# COSC 580 Project 2

Group Member: Yuming Wang, Ziang He, Jiancheng Sun
Group Number: 6

## 1. Source Identification

For this project, we used three datasets to build our Database System:
- Dataset_0: [United States COVID-19 Cases, Deaths, and Laboratory Testing (NAATs) by State, Territory, and Jurisdiction](#)
- Dataset_1: [Data Table for Vaccinations Equity (SVI)](#)
- Dataset_2: [Data Table for Vaccinations Equity (Metro/Non-Metro)](#)

### 1.1. Dataset_0

This dataset shows the number of COVID-19 **cases, deaths and laboratory testing** for every 100,000 people over the last 7 days, allowing us to compare areas with different population sizes.

### 1.2. Dataset_1 and Dataset_2

These two datasets provide a county-level view of COVID-19 vaccination coverage, social vulnerability and Metropolitan vs. Non-Metropolitan:

***Social vulnerability*** is measured by CDC Social Vulnerability Index (SVI), which uses U.S. Census data on categories like poverty, housing, and vehicle access to estimate a community's ability to respond to and recover from disasters or disease outbreaks.

***Metropolitan vs. Non-Metropolitan*** classification is based off an aggregation of the six 2013 National Center for Health Statistics (NCHS) Urban-Rural classifications, where "Metro" counties include Large Central Metro, Large Fringe Metro, Medium Metro, and Small Metro and "Non-Metro" counties include Micropolitan and Non-Core (Rural).

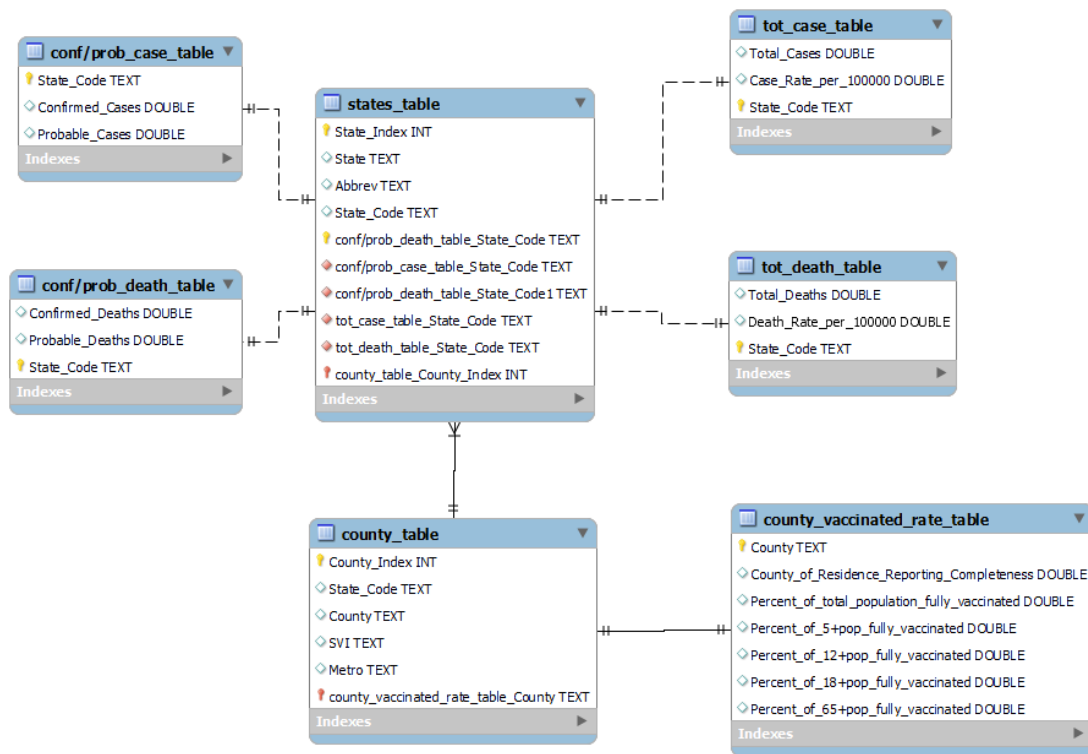### 1.3. Data Preprocessing and Table Building

For these datasets, we do some preprocessing to build our own database.
- For these datasets, we do some preprocessing to build our database. Although the datasets we collected contain ***regional information***, such as state names, some are abbreviated, and some are full. For this problem, we matched the abbreviation of the state name with the full name of the state name and constructed the "***states_table***".
- In addition, for the data of "Total_Cases", "Total_Deaths", "Confirmed_Deaths", "Probable_Deaths", "Confirmed_Cases", "Probable_Cases" of each state in Dataset_0, we have established different tables to store them to meet the requirements of ***3NF***. They all use "State_Code" as the ***primary key***.
- The data in Dataset_1 and Dataset_2 ***overlap with multiple attributes***, which provide a county-level view of COVID-19 vaccination coverage from social vulnerability and Metropolitan vs. Non-Metropolitan, respectively. We extracted "SVI", "Metro", "County" and "State_Code" from these two datasets and defined the "***County_Index***" as the

primary key to build the "***county_table***".

- In addition, for the ***overlapping parts*** of the two data: the proportion of fully vaccinated people in different age groups and the residence reporting completeness of each county, we extracted and constructed the "***county_vaccinated_rate_table***", and used "***county***" as the primary key.

## 2. Schema Definition



We used *Entity Relationship Diagram (ERD)* to present the relations and relationships. Every relation is in *3NF*, which means every non-key attribute is dependent on the key and nothing but the key.

## 3. Functional Dependencies

### 3.1. county_table

$$FD: \{County\_Index\# \rightarrow State\_Code, County, SVI, Metro\}$$

## 3.2. states_table

$$FD : \{State\_Index\# \rightarrow State\_Code, Abbrev, State\}$$



## 3.3. county_vaccinated_rate_table

$FD : \{County\# \rightarrow County\_of\_Residence\_Reporting\_Completeness,$

Percent_of_total_population_fully_vaccinated,

Percent_of_5+pop_fully_vaccinated,

Percent_of_12+pop_fully_vaccinated,

Percent_of_18+pop_fully_vaccinated,

$Percent\_of\_65+pop\_fully\_vaccinated\}$

**county_vaccinated_rate_table**
- County TEXT
- County_of_Residence_Reporting_Completeness DOUBLE
- Percent_of_total_population_fully_vaccinated DOUBLE
- Percent_of_5+pop_fully_vaccinated DOUBLE
- Percent_of_12+pop_fully_vaccinated DOUBLE
- Percent_of_18+pop_fully_vaccinated DOUBLE
- Percent_of_65+pop_fully_vaccinated DOUBLE

Indexes

### 3.4. tot_case_table

$FD : \{State\_Code\# \rightarrow Total\_Cases, Case\_Rate\_per\_100000\}$

**tot_case_table**
- Total_Cases DOUBLE
- Case_Rate_per_100000 DOUBLE
- State_Code TEXT

Indexes

### 3.5. tot_death_table

$FD : \{State\_Code\# \rightarrow Total\_Deaths, Death\_Rate\_per\_100000\}$

**tot_death_table**
- Total_Deaths DOUBLE
- Death_Rate_per_100000 DOUBLE
- State_Code TEXT

Indexes

### 3.6. conf/prob_case_table

$FD : \{State\_Code\# \rightarrow Confirmed\_Caes, Probable\_Cases\}$

| | |
|---|---|
| **conf/prob_case_table** ▼ | |
| 🔑 State_Code TEXT | |
| ◇ Confirmed_Cases DOUBLE | |
| ◇ Probable_Cases DOUBLE | |
| **Indexes** ► | |

### 3.7. conf/prob_death_table

$FD : \{State\_Code\# \rightarrow Confirmed\_Deaths, Probable\_Deaths\}$

| | |
|---|---|
| **conf/prob_death_table** ▼ | |
| ◇ Confirmed_Deaths DOUBLE | |
| ◇ Probable_Deaths DOUBLE | |
| 🔑 State_Code TEXT | |
| **Indexes** ► | |

# 4. Canned Queries

In this database, we can support some query functions.

Some examples of retrieving some **_extreme values_** are given below:

Eg1:

| Query | select * from tot_case_table where Total_Cases = (select max(Total_Cases) from tot_case_table); |
|---|---|
| Result | State_Code: CA, Total_Cases: 9015587, Case_Rate_per_100000: 22817 |

Eg2:

| Query | select * from tot_case_table where Total_Cases = (select min(Total_Cases) from tot_case_table); |
|---|---|
| Result | State_Code: VT, Total_Cases: 105475, Case_Rate_per_100000: 16903 |

Eg3:

| Query | select*from tot_death_table where Total_Deaths = (select max(Total_Deaths) from tot_death_table); |
|---|---|
| Result | State_Code: CA, Total_Deaths: 86185, Death_Rate_per_100000: 218 |

Eg4:

| Query | select*from tot_death_table where Total_Deaths = (select min(Total_Deaths) from tot_death_table); |
|---|---|
| Result | State_Code: VT, Total_Deaths: 578, Death_Rate_per_100000: 92 |

Some examples of **_range searches_** are given below:

Eg5:

| Query | SELECT State_Code, Case_Rate_per_100000 FROM tot_case_table order by Case_Rate_per_100000 desc limit 10; |
|---|---|

| Result | | State_Code | Case_Rate_per_100000 | |
|---|---|---|---|---|
| | ► | AK | 32105 | |
| | | RI | 32086 | |
| | | ND | 31368 | |
| | | TN | 29510 | |
| | | KY | 28979 | |
| | | UT | 28861 | |
| | | SC | 28444 | |
| | | WV | 27587 | |
| | | AR | 27337 | |
| | | AZ | 27303 | |

