# COVID-19
# Data Science Analysis

Mhealyssah Bustria & Anjelina Velazquez
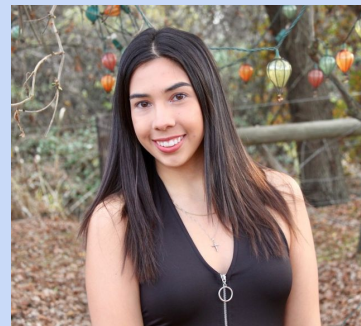Mentors: Dr. Niema Moshiri & Dr. Youwen Ouyang

UCSD DBMI Summer Internship 2020
August 12, 2020

# About Us

**Mhealyssah (Mhea) Bustria**

A Computer Science undergraduate at CSUSM. She is always happy to give back to her community, and she is especially-interested in advancing the fields of education and health.

**Anjelina Velazquez**

Currently a fourth year Computer Science student at CSUSM. Enjoys keeping busy by always learning new concepts and ideas. She is always encouraging others to do the same.

# Outline

**BACKGROUND**

- Motivation

- Our Project

- Methods

**RESULTS**

- Reading patient data

- Analyzing the dataset

**CONCLUSIONS**

- Lessons learned

- Next steps

# Background - Motivation

To combat the **COVID-19 pandemic**, researchers need to make inferences from data collected from patients.

# Background - Motivation

To combat the **COVID-19 pandemic**, researchers need to make inferences from data collected from patients.

**Are these inferences generalizable?**

# Background - Motivation

To combat the **COVID-19 pandemic**, researchers need to make inferences from data collected from patients.

**Are these inferences generalizable?**

Researchers need to account for...

**potential errors and inaccuracies** that are present in the data collected due to manual input

**potential confounding factors,** such as biases in patient sampling

# Background - Our project

Collection of over 75,000 SARS-CoV-2 Patient Records

**Manually-entered data**

Errors, inconsistencies, and missing data affects the data science analysis process.

**Large dataset**

What do the patient records in our dataset look like?

What are some potential confounding factors?

# Research Goals

**Reading patient data**

Demonstrate how
manually-entered data
affects
data science analysis.

**Analysis and visualization**

Show how
demographic information varies
in our large sample.

# Methods

# Methods

**Obtaining the dataset**

Global initiative on sharing all influenza data (GISAID)

- Extracted records were stored in a gzipped JSON file

# Methods



**Obtaining the dataset**

Global initiative on sharing all influenza data (GISAID)

- Extracted records were stored in a gzipped JSON file

**Data reading and analysis**

Python

- Identify **what information can be found** in the records
- Decide **what information to use** for analysis and visualization

# Methods



**Obtaining the dataset**

Global initiative on sharing all influenza data (GISAID)

- Extracted records were stored in a gzipped JSON file



**Data reading and analysis**

Python

- Identify **what information can be found** in the records
- Decide **what information to use** for analysis and visualization



**Data visualization**

**matplotlib** Python library

- pie charts

**seaborn** Python library (based on matplotlib)

- bar plots
- count plots
- violin plots

# RESULTS
---
# Reading Patient Data

# Results - Reading patient data

How were our analysis and visualization processes affected by manual data-entry practices?

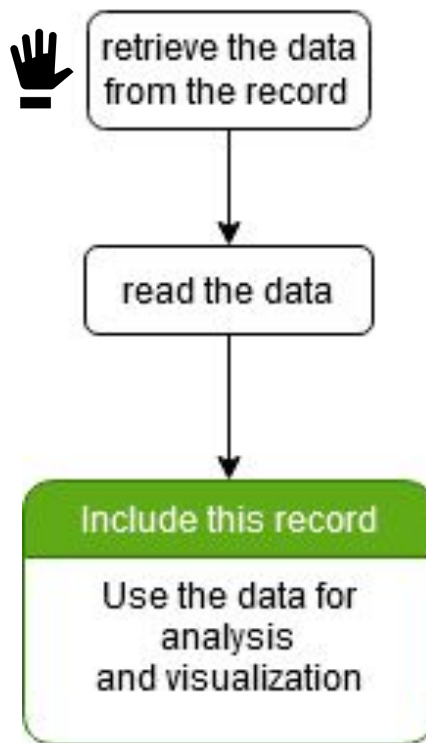How did we modify our data-reading strategies to account for unusable data?

- data correction
- data exclusion

What are potential solutions to address the issues caused by manual data entry?
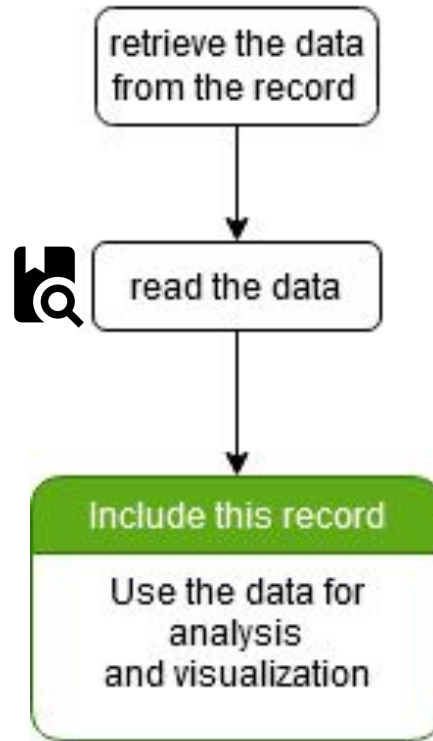
# Original strategy for reading the data
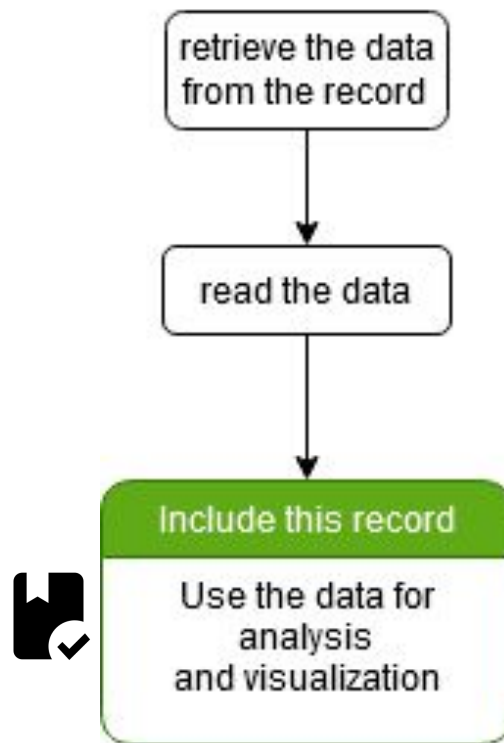
# Original strategy for reading the data

# Original strategy for reading the data

# Original strategy for reading the data

# 👤?Types of missing / unknown / invalid data

**Type 1**

Data for the attribute-of-interest was **not provided**.

# 👤?Types of missing / unknown / invalid data

**Type 1**

Data for the attribute-of-interest was **not provided**.

**Type 2**

Data for the attribute-of-interest was provided, but was entered as some variation of **"unknown" or "not applicable"**.

