

What's killing Americans?

Data Understanding and Preparation

Team Adams

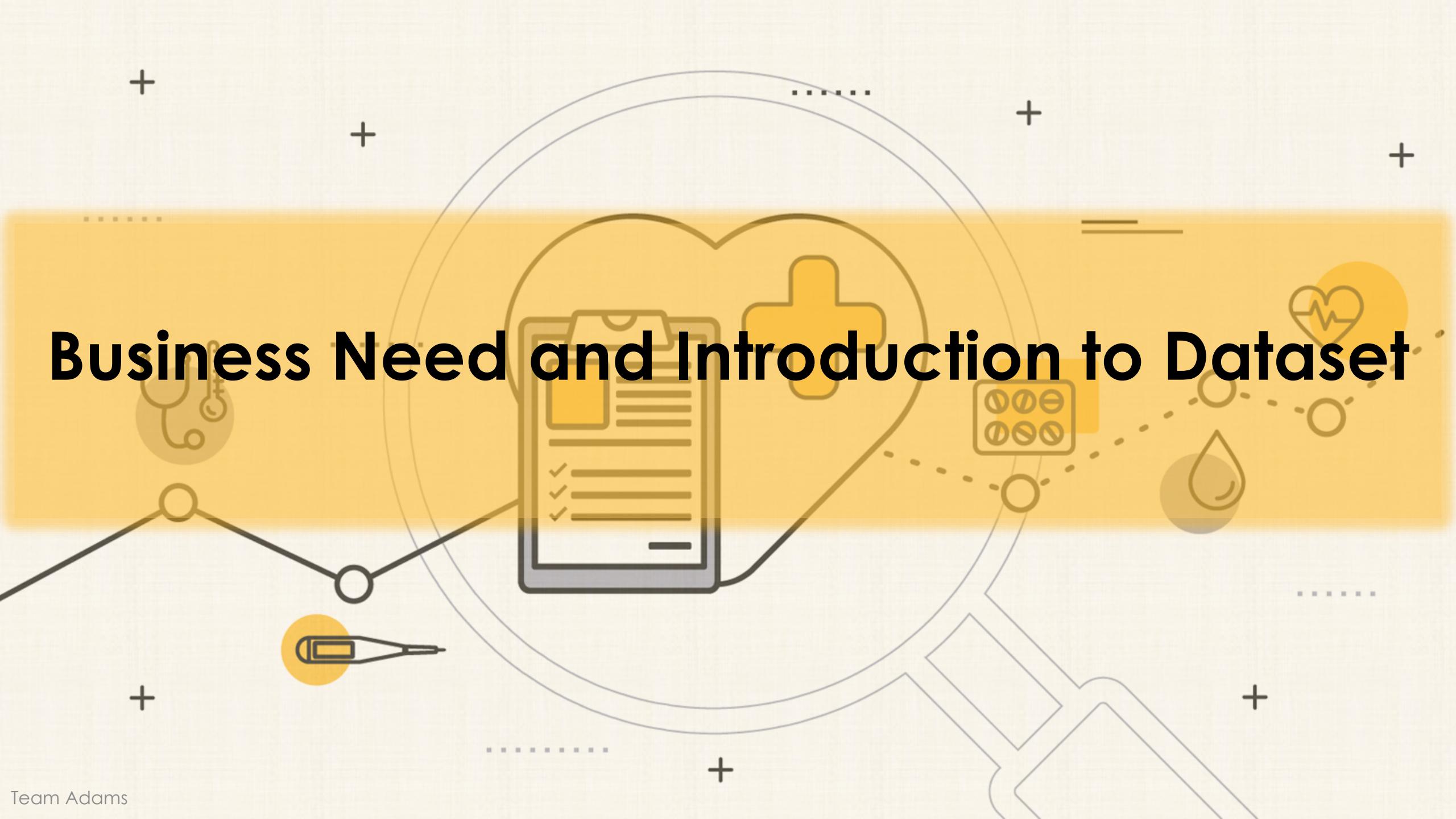
Team Members:

Pradeep Kumar M
Harsh Solanki
Priyanka Priya
Shiva R

Content

- Business Need and Introduction to Dataset
- Data Understanding:
 - Entity Relationship Diagram
 - Normalization
 - Uploading Data
- Analysis with SQL Queries
- Visualization
- Recommendations
- Lessons Learned

Business Need and Introduction to Dataset

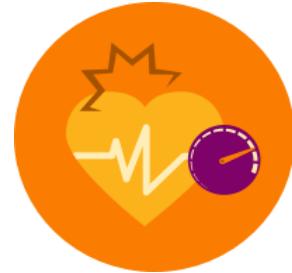


Business Need and Introduction to Dataset

This project is an attempt to analyze and recognize the correlation between various living factors and mortality due to major diseases in USA. These correlation will help governments in devising policies and initiatives at a County level



Comprehensive **Census Data** highlighting Demographics, Income, Nature of Employment, Mode of Transport, Population count granular to ethnicity at county level



Mortality Rate Data of Top 21 Categories corresponding to major diseases at County level



A collective Dataset of both Census and Mortality for years **2015 & 2017** to intersperse & capture the 'changing trends'



67,000+ records

A rich and complete collection of mortality rates corresponding to different categories for all counties across the USA

Data Understanding

Data Understanding

Census Data



Location



County, State

Demographic



Total Population, Men, Women, Hispanic, White, Black, Native, Asian, Citizen

Income



Income, Income Error, Income Per Capita, Income Per Capita Error, Poverty, Child Poverty

Mode of Transport



Drive, Carpool, Transit, Walk, Other Transport, Mean Commute, Work at Home

Employment Sector



Professional, Service, Office, Construction, Production

Employment Type



Private, Public, Family

Employment Status



Employed, Self Employed, Unemployed

Mortality Data



Diseases and other mortality categories



Cardiovascular, Chronic respiratory, Cirrhosis and other chronic liver diseases, Diabetes, urogenital, blood, and endocrine diseases, Diarrhea, lower respiratory, and other common infectious diseases, Digestive diseases, HIV/AIDS and tuberculosis



Maternal disorders, Mental and substance use disorders, Musculoskeletal disorders

Neglected tropical diseases and malaria

Neonatal disorders, Neoplasms, Neurological disorders



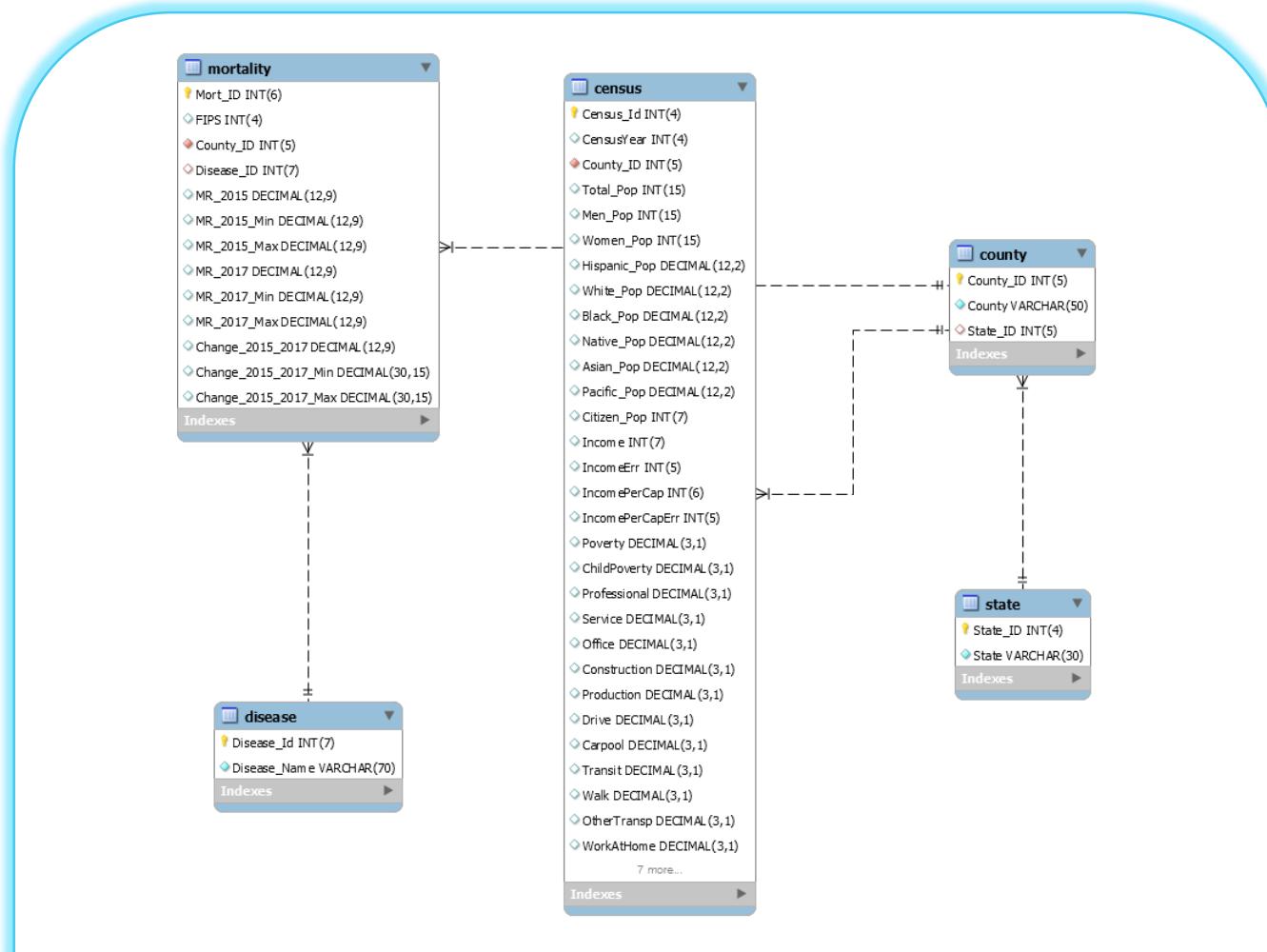
Nutritional deficiencies, Other communicable, maternal, neonatal, and nutritional diseases, Other non-communicable diseases



