

# Chapter 4: Readyng the Data

January 9, 2019

## 1 Concepts

### 1.1 Text is data

Speeches contains a wealth of information. As humans, we are taught to understand verbal and written communication – pick out the nouns, verbs, and adjectives, then combine the information to decipher meaning. Take the following excerpt from the 2010 State of the Union:

Now, one place to start is serious financial reform. Look, I am not interested in punishing banks. I'm interested in protecting our economy. A strong, healthy financial market makes it possible for businesses to access credit and create new jobs. It channels the savings of families into investments that raise incomes. But that can only happen if we guard against the same recklessness that nearly brought down our entire economy. We need to make sure consumers and middle-class families have the information they need to make financial decisions. We can't allow financial institutions, including those that take your deposits, to take risks that threaten the whole economy.

To many, text might not be considered data, but they are. Our minds analyze the text and arrive at inferences about what the author or speaker is signaling. Analyzing text is largely a matter of identifying key terms:

~~Now, one place to start is serious financial reform. Look, I am not interested in punishing banks. I'm interested in protecting our economy. A strong, healthy financial market makes it possible for businesses to access credit and create new jobs. It channels the savings of families into investments that raise incomes. But that can only happen if we guard against the same recklessness that nearly brought down our entire economy. We need to make sure consumers and middle-class families have the information they need to make financial decisions. We can't allow financial institutions, including those that take your deposits, to take risks that threaten the whole economy.~~

Much like the logic that guides keyword identification, text can be shaped from an unstructured dataset into a well-defined, structured dataset:

Table 1: Most frequent terms found in excerpt.

Terms	Frequency of Term	Number of Characters
financial	4	9
economy	3	7
families	2	8
interested	2	10

Of course, this process could be done manually, but imagine sorting through all 7,304 words in the 2010 address or scaling the process to the roughly *1.9 million words* in State of the Union addresses between 1790 and 2016. All the steps required to convert unstructured text into usable data can be done with a little bit of planning, technical imagination and data manipulation. Every little detail about the data needs to be considered and meticulously converted into a usable form. From a data format perspective, capitalized characters are not the same as lower case. Contractions are not the same as terms that are spelled out. Punctuation affects spacing. Carriage returns and new line markers, while not visible in reading mode, are recorded.

Let's take one line from above and dissect the changes that need to be made:

“We need to make sure consumers and middle-class families have the information they need to make financial decisions. We can’t allow financial institutions, including those that take your deposits, to take risks that threaten the whole economy.”

We then turn everything into lower case so all letters of the alphabet are read the same.

“we need to make sure consumers and middle-class families have the information they need to make financial decisions. we can’t allow financial institutions, including those that take your deposits, to take risks that threaten the whole economy.”

Then, we get rid of punctuation by substituting values with empty quotations ("").

“we need to make sure consumers and middleclass families have the information they need to make financial decisions we cant allow financial institutions including those that take your deposits to take risks that threaten the whole economy”

Each space between each word can be used as a *delimiter* that can be used as a symbol for a program to break apart words into elements in a list.

Table 2: Terms

we	families	financial	those	that
need	have	decisions	that	threaten
to	the	we	take	the
make	information	cant	your	whole
sure	they	allow	deposits	economy
consumers	need	financial	to	
and	to	institutions	take	
middleclass	make	including	risks	

There are words in there that don’t add much value as they are commonplace and filler. In text processing, these words are known as *stop words*. In each domain, the list of stop words likely differs, thus data scientists may need to build a customized list. For simplicity, we’ve used a stop words list that is used in MySQL – an open source relational database management system. The result is the list of remaining words.

Table 3: Terms after removing stop words

make	information	financial	risks
consumers	make	institutions	threaten
middleclass	financial	including	economy
families	decisions	deposits	

From that data, we can aggregate the data into a form that is meaningful to answer a research question. For example, the frequency of words may provide a clue as to what the text is about. In this case, each “financial” and “make” appear twice in the text, perhaps indicating that there is an orientation towards action (make) for financial considerations.

Table 4: Term Frequencies

Term	Freq	Term	Freq
financial	2	including	1
make	2	information	1
consumers	1	institutions	1
decisions	1	middleclass	1
deposits	1	risks	1

Term	Freq	Term	Freq
economy	1	threaten	1
families	1		

This is just the tip of the iceberg. Text processing is just one aspect of readying data for use.

There are entire texts dedicated to data engineering and manipulating data. This chapter is intended to be a simple primer, providing a brief review of programming paradigms that are necessary to shape data into usable form so it can be the engine that drives impact in the public and social sectors. The chapter begins with illustrating a workflow for retrieving and assembling data, then explores common methods of manipulating data values and formats, and closes with manipulating data structures.

## 2 Retrieval and Assembly

The retrieval and assembly of a data set can be one of the greatest barriers to launching a project. Here, we lay out a few practices that make the process far simpler.

There are a multitude of data storage formats in use such as extensible markup language (XML), comma separated values (CSV), Excel files among others. Each has its quirks. Some require special software to load while others simply need time and patience to crack the code. Fortunately, R is equipped to load virtually all data formats. Below is a recommended set of functions that are easy to use and flexible.

Table 5: Recommended functions to load an assortment of files.

Function	Package	Description
read_excel	readxl	Load Excel files into a data frame. Note that with excel files, generally need to specify which sheet to load using the <i>sheet</i> argument and may need to skip a number of rows at the beginning of worksheet to indicate the header row using <i>skip</i> .
read_csv	readr	Load any comma separated file into data frame format.
read_delim	readr	Load delimited file using any delimiter such as tab delimited (“\t”) using the <i>sep</i> argument.
readLines	base R	Read a text file, line by line. This is helpful for working with free form text.
xmlToDataFrame	XML	Read XML into a data frame if the structure is simple and flat. Note that XML files tend to have a more complex hierarchical structure.
xmlToList	XML	If the XML is complicated, read each element as a list.
fromJSON	rjson	Read JSON into list. If JSON is a flat, non-hierarchical file, then the resulting object can be simply converted into a data frame.
read_dta	haven	Load Stata data file.
read_sas	haven	Load SAS data files.
read_spss	haven	Load SPSS data files.
load	base R	Load saved R data set that can contain multiple objects. Note that this does not need to be assigned to another object.
readRDS	base R	Read individually saved R object.

In the course of writing this text, the authors have come across a number of misconceptions about how to work with data, leading to haphazardous manual editing of data files. In particular, we have found that:

- Analysts will often manually open files and delete unnecessary rows at the top of a file as functions

may not be able to read in files.

- The number of columns will differ, requiring columns to be deleted or empty columns added before loading.
- Column names differ, requiring manual editing of headers.
- The edits often overwrite the original file. This means that an edit were erroneous, there is no hope in restoring the original data.
- Manual edits are learned and shared through as an oral history.

These manual steps become tedious endeavor when presented with hundreds if not thousands of files. When loading and assembling data sets, the raw data should be stored as a copy. It then should be read into the programming environment where all processing is conducted. Since the processing is coded, it can be reproduced. The result is then stored as a new file.

In this section, we illustrate a simple workflow that is efficient and accurate. Our goal is to see if Energy Information Administration (EIA) spot gas prices recorded in New York Harbor and the US Gulf Coast are correlated. To do this, we will need to assemble a time series data set taking the following:

- Load a CSV of spot gas prices from the US Gulf Coast as a data frame.
- Load two Excel worksheets of spot gas prices from NY Harbor from two different time windows as data frames, then append one to the other.
- Combine the two into a single data frame as illustrated below, then calculate the Pearson Correlation coefficient.

