have at least some partial knowledge about their likely gains from treatment, i.e. their personal  $\delta_i$ . If those who expect to benefit a lot will select disproportionately into treatment, then our estimator  $\hat{\delta}$  will be biased upwards for the true average effect  $\delta^{ATE}$ . This is so because the average of observed outcomes in the treatment group, i.e.

$$\frac{1}{N_T} \sum_{i \in T} Y_i$$

will be **too high**. It represents the disproportionately high treatment outcome  $Y_i^1$  for all those who anticipated such a high outcome from treatment, and who therefore were particularly eager to get selected into treatment. It's not representative of the true population wide treatment outcome  $E[Y_i^1]$ .

Here is where randomization comes into play. Suppose we now flip a coin for each person to determine whether they obtain treatment or not. This takes away from them the possibility to select on expected gains into treatment. Crucially, the distribution of effects  $\delta_i$  is still the same in the study population, i.e. there are still people with high and people with low effects. But we have solved the missing data problem mentioned above, because whether  $Y_i^1$  or rather  $Y_i^0$  is observed for each  $i$  is now **random**, and no longer a function of any other factor that  $i$  could act upon! Hooray!

Notice how this links back to our initial discussion about DAGs above. Randomisation essentially cancels the links starting at confounder  $z$  in ??.

## 7.5 Omitted Variable Bias and DAGs

We want to revisit the underlying assumptions of the classical model outlined in ?? in the previous chapter, which is closely related to the previous discussion. Let's talk a bit more about assumption number 2 of the definition in ?. It said this:

The mean of the residuals conditional on  $x$  should be zero,  $E[\varepsilon|x] = 0$ . This means that  $Cov(\varepsilon, x) = 0$ , i.e. that the errors and our explanatory variable(s) should be uncorrelated. We want  $x$  to be **strictly exogenous** to the model.

Let us start again with

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \quad (7.5)$$

and imagine it represents the data generating process (DGP) of the impact of  $x$  on  $y$ . Writing down this equation is tightly linked to drawing this DAG from above:

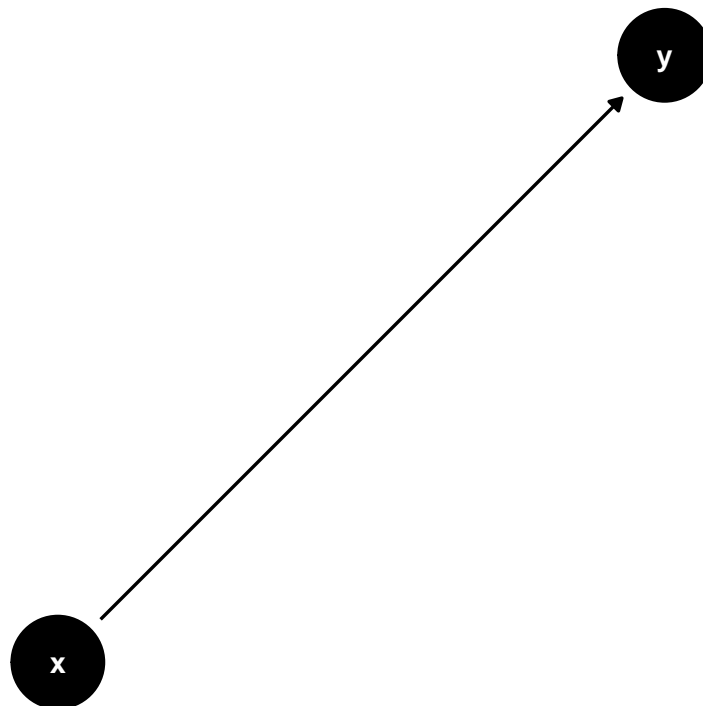


Figure 7.7: The same simple DAG showing the causal impact of  $x$  on  $y$ .

The role of  $\varepsilon_i$  in equation (??) is to allow for random variability in the data not captured by our model, almost as an acknowledgement that we would never be able to fully explain  $y_i$  with our necessarily simple model. However, assumption  $E[\varepsilon|x] = 0$  (or  $Cov(\varepsilon, x) = 0$ ) makes sure that those other factors are in **no systematic relationship** with our regressor  $x$ . Why? Well if it were the case that another factor  $z$  is related to  $x$ , we could never make our ceteris paribus statements of holding all other factors fixed, the impact of  $x$  on  $y$  is  $\beta$ . In other words, we'd have a confounder in our regression.