Self-harm and interpersonal violence, Transport injuries, Unintentional injuries, Forces of nature, war, and legal intervention

Data Understanding

1. Entity Relationship Diagram (EER)



Data Understanding

2. Normalization

1 st Normal form	
There are no repeating groups	Yes
Data atomic	Yes
Each field has unique name	Yes
Primary Key	Yes

2 nd Normal form	
First normal form	Yes
All non-key attributes are dependent	Yes
Eliminate partial dependencies	Yes

3 rd Normal form	
First normal form	Yes
2 nd normal form	Yes
Eliminate transitive dependencies	Yes

Partial Dependency

Census
CensusId
CensusYear
County_ID
County
State
Total_Pop

Morality
Mort_ID
FIPS
County_ID
Disease_ID
Disease
MR_2015

Transitive Dependency

County
County_ID
County
State
State_ID

State
State_ID
State
State

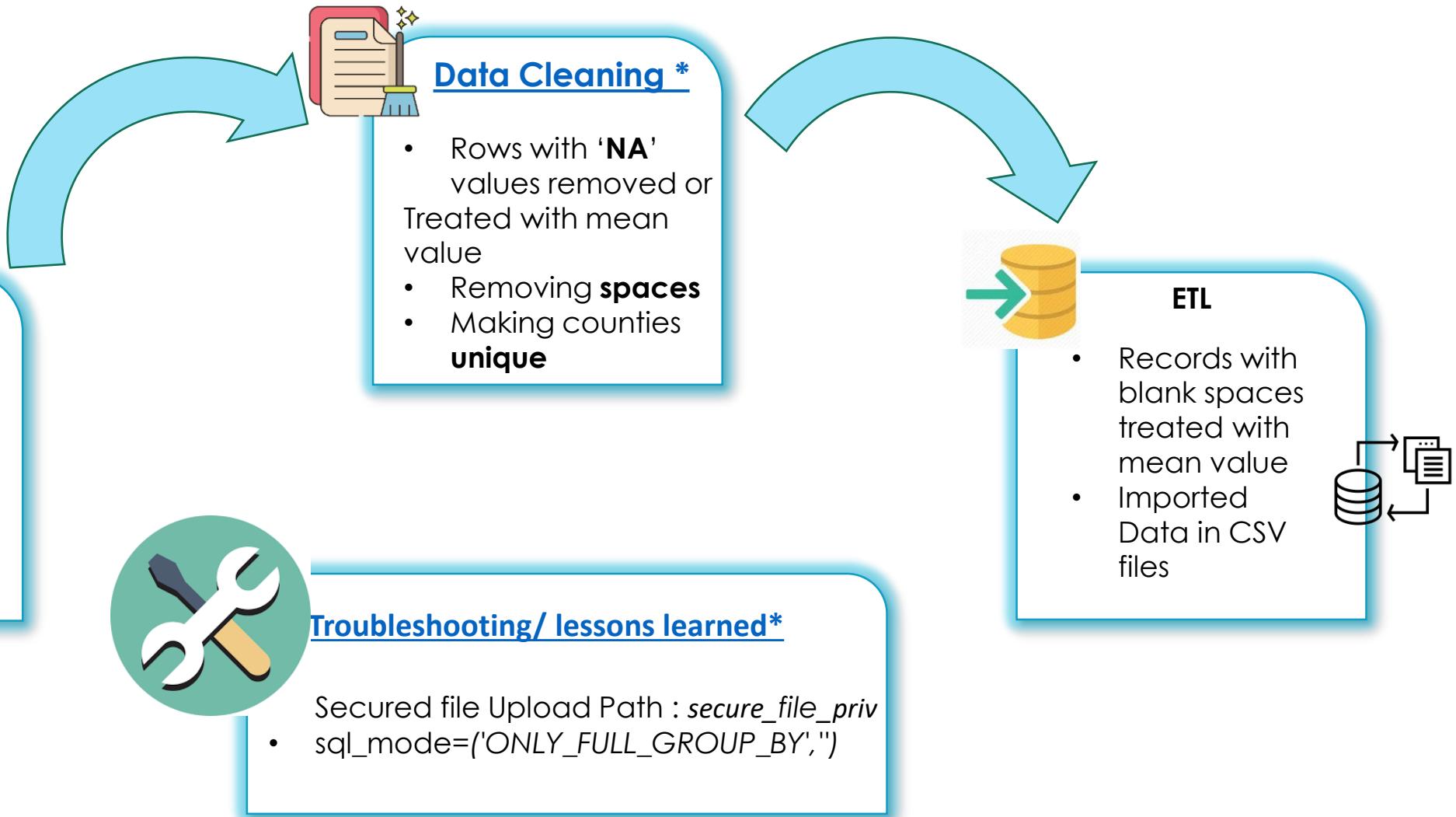
First Normal Form (1NF): Each row is Unique, and no columns contains multiple values. Table for Census was mostly downloaded in first normal form. For mortality table, we had to separate county and states. It was in following form (Autauga County, Alabama)

Second Normal Form (2NF): It is in 1NF and all non-key attributes are fully functional dependent on the primary key (does not contain partial functional dependencies)

Third Normal Form (3NF): It is in 2NF and non key columns do not determine other non key columns

Data Understanding

3. Uploading Data



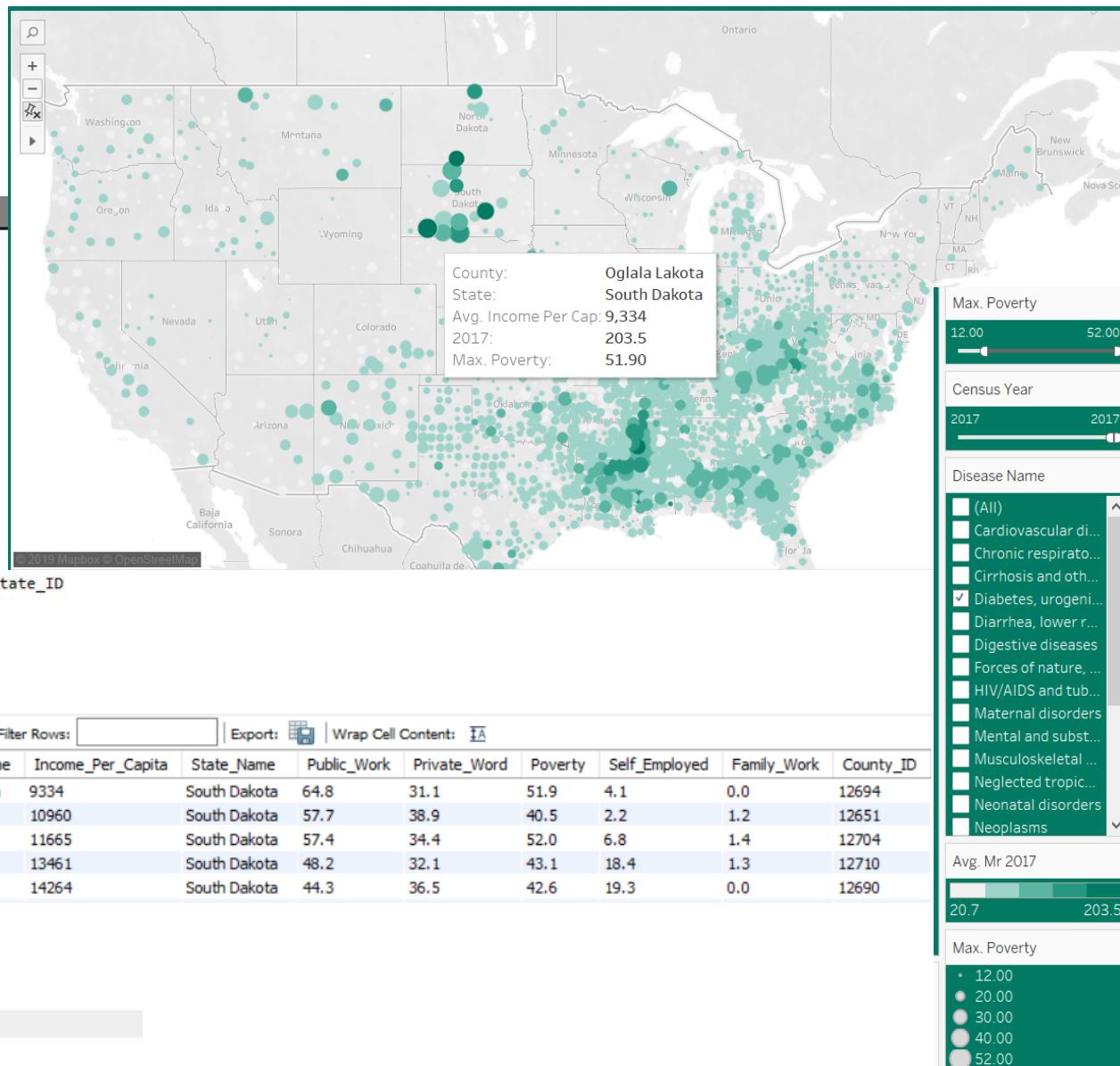
Analysis with SQL Queries

Public Jobs curse Oglala Lakota with Poverty and Diabetes!

