

# **Collecting, Classifying, and Analyzing Textual Data Using R**

Alison Blaine and Markus Wust  
NCSU Libraries

# Welcome!

# Schedule

**10:00-11:10**

**11:10-11:20**

**11:20-12:30**

**12:30-1:45**

**1:45-2:55**

**2:55-3:05**

**3:05-4:15**

**VM installation help**

**class  
break\***

**lunch**

**\*Please take additional breaks as needed**

# **Restrooms**

**single-stall: 1st floor, near elevators**

**multi-stall: this floor, near elevators**

**lactation room: 1st floor, rm 1202**

**All public restrooms throughout the building are ADA-accessible.**

# **Food & Coffee**

**Cafe (1st floor)**

**Dining Hall (On the Oval)**

**Food trucks**

**Thursday: Baguetteaboutit**

**Friday: American Meltdown**

**main campus drive**

# **Learning environment**

**Questions are welcomed**

**Requests to repeat/re-explain are  
welcomed**

**Collaborators in learning - help  
your neighbor if you can**

**We get imposter syndrome too**

# Agenda

## Thursday

Introduction

VM Installation

R Basics

Web scraping

Data parsing

Introduction to APIs

API data harvesting activity

## Friday

Intro to Textual Analysis

Textual Analysis Activities

Sentiment analysis

Word Frequencies

TF-IDF

Bigrams

Topic Modeling

# Don't Panic!

# Installation stuff

# Goals

**Learn how to harvest text  
data using R from the  
internet**

**Do some web scraping**

**Get some data into R using  
an Application Programming  
Interface (API)**

**Use an R library to harvest  
data**

**Learn & implement some R  
functions for preprocessing  
text data**

**Learn and implement some  
text analysis techniques  
using R**

**Know how to seek help &  
resources in the future**

**Tip: Getting code working is the first battle, understanding fully what it means comes later**

“Good bye”