# 👤?Types of missing / unknown / invalid data

**Type 1**

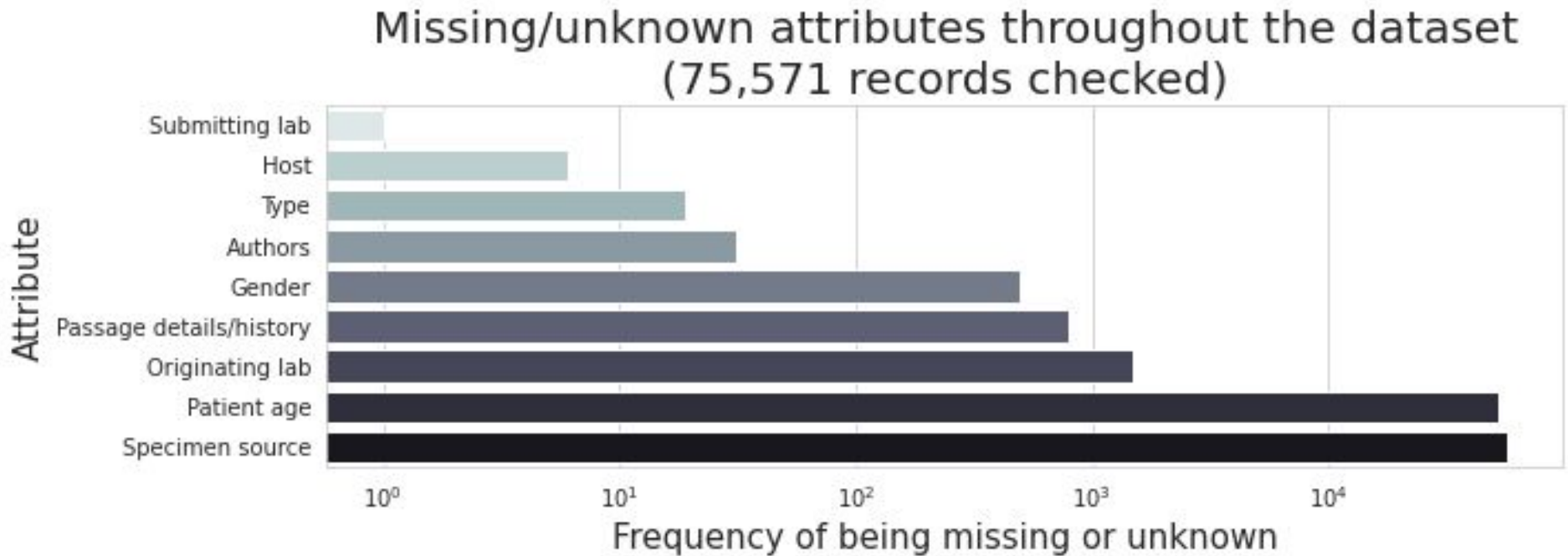Data for the attribute-of-interest was **not provided**.

**Type 2**

Data for the attribute-of-interest was provided, but was entered as some variation of **"unknown" or "not applicable"**.
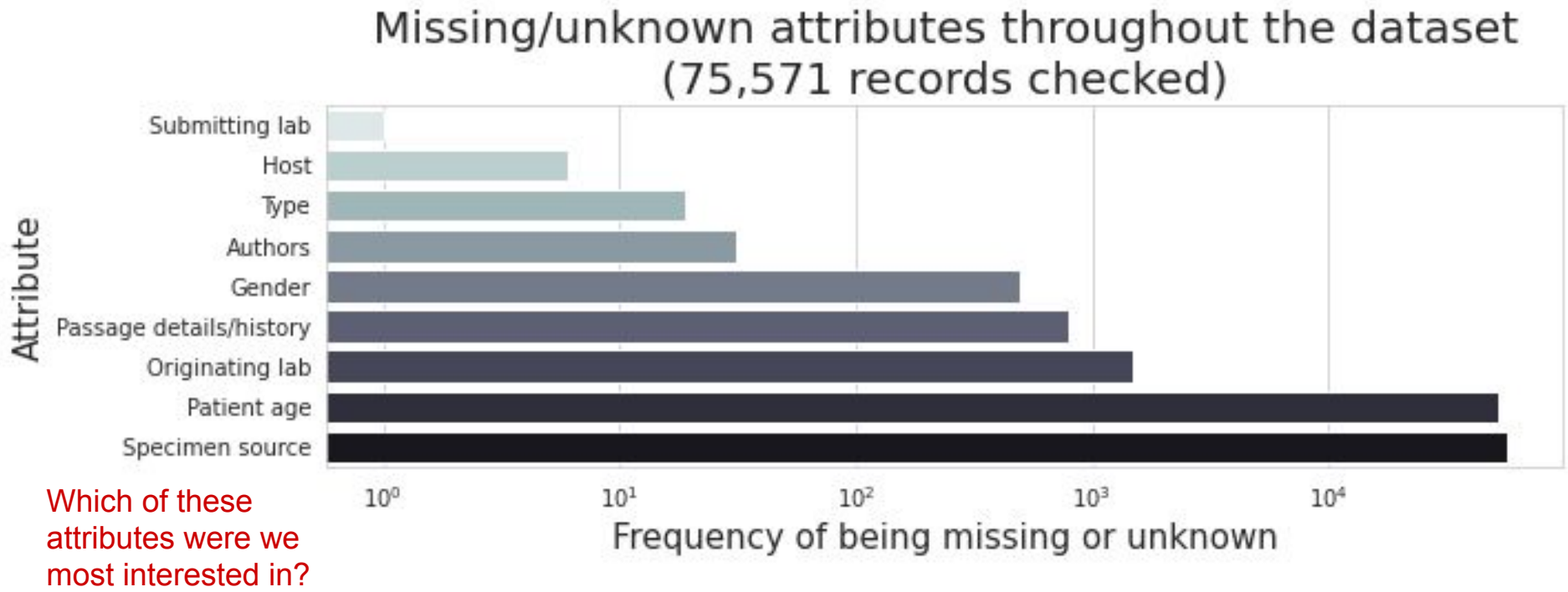
**Type 3**

Data for the attribute-of-interest contained an error such as **formatting inconsistencies or misspellings**.
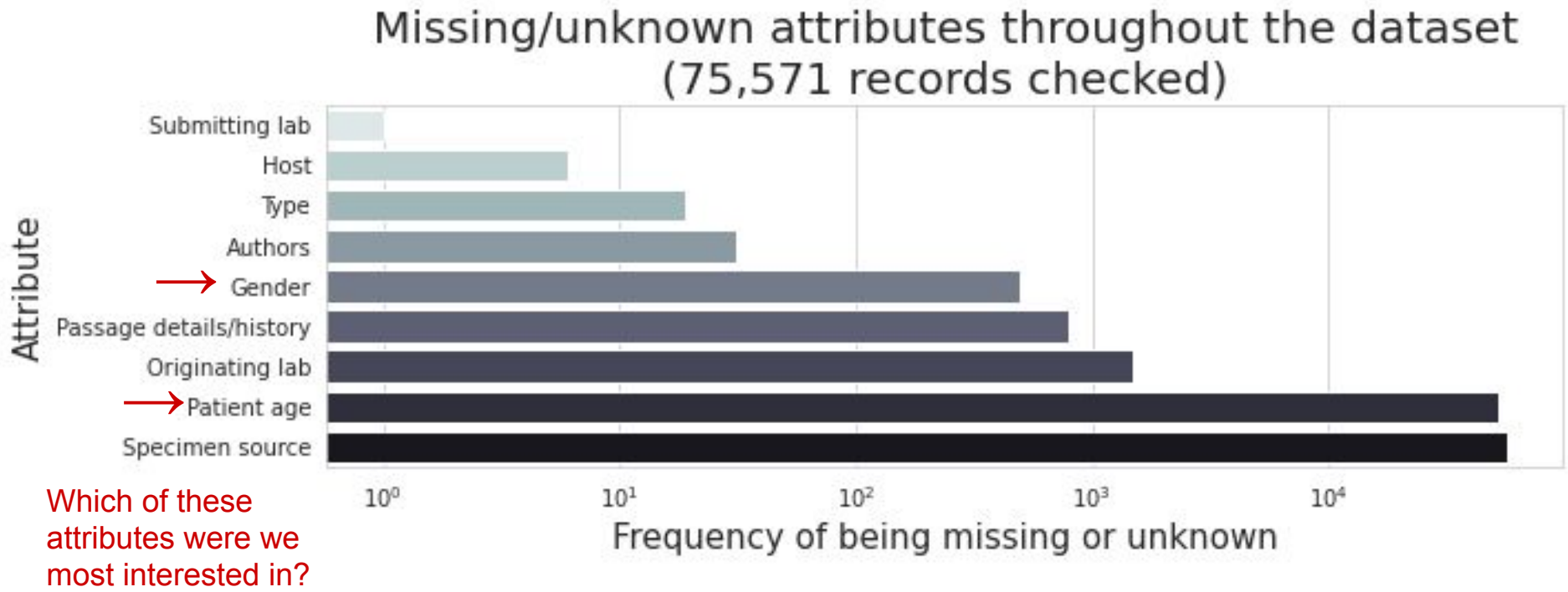
# 👤?How often was attribute data missing or invalid/unknown?