SQL Query-1

In 2017, Oglala Lakota County records the highest deaths due to Diabetes(across USA) having the least per capita income(across South Dakota) with half the population under poverty and majority working for public sector

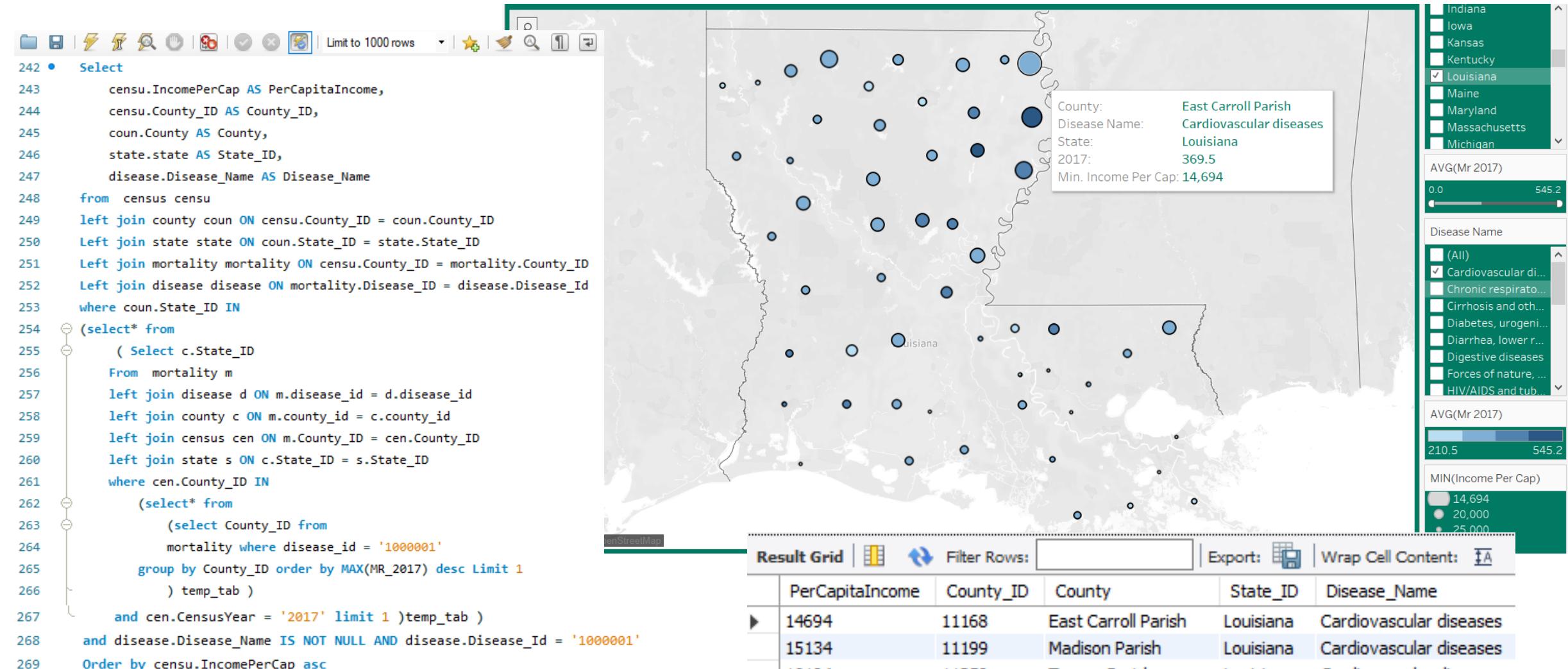
```
US_Final_Project_SQL_Prade... x
Limit to 2000 rows
Select censu.CensusYear AS Year,
       coun.County AS County_Name,
       censu.IncomePerCap AS Income_Per_Capita,
       state.state AS State_Name,
       censu.PublicWork AS Public_Work,
       censu.PrivateWork AS Private_Word,
       censu.poverty AS Poverty,
       censu.SelfEmployed AS Self_Employed,
       censu.FamilyWork AS Family_Work,
       censu.County_ID
  from census censu
 left join county coun ON censu.County_ID = coun.County_ID Left join state state ON coun.State_ID = state.State_ID
 where coun.State_ID IN
 (select* from ( Select c.State_ID
      From mortality m
      left join disease d ON m.disease_id = d.disease_id
      left join county c ON m.county_id = c.county_id
      left join census cen ON m.County_ID = cen.County_ID
      left join state s ON c.State_ID = s.State_ID
      where cen.County_ID IN
      (select* from (select County_ID
          from mortality
          where disease_id = '1000004'
          group by County_ID
          order by MAX(MR_2017) desc Limit 1 )
      temp_tab )and cen.CensusYear = '2017' limit 1 )temp_tab )
 and censu.CensusYear = '2017'
 Order by (censu.IncomePerCap) asc;
```



Franklin Parish suffers with Cardiovascular Diseases



In 2017, People in Franklin Parish, Louisiana saw the highest cardiovascular deaths. Interesting to note that nearby counties have similar mortality rate. East Carroll Parish, has the least per capita income in the same state

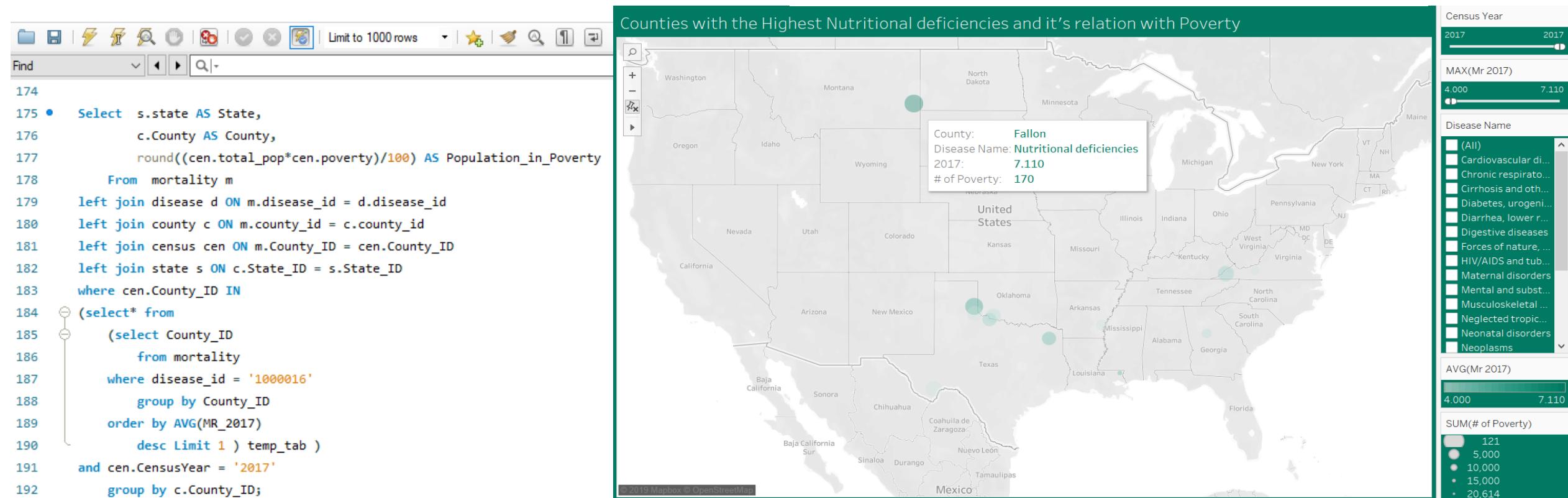


Fallon battles with Nutritional Deficiency



In the year 2017, Fallon from Montana has the greatest number of deaths due to 'Nutritional deficiencies' with 170 people more susceptible to the disease because of poverty

