

Data Transformation

Session 4
Patrick Mathias
September 26, 2021

September 26, 2021	Session	Instructor
8:30 am - 8:45 am	Course Introduction	Patrick Mathias
8:45 am – 9:30 am	Intro to R and Reproducible Reporting	Joseph Rudolf
9:45 am - 10:30 am	Coding Basics and Importing Data	Joseph Rudolf
10:45 am – 11:30 am	Data Visualization	Patrick Mathias
LUNCH		
12:30 pm - 1:30 pm	Data Transformation	Patrick Mathias
1:45 pm – 2:45 pm	Grouping and Summarizing Data	Joseph Rudolf
3:00 pm - 3:30 pm	Dashboard Demo and Course Wrap Up	Patrick Mathias

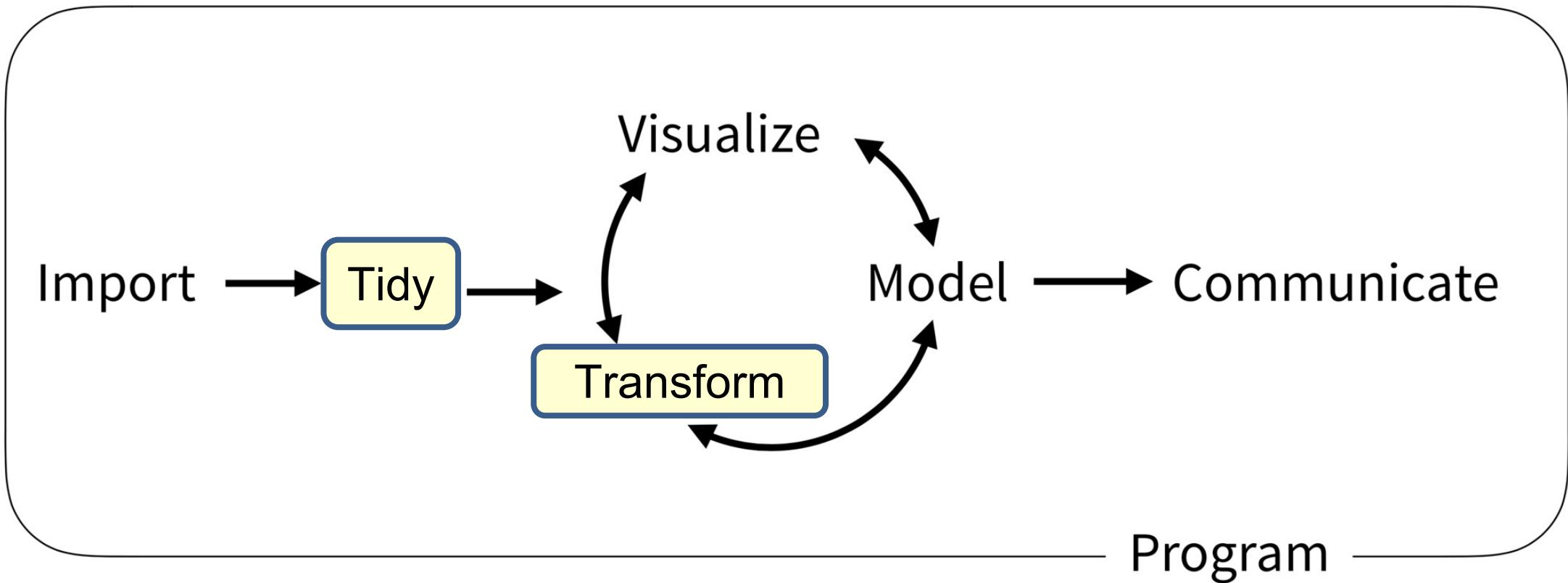
Goals

1. Learn how to use dplyr to transform data frames
2. Appreciate the role of piping in facilitating data transformation

Objectives

1. List the major forms of data transformation implemented in dplyr
2. Extract columns meeting inclusion criteria from a data frame
3. Create new calculated columns not found in the original data frame
4. Use the pipe operator to pass the output of one function as an input to the next function