Etymology: late 16th century: contraction of *God be with you!*

# Course Resources

**go.ncsu.edu/textr**

# Collecting Text Data



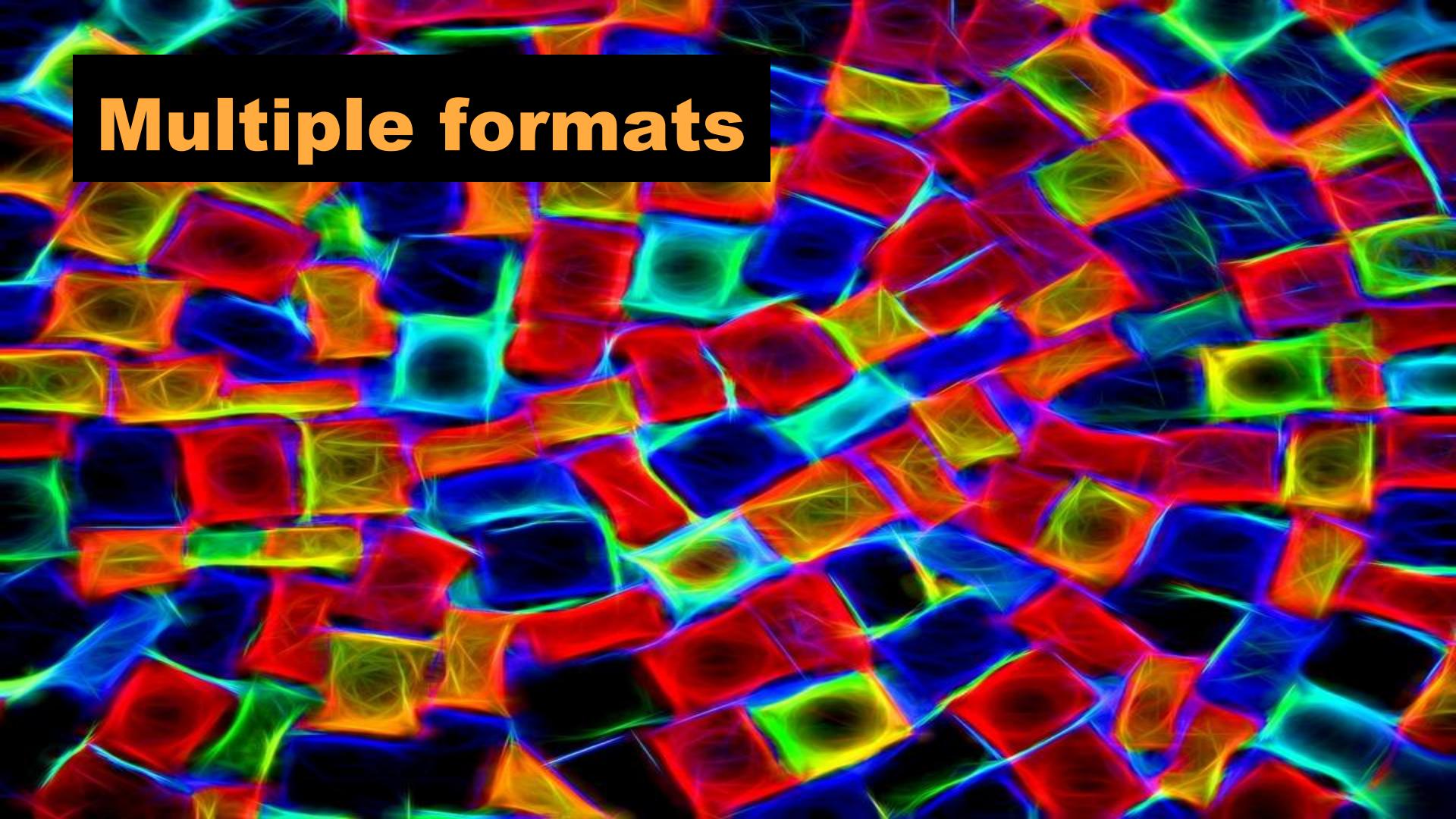
# Data on the Web



**Messy**



# Multiple formats



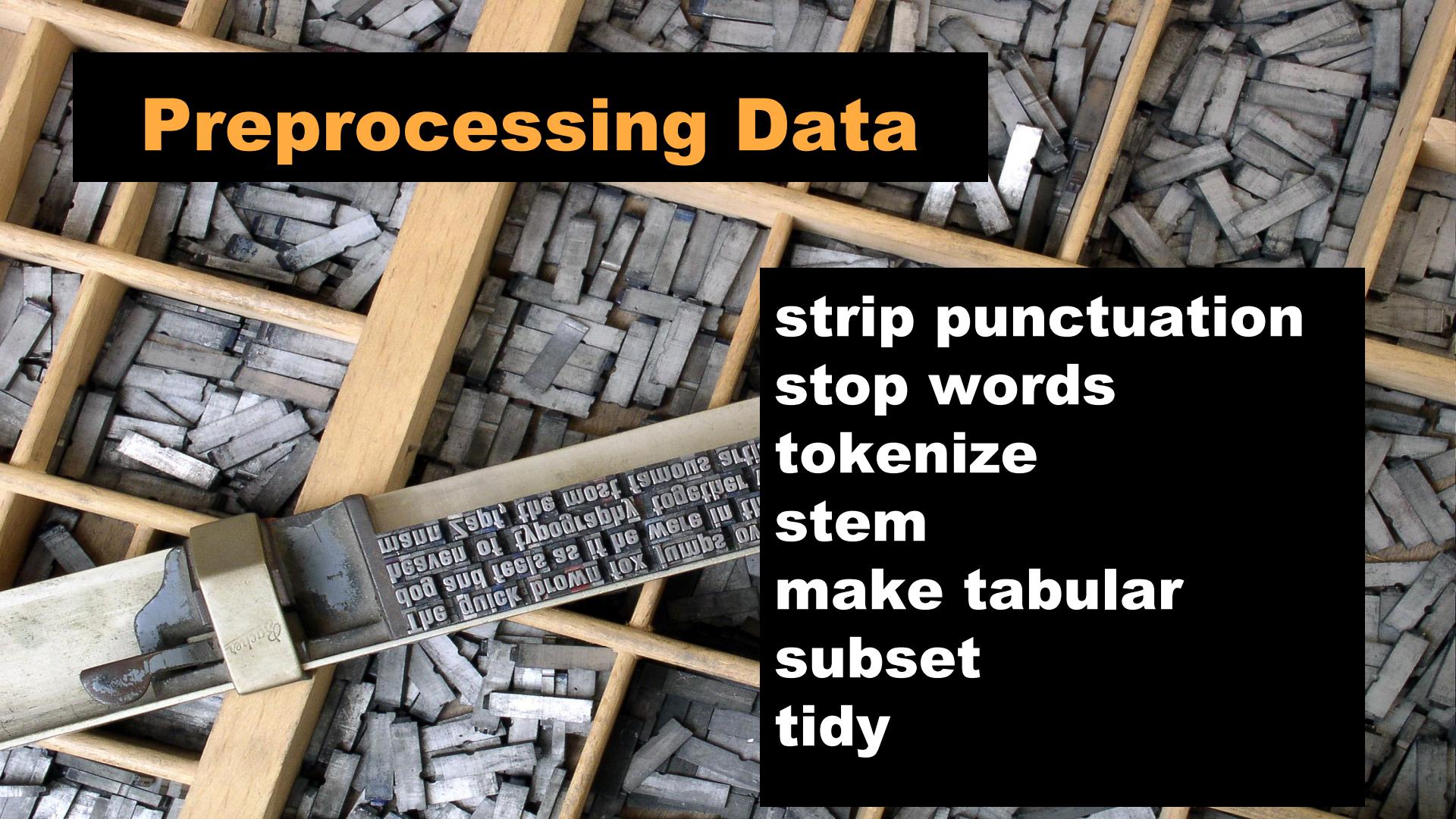
# Dynamic



# Unstructured



# Preprocessing Data

A close-up photograph of a wooden tray filled with metal letterpress type blocks. The blocks are rectangular with various letters and numbers stamped on them. A metal ruler is placed diagonally across the tray, showing markings from 1 to 6. The background consists of more trays of type.

**strip punctuation  
stop words  
tokenize  
stem  
make tabular  
subset  
tidy**

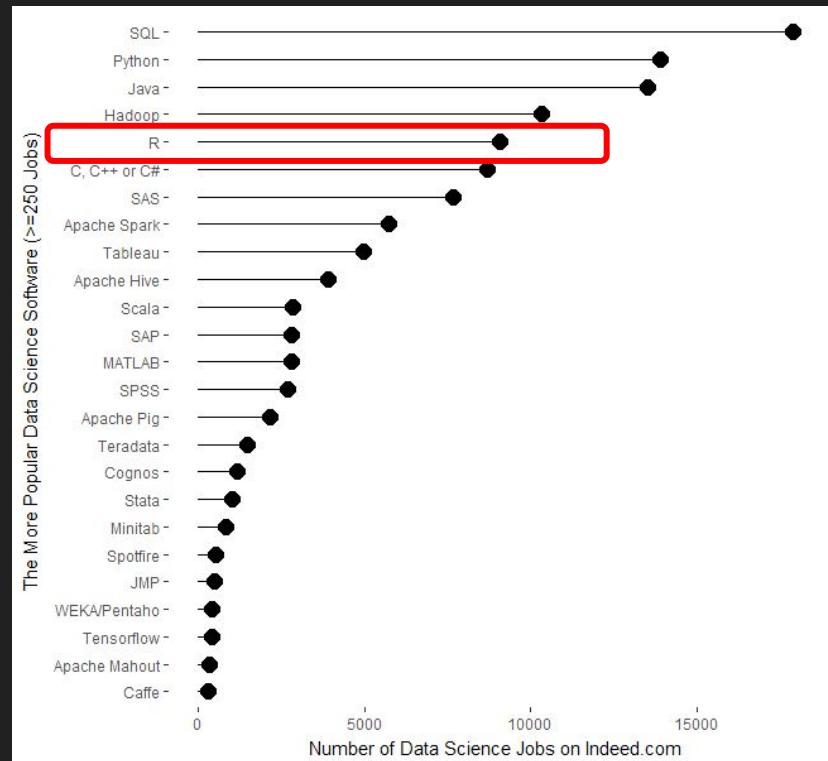
# Analyzing Data

**word frequencies  
topic modeling  
n-grams (bi-, tri-)  
part of speech tagging  
sentiment analysis**

# Popularity of R



[source](#)



# R Packages for Text Analysis

THERE ARE  
MORE THAN  
THIS, BUT  
THIS IS A  
PRETTY  
GOOD LIST!

## CRAN packages:

- [boilerpipeR](#)
- [corpora](#)
- [gsubfn](#)
- [gutenbergr](#)
- [hunspell](#)
- [kernlab](#)
- [KoNLP](#)
- [koRpus](#)
- [languageR](#)
- [lda](#)
- [lsa](#)
- [maxent](#)
- [monkeylearn](#)
- [movMF](#)
- [mscstexta4r](#)
- [mscswebm4r](#)
- [openNLP](#)
- [phonics](#)
- [qdap](#)
- [quanteda](#)
- [RcmdrPlugin.temis](#)
- [RKEA](#)
- [RTextTools](#)
- [RWeka](#)
- [skmeans](#)
- [SnowballC](#)
- [stringi](#)
- [tau](#)
- [tesseract](#)
- [text2vec](#)
- [textcat](#)
- [textir](#)
- [textruse](#)
- [tidytext](#)
- [tm \(core\)](#)
- [tm.plugin.alceste](#)
- [tm.plugin.dc](#)
- [tm.plugin.europresse](#)
- [tm.plugin.factiva](#)
- [tm.plugin.lexisnexis](#)
- [tm.plugin.mail](#)
- [tm.plugin.webmining](#)
- [tokenizers](#)
- [topicmodels](#)
- [wordcloud](#)
- [wordnet](#)
- [zipfR](#)

# R basics

# R-objects

# Vector

*ALL ELEMENTS  
ARE THE SAME  
DATA TYPE!*



[“raspberry”, “blackberry”,  
“strawberry”, “cherry”]

# List



[“raspberry”, “blackberry”,  
“strawberry”, 1, 2, 3, 5, 6, 7]

# Matrix

	col1	col2	col3
row1	2	4	3
row2	1	5	7

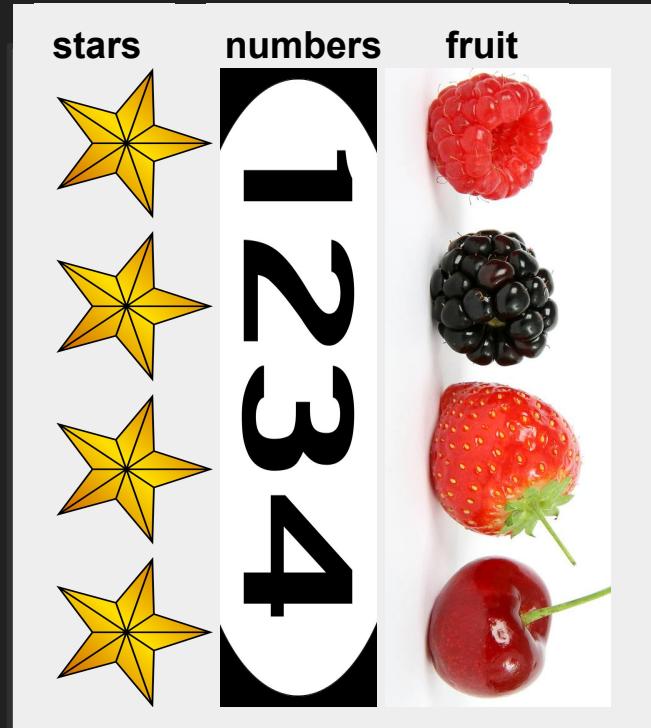
# Array



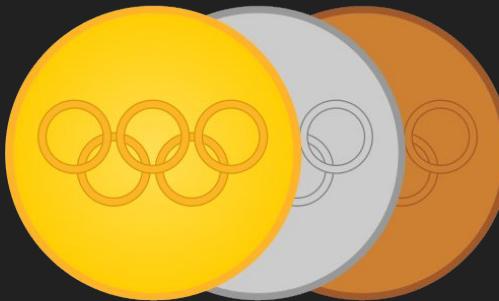
LIKE A  
MATRIX BUT  
CAN HAVE  
MULTIPLE  
DIMENSIONS!  
OR NOT.

# Data Frame

TABULAR. ALL ITEMS MUST  
HAVE THE SAME LENGTH!  
CAN BE VECTORS, ARRAYS  
OR MATRICES



# Factor



```
> x <- c(1, 1, 1, 1, 1, 2, 2, 2, 2, 9, 9, 9, 9)  
> factor(x)  
[1] 1 1 1 1 1 2 2 2 2 9 9 9 9  
Levels: 1 2 9
```

ORDERED CATEGORIES. FACTORS MAKE  
IT EASIER TO WORK WITH CATEGORICAL  
VARIABLES IN MODELING + GRAPHING!

# R packages

# Installation

ONLY HAVE TO  
INSTALL A PACKAGE  
ONCE!

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

# Loading

HAVE TO LOAD THE  
PACKAGE IN EACH R  
SESSION!

```
library(ggplot2)
```

# **Popular packages for working with data**

# dplyr

*dee-PLY-ur*

**select()** - select data

**filter()** - filter data

**arrange()** - order data

**mutate()** - create variables

**count()** - count

**group\_by()** - group

# **tidyr**

*TY-dee-ur*

**gather() - columns into rows**

**spread() - rows into columns**

**separate() - one column into  
several**

**unite() - several columns into  
one**

# Tidy data

country	year	cases	population
Afghanistan	1999	745	1637071
Afghanistan	2000	2666	20695360
Brazil	1999	37737	17206362
Brazil	2000	84488	174504898
China	1999	212258	1272915272
China	2000	216766	128042583

variables

country	year	cases	population
Afghanistan	1999	745	1637071
Afghanistan	2000	2666	20695360
Brazil	1999	37737	17206362
Brazil	2000	84488	174504898
China	1999	212258	1272915272
China	2000	216766	128042583

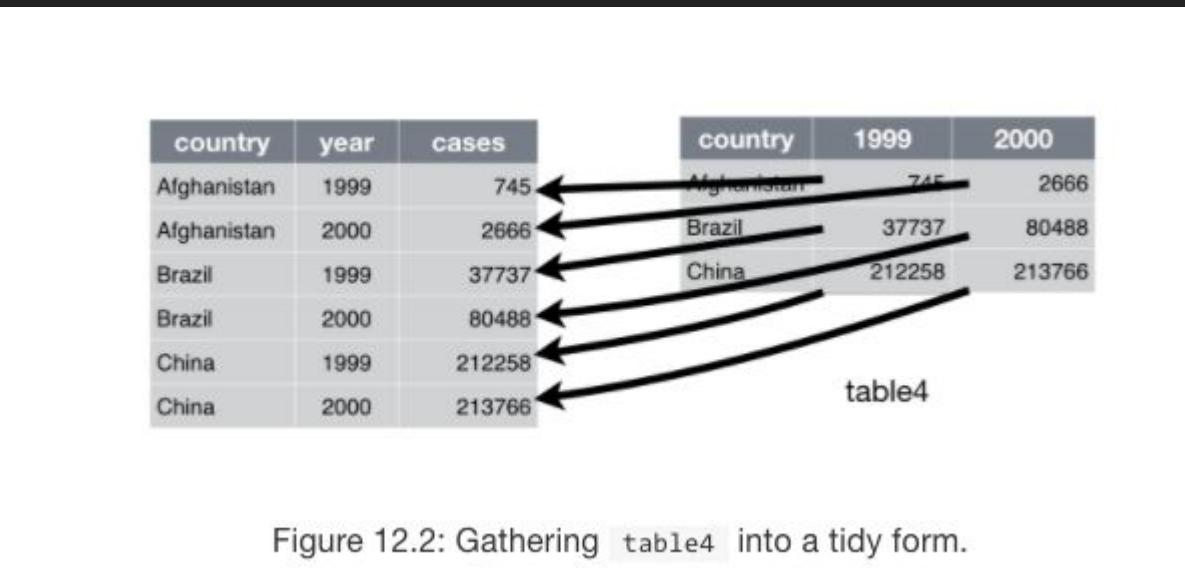
observations

country	year	cases	population
Afghanistan	99	745	1637071
Afghanistan	00	2666	20695360
Brazil	99	37737	17206362
Brazil	00	84488	174504898
China	99	212258	1272915272
China	00	216766	128042583

values

# tidyverse

## gather() - columns into rows



# spread() - rows into columns

tidyr

country	year	key	value	country	year	cases	population
Afghanistan	1999	cases	745	Afghanistan	1999	745	19987071
Afghanistan	1999	population	19987071	Afghanistan	2000	2666	20595360
Afghanistan	2000	cases	2666	Brazil	1999	37737	172006362
Afghanistan	2000	population	20595360	Brazil	2000	80488	174504898
Brazil	1999	cases	37737	China	1999	212258	1272915272
Brazil	1999	population	172006362	China	2000	213766	1280428583
Brazil	2000	cases	80488				
Brazil	2000	population	174504898				
China	1999	cases	212258				
China	1999	population	1272915272				
China	2000	cases	213766				
China	2000	population	1280428583				