	State	County	Population_in_Poverty
▶	Montana	Fallon	170

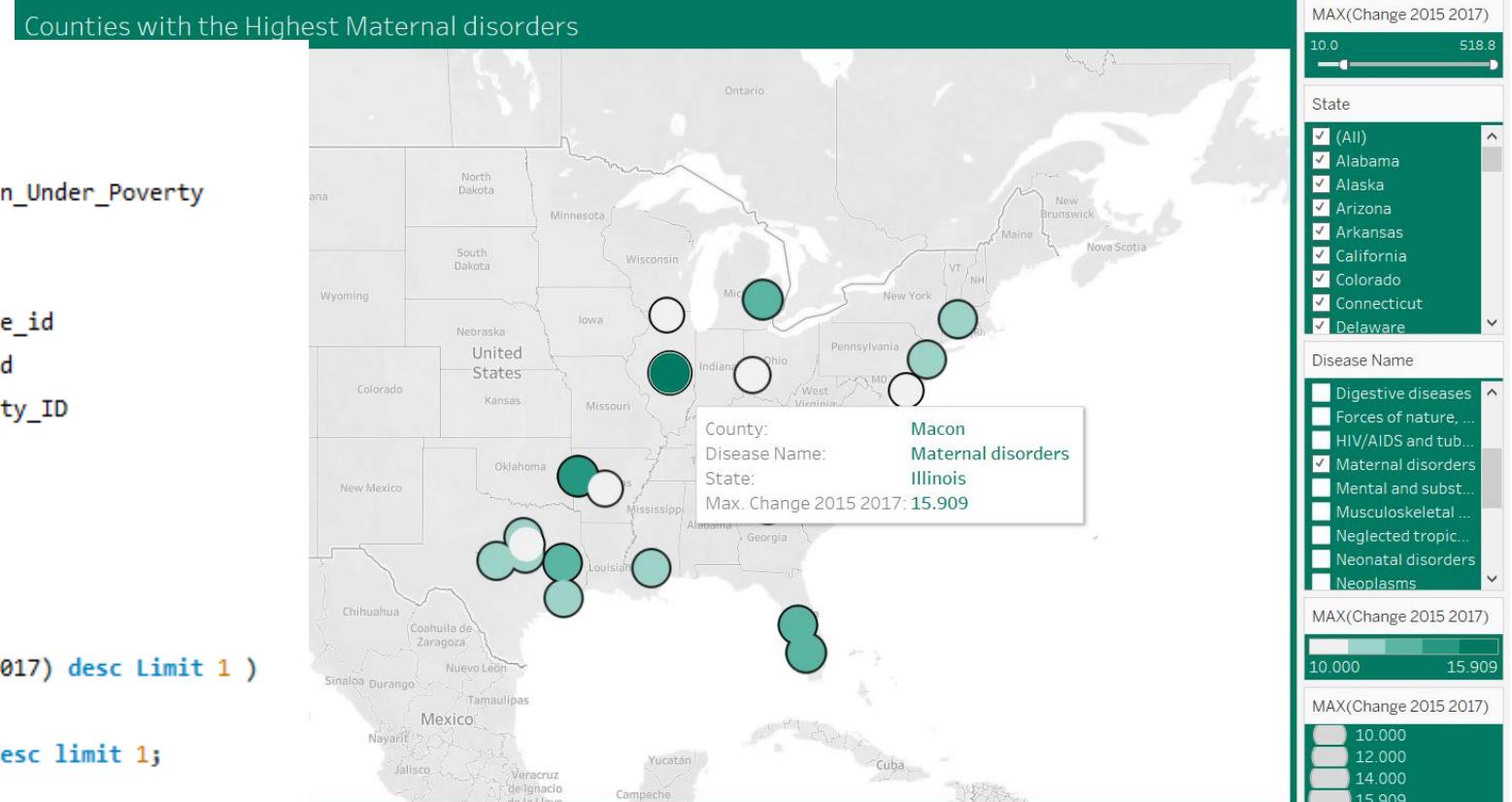


Macon's women in growing risk with Maternal Disorders



Women in Macon, Illinois have seen highest increase in Maternal disorders with around 10.5k women are under poverty

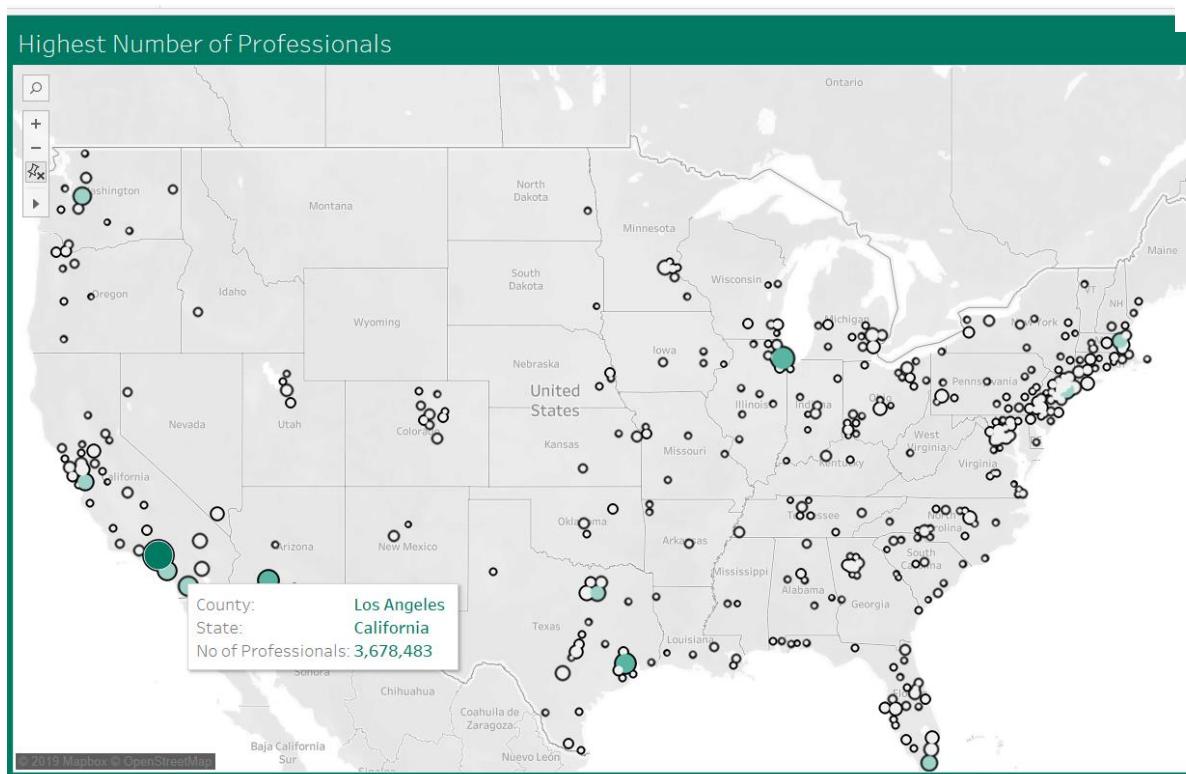
```
216 • Select c.County_id AS County_ID,  
217     s.state AS State,  
218     c.County AS County,  
219     round(cen.Women_Pop*(cen.poverty/100)) AS Women_Under_Poverty  
220  
221     From mortality m  
222         left join disease d ON m.disease_id = d.disease_id  
223         left join county c ON m.county_id = c.county_id  
224         left join census cen ON m.County_ID = cen.County_ID  
225         left join state s ON c.State_ID = s.State_ID  
226     where cen.County_ID IN  
227         (select* from  
228             (select County_ID from mortality  
229                 where disease_id = '1000009'  
230                 group by County_ID order by (Change_2015_2017) desc Limit 1 )  
231             temp_tab )  
232     group by cen.County_ID order by (m.Change_2015_2017) desc limit 1;  
233
```



Los Angeles embracing Immigrants



LA from CA is the most welcoming place in USA for opportunity seekers across the globe



```
200 • Select
201     state.state AS State,
202     coun.county AS County,
203     round(((censu.Total_Pop-censu.Citizen_Pop)*(100-censu.Unemployment))/100 ) AS EmployedImmigrants
204 FROM census censu
205     left join county coun ON censu.County_ID = coun.County_ID
206     Left join state state ON coun.State_ID = state.State_ID
207 where censu.CensusYear = '2017'
208 group by censu.County_ID
209 order by ((censu.Total_Pop-censu.Citizen_Pop)*(100-censu.Unemployment))/100 desc limit 1;
```

Result Grid | Filter Rows:

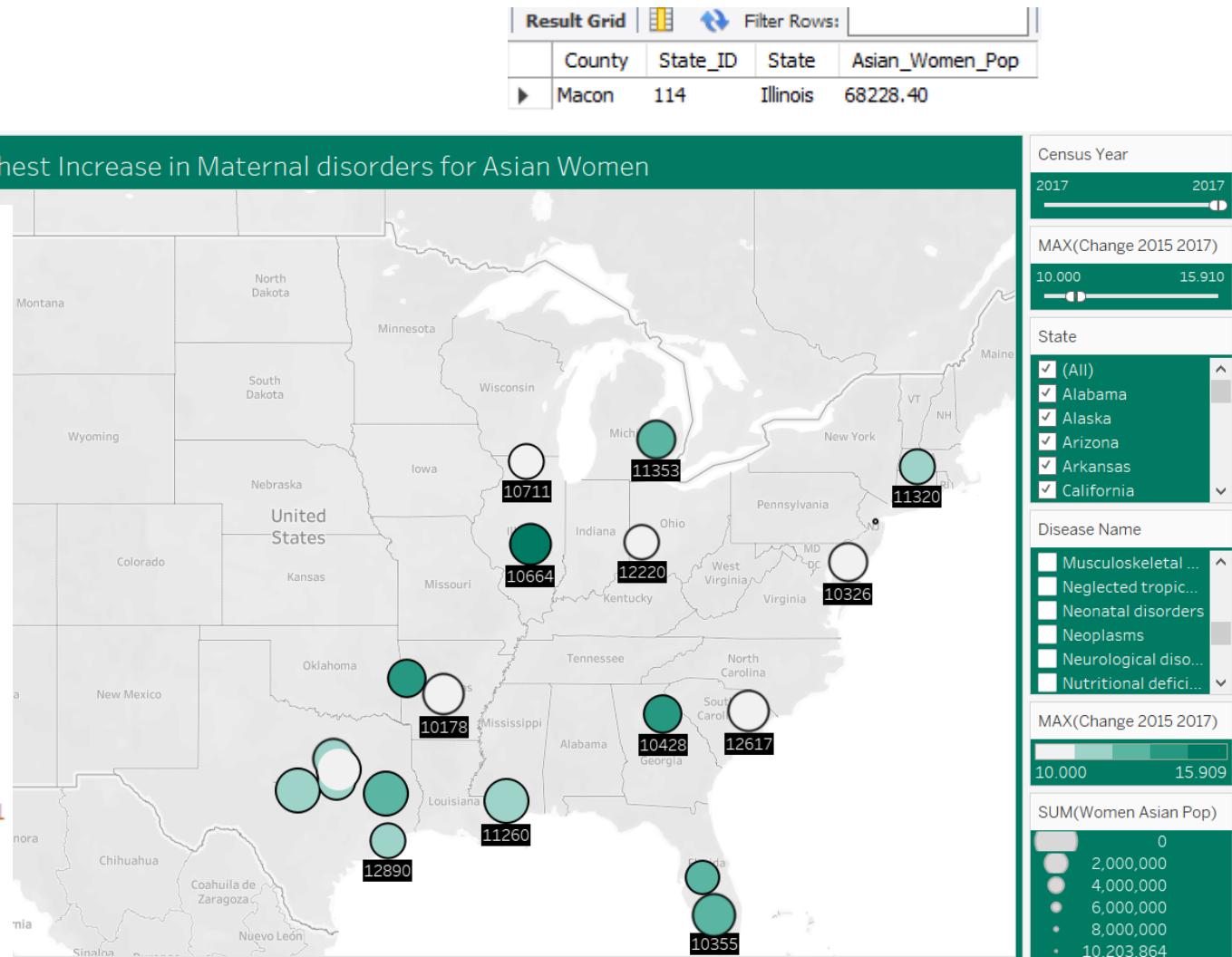
State	County	EmployedImmigrants
California	Los Angeles	3584222

Macon's Asian women suffer the most in maternity



Asian Women in Macon, Illinois are recorded with highest Maternal disorders

```
242
243 • Select county AS County,
244     s.state_id AS State_ID,
245     s.State AS State,
246     cen.Women_Pop*Asian_Pop AS Asian_Women_Pop
247 From mortality m
248 left join disease d ON m.disease_id = d.disease_id
249 left join county c ON m.county_id = c.county_id
250 left join state s ON c.state_id = s.state_id
251 left join census cen ON m.county_id = cen.county_id
252 where cen.county_id IN
253 (select* from (select county_id
254     from mortality
255     where disease_id = '10000009'
256     group by county_id order by MAX(Change_2015_2017)desc limit 1
257 ) temp_tab )
258 group by (cen.county_ID) order by cen.Women_Pop*Asian_pop desc
259 limit 3;
```



Visualization



Viz-1

Concentration of Child Poverty VS USA Average:

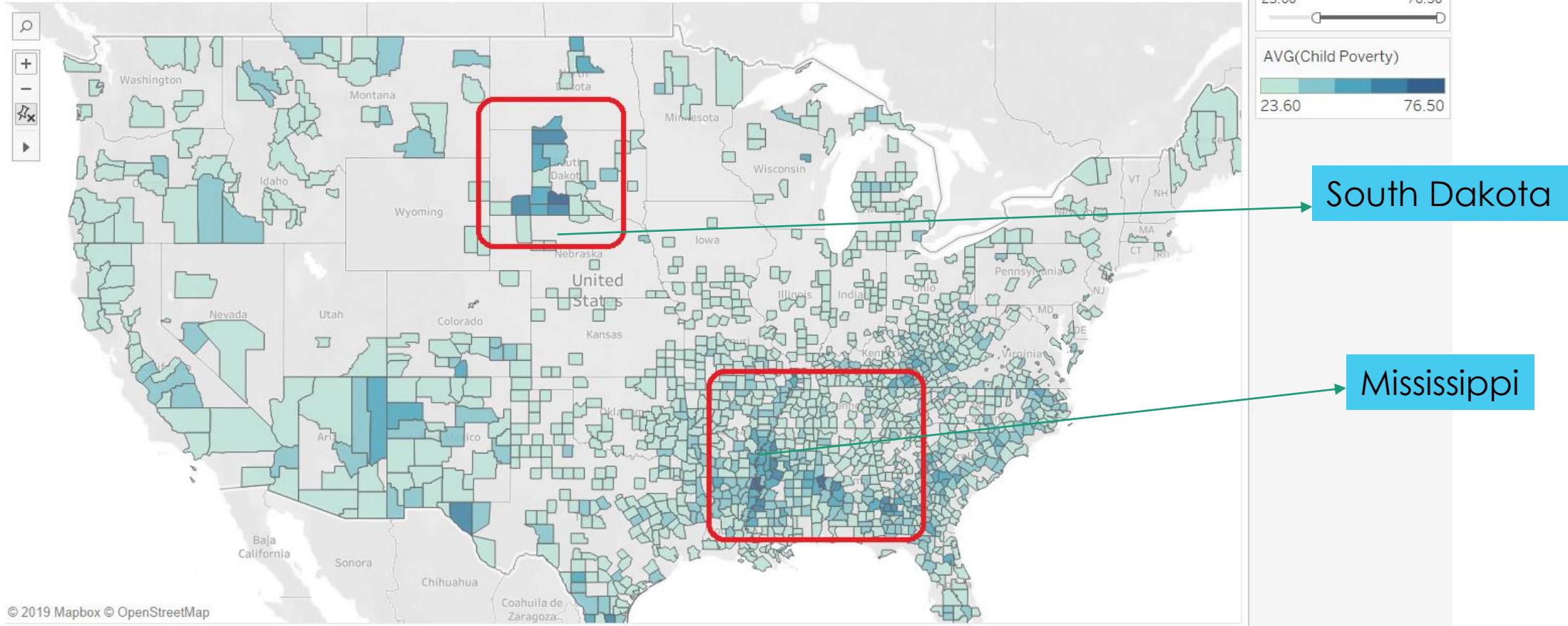
State wise:

Greatest concern : Mississippi

Least concern: Connecticut.

Heavily concentrated in the South Dakota and Mississippi. Both possess the largest number of census tracts where child poverty is above average.