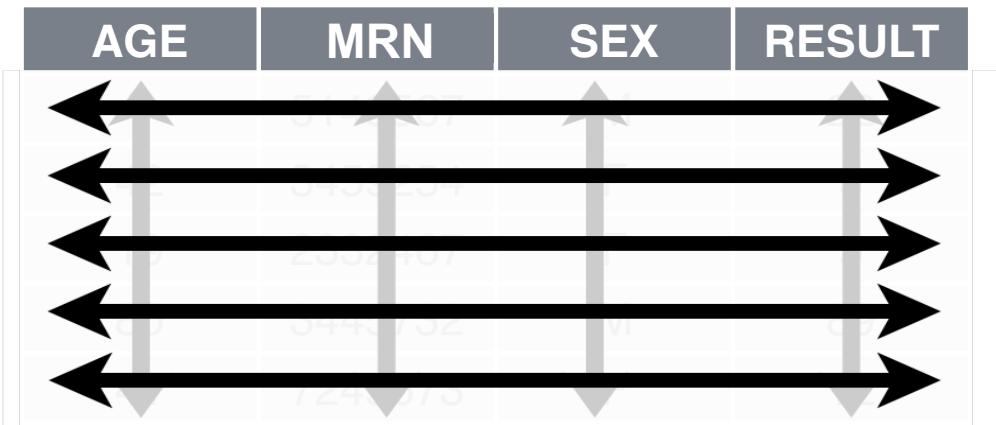
Typical Data Science Pipeline



What is a “Tidy” Data Frame

A data set is **tidy** if:

AGE	MRN	SEX	RESULT
1	100001	M	Normal
2	100002	M	Abnormal
3	100003	F	Normal
4	100004	F	Abnormal



1. Each **variable** is in its own **column**
2. Each **observation** is in its own **row**
3. Each **value** is in its own **cell**

Your Turn 1

Open "04-Transform.Rmd"

Run the setup chunk

```
```{r setup}
library(tidyverse) # Provides functions used throughout this session

covid_testing <- read_csv("covid_testing.csv")
```
```



Pop Quiz

How can you confirm that you have successfully loaded the data file into RStudio?

1. The code that imported the data did not yield an error
2. Code that references the `covid_testing` object runs without errors
3. The `covid_testing` object is present in the environment pane
4. All of the above

Transform Data with



dplyr



dplyr implements a *grammar* of data manipulation for transforming tabular data.

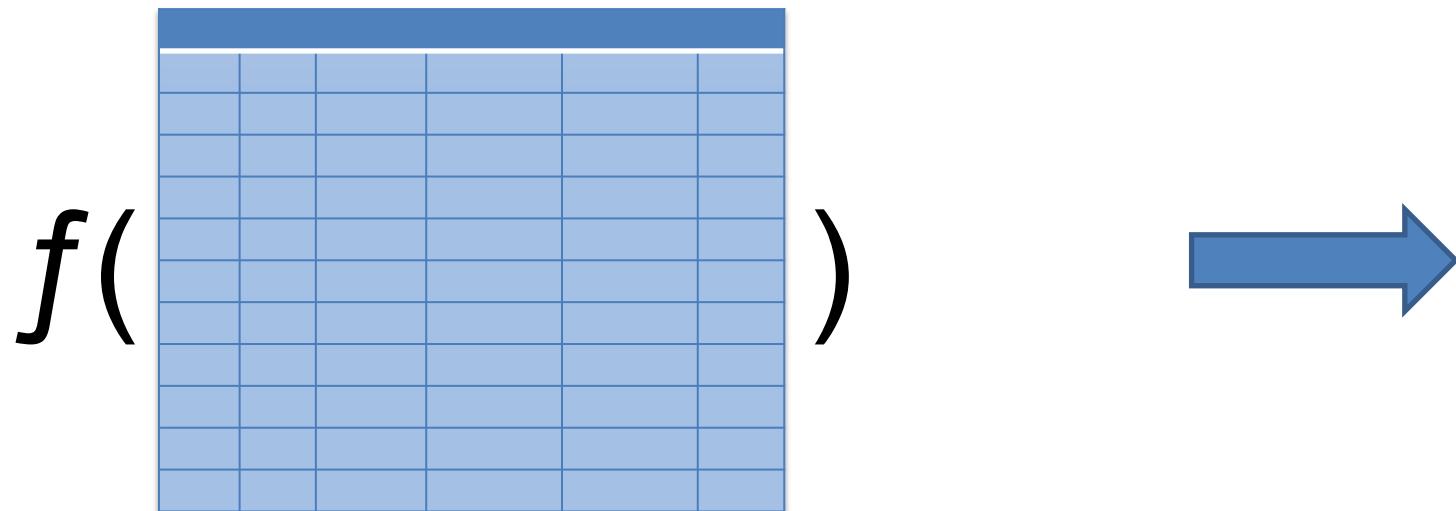


dplyr: a grammar for transforming data

- 1 Choose columns.** `select()`
- 2 Extract rows.** `filter()`
- 3 Derive new columns.** `mutate()`
- 4 Change the unit of analysis.** `summarize()`



