


# Data Wrangling

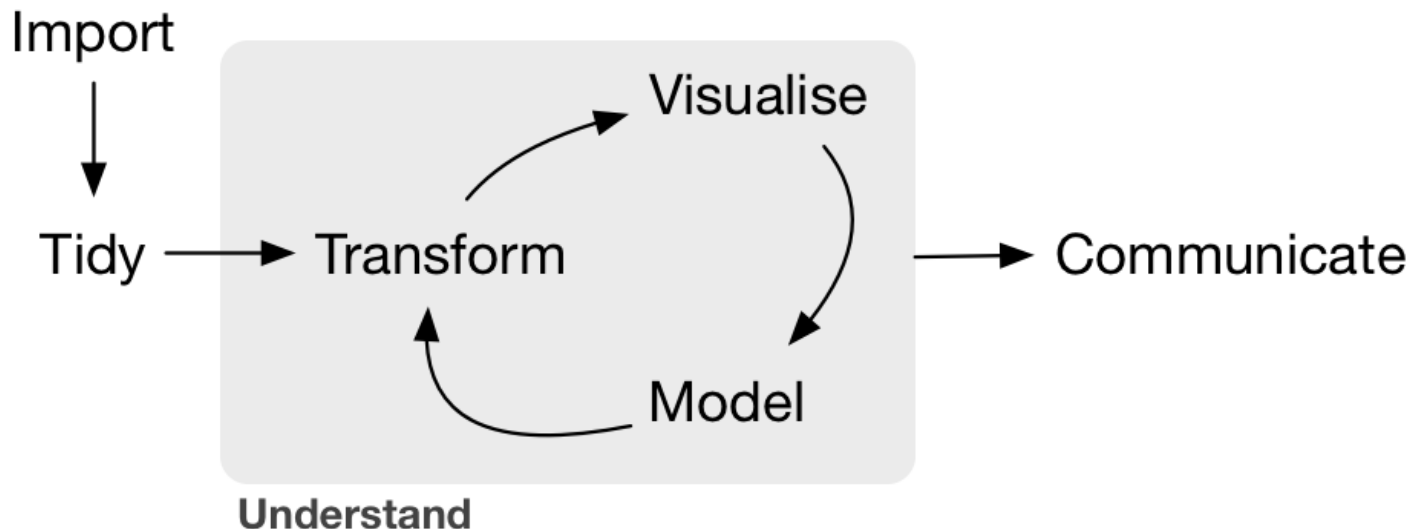
Country	2011	2012	2013
FR	7000	8000	9000
DE	5000	6000	6200
US	15000	16000	18000



# Data Wrangling: Two Goals

---

- 1 Make data suitable to use with a particular piece of software
- 2 Reveal information



Wrangling  
Munging  
Janitor Work  
Manipulation  
Transformation

50-80%

of your time?

# Two packages to help you work with the structure of data

---



tidyr



dplyr

# Data Wrangling with dplyr and tidy

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).

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)
```

RStudio® is a trademark of RStudio, Inc. • All rights reserved • info@rstudio.com • 844-448-1212 • rstudio.com

## Tidy Data - A foundation for wrangling in R



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

**dplyr: data\_frame(a = 1:3, b = 4:6)**  
Combine vectors into data frame (optimized).

**dplyr: arrange(mtcars, mpg)**  
Order rows by values of a column (low to high).

**dplyr: arrange(mtcars, desc(mpg))**  
Order rows by values of a column (high to low).

**dplyr: rename(tb, y = year)**  
Rename the columns of a data frame.

## Subset Observations (Rows)

**dplyr: filter(iris, Sepal.Length > 7)**  
Extract rows that meet logical criteria.

**dplyr: distinct(iris)**  
Remove duplicate rows.

**dplyr: sample\_frac(iris, 0.5, replace = TRUE)**  
Randomly select fraction of rows.

**dplyr: sample\_n(iris, 10, replace = TRUE)**  
Randomly select n rows.

**dplyr: slice(iris, 10:15)**  
Select rows by position.

**dplyr: top\_n(storms, 2, date)**  
Select and order top n entries (by group if grouped data).

## Subset Variables (Columns)

**dplyr: select(iris, Sepal.Width, Petal.Length, Species)**  
Select columns by name or helper function.

### 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("L"))
  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

devtools::install\_github("rstudio/EDAWR") for data sets Learn more with browseVignettes(package = c("dplyr", "tidyr")) • dplyr 0.4.0• tidyr 0.2.0 • Updated: 1/15

Also in Chinese...

## 利用dplyr和tidyr进行数据再加工



由 @supatcat 翻译

语法 - 有用的数据再加工规则

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

将数据转换为tbl类。相比数据框，tbl更易于查看，R只在显示适合屏幕大小的数据。

```
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)
```

## 整洁数据 - 在R中进行数据再加工的基础



## 重组数据 - 改变数据表的布局

**dplyr: data\_frame(a = 1:3, b = 4:6)**  
将向量合并成数据框 (已优化)

**dplyr: arrange(mtcars, mpg)**  
按指定列的值对数据进行排序 (从小到大)

**dplyr: arrange(mtcars, desc(mpg))**  
按指定列的值对数据进行排序 (从大到小)

**dplyr: separate(storms, date, c("y", "m", "d"))**  
将多列拆分为多列。

**dplyr: unite(data, col, ..., sep)**  
将多列拼为一列。

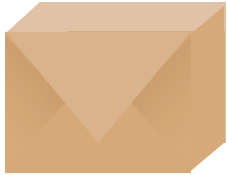
**dplyr: rename(tb, y = year)**  
重命名数据框中的列名。

Data sets come in many  
formats

...but R (often) prefers just one

# EDAWR

---



An R package with all of the data sets that shown in this lecture.

```
# install.packages("devtools")  
# devtools::install_github("rstudio/EDAWR")  
library(EDAWR)  
?storms  
?cases  
?pollution  
?tb
```

```
devtools::install_github("rstudio/EDAWR")
library(EDAWR)
```

---

storms

storm	wind	pressure	date
Arthur	45	100	2000-08-12
Alexander	45	100	1998-07-30
Allison	65	100	1995-06-04
Anne	40	101	1997-07-01
Andrew	50	101	1999-06-13
Arthur	45	101	1996-05-21

- Storm name
- Wind Speed (mph)
- Air pressure
- Date

pollution

Country	2000	2001	2002
FR	7000	6000	6000
DE	5000	6000	6200
US	15000	15000	15000

- Country
- Year
- Count

cases

city	particle size	amount ( $\mu\text{g}/\text{m}^3$ )
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56

- City
- Amount of Large Particles
- Amount of Small particles



