

Behavioral Risk Factor Surveillance



Abstract

- The "Behavioral Risk Factor Surveillance System (BRFSS)" dataset is a compilation of survey information on people's health and well being in the United States. This dataset is a precious resource for people researching the topic, and public health officials, who are passionate about understanding the behavioral health and factors of risk of the United States population. The Centers for Disease Control and Prevention (CDC) conducted phone conversations with people from different US states to gather the data. A variety of subjects including health practices, chronic illnesses, sleep patterns, and the utilization of preventive services are covered in the survey. The dataset contains data on demographic traits, smoking patterns, levels of physical activity, sleep cycles, and patterns. Insights show that certain demographic groups, such as Black/African American and Hispanic/Latino adults, are more likely to report poor health outcomes and behaviors such as smoking and obesity. We will utilize SQL queries to analyze data, and visualize it with Python and Tableau. This mental health topic is often taboo and it is one of the main reasons why we wanted to dive deeper into the statistics.

Motivation

- This dataset concentrates on the prevalence of health behaviors and chronic conditions by providing insights into health disparities over time. By analyzing this dataset, we hope to identify trends in health behaviors such as tobacco use, physical activity, and alcohol consumption, as well as chronic health conditions such as diabetes, hypertension, and heart disease by highlighting the effects of sleep negligence. This topic of interest caught our attention in mental health as it is often the elephant in the room and people tend to neglect it. We have come across many stories of mental breakdowns. From Influential celebrities to the general public, many have fallen into mental breakdown trauma and committed suicides. Hence one can be able to infer how important it is to value not just the physical, but also a healthy mindset/environment.
- The BRFSS dataset provides a unique opportunity to gain insights into the health behaviors and risk factors of the US population and inform evidence-based public health interventions to promote healthier lives. Over time, people have disregarded how important sleep is and have incurred a number of health problems. We intend to identify compelling causes for the decline in people's health and mentality. These findings can then be used to inform public health policies and interventions aimed at reducing the prevalence of these risk factors and improving overall health outcomes.

