ANT 6973: DATA VISUALIZATION AND EXPLORATION

DATA MANIPULATION, PART 1

PACKAGES FOR WORKING WITH DATA



tidyr

Both are part of core





dplyr

library("tidyverse")

Transform Data with



ACTIVITIES

• Open babynames-manip.Rmd and follow along.

babynames



Names of male and female babies born in the US from 1880 to 2015. 1.8M rows.

```
# install.packages("babynames")
library("babynames")
```

babynames

<pre>cdbl></pre>	n <int></int>	name <chr></chr>	<chr></chr>	<dbl></dbl>
7.238433e-02	7065	Mary	F	1880
2.667923e-02	2604	Anna	F	1880
2.052170e-02	2003	Emma	F	1880
1.986599e-02	1939	Elizabeth	F	1880
1.788861e-02	1746	Minnie	F	1880
1.616737e-02	1578	Margaret	F	1880
1.508135e-02	1472	Ida	F	1880
1.448711e-02	1414	Alice	F	1880
1.352404e-02	1320	Bertha	F	1880
1.319618e-02	1288	Sarah	F	1880

How to isolate?

year	sex	name	n	prop
1880	M	John	9655	0.0815
1880	M	William	9532	0.0805
1880	M	James	5927	0.0501
1880	M	Charles	5348	0.0451
1880	M	Fernando	8	0.0001
1881	M	John	8769	0.081
1881	M	William	8524	0.0787
1881	M	James	5442	0.0503
1881	M	Charles	4664	0.0431
1881	M	Fernando	6	0.0001
1881	M	Gideon	7	0.0001

year	sex	name	n	prop
1880	M	Fernando	8	0.0001
1881	M	Fernando	6	0.0001
• • •	• • •	Fernando	• • •	• • •

dplyr

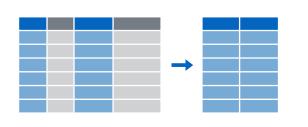


A package that transforms data.

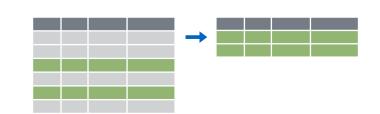
dplyr implements a *grammar* for transforming tabular data.

SINGLE TABLE VERBS





Extract variables with select()



Extract cases with filter()



Arrange cases with arrange()



Make new variables with mutate()



Make tables of summaries with summarise()

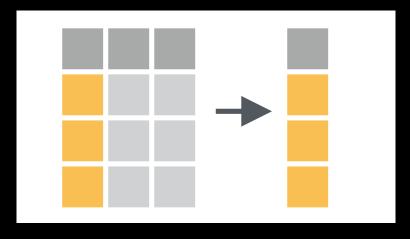
along with group_by()



select()



Extract columns by name.



```
select(.data, ...)
```

data frame to transform

name(s) of columns to extract (or a select helper function)

Select multiple variables by separating them with commas



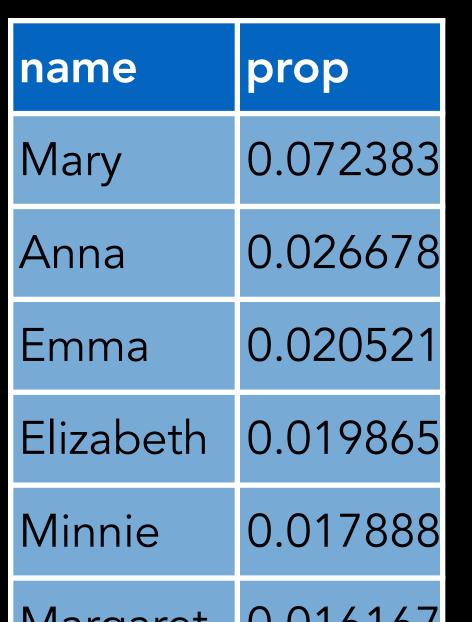
Extract columns by name.

select(babynames, name, prop)

Note how the order of columns is determined by the order of inputs

babynames

year	sex	name	n	prop
1880	F	Mary	7065	0.072383
1880	F	Anna	2604	0.026678
1880	F	Emma	2003	0.020521
1880	F	Elizabeth	1939	0.019865
1880	F	Minnie	1746	0.017888
1000	_	Margaret	1570	0.014147





Extract columns by name.



babynames %>% # Same but with pipe
select(name, prop)

babynames

year	sex	name	n	prop
1880	F	Mary	7065	0.072383
1880	F	Anna	2604	0.026678
1880	F	Emma	2003	0.020521
1880	F	Elizabeth	1939	0.019865
1880	F	Minnie	1746	0.017888
1000		Margaret	1570	0.014147



ACTIVITY 1

Alter the code to select just the n column:

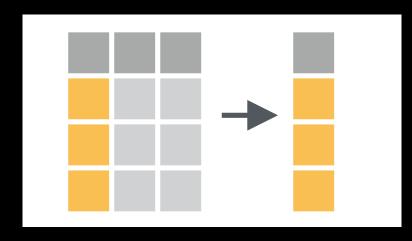
```
select(babynames, name, prop)
```

select(babynames, n)

```
<int>
# 1 7065
# 2 2604
# 3 2003
# 4 1939
# 5 1746
```



You can rename on the fly



select(babynames, name, popularity = prop)

babynames

year	sex	name	n	prop
1880	F	Mary	7065	0.072383
1880	F	Anna	2604	0.026678
1880	F	Emma	2003	0.020521
1880	F	Elizabeth	1939	0.019865
1880	F	Minnie	1746	0.017888
1000	F	Margaret	1570	0.014147



name	popularity
Mary	0.07238359
Anna	0.02667896
Emma	0.02052149
Elizabeth	0.01986579
Minnie	0.01788843
Margaret	0.0141472

: Select range of columns

```
select(penguins, species:sex)
```

! or – Negate selection (select every column but)

```
select(penguins, !c(island, sex))
```

starts_with() Select columns that start with...

```
select(penguins, starts_with("bill"))
```

ends_with() Select columns that end with...

```
select(penguins, ends_with("_mm"))
```



contains () Select columns whose names contain...

```
select(penguins, contains("length"))
```

matches() Select columns whose names match regular expression

```
select(penguins, matches("^.{4}$"))
```

any_of() Select columns whose names are in a set

```
select(penguins, any_of(c("year", "Year", "yEaR")))
```

num_range() Select columns named in prefix, number style

```
select(billboard, num_range("wk", 10:15))
```



everything() Select all variables

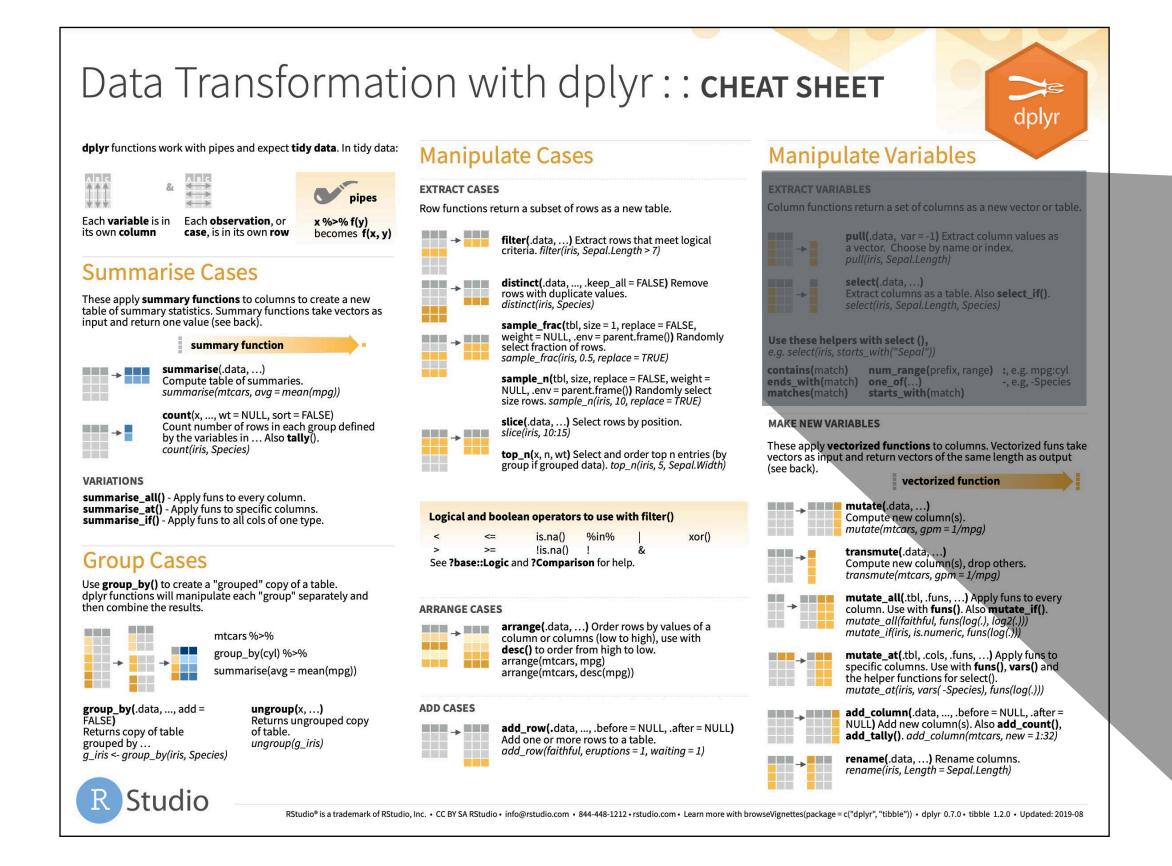
```
select(penguins, everything())
```

where () Select columns for which function returns TRUE

```
select(penguins, where(is.numeric) & !year)
```







EXTRACT VARIABLES Column functions return a set of columns as a new vector or table pull(.data, var = -1) Extract column values as a vector. Choose by name or index. pull(iris, Sepal.Length) select(.data, ...) Extract columns as a table. Also select_if(). select(iris, Sepal.Length, Species) Use these helpers with select (), e.g. select(iris, starts_with("Sepal")) num_range(prefix, range) :, e.g. mpg:cyl contains(match) ends_with(match) one_of(...) -, e.g, -Species starts_with(match) matches(match)