```
devtools::install_github("rstudio/EDAWR")
library(EDAWR)
```

---

storms

storm	wind	pressure	date
Arthur	45	100	2000-08-12
Alexander	45	100	1998-07-30
Allison	65	100	1995-06-04
Anna	40	101	1997-07-01
Anna	45	101	1999-06-13
Arthur	45	101	1996-05-21

```
storms$storm
storms$wind
storms$pressure
storms$date
```

pollution

Country	2013	2014	2015
FR	7000	6000	6000
DE	5000	6000	6200
US	15000	6000	6000

```
cases$country
names(cases)[-1]
unlist(cases[1:3, 2:4])
```

cases

city	particle size	amount (µg/m³)
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56

```
pollution$city[1,3,5]
pollution$amount[1,3,5]
pollution$amount[2,4,6]
```

# Adding/modifying columns

---

$$\text{ratio} = \frac{\text{pressure}}{\text{wind}}$$

storms

storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Alex	45	1009	1998-07-30
Allison	65	1005	1995-06-04
Ana	40	1013	1997-07-01
Arlene	50	1010	1999-06-13
Arthur	45	1010	1996-06-21

`storms$pressure / storms$wind`

950	/	110	→	8.6
1003	/	45	→	22.3
987	/	65	→	15.2
1004	/	40	→	25.1
1006	/	50	→	20.1
1000	/	45	→	22.2

# Tidy data

---

storms

storm	wind	pressure	date
Alfred	7	1012	2007-07-12
Alfred	4	100	1998-07-30
Alfred	8	100	1995-08-04
Alfred	4	101	1997-07-31
Alfred	7	10	1995-08-13
Alfred	43	1010	1998-08-21

1

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

2

Each **observation** is saved in its own **row**.

3

Each "type" of observation stored in a **single table** (here, storms).

# Recap: Tidy data

---

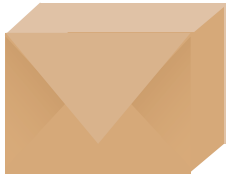
- 1 Variables in columns, observations in rows, each type in a table
- 2 Easy to access variables
- 3 Automatically preserves observations

tidyr

Tidyr: A package that reshapes the layout of tables.

---

## tidyr



Two main functions: **gather()** and **spread()**

```
# install.packages("tidyr")  
library(tidyr)  
?gather  
?spread
```

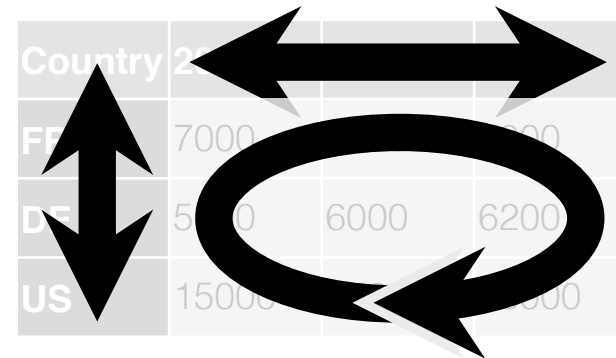
# Your Turn

---

Imagine how this data would look if it were tidy with three variables:  
*country, year, n*

cases

Country	2011	2012	2013
FR	7000	6900	7000
DE	5800	6000	6200
US	15000	14000	13000



---

Country	2011	2012	2013
FR	7000	6900	7000
DE	5800	6000	6200
US	15000	14000	13000



---

Country	2011	2012	2013
FR	7000	6900	7000
DE	5800	6000	6200
US	15000	14000	13000

Country	Year	n
---------	------	---

---

Country	2011	2012	2013
FR	7000	6900	7000
DE	5800	6000	6200
US	15000	14000	13000

Country	Year	n
FR	2011	7000

---

Country	2011	2012	2013
FR	7000	6900	7000
DE	5800	6000	6200
US	15000	14000	13000

Country	Year	n
FR	2011	7000
DE	2011	5800

---

Country	2011	2012	2013
FR	7000	6900	7000
DE	5800	6000	6200
US	15000	14000	13000

Country	Year	n
FR	2011	7000
DE	2011	5800
US	2011	15000

---

Country	2011	2012	2013
FR	7000	6900	7000
DE	5800	6000	6200
US	15000	14000	13000

Country	Year	n
FR	2011	7000
DE	2011	5800
US	2011	15000
FR	2012	6900

---

Country	2011	2012	2013
FR	7000	6900	7000
DE	5800	6000	6200
US	15000	14000	13000

Country	Year	n
FR	2011	7000
DE	2011	5800
US	2011	15000
FR	2012	6900
DE	2012	6000

---

Country	2011	2012	2013
FR	7000	6900	7000
DE	5800	6000	6200
US	15000	14000	13000

Country	Year	n
FR	2011	7000
DE	2011	5800
US	2011	15000
FR	2012	6900
DE	2012	6000
US	2012	14000

---

Country	2011	2012	2013
FR	7000	6900	7000
DE	5800	6000	6200
US	15000	14000	13000

Country	Year	n
FR	2011	7000
DE	2011	5800
US	2011	15000
FR	2012	6900
DE	2012	6000
US	2012	14000
FR	2013	7000



---

Country	2011	2012	2013
FR	7000	6900	7000
DE	5800	6000	6200
US	15000	14000	13000

Country	Year	n
FR	2011	7000
DE	2011	5800
US	2011	15000
FR	2012	6900
DE	2012	6000
US	2012	14000
FR	2013	7000
DE	2013	6200

---

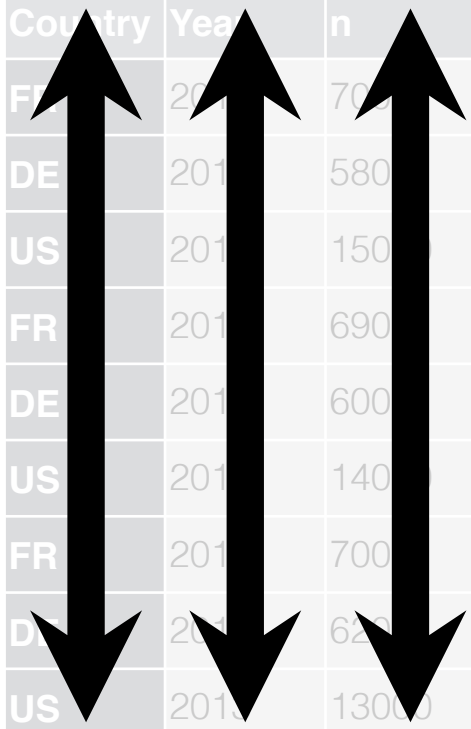
Country	2011	2012	2013
FR	7000	6900	7000
DE	5800	6000	6200
US	15000	14000	13000

Country	Year	n
FR	2011	7000
DE	2011	5800
US	2011	15000
FR	2012	6900
DE	2012	6000
US	2012	14000
FR	2013	7000
DE	2013	6200
US	2013	13000

---

Country	2011	2012	2013
FR	7000	6900	7000
DE	5800	6000	6200
US	15000	14000	13000

Country	Year	n
FR	2011	7000
DE	2011	5800
US	2011	15000
FR	2012	6900
DE	2012	6000
US	2012	14000
FR	2013	7000
DE	2013	6200
US	2013	13000



---

Country	2011	2012	2013
FR	7000	6900	7000
DE	5800	6000	6200
US	15000	14000	13000



Country	Year	n
FR	2011	7000
DE	2011	5800
US	2011	15000
FR	2012	6900
DE	2012	6000
US	2012	14000
FR	2013	7000
DE	2013	6200
US	2013	13000

---

Count	2011	2012	2013
FR	7000	6900	7000
DE	5800	6000	6200
US	15000	14000	13000

Country	Year	n
FR	2011	7000
DE	2011	5800
US	2011	15000
FR	2012	6900
DE	2012	6000
US	2012	14000
FR	2013	7000
DE	2013	6200
US	2013	13000

---

key (former column names)

Country	2011	2012	2013
FR	7000	6900	7000
DE	5800	6000	6200
US	15000	14000	13000

Country	Year	n
FR	2011	7000
DE	2011	5800
US	2011	15000
FR	2012	6900
DE	2012	6000
US	2012	14000
FR	2013	7000
DE	2013	6200
US	2013	13000

---

key value (former cells)

Country	2011	2012	2013
FR	7000	6900	7000
DE	5800	6000	6200
US	15000	14000	13000

Country	Year	n
FR	2011	7000
DE	2011	5800
US	2011	15000
FR	2012	6900
DE	2012	6000
US	2012	14000
FR	2013	7000
DE	2013	6200
US	2013	13000

# gather()

---

Collapses multiple columns into two columns:

1. a **key** column that contains the former column names
2. a **value** column that contains the former column cells