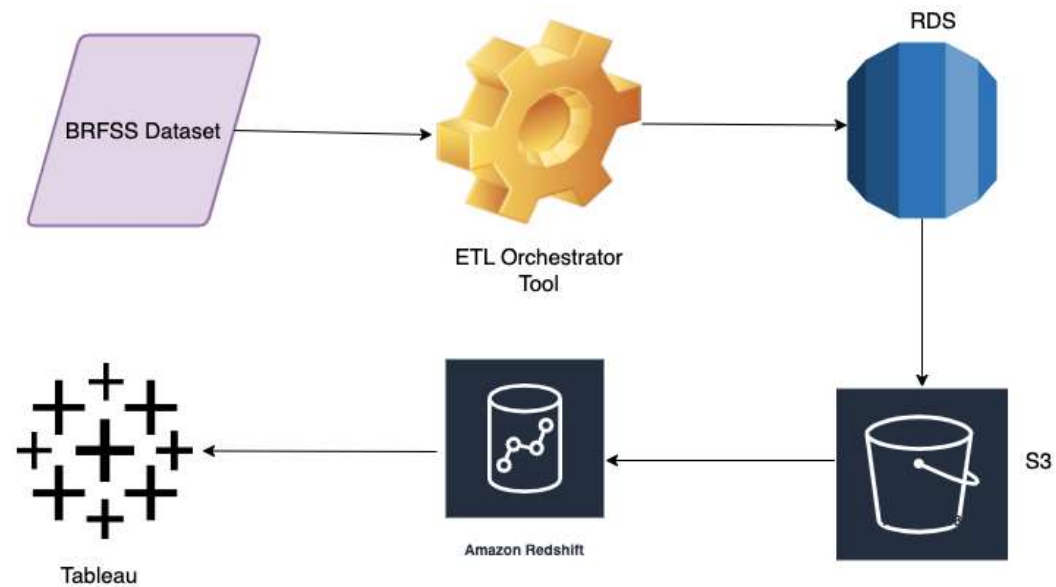
Methodology

- The dataset that we propose to use is called "Behavioral Risk Factor Surveillance System(BRFSS)" and it may be found on the "kaggle.com" website. Cleaning and preprocessing the data, so that it involves resolving inaccurate data and addressing outliers, and replacing duplicate entries with null values using Python would be our first step in the data engineering process. By doing so, we would be able to create appropriate data models and move on to the next stage., which would be storing the data in AWS Cloud. The BRFSS dataset includes information from numerous surveys that were carried out over the years. With data integration, we would compile the results of numerous surveys into a single dataset using Python. Data warehousing and ETL methods will be applied to the dataset. Transforming the data, by aggregating it with a common group and calculating new variables based on existing data would be the next objective of our deliverables. Finally, we would use exploratory analysis to locate anomalies, connected patterns, and summarize them to understand the dataset. Upon doing so, we would be able to draw out newer insights and key factors for the root cause. Finally, to explore and present the insights that we have discovered through data analysis, we would utilize data visualization tools like Tableau and PowerBi and python data libraries like Matplotlib, Seaborn, etc wherever required. This would enable us to draw conclusions from our data that would be clear even to the average person in terms of the underlying reason. Finally in order to provide reports, we would use tools like Grammarly and Latex to improve the phrases and establish a standard format for language use.