QUIZ

Which of these is NOT a way to select the name and n columns together?

```
select(babynames, -c(year, sex, prop))
select(babynames, name:n)
select(babynames, starts_with("n"))
select(babynames, ends_with("n"))
```

QUIZ

Which of these is NOT a way to select the name and n columns together?

```
select(babynames, -c(year, sex, prop))
select(babynames, name:n)
select(babynames, starts_with("n"))
select(babynames, ends_with("n"))
```



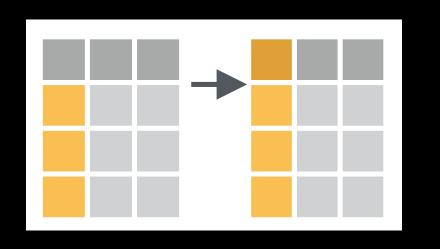
SELECT()'S COUSINS

For manipulating variables/columns

RENAME()



Rename columns without extracting.



rename(babynames, popularity = prop)

year	sex	name	n	prop
1880	F	Mary	7065	0.072383
1880	F	Anna	2604	0.026678
1880	F	Emma	2003	0.020521
1880	F	Elizabeth	1939	0.019865
1880	F	Minnie	1746	0.017888
1000	Г	N /1	1570	0 01/1/7

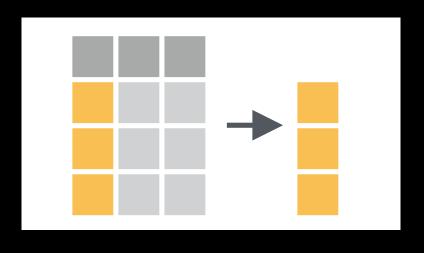


year	sex	name	n	popularity
1880	F	Mary	7065	0.07238359
1880	F	Anna	2604	0.02667896
1880	F	Emma	2003	0.02052149
1880	F	Elizabeth	1939	0.01986579
1880	F	Minnie	1746	0.01788843
1000	_	N /1	1570	0 01/1/72

PULL()



Extract column values only.



```
pull(babynames, n)
# Does same thing as babynames$n
```

year	sex	name	n	prop
1880	F	Mary	7065	0.072383
1880	F	Anna	2604	0.026678
1880	F	Emma	2003	0.020521
1880	F	Elizabeth	1939	0.019865
1880	F	Minnie	1746	0.017888
1000	Г	N /1	1570	0 01/1/7



[1] 7065 2604 2003 1939 ...



filter()



Extract rows that meet logical criteria.

```
data frame to transform

one or more logical tests (filter returns each row for which the test is TRUE)
```



Extract rows that meet logical criteria.

filter(babynames, name == "Fernando")

babynames

year	sex	name	n	prop
1880	M	John	9655	0.0815
1880	M	William	9532	0.0805
1880	M	James	5927	0.0501
1880	M	Charles	5348	0.0451
1880	M	Fernando	8	0.0001
1881	M	John	8769	0.081



year	sex	name	n	prop
1880	M	Fernando	8	0.0001
1881	M	Fernando	6	0.0001
• • •	• • •	Fernando	• • •	• • •



Extract rows that meet logical criteria.

filter(babynames, name == "Fernando")

babynames

year	sex	name	n	prop
1880	M	John	9655	0.0815
1880	M	William	9532	0.0805
1880	M	James	5927	0.0501
1880	M	Charles	5348	0.0451
1880	M	Fernando	8	0.0001
1881	M	John	8769	0.081

= sets
(returns nothing)
== tests if equal
(returns TRUE or

BASIC R: OPERATORS

COMPARISON OPERATORS

OPERATOR	DESCRIPTION
x < y	LESS THAN
x <= y	LESS THAN OR EQUAL TO
x > y	GREATER THAN
x >= y	GREATER THAN OR EQUAL TO
x == y	EXACTLY EQUAL TO
x != y	NOT EQUAL TO
x %in% y	GROUP MEMBERSHIP
is.na(x)	IS NA

?Comparison

ACTIVITY 2

- See if you can use the logical operators to manipulate babynames to show:
 - 1. All of the names where prop is greater than or equal to 0.08
 - 2. All of the children named "Sea"
 - 3. All of the names that have a missing value for **n** (Hint: this should return an empty data set).

filter(babynames, prop >= 0.08)

filter(babynames, name == "Sea")

```
# year sex name n prop
# 1 1982 F Sea 5 2.756771e-06
# 2 1985 M Sea 6 3.119547e-06
# 3 1986 M Sea 5 2.603512e-06
# 4 1998 F Sea 5 2.580377e-06
```

filter(babynames, is.na(n))

0 rows

TWO COMMON MISTAKES



1. Using = instead of ==

```
filter(babynames, name = "Sea")
filter(babynames, name == "Sea")
```

2. Forgetting quotes

```
filter(babynames, name == Sea)
filter(babynames, name == "Sea")
```

Multiple criteria to meet can be separated with commas...