Table 6: Target format

date	price1	price2
2016-12-28	1.739	1.702
2016-12-29	1.729	1.703
2016-12-30	1.722	1.699

**Loading CSVs.** Before reading in a file, it is generally a good idea to open the file in a code editor to take a peek at the structure, but “blind loading” the file can also give clues as well. Let’s start with the CSV `doe_usgulf.csv` using the `read_csv` function from the `readr` package, assign the data to the object `gulf`, then inspect the first few rows using `head`.

```
pacman::p_load(readr)
gulf <- read_csv("data/doe_usgulf.csv")
head(gulf, 3)
```

```
## # A tibble: 3 x 2
##   Date      `U.S. Gulf Coast Conventional Gasoline Regular Spot Price FOB (D~
##   <chr>                                     <dbl>
## 1 1/2/14                                     2.52
## 2 1/3/14                                     2.49
## 3 1/6/14                                     2.52
```

It appears that the data was successfully loaded. The data is read in as a *tibble*, which is a variant of a data frame designed for more efficient data manipulation. The headers require cosmetic changes. Generally, it is good practice to keep column names as the same case (lower case) and avoid the use of spaces and punctuation with the exception of periods. We can edit the column headers, which can be called as a vector using `colnames(gulf)`, by assigning a new string vector that turns all letters into lower case (`tolower`) and replaces substrings (`gsub`) that are either spaces or punctuation (`[[:space:]]` and `[[:punct:]]`) with periods. These string manipulation commands will be covered in more detail later in this chapter.

```
colnames(gulf) <- tolower(gsub("[:space:][:punct:]", ".", colnames(gulf)))
```

However, as these column names are quite long, we can opt to overwrite them with something more concise. Notice that we did not need to open the CSV and manually edit the column names.

```
colnames(gulf) <- c("date", "gulf.price")
```

**Loading an Excel file.** Excel files tend to have a few more curve balls than CSVs. While Excel files are more readable by aesthetic and cosmetic styling, parsing is challenging. Using `read_excel` from the `readxl` package, we take a first look at the first sheet (`sheet = 1`) of our Excel file `doe_ny.xlsx`, then use `head` to examine what was loaded. Unlike our CSV, our first attempt requires some tweaks: The `Date` header is treated as the first row of data and there are two additional empty columns.

```
pacman::p_load(readxl)
attempt1 <- read_excel("data/doe_ny.xlsx", sheet = 1)
head(attempt1, 3)
```

```
## # A tibble: 3 x 4
##   Sourcekey EER_EPMRU_PF4_Y35NY_DPG X__1 X__2
##   <chr>    <chr>                  <lg1> <lg1>
## 1 Date      New York Harbor Conventional Gasoline Regular Spo~ NA    NA
## 2 41641      2.718                      NA    NA
## 3 41642      2.6709999999999998    NA    NA
```

At this point, some may be tempted to manually delete the first row. Fortunately, we can skip rows when using most loading functions using a `skip` argument. Below, we read in the first and second sheets of `doe_ny.xlsx`, skipping the first row in each, then rename the headers.

```
#Load first two sheets
sheet1 <- read_excel("data/doe_ny.xlsx", sheet = 1, skip = 1)
sheet2 <- read_excel("data/doe_ny.xlsx", sheet = 2, skip = 1)

#Rename columns
colnames(sheet1) <- c("date", "ny.price")
colnames(sheet2) <- c("date", "ny.price")
```

**Appending data.** It turns out that `sheet1` and `sheet2` are from the same time series of spot covering different time windows. We can construct a new, more complete data frame by appending one data frame to the other using the `rbind` function.

```
ny <- rbind(sheet1, sheet2)
```

The above scenario is the ideal: two data frames with the same number of columns and the same column names. But what if one data frame has at least one more column than the other? The `rbind.fill` function in the `plyr` package appends the two data frames together, filling any additional missing columns with `NA` values. Below, rather than loading the entire `plyr` library, we can selectively load functions using the double colon operator `::`.

```
ny <- plyr::rbind.fill(sheet1, sheet2)
```

Note that if two columns in two different data frames represent the same concept but have different names, `rbind.fill` does not know to make the connection. Column names are interpreted literally, meaning that it is of the utmost importance that headers are exactly labeled.

**Combining data by rows.** The two data frames `gulf` and `ny` contain the same number of rows. If they are in the same exact desired order and have the same number of rows, we can use the `cbind` function to join the datasets by rows. The `cbind` function will join all columns in both data frames together, including the dates. To ensure the data is neatly organized, we keep all columns from the `gulf` dataset and only the second column from the `ny` data set. The resulting data frame contains three columns containing daily spot prices from two sources.

We calculate the correlation using the `cor` function finding a value of  $\rho = 0.989617$  – the two sets of prices are highly correlated. Note that `cor` only accepts numeric values, thus the date column is temporarily dropped.

```
#Combine data
prices <- cbind(gulf, ny[,2])

#Calculate correlation
cor(prices[, -1])

##           gulf.price ny[, 2]
## gulf.price  1.000000 0.989617
## ny[, 2]     0.989617 1.000000
```

### 3 Manipulating values

More often than not, working with raw data means working with string values. It is easy to imagine that necessary data are embedded in a sentence. For example, below is a vector of four string values:

```
budget <- c("Captain's Log, Stardate 1511.8. I have $10.20 for a big galactic mac.",
            "The ensign has $1,20 in her pocket.",
            "The ExO spent has $0.25 left after paying for overpriced warp core fuel.",
            "Chief medical officer is the high roller with $53,13.")
```

What if we need to extract the total available funds available to buy galactic big macs? All four elements contain dollar values, which can benefit from feature engineering. To do so, we use a combination of text manipulation functions and *regular expressions* or *regex* – a series of characters that describe a regularly occurring text pattern.

First, commas should be replaced with a period using `str_replace_all()`, assigning the result to a new object `new`. Note that in some regions, such as Europe, commas are used as decimals rather than periods.

```
require(stringr)
new <- str_replace_all(budget, ",", "\\.")
```

Second, find the elements that contain the following pattern: a dollar sign followed by one to two digits, followed by a period, then another two digits (`\\$\\d{1,2}\\d{2}`). Using the `str_extract_all` function

in the `stringr` package, we can extract all matching substrings, specifying `simplify = TRUE` to return a vector:

```
funds <- str_extract_all(new, "\\$\\d{1,2}\\\\.\\d{2}", simplify = TRUE)
print(funds)

##      [,1]
## [1,] "$10.20"
## [2,] "$1.20"
## [3,] "$0.25"
## [4,] "$53.13"
```

Third, we should replace dollar sign with blank, then strip out any leading white space using `str_trim()`.

```
funds <- str_replace_all(funds, "\\$", "")
funds <- str_trim(funds)
print(funds)

## [1] "10.20" "1.20" "0.25" "53.13"
```

Lastly, we convert the character vector to numeric values using `as.numeric`, then sum the vector.

```
money <- as.numeric(funds)
print(paste0("Total galactic big mac funds = $", sum(money)))

## [1] "Total galactic big mac funds = $64.78"
```

A number of observations. In steps one through three, you will have noticed that the characters "\$", ".", and "d" were preceded by double backslash. These are known as *escaped characters* as the double backslash preceding the characters changes their meanings. In step two, a sequence of unusual characters (`\\$\\d{1,2}\\\\.\\d{2}`) was used to find the `$x.xx` pattern, which can be broken into specific commands:

- `\\$` is a dollar sign.
- `\\d{1,2}` is a series of numerical characters that is between one to two digits long.
- `\\.` is a period.
- `\\d{2}` is a series of numerical characters that is exactly two digits long.

Mastering *regex* is a productivity multiplier, opening the possibility of ultra-precise text replacement, extraction, and other manipulation. Imagine scenarios where raw data is not quality controlled and mass errors plague the usefulness of the data. An analyst may spend days if not weeks or months cleaning data by hand (or rather through find and replace). With *regex*, haphazard cleaning is no longer an issue. To make the most of *regex* requires a command of both *text manipulation functions* that are designed to interpret *regex* as well as *regex* itself.

### 3.1 Text manipulation functions

*Find and replace* functionality are common in word processing and spreadsheet softwares. But what does it take to do apply it to a sizable amount of data? The following text manipulation functions simplify string manipulation, whether to extracting or editing specific values. We provide two sets of functions. The first are core functions that are furnished with R – these closely resemble string manipulation functions implemented in most programming languages. While they are powerful, they at times may require some additional trial

and error to use. Alternatively, the **stringr** package simplifies text manipulation and syntax, and in some case extends functionality. For the most part, we will use **stringr** in this text.

Table 7: Recommended text manipulation functions.