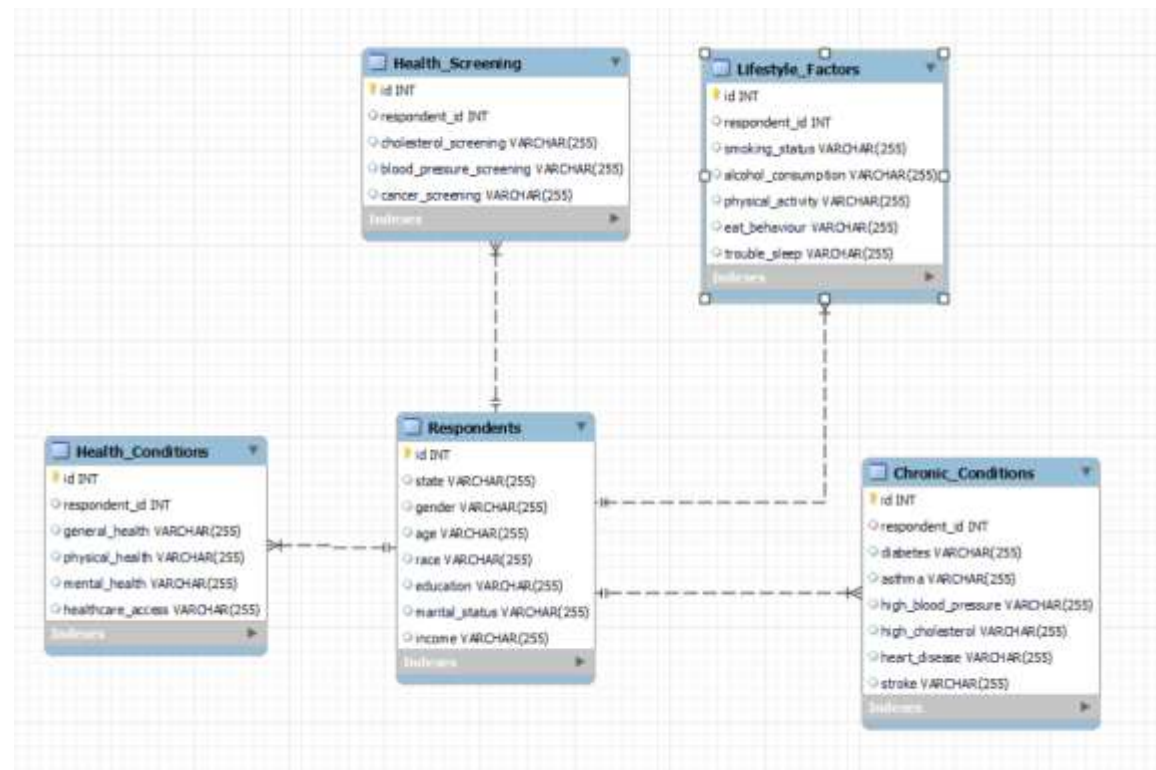
DATA PREPARATION AND EXPLORATION

- The crash data is provided as a comma-separated document (.csv). We took the 5 csv and loaded it into Amazon S3 cloud storage, then put it through AWS Glue to create a data warehouse, where we divided the 5 csv files as a Data Model.
- For preliminary exploration we used python to understand the data structure and for removal of noisy / null data. Based on the column the null values were replaced by the measures of central tendencies such as mean, mode.
- We connected Amazon Redshift to tableau for data visualization and performed operations for meaningful insights.

Architecture Diagram



ER



Insights

```
Run Limit 100 Explain Isolated session Serverless: de... brfss_redshift
```

```
--3. Query to get the percentage of respondents with diabetes:
20 SELECT
21 COUNT(CASE WHEN diabetes = 'yes' THEN 1 END) * 100.0 / COUNT(*) AS percentage
22 FROM brfss.chronic_condition;
23
```

Result 1 (1)

percentage
13.085379238285349

```
Run Limit 100 Explain Isolated session Serverless: de... brfss_redshift
```

```
--3 Query to get the percentage of respondents with high blood pressure:
30 SELECT
31 COUNT(CASE WHEN high_blood_pressure = 'yes' THEN 1 END) * 100.0 / COUNT(*) AS percentage
32 FROM brfss.chronic_condition;
33
```

Result 1 (1)

percentage
40.935315557485869

```
Run Limit 100 Explain Isolated session Serverless: de... brfss_redshift
```

```
--4. Query to get the percentage of respondents with high cholesterol:
35 SELECT
36 COUNT(CASE WHEN high_cholesterol = 'yes' THEN 1 END) * 100.0 / COUNT(*) AS percentage
37 FROM brfss.chronic_condition;
38
```

Result 1 (1)

percentage
37.005298886724458

```
Run Limit 100 Explain Isolated session Serverless: de... brfss_redshift
```

```
--4 query to get the percentage of respondents with asthma:
25 SELECT
26 COUNT(CASE WHEN asthma = 'yes' THEN 1 END) * 100.0 / COUNT(*) AS percentage
27 FROM brfss.chronic_condition;
28
```

Result 1 (1)

percentage
13.827245995738961


```
Run Limit 100 Explain Isolated session Serverless: de... brfss_redshift
```

```
44 --8.Query to get the percentage of respondents who have had a stroke:
45 SELECT
46 COUNT(CASE WHEN stroke = 'Yes' THEN 1 END) * 100.0 / COUNT(*) AS percentage
47 FROM brfss.chronic_condition;
48
```

Result 1 (1)	
percentage	
	4.287448881713277

```
Run Limit 100 Explain Isolated session Serverless: de... brfss_redshift
```

```
38
39 --9.Query to get the percentage of respondents with heart disease:
40 SELECT
41 COUNT(CASE WHEN heart_disease = 'Yes' THEN 1 END) * 100.0 / COUNT(*) AS percentage
42 FROM brfss.chronic_condition;
43
```

Result 1 (1)	
percentage	
	5.781161044862195

```
Run Limit 100 Explain Isolated session Serverless: de... brfss_redshift
```

```
60
61 --10.Query to get the average score for physical health:
62 SELECT AVG(physical_health) FROM brfss.Health_conditions
63 WHERE physical_health > 0;
64
```

Result 1 (1)	
avg	
	11

```
Run Limit 100 Explain Isolated session Serverless: de... brfss_redshift
```

```
64
65 --11.Query to get the average score for mental health:
66 SELECT AVG(mental_health) FROM brfss.Health_conditions
67 WHERE mental_health > 0;
68
```