Dplyr Approach



Dplyr Approach

$f($

| | | | | | |
|--|--|--|--|--|--|
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |

6 columns x 11 rows

) ↴

$f($

| | | | |
|--|--|--|--|
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |

4 columns x 11 rows

↗

$f($

| | | | |
|--|--|--|--|
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |

4 columns x 7 rows

↳

| | | | | |
|--|--|--|--|--|
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |

5 columns x 7 rows



Common syntax

Each function takes a data frame as its first argument and returns a data frame as its output.

```
function(data, ...)
```

dplyr
function

data frame to
transform

specific
arguments



Filtering a Subset of Rows

filter()

Extract rows that meet logical criteria

| | | | |
|--|--|--|--|
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |



| | | | |
|--|--|--|--|
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |

↓ Number of rows
= Number of Columns



Common syntax

Each function takes a data frame as its first argument
and returns a data frame as its output.

```
function(data, ...)
```

dplyr
function

data frame to
transform

specific
arguments



filter()

Extract rows that meet logical criteria

```
filter(data, ...)
```

**data frame to
transform**

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

| | | | | |
|--|--|--|--|-------|
| | | | | FALSE |
| | | | | FALSE |
| | | | | TRUE |
| | | | | FALSE |
| | | | | TRUE |
| | | | | FALSE |



| | | | |
|--|--|--|--|
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |



filter()

Extract rows that meet logical criteria

```
filter(data, column_name == criteria )
```

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

| | | | | |
|--|--|--|--|-------|
| | | | | FALSE |
| | | | | FALSE |
| | | | | TRUE |
| | | | | FALSE |
| | | | | TRUE |
| | | | | FALSE |



| | | | |
|--|--|--|--|
| | | | |
| | | | |
| | | | |
| | | | |



filter()

Extract rows that meet logical criteria

```
filter(covid_testing, mrn == 5000083)
```

| | mrn | first_name | last_name |
|-------|---------|------------|------------|
| FALSE | 5000876 | sarella | stark |
| FALSE | 5006017 | alester | stark |
| FALSE | 5001412 | jhezane | westerling |
| TRUE | 5000083 | lollys | clegane |



| | mrn | first_name | last_name |
|--|---------|------------|-----------|
| | 5000083 | lollys | clegane |



filter()

Extract rows that meet logical criteria

```
filter(covid_testing, mrn == 5000083)
```

| mrn | first_name | last_name |
|---------|------------|------------|
| 5000876 | sarella | stark |
| 5006017 | alester | stark |
| 5001412 | jhezane | westerling |
| 5000083 | lollys | clegane |

= sets

(returns nothing)

== tests if equal

(returns TRUE or FALSE)



filter()

Values coded as character strings must be surrounded by quotes

Extract rows that meet logical criteria.

```
filter(covid_testing, last_name == "stark")
```

| mrn | first_name | last_name | |
|---------|------------|------------|-------|
| 5000876 | sarella | stark | TRUE |
| 5006017 | alester | stark | TRUE |
| 5001412 | jhezane | westerling | FALSE |
| 5000083 | lollys | clegane | FALSE |



| mrn | first_name | last_name |
|---------|------------|-----------|
| 5000876 | sarella | stark |
| 5006017 | alester | stark |



filter()

Extract rows that meet logical criteria

```
filter(data, ... )
```

**data frame to
transform**

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



Logical tests

| | |
|------------------------|--------------------------|
| <code>x < y</code> | Less than |
| <code>x > y</code> | Greater than |
| <code>x == y</code> | Equal to |
| <code>x <= y</code> | Less than or equal to |
| <code>x >= y</code> | Greater than or equal to |
| <code>x != y</code> | Not equal to |
| <code>x %in% y</code> | Group membership |
| <code>is.na(x)</code> | Is NA |
| <code>!is.na(x)</code> | Is not NA |



Pop Quiz

What is the result?