AVG Child Poverty

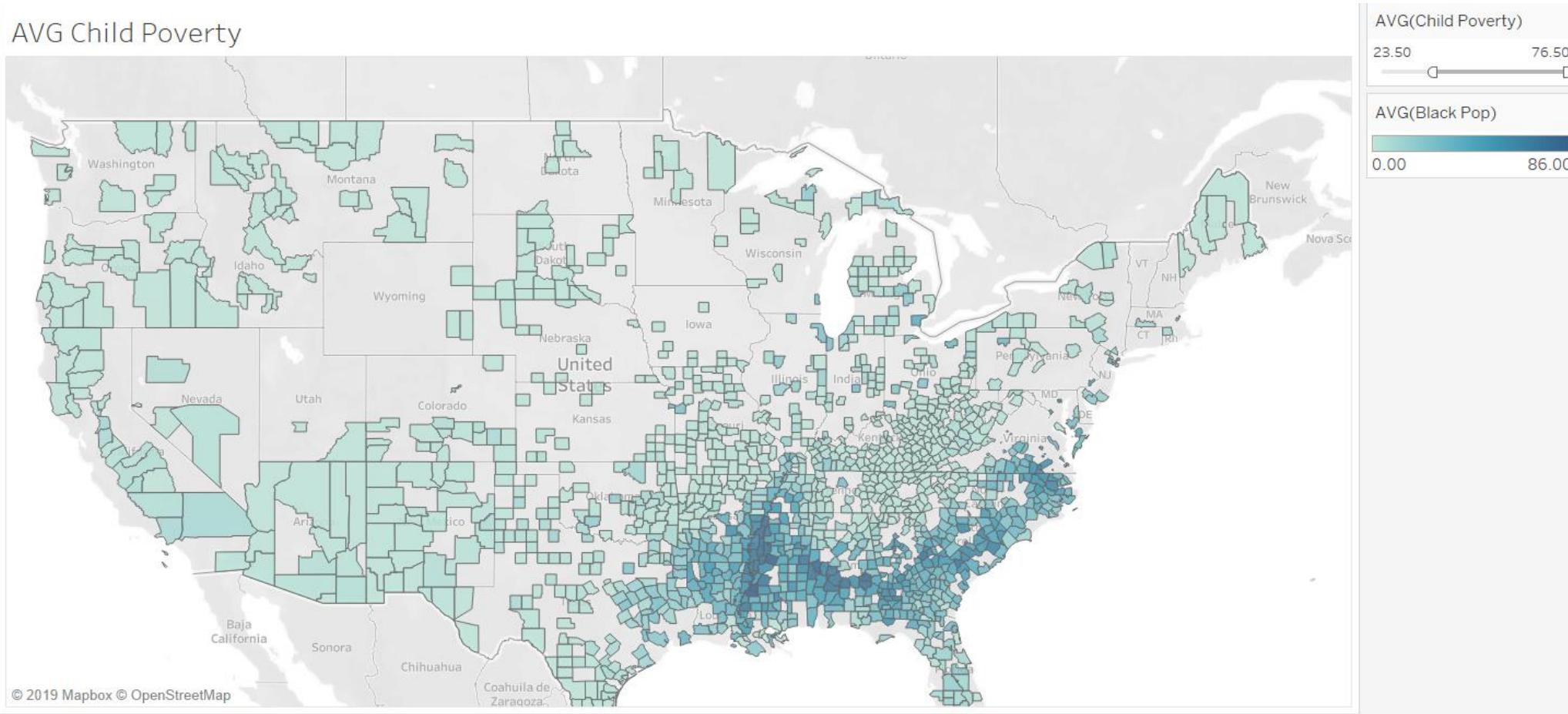




Viz-2

Child Poverty VS Black Population Concentration:

Child Poverty is weakly correlated with the proportion of black residents, showing a positive relationship.

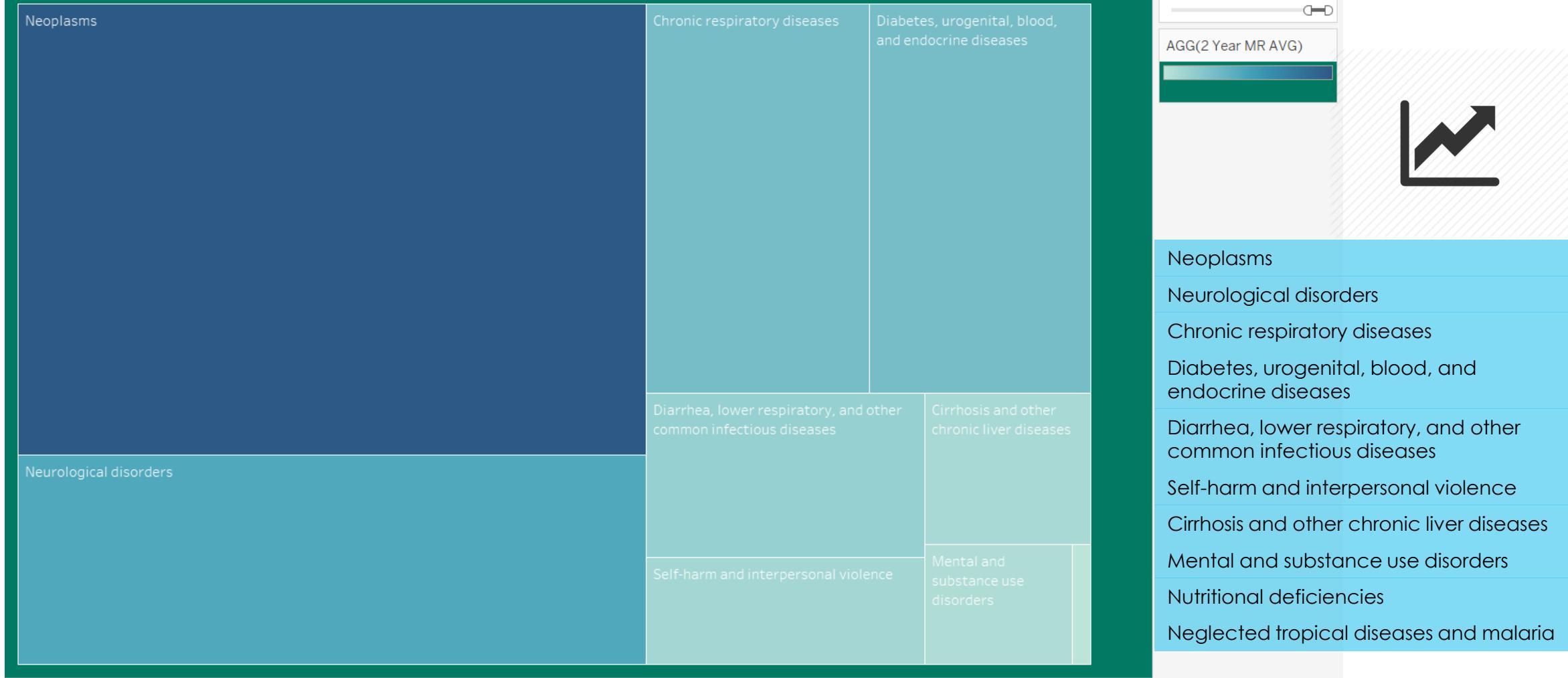




Viz-3

What increased the deaths?

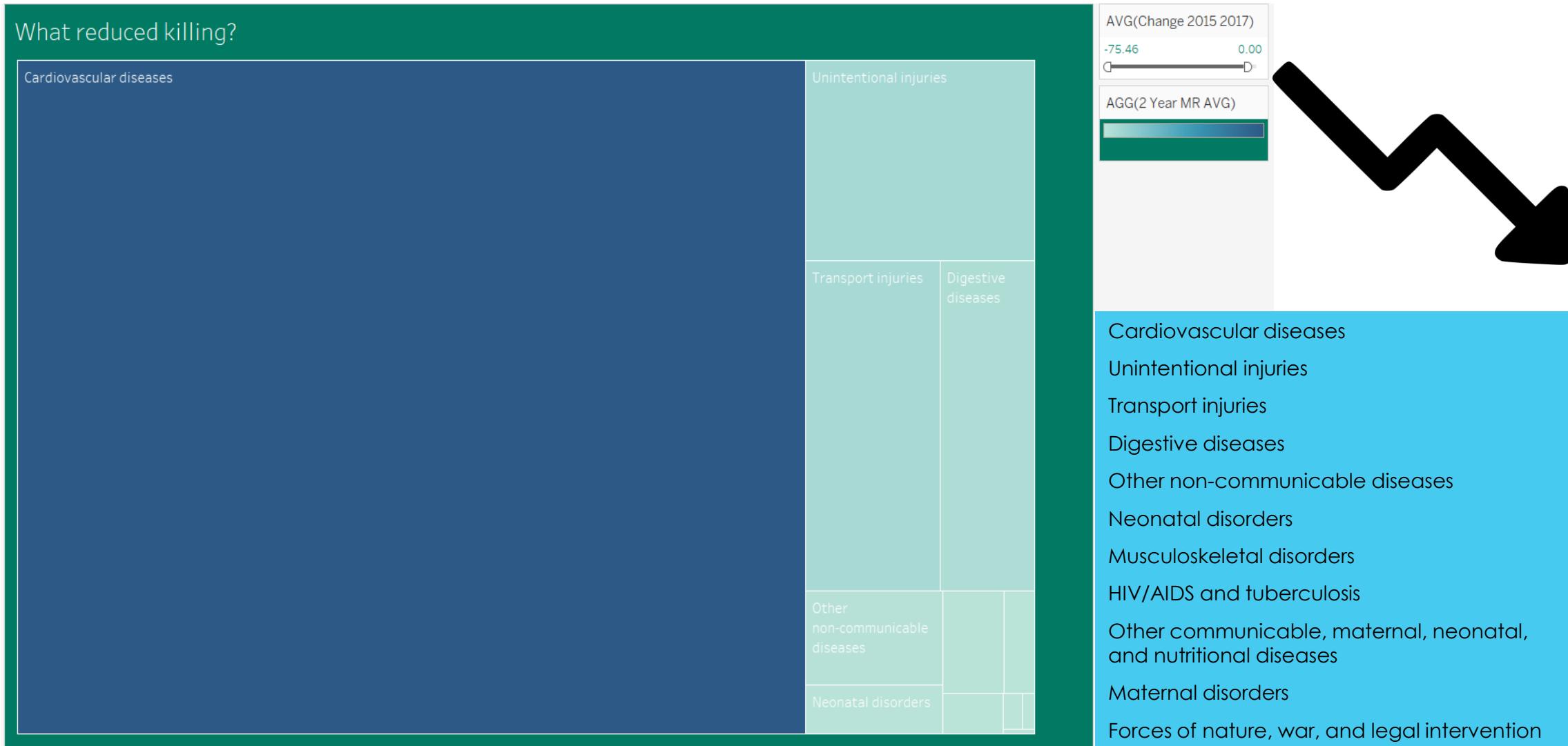
What increased killing?