Result 1 (1)	
avg	
	10

Insights

Run Limit 100 Explain Isolated session Serverless: de... brfss_redshift

```

154 --24.Query to find the percentage of respondents with high blood pressure by age group
155 SELECT
156 CASE
157 WHEN age BETWEEN 18 AND 24 THEN 18-24
158 WHEN age BETWEEN 25 AND 34 THEN 25-34
159 WHEN age BETWEEN 35 AND 44 THEN 35-44
160 WHEN age BETWEEN 45 AND 54 THEN 45-54
161 WHEN age BETWEEN 55 AND 64 THEN 55-64
162 WHEN age BETWEEN 65 AND 74 THEN 65-74
163 WHEN age >= 75 THEN 75+
164 END AS age_group,
165 COUNT(CASE WHEN high_blood_pressure = 1 THEN 1 END) * 100.0 / COUNT(*) AS percentage
166 FROM brfss.chronic_condition
167 JOIN brfss.respondents ON brfss.chronic_condition.respondent_id = brfss.respondents.id
168 GROUP BY age_group
169

```

Result 1 (6)

percentage	age_group +
7.046427385456042	18-24
12.48296567248171	25-34
20.400552857290482	35-44
32.410466238727632	45-54
46.415421838838628	55-64
60.483232108849486	65-74

Run Limit 100 Explain Isolated session Serverless: de... brfss_redshift

```

171 --25.Query to find the percentage of respondents who consume alcohol by age group
172 SELECT
173 CASE
174 WHEN age BETWEEN 18 AND 24 THEN 18-24
175 WHEN age BETWEEN 25 AND 34 THEN 25-34
176 WHEN age BETWEEN 35 AND 44 THEN 35-44
177 WHEN age BETWEEN 45 AND 54 THEN 45-54
178 WHEN age BETWEEN 55 AND 64 THEN 55-64
179 WHEN age BETWEEN 65 AND 74 THEN 65-74
180 WHEN age >= 75 THEN 75+
181 END AS age_group,
182 COUNT(CASE WHEN alcohol_consumption = 1 THEN 1 END) * 100.0 / COUNT(*) AS percentage
183 FROM brfss.lifestyle_factors
184 JOIN brfss.respondents ON brfss.lifestyle_factors.respondent_id = brfss.respondents.id
185 GROUP BY age_group
186

```

Result 1 (6)

age_group +	percentage
18-24	5.872022828946701
25-34	6.183855480522728
35-44	5.38878642676311
45-54	6.704300886252106
55-64	5.186198190616462
65-74	3.965804630897968

Run Limit 100 Explain Isolated session Serverless: de... brfss_redshift

```

187 --22.Query to find the percentage of respondents with asthma by gender
188 SELECT gender, COUNT(CASE WHEN asthma = 1 THEN 1 END) * 100.0 / COUNT(*) AS percentage
189 FROM brfss.chronic_condition
190 JOIN brfss.respondents ON brfss.chronic_condition.respondent_id = brfss.respondents.id
191 GROUP BY gender
192

```

Result 1 (2)

gender	percentage
Female	15.004583477499647
Male	11.404186628546698

Run Limit 100 Explain Isolated session Serverless: de... brfss_redshift

```

193 --23.Query to find the percentage of respondents with diabetes by gender
194 SELECT gender, COUNT(CASE WHEN diabetes = 1 THEN 1 END) * 100.0 / COUNT(*) AS percentage
195 FROM brfss.chronic_condition
196 JOIN brfss.respondents ON brfss.chronic_condition.respondent_id = brfss.respondents.id
197 GROUP BY gender
198

```

Result 1 (2)

gender	percentage
Male	13.585434173669467
Female	12.718538474487616

insights

Insights

```
Run Limit 100 Explain Isolated session Serverless: de... brfss_redshift
```

```
--29.Query to get the percentage of respondents who had a stroke and have high blood pressure:  
220 SELECT COUNT(*) * 100.0 / (SELECT COUNT(*) FROM brfss.chronic_condition WHERE stroke = 'Yes') AS percent  
221 FROM brfss.chronic_condition  
222 WHERE stroke = 'Yes' AND high blood pressure = 'Yes';  
223  
224
```