```
gather(cases, "year", "n", 2:4)
```

data frame  
to reshape

name of the new key  
column  
(a character string)

name of the new  
value column  
(a character string)

names or numeric  
indexes of columns to  
collapse



# gather(cases, "year", "n", 2:4)

---

##	country	2011	2012	2013
## 1	FR	7000	6900	7000
## 2	DE	5800	6000	6200
## 3	US	15000	14000	13000



##	country	year	n
## 1	FR	2011	7000
## 2	DE	2011	5800
## 3	US	2011	15000
## 4	FR	2012	6900
## 5	DE	2012	6000
## 6	US	2012	14000
## 7	FR	2013	7000
## 8	DE	2013	6200
## 9	US	2013	13000

# Your Turn

---

Imagine how the pollution data set would look tidy with three variables:  
*city, large, small*

pollution

city	size	amount
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56

city	particle size	amount ( $\mu\text{g}/\text{m}^3$ )
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56

---

city	size	amount
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56

---

city	size	amount
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56

city	large	small
------	-------	-------

---

city	size	amount
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56

city	large	small
New York	23	

---

city	size	amount
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56

city	large	small
New York	23	14

---

city	size	amount
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56

city	large	small
New York	23	14
London	22	

---

city	size	amount
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56

city	large	small
New York	23	14
London	22	16



---

city	size	amount
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56

city	large	small
New York	23	14
London	22	16
Beijing	121	

---

city	size	amount
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56

city	large	small
New York	23	14
London	22	16
Beijing	121	56

---

city	size	amount
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56

city	large	small
New York	23	14
London	22	16
Beijing	121	56

---

city	size	amount
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56



city	large	small
New York	23	14
London	22	16
Beijing	121	56

---

**key (new column names)**

city	size	amount
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56

city	large	small
New York	23	14
London	22	16
Beijing	121	56

---

key    value (new cells)

city	size	amount
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56

city	large	small
New York	23	14
London	22	16
Beijing	121	56

# spread()

---

Generates multiple columns from two columns:

1. each unique value in the **key** column becomes a column name
2. each value in the **value** column becomes a cell in the new columns

`spread(pollution, size, amount)`

data frame to  
reshape

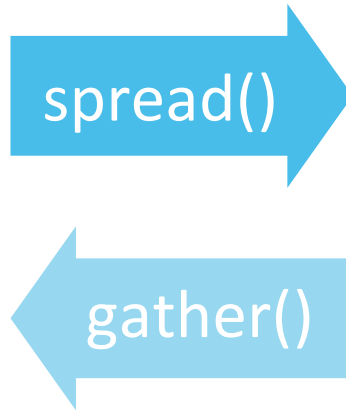
column to use for  
keys (new columns  
names)

column to use for  
values (new column  
cells)

# spread(pollution, size, amount)

---

city	size	amount
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56



city	large	small
New York	23	14
London	22	16
Beijing	121	56

Separate **all variables** *implied by law, formula or goal*



# unite() and separate()

---

There are three more variables hidden in storms:

storms

storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Alex	45	1009	1998-07-30
Allison	65	1005	1995-06-04
Ana	40	1013	1997-07-01
Arlene	50	1010	1999-06-13
Arthur	45	1010	1996-06-21

- Year
- Month
- Day

# separate()

---

Separate splits a column by a character string separator.

```
separate(storms, date, c("year", "month", "day"), sep = "-")
```

storms

storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Alex	45	1009	1998-07-30
Allison	65	1005	1995-06-04
Ana	40	1013	1997-07-01
Arlene	50	1010	1999-06-13
Arthur	45	1010	1996-06-21



storms2

storm	wind	pressure	year	month	day
Alberto	110	1007	2000	08	12
Alex	45	1009	1998	07	30
Allison	65	1005	1995	06	04
Ana	40	1013	1997	07	1
Arlene	50	1010	1999	06	13
Arthur	45	1010	1996	06	21

# unite()

---

Unite unites columns into a single column.

