

# Model

## Dataset Manipulation:

- The data on total vaccines, people vaccinated, and total booster given was deleted on all the specific country datasets because this same data can be found on the vaccinations dataset.
- A new spreadsheet of vaccines used per country was created from the locations csv because there was list of multiple vaccines in each cell. The vaccines column was then deleted from the locations csv
- All the specific country csvs are combined into one called countries\_with\_source. The list of vaccines in the vaccines column will be made into a new spreadsheet called countries\_vaccines\_by\_day and will be deleted from the countries\_with\_source csv. This new csv will contain the location date and vaccine used on that date with multiple rows for each location and date if there are multiple vaccines used on each day.
- All the per\_capita columns except the one without missing data was deleted from the vaccinations and us\_state\_vaccinations as per capita information for other information can be calculated using the kept per capita and the information and the other columns
- Iso\_code is removed from all csv except the locations csv because location and iso\_code is not needed as they both describe the country or location.
- For us\_state\_vaccinations, changed existing location column name to state
- Added missing locations in locations because some are in vaccinations that are not in locations
- Add location US in state vaccinations csv to be used as foreign key

## Initial Schema (after initial dataset manipulation):

**Vaccinations**(location\*, date, total\_vaccination, people\_vaccinated, people\_fully\_vaccinated, total\_boosters, daily\_vaccinations\_raw, daily\_vaccinations, daily\_vaccinations\_per\_million)

**Vaccinations-by-age-group**(location\*, date, age\_group, people\_vaccinated\_per\_hundred, people\_fully\_vaccinated\_per\_hundred, people\_with\_booster\_per\_hundred)

**Vaccinations-by-manufacturer**(location\*, date, vaccine, total\_vaccinations)

**Locations**(location\*, iso\_code, last\_observation\_date, source\_name, source\_website)

**Locations\_vaccines**(location\*, vaccine)

**Countries\_vaccines\_by\_day**(location\*, date, vaccine)

**countries\_source\_by\_day**(location\*, date, source\_url)

**Us\_state\_vaccinations**(date\*, state, total\_vaccinations, total\_distributed, people\_vaccinated, people\_fully\_vaccinated, daily\_vaccinations\_raw, daily\_vaccinations, share\_doses\_used, total\_boosters, location\*)

### Normalisation:

Vaccinations is in third normal form

Vaccinations\_by\_age\_group is in third normal form

Locations:

Location → iso\_code, last\_observation, source\_url

source\_url → source\_name

Locations\_vaccines is in third normal form

Countries\_vaccines\_by\_day is in third normal form

countries\_source\_by\_day is in third normal form

Us\_state\_vaccinations is in third normal form

Vaccinations\_by\_manufacturer is in third normal form

### Final Schema:

**Vaccinations**(location\*, date\*, total\_vaccination, people\_vaccinated, people\_fully\_vaccinated, total\_boosters, daily\_vaccinations\_raw, daily\_vaccinations, daily\_vaccinations\_per\_million)

**Vaccinations-by-age-group**(location\*, date\*, age\_group, people\_vaccinated\_per\_hundred, people\_fully\_vaccinated\_per\_hundred, people\_with\_booster\_per\_hundred)

**Vaccinations-by-manufacturer**(location\*, date\*, vaccine, total\_vaccinations)

**Locations**(location\*, iso\_code, last\_observation\_date, source\_url\*)

**Source**(location\_url\*, source\_name)

**Locations\_vaccines\_used**(location\*, vaccine)

**countries\_vaccines\_by\_day**(location\*, date\*, vaccine)

**countries\_source\_by\_day**(location\*, date\*, source\_url\*)

**Us\_state\_vaccinations**(date, state, total\_vaccinations, total\_distributed, people\_vaccinated, people\_fully\_vaccinated, daily\_vaccinations\_raw, daily\_vaccinations, share\_doses\_used, total\_boosters, location\*)

## ER Model:

