

# DSBA 5122: Visual Analytics

## Why R/RStudio and tidyverse?

Ryan Wesslen

August 26, 2019

# What is R?

Terminal

```
R version 3.4.4 (2018-03-15) -- "Someone to Lean On"  
Copyright (C) 2018 The R Foundation for Statistical Computing  
Platform: x86_64-pc-linux-gnu (64-bit)
```

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

```
  Natural language support but running in an English locale
```

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

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

```
[Previously saved workspace restored]
```

```
> █
```

# What is RStudio?

The screenshot displays the RStudio interface with the following components:

- Source Editor:** Contains R code for data processing and plotting. The code includes a `count()` function to aggregate data by date and verb, a `mutate()` function to adjust for sampling, and a `ggplot()` function to create a line plot with points.
- Environment Pane:** Shows the current environment with two data objects: `counts` (14 observations, 3 variables) and `protestData` (135847 observations, 11 variables).
- Console:** Shows the execution of the R script, including library loading, data reading, and the execution of the `count()` and `mutate()` functions.

```
9   count(Date = as.Date(postedTime), Type = verb) %>%
10  mutate(n = n * 10) # adjust for sampling
11
12  # plot
13  ggplot(counts, aes(x = Date, y = n, color = Type)) +
14    geom_line() +
15    geom_point() +
16    labs(x = "Date", y = " ",
17         title = "Charlotte Protest Tweets",
18         subtitle = "Adjusted for 10% Sample") +
19    scale_y_continuous(label = scales::comma) +
20    theme(legend.position = c(0.8,0.8))
21
```

```
> # Project 1: URLs
> library(tidyverse)
> # 10% sample
> protestData <- readRDS("Protest.RData")
> # get daily counts
> counts <- protestData %>%
+   count(Date = as.Date(postedTime), Type = verb) %>%
+   mutate(n = n * 10) # adjust for sampling
> |
```

# 1. Why R / RStudio?

It's free, as in ... free beer?

## 2. Why R / RStudio?



**JD Long**

@CMastication



well this R=Batman, Python=Superman apparently showed up in a vendor pitch to one of my colleagues. That must make it official. [#rstats](#) [#python](#)

| Analysis Tool   | Similar Superhero   | Super Powers in Common   |
|---|---|--|
| <br>R      | <br>Batman   | <ul style="list-style-type: none"><li>• Detective Work</li><li>• Intelligence</li><li>• Cunning</li><li>• Usage of Tools</li><li>• More Brain than Muscles</li></ul> |
| <br>Python | <br>Superman | <ul style="list-style-type: none"><li>• Muscle Power</li><li>• Super Strength</li><li>• Elegance</li><li>• Wide Range</li><li>• More Muscles than Brain</li></ul>    |

♥ 415 7:44 PM - Sep 5, 2018



💬 172 people are talking about this





Front end



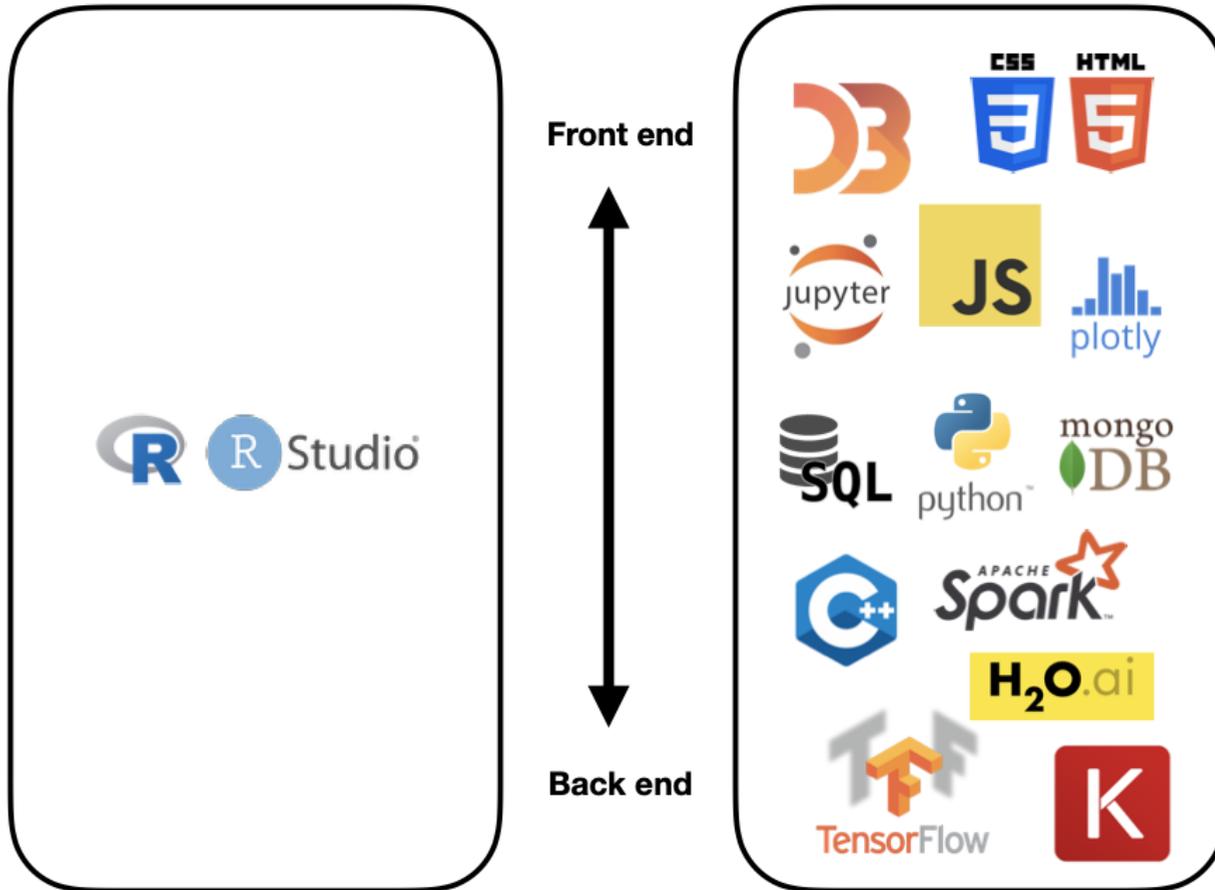


Front end



Back end





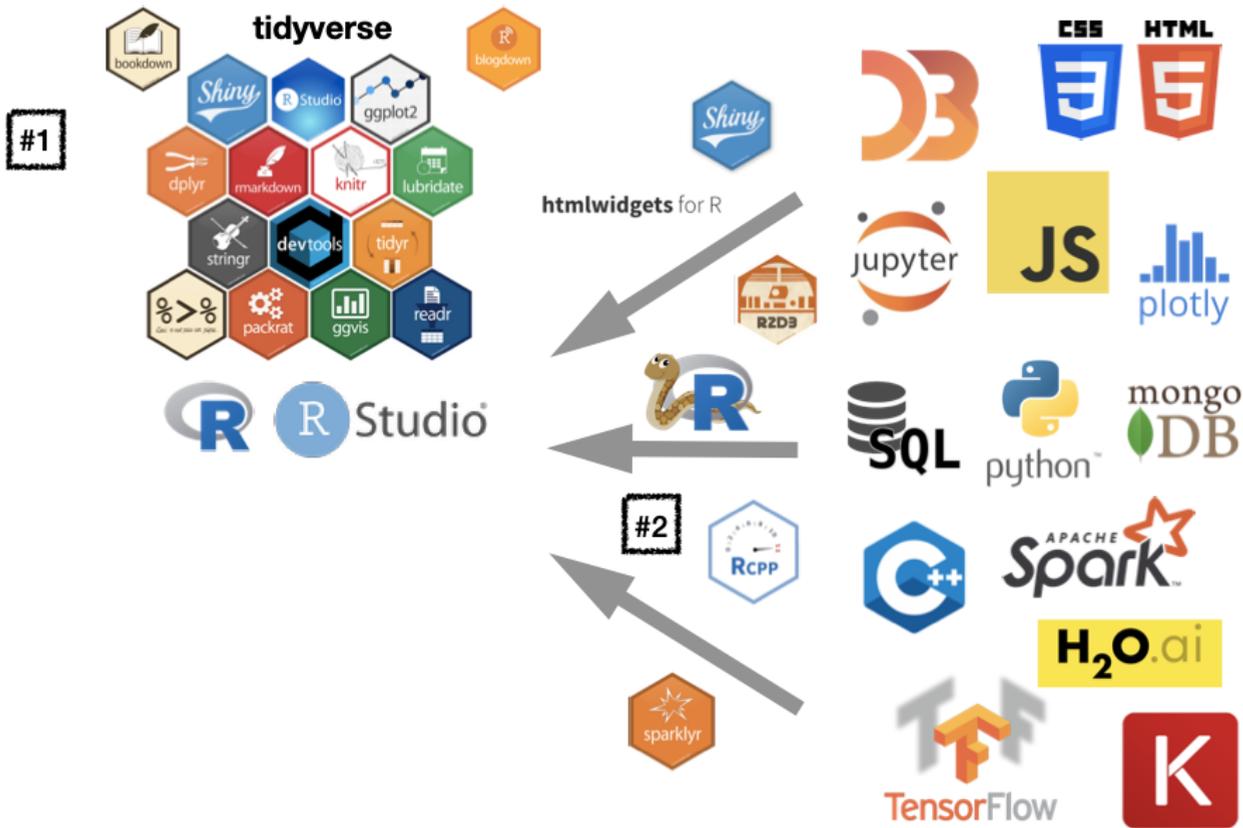
#1

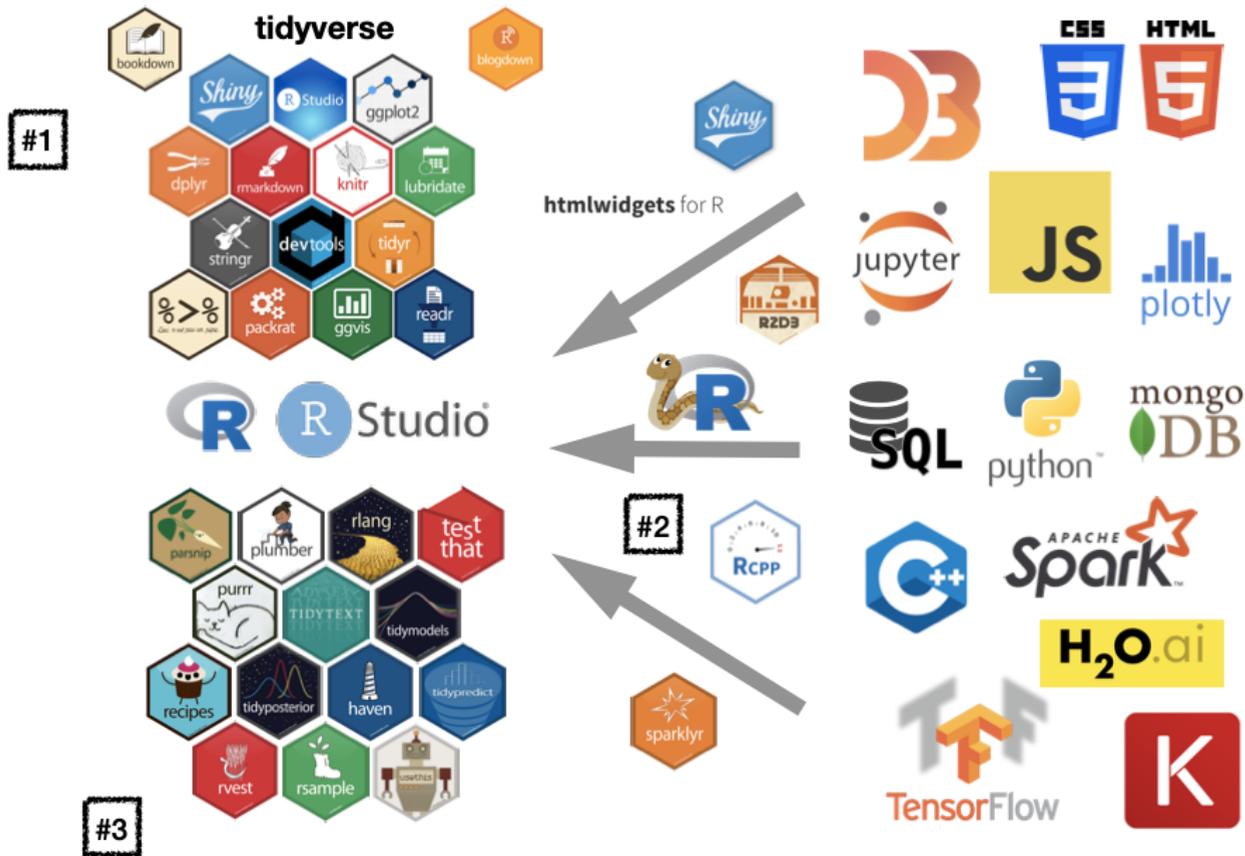


Front end



Back end





### 3. Why R / RStudio?

The `#rstats` community (h/t `@hrbrmstr`)

```
library(rtweet)
library(igraph)
library(ggraph)
library(tidyverse)

rt_g <- search_tweets("#rstats", n=3200) %>% # twitter api
  filter(retweet_count > 0) %>% # keep tweets with RT's
  select(screen_name, mentions_screen_name) %>% # select column from/to
  unnest(mentions_screen_name) %>% # unnest json
  filter(!is.na(mentions_screen_name)) %>% # include mentions
  graph_from_data_frame() # convert to ggraph format
```

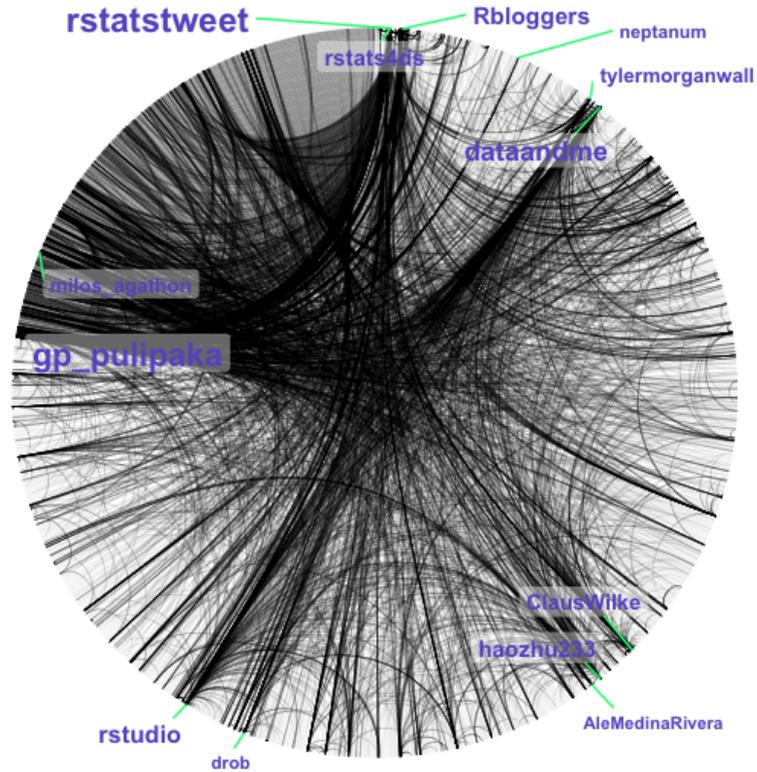
### 3. Why R / RStudio?

The `#rstats` community (h/t `@hrbrmstr`)

```
# ggplot inspired networks: ggraph
ggraph(rt_g, layout = 'linear', circular = TRUE) +
  geom_edge_arc(edge_width=0.125, aes(alpha=..index..)) +
  geom_node_label(aes(label=node_label, size=node_size),
                  label.size=0, fill="#ffffff66", segment.colour="springgreen",
                  color="slateblue", repel=TRUE, fontface="bold") +
  coord_fixed() +
  scale_size_area(trans="sqrt") +
  labs(title="Retweet Relationships",
        subtitle= subt) +
  theme_graph() +
  theme(legend.position="none")
```

## Retweet Relationships

Most retweeted screen names labeled. Darker edges == more retweets. Node size == larger degree



# And it's just fun!

 **Tyler Morgan-Wall**  
@tylermorganwall 

Replying to @ClausWilke  
I can and it's HORRIFYING 🤢🤢🤢  
🤢 #rstats



🍷 183 9:54 PM - Jan 25, 2019 

💬 45 people are talking about this 



**David Schoppik**

@schoppik



Replying to @tylormorganwall



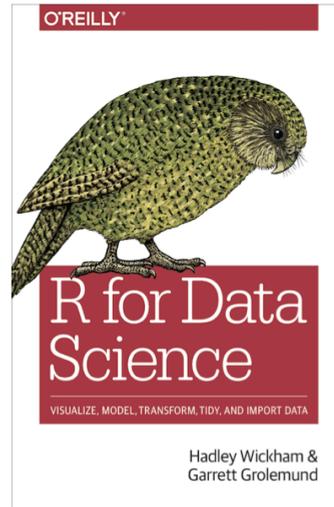
♡ 20 10:14 PM - Jan 25, 2019



 [See David Schoppik's other Tweets](#)



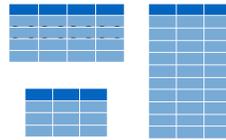
# tidyverse



# What are R packages?



function1()  
function2()  
function3()  
function4()



function5()  
function6()  
function7()  
function8()



function9()  
functionA()  
functionB()  
functionC()

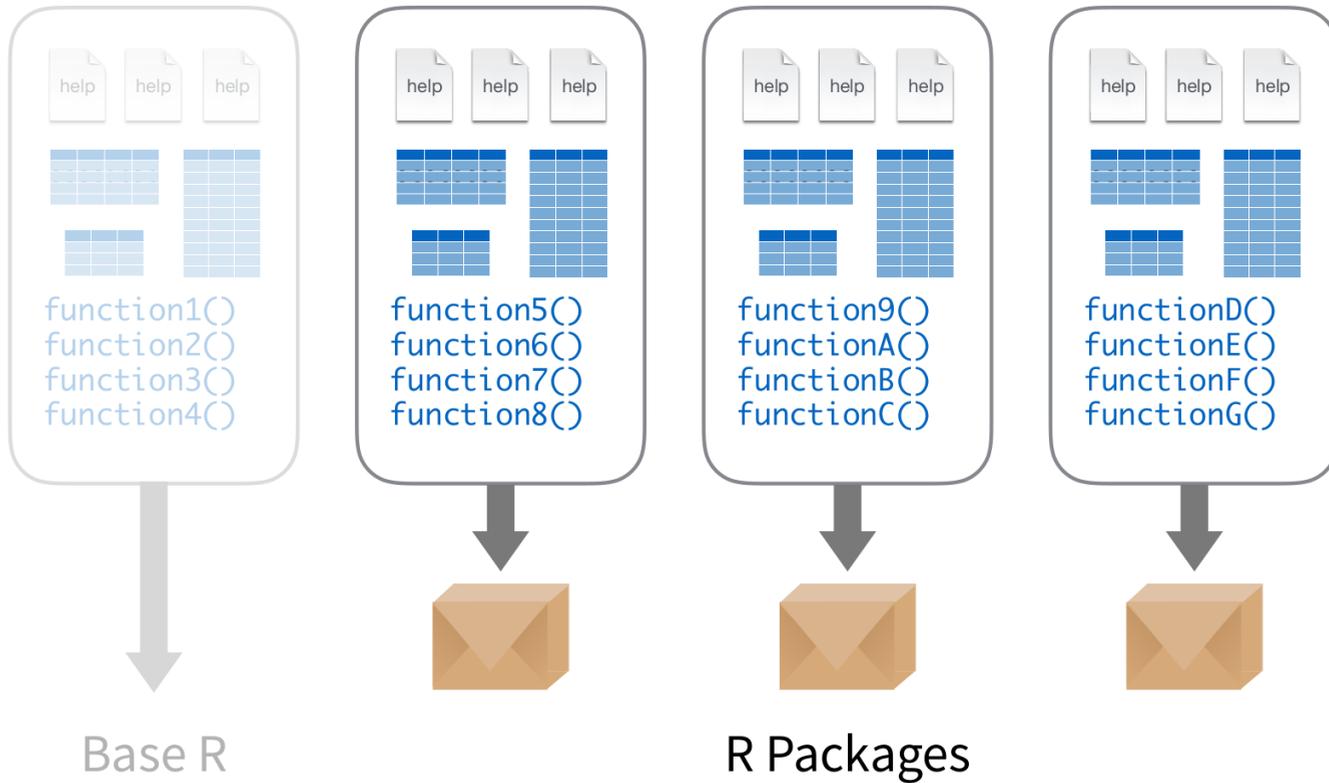


functionD()  
functionE()  
functionF()  
functionG()

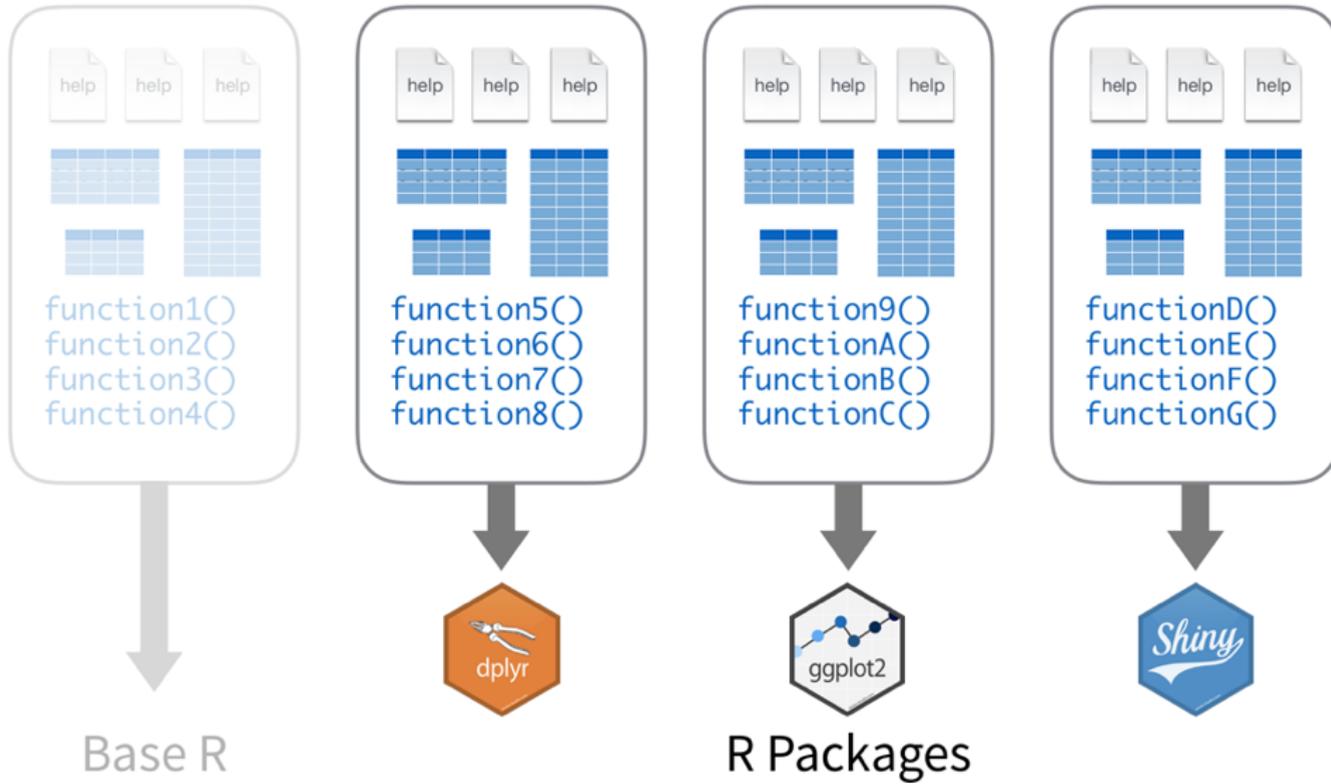
# What are R packages?



# What are R packages?



# What are R packages?



# How to install and run packages:

## 1

```
install.packages("foo")
```

Downloads files to computer

1 x per computer

# How to install and run packages:

**1**

```
install.packages("foo")
```

Downloads files to computer

1 x per computer

**2**

```
library("foo")
```

Loads package

1 x per R Session

# How to install and run packages:

**1**

```
install.packages("foo")
```

Downloads files to computer

1 x per computer

**2**

```
library("foo")
```

Loads package

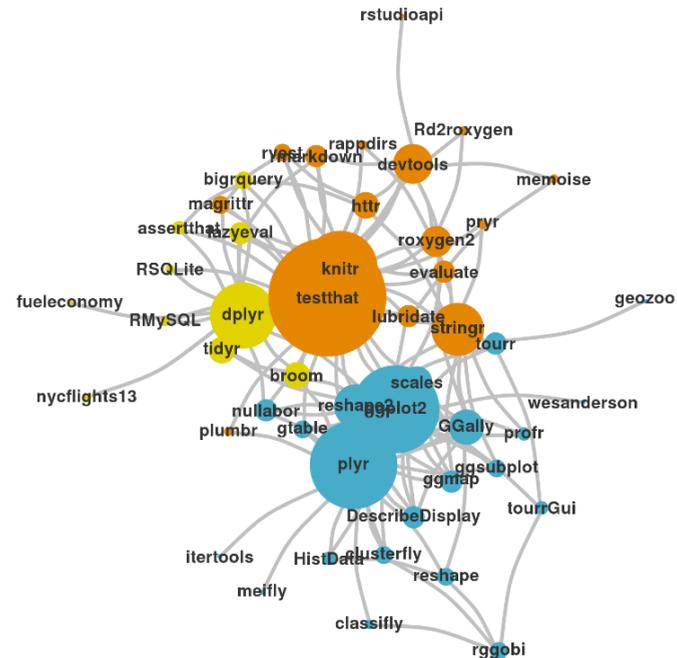
1 x per R Session

# How to install and run packages:

```
install.packages("tidyverse")
```

does the equivalent of

```
install.packages("ggplot2")
install.packages("dplyr")
install.packages("tidyr")
install.packages("readr")
install.packages("purrr")
install.packages("tibble")
install.packages("hms")
install.packages("stringr")
install.packages("lubridate")
install.packages("forcats")
install.packages("DBI")
install.packages("haven")
install.packages("httr")
install.packages("jsonlite")
install.packages("readxl")
install.packages("rvest")
install.packages("xml2")
install.packages("modelr")
install.packages("broom")
```



# How to install and run packages:

```
install.packages("tidyverse")
```

does the equivalent of

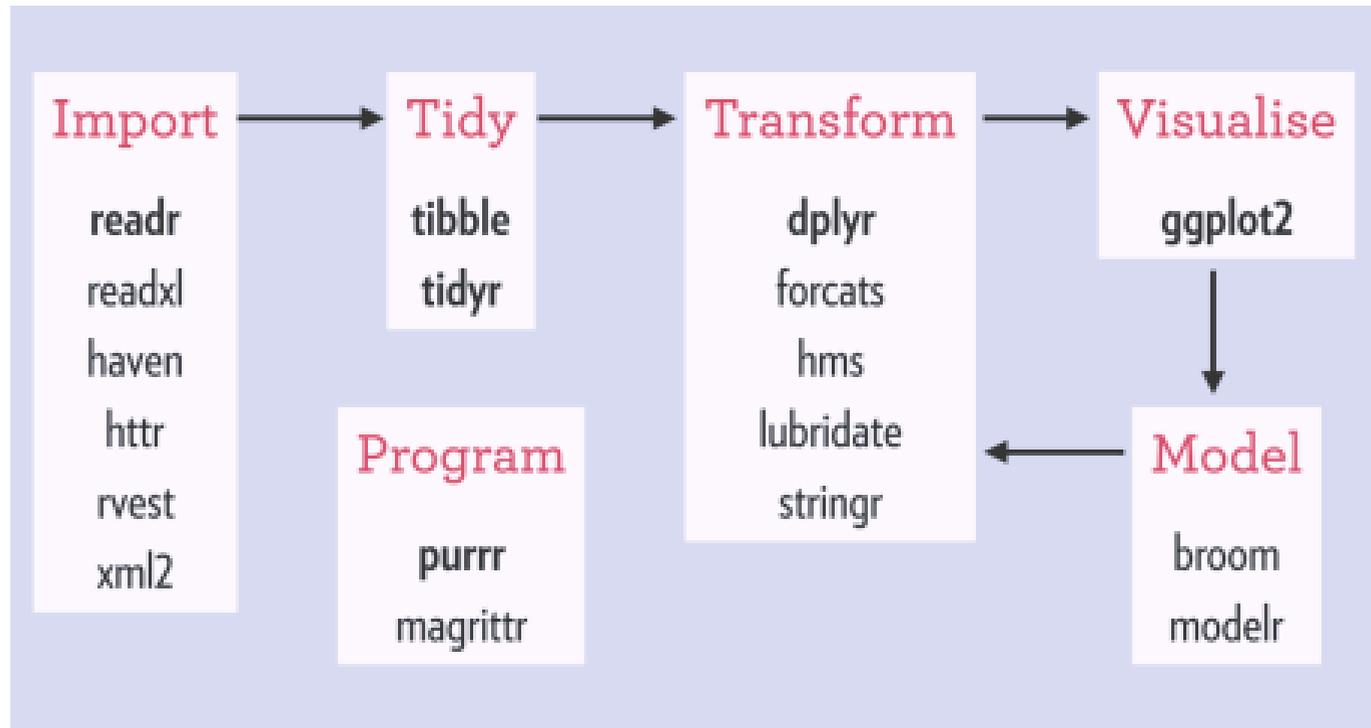
```
install.packages("ggplot2")  
install.packages("dplyr")  
install.packages("tidyr")  
install.packages("readr")  
install.packages("purrr")  
install.packages("tibble")  
install.packages("hms")  
install.packages("stringr")  
install.packages("lubridate")  
install.packages("forcats")  
install.packages("DBI")  
install.packages("haven")  
install.packages("httr")  
install.packages("jsonlite")  
install.packages("readxl")  
install.packages("rvest")  
install.packages("xml2")  
install.packages("modelr")  
install.packages("broom")
```

```
library("tidyverse")
```

does the equivalent of

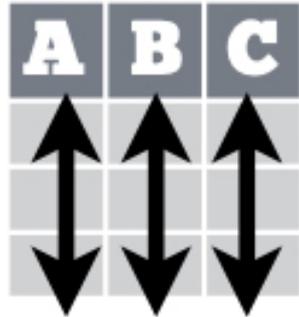
```
library("ggplot2")  
library("dplyr")  
library("tidyr")  
library("readr")  
library("purrr")  
library("tibble")
```

# Data science workflow



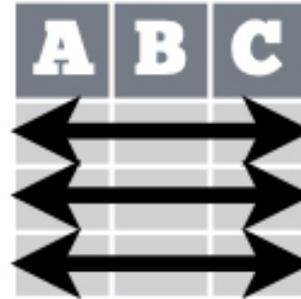
# tidy data

A table is tidy if:



Each **variable** is in its own **column**

&



Each **observation**, or **case**, is in its own **row**

“Tidy datasets are all alike but every messy dataset is messy in its own way.” – Hadley Wickham

# tidy data: "pivoting"

**gather**(data, key, value, ..., na.rm = FALSE, convert = FALSE, factor\_key = FALSE)

gather() moves column names into a **key** column, gathering the column values into a single **value** column.

table4a

| country | 1999 | 2000 |
|---------|------|------|
| A       | 0.7K | 2K   |
| B       | 37K  | 80K  |
| C       | 212K | 213K |

→

| country | year | cases |
|---------|------|-------|
| A       | 1999 | 0.7K  |
| B       | 1999 | 37K   |
| C       | 1999 | 212K  |
| A       | 2000 | 2K    |
| B       | 2000 | 80K   |
| C       | 2000 | 213K  |

key value

*gather(table4a, `1999`, `2000`,  
key = "year", value = "cases")*

**spread**(data, key, value, fill = NA, convert = FALSE, drop = TRUE, sep = NULL)

spread() moves the unique values of a **key** column into the column names, spreading the values of a **value** column across the new columns.

table2

| country | year | type  | count |
|---------|------|-------|-------|
| A       | 1999 | cases | 0.7K  |
| A       | 1999 | pop   | 19M   |
| A       | 2000 | cases | 2K    |
| A       | 2000 | pop   | 20M   |
| B       | 1999 | cases | 37K   |
| B       | 1999 | pop   | 172M  |
| B       | 2000 | cases | 80K   |
| B       | 2000 | pop   | 174M  |
| C       | 1999 | cases | 212K  |
| C       | 1999 | pop   | 1T    |
| C       | 2000 | cases | 213K  |
| C       | 2000 | pop   | 1T    |

key value

*spread(table2, type, count)*

# Core packages: dplyr, ggplot2, piping

## Data Transformation with dplyr : : CHEAT SHEET

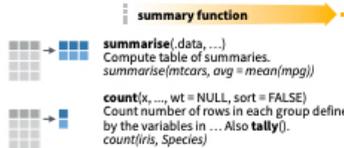


dplyr functions work with pipes and expect tidy data. In tidy data:



### Summarise Cases

These apply **summary functions** to columns to create a new table of summary statistics. Summary functions take vectors as input and return one value (see back).



#### VARIATIONS

- summarise\_all()** - Apply funs to every column.
- summarise\_at()** - Apply funs to specific columns.
- summarise\_if()** - Apply funs to all cols of one type.

### Group Cases

Use **group\_by()** to create a "grouped" copy of a table. dplyr functions will manipulate each "group" separately and then combine the results.



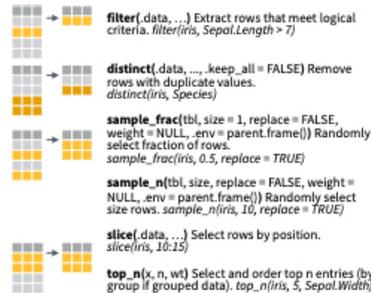
**group\_by(data, ..., add = FALSE)**  
Returns copy of table grouped by ...  
`g_iris <- group_by(iris, Species)`

**ungroup(x, ...)**  
Returns ungrouped copy of table.  
`ungroup(g_iris)`

### Manipulate Cases

#### EXTRACT CASES

Row functions return a subset of rows as a new table.

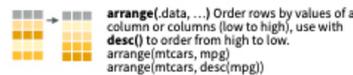


#### Logical and boolean operators to use with filter()

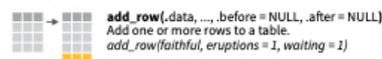
|   |    |          |      |   |       |
|---|----|----------|------|---|-------|
| < | <= | is.na()  | %in% |   | xor() |
| > | >= | !is.na() | !    | & |       |

See ?base::logic and ?Comparison for help.

#### ARRANGE CASES



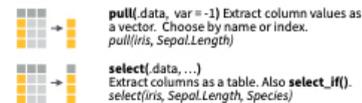
#### ADD CASES



### Manipulate Variables

#### EXTRACT VARIABLES

Column functions return a set of columns as a new vector or table.

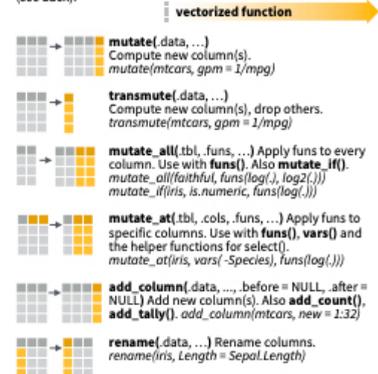


Use these helpers with **select()**, e.g. `select(iris, starts_with("Sepal"))`

**contains(match)** **num\_range(prefix, range)** ; e.g. `mpg:cyl`  
**ends\_with(match)** **one\_of(...)** ; e.g. `Species`  
**matches(match)** **starts\_with(match)**

#### MAKE NEW VARIABLES

These apply **vectorized functions** to columns. Vectorized funs take vectors as input and return vectors of the same length as output (see back).



# Core packages: dplyr, ggplot2, piping

## Data Visualization with ggplot2

Cheat Sheet



### Basics

ggplot2 is based on the **grammar of graphics**, the idea that you can build every graph from the same few components: a **data set**, a set of **geoms**—visual marks that represent data points, and a **coordinate system**.



To display data values, map variables in the data set to aesthetic properties of the geom like **size**, **color**, and **x** and **y** locations.



Build a graph with **qplot()** or **ggplot()**

**aes** aesthetic mappings **data** data **geom** geom

`qplot(x = cty, y = hwy, color = cyl, data = mpg, geom = "point")`

Creates a complete plot with given data, geom, and mappings. Supplies many useful defaults.

`ggplot(data = mpg, aes(x = cty, y = hwy))`

Begins a plot that you finish by adding layers to. No defaults, but provides more control than `qplot()`.

**data** **geom** **geom**

`ggplot(mpg, aes(hwy, cty)) + geom_point(aes(color = cyl)) + geom_smooth(method = "lm") + coord_cartesian() + scale_color_gradient() + theme_bw()`

**add layers, elements with +**  
**layer = geom + default stats + layer specific mappings**  
**additional elements**

Add a new layer to a plot with a **geom\_\*()** or **stat\_\*()** function. Each provides a geom, a set of aesthetic mappings, and a default stat and position adjustment.

**last\_plot()**  
Returns the last plot

**ggsave("plot.png", width = 5, height = 5)**  
Saves last plot as 5 x 5 file named "plot.png" in working directory. Matches file type to file extension.

RStudio is a trademark of RStudio, Inc. • [RStudio.com](http://RStudio.com) • info@rstudio.com • 844-468-2322 • [rstudio.com](http://rstudio.com)

## Geoms - Use a geom to represent data points, use the geom's aesthetic properties to represent variables. Each function returns a layer.