Description	Base R	stringr
Returns either the index position of a matched string or the string containing the matched portion.	grep	str_which
Returns a logical vector indicating if a matched string was found.	grepl	str_detect
Searches for a specified pattern and replaces with user-specified substring.	gsub	str_replace_all
Remove matched pattern from string.	gsub	str_remove_all
Returns the first position of matched pattern in a string.	regexpr	str_locate
Returns the position of all matched patterns in a string.	gregexpr	str_locate_all
Extract substring based on matched pattern. Note that regmatches is used in conjunction with regexpr or gregexpr.	regmatches	str_extract_all
Splits strings into a list of values based on a delimiter.	strsplit	str_split
Extract substring based on start and end position.	substr	str_sub
Trim whitespace on either end of string (excessive spaces)	trimws	str_trim
Returns number of characters in string	nchar	str_length
Returns the number of matched patterns.		str_count
Convert to upper case	toupper	str_to_upper
Convert to lower case.	tolower	str_to_lower
Convert to title case.		str_to_title
Pad string (e.g. add leading zeros to string)	sprintf	str_pad

Some basic tasks can be accomplished such as exact matches of specific text. As will be seen later, these functions combined with regex are quite powerful. To illustrate the basic functionality, let's assume we have four sentences that indicate when four US laws were signed.

```
laws <- c(". Dodd-Frank Act was signed into federal law on July 21, 2010.",
  "Glass-Steagall Act was signed into federal law by FDR on June 16, 1933",
  "Hatch Act went into effect on August 2, 1939",
  "Sarbanes-Oxley Act was signed into law on July 30, 2002")
```

Suppose we need to find acts that are named for two congressmen. The **str\_detect()** function can be used to produce a boolean vector if the regex pattern **[A-z]{1,}-[A-z]{1,}** is matched (more on this pattern in the next section).

```
str_detect(laws, "[A-z]{1,}-[A-z]{1,}")
```

```
## [1] TRUE TRUE FALSE TRUE
```

We can extract the specific senator names using **str\_extract\_all**, returning a vector of string values.

```
str_extract_all(laws, "[A-z]{1,}-[A-z]{1,}", simplify = TRUE)
```

```
##      [,1]
## [1,] "Dodd-Frank"
## [2,] "Glass-Steagall"
```



```
## [3,] ""
## [4,] "Sarbanes-Oxley"
```

## 3.2 Regular Expressions (regex)

As we can see, much of the secret ingredient of the text manipulation functions is the regex, which are powerful commands that give coders the flexibility to search data and surface results following a substring pattern. Before proceeding into more complex string combinations, knowledge of a few cleverly designed capabilities may go a long way:

- (1) Alternatives (e.g. “OR” searches) can be surfaced by using a pipe “|”. For example, a string search for “Bob or Moe” would be represented as “Bob|Moe”.
- (2) The extent of a search should be denoted by parentheses (). For example, a string search for “Jenny” or an alternative spelling like Jenn would be represented as “Jenn(y|i)”.
- (3) A search for one specific character should be placed between square brackets [].
- (4) The length of a match is specified using curly brackets {}.

In New York City, the famed avenue *Broadway* is may be written and abbreviated in a number of ways. The vector `streets` contains a few instances of spellings of Broadway mixed in with other streets that start with the letter B.

```
#A sampling of street names
streets <- c("Bruckner Blvd", "Bowery", "Broadway", "Bway", "Bdway",
            "Broad Street", "Bridge Street", "B'way")

#Search for two specific options
str_detect(streets, "Broadway|Bdway")

## [1] FALSE FALSE  TRUE FALSE  TRUE FALSE FALSE FALSE

#Search for two variations of Broadway
str_detect(streets, "B(road|')way")

## [1] FALSE FALSE  TRUE FALSE FALSE FALSE FALSE  TRUE

#Search for cases where either d or apostrophe are between B and way
str_detect(streets, "B[d']way")

## [1] FALSE FALSE FALSE FALSE  TRUE FALSE FALSE  TRUE
```

### 3.2.1 Escaped Characters

Quite a few single characters hold a special meaning in addition to the literal meaning. To disambiguate their meaning, a backslash precedes these characters to denote the alternative meaning. A few include:

- `\n`: new line
- `\r`: carriage return
- `\t`: tab
- `\'`: single quote when in a string enclosed in single quotes ('Nay, I can\'t')
- `\"`: double quote when in a string enclosed in double quotes ("I have a \"guy\".")

In other cases, double backslashes should be used:

- `\\.`: period. Otherwise, un-escaped periods indicate searches for *any* single character.
- `\\$`: dollar sign. A dollar sign without backslashes indicates to find patterns at the end of a string.

### 3.2.2 Character Classes

A *character class* or *character set* is used to identify specific characters within a string. How would one represent “12.301.1034” or “?!?!?”? One or more of the following character classes can do the job:

- `[:punct:]`: Any and all punctuation such as periods, commas, semicolons, etc. For specific punctuation, simply enclose the characters between two brackets. For example, to find only commas and carots, use `[<>,]`.
- `[:alpha:]`: Alphabetic characters such as a, b, c, etc. With other languages including R, it is commonly written as `[a-z]` for lower case, `[A-Z]` for upper case, and `[A-Za-z]` for mixed case.
- `[:digit:]`: Numerical values. With other languages including R, it is commonly written as `\\d` or `[0-9]`. For any non-digit, write `\\D`.
- `[:alnum:]`: Alphanumeric characters (mix of letters and numbers). With other languages including R, it is indicated using to as `[0-9A-Za-z]` or `\\w`. For any non-alphanumeric character, use `\\W`.
- `[:space:]`: Spaces such as tabs, carriage returns, etc. For any white space, use `\\s`. For any non-whitespace character, use `\\S`.
- `[:graph:]`: Human readable characters including `[:alnum:]` and `[:punct:]`.
- `\\b`: Used to denote “whole words”. `\\b` should be placed before and after a regex pattern. For example, `\\b\\w{10}\\b` indicates a 10 letter word.

These are only a few of the most important character classes. It is worth keeping in mind that these are R-specific character classes and implementations may differ from one programming language to another.

### 3.2.3 Quantifiers

Each character class on its own indicates a search for any character of that class. For example, a search for the pattern `[:alpha]` will return any element with an alphabetic character – not very specific. In practice, most character searches will involve a search for more than just one character. To indicate such a search, regex relies on *quantifiers* to indicate the length of patterns. For example, a search for a year between the year 1980 and 2000 will require exactly four digits, but a search for the speed of a gust of wind will likely vary between 1 and 3 digits. The following six quantifiers provide flexibility and specificity to effectively accomplish search tasks:

- `{n}`: match pattern n times for a preceding character class. For example `\\d{4}` looks for a four digit number.
- `{n, m}`: match pattern at least n-times and not more than m times for a preceding character class. For example `\\d{1,4}` looks for one to four digit number.
- `{n, }`: match at least n times for a preceding character class. For example `\\d{4,}` searches for a number that has at least four digits.
- `*`: Wildcard, or match at least 0 times.
- `+`: Match at least once.
- `?`: Match at most once.

In the example below, quantifiers are used to extract specific number patterns with a high degree of accuracy.

```
big_dates <- c("Octavian became Augustus on 16 Jan 27 BCE",
              "In the year 2000, a computer bug was expected to topple society.",
              "In the 5400000000 years, our sun will become a red dwarf.")

#Return a 9 digit number
str_extract(big_dates, "\\d{9}")

## [1] NA          NA          "5400000000"

#Return a 4 digit substring that is flanked by empty value at either end
str_extract(big_dates, "\\b\\d{4}\\b")

## [1] NA          "2000" NA

#Match a date that follows 16 January 27 BCE
str_extract(big_dates, "\\d{2}\\s\\w{3}\\s\\d{2}\\s\\w{3}")

## [1] "16 Jan 27 BCE" NA          NA
```

### 3.3 Positions

Regex builds in functionality to search for patterns based on location of a substring, such as at the start or end of a string. There are quite a few other position matching patterns, but the following two are the main workhorses:

- `$`: Search at the end of a string.
- `^`: Start of string when placed at the beginning of a regex pattern.

To demonstrate these patterns, we'll apply `str_extract` to three headlines from the BBC.

```
headlines <- c("May to deliver speech on Brexit",
              "Pound falls with May's comments",
              "May: Brexit plans to be laid out in new year")
print(headlines)

## [1] "May to deliver speech on Brexit"
## [2] "Pound falls with May's comments"
## [3] "May: Brexit plans to be laid out in new year"

#Find elements that contain May at the beginning of the string
str_extract(headlines, "^May")

## [1] "May" NA          "May"

#Find elements that contain Brexit at the beginning of the string
str_extract(headlines, "Brexit$")

## [1] "Brexit" NA          NA
```

### 3.3.1 DIY: Redact PII

In an increasingly digital world, data privacy is a sensitive issue that has taken center stage. At the center of it is the safeguarding of Personally identifiable information (PII). Legislation in the European Union, namely the General Data Protection Regulation or GDPR, requires companies to protect the personal data of European Union (EU) citizens associated with transactions conducted in the EU ((“EU Gdpr.org,” n.d.)). The US Census Bureau, which administers the decennial census, must apply disclosure avoidance practices in order so that individuals cannot be identified (( ???)). Anonymization has become a common task when working with sensitive PII data, spanning complex probabilistic methods to simple redaction. We focus here on the latter.

