

# **Loading and saving data**

**(a refresher)**

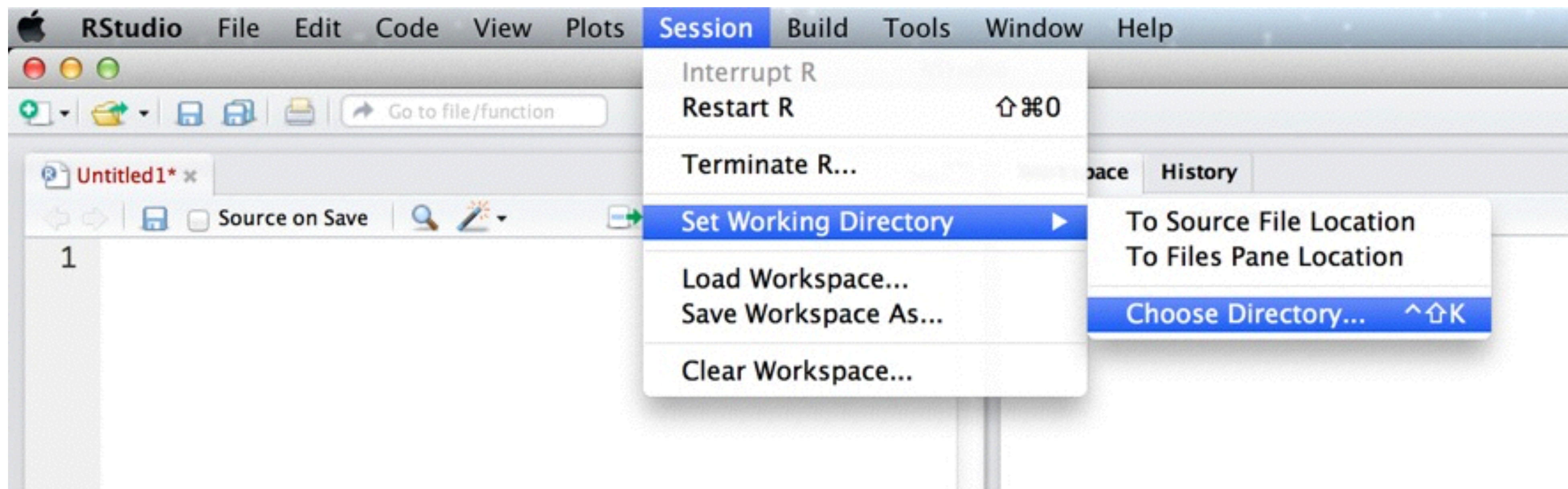
# Working directory

When you start R, it associates itself with a folder (i.e, directory) on your computer.

- This folder is known as your "**working directory**"
- When you save files, R will save them here
- When you load files, R will look for them here
- To see which folder is your working directory, run **getwd()**

# To change the working directory

In the toolbar, go to  
Session>Set Working Directory>Choose Directory...



# Loading data

Load csv files with `read.csv()`

```
read.csv("data/my-data.csv", stringsAsFactors = FALSE)
```