### One Variable

#### Continuous

`a <- ggplot(mpg, aes(hwy))`

**a** `geom_area(stat = "bin")`  
x, y, alpha, color, fill, linetype, size  
b `geom_area(aes(y = ..density.., stat = "bin"))`  
**a** `geom_density(kernel = "gaussian")`  
x, y, alpha, color, fill, linetype, size, weight  
b `geom_density(aes(y = ..density..))`  
**a** `geom_dotplot()`  
x, y, alpha, color, fill

**a** `geom_freqpoly()`  
x, y, alpha, color, linetype, size  
b `geom_freqpoly(aes(y = ..density..))`  
**a** `geom_histogram(binwidth = 5)`  
x, y, alpha, color, fill, linetype, size, weight  
b `geom_histogram(aes(y = ..density..))`

#### Discrete

b `ggplot(mpg, aes(fit))`

**a** `geom_bar()`  
x, alpha, color, fill, linetype, size, weight

### Two Variables

#### Continuous X, Continuous Y

f `ggplot(mpg, aes(cty, hwy))`

**f** `geom_blank()`

**f** `geom_jitter()`  
x, y, alpha, color, fill, shape, size

**f** `geom_point()`  
x, y, alpha, color, fill, shape, size

**f** `geom_quantile()`  
x, y, alpha, color, linetype, size, weight

**f** `geom_rug(sides = "bl")`  
alpha, color, linetype, size

**f** `geom_smooth(model = lm)`  
x, y, alpha, color, fill, linetype, size, weight

**f** `geom_text(aes(label = cty))`  
x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust

#### Discrete X, Continuous Y

g `ggplot(mpg, aes(class, hwy))`

**g** `geom_bar(stat = "identity")`  
x, y, alpha, color, fill, linetype, size, weight

**g** `geom_boxplot()`  
lower, middle, upper, x, ymax, ymin, alpha, color, fill, linetype, shape, size, weight

**g** `geom_dotplot(binaxis = "y", stackdir = "center")`

**g** `geom_violin(scale = "area")`  
x, y, alpha, color, fill, linetype, size, weight

#### Discrete X, Discrete Y

h `ggplot(diamonds, aes(cut, color))`

**h** `geom_jitter()`  
x, y, alpha, color, fill, shape, size

**h** `geom_rect(aes(xmin = long, ymin = lat, xmax = long + delta_long, ymax = lat + delta_lat))`  
x, xmax, ymin, ymax, alpha, color, fill, linetype, size

**e** `geom_segment(aes(xend = long + delta_long, yend = lat + delta_lat))`  
x, xend, y, yend, alpha, color, linetype, size

**e** `geom_rect(aes(xmin = long, ymin = lat, xmax = long + delta_long, ymax = lat + delta_lat))`  
x, xmax, ymin, ymax, alpha, color, fill, linetype, size

#### Continuous Bivariate Distribution

l `ggplot(movies, aes(year, rating))`

**i** `geom_bin2d(binwidth = c(5, 0.5))`  
xmax, xmin, ymax, ymin, alpha, color, fill, linetype, size, weight

**i** `geom_density2d()`  
x, y, alpha, colour, linetype, size

**i** `geom_hex()`  
x, y, alpha, colour, fill, size

#### Continuous Function

j `ggplot(economics, aes(date, unemployment))`

**j** `geom_area()`  
x, y, alpha, color, fill, linetype, size

**j** `geom_line()`  
x, y, alpha, color, linetype, size

**j** `geom_step(direction = "hv")`  
x, y, alpha, color, linetype, size

#### Visualizing error

df `<- data.frame(grp = c("A", "B"), fit = 4.5, se = 1.2)`  
k `ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))`

**k** `geom_crossbar(fatten = 2)`  
x, y, ymax, ymin, alpha, color, fill, linetype, size

**k** `geom_errorbar()`  
x, ymax, ymin, alpha, color, linetype, size, width (also `geom_errorbarh()`)

**k** `geom_linerange()`  
x, ymin, ymax, alpha, color, linetype, size

**k** `geom_pointrange()`  
x, y, ymin, ymax, alpha, color, fill, linetype, shape, size

#### Maps

data `<- data.frame(murder = USArrests$Murder, state = tolower(row.names(USArrests)))`  
map `<- map_data("state")`  
l `ggplot(data, aes(fill = murder))`

**l** `geom_map(aes(map_id = state), map = map) + expand_limits(x = map$long, y = map$lat, map_id = alpha, color, fill, linetype, size)`

### Three Variables

seals\$z <- with(seals, sqrt(delta\_long^2 + delta\_lat^2))  
m < ggplot(seals, aes(long, lat))

**m** `geom_raster(aes(fill = z), hjust = 0.5, vjust = 0.5, interpolate = FALSE)`  
x, y, alpha, fill

**m** `geom_tile(aes(fill = z))`  
x, y, alpha, color, fill, linetype, size, weight

Learn more at [docs.ggplot2.org](https://docs.ggplot2.org) • ggplot2 0.9.3.1 • Updated 3/15

# Core packages: dplyr, ggplot2, piping

What are the top 3 most popular years for males born "Taylor"?

```
# Load the data
library(babynames)
data(babynames)
```

```
df <- filter(babynames, sex=="M")
df <- filter(df, name=="Taylor")
df <- arrange(df, desc(n))
df <- select(df, year, n, prop)

head(df, n=3)
```

```
## # A tibble: 3 x 3
##   year     n  prop
##   <dbl> <int> <dbl>
## 1  1992  8240 0.00393
## 2  1991  7967 0.00376
## 3  1993  7688 0.00372
```

```
babynames %>%
  filter(sex=="M") %>%
  filter(name=="Taylor") %>%
  arrange(desc(n)) %>%
  select(year, n, prop) %>%
  head(n=3)
```

```
## # A tibble: 3 x 3
##   year     n  prop
##   <dbl> <int> <dbl>
## 1  1992  8240 0.00393
## 2  1991  7967 0.00376
## 3  1993  7688 0.00372
```

# Happy R programming!