1 == 1

Pop Quiz

What is the result?

$$3 != 1$$

Your Turn 2

Use filter() with the logical operators to find:

- Every test for patients **over age 80**
- All of the covid testing where the demographic group (demo_group) is **equal to "client"**

filter()

Filter to multiple matches

```
covid_testing %>%  
  filter(first_name %in% c("jon", "daenerys"))
```

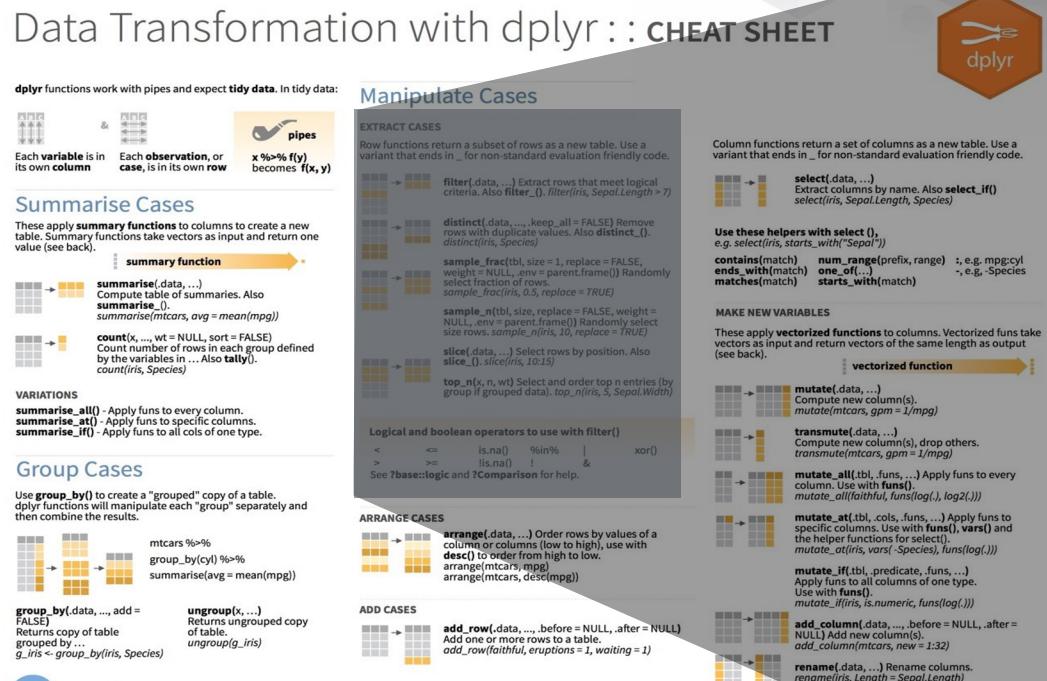
| mrn | first_name |
|---------|------------|
| <dbl> | <chr> |
| 5001412 | jhezane |
| 5000533 | penny |
| 5009134 | grunt |
| 5008518 | melisandre |
| 5008967 | rolley |



| mrn | first_name |
|---------|------------|
| <dbl> | <chr> |
| 5002427 | daenerys |
| 5011120 | jon |
| 5001092 | jon |
| 5004082 | jon |
| 5005197 | daenerys |



filter() variants



EXTRACT CASES

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

filter(.data, ...) Extract rows that meet logical criteria. `filter(iris, Sepal.Length > 7)`

distinct(.data, ..., .keep_all = FALSE) Remove rows with duplicate values. `distinct(iris, Species)`

sample_frac(tbl, size = 1, replace = FALSE, weight = NULL, env = parent.frame()) Randomly select fraction of rows. `sample_frac(iris, 0.5, replace = TRUE)`

sample_n(tbl, size, replace = FALSE, weight = NULL, env = parent.frame()) Randomly select size rows. `sample_n(iris, 10, replace = TRUE)`

slice(.data, ...) Select rows by position. `slice(iris, 10:15)`

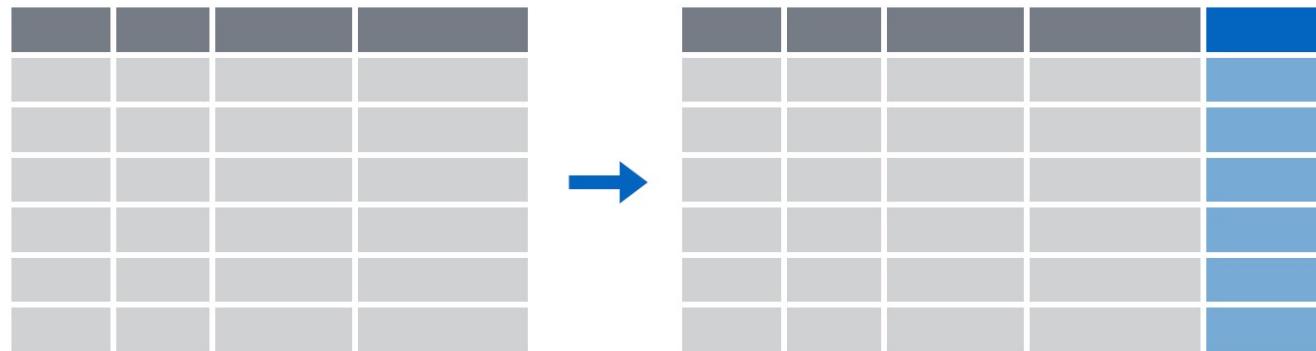
top_n(x, n, wt) Select and order top n entries (by group if grouped data). `top_n(iris, 5, Sepal.Width)`