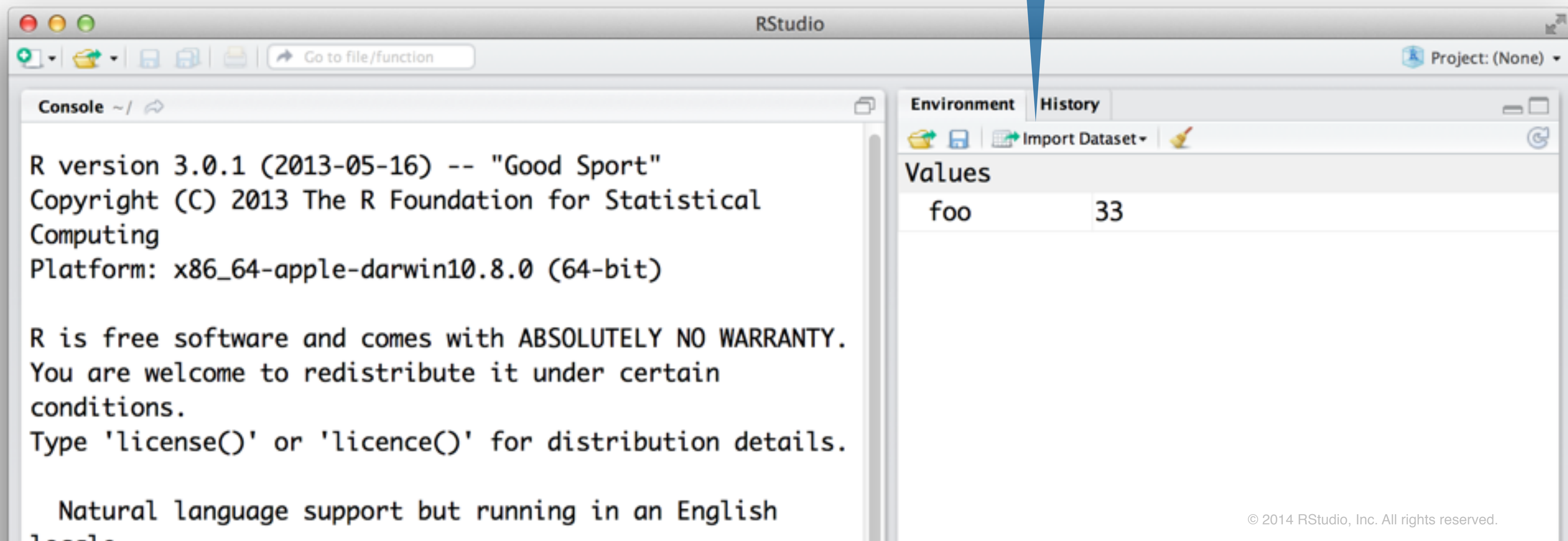
A file path from your working directory to the csv file

To prevent headaches



# Import Dataset

Click to select and  
read in a file



# Saving data

Save data frames with csv files with **write.csv()**

```
write.csv(df, file = "data/my-data.csv", row.names = FALSE)
```