Missing/unknown attributes throughout the dataset (75,571 records checked)

# ?How often was attribute data missing or invalid/unknown?



Missing/unknown attributes throughout the dataset
(75,571 records checked)

Which of these attributes were we most interested in?

# 👤?How often was attribute data missing or invalid/unknown?



Missing/unknown attributes throughout the dataset
(75,571 records checked)
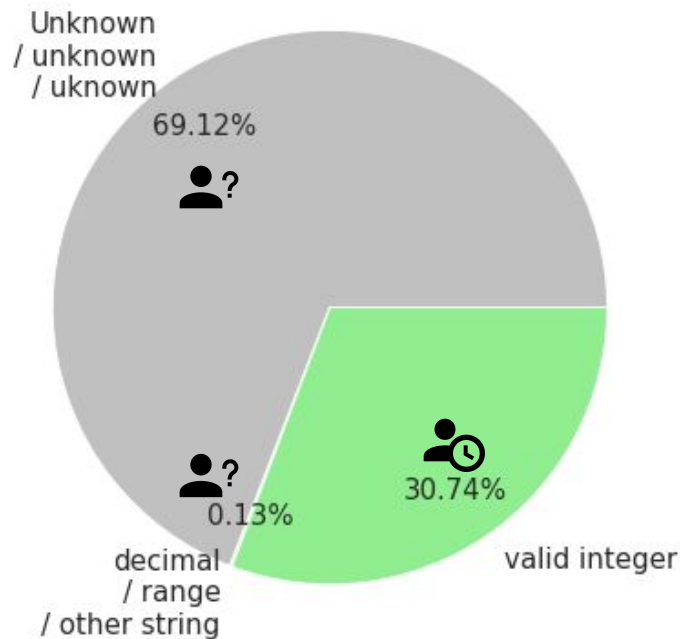
Which of these attributes were we most interested in?

# 👫 How often did we encounter unknown gender data?



Amounts of known and unknown
gender data from 75,571 records

male
16.96%

female
14.23%
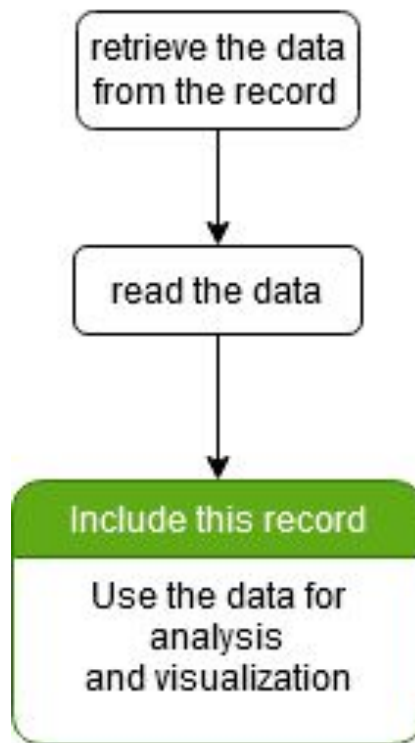
unknown
68.81%

# How often did we encounter unknown or unusable age data?



Formats of age data of
75,571 COVID-19 patient records

Unknown / unknown / uknown
69.12%

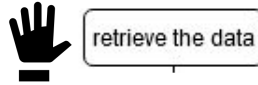decimal / range / other string
0.13%

valid integer
30.74%

# Original strategy for reading the data



Throughout the study,
how did we
revise our strategies
to handle the issues
involved with
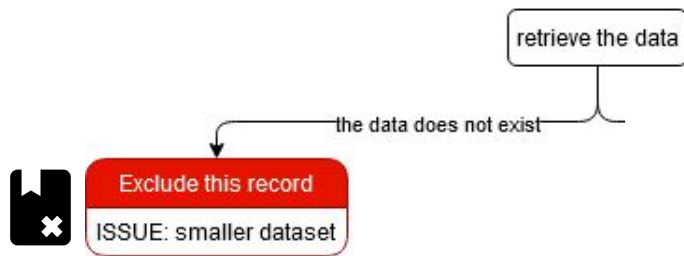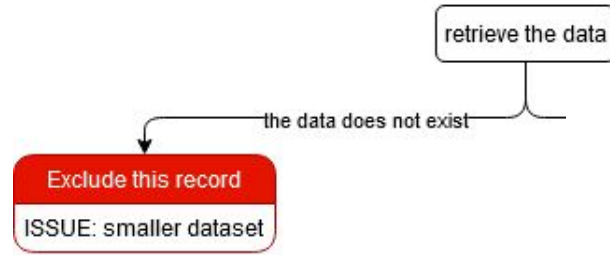manually-entered data?

# Handling manually-entered data


retrieve the data

# Handling manually-entered data

retrieve the data

the data may or may not
be missing

# Handling manually-entered data



retrieve the data
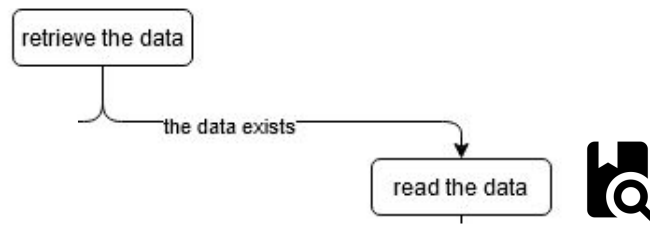
the data does not exist

Exclude this record

ISSUE: smaller dataset

# Handling manually-entered data - example of data exclusion



retrieve the data
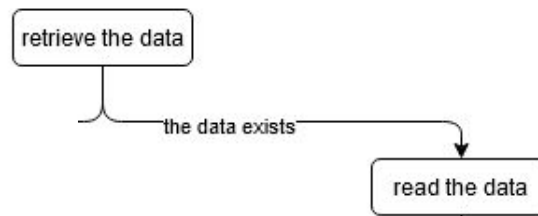
the data does not exist

Exclude this record

ISSUE: smaller dataset

If age is an attribute of interest but there is no age provided in that record → **data exclusion**