Creating a New Column

mutate()

Creating new calculated columns



= Number of rows
↑ Number of Columns



mutate()

Creating new calculated columns

```
Covid_testing %>%  
  mutate(new_column = calculation)
```

name for new column

equals

function whose results will populate columns



mutate()

Creating new calculated columns

```
covid_testing %>%  
  mutate(c_r_tat_mins = col_rec_tat * 60)
```

| mrn | col_rec_tat | rec_ver_tat |
|---------|-------------|-------------|
| 5000876 | 29.5 | 11.5 |
| 5006017 | 3.6 | 5 |
| 5001412 | 1.4 | 5.2 |
| 5000533 | 2.3 | 5.8 |



| mrn | col_rec_tat | rec_ver_tat | c_r_tat_mins |
|---------|-------------|-------------|--------------|
| 5000876 | 29.5 | 11.5 | 1770 |
| 5006017 | 3.6 | 5 | 216 |
| 5001412 | 1.4 | 5.2 | 84 |
| 5000533 | 2.3 | 5.8 | 138 |



Your Turn 3

Create a new column using the `mutate()` function that contains the total TAT (sum of `col_rec_tat` and `rec_ver_tat`)

Functions to use in mutate()

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

dplyr::cumall() - Cumulative all()
dplyr::cumany() - Cumulative any()
 cummax() - Cumulative max()
dplyr::cummean() - Cumulative mean()
 cummin() - Cumulative min()
 cumprod() - Cumulative prod()
 cumsum() - Cumulative sum()

RANKINGS

dplyr::cume_dist() - Proportion of all values <=
dplyr::dense_rank() - rank with 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"

MATH

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

MISC

dplyr::case_when() - multi-case if_else()
dplyr::coalesce() - first non-NA values by element across a set of vectors
dplyr::if_else() - element-wise if() + else()
dplyr::na_if() - replace specific values with NA
 pmax() - element-wise max()
 pmin() - element-wise min()
dplyr::recode() - Vectorized switch()
dplyr::recode_factor() - Vectorized switch() for factors

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

dplyr::cumall() - Cumulative all()
dplyr::cumany() - Cumulative any()
 cummax() - Cumulative max()
dplyr::cummean() - Cumulative mean()
 cummin() - Cumulative min()
 cumprod() - Cumulative prod()
 cumsum() - Cumulative sum()

RANKINGS

dplyr::cume_dist() - Proportion of all values <=
dplyr::dense_rank() - rank with 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"

MATH

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

MISC

dplyr::case_when() - multi-case if_else()
dplyr::coalesce() - first non-NA values by element across a set of vectors
dplyr::if_else() - element-wise if() + else()
dplyr::na_if() - replace specific values with NA
 pmax() - element-wise max()
 pmin() - element-wise min()
dplyr::recode() - Vectorized switch()
dplyr::recode_factor() - Vectorized switch() for factors

Summary Functions

TO USE WITH SUMMARISE ()

summarise() applies summary functions to columns to create a new table. Summary functions take vectors as input and return single values as output.

summary function

COUNTS

dplyr::n() - number of values/rows
dplyr::n_distinct() - # of uniqueness
 sum(is.na()) - # of non-NA's

LOCATION

mean() - mean, also **mean(is.na())**
median() - median

LOGICALS

mean() - proportion of TRUE's
sum() - # of TRUE's

POSITION/ORDER

dplyr::first() - first value
dplyr::last() - last value
dplyr::nth() - value in nth location of vector

RANK

quantile() - nth quantile
min() - minimum value
max() - maximum value

SPREAD

IQR() - Inter-Quartile Range
mad() - median absolute deviation
sd() - standard deviation
var() - variance

Row Names

Tidy data does not use rownames, which store a variable outside of the columns. To work with the rownames, first move them into a column.

rownames_to_column()
Move row names into col.
o <- rownames_to_column(iris, var = "C")

column_to_rownames()
Move col in row names.
o <- column_to_rownames(a, var = "C")

Also has **rownames()**, **remove_rownames()**

Combine Tables

COMBINE VARIABLES

x + **y** = **bind_cols()**
Use **bind_cols()** to paste tables beside each other as side-by-side.

x + **y** = **bind_col()**
Use **bind_col()** to paste tables placed side by side as a single table.
BE SURE THAT ROWS ALIGN.