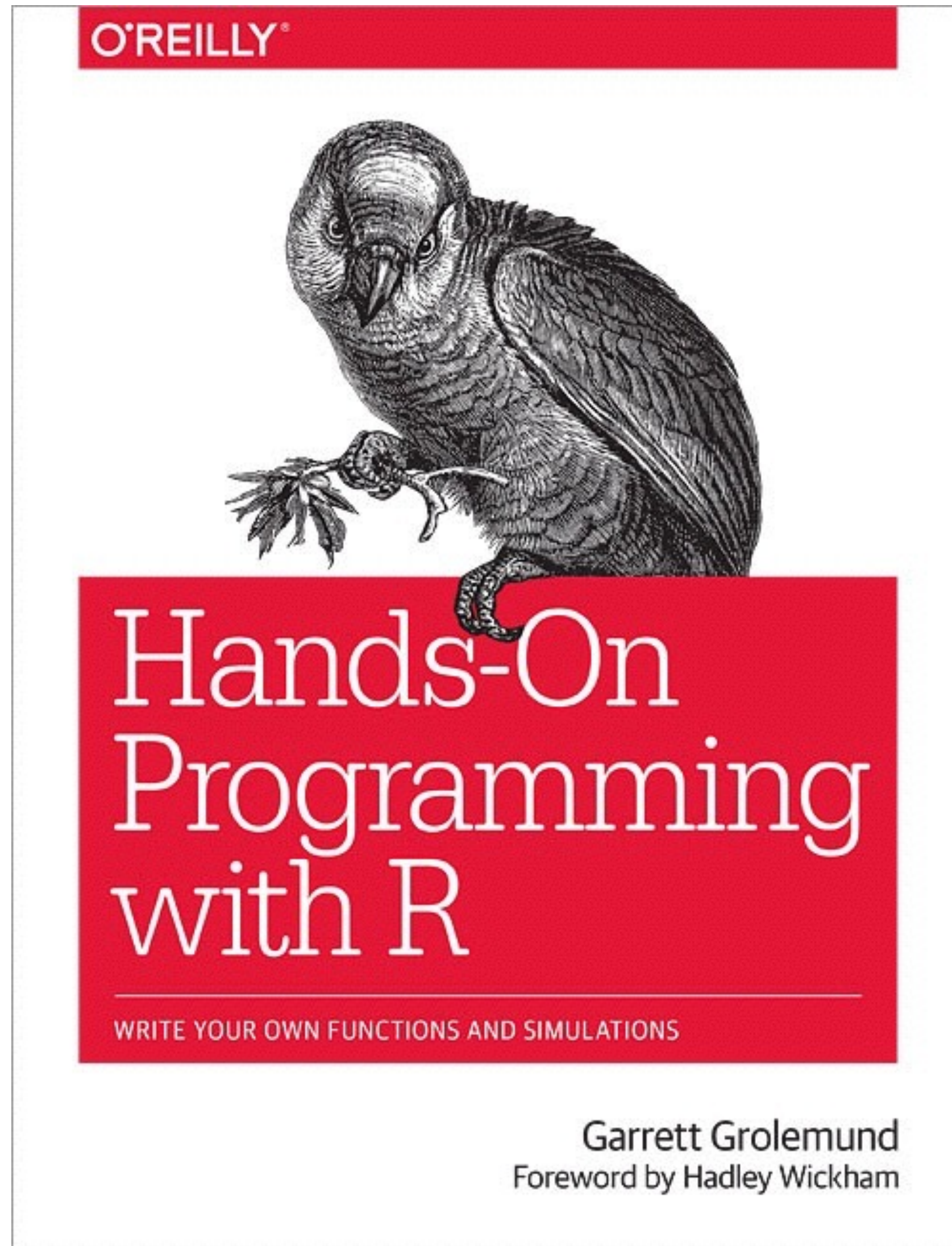
Data frame  
to save

Path to file you wish  
to create

To prevent  
headaches

**How to  
learn more**





# The R language

Three hands-on projects that teach the breadth of the R language.

- Storing, saving, and retrieving data.
- Writing programs and simulations.
- Vectorized programming.

[bit.ly/hands-on-programming-with-R](http://bit.ly/hands-on-programming-with-R)