The first element below, for example, contains hypothetical PII and sensitive information – John’s social security number and balance in his savings account are shown. When presented with many lines of sensitive information, one could review each sentence and manually redact sensitive information, but given thousands if not millions of sensitive information, this is simply not feasible.

```
statement <- c( "John Doe (SSN: 012-34-5678) has $2303 in savings in his account.",
               "Georgette Smith (SSN: 000-99-0000) owes $323 to the IRS.",
               "Alexander Doesmith (SSN: 098-76-5432) was fined $14321 for overdue books.")
```

Using a combination of regex and `stringr` we can redact sensitive information with placeholders. To remove the SSN, we need a regex pattern that captures a pattern with three digits (`\\d{3}`), a hyphen, two digits (`\\d{2}`), a hyphen, then four digits (`\\d{4}`), or when combined: `\\d{3}-\\d{2}-\\d{4}`. The matched pattern is then replaced with `XXXXX`.

```
pacman::p_load(stringr)
new_statement <- str_replace(statement, "\\d{3}-\\d{2}-\\d{4}", "XXXXX")
print(new_statement)
```

```
## [1] "John Doe (SSN: XXXXX) has $2303 in savings in his account."
## [2] "Georgette Smith (SSN: XXXXX) owes $323 to the IRS."
## [3] "Alexander Doesmith (SSN: XXXXX) was fined $14321 for overdue books."
```

Next, we replace the dollar value by matching a string that starts with the dollar sign (`\\$`) followed by at least one digit (`\\d{1,}`). And finally, the John Doe’s first and last name are replaced by looking for two substrings that each have at least one uppercase letter with an unspecified length (`[A-z]{1,}` `[A-z]{1,}`) and are found at the beginning of the string (`^`). The resulting sentence has little to no information about the individual in question.

```
#Find a replace dollar amount
new_statement <- str_replace(new_statement, "\\$(\\d{1,})", "XXXXX")

#Find and replace first and last name
new_statement <- str_replace(new_statement, "^[A-z]{1,} [A-z]{1,}", "XXXXX")

print(new_statement)
```

```
## [1] "XXXXX (SSN: XXXXX) has XXXXX in savings in his account."
## [2] "XXXXX (SSN: XXXXX) owes XXXXX to the IRS."
## [3] "XXXXX (SSN: XXXXX) was fined XXXXX for overdue books."
```

### 3.3.2 DIY: Extract PII

The inverse of this task is to extract all information and create a data set. Using the same regex patterns, we can apply `str_extract` to create a three variable data frame containing person name, SSN and money – all done with minimal effort.

```
ref_table <- data.frame(name = str_extract(statement, "[A-z]{1,} [A-z]{1,}"),
                        ssn = str_extract(statement, "\\d{3}-\\d{2}-\\d{4}"),
                        money = str_extract(statement, "\\$(\\d{1,})"))
print(ref_table)
```

```
##           name      ssn  money
## 1      John Doe 012-34-5678 $2303
## 2  Georgette Smith 000-99-0000 $323
## 3 Alexander Doesmith 098-76-5432 $14321
```

## 3.4 Working with Dates

Working dates can be challenging. Most programming languages are unable to “auto-magically” recognize a date variable, requiring coders to specify the format using a set of date symbols. The `lubridate` package significantly lowers the bar for R users, making the process of working and manipulating dates more seamless and less prone to error.

```
require(lubridate)
```

Upon loading `lubridate`, we can convert string and numeric objects with values that represent dates into date objects. The `as_date` function is most intuitive, automatically detecting the date format and assumes that the date is recorded in UTC.

```
d0 <- as_date("2010-08-20")
```

In case `as_date` is unable to detect the date, the user can more formally define the date format using functions such as `mdy` and `ymd`. Both are able to accommodate both string and numeric formats.

```
d1 <- mdy("01/20/2010")
d2 <- mdy_hm("01/20/2010 00:00 AM")
d3 <- ymd(20100101)
```

Once data are converted into date objects, it is easy to process derivative information. To calculate duration between two dates is as simple as subtraction.

```
d1 - d3
```

```
## Time difference of 19 days
```

Often times will need to extract parts of the date. With `lubridate`, each `year`, `month`, `quarter`, `day`, `hour`, `minute`, and `second` can be extracted using a dedicated function. For example, we can extract the `year` and `quarter` and concatenate into a string using `paste0`.

```
y1 <- year(d0)
q1 <- quarter(d0)
paste0(y1, "Q", q1)
```

```
## [1] "2010Q3"
```

In base R, we can alternatively use the `format` function to specify an output format using date symbols. The example below outputs a concatenation of year and month.

```
format(d0, "%Y-%m")
```

```
## [1] "2010-08"
```

## 4 Sculpting the data

The ultimate goal of data retrieval and processing is to construct a data set that is ready for analysis, modeling and visualization. This process, however, can occupy between 50% to 80% of a data scientist's time (( Lohr 2014)). Data needs to be processed into different forms depending on the use case. The number of tools available to *wrangle* data are many. But what if a standard framework can be applied make data more actionable sooner. As laid out in Wickham (2014), we can follow the principles of *tidy data*:

- Each variable is a column.
- Each observation is a row.
- Each type of observational unit forms a table.

Otherwise, a *messy* data set is one that does not follow the described structure. Let's walk through a hypothetical company's financial data. Below, each month's data is a column and each row represents a different financial concept. It is a standard way of displaying balance sheets, optimized for readability but not empirical analysis. Working with this data in the current form to produce a time series forecast can prove challenging – it's messy data. The data should be reshaped into a multivariate time series: each row contains a month's data months and each financial concept as a column.

Table 8: Example of a messy data set.

value	Dec.1.2017	Jan.1.2018	Feb.1.2018
revenue	46188	13009	8223
cost	32000	41277	1900
heads	73	36	37

Fortunately, this is a relatively simple task in this case – we simply need to transpose the data into tidy form so the data can be more easily manipulated and used for data science purposes. The tidied table can not only be used for time series and cross sectional regression analysis, but visualized in a variety of ways.

Table 9: Example of a tidy data set.

date	revenue	cost	heads
12/1/17	46188	32000	73
1/1/18	13009	41277	36
2/1/18	8223	1900	37

For the much of this text, we will tend towards tidy data and use them as the preferred data structure that to enable analysis and modeling. Tidying involves moving columns, rotating, slicing, deleting – all tasks that change the shape of the data, requiring analysts to become intimately familiar with data structures. In the remainder of this section, we introduce core concepts for sculpting data into the right shape.

## 4.1 Matrix or Data Frame?

Last chapter, we introduced matrices and data frames as options for storing multivariate data. Matrices are vectors of the same variable type, but in two dimensions ( $n \times m$  dimensions). Data frames are a generalization of matrices that allow for each column of data to hold different data types as well as refer to individual columns by a user-specified name.

*When should each matrices and data frames be used?* From a pure logistical perspective, data frames are more flexible with respect to its ability to store multiple data types. Some code libraries are built specifically for matrices and others for data frames. Ultimately, data frames are generally a safe bet. They can be subsetted, sorted, reshaped, collapsed, and merged, retaining the variable types.

Table 10: Matrices versus Data Frames.

	Matrix	Data Frame
Pros	Memory efficient. Good for advanced mathematical operations.	Store mixed types of data types. Allows user to refer to columns by an explicit name.
Cons	Able to store one data type at a time – leads to slightly more work required to manage multiple matrices. Columns can only be referred to by index number.	Not as memory efficient.

To illustrate data in action, we consider a data set containing all reported thefts, burglaries and robberies in Chicago between 2014 through 2018. The data are available on the Chicago Open Data Portal. A cleaned version in .Rda is available at [link goes here](#). The file contains a data frame `crimes` with  $n = 428,250$  with 17 variables.

```
load("data/chicago_crime.Rda")
```

Table 11: Example from crimes data set.

case.number	date	primary.type	description
HY189976	03/18/2015 09:00:00 PM	ROBBERY	ARMED: HANDGUN
HY190019	03/18/2015 10:31:00 PM	THEFT	RETAIL THEFT
HY189725	03/18/2015 12:55:00 PM	BURGLARY	FORCIBLE ENTRY

## 4.2 Array Indexes

We can think of data frames and matrices as data organized along a grid – essentially a spreadsheet. Each cell in a data frame can be located using a coordinate set of *array indexes*, or a number that indicate the location along a row or column. The first, fifth, and millionth rows have index values of 1, 5, and 1000000. The second, ninth and one-thousandth columns have values of 2, 9, 1000. This is similar to languages like FORTRAN, MATLAB, and Julia. Note that array indexes in programming languages such as Python, Perl, C, Java and Ruby start with 0 – often a point of confusion when transitioning between languages.

*Why do array indexes matter?* Understanding this is the gateway to subsetting data, or slicing and dicing as one pleases. To extract any part of the data frame requires referring to the array indexes in a specific format:

```
crimes[7, 6]
```

```
## # A tibble: 1 x 1
```

```
## primary.type
## <chr>
## 1 ROBBERY
```

where `crimes` is the name of a data frame. 7 is a row index and 6 is a column index, both of which are sandwiched in square brackets and separated by a comma. Executing the above state will extract the value in the seventh row and sixth column.