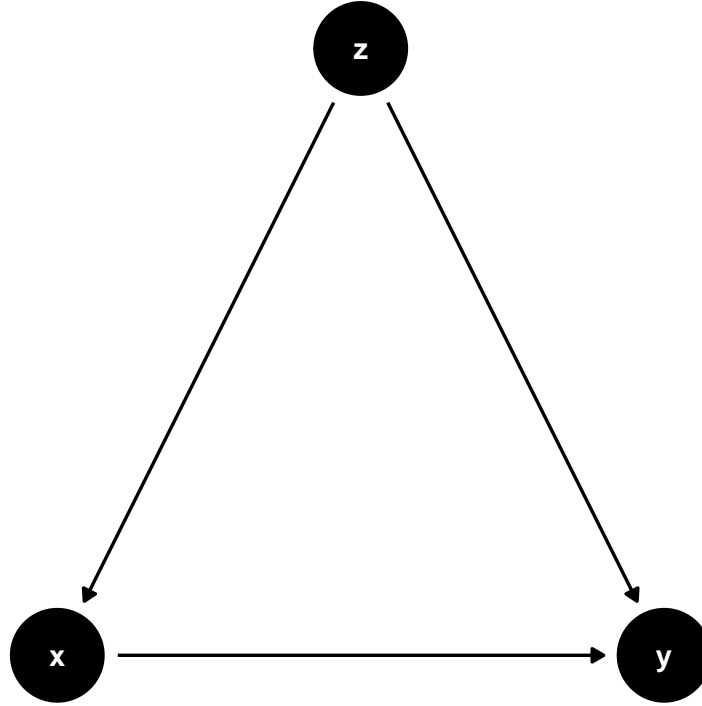


Figure 7.8: The same simple DAG where  $z$  is a confounder that needs to be controlled for.

Notice, again, that the key here is that if we don't control for  $z$ , it will form part of the error term  $\varepsilon$ . Given the causal link from  $z$  to  $x$ , we will then observe that  $Cov(x, u) = Cov(x, \varepsilon + z) \neq 0$ , invalidating our assumption.

### 7.5.1 House Prices and Bathrooms

Let's imagine that equation (??) represents the impact of number of bathrooms ( $x$ ) on the sales price of houses ( $y$ ). We run OLS as

$$y_i = b_0 + b_1 x_i + e_i$$

and find a positive impact of bathrooms on houses:

```
data(Housing, package="Ecdat")
hlm = lm(price ~ bathrms, data = Housing)
summary(hlm)

#OUT>
#OUT> Call:
#OUT> lm(formula = price ~ bathrms, data = Housing)
#OUT>
#OUT> Residuals:
#OUT>      Min       1Q   Median       3Q      Max
#OUT> -77225 -15271  -2510   11704  102729
#OUT>
#OUT> Coefficients:
#OUT>              Estimate Std. Error t value Pr(>|t|)
#OUT> (Intercept)    32794         2694   12.17  <2e-16 ***
#OUT> bathrms       27477         1952   14.08  <2e-16 ***
#OUT> ---
#OUT> Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#OUT>
#OUT> Residual standard error: 22880 on 544 degrees of freedom
#OUT> Multiple R-squared:  0.267, Adjusted R-squared:  0.2657
#OUT> F-statistic: 198.2 on 1 and 544 DF,  p-value: < 2.2e-16
```

In fact, from this you conclude that each additional bathroom increases the sales price of a house by 27477 dollars. Let's see if our assumption  $E[\varepsilon|x] = 0$  is satisfied:

```
library(dplyr)
# add residuals to the data
Housing$resid <- resid(hlm)
Housing %>%
  group_by(bathrms) %>%
  summarise(mean_of_resid=mean(resid))

#OUT> # A tibble: 4 x 2
#OUT>   bathrms mean_of_resid
#OUT>   <dbl>      <dbl>
#OUT> 1     1        -118.
#OUT> 2     2         955.
#OUT> 3     3       -11195.
#OUT> 4     4       32298.
```

Oh, that doesn't look good. Even though the unconditional mean  $E[e] = 0$  is very close to zero (type `mean(resid(hlm))`!), this doesn't seem to hold at all by categories of  $x$ . This indicates that there is something in the error term  $e$  which is correlated with `bathrms`. Going back to our discussion about ceteris paribus

in section ??, we stated that the interpretation of our OLS slope estimate is that

Keeping everything else fixed at the current value, what is the impact of  $x$  on  $y$ ? Everything also includes things in  $\varepsilon$  (and, hence,  $e$ )!

It looks like our DGP in (??) is the wrong model. Suppose instead, that in reality sales prices are generated like this:

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 z_i + \varepsilon_i \quad (7.6)$$

This would now mean that by running our regression, informed by the wrong DGP, what we estimate is in fact this:

$$y_i = b_0 + b_1 x_i + (b_2 z_i + e_i) = b_0 + b_1 x_i + u_i.$$

This is to say that by omitting variable  $z$ , we relegate it to a new error term, here called  $u_i = b_2 z_i + e_i$ . Our assumption above states that all regressors need to be uncorrelated with the error term - so, if  $Corr(x, z) \neq 0$ , we have a problem. Let's take this idea to our running example.

### 7.5.2 Including an Omitted Variable

What we are discussing here is called Omitted Variable Bias. There is a variable which we omitted from our regression, i.e. we forgot to include it. It is often difficult to find out what that variable could be, and you can go a long way by just reasoning about the data-generating process. In other words, do you think it's reasonable that price be determined by the number of bathrooms only? Or could there be another variable, omitted from our model, that is important to explain prices, and at the same time correlated with `bathrms`?