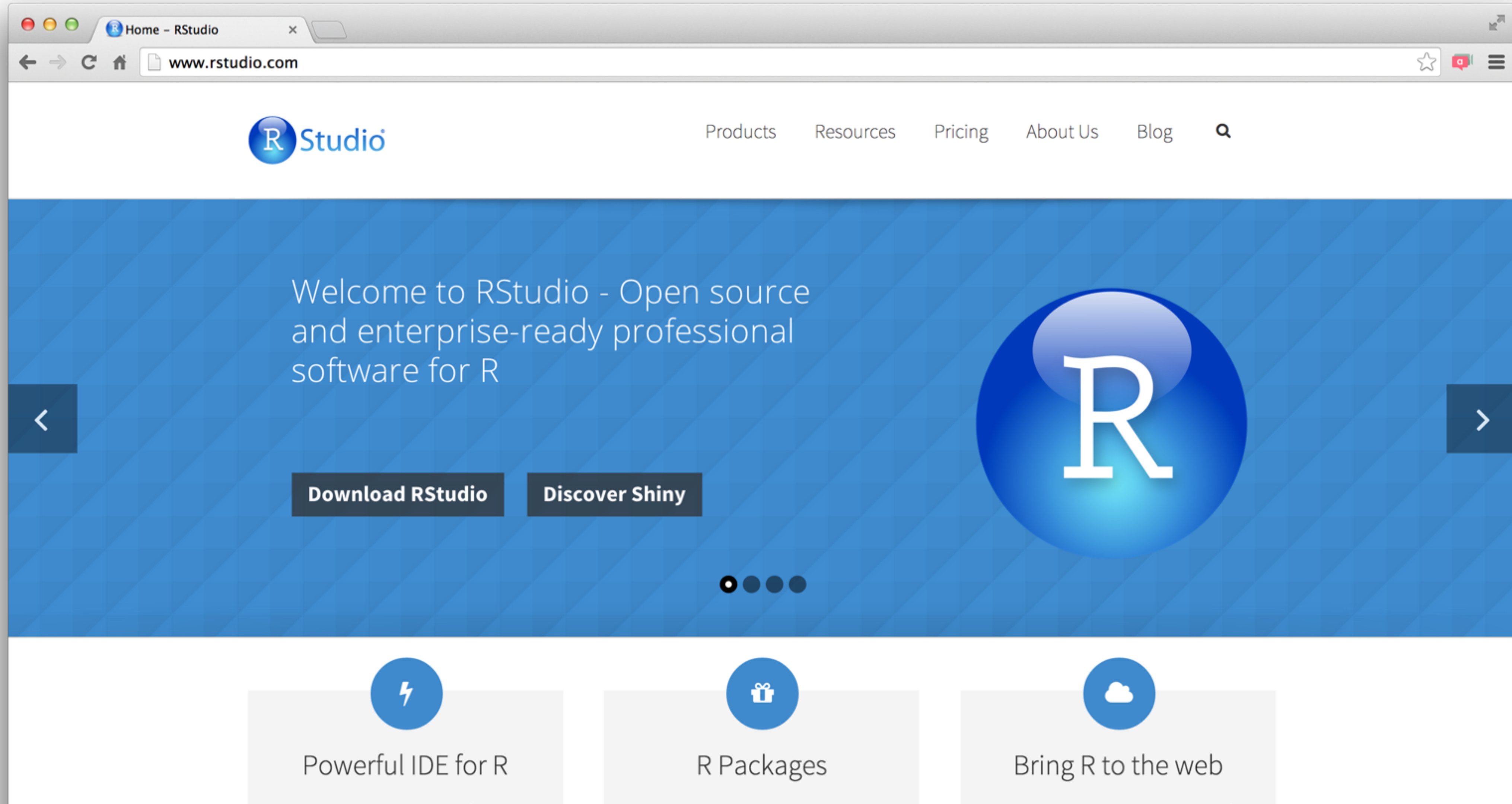
# Data Science with R

R's tools for data science.  
Reshape2, dplyr, and ggplot2  
packages.

- Tidy data
- Data visualization and customizing graphics
- Statistical modeling with R

[bit.ly/intro-to-data-science-with-R](http://bit.ly/intro-to-data-science-with-R)

# [www.rstudio.com](http://www.rstudio.com)





# Data Wrangling with dplyr and tidyr

Cheat Sheet



## Syntax - Helpful conventions for wrangling

### dplyr::tbl\_df(iris)

Converts data to tbl class. tbl's are easier to examine than data frames. R displays only the data that fits onscreen:

```
Source: local data frame [150 x 5]
  Sepal.Length Sepal.Width Petal.Length
1           5.1           3.5           1.4
2           4.9           3.0           1.4
3           4.7           3.2           1.3
4           4.6           3.1           1.5
5           5.0           3.6           1.4
..          ...           ...           ...
Variables not shown: Petal.Width (dbl),
Species (fctr)
```

### dplyr::glimpse(iris)

Information dense summary of tbl data.

### utils::View(iris)

View data set in spreadsheet-like display (note capital V).