*What if we want to extract an entire row?* Simply specify the row index of interest and leave the column blank. The result is a single row data frame describing a robbery that occurred on March 18, 2015 at 11:55 PM.

```
crimes[7, ]
```

To *extract multiple records by row index* depends on whether the request is sequential or piecemeal. Below, the first line extracts a range of rows from the 2nd through 4th rows in `crimes`, whereas the second extracts two non-overlapping ranges that are included in a vector.

```
crimes[2:4, ] # apply the index range to extract rows
crimes[c(1:2, 10:11), ] # specific indexes
```

The same notation can be used *to extract all street blocks that had a crime reported* by typing the number 4 after the comma in square brackets. The number 4 is the column index that contains `block` names. Rather than returning a data frame, the result is a vector containing the partially redacted street block.

```
crimes[,4]
```

In addition, data frames provide a few additional options for extracting the `block` column.

```
crimes[, "block"] # extract column with "block" label
crimes[["block"]] # list syntax to extract column from data frame
crimes$block # compact version of data column manipulation
```

To *extract two or more columns* follows a familiar pattern, making use of either a range of column indexes or a vector of column names.

```
crimes[, 3:4] # extract multiple columns
crimes[, c("date", "block")] # multiple column labels
```

## 4.3 Subsetting

Extracting specific *subsets* is perhaps one of the most essential skills of data processing as it allows one to focus on specific aspects of the data. In addition to referring to specific array indexes to subset data, we can extract specific rows based on user-specified criteria, such as “crimes that happened in 2015” or “robberies that involved a handgun”. Here, we describe three approaches using logical vectors, the `which` command and



the `filter` command in the `dplyr` library. The first two options are possible using the pre-built functionality of R.

The **logical vector** approach involves creating a vector that describes some comparison, then passing the resulting logical vector to the data frame.

```
#Logical vector
group1 <- crimes$year == 2015

#Select based on vector
out <- crimes[group1, ]
```

The **which** approach is largely the same as above, but uses the `which` function to retrieve the row index of each TRUE value in the logical vector. The relevant row indexes are then passed to the data frame, returning a result far faster than the first option.

```
#Logical vector
group1 <- which(crimes$year == 2015)

#Select based on vector
out <- crimes[group1, ]
```

Lastly, and perhaps more convenient, is the `filter` function in the `dplyr` package. Not only is the syntax concise, the processing time is fast. The command only requires the data frame to be specified once and variables can be specified without the dollar sign notation.

```
pacman::p_load(dplyr)
out <- filter(crimes, year == 2015)
```

The `dplyr` allows for *piping* of operations, meaning that multiple steps can be combined into one function call. The pipe operator `%>%` allows the output of one step to be input into another, thereby creating a pipeline. For example, the statement below is the same as the above.

```
out <- crimes %>% filter(year == 2015)
```

To get comfortable, let's walk through a few different scenarios.

**Example 1:** Find all crimes that were armed. To do this, we use `str_detect` function from `stringr` to search for the term `ARMED` from the `description` field in `crimes`, then use the result to subset using `filter`.

```
out <- filter(crimes, str_detect(crimes$description, "ARMED"))
```

**Example 2:** Find all cases that where the `primary.type` of crime was either a `THEFT` or `BURGLARY`. So that the statement can be one compact line, we use the `%in%` operator to search if any the elements in the vector `c("THEFT", "BURGLARY")` are found in `primary.type`.

```
out <- filter(crimes, primary.type %in% c("THEFT", "BURGLARY"))
```

*Extract rows that meet two criteria* is as simple as using logical operators. The pipe operator "`|`" should be used to represent *OR*, whereas "`&`" represents *AND*. To find all cases where an `arrest` occurred after the year 2016 :

```
out <- filter(crimes, arrest == TRUE & year > 2016)
```

## 4.4 Sorting and Re-ordering

At times, data frames need to be re-arranged to make it easier to navigate the data.

*Re-arranging rows and columns.* For one thing, columns and rows can be re-arranged or even removed by listing column indexes in a desired order. By passing a vector of column indexes to a data frame, we can not only reorder columns but also remove them. In the example below, we first re-order columns, moving columns 10 through 17 to the beginning of the data set and moving one through nine to the end. Then, we create a new data set `crime_new` dropping columns 10 and 13 by placing a negative (`-`) before the column indexes.

```
#Re-order
crimes <- crimes[, c(10:17, 1:9)]

#Drop columns
crime_new <- crimes[, -c(10, 13)]
```

It is more likely that we would organize our data calling specific columns by name. The standard approach is to pass a string vector of column names, requiring each name to be placed in quotations. A more concise approach using the `select` function from `dplyr` lists the names of desired columns without quotations. Below, we illustrate how to create a data frame of geographic coordinates.

```
#base R
crimes_geo <- crimes[, c("latitude", "longitude")]

#dplyr approach
crimes_geo <- dplyr::select(crimes, latitude, longitude)
```

To *sort a data set* can be done using `order`, which returns a vector of row indices in ascending order based on the variable of interest. To sort descending, add we can either `"-"` before the variable if it is a numeric or specify `decreasing = T`. The `order` function should not be confused with its cousin the `sort` function, which is designed to order and return a vector of values.

```
#base R - sort ascending
crimes[order(crimes$latitude, crimes$longitude), ]

#base R - sort descending (2 options)
crimes[order(crimes$latitude, crimes$longitude, decreasing = TRUE), ]
crimes[order(-crimes$latitude, -crimes$longitude), ]
```

The base R options are somewhat clunky, requiring repetition of the data frame name if many variables are to be sorted. The `arrange` function in `dplyr` is more concise. Similar to `order`, the `arrange` function sorts ascending by default.

```
#dplyr
crimes <- arrange(crimes, latitude, longitude)
crimes <- arrange(crimes, desc(latitude), desc(longitude))
```

## 4.5 Aggregating data

Seeing the pattern in the data can be challenging when the units are too granular. Aggregating data into coarser units is a common strategy to see the forest rather than the individual tree. For example, it is only possible to determine if crime has become more prevalent by aggregating events.

Here, we introduce the `summarise` and `group_by` functions from `dplyr`. The former `summarise` is used to calculate summary statistics on a set of variables, such as the count of crimes or the average latitude of crimes. The latter `group_by` indicates that manipulations will be applied to each level of a variable. By applying these functions, we can convert a highly granular data set into a more compacted summarized set.

Suppose we would like to see the total number of each type of crime for each ward of Chicago. We can easily send the `crimes` data set into a pipeline like the one below. A `group_by` is specified up immediately after the data set, indicating that all subsequent calculations will be applied groupwise for each type of crime (`primary.type`) and ward of the city (`ward`). The data is then summarized by calculating two variables: the number of crimes `num.crimes` using the count function `n` and the `arrest.rate` by taking the `mean` of the `arrest` variable. The result of this process is a new data frame with only  $n = 151$  and  $k = 4$ . It is clear how this can be used – turning raw data into aggregate, communicable insights.

```
crimes %>%
  group_by(primary.type, ward) %>%
  summarise(num.crimes = n(),
            arrest.rate = 100 * mean(arrest))
```

Any variable can be used in the `group_by` statement, separating each dimension by a comma. The `summarise` function is quite flexible encompassing a large number of summary statistics:

- *Central tendency*: `mean`, `median`
- *Dispersion*: `sd`, `IQR`, `mad`
- *Range*: `min`, `max`, `quantile`
- *Position*: `first`, `last`, `nth`,
- *Count*: `n`, `n_distinct`
- *Logical*: `any`, `all`

**Constructing Time Series.** Using these functions, we can aggregate the events data into a *time series* to uncover patterns and trends over time. Are crimes increasing? Do they follow a predictable seasonal pattern? Suppose we would like to visualize arrest rates in Chicago. We will follow a simple three discrete step:

- *Create quarterly time variables* using a combination of `mdy_hms` and `quarter` functions from `lubridate`. The resulting time variable will be used in `group_by`.
- *Summarize the data* into a new data frame.
- *Plot the new time series*.