Let's try with `lotsize`, i.e. the size of the area on which the house stands. Intuitively, larger lots should command a higher price; At the same time, however, larger lots imply more space, hence, you can also have more bathrooms! Let's check this out:

```
#OUT>
#OUT> Call:
#OUT> lm(formula = price ~ bathrms + lotsize, data = Housing)
#OUT>
#OUT> Residuals:
#OUT>      Min       1Q   Median       3Q      Max
#OUT> -60752 -12532  -1674   10514   92931
#OUT>
#OUT> Coefficients:
```



```

#OUT>               Estimate Std. Error t value Pr(>|t|)
#OUT> (Intercept) 1.008e+04  2.810e+03   3.588 0.000364 ***
#OUT> bathrms      2.281e+04  1.703e+03  13.397 < 2e-16 ***
#OUT> lotsize      5.575e+00  3.944e-01  14.136 < 2e-16 ***
#OUT> ---
#OUT> Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#OUT>
#OUT> Residual standard error: 19580 on 543 degrees of freedom
#OUT> Multiple R-squared:  0.4642, Adjusted R-squared:  0.4622
#OUT> F-statistic: 235.2 on 2 and 543 DF,  p-value: < 2.2e-16

```

Here we see that the estimate for the effect of an additional bathroom decreased from 27477 to 22811.5 by almost 5000 dollars! Well that's the problem then. We said above that one more bathroom is worth 27477 dollars - if **nothing else changes**! But that doesn't seem to hold, because we have seen that as we increase `bathrms` from 1 to 2, the mean of the resulting residuals changes quite a bit. So there **is something in  $\epsilon$  which does change**, hence, our conclusion that one more bathroom is worth 27477 dollars is in fact invalid!

The way in which `bathrms` and `lotsize` are correlated is important here, so let's investigate that:

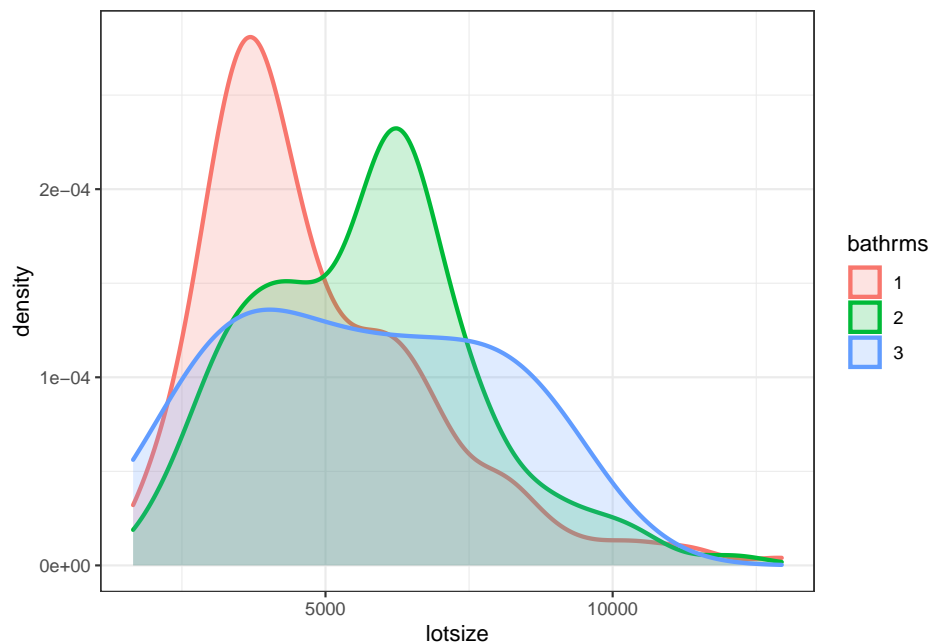


Figure 7.9: Distribution of 'lotsize' by 'bathrms'

This shows that `lotsize` and the number of bathrooms is indeed positively related. Larger lot of the house, more bathrooms. This leads to a general result:

**Direction of Omitted Variable Bias**

If the direction of correlation between omitted variable  $z$  and  $x$  is the same as that between  $x$  and  $y$ , we will observe upward bias in our estimate of  $b_1$ , and vice versa if the correlations go in opposite directions. In other words, we have positive bias if  $b_2 z_i > 0$  and vice versa.

## 7.6 STAR Experiment

tbd

## Chapter 8

# STAR Experiment

How to best allocate spending on schooling is an important question. What's the impact of spending money to finance smaller classrooms on student performance and outcomes, both in the short and in the long run? A vast literature in economics is concerned with this question, and for a long time there was no consensus.

The big underlying problem in answering this question is that we do not really know how student outcomes are produced. In other words, what makes a successful student? Is it the quality of their teacher? Surely matters. Is it quality of the school building? Could be. Is it that the other pupils are of high quality and this somehow rubs off to weaker pupils? Also possible. What about parental background? Sure. You see that there are many potential channels that could determine student outcomes. What is more, there could be several interdependencies amongst those factors. Here's a DAG!

We will look at an important paper in this literature now, which used a randomized experiment to make some substantial progress in answering the question what is the production function for student outcomes. We will study ?, which analyses the Tennessee Student/Teacher Achievement Ratio Experiment, STAR in short.