# Handling manually-entered data



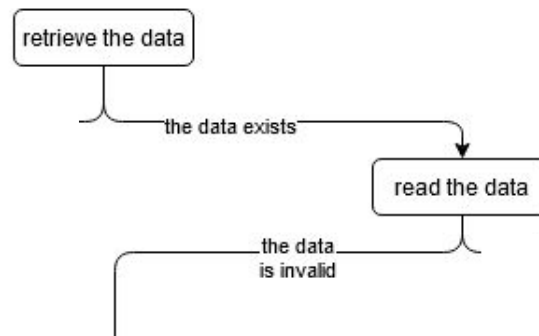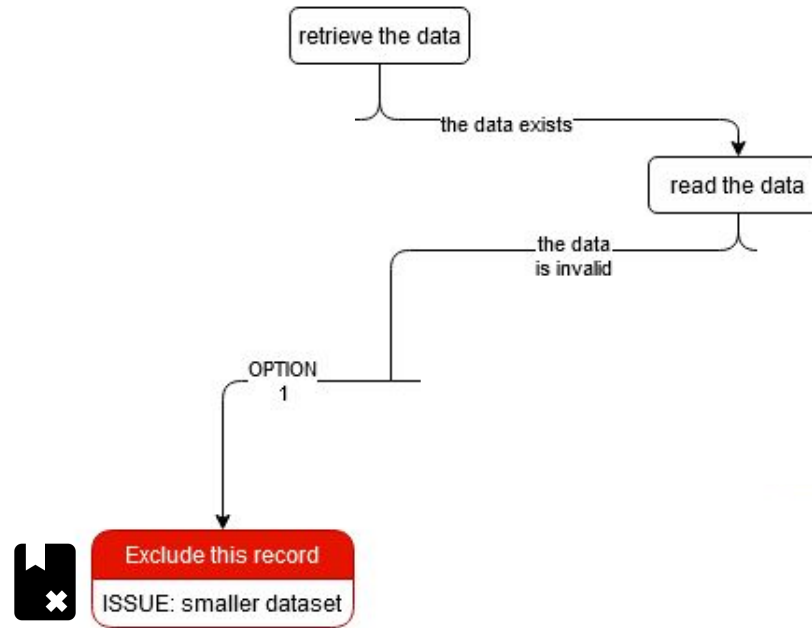retrieve the data

the data exists

read the data

# Handling manually-entered data

retrieve the data

the data exists

read the data

a data-entry error
may or may not
be discovered

# Handling manually-entered data



retrieve the data

the data exists

read the data
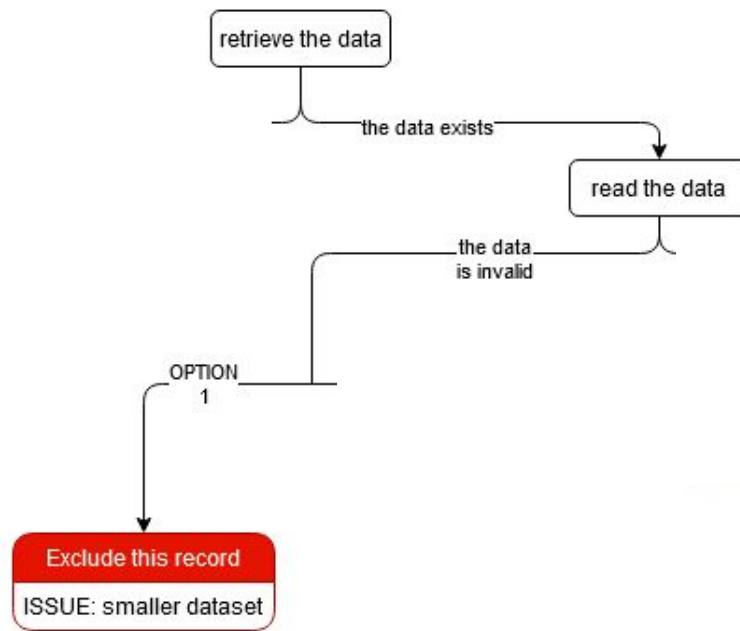
the data
is invalid

OPTIONS:
data-exclusion
or
data-correction
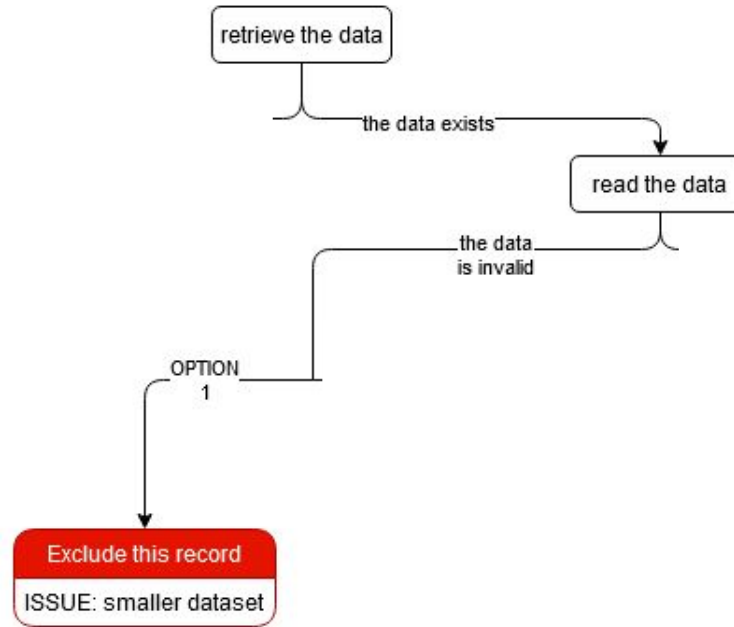
# Handling manually-entered data

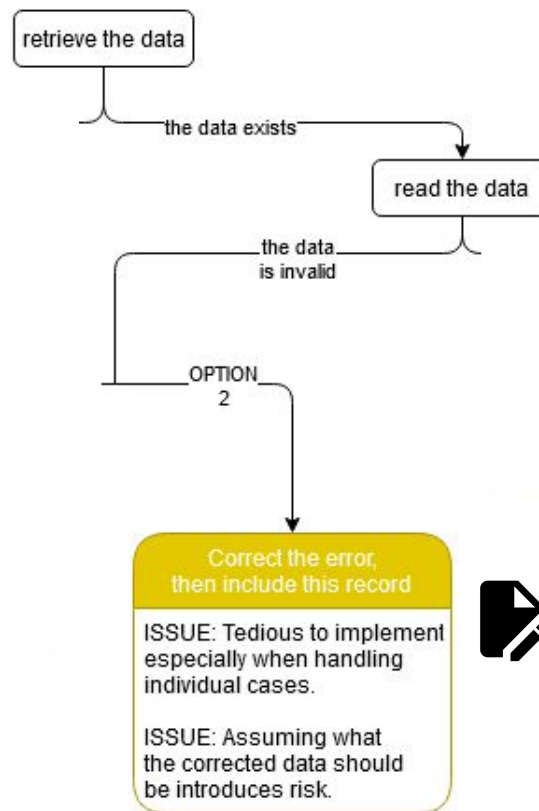# Handling manually-entered data - example of data exclusion



**"North America / USA / LA"**
could refer to Los Angeles or Louisiana. It is not specified whether the third piece is the state or the city.
**ambiguous → exclude**