Result 1 (1)

percentage
73.319644079397672

```
Run Limit 100 Explain Isolated session Serverless: de... brfss_redshift
```

```
--30.Query to get the percentage of respondents who smoke and have heart disease:  
225 SELECT COUNT(*) * 100.0 / (SELECT COUNT(*) FROM brfss.chronic_condition WHERE heart_disease = 'Yes') AS percentage  
226 FROM brfss.chronic_condition  
227 JOIN brfss.lifestyle_factors ON brfss.chronic_condition.respondent_id = brfss.lifestyle_factors.respondent_id  
228 WHERE heart_disease = 'Yes' AND smoking_status = 'Current smoker' OR 'Former smoker';  
229  
230
```

Result 1 (1)

percentage
9.421319796954314

```
Run Limit 100 Explain Isolated session Serverless: de... brfss_redshift
```

```
--31.Query to get the percentage of respondents who have had a cancer screening by age group and gender:  
233 SELECT  
234 CASE  
235 WHEN age BETWEEN 18 AND 24 THEN '18-24'  
236 WHEN age BETWEEN 25 AND 34 THEN '25-34'  
237 WHEN age BETWEEN 35 AND 44 THEN '35-44'  
238 WHEN age BETWEEN 45 AND 54 THEN '45-54'  
239 WHEN age BETWEEN 55 AND 64 THEN '55-64'  
240 WHEN age BETWEEN 65 AND 74 THEN '65-74'  
241 WHEN age >= 75 THEN '75+'  
242 END AS age_group,  
243 gender,  
244 COUNT(CASE WHEN cancer_screening = 'Yes' THEN 1 END) * 100.0 / COUNT(*) AS percentage  
245 FROM brfss.Health_Screening  
246 JOIN brfss.Respondents ON brfss.Health_Screening.respondent_id = brfss.Respondents.id  
247 GROUP BY age_group, gender;  
248
```

Result 1 (12)

age_group ↑	gender	percentage
18-24	Male	0
18-24	Female	0
25-34	Male	0
25-34	Female	0.011827321111768
35-44	Male	0.02348796241926
35-44	Female	0.036903773410831
45-54	Female	3.101534443356186

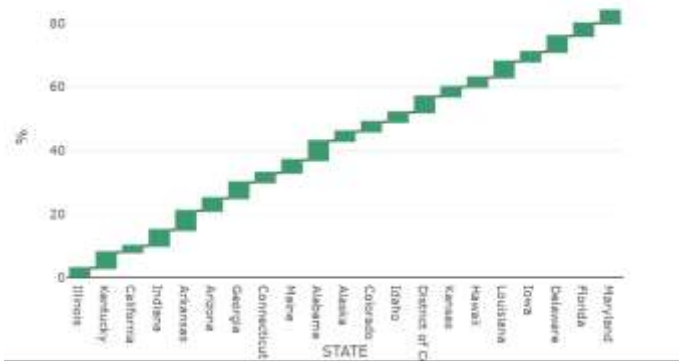


Trends based on Insights

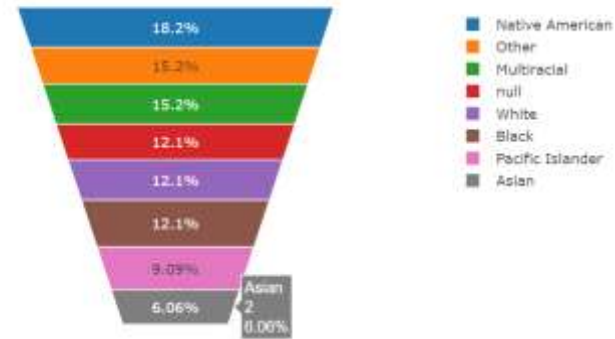
- Chronic conditions such as high blood pressure, high cholesterol, and diabetes are prevalent among respondents, with high blood pressure being the most common.
- Females have a higher prevalence of asthma, while males have a higher prevalence of diabetes.
- Age is a significant factor in the prevalence of chronic conditions, with older respondents having higher rates of high blood pressure, high cholesterol, and stroke.
- Education and marital status both influence respondents' average income.
- There is a wide variation in the distribution of respondents across states, as well as in healthcare access by state.
- Lifestyle factors such as smoking, alcohol consumption, and physical activity vary among respondents, with smoking being linked to heart disease and alcohol consumption varying by age group.