	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
1	5.1	3.5	1.4	0.2	setosa
2	4.9	3.0	1.4	0.2	setosa
3	4.7	3.2	1.3	0.2	setosa
4	4.6	3.1	1.5	0.2	setosa
5	5.0	3.6	1.4	0.2	setosa
6	5.4	3.9	1.7	0.4	setosa
7	4.6	3.4	1.4	0.3	setosa
8	5.0	3.4	1.5	0.2	setosa

### dplyr::%>%

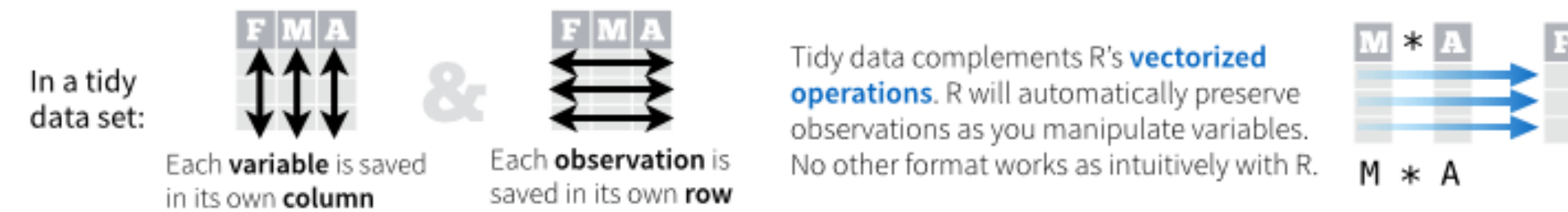
Passes object on left hand side as first argument (or argument) of function on righthand side.

`x %>% f(y)` is the same as `f(x, y)`  
`y %>% f(x, ., z)` is the same as `f(x, y, z)`

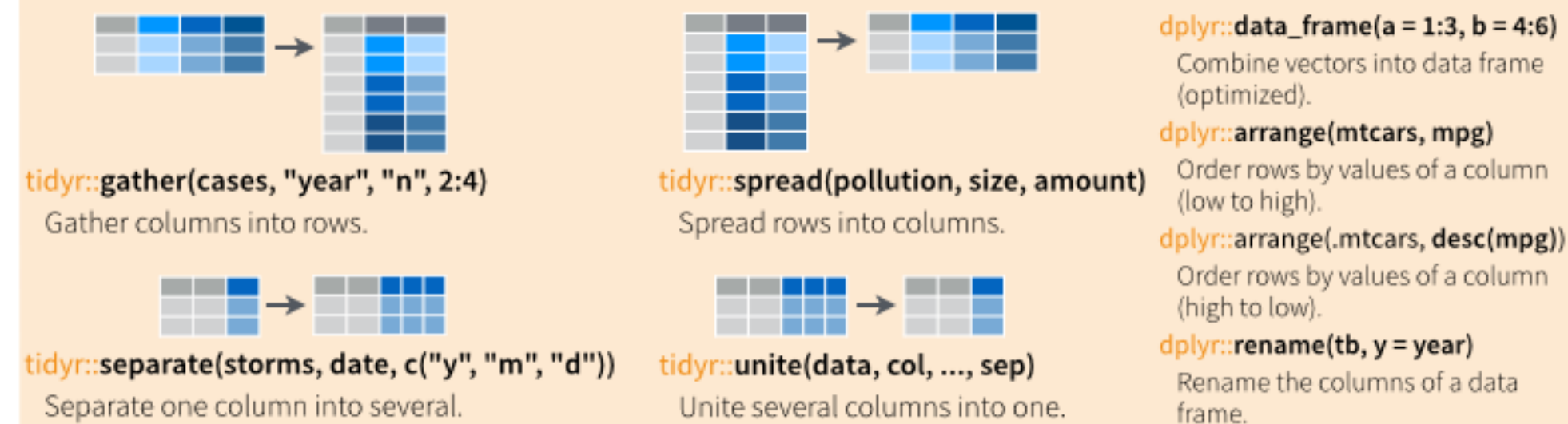
"Piping" with %>% makes code more readable, e.g.

```
iris %>%
  group_by(Species) %>%
  summarise(avg = mean(Sepal.Width)) %>%
  arrange(avg)
```