### 8.1 The STAR Experiment

Starting in 1985-1986 and lasting for four years, young pupils starting Kindergarten and their teachers were randomly allocated to to several possible groups:

1. small classes with 13-17 students
2. regular classes with 22-25 students

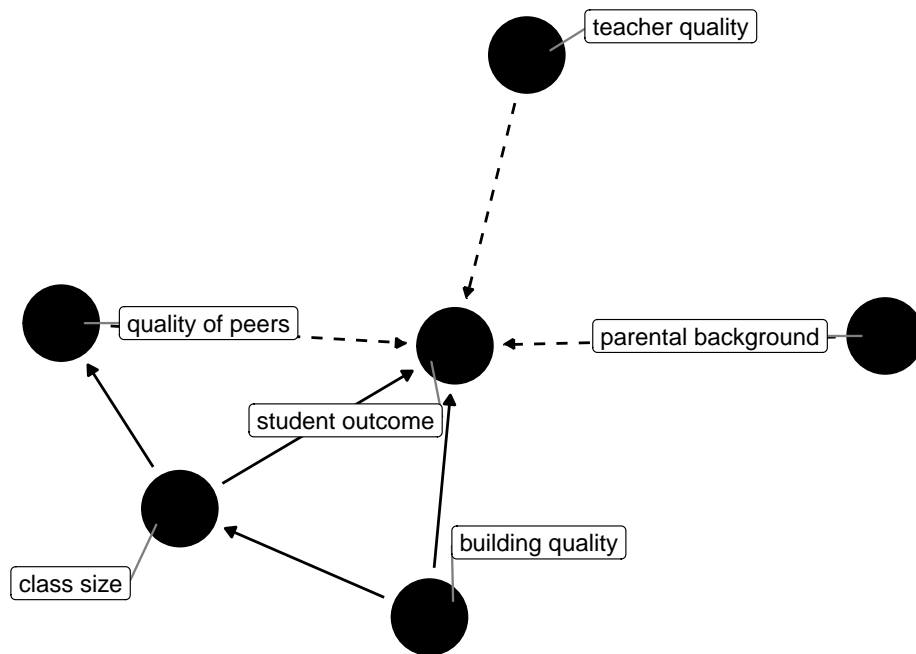


Figure 8.1: Possible Channels determining student outcomes. Dashed arrows represent potentially unobserved links.

3. regular classes with 22-25 students but with an additional full-time teaching aide.

The experiment involved about 6000 students per year, for a total of 11,600 students from 80 schools. Each school was required to have at least on class of each size-type above, and random assignment happened at the school level. At the end of each school grade (kindergarden and grades 1 thru 3) the pupils were given a standardized test. Now, looking back at figure ??, what are the complications when we'd like to assess the impact of class size on student outcome? Put differently, why can't we just look at observational data of all schools (absent any experiment!), group classes by their size, and compute the mean outcomes for each group? Here is a short list:

1. There is selection into schools with different sized classes. Suppose parents have a prior that smaller classes are better - they will try to get their kids into those schools.
2. Relatedly, who ends up being in the classroom with a child could matter (peer effects). So, if high quality kids are sorting into schools with small classes, and if peer effects are strong, we could concluded that small classes improved student outcomes when in reality this was due to the high quality of peers in class.
3. Also related, teachers could sort towards schools with smaller classes because it's easier to teach a small rather than a large class, and if there is competition for those places, higher quality teachers will have an advantage.

Now, what can STAR do for us here? There will still be selection into schools, however, once selected a school it is random whether one ends up in a small or a large class. So, the quality of peers present in the school (determined before the experiment through school choice) will be similar across small and big groups. In figure ??, you see that some factors are drawn as unobserved (dashed arrow), and some are observed (solid). In any observational dataset, the dashed arrows would be really troubling. Here, given randomisation into class sizes, we don't care whether those factors are unobserved or not: It's reasonable to assume that across randomly assigned groups, the distributions of each of those factors is roughly constant! If we can in fact proxy some of those factors (suppose we had data on teacher qualifications), even better, but not necessary to identify the causal effect of class size.

## 8.2 PO as Regression

Before we start replicating the findings in ?, let's augment our potential outcomes (PO) notation from the previous chapter. To remind you, we had defined the PO model in equation (??):

$$Y_i = D_i Y_i^1 + (1 - D_i) Y_i^0$$

and we had defined the treatment effect of individual  $i$  as in (??):

$$\delta_i = Y_i^1 - Y_i^0.$$

Now, as a start, let's assume that the treatment effect of small class is identical for all  $i$ : in that case we have

$$\delta_i = \delta, \forall i$$

Next, let's distribute the  $Y_i^0$  in (??) as follows:

$$\begin{aligned} Y_i &= Y_i^0 + D_i(Y_i^1 - Y_i^0) \\ &= Y_i^0 + D_i\delta \end{aligned}$$

finally, let's add  $E[Y_i^0] - E[Y_i^0] = 0$  to the RHS of that last equation to get