Extract rows that meet every logical criterion.

filter(babynames, name == "Fernando", year == 1880)

babynames

year	sex	name	n	prop
1880	M	John	9655	0.0815
1880	M	William	9532	0.0805
1880	M	James	5927	0.0501
1880	M	Charles	5348	0.0451
1880	Μ	Fernando	8	0.0001
1881	M	John	8769	0.081



year	sex	name	n	prop
1880	M	Fernando	8	0.0001

... which is equivalent to using the & operator



Extract rows that meet every logical criterion.

filter(babynames, name == "Fernando" & year == 1880)

babynames

year	sex	name	n	prop
1880	M	John	9655	0.0815
1880	M	William	9532	0.0805
1880	M	James	5927	0.0501
1880	M	Charles	5348	0.0451
1880	M	Fernando	8	0.0001
1881	M	John	8769	0.081



year	sex	name	n	prop
1880	M	Fernando	8	0.0001

BASIC R: OPERATORS

LOGICAL OPERATORS

OPERATOR	DESCRIPTION
! X	NOT
x	OR
x & y	AND
xor(x, y)	EXACTLY OR (FALSE IF BOTH ARE TRUE)

?base::Logic

ACTIVITY 3

- Use logical operators to alter the code below to return only the rows that contain:
 - 1. Girls named Sea
 - 2. Names that were used by exactly 5 or 6 children in 1880
 - 3. Names that are one of Acura, Lexus, or Yugo

```
filter(babynames, name == "Sea" | name == "Anemone")
```

TWO MORE COMMON MISTAKES



3. Collapsing multiple tests into one

```
filter(babynames, 10 < n < 20)
filter(babynames, 10 < n, n < 20)</pre>
```

4. Stringing together many tests (when you could use %in%)

```
filter(babynames, n == 5 \mid n == 6 \mid n == 7 \mid n == 8)
filter(babynames, n \% in\% c(5, 6, 7, 8))
```



FILTER()'S COUSINS

For manipulating cases/rows

SLICE()

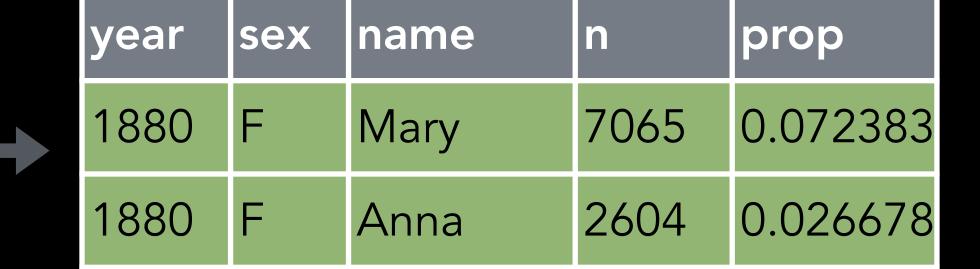


Extract rows based on position rather than criteria.

slice(babynames, 1:2)

babynames

year	sex	name	n	prop
1880	F	Mary	7065	0.072383
1880	F	Anna	2604	0.026678
1880	F	Emma	2003	0.020521
1880	F	Elizabeth	1939	0.019865
1880	F	Minnie	1746	0.017888
1880	F	Margaret	1578	0.016167



SLICE() — USEFUL VARIANTS



Variable to order by

slice_min() Subset rows v th lowest values of variable

```
slice_min(penguins, bill_length_mm, n = 5)
```

Number of rows to extract

slice_max() Subset rows with highest values of variable

```
slice_max(penguins, bill_length_mm, n = 5)
```

slice_sample() Subset rows randomly

```
slice_sample(penguins, n = 5)
```

SLICE() – USEFUL VARIANTS



slice_min() Subset rows with lowest values of variable

Watch for ties!

```
slice_min(penguins, bill_length_mm, prop = 0.1)
```

Or proportion of rows to extract

slice_max() Subset rows with highest values of variable

```
slice_max(penguins, bill_length_mm, prop = 0.1)
```

slice_sample() Subset rows randomly

```
slice_sample(penguins, prop = 0.1)
```

slice_min(penguins, body_mass_g, n = 5) species island bill_length_mm bill_depth_mm flipper_length_mm body_mass_g sex year <int> <fct> <int> <fct> <fct> <dbl> <dbl> <int> 1 Gentoo Biscoe 49.2 6300 male 2007 15.2 221 2 Gentoo 6050 male Biscoe 59.6 230 17 2007 3 Gentoo Biscoe 51.1 16.3 220 6000 male 2008 4 Gentoo Biscoe 48.8 6000 male 16.2 222 2009 Biscoe 5 Gentoo 16.4 5950 male 45.2 2008 223 6 Gentoo Biscoe 49.8 15.9 5950 male 2009 229

slice_	min(pe	nguins, bod	ly_mass_g,	$n = 5$, with_t	ies = FAL	SE)	
species	sisland	bill_length_mm	bill_depth_mm	flipper_length_mm	body_mass_g	sex	year
<fct></fct>	<fct></fct>	<dbl></dbl>	<dbl></dbl>	<int></int>	<int></int>	<fct></fct>	<int></int>
1 Gentoo	Biscoe	49.2	15.2	221	6300	male	2007
2 Gentoo	Biscoe	59.6	17	230	6050	male	2007
3 Gentoo	Biscoe	51.1	16.3	220	6000	male	2008
4 Gentoo	Biscoe	48.8	16.2	222	6000	male	2009
5 Gentoo	Biscoe	45.2	16.4	223	5950	male	2008