Use a "Mutating Join" to join one table to columns from another, matching values with the rows that they correspond to. Each join retains a different combination of values from the tables.

left_join(x, y, by = NULL, copy = FALSE, suffix = c("x", "y"))
Join matching values from y to x.

right_join(x, y, by = NULL, copy = FALSE, suffix = c("x", "y"))
Join matching values from x to y.

inner_join(x, y, by = NULL, copy = FALSE, suffix = c("x", "y"))
Join data. Retain only rows with matches.

full_join(x, y, by = NULL, copy = FALSE, suffix = c("x", "y"))
Join all values, all rows.

use_by = c("col1", "col2", ...) to specify one or more common columns to match on.
left_join(x, y, by = "A")

use_by = c("col1", "col2", ...) to match on columns that have different names in each table.
left_join(x, y, by = "C")

use_by = c("col1", "col2", ...) to give to unmatched columns that have the same name in both tables.
left_join(x, y, by = "C" ~ "D", suffix = c("1", "2"))

semi_join(x, y, by = NULL, ...)
Returns rows of x that have a match in y.
USEFUL TO SEE WHAT WILL BE JOINED.

anti_join(x, y, by = NULL, ...)
Return rows of x that do not have a match in y.
USEFUL TO SEE WHAT WILL NOT BE JOINED.



on back



dplyr

RStudio® is a trademark of RStudio, Inc. • CC BY SA RStudio • info@rstudio.com • 844-448-1212 • rstudio.com • Learn more with browseVignettes(package = c("dplyr", "tidyverse")) • dplyr 0.7.0 • tidyverse 1.2.0 • Updated: 2017-03

mutate()

Replacing columns

```
covid_testing %>%  
  mutate(mrn = as.character(mrn))
```

Function to "coerce" one type of data into another type of data

| mrn
<dbl> | first_name
<chr> | last_name
<chr> |
|--------------|---------------------|--------------------|
| 5000876 | sarella | stark |
| 5006017 | alester | stark |
| 5001412 | jhezane | westerling |
| 5000533 | penny | targaryen |

| mrn
<chr> | first_name
<chr> | last_name
<chr> |
|--------------|---------------------|--------------------|
| 5000876 | sarella | stark |
| 5006017 | alester | stark |
| 5001412 | jhezane | westerling |
| 5000533 | penny | targaryen |



mutate()

Conditionally replacing values

```
covid_testing %>%  
  mutate(last_name = if_else(last_name=="targaryen",  
                            "TARGARYEN",last_name))
```

| mrn
<dbl> | first_name
<chr> | last_name
<chr> |
|--------------|---------------------|--------------------|
| 5001412 | jhezane | westerling |
| 5000533 | penny | targaryen |
| 5009134 | grunt | rivers |
| 5008518 | melisandre | swyft |



| mrn
<dbl> | first_name
<chr> | last_name
<chr> |
|--------------|---------------------|--------------------|
| 5001412 | jhezane | westerling |
| 5000533 | penny | TARGARYEN |
| 5009134 | grunt | rivers |
| 5008518 | melisandre | swyft |



Piping Data from One Row to the Next

Answering the age old lab question

There have been some complaints from a few inpatient units that turnaround times are longer than expected. You would like to take a closer look at the tests that took longer than expected. For these inpatient samples turnaround times should rarely be longer than 18 hours, so you would specifically like to filter rows that are from inpatients and have a total turnaround time (collect to verify) of more than 18 hours.