$$Y_i = E[Y_i^0] + D_i\delta + Y_i^0 - E[Y_i^0]$$

which we can rewrite in our well-known regression format

$$Y_i = b_0 + \delta D_i + u_i \tag{8.1}$$

In that formulation, the first  $E[Y_i^0]$  is the average non-treatment outcome, which we could regard as some sort of baseline - i.e. our intercept.  $\delta$  is the coefficient on the binary treatment indicator. The random deviation  $Y_i^0 - E[Y_i^0]$  is the residual  $u$ . Under only very specific circumstances will the OLS estimator  $\hat{\delta}$  identify the true Average Treatment Effect  $\delta^{ATE}$ . Random assignment ensures that the crucial assumption  $E[u|D] = E[Y_i^0 - E[Y_i^0]|D] = E[Y_i^0|D] - E[Y_i^0] = 0$ , in other words, there is no difference in nontreatment outcomes across treatment groups. Additionally, we could easily include regressors  $X_i$  in equation (??) to account for additional variation in the outcome.

With that out of the way, let's write down the regression that ? wants to estimate. Equation (2) in his paper reads like this:

$$Y_{ics} = \beta_0 + \beta_1 \text{small}_{cs} + \beta_2 \text{REG/A}_{cs} + \beta_3 X_{ics} + \alpha_s + \varepsilon_{ics} \tag{8.2}$$

where  $i$  indexes pupil,  $c$  is class id and  $s$  is the school id.  $\text{small}_{cs}$  and  $\text{REG/A}_{cs}$  are both dummy variables equal to one if class  $c$  in school  $s$  is either small,

or regular with aide.  $X_{ics}$  contains student specific controls (like gender). Importantly, given that randomization was at the school level, we control for the identify of the school with a school fixed effect  $\alpha_s$ .

Before we proceed to run this regression, we need to define the outcome variable  $Y_{ics}$ . ? combines the various SAT test scores in an average score for each student in each grade. However, given that the SAT scores are on different scales, he first computes a ranking of all scores for each subject (reading or math), and then assigns to each student their percentile in the rank distribution. The highest score is 100, the lowest score is 0.

## 8.3 Implementing STAR

Let's start with computing the ranking of grades. Let's load the data and the `data.table` package:

```
data("STAR", package = "AER")
library(data.table)
x = as.data.table(STAR)
x
```

```
#OUT>      gender ethnicity  birth      stark      star1      star2      star3 readk read1 read2 read3
#OUT>  1: female    afam 1979.50      <NA>      <NA>      <NA>    regular  NA   NA   NA  580  NA
#OUT>  2: female    cauc 1980.00    small    small    small    small  447  507  568  587  4
#OUT>  3: female    afam 1979.75    small    small regular+aide regular+aide  450  579  588  6
#OUT>  4: male      cauc 1979.75      <NA>      <NA>      <NA>    small  NA   NA   NA  686  NA
#OUT>  5: male      afam 1980.00 regular+aide      <NA>      <NA>      <NA>  439  NA   NA   NA  4
#OUT>    ---
#OUT> 11594: male    cauc 1979.50    small    small    small    small  483  590  650  675
#OUT> 11595: female    cauc 1980.50    regular    regular    regular    regular  437  533  586  65
#OUT> 11596: male    cauc 1980.00      <NA>    regular    regular    regular  NA  571  604  595
#OUT> 11597: female    afam 1980.00    regular regular+aide    regular regular+aide  431  475  542
#OUT> 11598: male    afam 1980.25 regular+aide regular+aide regular+aide regular+aide  421  468  5
#OUT>      lunchk lunch1 lunch2 lunch3 schoolk school1 school2 school3 degreek degree1
#OUT>  1:  <NA>  <NA>  <NA>  free  <NA>  <NA>  <NA>  suburban  <NA>  <NA>  <NA>
#OUT>  2: non-free  free non-free  free  rural  rural  rural  rural  bachelor bachelor
#OUT>  3: non-free  <NA> non-free non-free  suburban  suburban  suburban  suburban  bachelor ma
#OUT>  4:  <NA>  <NA>  <NA> non-free  <NA>  <NA>  <NA>  rural  <NA>  <NA>  <NA>
#OUT>  5:  free  <NA>  <NA>  <NA> inner-city  <NA>  <NA>  <NA>  bachelor  <NA>  <
#OUT>    ---
#OUT> 11594: non-free non-free non-free non-free  rural  rural  rural  rural  bachelor mast
#OUT> 11595:  free  free  free  free  rural  rural  rural  rural  bachelor bachelor ba
#OUT> 11596:  <NA> non-free non-free non-free  <NA>  suburban  suburban  suburban  <NA> bache
#OUT> 11597:  free  free  free  free inner-city inner-city inner-city inner-city  bachelor bache
#OUT> 11598: non-free  free  free non-free inner-city inner-city inner-city inner-city  bachelor b
```

```

#OUT>      ladder1  ladder2  ladder3 experiencek experience1 experience2 experience3 t
#OUT>  1:   <NA>    <NA>   level1      NA      NA      NA      30   <NA>    <NA>
#OUT>  2:  level1 apprentice apprentice      7      7      3      1   cauc
#OUT>  3: probation  level1  level1      21     32      4      4   cauc
#OUT>  4:   <NA>    <NA>   level1      NA     NA      NA     10   <NA>    <NA>
#OUT>  5:   <NA>    <NA>    <NA>      0     NA     NA     NA   cauc    <NA>
#OUT>    ---
#OUT> 11594:  level1  level3  level1      8     13     15     17   cauc
#OUT> 11595:  level1 apprentice notladder      0      7      1      7   cauc
#OUT> 11596: probation  level1  level1      NA      0      8     22   <NA>
#OUT> 11597:  level1  level1  probation      24     27      7     12   afam
#OUT> 11598:  level1  level1  level1      2     10     14     33   cauc
#OUT>      tethnicity3 systemk system1 system2 system3 schoolidk schoolid1 schoolid2 schoo
#OUT>  1:   cauc  <NA>  <NA>  <NA>      22   <NA>  <NA>  <NA>      54
#OUT>  2:   cauc   30   30   30   30     63   63   63   63
#OUT>  3:   cauc   11   11   11   11     20   20   20   20
#OUT>  4:   cauc  <NA>  <NA>  <NA>      6   <NA>  <NA>  <NA>      8
#OUT>  5:  <NA>   11  <NA>  <NA>  <NA>     19  <NA>  <NA>  <NA>
#OUT>    ---
#OUT> 11594:   cauc   21   21   21   21     49   49   49   49
#OUT> 11595:   cauc   33   33   33   33     67   67   67   67
#OUT> 11596:   cauc  <NA>   25   25   25   <NA>   58   58   58
#OUT> 11597:   cauc   11   11   11   11     22   22   22   22
#OUT> 11598:  afam   11   11   11   11     32   32   32   32