DISTINCT()



Select only unique/distinct row

When using multiple variables to determine uniqueness, separate with commas

```
distinct(penguins, species, body_mass_g)
```

distinct(penguins, species, body_mass_g, .keep_all = TRUE)

Keep all variables in the data? By default, only the columns used to determine uniqueness are retained



arrange()

ARRANGE()



Order rows from smallest to largest values.

```
arrange(.data, ...)
```

data frame to transform

one or more columns to order by (additional columns will be used as tie breakers)

ARRANGE()



Order rows from smallest to largest values.

arrange(babynames, n)

babynames

year	sex	name	n	prop
1880	М	John	9655	0.0815
1880	M	William	9532	0.0805
1880	M	James	5927	0.0501
1880	M	Charles	5348	0.0451
1880	M	Fernando	8	0.0001
1881	M	John	8769	0.081



year	sex	name	n	prop
1880	M	Fernando	8	0.0001
1880	M	Charles	5348	0.0451
1880	M	James	5927	0.0501
1881	M	John	8769	0.081
1880	M	William	9532	0.0805
1880	M	John	9655	0.0815

ACTIVITY 4

- Arrange babynames by n, and add name as a second (tie breaking)
 variable by which to arrange.
- Can you tell what the smallest value of n is?

arrange(babynames, n, name)

```
year
            sex
                       name
                                           prop
     2007
               M
                       Aaban
                                 5 2.259872e-06
               M
  2
     2007
                                 5 2.259872e-06
                      Aareon
                      Aaris
  3 2007
                                 5 2.259872e-06
  4 2007
                         Abd
                                 5 2.259872e-06
#
                  Abdulazeez
  5
     2007
                                 5 2.259872e-06
                  Abdulhadi
     2007
                                 5 2.259872e-06
     2007
                 Abdulhamid
                                 5 2.259872e-06
                 Abdulkadir
#
     2007
  8
                                 5 2.259872e-06
               M Abdulraheem
  9
     2007
                                 5 2.259872e-06
# 10
               M Abdulrahim
     2007
                                 5 2.259872e-06
        # ... with 1,858,679 more rows
```

DESC()



Changes ordering to largest to smallest.

arrange(babynames, desc(n))

babynames

year	sex	name	n	prop
1880	Μ	John	9655	0.0815
1880	M	William	9532	0.0805
1880	M	James	5927	0.0501
1880	M	Charles	5348	0.0451
1880	M	Fernando	8	0.0001
1881	M	John	8769	0.081



year	sex	name	n	prop
1880	М	John	9655	0.0815
1880	M	William	9532	0.0805
1881	M	John	8769	0.081
1880	M	James	5927	0.0501
1880	M	Charles	5348	0.0451
1880	M	Fernando	8	0.0001

ACTIVITY 5

- Use desc() to find the names with the highest prop.
- Then, use desc() to find the names with the highest n.

arrange(babynames, desc(prop))

```
#
      year
             sex
                     name
                                       prop
                           9655 0.08154630
      1880
                     John
                           8769 0.08098299
      1881
                 William
#
      1880
                           9531 0.08049899
                     John
                           8894 0.07907324
      1883
                M William
#
                           8524 0.07872038
      1881
                     John
      1882
                           9557 0.07831617
                     John
                           9388 0.07648751
      1884
                 William
                           9298 0.07619375
      1882
                     John
                           9026 0.07582198
      1886
                           8756 0.07551791
#
  10
      1885
                     John
     # ... with 1,858,679 more rows
```

arrange(babynames, desc(n))

```
year
           sex
                   name
                                     prop
                  Linda 99680 0.05483609
    1947
                  Linda 96211 0.05521159
    1948
    1947
                 James 94763 0.05102057
             M Michael 92726 0.04238659
    1957
                 Robert 91646 0.04934237
    1947
                  Linda 91010 0.05184281
    1949
    1956
             M Michael 90623 0.04225479
             M Michael 90517
    1958
                              0.04203881
                  James 88588 0.04969679
    1948
10
             M Michael 88493 0.04279403
    1954
   # ... with 1,858,679 more rows
```

PIPES



Here's where they become really useful!



```
boys_2015 <- filter(babynames, year == 2015, sex == "M")
boys_2015 <- select(boys_2015, name, n)
boys_2015 <- arrange(boys_2015, desc(n))
boys_2015</pre>
```

- 1. Filter babynames to just boys born in 2015
- 2. Select the name and n columns from the result
- 3. Arrange those columns so that the most popular names appear near the top.