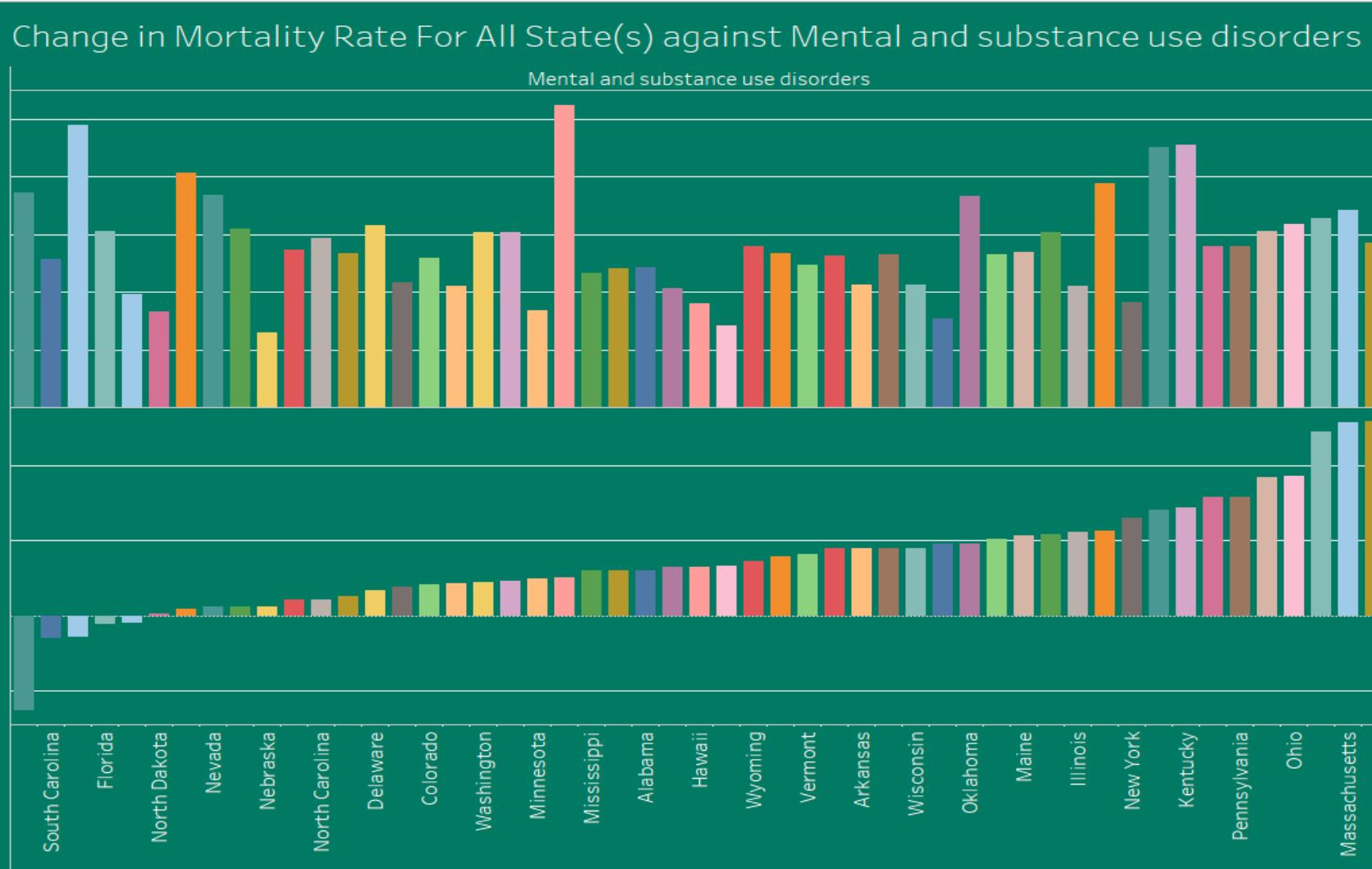


Viz-4

What reduced the deaths?

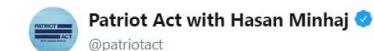


It is interesting to observe that the variation of "Mental and substance use disorders" has almost no negative values. That means that almost no State except DOC, South Carolina, Florida, Alaska, South Dakota has decreased their consumption of substances:



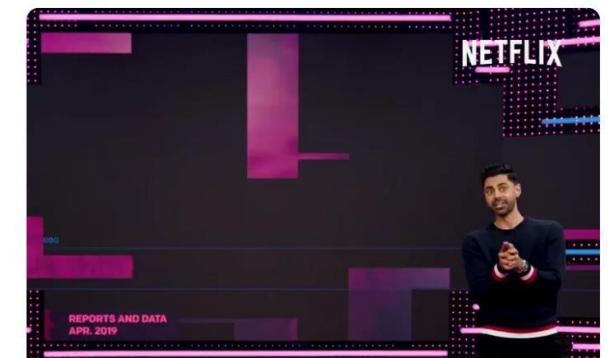
This is a pretty serious problem in USA. Fentanyl-Involved Deaths are Soaring In U.S.
Do watch how serious this is here:

<https://www.youtube.com/watch?v=vkhMfgOT9Xs>

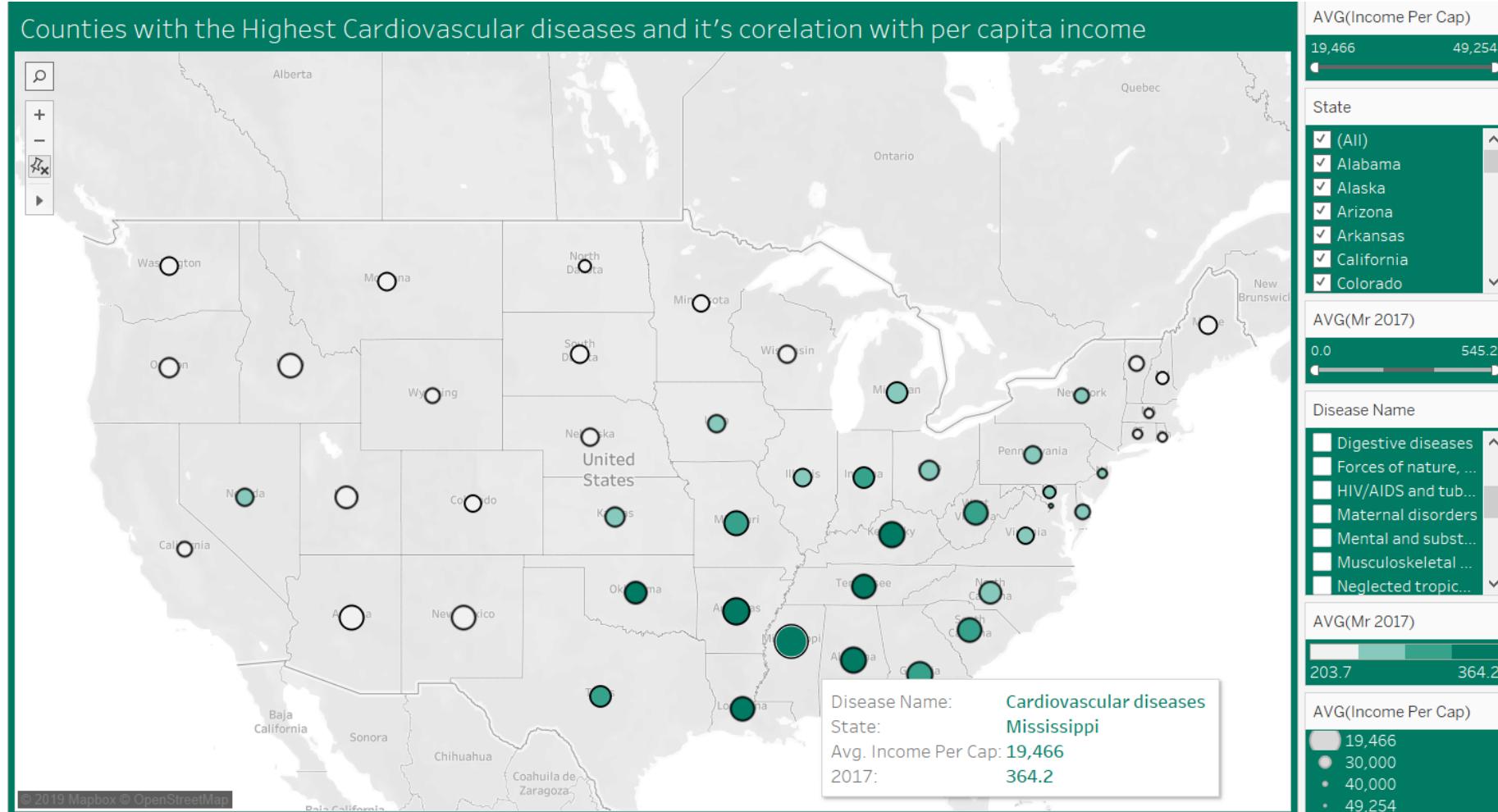


Follow

In the last five years, fentanyl became the reason for nearly two-thirds of all opioid deaths, making it the third wave of the opioid crisis.