This process, as illustrated below, creates a 20-observation data frame `crimes_qtr` that shows that arrest rates have fallen over time. Tasks such as these are core processing steps that underlie most data science projects.

```
#Convert to date form and extract the the quarter
crimes$date1 <- mdy_hms(crimes$date)
```

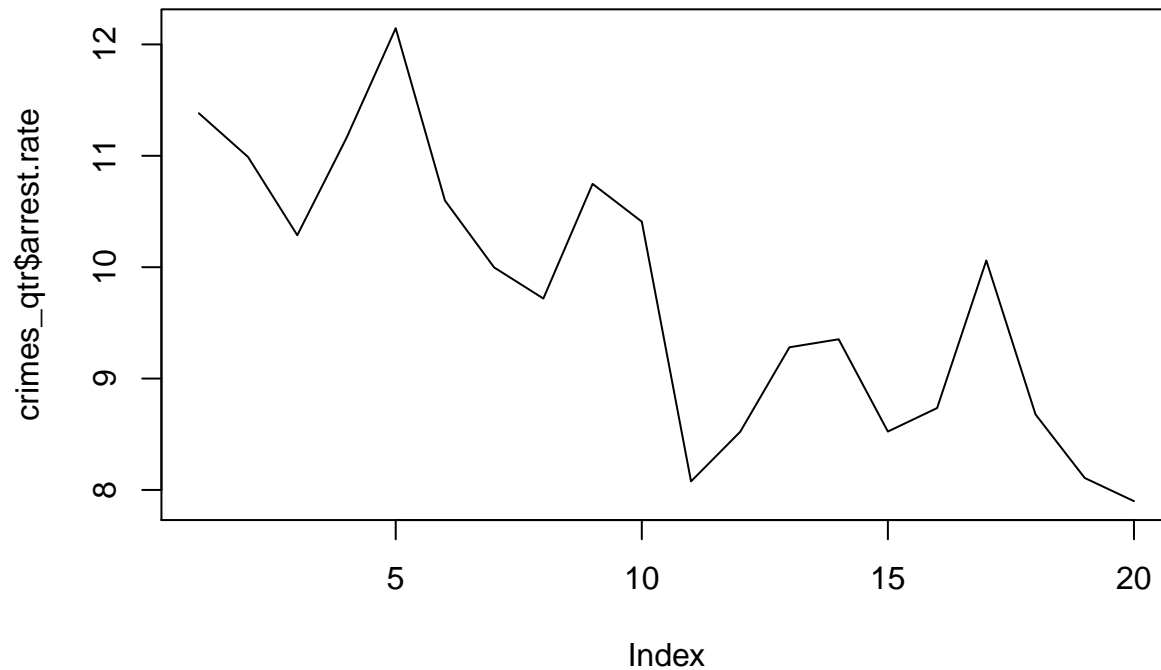
```

crimes$qtr <- quarter(crimes$date1, with_year = TRUE)

#Summarized data
crimes_qtr <- crimes %>%
  group_by(qtr) %>%
  summarise(arrest.rate = 100 * mean(arrest))

#Plot a line chart
plot(crimes_qtr$arrest.rate, type = "l")

```



## Exercises

1. Load the `iris` dataset using the `datasets` package. Calculate the average and maximum sepal length for each Iris species. (**Bonus:** Write the commands without error messages.)
2. Use the `summarise()` function to count the number of observations of sepal width for each species.
3. Create a data frame with the 35th observation of sepal width for each species.

## 4.6 Reshape

Data are generally stored in either *long* form and *wide* form. Long form represents each subject (or individual unit of analysis) at one point in time as a row and each variable as a column. This format also is the most analytically convenient format. In the social sciences, long format is analogous to panel data or stacked data. For example, the following is a long format data set with arrest rates associated with burglary, robberies and thefts for each ward and quarter.

Table 12: Example of long data.

qtr	ward	arrest.rate
2014.1	1	7.17
2014.1	13	24.51
2014.1	25	11.49
2014.1	37	30.30
2014.1	49	7.59

Wide format, in contrast, represents each subject as one row and each column is a variable. If there are multiple dimensions such as time, a variable is spread across multiple time columns. This format is more common when presenting data as it is easier to visually inspect data in a single viewport or sheet of paper. The long format crime data above can be restructured into wide form, as shown below.

The shape of data has a direct implications for how analyses can be conducted and depend on the use case. Panel analyses, or analyses using data for each subject over time, require data in long form. Cross sectional analyses may require data to be in wide form. The process of converting data from one form to another is known as *reshaping*. For those who have familiarity with spreadsheet software, reshaping data essentially is the same as working with a pivot table.

Table 13: Example of wide data.

qtr	ward. 1	ward. 10	ward. 11	ward. 12
2014.1	7.17	14.96	11.50	12.22
2014.2	7.42	15.23	9.12	14.19
2014.3	6.69	9.41	9.76	13.25
2014.4	7.56	13.33	16.99	12.12
2015.1	10.88	21.84	10.44	13.82

To reshape data, we rely on a package called **reshape2**, which simplifies the process of reshaping into two functions: **melt** and **dcast**. The former **melt** is used to convert data into long form, whereas the latter **dcast** converts data into wide form.

```
pacman::p_load(reshape2)
```

**Long to Wide.** Suppose we summarized the crimes data into a quarterly data set of arrest rates by type of crime and ward. This data would be returned in long form.

```
crimes_long <- crimes %>%
  group_by(qtr, ward, primary.type) %>%
  summarise(arrest.rate = 100 * mean(arrest),
            num.crimes = n())
```

The **dcast** function in **reshape2** converts data from long form to wide form. The function expects a data frame, a formula, and at least one value variable. Let's walk through the example below that produces a new data frame where each row is a type of crime in a ward and each column is the arrest rate for each time

period.

- The formula dictates the desired structure. In the formula `ward + primary.type ~ period`, for example, variables to the left of the tilde `~` are treated as row identifiers and variables to the right are dimensions that will be spread across the columns. Multiple dimensions can be specified, each separated by an addition symbol `+`.
- The value variable `value.var`, which in this case is the `arrest.rate`, will fill the new data structure.

The resulting data frame has  $n = 150$  rows and  $k = 22$  columns.

```
crimes_wide <- dcast(crimes_long,
                    ward + primary.type ~ qtr,
                    value.var = "arrest.rate")
```

Suppose we would like to reshape both the `arrest.rate` as well as the `num.crimes`. This requires two modifications to the above code:

- First, we will need to install and load the `data.table` package, which converts data frames into a class of objects for memory efficient and speedy data manipulation. From that package, we will need to use `setDT` to convert our data frame into the appropriate format.
- Second, specify a string vector of variable names for the `value.var`.

The resulting data frame has  $n = 150$  rows and  $k = 42$  columns.

```
pacman::p_load(data.table)
crimes_widetoo <- dcast(setDT(crimes_long),
                      ward + primary.type ~ qtr,
                      value.var = c("arrest.rate", "num.crimes"))
```

**Wide to Long.** It is just as simple to roll the data back into long form using the `melt` function. The example below re-converts the `crimes_wide` data frame into long form. Remember that all but two variables contain `arrest.rate` for a each quarter. Moving from wide to long means that all the `arrest.rate` columns will be consolidated into two variables: column headers will be captured as a new time variable and column values will be arrest rates. The function expects:

- a data frame (`crimes_wide`).
- `id.vars`: a vector of variable names or column indexes that indicate row identifiers.
- `variable.name`: a name of a new variable to which the wide column names will be stored.
- `value.var`: a name of a new variable to store the values.

```
crimes_longtoo <- melt(crimes_wide,
                      id.vars = c("ward", "primary.type"),
                      variable.name = "qtr",
                      value.name = "arrest.rate")
```

## 5 Control Structures

Much of data science requires developing specialized code to handle the eccentricities of a dataset. Re-running blocks of code is required, often times on multiple data samples and subpopulations. It's simply not scalable to manually change variables and assumptions of the code everytime.

Variables are typically treated differently based on their quality and characteristics. In order to accomplish analytical and programming tasks, control structures are used to determine how a program will treat a given variable given conditions and parameters. In this section, we will cover two commonly used control structures: if...else statements and for loops.

## 5.1 If and If...Else Statement

If statements evaluate a logical statement, then execute a script based on whether the evaluated statement is true or false. If the statement is TRUE, then the code block is executed.

```
budget <- 450
if(budget > 400){
  #If statement true, run script goes here
  print("You're over budget. Cut back.")
}
```

```
## [1] "You're over budget. Cut back."
```

In cases where there are two or more choices, if...else statements would be appropriate. In addition to the if() statement, an else statement is included to handle cases where the logical statement is FALSE.

```
budget <- 399
if(budget >= 400){
  #If statement true, run script goes here
  print("You're over budget. Cut back.")
} else {
  #else, run script goes here
  print("You're under budget, but watch it.")
}
```

The complexity of these statements can be as simple as if(x > 10){ print("Hello")} more complex trees:

```
age <- 23

if(age <= 12){
  print("kid")
} else if(age >12 && age <20) {
  print("teenager")
} else if(age >=20 && age <65) {
  print("adult")
} else{
  print("senior")
}
```

```
## [1] "adult"
```

## 5.2 For-loops

Loops can be used to run the a given statement of code multiple times for a specified number of times or a list of index value. This is a functionality that is available in most programming languages, but the programming

syntax will be different. Conceptually, for loops can be likened to an assembly line in a car factory. In order to build a car, a series of well-defined, well-timed processes need to be coordinated in a serial fashion. To build 500 cars, the process needs to be executed 500 times. For-loops are essentially the same: Given a well-defined, self-contained process, a process can be iteratively applied to address repetitive tasks.