table2

# magrittr

*ma-GRIT-ur*

`%>%`

**pipe operator**

**you can put one variable  
through multiple  
functions**

# Function for washing a car using **magrittr**

```
dirty_car %>%  
  soap() %>%  
  rinse() %>%  
  dry()
```



# Example using **dplyr** and **magrittr**

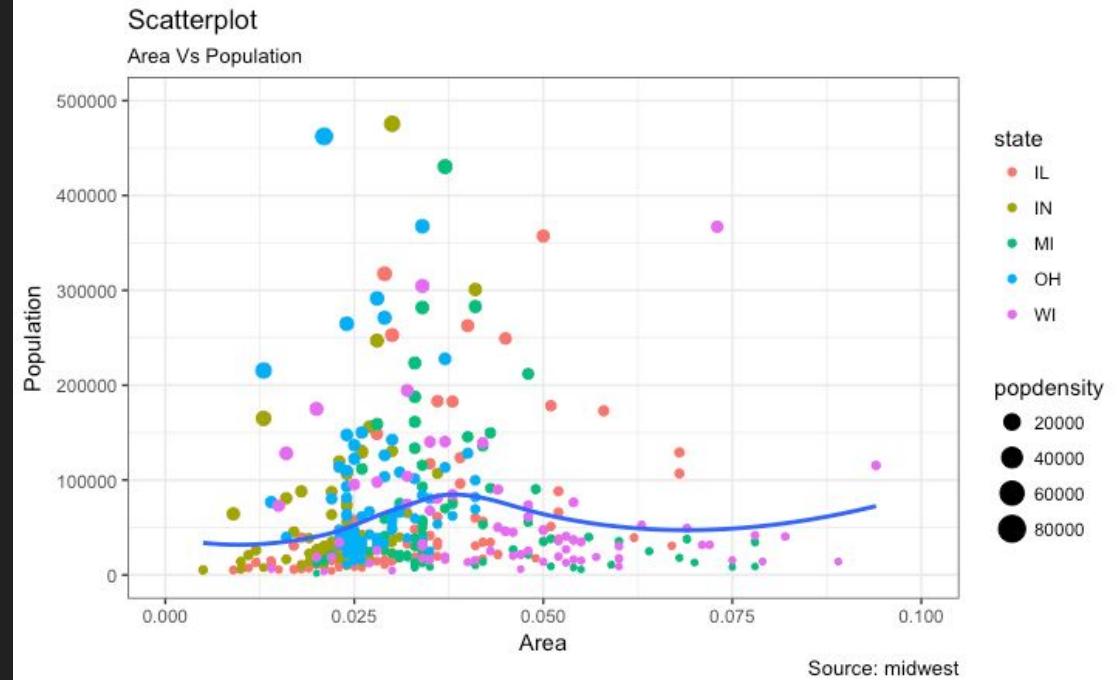
“books” is a data frame with the following information:

title	author	edition	date_published	type
A Day in the Life of Squid	O. Isles	first	10/4/15	not banned
Recipes for Disaster	Y. Not	second	1/28/17	banned

```
s  
banned_books <- books %>%  
  filter(type == "banned")
```

**What does `banned_books` contain?**

# ggplot2



# The tidyverse

## Components



Cheat sheets

data wrangling

ggplot2

# R Studio

~/machine-learning-test - RStudio

File Edit View Insert Cell Run Source Environment History Global Environment Data Compare Result Emp Data Earnings S Rating InTrain Test Testing Training Wage Files Plots Packages Help Viewer Zoom Export Publish

Session info

rf\_model

plots, help, file menu

education education age age jobclass %IncMSE IncNodePurity

output

```
1 install.packages("ISLR")
2 install.packages("randomForest")
3 library(ISLR)
4 library("randomForest")
5 library(ggplot2)
6 library(caret)
7
8 # Know more about the Wage dataset
9 data(Wage)
10 summary(Wage)
11 dim(Wage)
12
13
14 # Some feature engineering, we don't need logwage
15 Wage<- subset(Wage, select=- c(logwage))
16
17 # Some Exploratory data analysis a.k.a. visualization
18 qplot(age, wage, data=Wage, colour = race)
19 qplot(age, wage, data=Wage, colour=education)
20
21 Wage_x <- data.frame(as.numeric(Wage$age), as.numeric(Wage$education), as.numeric(Wage$jobclass),as.numeric(Wage$wage))
22 cor(Wage_x, Wage$wage)
23
24
25 inTrain<- createDataPartition(y = Wage$wage, p=0.7, list= FALSE)
26
27 training<- Wage[inTrain, ]
```

# Web scraping

# Web Scraping

- Accessing one or more webpages and extracting data
- What you do if there is no API (Application Programming Interface) to get data
- For this to work, webpages should use consistent HTML markup