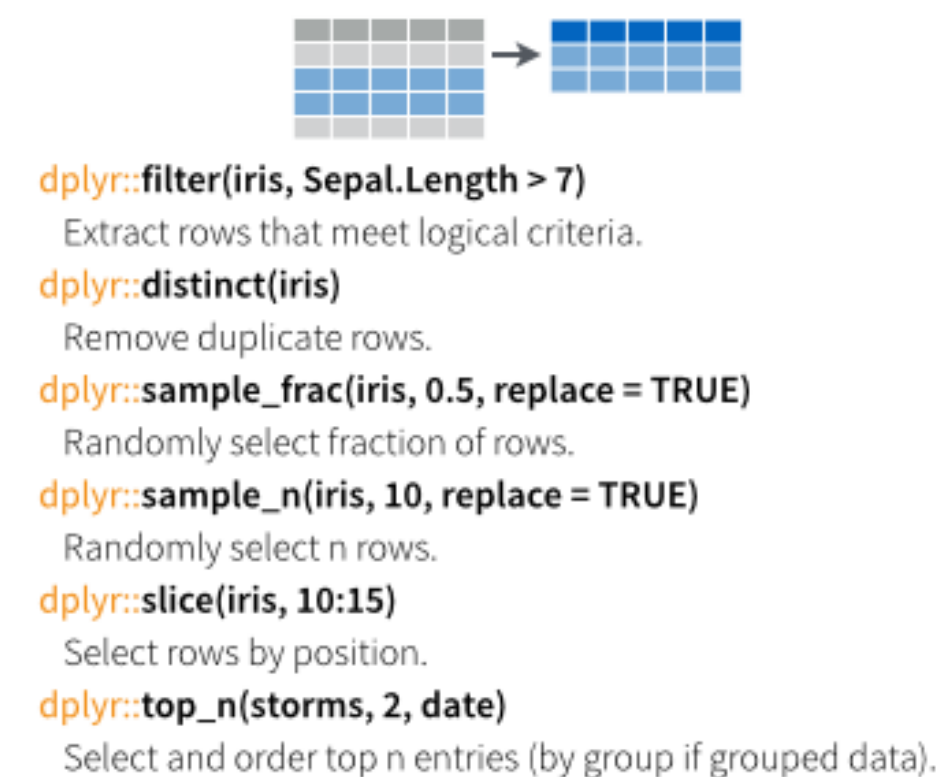
## Tidy Data - A foundation for wrangling in R



## Reshaping Data - Change the layout of a data set



## Subset Observations (Rows)



## Subset Variables (Columns)



### Helper functions for select - ?select

```
select(iris, contains("."))
  Select columns whose name contains a character string.
select(iris, ends_with("Length"))
  Select columns whose name ends with a character string.
select(iris, everything())
  Select every column.
select(iris, matches(".t."))
  Select columns whose name matches a regular expression.
select(iris, num_range("x", 1:5))
  Select columns named x1, x2, x3, x4, x5.
select(iris, one_of(c("Species", "Genus")))
  Select columns whose names are in a group of names.
select(iris, starts_with("Sepal"))
  Select columns whose name starts with a character string.
select(iris, Sepal.Length:Petal.Width)
  Select all columns between Sepal.Length and Petal.Width (inclusive).
select(iris, -Species)
  Select all columns except Species.
```

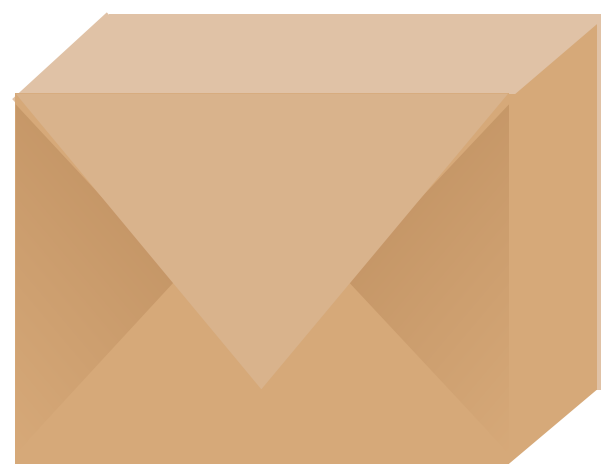
### Logic in R - ?Comparison, ?base::Logic