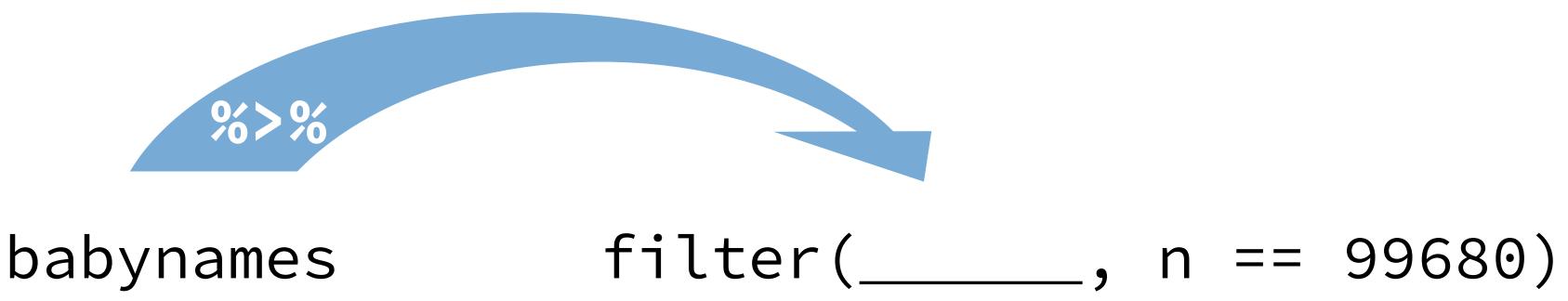
```
boys_2015 <- filter(babynames, year == 2015, sex == "M")
boys_2015 <- select(boys_2015, name, n)
boys_2015 <- arrange(boys_2015, desc(n))
boys_2015</pre>
```



```
arrange(select(filter(babynames, year == 2015,
    sex == "M"), name, n), desc(n))
```

RECALL: THE PIPE OPERATOR %>%





Passes result on left into first argument of function on right. So, for example, these do the same thing.

```
filter(babynames, n == 99680)
babynames %>% filter(n == 99680)
```



```
babynames
boys_2015 <- filter(babynames, year == 2015, sex == "M")
boys_2015 <- select(boys_2015, name, n)
boys_2015 <- arrange(boys_2015, desc(n))
boys_2015</pre>
```

```
babynames %>%
  filter(year == 2015, sex == "M") %>%
  select(name, n) %>%
  arrange(desc(n))
```

PIPES



Shortcut to type %>%

ACTIVITY 6

- Use %>% to write the following sequence of functions:
 - 1. Filter babynames to just the girls that were born in 2015
 - 2. Select the name and n columns
 - 3. Arrange the results so that the most popular names are near the top.

```
babynames %>%
  filter(year == 2015, sex == "F") %>%
  select(name, n) %>%
  arrange(desc(n))
                  # name n
                  # 1 Emma 20355
                  # 2 Olivia 19553
                  # 3 Sophia 17327
                  # 4 Ava 16286
                  # 5 Isabella 15504
                  # 6 Mia 14820
                  # 7 Abigail 12311
                  # 8 Emily 11727
                  # 9 Charlotte 11332
                  # 10 Harper 10241
                # ... with 18,983 more rows
```



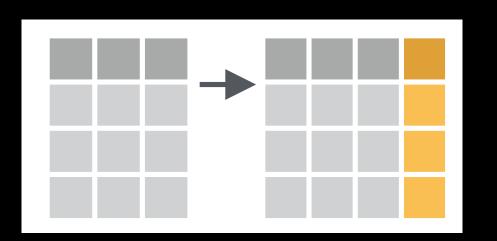
mutate()

MUTATE()

You can do basic arithmetic



Create new columns.



babynames %>%

mutate(percent = prop * 100)

babynames

year	sex	name	n	prop
1880	F	Mary	7065	0.072383
1880	F	Anna	2604	0.026678
1880	F	Emma	2003	0.020521
1880	F	Elizabeth	1939	0.019865
1880	F	Minnie	1746	0.017888

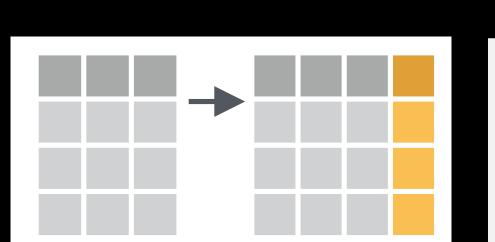


year	sex	name	n	prop	percent
1880	F	Mary	7065	0.072383	7.238359
1880	F	Anna	2604	0.026678	2.667896
1880	F	Emma	2003	0.020521	2.052149
1880	F	Elizabeth	1939	0.019865	1.986579
1880	F	Minnie	1746	0.017888	1.788843

MUTATE()



Create new columns.



babynames %>%

mutate(percent = round(prop * 100, 2))

babynames

year	sex	name	n	prop
1880	F	Mary	7065	0.072383
1880	F	Anna	2604	0.026678
1880	F	Emma	2003	0.020521
1880	F	Elizabeth	1939	0.019865
1880	F	Minnie	1746	0.017888



year	sex	name	n	prop	percent
1880	F	Mary	7065	0.072383	7.24
1880	F	Anna	2604	0.026678	2.67
1880	F	Emma	2003	0.020521	2.05
1880	F	Elizabeth	1939	0.019865	1.99
1880	F	Minnie	1746	0.017888	1.79

Operations can be

nested

ASIDE: ROUND()

Round a number to a specified number of decimal digits (0 by default)

```
x <- c(1.2, 1/3, 10.01)
round(x)
[1] 1 0 10
round(x, digits = 2)
[1] 1.20 0.33 10.01</pre>
```

When using a function in mutate() the argument will be a column name from the data

```
babynames %>%
mutate(percent = round(prop * 100, 2))
```

MUTATE()