# Web Scraping: How Does It Work?

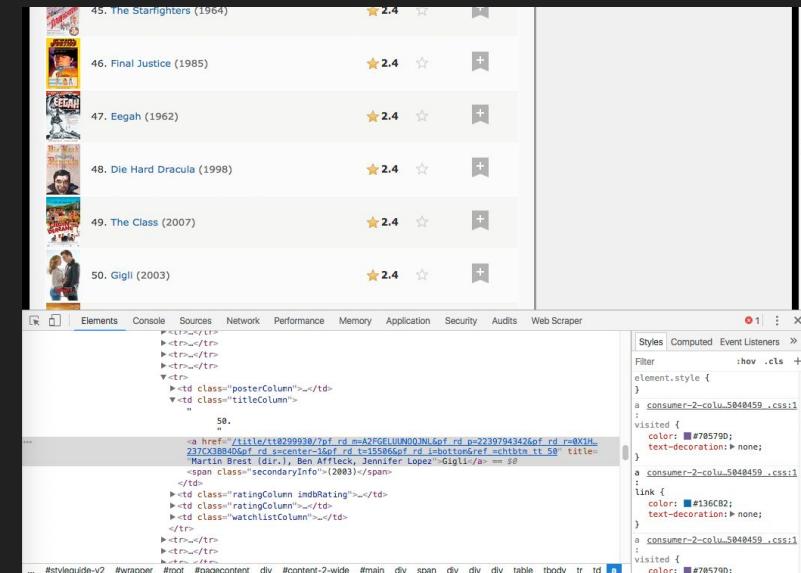
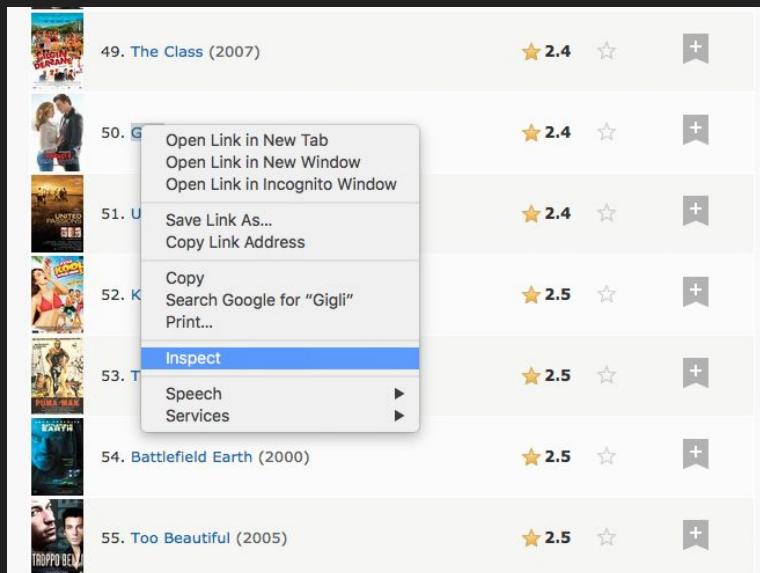
Webpages consist of content (text, images, etc.) and HTML markup:

```
<html>
<body>
<div id="content">
<p id="descripton">Lorem ipsum...</p>
<p id="summary">Lorem ipsum...</p>
</div>
</body>
</div>
```

# Web Scraping: How Does It Work?

- You have to identify the HTML elements that contain the data you are interested in
- You can select elements by:
  - The elements' names
  - IDs
  - Classes
  - Their parents (the elements that they are nested in)
  - Their siblings (elements that are at the same level)
  - All of the above!

# Web Scraping: How Does It Work?



Note:

- This example uses Google Chrome, but other browsers offer similar functionality out of the box.
- This is different from the browser's "View Page Source" function.

# Web Scraping: How Does It Work?

- R has library (`rvest`) to support web scraping:
  - Loads webpage(s)
  - Selects content based on identifiers you provide
- Sometimes, you cannot clearly identify certain data points based on a webpage's HTML markup
- If you are looking for information that follows a particular format, such as telephone numbers, you can instead use **Regular Expressions**

# Regular Expressions

Regular Expression:

"... a sequence of characters that define a search pattern" (Wikipedia)

- They help you find text that follows a specific pattern:
  - Phone numbers / Social Security Numbers
  - URLs
  - Addresses
- You can extract those data or modify them
- Note: There are different flavors of regular expressions

# Regular Expressions

- All modern programming languages support regular expressions
- Regular expressions can be used in other environments:
  - Google Sheets
  - Many text editors that you would use for programming
- Since we could do a workshop just on regular expressions:
  - Excellent free tutorial:  
<http://www.regular-expressions.info/tutorial.html>
  - Free sandbox for testing: <https://regex101.com>

# Regular Expressions: Examples

- `\w <-` All word characters. That includes all letters, all digits, and the underscore (`_`) character.
- `\d <-` All digits.
- `\s <-` Any whitespace character, i.e., a space, a tab, a line break, or a form feed.
- `\W <-` Anything but a word character.
- `\D <-` Anything but a digit.
- `\S <-` Anything but a whitespace character.
- `\d{1,3} <-` Between 1 and 3 digits.
- `\d* <-` 0 or more digits.

# Regular Expressions

What does this regular expression find:

```
(?:[a-zA-Z!#$%&'*+=?^_`{|}~-]+(?:\.[a-zA-Z!#$%&'*+=?^_`{|}~-]+)*|"(?:[\x01-\x08\x0b\x0c\x0e-\x1f\x21-\x23-\x5b\x5d-\x7f]|\\"[\x01-\x09\x0b\x0c\x0e-\x7f])*")@(?:(:?[a-zA-Z](?:[a-zA-Z-]*[a-zA-Z])?\.)+[a-zA-Z](?:[a-zA-Z-]*[a-zA-Z])?)|(?:(:?25[0-5]|2[0-4][0-9])[01]?[0-9][0-9]?)\.\{3\}(?:25[0-5]|2[0-4][0-9])[01]?[0-9][0-9]?|[a-zA-Z-]*[a-zA-Z](?:(?:[\x01-\x08\x0b\x0c\x0e-\x1f\x21-\x5a\x53-\x7f]|\\"[\x01-\x09\x0b\x0c\x0e-\x7f])+)\])
```

# Regular Expressions

Answer: Email addresses!

```
(?:[a-zA-Z0-9!#$%&'*+/=?^_`{|}~-]+(?:\.[a-zA-Z0-9!#$%&'*+/=?^_`{|}~-]+)*|"(?:[\x01-\x08\x0b\x0c\x0e-\x1f\x21-\x23-\x5b\x5d-\x7f]|\\"[\x01-\x09\x0b\x0c\x0e-\x7f])*")@(?:(:[a-zA-Z](?:[a-zA-Z-]*[a-zA-Z])?\.\.)+[a-zA-Z](?:[a-zA-Z-]*[a-zA-Z])?|\[(?:(:25[0-5]|2[0-4][0-9])[01]?[0-9][0-9]?)\]\{3\}(?:25[0-5]|2[0-4][0-9])[01]?[0-9][0-9]?|[a-zA-Z-]*[a-zA-Z](?:[\x01-\x08\x0b\x0c\x0e-\x1f\x21-\x5a\x53-\x7f]|\\"[\x01-\x09\x0b\x0c\x0e-\x7f])+)\])
```

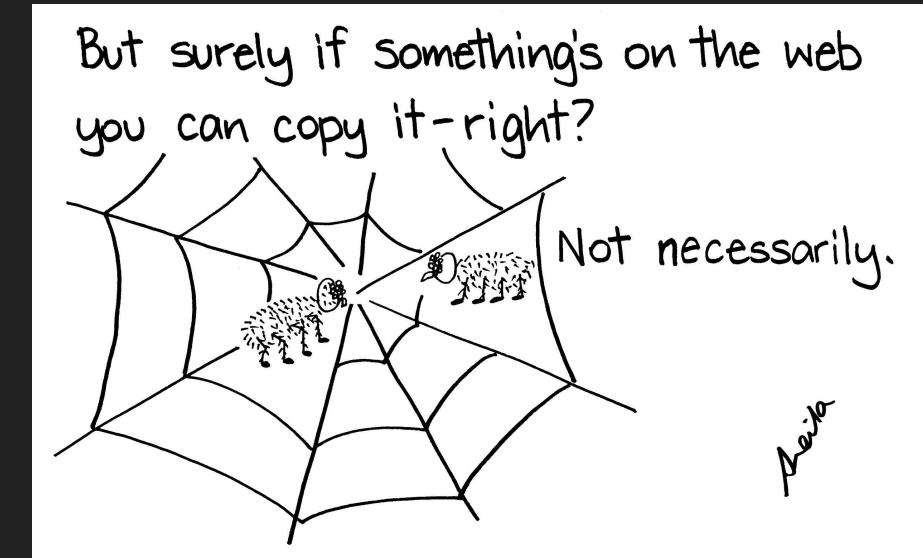
# Web Scraping: Issues

- Although data display is created with HTML markup, data is not really structured
- You can't be 100% certain that elements will always contain the same data
- Scraping becomes more difficult if a page does not use unique IDs for HTML elements



# Web Scraping: Issues

- Web scraping is a quick way to get data
- Problem:
  - Just because it is made public does not mean that you can just take and use it for your own work

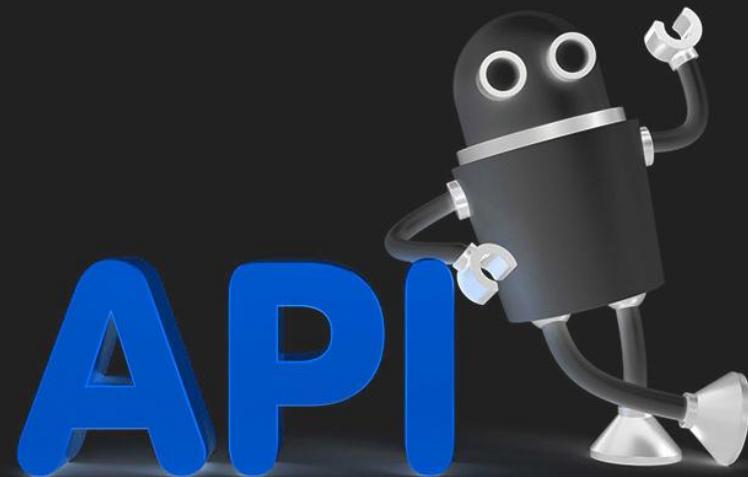


# Web Scraping: Issues

- You should ask for permission if there is no explicit permission for data reuse
- That way, you can also ask whether the data are available in an alternative format, e.g., as a spreadsheet; this makes them easier to work with

# API: Better Than Web Scraping

- A better way to get data is to use an **API** if one is available and accessible:
  - Legal situation is usually clear
  - Reliable data structure
  - Faster
  - Easier to use



# API

API <- Application Programming Interface

Web APIs:

- Allow you to send data requests to many web-based services
- Let you request individual data points or whole data sets
- Return data in a structured format that can easily be processed

The screenshot shows the ProgrammableWeb website interface. At the top, there's a navigation bar with links for 'API NEWS' and 'API DIRECTORY'. Below the navigation is a search bar and social media sharing icons. The main content area features a large banner for the 'SEMANTIC TIME CAPABILITY AWARENESS API' with a sub-headline 'Google Adds Semantic Time Capability to Awareness API'. Below the banner are three smaller news cards: 'Mapbox Launches Optimization API, Mapbox Navigation SDK Updated' (by Janet Wagner), 'Mozilla Announces WebExtensions API for Cross-Browser Extensions' (by Eric Carter), and 'Google Releases Final APIs for Android O' (by Eric Zeman). At the bottom of the page, there's a section titled 'API NEWS' with a card for 'MapQuest Announces Release of MapQuest.js: An Interactive JavaScript Mapping Library'.