# Handling manually-entered data - example of data exclusion

retrieve the data

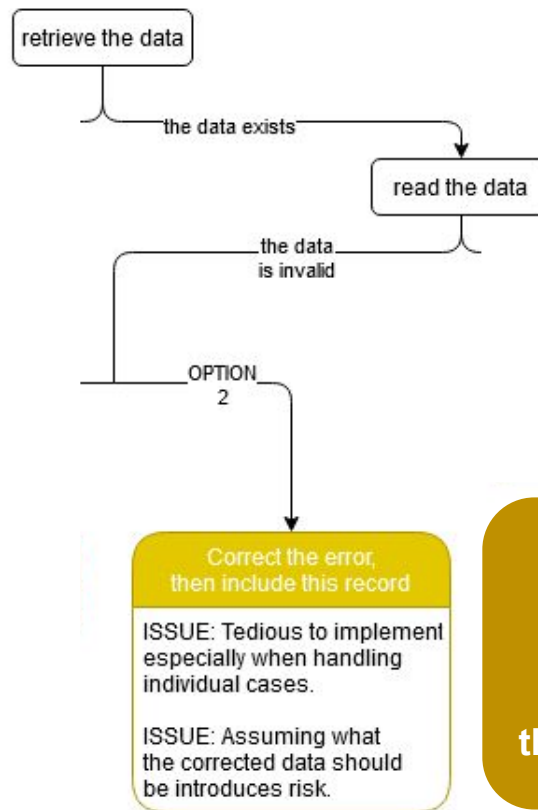the data exists

read the data

the data is invalid

OPTION 1

Age data is present in the record, but a value such as "unknown" or "not applicable" was provided.
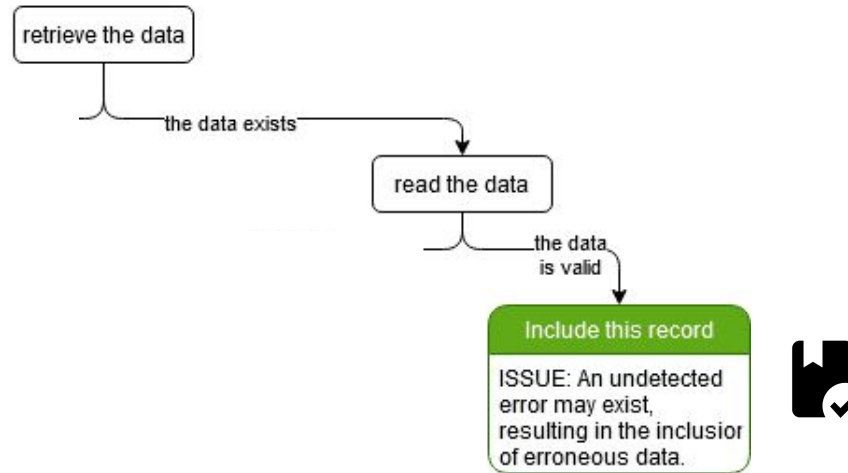→ **exclude**

Exclude this record

ISSUE: smaller dataset

# Handling manually-entered data

# Handling manually-entered data - example of data correction



retrieve the data

the data exists

read the data

the data is invalid

OPTION 2

Correct the error, then include this record

ISSUE: Tedious to implement especially when handling individual cases.

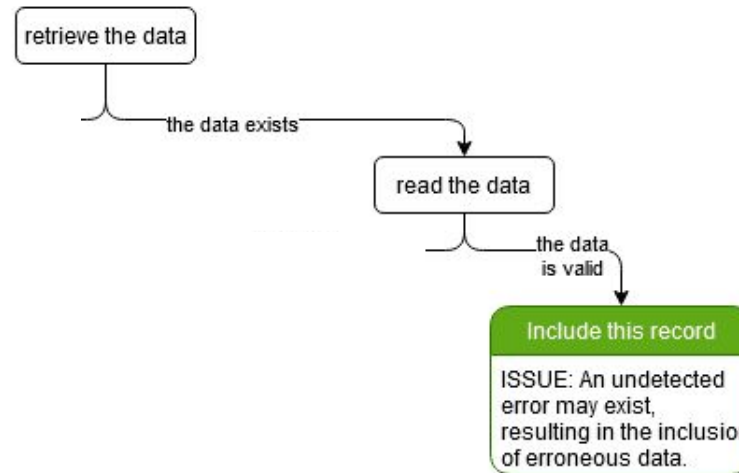ISSUE: Assuming what the corrected data should be introduces risk.

Location is provided as **"Oceania / Australia / NSW"** → **assume "NSW" means "New South Wales", correct the value and include the data**

# Handling manually-entered data

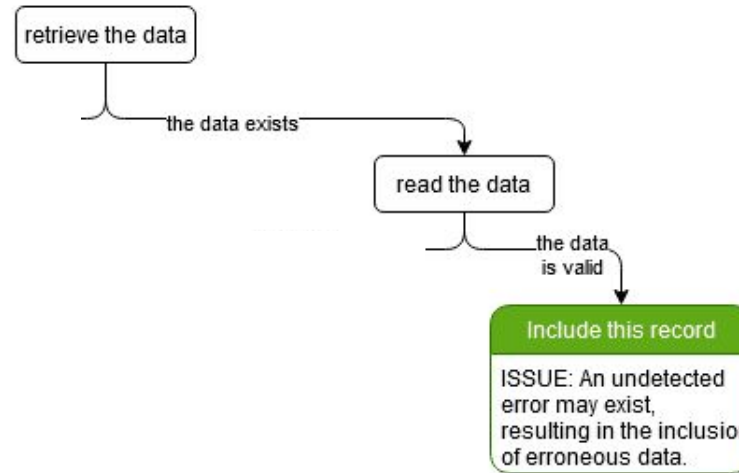# Handling manually-entered data - example of data inclusion

retrieve the data

the data exists

read the data

the data is valid

Include this record

ISSUE: An undetected error may exist, resulting in the inclusion of erroneous data.

**"Asia / China / Wuhan"**
Wuhan gets categorized as a state in China, but is really a city in Hubei, China.
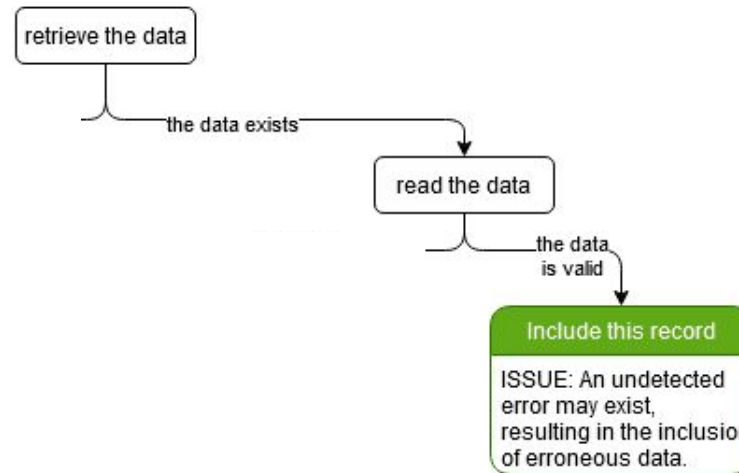→ **incorrect data was included due to unclear formatting**

# Handling manually-entered data - example of data inclusion

retrieve the data

the data exists

read the data

the data is valid

Include this record

ISSUE: An undetected error may exist, resulting in the inclusion of erroneous data.