It is interesting to observe that the Mississippi has highest deaths due to Cardiovascular disease and lowest per capita income



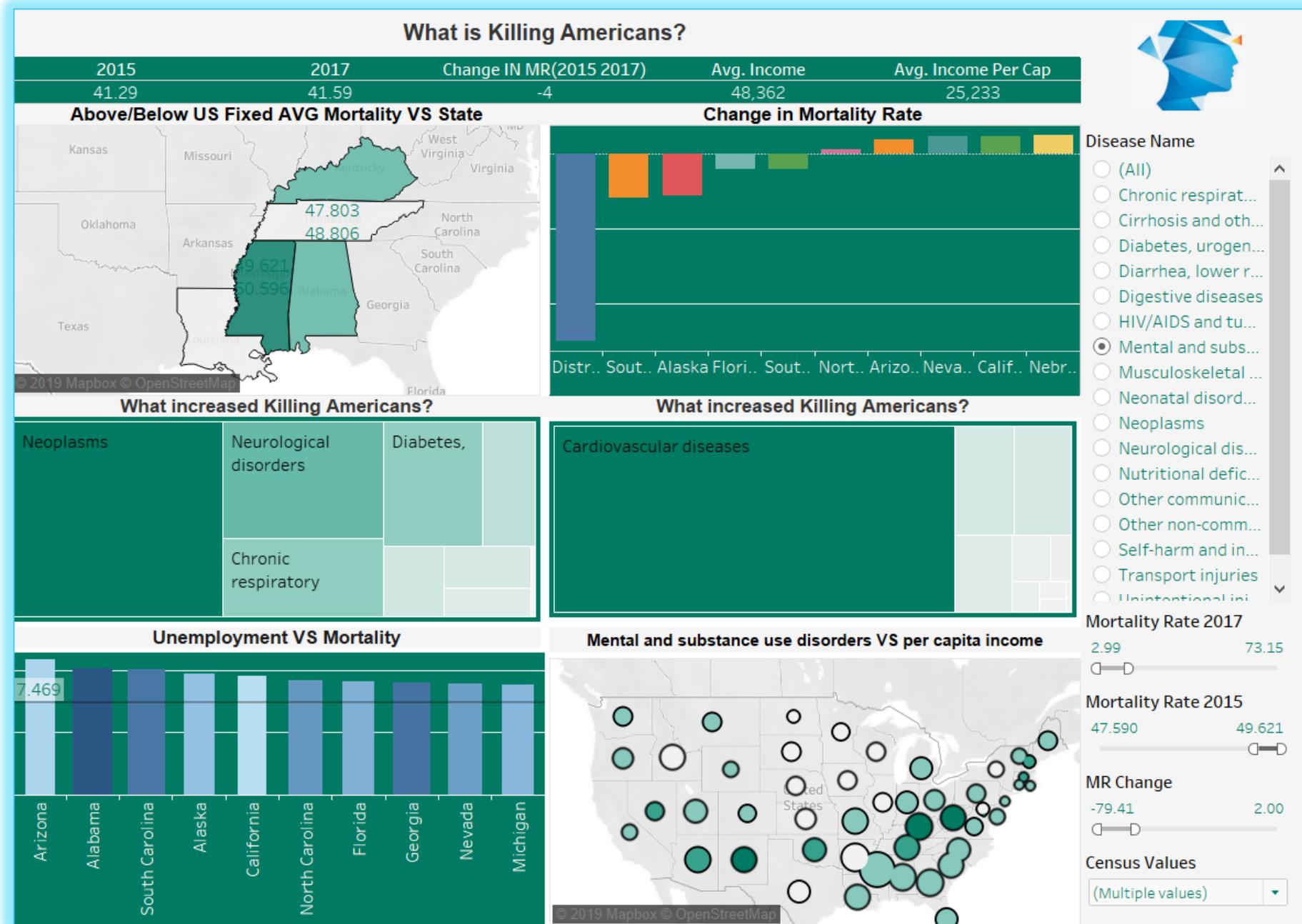
Less income
equals
more junk food



Dashboard



Viz-7



Recommendations

Interesting insights Gained from the SQL Analysis and Reports:

Growing unemployment leads to poverty:

Corson, **Unemployment**: 29.4, **Poverty**: 45.6
Kusilvak, Census Area: **Unemployment**: 28.8, **Poverty**: 39.1
Oglala Lakota, **Unemployment**: 28.7, **Poverty**: 53.3

Poverty decreases with higher income

Counties with 6 figure income has poverty within national average

Poverty encourages people to carpool:

Oglala Lakota with highest poverty has the maximum number of people using Carpool

Communities that take Transit have longer commutes:

Bronx, New York and Queens are topping the chart

Longest and Shortest County names:

Prince of Wales-Hyder Census Area: 33
Lee: 3

Natives and Poverty, Unemployment:

Oglala Lakota has highest Native Population and it is 3rd in unemployment and 1st in Poverty.

Following Counties have 100% White Population:

Jones, Webster, Jackson, Magoffin, Wheeler, Highland, Garfield, Loup

Can two states have the same county name?

The answer is a resounding, patriotic yes and coincidentally there are 51 counties with names similar to the State

Commute:

More people walk in New York City than people who drive, but they do take transit. Clay, Georgia has the highest rate of carpooling.
Alaska likes to walk and, it has the lowest mean commute time as well.

Career and Income:

The highest fraction of people in a 'Professional' career live in Fall Church City, Virginia. This is the same county with third highest income.
Counties with the lowest portion of Construction jobs include Washington D.C., Florida, Rhode Island, New York and California. This is a good sanity check as they are already developed, and any level of development jobs are diluted by their large population.

Population:

The most populated County is Los Angeles, California with a population of 10105722 people The least populated County is: Loving, Texas with a population of: 74 people.
The most populated State is California with a population of 77404311 people The least populated State is: North Dakota with a population of: 1467115 people

A well thought idea of collecting Census and Mortality data is primarily to understand and improve livelihood. The results/observations shown earlier and on this page are vital inputs to the county/state/federal administrators in the policy making/ driving initiatives around healthcare and medical sectors

Because as θ is between $\pi/2$ and 2π , $r = \sin \theta$
and $2\pi - \theta$ results from it retraces its steps.

$$r = |\sin \theta|$$



$0 \leq \theta \leq \pi$

$\pi \leq \theta \leq 2\pi$

$2\pi - \theta$ results from it retraces its steps.

$$Q(T_0 T_2) = -\pi R \cdot \int_{\frac{V_2}{R}}^{\frac{V_1}{R}} - \int_{\frac{V_2}{R}}^{\frac{V_1}{R}} (V_2 - V_1)$$

$$\pi \theta = \pi/3$$

Lessons Learned

$$r = \cos \theta \text{ for } 0 \leq \theta \leq \pi/2$$



$(1, 0) = \text{begin}$
 $(0, 1/2) \text{ end}$

$V_1 - V_2 = -\frac{V_1}{2}$
 $r = |\cos \theta|$

$$r = \cos \theta \text{ for } \pi/2 \leq \theta \leq \pi$$



$5 = A + B \cos \pi$

$5 = A - B$

$$\cos \pi = -1$$

$$5 = A + B$$

$$5 = A - B$$

Lessons Learned

1. Complete Data wrangling process to bring necessary data required (from multiple sources) for the analysis
2. Data cleaning to maintain accurate records for efficient querying and visualizing processes
3. Data Normalization and appropriate categorization of data into respective tables
4. Complex queries and nuisance of SQL with the knowledge of multiple paths of getting results
5. Sequential process in querying, with the quest of finding more business implications from the data
6. Learning from encountered errors like knowing to insert records from correct folder (secure_file_priv), aggregation function (sql_mode=('ONLY_FULL_GROUP_BY',''))
7. Garbage 'IN' equals Garbage 'Out': We initially started off with very little cleaning of the data and we realized that we are not getting correct result/visualization.
Example: Similar County Names leading to incorrect results. We were not aware that USA has various counties with similar names and hence, after realization concatenated the county name with the state name to get unique record
8. Hands on with Tableau taught usage of various kinds of charts that can be associated with a particular form a data

End of Presentation