# API: Openness

APIs differ in their degree of openness:

- Subscription only
- Open but require API key
- Open without API key but rate-limited
  - How many requests can you send per hour?
- Open without limits

# API: How To Get Data

- To get data through an API you first have to send a request
- Different technical approaches
- Some require you to first authenticate
- Some ask you to turn your request into a small data package and send that to the service
- Easiest way:
  - Construct a request URL that includes all the required information
  - Have your script send a GET request to that URL
  - API sends back data

# API: How To Get Data

Example: arxiv.org (preprint archive [Physics, Math, Comp Sci, Quantitative Biology/Finance, Statistics])

`http://export.arxiv.org/api/query?search_query=all:electron&start=0&max_results=20`

Example: Wikipedia

`https://en.wikipedia.org/w/api.php?action=opensearch&search=electron&namespace=0&format=json`

To see what the data returned by these services looks like open the URLs in your browser!

# API: Data Formats

- APIs also differ in the data formats they provide
- Some services only offer one format, others let you choose between formats
- The most common formats are:
  - XML <- eXtensible Markup Language
  - JSON <- JavaScript Object Notation

# API: Data Formats

XML

```
<?xml version="1.0"?>  
<book>  
  <id>12345</id>  
  <title>Beginning R</title>  
  <author>Larry Pace</author>  
  <published>  
    <publisher>Apress</publisher>  
    <year>2012</year>  
  </published>  
</book>
```

JSON

```
{  
  "book": {  
    "id": "12345",  
    "title": "Beginning R",  
    "author": "Larry Pace",  
    "published": {  
      "publisher": "Apress",  
      "year": "2012"  
    }  
  }  
}
```

# API

Once the API response is received, you can then parse it and directly address the data elements you are looking for:

```
temp <- data["properties"]["periods"]["temperature"]
```

How to find correct elements:

- Look at data returned by API
- Read API documentation

```
{
  "@context": [
    "https://raw.githubusercontent.com/geojson/geojson-ld/master-contexts/geojson-ld.jsonld"
  ],
  "wx": "https://api.weather.gov/ontology#",
  "geo": "http://www.opengis.net/ont/geosparql#",
  "unit": "http://codes.wmo.int/common/unit/",
  "@vocab": "https://api.weather.gov/ontology#"
},
{
  "type": "Feature",
  "geometry": {
    "type": "Point",
    "coordinates": [
      -78.6382,
      35.7796
    ]
  },
  "properties": {
    "updated": "2017-07-17T16:30:46+00:00",
    "units": "us",
    "forecastGenerator": "HourlyForecastGenerator",
    "generatedAt": "2017-07-17T20:25:18+00:00",
    "periods": [
      {
        "number": 1,
        "name": "",
        "startTime": "2017-07-17T16:00:00-04:00",
        "endTime": "2017-07-17T17:00:00-04:00",
        "temperature": 84,
        "temperatureUnit": "F",
        "temperatureTrend": null,
        "windSpeed": "5 mph",
        "windDirection": "S",
        "icon": "https://api.weather.gov/icons/land/day/tsra,20?size=small",
        "shortForecast": "Slight Chance Showers And Thunderstorms",
        "detailedForecast": ""
      }
    ]
  }
}
```

# API: Tools

Browser Extension:

- JSON Formatter (Chrome): Make JSON more human-readable

R Libraries

- httr: Handling API requests
- jsonlite: Parsing JSON data, creating R objects

# Text Analysis

**Train computers to find  
similarities between texts,  
extract themes, and analyze  
expressed sentiments**

# **Text Analysis Methods**

# Word frequency

**The mouse ate the cheese.**

**The 2**

**mouse 1**

**ate 1**

**cheese 1**

# **Collocation**

**“love”**

**true love**

**my love**

**with love**

# **Concordance** - key word in context

**‘Will you be my own true love?’ the musician  
sang...’**

**“took the dusty books with love and dusted  
them...”**

# N-grams (bi, tri-, etc)

Example: “The merry month of May”

Where n = 2:

The merry  
merry month  
month of  
of May

THESE ARE  
BIGRAMS!

# **Natural Language Processing**

**Part of speech tagging**

**Topic modeling**

**Sentiment analysis**

**etc**

# Text Analysis: Glossary

**stopwords** - [the, a, an, and, not, for, but, so...]

**tokenize** - the/cat/sat/on/the/mat/

**stem** - playing, played, player → play\*

**lemmatize** – good, better, best → good

trying, tried, tries → try

# **Text Analysis: Glossary**

**tf-idf - how important a term is in distinguishing a text from others in a collection**

**document-term matrix - rows: documents, columns: terms**

**bigram - pairs of contiguous words in a document**

# **Workshop activities**

**sentiment analysis**

**tokenization**

**word frequency**

**document-term matrix**

**tf-idf**

**bigrams**

**visualization**

**part of speech tagging**

**topic modeling**

# **Get the Data**

**Download “Day 2 files” folder**

**Unzip/extract folder**

[go.ncsu.edu/textr](http://go.ncsu.edu/textr)

# Resources

**Julia Silge and David Robinson (2017), Text**

**Mining with R: A Tidy Approach**

**Matthew Jockers (2014), Text Analysis with R for Students of Literature**

**stackoverflow.com**