```
unite(storms2, "date", year, month, day, sep = "-")
```

storms2

storm	wind	pressure	year	month	day
Alberto	110	1007	2000	08	12
Alex	45	1009	1998	07	30
Allison	65	1005	1995	06	04
Ana	40	1013	1997	07	1
Arlene	50	1010	1999	06	13
Arthur	45	1010	1996	06	21



storms

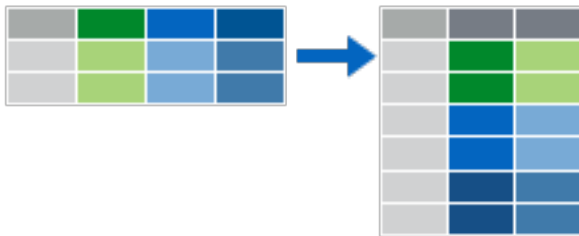
storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Alex	45	1009	1998-07-30
Allison	65	1005	1995-06-04
Ana	40	1013	1997-07-01
Arlene	50	1010	1999-06-13
Arthur	45	1010	1996-06-21

# Recap: tidyr

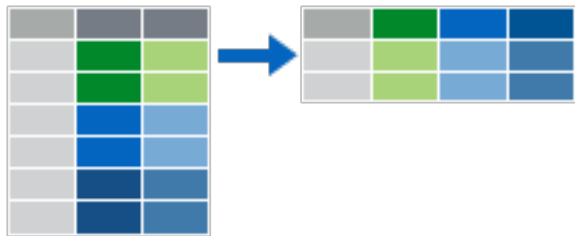
---



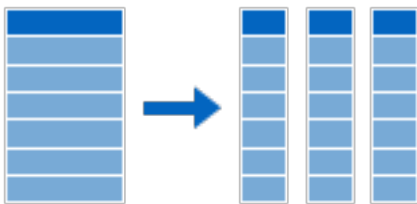
A package that reshapes the layout of data sets.



Make observations from variables with `gather()`



Make variables from observations with `spread()`



Split and merge columns with `unite()` and `separate()`

Also `reshape2` package

Data sets  
contain more  
information than  
they display

# dplyr: A package that helps transform tabular data.

---

## dplyr



```
# install.packages("dplyr")  
library(dplyr)  
?select                ?mutate  
?filter                 ?summarise  
?arrange                ?group_by
```

# Ways to access information

---

- 1 **Extract** existing variables. `select()`
- 2 **Extract** existing observations. `filter()`
- 3 **Derive** new variables  
(from existing variables) `mutate()`
- 4 **Change** the unit of analysis `summarise()`

# select()

---

storms

storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Alex	45	1009	1998-07-30
Allison	65	1005	1995-06-04
Ana	40	1013	1997-07-01
Arlene	50	1010	1999-06-13
Arthur	45	1010	1996-06-21



storm	pressure
Alberto	1007
Alex	1009
Allison	1005
Ana	1013
Arlene	1010
Arthur	1010

`select(storms, storm, pressure)`



# select()

---

storms

storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Alex	45	1009	1998-07-30
Allison	65	1005	1995-06-04
Ana	40	1013	1997-07-01
Arlene	50	1010	1999-06-13
Arthur	45	1010	1996-06-21



wind	pressure	date
110	1007	2000-08-12
45	1009	1998-07-30
65	1005	1995-06-04
40	1013	1997-07-01
50	1010	1999-06-13
45	1010	1996-06-21

```
select(storms, -storm)  
# see ?select for more
```

# select()

---

storms

storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Alex	45	1009	1998-07-30
Allison	65	1005	1995-06-04
Ana	40	1013	1997-07-01
Arlene	50	1010	1999-06-13
Arthur	45	1010	1996-06-21



wind	pressure	date
110	1007	2000-08-12
45	1009	1998-07-30
65	1005	1995-06-04
40	1013	1997-07-01
50	1010	1999-06-13
45	1010	1996-06-21

```
select(storms, wind:date)
```

```
# see ?select for more
```

# Useful select functions

---

-	Select everything but
:	Select range
contains()	Select columns whose name contains a character string
ends_with()	Select columns whose name ends with a string
everything()	Select every column
matches()	Select columns whose name matches a regular expression
num_range()	Select columns named x1, x2, x3, x4, x5
one_of()	Select columns whose names are in a group of names
starts_with()	Select columns whose name starts with a character string

\* Blue functions come in dplyr

# filter()

---

storms

storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Alex	45	1009	1998-07-30
Allison	65	1005	1995-06-04
Ana	40	1013	1997-07-01
Arlene	50	1010	1999-06-13
Arthur	45	1010	1996-06-21



storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Allison	65	1005	1995-06-04
Arlene	50	1010	1999-06-13

```
filter(storms, wind >= 50)
```

# filter()

---

storms

storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Alex	45	1009	1998-07-30
Allison	65	1005	1995-06-04
Ana	40	1013	1997-07-01
Arlene	50	1010	1999-06-13
Arthur	45	1010	1996-06-21



storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Allison	65	1005	1995-06-04

```
filter(storms, wind >= 50,  
       storm %in% c("Alberto", "Alex", "Allison"))
```

# logical tests in R

---

## ?Comparison

<	Less than
>	Greater than
==	Equal to
<=	Less than or equal to
>=	Greater than or equal to
!=	Not equal to
%in%	Group membership
is.na	Is NA
!is.na	Is not NA

## ?base::Logic

&	boolean and
	boolean or
xor	exactly or
!	not
any	any true
all	all true

# mutate()

---

storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Alex	45	1009	1998-07-30
Allison	65	1005	1995-06-04
Ana	40	1013	1997-07-01
Arlene	50	1010	1999-06-13
Arthur	45	1010	1996-06-21



storm	wind	pressure	date	ratio
Alberto	110	1007	2000-08-12	9.15
Alex	45	1009	1998-07-30	22.42
Allison	65	1005	1995-06-04	15.46
Ana	40	1013	1997-07-01	25.32
Arlene	50	1010	1999-06-13	20.20
Arthur	45	1010	1996-06-21	22.44

```
mutate(storms, ratio = pressure / wind)
```

# mutate()

---

storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Alex	45	1009	1998-07-30
Allison	65	1005	1995-06-04
Ana	40	1013	1997-07-01
Arlene	50	1010	1999-06-13
Arthur	45	1010	1996-06-21



storm	wind	pressure	date	ratio	inverse
Alberto	110	1007	2000-08-12	9.15	0.11
Alex	45	1009	1998-07-30	22.42	0.04
Allison	65	1005	1995-06-04	15.46	0.06
Ana	40	1013	1997-07-01	25.32	0.04
Arlene	50	1010	1999-06-13	20.20	0.05
Arthur	45	1010	1996-06-21	22.44	0.04

```
mutate(storms, ratio = pressure / wind, inverse = ratio^-1)
```



# Useful mutate functions

---

\* All take a vector of values and return a vector of values

\*\* Blue functions come in dplyr

<code>pmin()</code> , <code>pmax()</code>	Element-wise min and max
<code>cummin()</code> , <code>cummax()</code>	Cumulative min and max
<code>cumsum()</code> , <code>cumprod()</code>	Cumulative sum and product
<code>between()</code>	Are values between a and b?
<code>cume_dist()</code>	Cumulative distribution of values
<code>cumall()</code> , <code>cumany()</code>	Cumulative all and any
<code>cummean()</code>	Cumulative mean
<code>lead()</code> , <code>lag()</code>	Copy with values one position
<code>ntile()</code>	Bin vector into n buckets
<code>dense_rank()</code> , <code>min_rank()</code> , <code>percent_rank()</code> , <code>row_number()</code>	Various ranking methods

# "Window" functions

---

\* All take a vector of values and return a vector of values

pmin(), pmax()

cummin(), cummax()

cumsum(), cumprod()

between()

cume\_dist()

cumall(), cumany()

cummean()

lead(), lag()

ntile()

dense\_rank(), min\_rank(),  
percent\_rank(), row\_number()

1
2
3
4
5
6

cumsum()

1
3
6
10
15
21

# summarise()

---

city	particle size	amount ( $\mu\text{g}/\text{m}^3$ )
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56



median	variance
22.5	1731.6