<	Less than	!=	Not equal to
>	Greater than	%in%	Group membership
==	Equal to	is.na	Is NA
<=	Less than or equal to	!is.na	Is not NA
>=	Greater than or equal to	&,  , !, xor, any, all	Boolean operators



# **Final Project**

# nycflights13



A package of data sets about airline travel in 2013 in New York City.

```
library(nycflights13)
```

```
?airlines
```

```
?airports
```

```
?flights
```

```
?planes
```

```
?weather
```

# flights

(336,776 x 16)

On-time data for all flights that departed NYC (i.e. JFK, LGA or EWR) in 2013.

year	month	day	dep_time	dep_delay	arr_time	arr_delay	carrier	tailnum	flight	origin	dest	air_time	distance	hour	minute
2013	1	1	517	2	830	11	UA	N14228	1545	EWR	IAH	227	1400	5	17
2013	1	1	533	4	850	20	UA	N24211	1714	LGA	IAH	227	1416	5	33
2013	1	1	542	2	923	33	AA	N619AA	1141	JFK	MIA	160	1089	5	42
2013	1	1	544	-1	1004	-18	B6	N804JB	725	JFK	BQN	183	1576	5	44
2013	1	1	554	-6	812	-25	DL	N668DN	461	LGA	ATL	116	762	5	54
2013	1	1	554	-4	740	12	UA	N39463	1696	EWR	ORD	150	1065	5	54



# flights

(336,776 x 16)

On-time data for all flights that departed NYC (i.e. JFK, LGA or EWR) in 2013.

year	month	day	dep_time	dep_delay	arr_time	arr_delay	carrier	tailnum	flight	origin	dest	air_time	distance	hour	minute
2013	1	1	517	2	830	11	UA	N14228	1545	EWR	IAH	227	1400	5	17
2013	1	1	533	4	850	20	UA	N24211	1714	LGA	IAH	227	1416	5	33
2013	1	1	542	2	923	33	AA	N619AA	1141	JFK	MIA	160	1089	5	42
2013	1	1	544	-1	1004	-18	B6	N804JB	725	JFK	BQN	183	1576	5	44
2013	1	1	554	-6	812	-25	DL	N668DN	461	LGA	ATL	116	762	5	54
2013	1	1	554	-4	740	12	UA	N39463	1696	EWR	ORD	150	1065	5	54

# airlines

(16 x 2)