Let's take the following example. The code block essentially says "print values for the range of 1 through 5", where *i* is an *index value*. When executing the statement, R will push the first value in the sequence of 1:5 into the index (in this case, it's the number 1), then the code block in between the {} (curly brackets) will be executed, treating *i* as if it's the number 1. Upon executing the code without error, R will advance to the next value in the sequence and repeat the process until all values in the sequence have been completed.

```
for(i in 1:5){  
  print(paste0("Car #", i))  
}
```

```
## [1] "Car #1"  
## [1] "Car #2"  
## [1] "Car #3"  
## [1] "Car #4"  
## [1] "Car #5"
```

We can do the same for a vector or list of values. In the example below, the vector `news` contains six terms. Using a for-loop, we can print out each word in the vector.

```
news <- c("The", "Dow", "Is", "Up", "By", "400pts")  
for(i in news){  
  print(i)  
}
```

```
## [1] "The"  
## [1] "Dow"  
## [1] "Is"  
## [1] "Up"  
## [1] "By"  
## [1] "400pts"
```

For-loops has a few qualities that users should be aware. First, what happens within the for-loop is written to the R environment as *global variables*. That means that any object (e.g. calculations, models) that is created in the loop will be accessible in the programming environment even after the loop ends. This may be a good or bad, depending on the use case: Good if one wants to keep copies of the intermediate results of a loop iteration, but bad if the user is not careful to take note of the potential flood of extraneous objects that may effect downstream calculations. Second, one of the most common mistakes when using loops is failing to record the result of the loop. There are functions in R that are designed to log and package results from loops, but in plain vanilla loops, this is not the case.

**A common paradigm** with for-loops is to iteratively execute repetitive tasks. For example, if a calculation needed to be applied to each of one million files and the results need to be logged, then for-loops are a good option. Typically, the paradigm proceeds as follows:

1. Create placeholder object (e.g. a vector, matrix, or data frame);
2. Initialize loop; and
3. Add outputs to placeholder at the end of each loop iteration.

This may be applied in a broad variety of cases such as processes each data set in a repository of many large data sets, calculating complex statistics for various strata and subsets within the data, among others.



Best practices with loops start with initializing new placeholder objects to full length before the loop rather than increasing the object size within the loop<sup>1</sup>. In R, this is particularly important issue for efficient data processing.

In the example below, we would like to calculate the minimum and maximum of each of 1000 randomly generated normal distributions with  $\mu = 1000$  and  $\sigma = 10$ . To do this, a placeholder data frame `x` with three columns (iteration, min and max) is created with  $n = 1000$  rows for each of the random distributions to be generated. Then, we use `Sys.time()` to capture when the loop starts and end – a common practice for optimizing code. The loop is initiated for 1 to 1000 iterations to calculate the minimum and maximum. At the end of each iteration, the min and max results are overwritten to the row that corresponds to the iteration in the placeholder `x`.

```
#Set placeholder data frame with n rows
```

```
n <- 1000
x <- data.frame(iteration = 1:n,
               min = numeric(n),
               max = numeric(n))
```

```
#Loop
```

```
start <- Sys.time()
for(i in 1:n){
  y <- rnorm(10000, 1000, 10)
  x$min[i] <- min(y)
  x$max[i] <- max(y)
}
Sys.time() - start
```

```
## Time difference of 0.9130011 secs
```

The above process required roughly 0.8 seconds to process. *What happens if the placeholder length were not pre-specified?* For the given parameters, the task normally may last between 1.2 and 1.5 seconds. This may not seem to be much time, but at scale with millions if not billions of records and iterations, the time does tend to add up.

```
#Set placeholder data frame without dimensions
```

```
n <- 1000
x <- data.frame()

#Loop
start <- Sys.time()
for(i in 1:n){
  set.seed(i)
  y <- rnorm(10000, 1000, 10)
  x <- rbind(x, cbind(iteration = i,
                     min = min(y),
                     max = max(y)))
}
Sys.time() - start
```

```
## Time difference of 1.045437 secs
```

---

<sup>1</sup>[https://www.r-project.org/doc/Rnews/Rnews\\_2008-1.pdf](https://www.r-project.org/doc/Rnews/Rnews_2008-1.pdf)

### 5.2.1 R-specific: apply

For-loops are common across all languages, but the efficiency of their implementation will vary. As was described in the previous chapter, R is an interpreted language optimized for mathematical and statistical calculation – quite different than other languages. This means that programming in R is most optimal when vectorizing calculation – linear algebra calculations of vectors and matrices using operations such as `+`, `-`, `*`, `%*%`, among others.

In R, the speed of for-loops may be improved using `lapply()` under certain circumstances. `lapply()`, or *list apply* Whereas the intermediate objects in for-loops are global variables, `lapply()` creates temporary *local variables*.

```
#Set n
n <- 1000

#Loop
start <- Sys.time()
x <- lapply(1:n, function(i){
  y <- rnorm(10000, 1000, 10)
  return(cbind(iteration = i,
               min = min(y),
               max = max(y)))
})
x <- do.call(rbind, x)
Sys.time() - start

## Time difference of 0.7808199 secs
```

## 5.3 While

Whereas for loops require a range or list of values through which to iterate, `while()` statements keep iterating until some condition is met. The `while()` statement is formulated as follows:

```
while([condition is true]){

  [execute this statement]

}
```

A simple case may involve drawing a random value  $x$  from a normal distribution ( $\mu = 1.0$ ,  $\sigma = 0.5$ ) while  $x$  is greater than 0.01.

```
x <- 1
while(x > 0.01){
  x <- rnorm(1, 1, 0.5)
  print(x)
}
print("done!")
```

## 6 Functions

Data manipulation tasks are often repeated for many different projects and it is not uncommon for two or more scripts to contain the same exact steps, but the code is hardcoded. Same logic and different variables names equates to a significant amount of time spent editing and modifying programs.

Rather than tediously modifying programs, try to write your code once, then never again. Each set of code can serve as re-usable tools that can be re-applied to similar problems, but only if it is standardized with well-laid logic. This is the basis of *user-defined functions*: a coder can define some set of standard required inputs on which a set of steps can be applied to produce a standard output.

A typical function is constructed as follows. Using `function`, a set of input parameters are specified as placeholders for any kind of object. For example, `df1` represents a data frame and `var1` is a variable name in string format. Within the curly brackets, we insert code treating the parameters of actual data. In the example below, we calculate the mean of `var1` in data frame `df1`, then assign it to a `temp.mean`. These calculations are executed in a *local environment*, meaning that any calculations steps within the function are temporary. Thus, `temp.mean` is wiped once the function finishes. The `temp.mean` object can be extracted by passing it to `return`. All of the above steps are assigned to the `meanToo` object that is treated like any other function.

It is good form to include commentary about how to use the function. At a minimum, there should be comments containing what the function is, the arguments, and the output. In the open source tradition, you should be writing the code as if others will read and use it.

```
meanToo <- function(df1, var1, ...){  
  #  
  # Calculate mean of a variable in a data frame  
  #  
  # Args:  
  #   df1 = data frame  
  #   var1 = variable name (character)  
  #  
  # Returns:  
  #   A numeric value  
  
  #Code goes here  
  temp <- mean(df1[[var1]])  
  
  #Return desired result  
  return(temp)  
}
```

To execute the function, we will simply need to call the function with a data frame and a variable name. Basically any script can be genericized into a standardized function.

```
meanToo(data, "x1")
```

### 6.1 DIY: Tracking the supply of online content

Services like Google Trends illustrate the demand for online content, giving a sense of what the public are interested in and what people are reacting to at a given moment. But how about the supply of online content?

The supply provides a sense if content producers feel it is worth spending any time on creating new materials.

The opioid epidemic has been a social issue that has long been brewing, but has only become the centered of public attention in recent years. As the crisis deepens, we'd expect more content to be generated. For each year, we could tediously use the web browser to manually copy the number of search results from a platform like Bing, but what if we need to track a large number of search terms. We can package searches on Bing into a simple function that serves as a wrapper that extracts the number of search results. For example, the search query below yields approximately 18,000 search results for the year 2010 (the number may change as this is only an estimate):

```
https://www.bing.com/search?q=opioid%20epidemic&filters=ex1%3a%22ez5_12810_13893%22
```

Breaking down the URL, we see that the substring `q=opioid%20epidemic` indicates that the search `q=` will be based on the terms that follow `q=` and spaces are encoded as `%20`. Next, the substring `filters=ex1%3a%22ez5_12810_13893%22` indicates a time range filter is applied, specifically `12810_13893` are indexes that represent the range 1/27/2005 to 1/15/2008.