```

It's a bit unfortunate to switch to `data.table`, but I haven't been able to do what I wanted in `dplyr` :-(. Ok, here goes. First thing, you can see that this data set is wide. First thing we want to do is to make it long, i.e. reshape it so that it has 4 ID columns, and several measurements columns thereafter. First, let's add a student ID:

```

x[, ID := 1:nrow(x)] # add student ID

# `melt` a data.table means to dissolve it and reassemble for some ID variables
mx = melt.data.table(x,
  id = c("ID", "gender", "ethnicity", "birth"),
  measure.vars = patterns("star*", "read*", "math*", "schoolid*",
    "degree*", "experience*", "tethnicity*", "lun"),
  value.name = c("classtype", "read", "math", "schoolid", "degree",
    "experience", "tethniticy", "lunch"),
  variable.name = "grade")

levels(mx$grade) <- c("stark", "star1", "star2", "star3") # reassign levels to grade fac
mx[, 1:8] # show first 8 cols

```

```

#OUT>      ID gender ethnicity  birth grade  classtype read math
#OUT>  1:    1 female      afam 1979.50 stark    <NA>   NA   NA

```



```
#OUT> 2: 2 female cauc 1980.00 stark small 447 473
#OUT> 3: 3 female afam 1979.75 stark small 450 536
#OUT> 4: 4 male cauc 1979.75 stark <NA> NA NA
#OUT> 5: 5 male afam 1980.00 stark regular+aide 439 463
#OUT> ---
#OUT> 46388: 11594 male cauc 1979.50 star3 small 675 678
#OUT> 46389: 11595 female cauc 1980.50 star3 regular 654 651
#OUT> 46390: 11596 male cauc 1980.00 star3 regular 595 672
#OUT> 46391: 11597 female afam 1980.00 star3 regular+aide 624 610
#OUT> 46392: 11598 male afam 1980.25 star3 regular+aide 580 577
```

You can see here that for example pupil ID=1 was not present in kindergarden, but joined later. We will only keep complete records, hence we drop those NAs:

```
mx = mx[complete.cases(mx)]
mx[ID==2] # here is pupil number 2
```

```
#OUT> ID gender ethnicity birth grade classtype read math schoolid degree experience tethnity
#OUT> 1: 2 female cauc 1980 stark small 447 473 63 bachelor 7 cauc non-free
#OUT> 2: 2 female cauc 1980 star1 small 507 538 63 bachelor 7 cauc free
#OUT> 3: 2 female cauc 1980 star2 small 568 579 63 bachelor 3 cauc non-free
#OUT> 4: 2 female cauc 1980 star3 small 587 593 63 bachelor 1 cauc free
```

Ok, now on to standardizing those read and math scores. you can see they are on their kind of arbitrary SAT scales