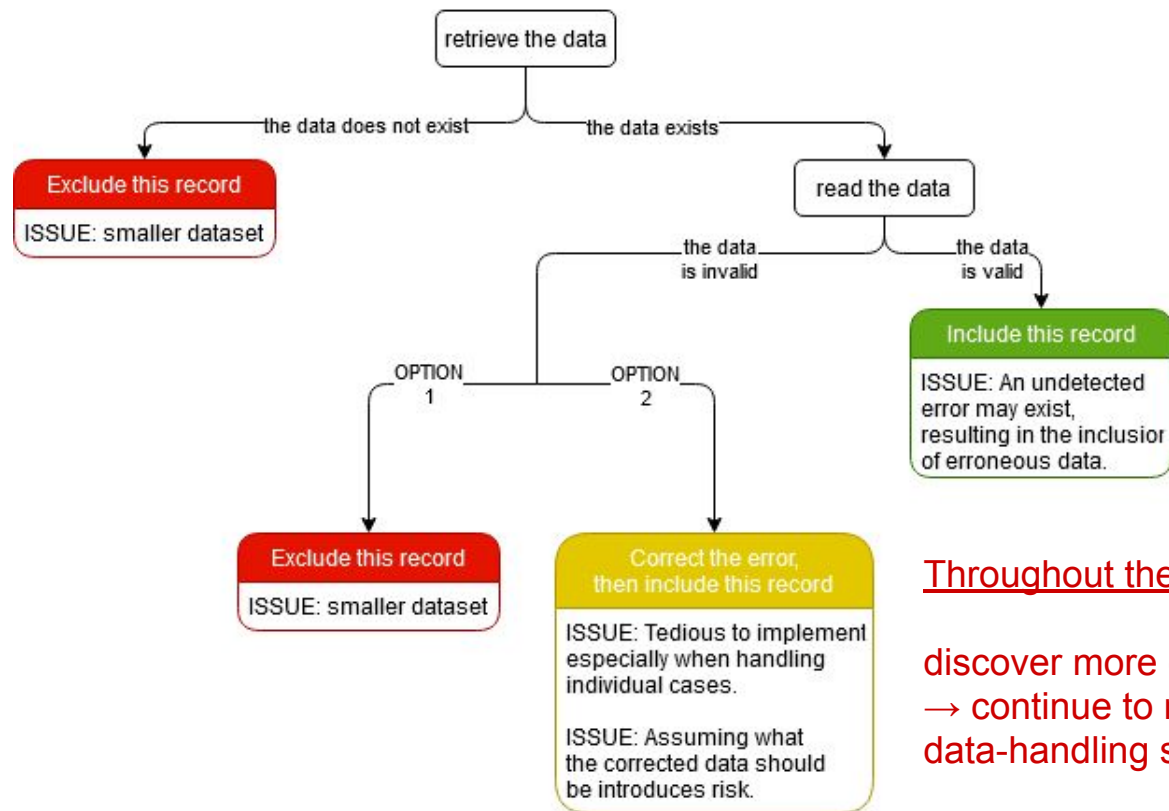**"Washington DC" and "District of Columbia"** are stored as two separate categories.
→ **dataset becomes misleading due to inconsistent naming conventions**

# Handling manually-entered data - example of data inclusion



retrieve the data

the data exists

read the data

the data is valid

**Include this record**

ISSUE: An undetected error may exist, resulting in the inclusion of erroneous data.

**"Oceania / <u>Australia</u> /
New South <u>Wales</u> / Sydney"**
Sydney gets categorized into both
"Australia" and "Wales".
→ **incorrect data was included due to
the categorization methods used**

# Handling manually-entered data



**Throughout the study**

discover more error cases
→ continue to modify the
data-handling strategies

# Potential solutions for improving data-entry practices

**Improving Manual Data Entry**

- The use of dropdown menus instead of allowing the user to input information
- Adding Data Validation will ensure that all information being added is formatted correctly

**Automated Data Capture instead of Manual Data Entry.**

Benefits of automation:
- Significantly reduces errors
- Improved efficiency and cleaner data

# RESULTS
## ---
# Analyzing the dataset

# Results - Analyzing the dataset

What did our collection of 75,571 patient records look like?