Data Analysis Steps

```
tat_outliers <- filter(covid_testing, patient_class == 'inpatient')

tat_outliers <- mutate(tat_outliers,
                      total_tat = col_rec_tat + rec_ver_tat)

tat_outliers <- filter(tat_outliers, total_tat > 18)
```

1. Filter to extract a subset of inpatient results
2. Mutate to add a column for total turnaround time
3. Filter to select only the results with a total turnaround time greater than 18 hours



The Pipe Operator %>%

Passes result on left into first argument of function on right.



```
covid_testing %>% filter(____, patient_class == 'inpatient')
```

```
filter(covid_testing, patient_class == 'inpatient')
```

```
covid_testing %>% filter(patient_class == 'inpatient')
```



Data Analysis Steps

```
tat_outliers <- covid_testing %>%  
  filter(patient_class == 'inpatient') %>%  
  mutate(total_tat = col_rec_tat + rec_ver_tat) %>%  
  filter(total_tat > 18)
```

1. Filter to extract a subset of inpatient results
2. Mutate to add a column for total turnaround time
3. Filter to select only the results with a total turnaround time greater than 18 hours



Shortcut to type %>%

Cmd + Shift + M (Mac)

Ctrl + Shift + M (Windows)



Scene

The PICU would like a word with you because of a recent incident involving a delay in results for a patient who required a AGP

They had to wait over 10 hours before the procedure could begin

You decide to investigate... WITH DATA



Your Turn 4

Use `%>%` to write a sequence of three functions that:

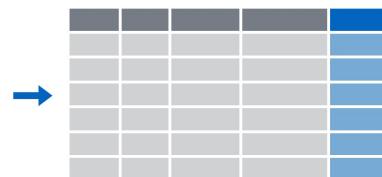
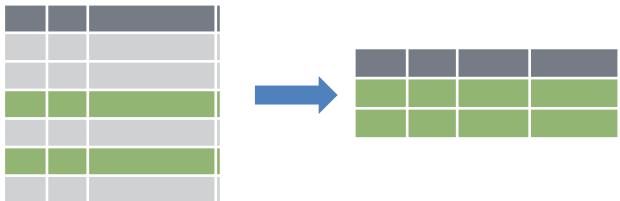
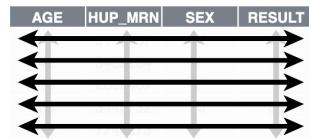
1. Filters to tests from the clinic (`clinic_name`) of "picu"
2. Adds a column with the total turnaround time (sum of `col_ver_tat` and `rec_ver_tat`)
3. Filters to include only the turnaround times greater than 10 hours

Using `<-`, assign the result to a new variable, call it whatever you want.

Recap



dplyr is a package that provides a **grammar of data manipulation**. Common operations include extracting a subset of your data set and adding new columns.



1. Fully utilizing **dplyr** requires a **tidy data frame**, in which each **variable** is in its own **column**, each **observation** is in its own **row**, each **value** is in its own **cell**;
2. The **filter** function allows you to extract a subset of your data using logical conditions
3. The **mutate** function allows you to add new columns to your data frame that can apply calculations across columns in your data frame



What Else?

select()

Extract columns from a data frame

| | | | |
|------|------|------|------|
| Blue | Grey | Blue | Grey |
| Blue | Grey | Blue | Grey |
| Blue | Grey | Blue | Grey |
| Blue | Grey | Blue | Grey |
| Blue | Grey | Blue | Grey |



| | |
|------|--|
| Blue | |

= Number of rows

↓ Number of Columns



select()

Extract columns from a data frame

```
select(covid_testing, mrn, last_name)
```

dplyr
function

data frame to
transform

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



select()

Extract columns from a data frame **by name**

```
select(covid_testing, mrn, last_name)
```

covid_testing

| mrn | first_name | last_name | gender |
|---------|------------|------------|--------|
| 5000876 | sarella | stark | female |
| 5006017 | alester | stark | male |
| 5001412 | jhezane | westerling | female |
| 5000533 | penny | targaryen | female |



| mrn | last_name |
|---------|------------|
| 5000876 | stark |
| 5006017 | stark |
| 5001412 | westerling |
| 5000533 | targaryen |

...



select()

Extract columns from a data frame **by name**

```
select(covid_testing, -mrn, -last_name)
```

covid_testing

| mrn | first_name | last_name | gender |
|---------|------------|------------|--------|
| 5000876 | sarella | stark | female |
| 5006017 | alester | stark | male |
| 5001412 | jhezane | westerling | female |
| 5000533 | penny | targaryen | female |
| ... | | | |



| first_name | gender |
|------------|--------|
| sarella | female |
| alester | male |
| jhezane | female |
| penny | female |



select() helpers

Data Transformation with dplyr :: CHEAT SHEET

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

- Each variable is in its own column
- Each observation, or case, is in its own row
- x %>% f(y) becomes f(x, y)

Manipulate Cases

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

- filter(data, ...)
- distinct(..., keep_all = FALSE)
- sample_frac(tbl, size = 1, replace = FALSE, weight = NULL, prob = parent.frame())
- sample_n(tbl, size, replace = FALSE, weight = NULL, env = parent.frame())
- slice(...)
- top_n(..., n, wt)

VARIATIONS

- summarise_all()
- summarise_att()
- summarise_if()

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.