```
pollution %>% summarise(median = median(amount), variance = var(amount))
```

# summarise()

---

city	particle size	amount ( $\mu\text{g}/\text{m}^3$ )
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56



mean	sum	n
42	252	6

```
pollution %>% summarise(mean = mean(amount), sum = sum(amount), n = n())
```

# Useful summary functions

---

\* All take a vector of values and return a single value

\*\* Blue functions come in dplyr

min(), max()	Minimum and maximum values
mean()	Mean value
median()	Median value
sum()	Sum of values
var, sd()	Variance and standard deviation of a vector
first()	First value in a vector
last()	Last value in a vector
nth()	Nth value in a vector
n()	The number of values in a vector
n_distinct()	The number of distinct values in a vector

# "Summary" functions

---

\* All take a vector of values and return a single value

min(), max()

mean()

median()

sum()

var, sd()

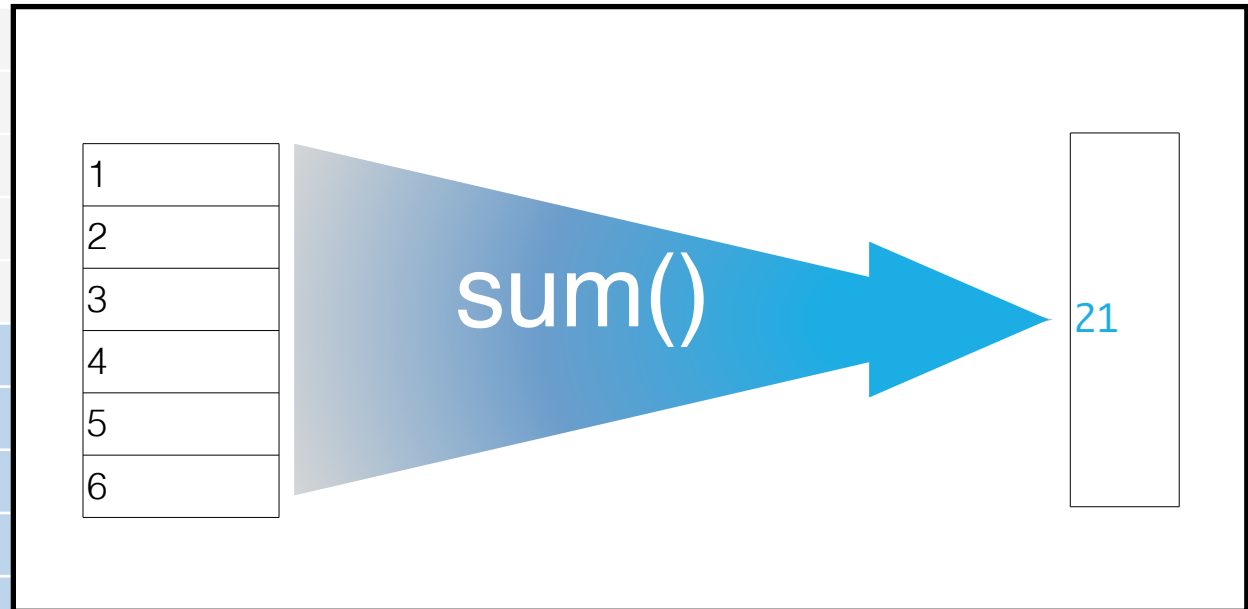
first()

last()

nth()

n()

n\_distinct()



The number of distinct values in a vector

# arrange()

---

storms

storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Alex	45	1009	1998-07-30
Allison	65	1005	1995-06-04
Ana	40	1013	1997-07-01
Arlene	50	1010	1999-06-13
Arthur	45	1010	1996-06-21



storm	wind	pressure	date
Ana	40	1013	1997-07-01
Alex	45	1009	1998-07-30
Arthur	45	1010	1996-06-21
Arlene	50	1010	1999-06-13
Allison	65	1005	1995-06-04
Alberto	110	1007	2000-08-12

`arrange(storms, wind)`

# arrange()

---

storms

storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Alex	45	1009	1998-07-30
Allison	65	1005	1995-06-04
Ana	40	1013	1997-07-01
Arlene	50	1010	1999-06-13
Arthur	45	1010	1996-06-21



storm	wind	pressure	date
Ana	40	1013	1997-07-01
Alex	45	1009	1998-07-30
Arthur	45	1010	1996-06-21
Arlene	50	1010	1999-06-13
Allison	65	1005	1995-06-04
Alberto	110	1007	2000-08-12

`arrange(storms, wind)`



# arrange()

---

storms

storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Alex	45	1009	1998-07-30
Allison	65	1005	1995-06-04
Ana	40	1013	1997-07-01
Arlene	50	1010	1999-06-13
Arthur	45	1010	1996-06-21



storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Allison	65	1005	1995-06-04
Arlene	50	1010	1999-06-13
Arthur	45	1010	1996-06-21
Alex	45	1009	1998-07-30
Ana	40	1013	1997-07-01

`arrange(storms, desc(wind))`

# arrange()

---

storms

storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Alex	45	1009	1998-07-30
Allison	65	1005	1995-06-04
Ana	40	1013	1997-07-01
Arlene	50	1010	1999-06-13
Arthur	45	1010	1996-06-21



storm	wind	pressure	date
Ana	40	1013	1997-07-01
Alex	45	1009	1998-07-30
Arthur	45	1010	1996-06-21
Arlene	50	1010	1999-06-13
Allison	65	1005	1995-06-04
Alberto	110	1007	2000-08-12

`arrange(storms, wind)`

# arrange()

---

storms

storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Alex	45	1009	1998-07-30
Allison	65	1005	1995-06-04
Ana	40	1013	1997-07-01
Arlene	50	1010	1999-06-13
Arthur	45	1010	1996-06-21



storm	wind	pressure	date
Ana	40	1013	1997-07-01
Arthur	45	1010	1996-06-21
Alex	45	1009	1998-07-30
Arlene	50	1010	1999-06-13
Allison	65	1005	1995-06-04
Alberto	110	1007	2000-08-12



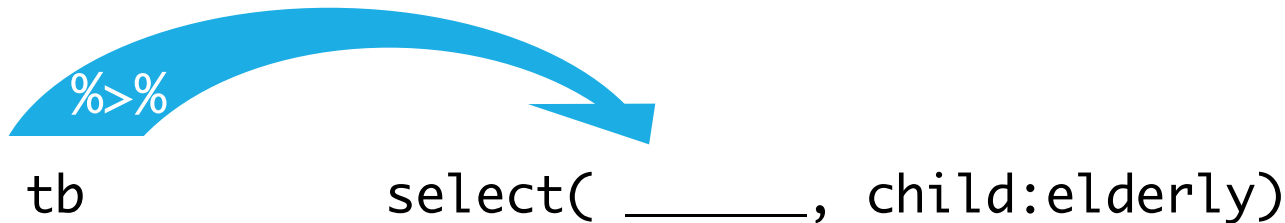
`arrange(storms, wind, date)`

# The pipe operator

`%>%`