Attributes that we looked at

**Location**
Where the data was submitted from
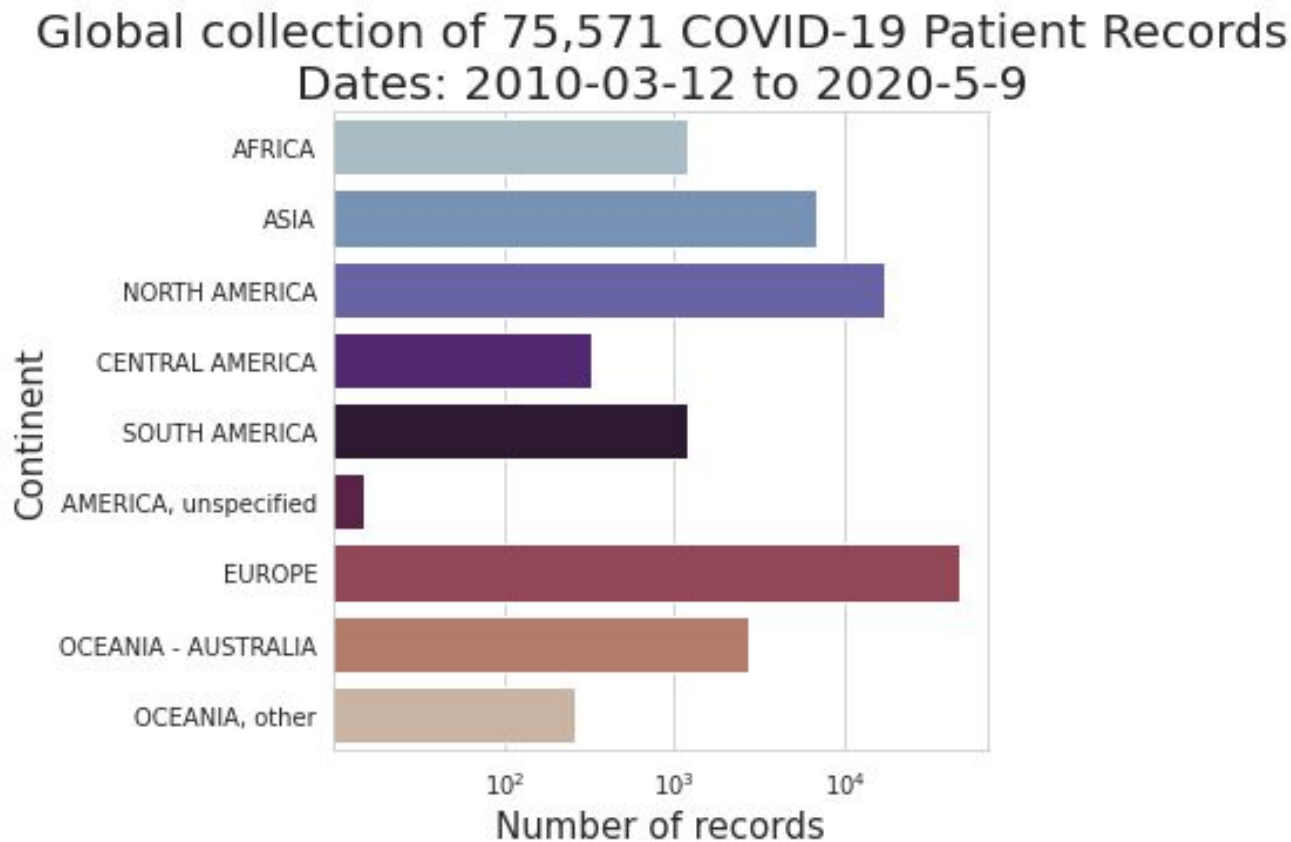
→ Continent
→ Country
→ State

**Gender**

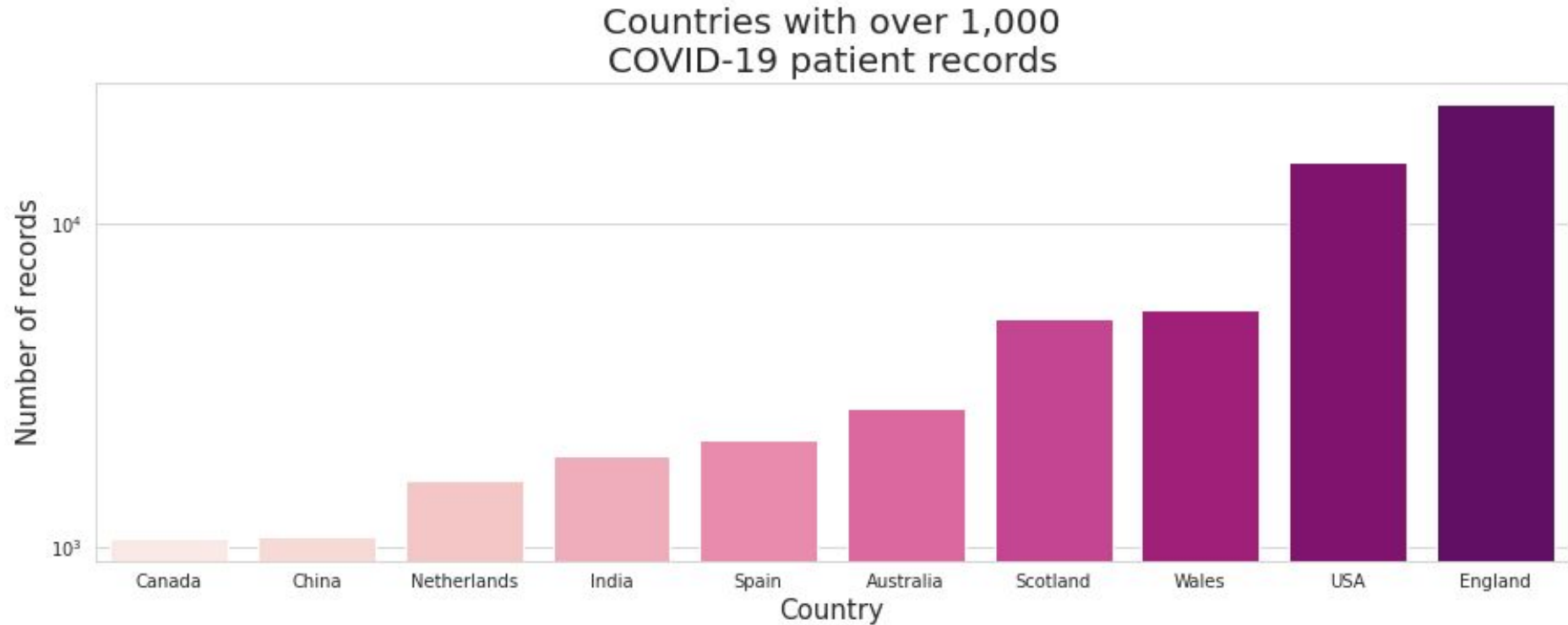→ The gender data that was provided

**Age**

→ The age in years of the patient

Global collection of 75,571 COVID-19 Patient Records
Dates: 2010-03-12 to 2020-5-9

**slide 48**
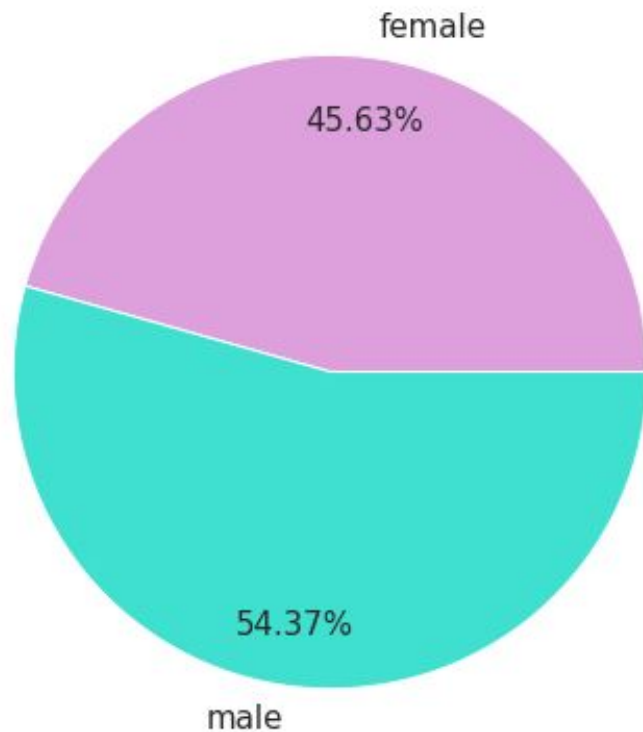
# 📍 Which countries have submitted over 1,000 records?



Countries with over 1,000 COVID-19 patient records

# What genders were found in the dataset?



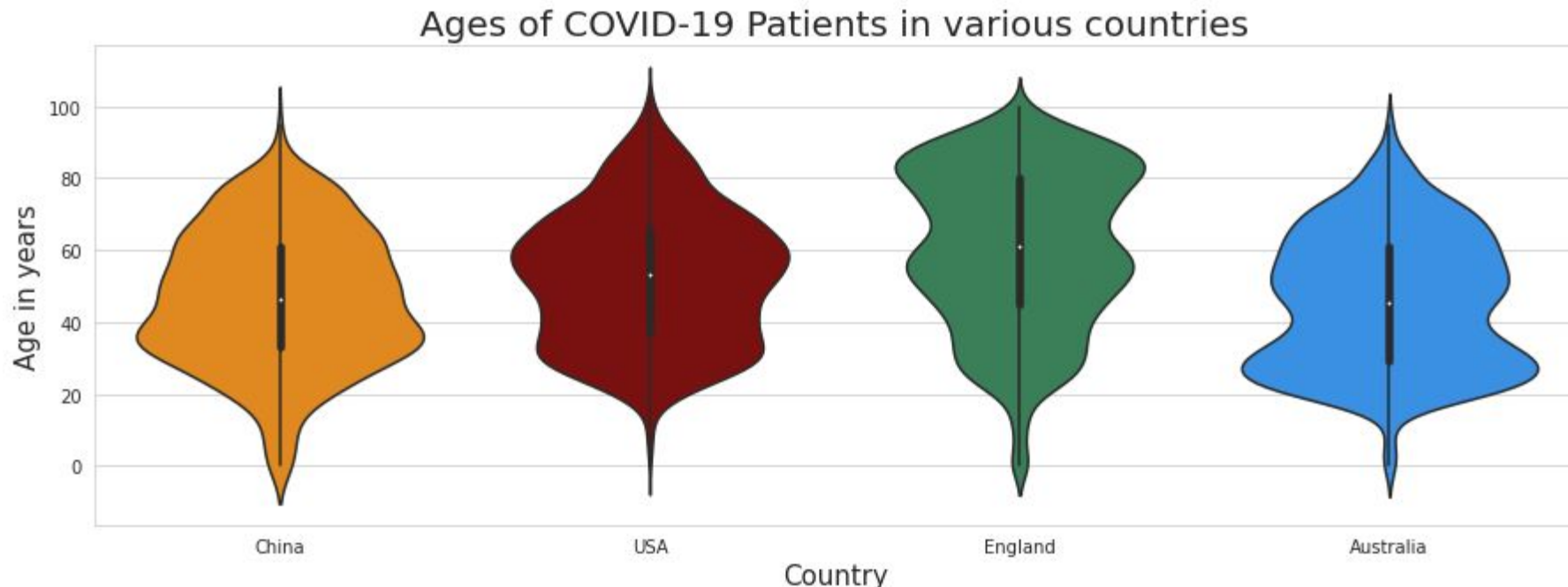Gender data found in 23,572
out of 75,571 records (31.19%)