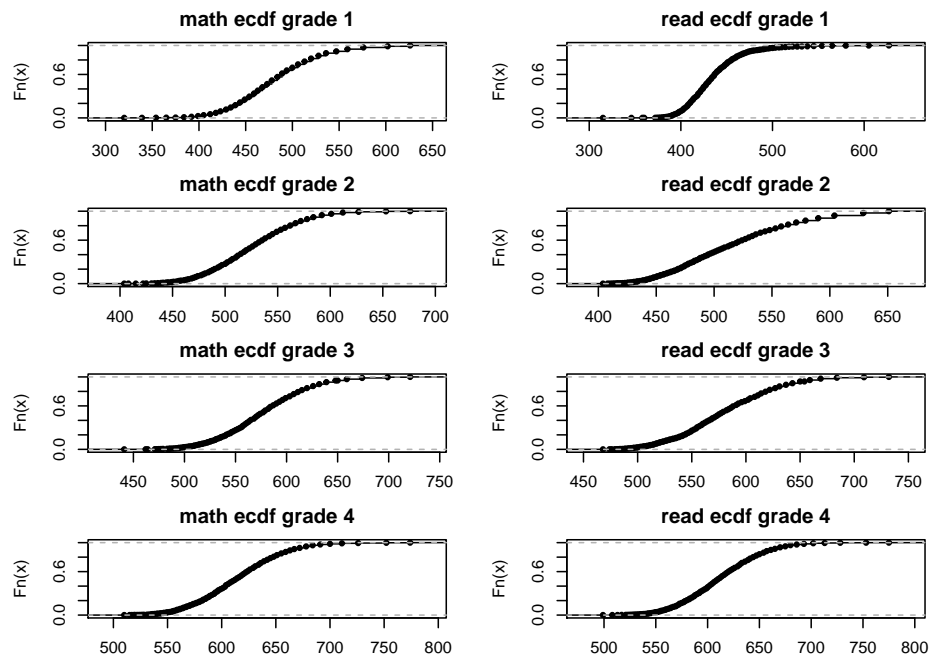
```
mx[,range(read)]
```

```
#OUT> [1] 315 775
```

First thing to do is to create an empirical cdf of each of those scores within a certain grade. That is the ranking of scores from 0 to 1:

```
setkey(mx, classtype) # key mx by class type
ecdfs = mx[classtype != "small", # subset data.table to this
  list(readcdf = list(ecdf(read)), # create cols readcdf and mathcdf
    mathcdf = list(ecdf(math))
  ),
  by = grade] # by grade

# let's look at those cdf!
om = par("mar")
par(mfcol=c(4,2),mar = c(2,om[2],2.5,om[4]))
ecdfs[,.SD[,plot(mathcdf[[1]],main = paste("math ecdf grade",.BY))],by = grade]
ecdfs[,.SD[,plot(readcdf[[1]],main = paste("read ecdf grade",.BY))],by = grade]
```



```
par(mfcol=c(1,1),mar = om)
```

You can see here how the cdf maps SAT scores (650, for example), into the interval  $[0, 1]$ . Now, in the `ecdfs` `data.table` object, the `readcdf` column contains a function (a cdf) for each grade. We can evaluate the observed test scores for each student in that function to get their ranking in  $[0, 1]$ , by grade:

```
setkey(ecdfs, grade) # key ecdfs according to `grade`
setkey(mx, grade)

z = mx[, list(ID, perc_read = ecdfs[.(BY), readcdf][[1]](read),
             perc_math = ecdfs[.(BY), mathcdf][[1]](math)),
       by=grade] # stick `grade` into `ecdfs` as `.BY`

z[, score := rowMeans(.SD)*100, .SDcols = c("perc_read", "perc_math")] # take average of
# and multiply by 100, so it's comparable to Krueger

# merge back into main data
mxz = merge(mx, z, by = c("grade", "ID"))

# make a plot
ggplot(data = mxz, mapping = aes(x = score, color=classtype)) + geom_density() + facet_v
```

You can compare figure ?? to ? figure 1. You can see that the density estimates are almost identical, the discrepancy comes mainly from the fact that we split the regular classes also by with/without aide.