```
library(dplyr)
```

```
select(tb, child:elderly)  
tb %>% select(child:elderly)
```



# Little Bunny FooFoo (a nursery rhyme)

---

Little bunny Foo Foo  
Went hopping through the forest  
Scooping up the field mice  
And bopping them on the head

# Little Bunny FooFoo (a nursery rhyme)

---

Little bunny Foo Foo  
Went hopping through the forest  
Scooping up the field mice  
And bopping them on the head

Using temporary objects:

```
T1=hop_through(foo_foo, forest)
T2=scoop_up(T1, field_mice)
T3=bop_on(T2, head)
```

# Little Bunny FooFoo (a nursery rhyme)

---

Little bunny Foo Foo  
Went hopping through the forest  
Scooping up the field mice  
And bopping them on the head

## Using nested function calls:

```
bop_on(  
    scoop_up(  
        hop_through(  
            foo_foo,  
            forest  
        ),  
        field_mice),  
    head)
```

# Little Bunny FooFoo (a nursery rhyme)

---

Little bunny Foo Foo  
Went hopping through the forest  
Scooping up the field mice  
And bopping them on the head

Using dplyr pipes:

```
foo_foo %>%  
  hop_through(forest) %>%  
  scoop_up(field_mice) %>%  
  bop_on(head)
```

Using pipes usually leads to more transparent code...

- No temporary objects to remember / mess up
- Reads chronologically



# select()

storms

storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Alex	45	1009	1998-07-30
Allison	65	1005	1995-06-04
Ana	40	1013	1997-07-01
Arlene	50	1010	1999-06-13
Arthur	45	1010	1996-06-21



storm	pressure
Alberto	1007
Alex	1009
Allison	1005
Ana	1013
Arlene	1010
Arthur	1010

```
select(storms, storm, pressure)
```

# select()

storms

storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Alex	45	1009	1998-07-30
Allison	65	1005	1995-06-04
Ana	40	1013	1997-07-01
Arlene	50	1010	1999-06-13
Arthur	45	1010	1996-06-21



storm	pressure
Alberto	1007
Alex	1009
Allison	1005
Ana	1013
Arlene	1010
Arthur	1010

```
storms %>% select(storm, pressure)
```

# filter()

storms

storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Alex	45	1009	1998-07-30
Allison	65	1005	1995-06-04
Ana	40	1013	1997-07-01
Arlene	50	1010	1999-06-13
Arthur	45	1010	1996-06-21



storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Allison	65	1005	1995-06-04
Arlene	50	1010	1999-06-13

```
filter(storms, wind >= 50)
```

# filter()

storms

storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Alex	45	1009	1998-07-30
Allison	65	1005	1995-06-04
Ana	40	1013	1997-07-01
Arlene	50	1010	1999-06-13
Arthur	45	1010	1996-06-21



storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Allison	65	1005	1995-06-04
Arlene	50	1010	1999-06-13

```
storms %>% filter(wind >= 50)
```

## storms

storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Alex	45	1009	1998-07-30
Allison	65	1005	1995-06-04
Ana	40	1013	1997-07-01
Arlene	50	1010	1999-06-13
Arthur	45	1010	1996-06-21



storm	pressure
Alberto	1007
Allison	1005
Arlene	1010

```
storms %>%
```

```
  filter(wind >= 50) %>%
```

```
  select(storm, pressure)
```

# mutate()

storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Alex	45	1009	1998-07-30
Allison	65	1005	1995-06-04
Ana	40	1013	1997-07-01
Arlene	50	1010	1999-06-13
Arthur	45	1010	1996-06-21



storms %>%

```
mutate(ratio = pressure / wind) %>%
```

```
select(storm, ratio)
```

# mutate()

storm	wind	pressure	date
Alberto	110	1007	2000-08-12
Alex	45	1009	1998-07-30
Allison	65	1005	1995-06-04
Ana	40	1013	1997-07-01
Arlene	50	1010	1999-06-13
Arthur	45	1010	1996-06-21



storm	ratio
Alberto	9.15
Alex	22.42
Allison	15.46
Ana	25.32
Arlene	20.20
Arthur	22.44

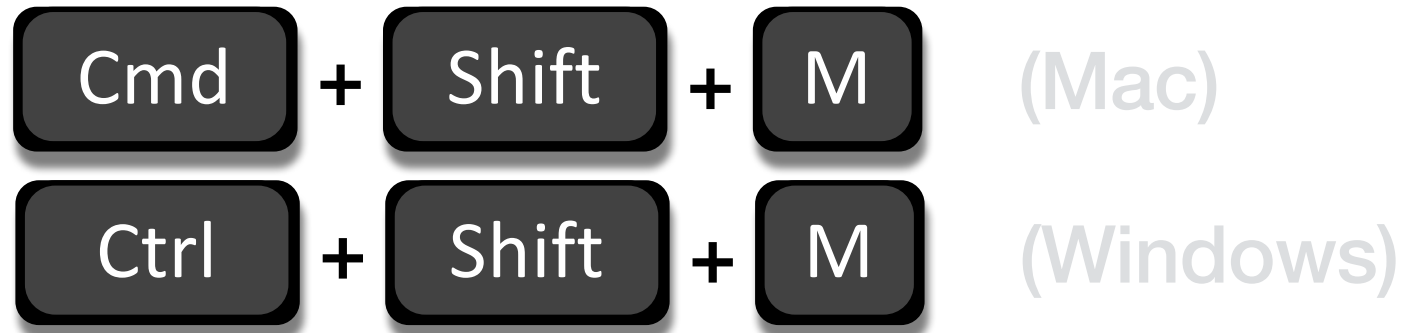
storms %>%

mutate(ratio = pressure / wind) %>%

select(storm, ratio)

# Shortcut to type %>%

---





Unit of  
analysis

---

city	particle size	amount (ug/m <sup>3</sup> )

n	sum	n

`summarize()`

---

city	particle size	amount ( $\mu\text{g}/\text{m}^3$ )
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56

mean	sum	n
42	252	6

city	particle size	amount ( $\mu\text{g}/\text{m}^3$ )
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56



mean	sum	n
18.5	37	2



19.0	38	2
------	----	---



88.5	177	2
------	-----	---

`group_by() + summarise()`

# group\_by()

---

city	particle size	amount ( $\mu\text{g}/\text{m}^3$ )
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56



city	particle size	amount ( $\mu\text{g}/\text{m}^3$ )
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56

