# Data Transformation: Slicing Your Data

**Patrick Mathias**

December 13, 2020

*Presentation adapted from…*

# Amrom Obstfeld

Assistant Professor of Clinical Pathology and Laboratory Medicine

University of Pennsylvania Perelman School of Medicine

Director of Hematology and Coagulation Laboratories
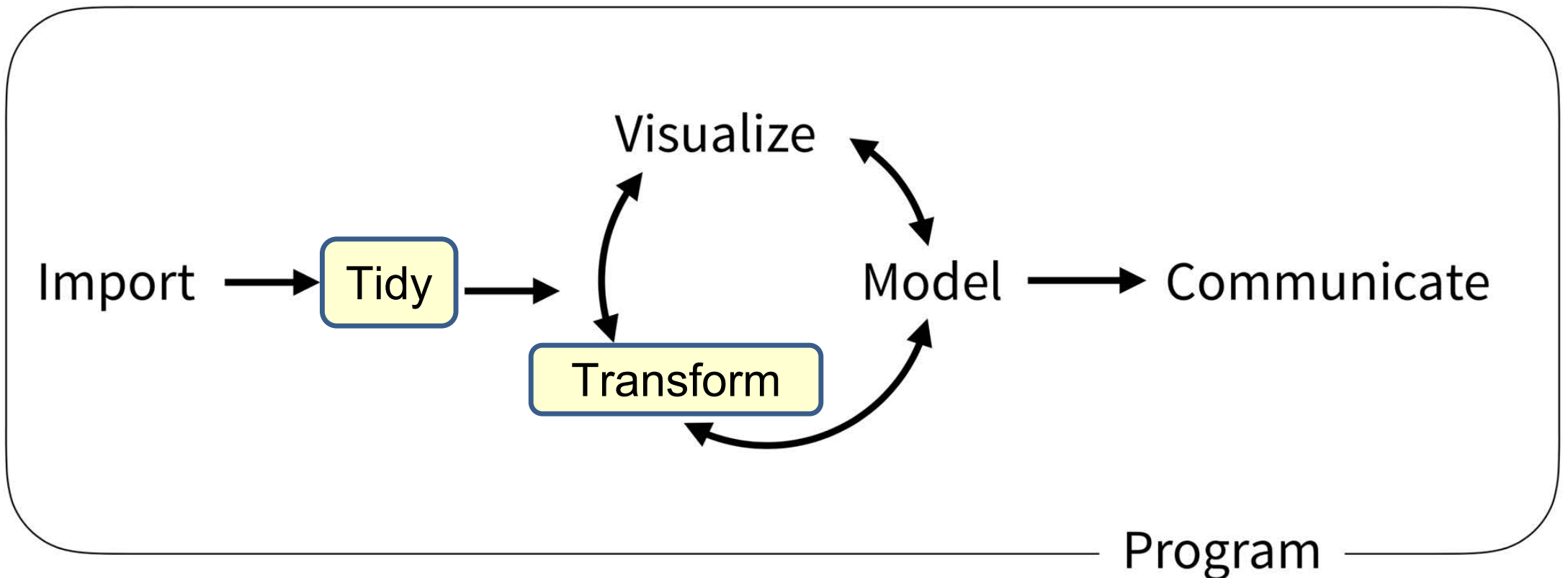
Children's Hospital of Philadelphia

# Goal

1. Learn how to use dplyr to transform data frames

# Objectives

1. List the major forms of data transformation implemented in dplyr
2. Use code templates with dplyr functions to tidy a raw data set

# Typical Data Science Pipeline

# What is a "Tidy" Data Frame

| AGE | MRN | SEX | RESULT |
|-----|-----|-----|--------|

A data set is **tidy** if:

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

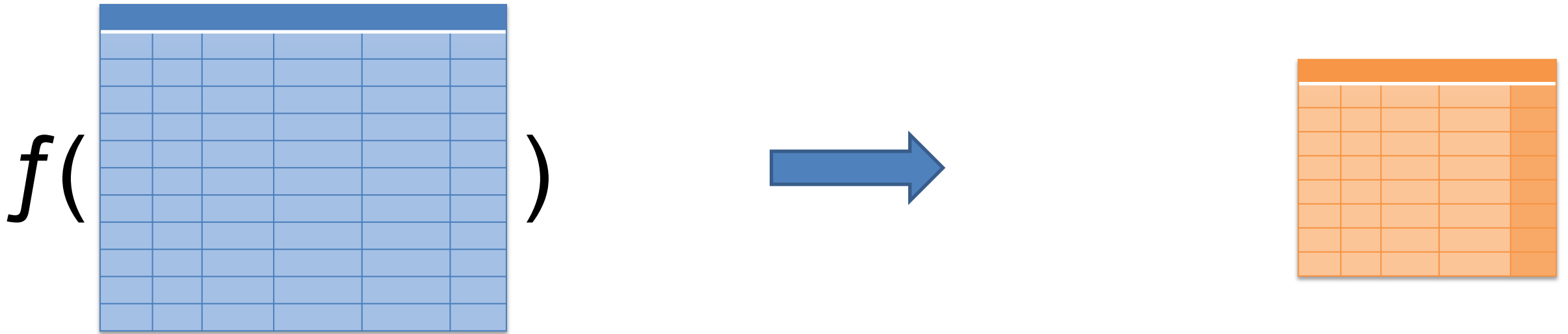covid_testing <- read_csv("data/covid_testing.csv")
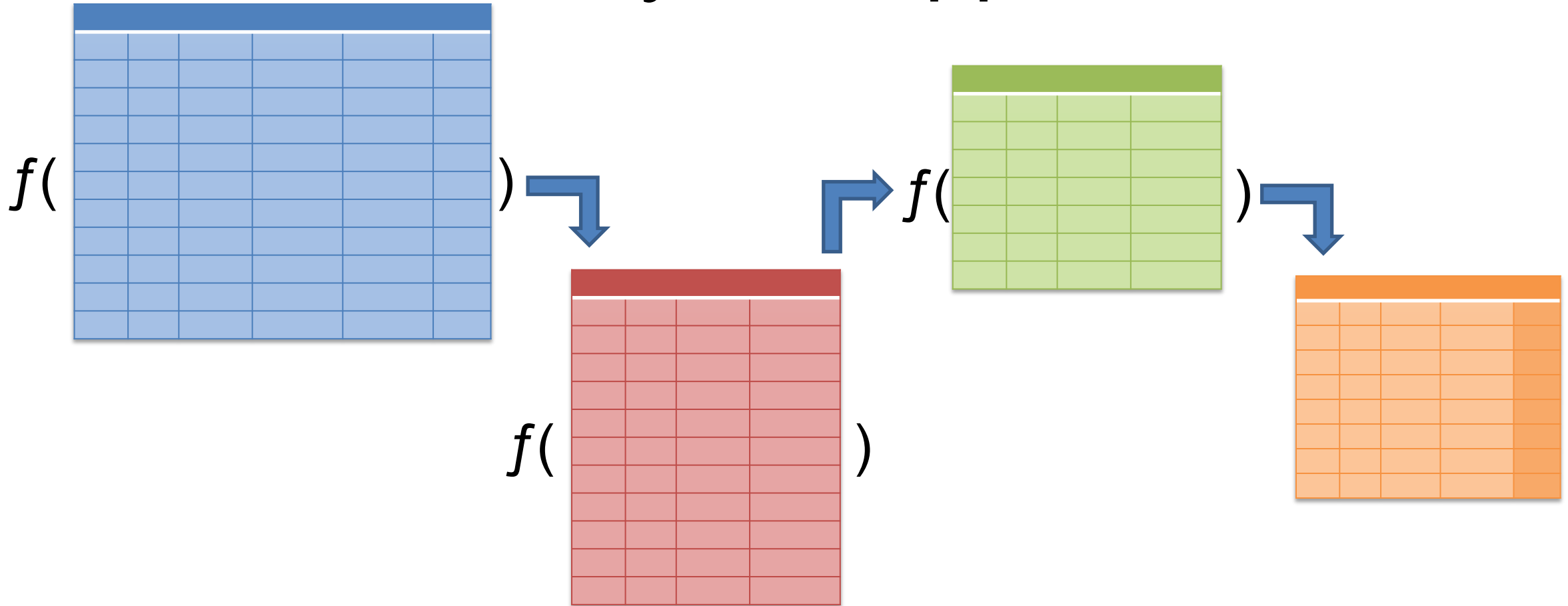```
```

# Transform Data with

# dplyr



dplyr implements a *grammar* for transforming tabular data.

# Analytical Approach

$f($  $)$ →

# Analytical Approach

# 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()**

# 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**

# Pulling specific columns out of your data frame

# select()

# select()

Extract columns from a data frame



= 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 |

dplyr

# 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 |

# Your Turn 2

- Alter the code to select just the `first_name` column from `covid_testing`

- If you have time, try to remove the `first_name` column

`covid_testing_2 <- select(covid_testing, _____)`

# select() helpers

# Pulling specific rows out of your data frame

# filter()

# 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)

# 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)

# 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 |

FALSE
FALSE
FALSE
TRUE

→

| 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)

dplyr

# filter()

Extract rows that meet logical criteria.

Values coded as character strings must be surrounded by quotes

```
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 |

dplyr

# 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)

dplyr

# Logical tests

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

# Pop Quiz

What is the result?

**1 == 1**

# Pop Quiz

## What is the result?

### 3 != 1

# Your Turn 3

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"**
- CHALLENGE:
  - All of the covid testing where the patient class (patient_class) **is NA** [Hint: See slide titled "Logical Tests"]

# filter() variants

# What else?

# arrange()

# 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 |

dplyr

# 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 |

dplyr

# Your Turn 4

The column `ct_result` contains the cycle threshold (Ct) for the real-time PCR that generated the final result.

How might you use `arrange()` to determine the highest and lowest Ct result in the dataset?

# Pop Quiz

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

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

# Goal

1.  Learn how to use dplyr to transform data frames

# Objectives

1.  List the major forms of data transformation implemented in dplyr
2.  Use code templates with dplyr functions to tidy a raw data set