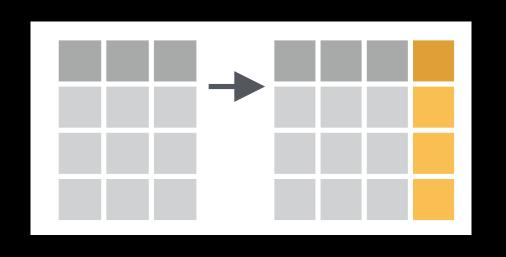


Newly created variables are

available immediately for

further manipulation

Create new columns.



```
babynames %>%
```

babynames

sex	name	n	prop
F	Mary	7065	0.072383
F	Anna	2604	0.026678
F	Emma	2003	0.020521
F	Elizabeth	1939	0.019865
F	Minnie	1746	0.017888
	F F	F Mary F Anna F Emma F Elizabeth	F Mary 7065 F Anna 2604 F Emma 2003 F Elizabeth 1939



year	sex	name	n	prop	percent	pcnt_rnd
1880	F	Mary	7065	0.072383	7.238359	7.2
1880	F	Anna	2604	0.026678	2.667896	2.7
1880	F	Emma	2003	0.020521	2.052149	2.1
1880	F	Elizabeth	1939	0.019865	1.986579	2
1880	F	Minnie	1746	0.017888	1.788843	1.8

Vector Functions

TO USE WITH MUTATE ()

mutate() and transmute() apply vectorized functions to columns to create new columns. Vectorized functions take vectors as input and return vectors of the same length as output.

vectorized function



OFFSETS

dplyr::lag() - Offset elements by 1 dplyr::lead() - Offset elements by -1

CUMULATIVE AGGREGATES

RANKINGS

dplyr::cume_dist() - Proportion of all values <=
dplyr::dense_rank() - rank w ties = min, no gaps
dplyr::min_rank() - rank with ties = min
dplyr::ntile() - bins into n bins
dplyr::percent_rank() - min_rank scaled to [0,1]
dplyr::row_number() - rank with ties = "first"</pre>

MATH

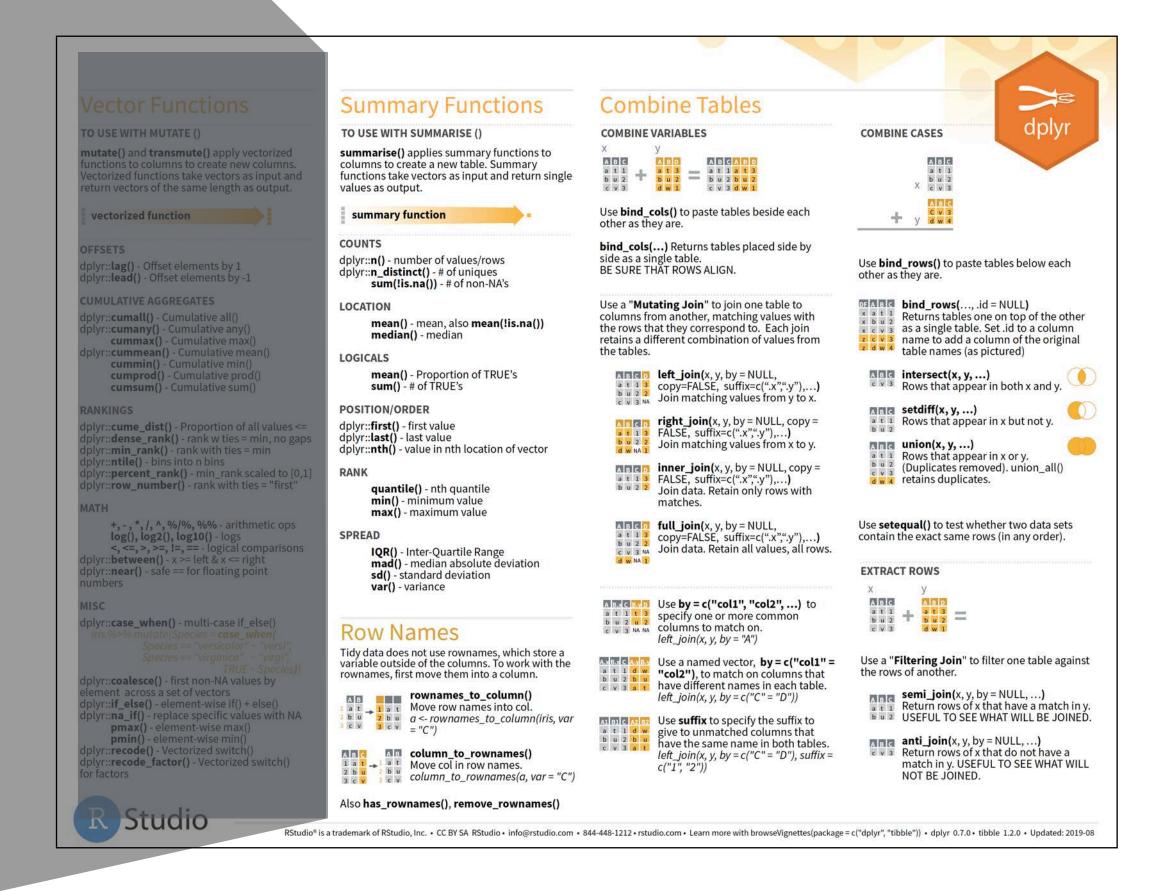
for factors

+,-,*,/,^,%/%, %% - arithmetic ops log(), log2(), log10() - logs <, <=, >, >=, !=, == - logical comparisons dplyr::between() - x >= left & x <= right dplyr::near() - safe == for floating point numbers

Vectorized functions

Take a vector as input.