```
pollution %>% group_by(city)
```

```
pollution %>% group_by(city)
## Source: local data frame [6 x 3]
## Groups: city
##
##   city  size amount
## 1 New York large    23
## 2 New York small   14
## 3 London large    22
## 4 London small   16
## 5 Beijing large  121
## 6 Beijing small   56
```

# group\_by() + summarise()

city	particle size	amount ( $\mu\text{g}/\text{m}^3$ )
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56

```
pollution %>% group_by(city) %>%  
  summarise(mean = mean(amount), sum = sum(amount), n = n())
```

---

city	particle size	amount ( $\mu\text{g}/\text{m}^3$ )
New York	large	23
New York	small	14



city	mean	sum	n
New York	18.5	37	2

London	large	22
London	small	16



London	19.0	38	2
--------	------	----	---

Beijing	large	121
Beijing	small	56



Beijing	88.5	177	2
---------	------	-----	---

```
pollution %>% group_by(city) %>%  
  summarise(mean = mean(amount), sum = sum(amount), n = n())
```



city	particle size	amount ( $\mu\text{g}/\text{m}^3$ )
New York	large	23
New York	small	14

London	large	22
London	small	16

Beijing	large	121
Beijing	small	56



city	mean	sum	n
New York	18.5	37	2
London	19.0	38	2
Beijing	88.5	177	2

```
pollution %>% group_by(city) %>%  
  summarise(mean = mean(amount), sum = sum(amount), n = n())
```

city	size	amount
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56



city	size	amount
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56



city	mean
New York	18.5
London	19.0
Beijing	88.5

```
pollution %>% group_by(city) %>% summarise(mean = mean(amount))
```

city	size	amount
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56



city	size	amount
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56



size	mean
large	55.3
small	28.6

```
pollution %>% group_by(size) %>% summarise(mean = mean(amount))
```

# ungroup()

city	particle size	amount ( $\mu\text{g}/\text{m}^3$ )
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56



city	particle size	amount ( $\mu\text{g}/\text{m}^3$ )
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56

```
pollution %>% ungroup()
```

# Hierarchy of information

## Larger units of analysis

country	year	sex	cases
Afghanistan	1999	female	1
Afghanistan	1999	male	1
Afghanistan	2000	female	1
Afghanistan	2000	male	1
Brazil	1999	female	2
Brazil	1999	male	2
Brazil	2000	female	2
Brazil	2000	male	2
China	1999	female	3
China	1999	male	3
China	2000	female	3
China	2000	male	3



country	year	sex	cases
Afghanistan	1999	female	1
Afghanistan	1999	male	1
Afghanistan	2000	female	1
Afghanistan	2000	male	1
Brazil	1999	female	2
Brazil	1999	male	2
Brazil	2000	female	2
Brazil	2000	male	2
China	1999	female	3
China	1999	male	3
China	2000	female	3
China	2000	male	3



country	year	cases
Afghanistan	1999	2
Afghanistan	2000	2
Brazil	1999	4
Brazil	2000	4
China	1999	6
China	1999	6



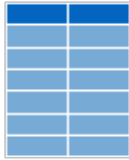
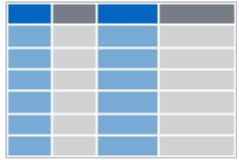
country	cases
Afghanistan	4
Brazil	8
China	12

```
tb %>%
```

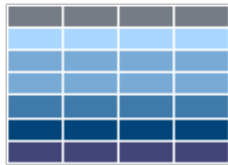
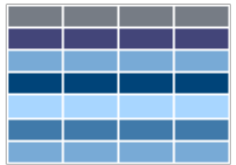
```
  group_by(country, year) %>%  
  summarise(cases = sum(cases)) %>%  
  summarise(cases = sum(cases))
```

# Recap: Information

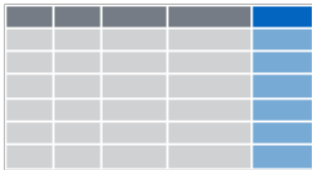
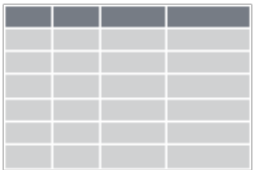
---



Extract variables and observations with `select()` and `filter()`



Arrange observations, with `arrange()`.



Make new variables, with `mutate()`.



Group observations with `group_by()` and `summarise()`.

# Joining data

# dplyr::bind\_cols()

y

x1	x2
A	1
B	2
C	3

+

z

x1	x2
B	2
C	3
D	4

=

x1	x2	x1	x2
A	1	B	2
B	2	C	3
C	3	D	4

`bind_cols(y, z)`



# dplyr::bind\_rows()

y

x1	x2
A	1
B	2
C	3

+

z

x1	x2
B	2
C	3
D	4

=

x1	x2
A	1
B	2
C	3
B	2
C	3
D	4

`bind_rows(y, z)`

# dplyr::union()

y

x1	x2
A	1
B	2
C	3

+

z

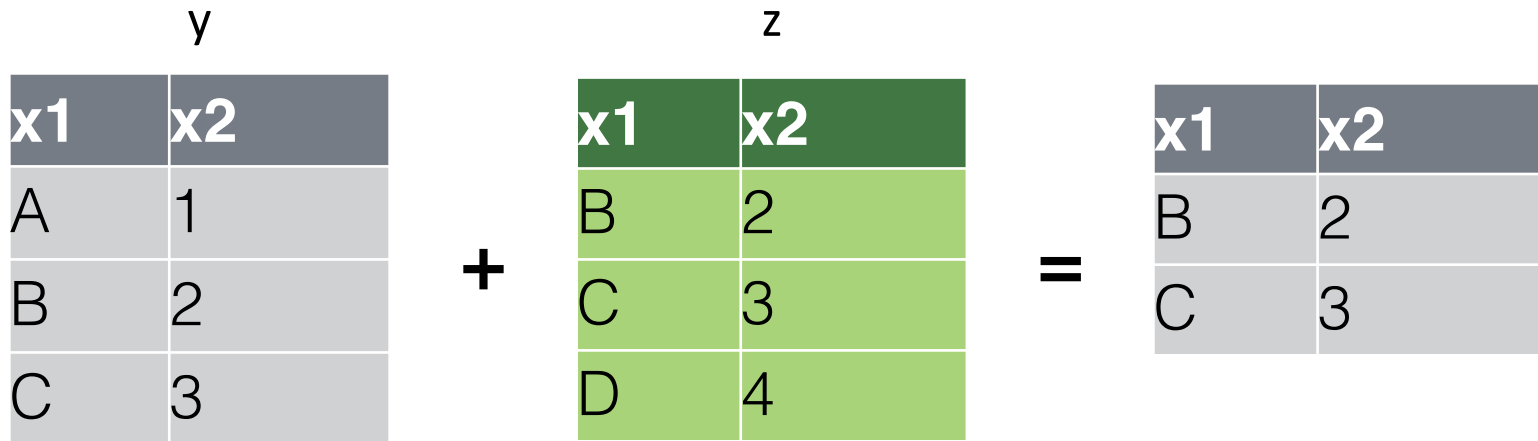
x1	x2
B	2
C	3
D	4

=

x1	x2
A	1
B	2
C	3
D	4

`union(y, z)`

# dplyr::intersect()



`intersect(y, z)`

---

## dplyr::setdiff()

y			z			y	
x1	x2		x1	x2		x1	x2
A	1	+	B	2	=	A	1
B	2		C	3		D	4
C	3		D	4			

setdiff(y, z)

*left\_join(x, y): Return all rows from x, and all columns from x and y. If there are multiple matches between x and y, all combination of the matches are returned. This is a mutating join.*

## dplyr::left\_join()

songs			artists			result		
song	name		name	plays	=	song	name	plays
Across the Universe	John	+	George	sitar		Across the Universe	John	guitar
Come Together	John		John	guitar		Come Together	John	guitar
Hello, Goodbye	Paul		Paul	bass		Hello, Goodbye	Paul	bass
Peggy Sue	Buddy		Ringo	drums		Peggy Sue	Buddy	<NA>