- mtcars %>% group_by(cyl) %>% summarise(avg = mean(mpg))
- ungroup(x, ...)
- group_by(data, ..., add = FALSE)
- group_by_copy(...)
- g_iris <- group_by(iris, Species)

Manipulate Variables

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

- pull(data, var = -1)
- select(...)
- select_if(...)

MAKE NEW VARIABLES
These apply vectors as inputs to functions (see back).

- contains(match)
- ends_with(match)
- matches(match)

Logical and boolean operators to use with filter()

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

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

ARRANGE CASES

- arrange(data, ...)
- arrange(mtcars, mpg)
- arrange(mtcars, desc(mpg))

ADD CASES

- add_row(data, ..., before = NULL, after = NULL)
- add_rows(faithful, eruptions = 1, waiting = 1)

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)**



select()

Renaming columns

```
covid_testing %>%  
  select(MRN = mrn, first_name, last_name)
```

| mrn
<dbl> | first_name
<chr> | last_name
<chr> |
|--------------|---------------------|--------------------|
| 5001412 | jhezane | westerling |
| 5000533 | penny | targaryen |
| 5009134 | grunt | rivers |
| 5008518 | melisandre | swyft |



| MRN
<dbl> | first_name
<chr> | last_name
<chr> |
|--------------|---------------------|--------------------|
| 5001412 | jhezane | westerling |
| 5000533 | penny | targaryen |
| 5009134 | grunt | rivers |
| 5008518 | melisandre | swyft |



arrange()

Order rows by values in a column

| | | | |
|--|--|--|--|
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |



| | | | |
|--|--|--|--|
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |

- = Number of rows
- = Number of Columns



arrange()

Order rows by values in a column

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

data frame to
transform

name(s) of columns to
arrange by



arrange()

Order rows by values in a column

```
arrange(covid_testing, first_name)
```

| mrn | first_name | last_name |
|---------|------------|-----------|
| 5000876 | sarella | stark |
| 5006017 | alester | stark |
| 5001412 | jhezane | targaryen |
| 5000533 | penny | targaryen |



| mrn | first_name | last_name |
|---------|------------|-----------|
| 5006017 | alester | stark |
| 5001412 | jhezane | targaryen |
| 5000533 | penny | targaryen |
| 5000876 | sarella | stark |



arrange()

Order rows by values in a column

```
arrange(covid_testing, desc(mrn))
```

| mrn | first_name | last_name |
|---------|------------|-----------|
| 5000876 | sarella | stark |
| 5006017 | alester | stark |
| 5001412 | jhezane | targaryen |
| 5000533 | penny | targaryen |



| mrn | first_name | last_name |
|---------|------------|-----------|
| 5006017 | alester | stark |
| 5001412 | jhezane | targaryen |
| 5000876 | sarella | stark |
| 5000533 | penny | targaryen |

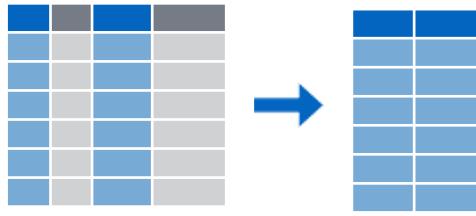


Pop Quiz

The default behavior of `arrange()` is to order from lower to higher values.

When might `arrange()` place "1000" before "50"?

Transforming data



Extract variables with `select()`



Extract rows with `filter()`



Arrange rows, with `arrange()`.



Add calculated columns with
`mutate()`



Example: Applying the pipe

```
covid_testing %>%  
  filter(pan_day <= 10) %>%  
  select(clinic_name) %>%  
  arrange(clinic_name)
```



Goal

1. Learn how to use dplyr to transform data frames
2. Appreciate the role of piping in facilitating data transformation

Objectives

1. List the major forms of data transformation implemented in dplyr
2. Extract columns meeting inclusion criteria from a data frame
3. Create new calculated columns not found in the original data frame
4. Use the pipe operator to pass the output of one function as an input to the next function