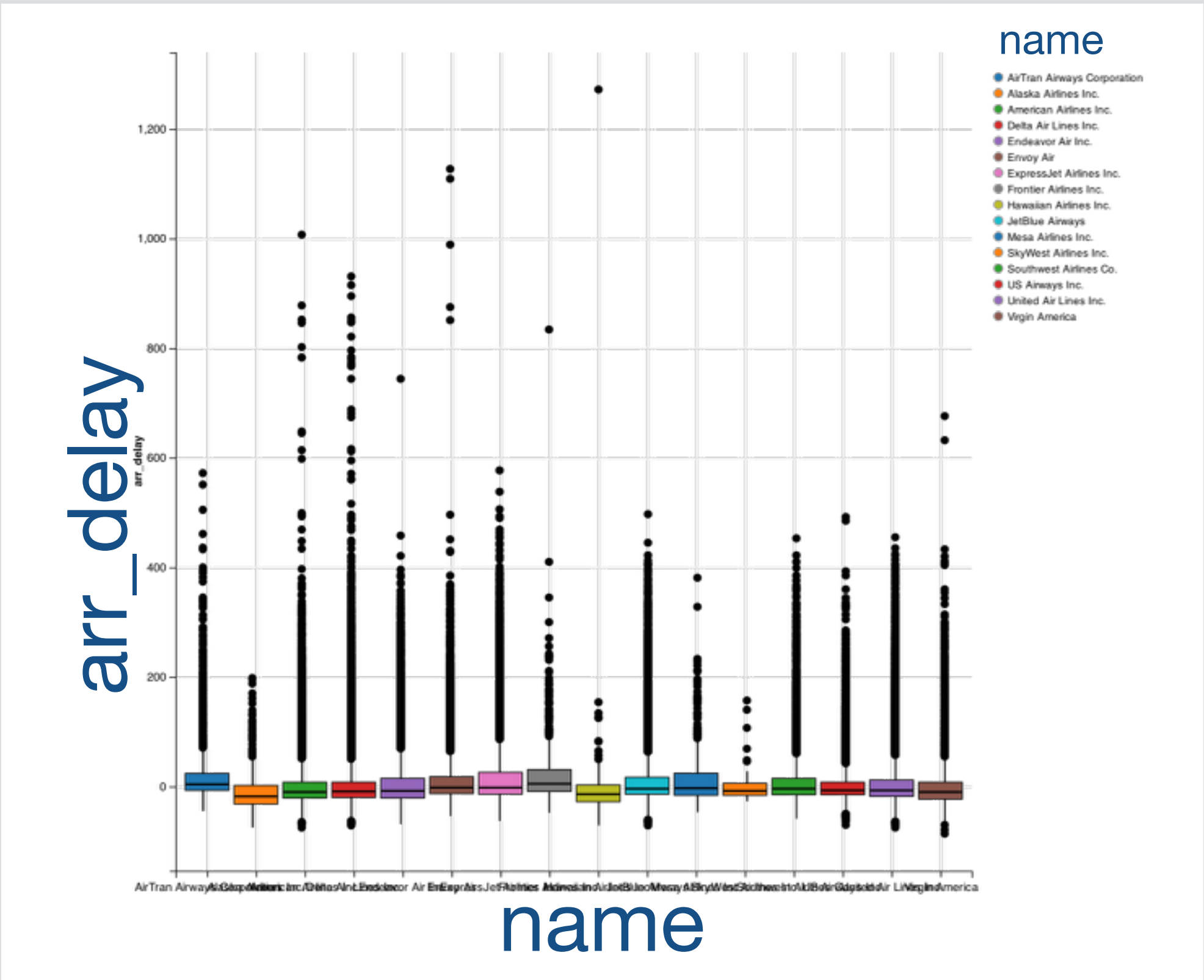
Airline names by carrier code

carrier	name
9E	Endeavor Air Inc.
AA	American Airlines Inc.
AS	Alaska Airlines Inc.
B6	JetBlue Airways
DL	Delta Air Lines Inc.
EV	ExpressJet Airlines Inc.

# Final project

Use the flights and airlines datasets to recreate the plot on the left and the table on the right.

(16 x 2)



name	median (arrival delay)
Frontier Airlines Inc.	6
AirTran Airways Corporation	5
Envoy Air	-1
ExpressJet Airlines Inc.	-1
Mesa Airlines Inc.	-2
JetBlue Airways	-3
Southwest Airlines Co.	-3
United Air Lines Inc.	-6
US Airways Inc.	-6
Endeavor Air Inc.	-7



# Plot

**1 Join**

**2 Filter** (NA's)

**3 Visualize**

# Table

**1 Join**

**2 Filter** (NA's)

**3 Group**

**4 Summarise**

**5 Arrange**

# Plot

```
flights %>%  
  left_join(airlines, by = "carrier") %>%  
  filter(!is.na(arr_delay)) %>%  
  ggvis(x = ~name, y = ~arr_delay, fill = ~name) %>%  
  layer_boxplots()
```

# Table

```
flights %>%
```

```
  left_join(airlines, by = "carrier") %>%
```

```
  filter(!is.na(arr_delay)) %>%
```

```
  group_by(name) %>%
```

```
  summarise(median = median(arr_delay)) %>%
```

```
  arrange(desc(median))
```