Eg6:

| Query | SELECT State_Code, Case_Rate_per_100000 FROM tot_case_table where Case_Rate_per_100000 > 27000; |
|---|---|

| Result | | State_Code | Case_Rate_per_100000 | |
|---|---|---|---|---|
| | ► | AK | 32105 | |
| | | AZ | 27303 | |
| | | AR | 27337 | |
| | | FL | 27120 | |
| | | KY | 28979 | |
| | | ND | 31368 | |
| | | RI | 32086 | |
| | | SC | 28444 | |
| | | TN | 29510 | |
| | | UT | 28861 | |
| | | WV | 27587 | |
| | | WI | 27088 | |

Eg7:

| Query | SELECT State_Code, Death_Rate_per_100000 FROM tot_death_table where Death_Rate_per_100000 > 350; |
|---|---|

| Result | | State_Code | Death_Rate_per_100000 | |
|---|---|---|---|---|
| | ► | AL | 384 | |
| | | AZ | 385 | |
| | | AR | 359 | |
| | | LA | 362 | |
| | | MI | 352 | |
| | | MS | 411 | |
| | | NJ | 372 | |
| | | TN | 363 | |
| | | WV | 365 | |

Here are some examples of *sorting by value*:

Eg8:

| Query | select State_Code, (count(case when SVI='High' then State_Code end)/count(State_Code)) as rate from county_table group by State_Code order by rate desc limit 10; |
|---|---|

| Result | |
|---|---|
| | State_Code | rate |
| | AZ | 0.7333 |
| | LA | 0.7031 |
| | MS | 0.6707 |
| | NM | 0.6364 |
| | GA | 0.5786 |
| | SC | 0.5435 |
| | AR | 0.5333 |
| | OK | 0.4805 |
| | AL | 0.4478 |
| | NC | 0.4400 |

Eg9:

| Query | select State  Code,rate<br>from(select  State_Code,  (count(case  when  SVI='High'  then  State_Code end)/count(State_Code)) as rate from county_table<br>group by State_Code<br>order by rate desc) as a<br>where rate > 0.5; |
|---|---|

| Result | |
|---|---|
| | State_Code | rate |
| | AZ | 0.7333 |
| | LA | 0.7031 |
| | MS | 0.6707 |
| | NM | 0.6364 |
| | GA | 0.5786 |
| | SC | 0.5435 |
| | AR | 0.5333 |

Eg10:

| Query | select State_Code,rate<br>from(select  State Code,  (count(case  when  SVI='High'  then  State Code end)/count(State_Code)) as rate from county_table<br>group by State_Code<br>order by rate desc) as a<br>where rate > 0.3; |
|---|---|

| Result | | State_Code | rate |
|--------|--|------------|------|
| | | AZ | 0.7333 |
| | | LA | 0.7031 |
| | | MS | 0.6707 |
| | | NM | 0.6364 |
| | | GA | 0.5786 |
| | | SC | 0.5435 |
| | | AR | 0.5333 |
| | | OK | 0.4805 |
| | | AL | 0.4478 |
| | | NC | 0.4400 |
| | | FL | 0.4328 |
| | | TX | 0.4291 |
| | | CA | 0.3966 |
| | | KY | 0.3583 |

Eg11:

| Query | select State_Code,Total_Deaths/Total_Cases as death_rate from tot_case_table join tot_death_table using(State_Code) order by death_rate desc limit 10; |
|-------|------|

| Result | | State_Code | death_rate |
|--------|--|------------|------------|
| | | PA | 0.015815759965569366 |
| | | MS | 0.015455774332345899 |
| | | NJ | 0.0151754769212694441 |
| | | MI | 0.014831186743799074 |
| | | AL | 0.014609718348648763 |
| | | GA | 0.0145946922002077716 |
| | | CT | 0.014553975161581344 |
| | | LA | 0.01446299422235031 |
| | | NV | 0.014392601664567291 |
| | | MD | 0.014151215332653455 |

Eg12:

| Query | select Metro,SVI,count(*) as group_number from county_table group by Metro,SVI; |
|-------|------|

| Result | | Metro | SVI | group_number |
|--------|--|-------|-----|--------------|
| | | Metro | Low-Mod | 344 |
| | | Metro | Low | 341 |
| | | Non-metro | High | 584 |
| | | Metro | Mod-High | 329 |
| | | Metro | High | 221 |
| | | Non-metro | Mod-High | 475 |
| | | Non-metro | Low-Mod | 459 |
| | | | | 2 |
| | | Non-metro | Low | 465 |
| | | | Low-Mod | 1 |
| | | Non-metro | | 1 |