female
45.63%

male
54.37%

Country categorized by Gender

Out of the 42,615 patient records from these countries, 55.0% held either female of male gender data.
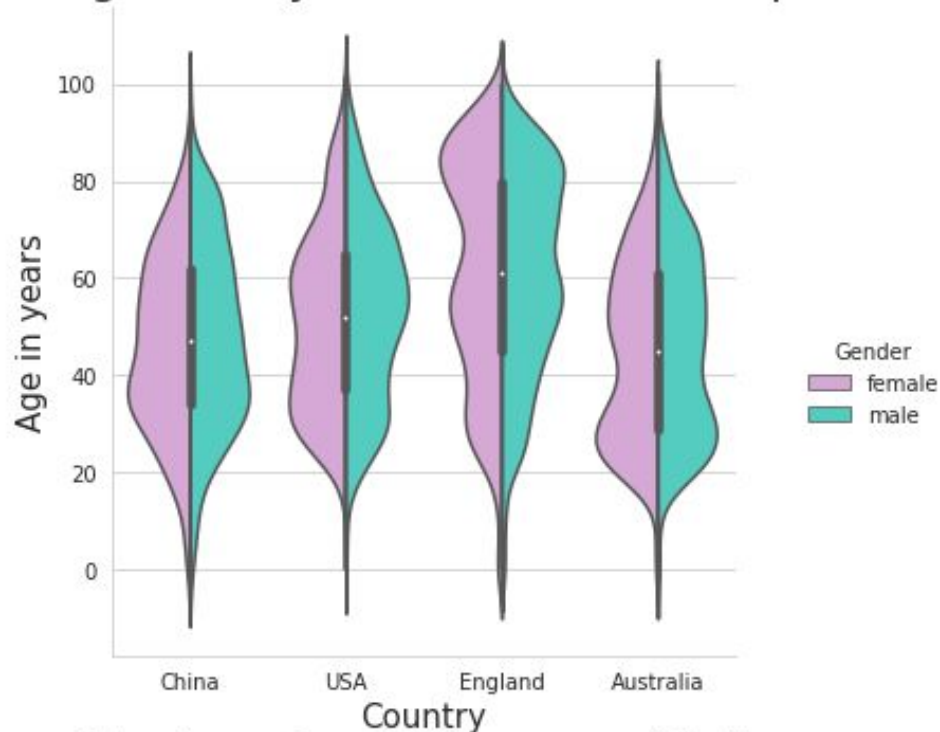
# What was the age data in various countries?



Ages of COVID-19 Patients in various countries

Country
Note: Age data was obtained by
9,689 out of 42,830 records in these countries (22.62%)

Age categorized by Gender of COVID-19 patients
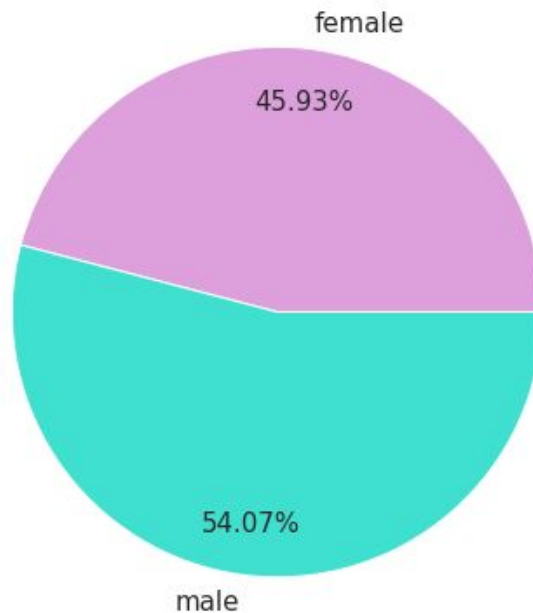
Note: Age and gender data was provided by
9,457 out of 42,830 records in these countries (22.08%)
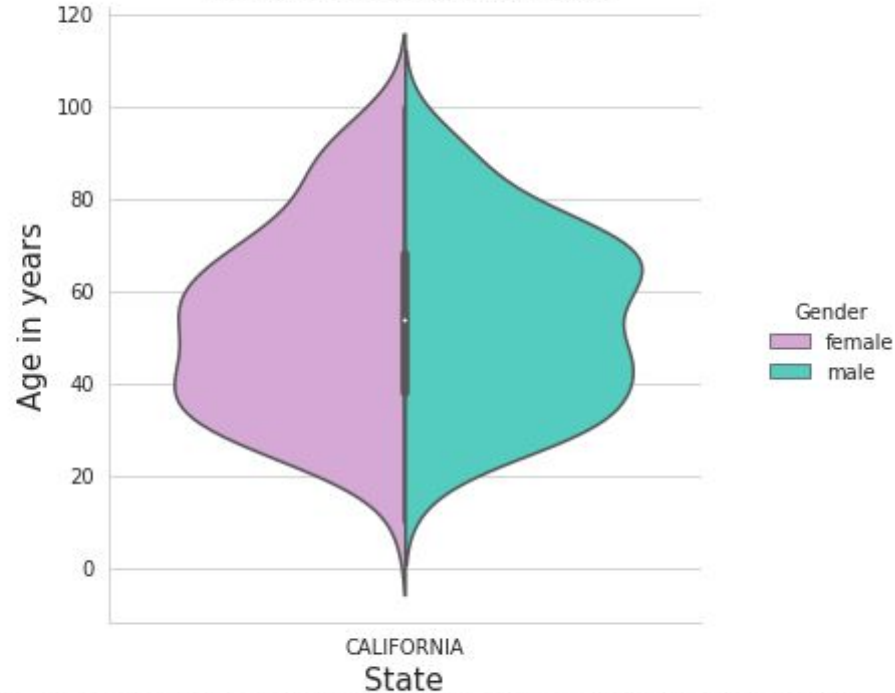
Genders of Patients in CALIFORNIA, USA
Note: Gender data was provided by
246 out of 2,176 records
in CALIFORNIA, USA (11.31%)

female 45.93%

male 54.07%

Age categorized by Gender of COVID-19 patients in CALIFORNIA, USA

Note: Age and gender provided by 246 out of 2,176 records in CALIFORNIA, USA (11.31%)

# Lessons Learned - Biomedical Research

**Reading patient data**

Data analysis is affected by the data-entry methods that were used to create the data.

Confounding errors such as the lack of correct data entry can cause misinterpretation.

**Analysis and visualization**

In order to generalize inferences about a sample of patients, the differences in the patients should be accounted for instead of treating each patient as if they are the same person.

# Lessons Learned - Data Science

**Reading patient data**

Loading and reading a collection of data from a file.

Identifying what information should be included in analysis and visualization.

**Analysis and visualization**

Choosing the appropriate visualization types to show trends in data.

Using tools to generate the visualizations.

# Next Steps

**Improve data-handling strategies in our code**

Current focus: error-handling for location data

Example: Use pre-defined lists of all states in a country to prevent misrepresentation

Improve **efficiency**, **readability**, and **reusability** of our code so that our methods can be used to assist with the study other datasets

Advocate for **improvement in data-entry practices**

# Acknowledgements

**DBMI Summer 2020 Internship Program Leaders**

Nancy Herbst

Dr. Jejo Koola

Dr. Tsung-Ting (Tim) Kuo

Dr. Lucila Ohno-Machado

Elizabeth Santillanez

**Project Mentors**

Dr. Niema Moshiri | Dr. Youwen Ouyang

# Thank you for your time!

## Any questions are welcome.

## Mhealyssah Bustria
bustr003@cougars.csusm.edu

## Anjelina Velazquez
velaz035@cougars.csusm.edu