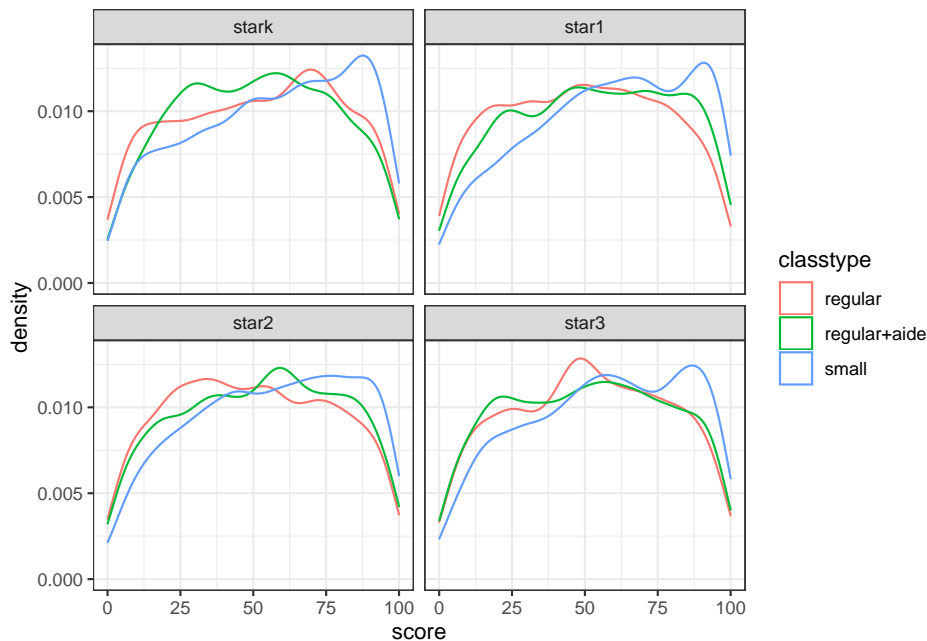


Figure 8.2: Reproducing Figure I in @krueger1999

So far, so good! Now we can move to run a regression and estimate (??).

```
# create Krueger's dummy variables
mxz = as_tibble(mxz) %>%
  mutate(small = classtype == "small",
         rega = classtype == "regular+aide",
         girl = gender == "female",
         freelunch = lunch == "free")

# reproduce columns 1-3
m1 = mxz %>%
  group_by(grade) %>%
  do(model = lm(score ~ small + rega, data = .))
m2 = mxz %>%
  group_by(grade) %>%
  do(model = lm(score ~ small + rega + schoolid, data = .))
m3 = mxz %>%
  group_by(grade) %>%
  do(model = lm(score ~ small + rega + schoolid + girl + freelunch, data = .))

# get school id names to omit from regression tables
school_co = grep(names(coef(m2[1,]$model[[1]])), pattern = "schoolid*", value=T)
school_co = c(unique(school_co), grep(names(coef(m3[1,]$model[[1]])), pattern = "schoolid*", value=T))
```

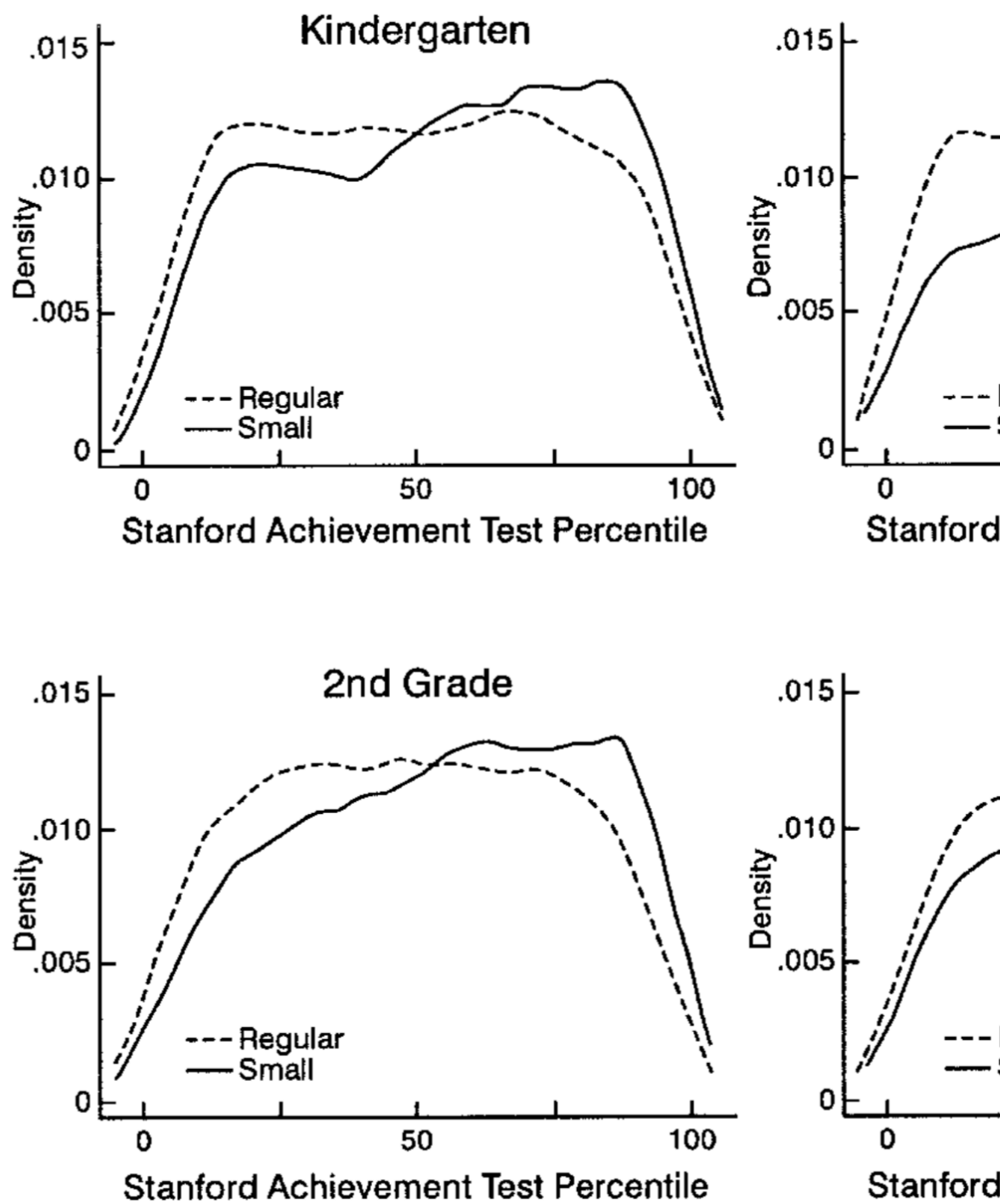


Figure 8.3: Outcome densities, @krueger1999 figure 1.