To make this a seamless process, we should construct a function `bingCounts` that will retrieve the number of search results for a given `search.term` in a specified calendar year (`year.filter`). First, we will need to load packages using `require`, which is specifically designed for within-function package loading. In particular, we will use:

- To convert the year into indexes, we use the `lubridate` package to work with date objects.
- To efficiently construct a search URL, we rely on the `stringr` package.
- `RCurl` is used to make use of the Client URL package to make a request to the web for a webpage.
- As the webpage is in HTML format, we use the `XML` package to parse the information.

```
bingCounts <- function(search.term, year.filter){  
  #  
  # Retrieve number of search results for exact query  
  #  
  # Args:  
  #   search.term = search query (character)  
  #   year.term = search year.term (int)  
  #  
  # Returns:  
  #   A numeric value  
  
  #Load package  
  require(lubridate)  
  require(stringr)  
  require(RCurl)  
  require(XML)  
  
  #Get date indexes  
  origin <- as_date("1970-01-01")  
  start.index <- as_date(paste0(year.filter, "01-01")) - origin  
  end.index <- as_date(paste0(year.filter, "12-31")) - origin  
  
  #Construct search URL by replacing placeholder terms  
  search.url <- "https://www.bing.com/search?q=term&filters=ex1%3a%22ez5_tstart_tend%22"  
  search.url <- str_replace(search.url, "term", gsub(" ", "%20", search.term))  
  search.url <- str_replace_all(search.url, "tstart", as.character(start.index))  
  search.url <- str_replace_all(search.url, "tend", as.character(end.index))  
  
  #Get URL content as HTML
```

```

search.html <- getURL(search.url)

#Parse HTML
parse.search <- htmlTreeParse(search.html, useInternalNodes = TRUE)

#Extract result statistics
nodes <- getNodeSet(parse.search, "//*[@id='b_tween']/span[1]")
value <- strsplit(xmlValue(nodes[[1]]), " ", fixed = TRUE)[[1]][1]
return(as.numeric(gsub(",", "", value, fixed = TRUE)))
}

```

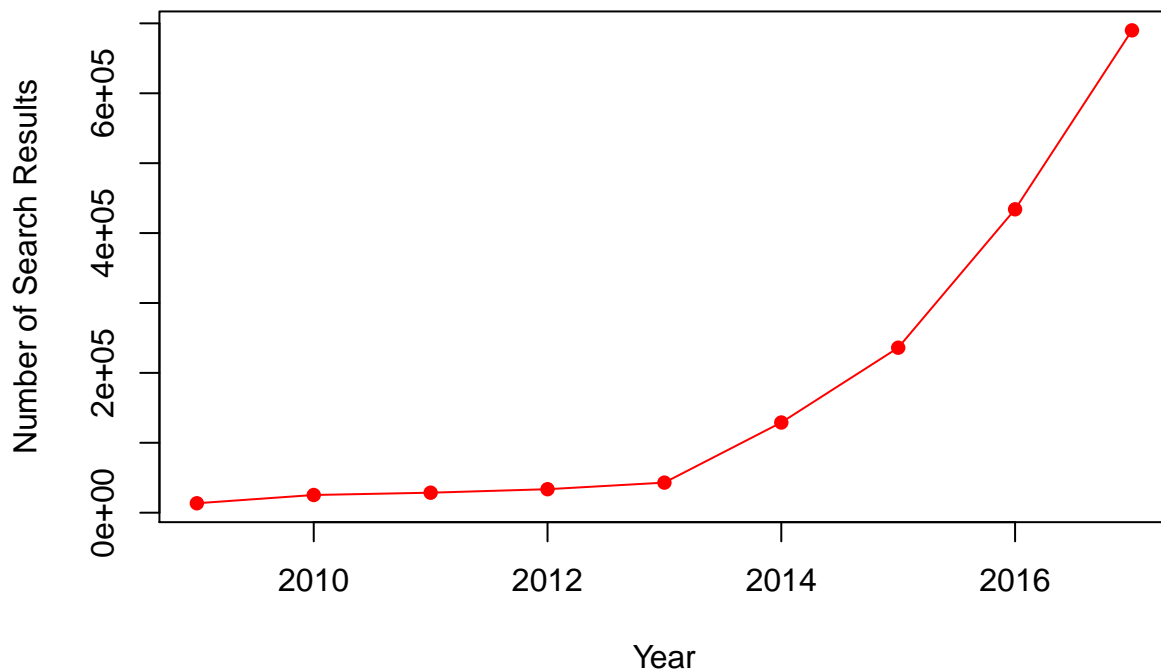
Rather than copying the above function multiple times, the resulting function can then be neatly wrapped into a loop to extract the number of search results for each year in a nine-year period between 2009 and 2017. We find the number of search results has grown 44-times over the period. Based on CDC data, the number of opioid-related deaths increased 2.4-times from approximately 20,400 to 49,000.<sup>2</sup> Thus, we can see the search queries give a measure of direction of growth, but are misleading in their magnitude.

```

#Set parameters
term <- "opioid epidemic"
results <- data.frame()

#Loop through
for(i in 2009:2017){
  results <- rbind(results,
                    data.frame(year = i,
                                cnt = bingCounts(term, i)))
}

```



<sup>2</sup><https://www.drugabuse.gov/related-topics/trends-statistics/overdose-death-rates>

## 7 Best Practices

*Writing just any piece of code is easy. The challenging thing is to ensure that it works well and is repeatable.*

To be successful in the wild requires the appropriate mindset and embracing of data etiquette. Data processing is not like a vampiric feeding frenzy of the sort that might be seen on Buffy the Vampire slayer (the TV show): one simply does not arbitrarily chomp at and slice data without any thought as to what comes next. It requires discipline and finesse to acquire a raw piece of information and extract the empirical substance.

Hypothetically, imagine that you have conducted a three-month analytics project, rushing to obtain to a result. If asked to provide backup, documentation, and a process, it may very well require an additional few months to decipher your passionate but hasty explorations. Where did this data come from? Why did I take a dot product in this equation? What happened to the raw file? Who is Bob mentioned in this email? Ultimately, the data scientist is in control and should be empowered to conduct data processing in a structured, scientific manner. *So, what exact does that mean?* Here are five guidelines that make for successful, scalable data science projects:

- *Define a file structure and stick to it.* Create a folder for *scripts* to store all code, another for *raw* data, a folder for *outputs* from data processing, and another for *documentation*. Choose any file structure, but make sure it is consistent and always followed.
- *Think modularly.* One of the greatest human inventions is the assembly line. Each step in the process is discrete, requiring workers to specialize in one set of well-defined tasks. As a worker becomes specialized, production speed increases as workers focus on a repetitive task as opposed to switching between tasks. Also, if any step in the process needs to be improved, revised, or replaced, all other steps in the process are not disrupted – only one step needs to be addressed. This standardization of processes is the only way to scale. Coding is the same. Each group of similar tasks can be rolled up into a function that is comprised of generic actions. For example, a function could transform involve downloading a file, unzipping it, loading into memory, and extract a specific set of fields, then saving out the extracted data.
- *Adopt and stick to a coding style.* Each file, data field, and function should follow a style so that the logic is cohesive, readable and traceable. The Google R Style Guide, for example, sets easy to follow rules. For example, variable names should be labeled as lower case characters with a "." as a separator: `variable.name`. Functions should be named in proper case without spaces: `NewFunction` or `MyFunction`. Regardless of the style guide, file names should be meaningful and the only punctuation used should be "\_" for spaces and "." to accompany the file type: `gdp_2014_10_02.json`. Be vigilant with your code style as it may be the difference between a successful and a failed project.
- *Work forward, not in circles.* Raw meats should never be placed with cooked meats. It's a one way street to keep things clean and sanitary. Data is also a one way street Start from the raw data, make a copy when working with the raw, do thy bidding on the data, then output results or processed data as a new file. Never overwrite the raw file as this is the equivalent to repeatedly re-fridgerating and re-cooking already cooked meat – the result may be less than desirable in the future.
- *Version up.* Each time code is written, modified, amended, or updated, a new version should be saved. Saving can be done in one many ways. If on a local computer, simply adding a version number to the code file name would suffice with an entry to a log file. The log file can be as simple as a text file with time stamps indicating the date fo the edit, the file name and the nature of the edit. Services such as Github are particularly useful for version control and working on code in groups.

Following these basic guidelines is the difference between successful, scalable projects and ones that get stuck in the mud.

## Chapter References

“EU Gdpr.org.” n.d. <https://eugdpr.org>.

Lohr, Steve. 2014. “For Big-Data Scientists, ‘Janitor Work’ Is Key Hurdle to Insights.” *The New York Times*, August. <https://www.nytimes.com/2014/08/18/technology/for-big-data-scientists-hurdle-to-insights-is-janitor-work.html>.

Wickham, Hadley. 2014. “Tidy Data.” *Journal of Statistical Software* 59 (10).