```
left_join(songs, artists, by = "name")
```

# dplyr::left\_join()

songs			artists			result		
song	name		name	plays	=	song	name	plays
Across the Universe	John	+	George	sitar		Across the Universe	John	guitar
Come Together	John		John	guitar		Come Together	John	guitar
Hello, Goodbye	Paul		Paul	bass		Hello, Goodbye	Paul	bass
Peggy Sue	Buddy		Ringo	drums		Peggy Sue	Buddy	<NA>

```
left_join(songs, artists, by = "name")
```

# dplyr::left\_join()

songs2

song	first	last
Across the Universe	John	Lennon
Come Together	John	Lennon
Hello, Goodbye	Paul	McCartney
Peggy Sue	Buddy	Holly

artists2

first	last	plays
George	Harrison	sitar
John	Lennon	guitar
Paul	McCartney	bass
Ringo	Starr	drums
Paul	Simon	guitar
John	Coltrane	sax

+

=

song	first	last	plays
Across the Universe	John	Lennon	guitar
Come Together	John	Lennon	guitar
Hello, Goodbye	Paul	McCartney	bass
Peggy Sue	Buddy	Holly	<NA>

```
left_join(songs2, artists2, by = c("first", "last"))
```

# dplyr::left\_join()

songs2

song	first	last
Across the Universe	John	Lennon
Come Together	John	Lennon
Hello, Goodbye	Paul	McCartney
Peggy Sue	Buddy	Holly

+

artists2

first	last	plays
George	Harrison	sitar
John	Lennon	guitar
Paul	McCartney	bass
Ringo	Starr	drums
Paul	Simon	guitar
John	Coltrane	sax

=

song	first	last	plays
Across the Universe	John	Lennon	guitar
Come Together	John	Lennon	guitar
Hello, Goodbye	Paul	McCartney	bass
Peggy Sue	Buddy	Holly	<NA>

```
left_join(songs2, artists2, by = c("first", "last"))
```



# left\_join()

songs

song	name
Across the Universe	John
Come Together	John
Hello, Goodbye	Paul
Peggy Sue	Buddy

+

artists

name	plays
George	sitar
John	guitar
Paul	bass
Ringo	drums

=

song	name	plays
Across the Universe	John	guitar
Come Together	John	guitar
Hello, Goodbye	Paul	bass
Peggy Sue	Buddy	<NA>

```
left_join(songs, artists, by = "name")
```

`inner_join(x, y)`: Return all rows from `x` where there are matching values in `y`, and all columns from `x` and `y`. If there are multiple matches between `x` and `y`, all combination of the matches are returned. This is a mutating join.

## `inner_join()`

songs			artists			result		
song	name		name	plays	=	song	name	plays
Across the Universe	John	+	George	sitar	=	Across the Universe	John	guitar
Come Together	John		John	guitar		Come Together	John	guitar
Hello, Goodbye	Paul		Paul	bass		Hello, Goodbye	Paul	bass
Peggy Sue	Buddy		Ringo	drums				

```
inner_join(songs, artists, by = "name")
```

`semi_join(x, y)`: Return all rows from `x` where there are matching values in `y`, keeping just columns from `x`. A semi join differs from an inner join because an inner join will return one row of `x` for each matching row of `y`, where a semi join will never duplicate rows of `x`. This is a filtering join.

## semi\_join()

songs			artists			songs	
song	name		name	plays	=	song	name
Across the Universe	John	+	George	sitar		Across the Universe	John
Come Together	John		John	guitar		Come Together	John
Hello, Goodbye	Paul		Paul	bass		Hello, Goodbye	Paul
Peggy Sue	Buddy		Ringo	drums			

```
semi_join(songs, artists, by = "name")
```

`anti_join(x, y)`: Return all rows from x where there are not matching values in y, keeping just columns from x. This is a filtering join.

## anti\_join()

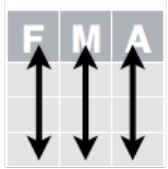
songs			artists			songs	
song	name		name	plays	=	song	name
Across the Universe	John	+	George	sitar	=	Peggy Sue	Buddy
Come Together	John		John	guitar			
Hello, Goodbye	Paul		Paul	bass			
Peggy Sue	Buddy		Ringo	drums			

`anti_join(songs, artists, by = "name")`

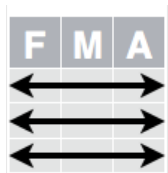
Great Join Cheatsheet: [http://stat545.com/bit001\\_dplyr-cheatsheet.html](http://stat545.com/bit001_dplyr-cheatsheet.html)

# Recap: Best format for analysis

---



**Variables** in columns



**Observations** in rows



Separate **all variables** *implied by law, formula or goal*



**Unit of analysis matches** the unit of analysis *implied by law, formula or goal*



**Single** table

# Interactive Exercises