Return a vector of the *same length* as output.



min_rank()



A convenient ranking function (ties share the lowest rank)

```
min_rank(c(50, 100, 1000))
# [1] 1 2 3
```

```
min_rank(c(50, 100, 100, 1000))
# [1] 1 2 2 4
```

```
min_rank(desc(c(50, 100, 1000)))
# [1] 3 2 1
```

ACTIVITY 7

- Use mutate() and min_rank() to rank each row in babynames from largest n to smallest n, and filter to extract the top 10 ranking names in the entire data set.
- Do the same but use prop rather than n to rank the names.
- What differences do you see in these rankings?
- How could you accomplish the same thing using some version of slice()?


```
rank
   <dbl> <chr> <int> <dbl> <int> <dbl> <int>
   1947 F
               Linda
                       99686 0.0548
                                        3
               James
                       94756 0.0510
   1947 M
              Robert
                                        5
   1947 M
                       91642 0.0493
               Linda
   1948 F
                       96209 0.0552
                                        9
   1948 M
               James
                       88588 0.0497
               Linda
   1949 F
                       91016 0.0518
                                        6
               Michael 88514 0.0428
   1954 M
                                       10
               Michael 90620 0.0423
   1956 M
   1957 M
               Michael 92695 0.0424
               Michael 90520 0.0420
   1958 M
                                        8
10
```

1880 M

1881 M

1881 M

1882 M

1882 M

1883 M

1884 M

1885 M

1886 M

10

William

John

William

John

William

John

John

John

John

3

2

5

6

8

4

10

9

9532 0.0805

8769 0.0810

8524 0.0787

9557 0.0783

9298 0.0762

8894 0.0791

9388 0.0765

8756 0.0755

9026 0.0758

babynames %>% slice_max(prop, n = 10)

```
year sex
              name
                               prop
   <dbl> <chr> <chr>
                       <int>
                              <dbl>
   1880 M
              John
                       9655 0.0815
   1881 M
              John
                       8769 0.0810
              William
   1880 M
                       9532 0.0805
   1883 M
              John
                       8894 0.0791
              William
   1881 M
                       8524 0.0787
   1882 M
               John
                       9557 0.0783
   1884 M
              John
                       9388 0.0765
   1882 M
              William
                       9298 0.0762
   1886 M
               John
                       9026 0.0758
              John 8756 0.0755
   1885 M
10
```

lead() and lag()



Find the "next" or "previous" values.

```
1:10
[1] 1 2 3 4 5 6 7 8 9 10
lead(1:10)
[1] 2 3 4 5 6 7 8 9 10 NA
lag(1:10)
[1] NA 1 2 3 4 5 6 7 8 9
```

ACTIVITY 7

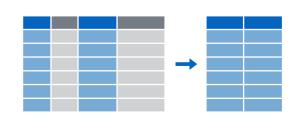
 Use mutate() and lag() to calculate year-to-year change in popularity (using the n column), then arrange the data set by this new variable with the year of greatest increase at the top. babynames %>%

filter(name == "Fernando" & sex == "M") %>%
mutate(n_change = n - lag(n)) %>%
arrange(desc(n))

	year	sex	name	n	prop	n_change
	<dbl></dbl>	<chr></chr>	<chr></chr>	<int></int>	<dbl></dbl>	<int></int>
1	1997	M	Fernando	2315	0.00116	439
2	2000	M	Fernando	2601	0.00125	426
3	2006	M	Fernando	2758	0.00126	318
4	1977	M	Fernando	1135	0.000664	249
5	1989	M	Fernando	1581	0.000754	221
6	1991	M	Fernando	1897	0.000895	167
7	1988	M	Fernando	1360	0.000680	157
8	1990	M	Fernando	1730	0.000804	149
9	1980	M	Fernando	1264	0.000681	134
10	2003	M	Fernando	2551	0.00121	128

RECAP: SINGLE TABLE VERBS

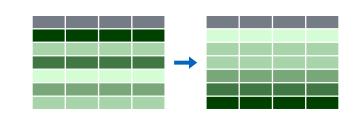




Extract variables with select()



Extract cases with filter()



Arrange cases, with arrange().



Make new variables, with mutate().



Make tables of summaries with summarise().

THE GLUE PACKAGE



```
library("glue")
my_name <- "Fernando"</pre>
glue("My name is {my_name}.")
my_nums <- 1:5
glue("This is row {my_nums}.")
glue("The price is now ${1:5}.00 higher.")
glue("The average of 10, 14, and 33 is: {mean(c(10, 14, 33))}")
name_sequence <- c("first name", "middle name", "last name")</pre>
my_names <- c("Fernando", "Alonso", "Campos")</pre>
glue("My {name_sequence} is {my_names}.")
```