Visualizations

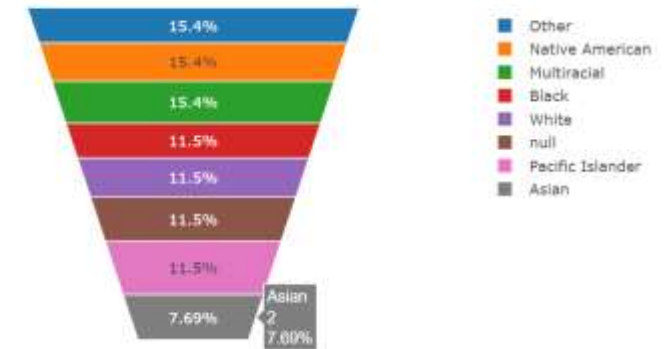
Diabetes percentages across the states



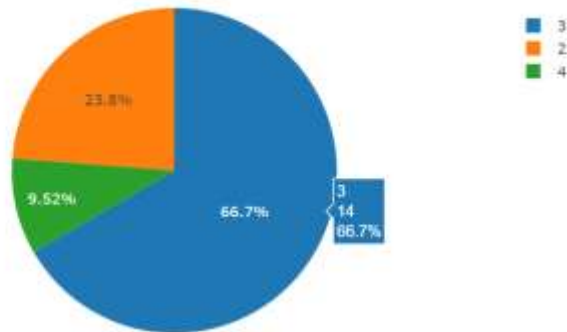
Average Physical Health Score by Race



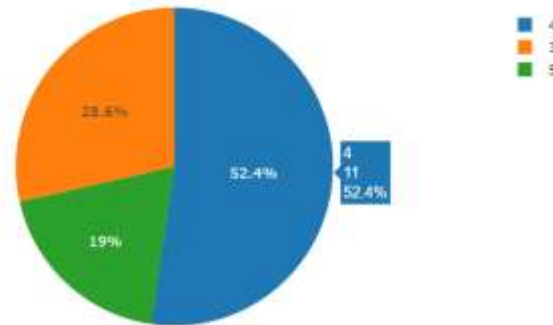
Average Mental Health Score by Race



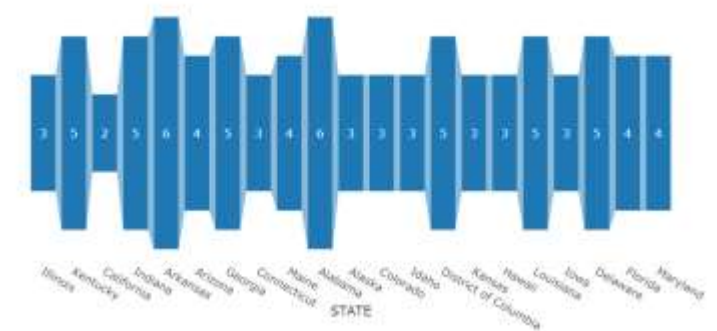
Average Mental Health Percentage



Average Physical Health Percentage



Heart Stroke Rates by State





Impact on Society

- By enhancing our knowledge of behavioral health, factors of risk, and outcomes of the US population, this project has the potential to have a positive impact on society. We can spot trends and patterns in health-related behaviors and results, which can enable public health policies and programs to create focused interventions that can help people and communities live healthier lives by preventing chronic diseases.



Project Development Methodology



CONCLUSION

1. The project analyzed a diverse sample and found key trends and insights that can inform public health policies and interventions.
2. Chronic conditions such as high blood pressure, high cholesterol, and diabetes are highly prevalent among respondents, underscoring the need for effective public health measures.
3. Gender differences exist in the prevalence of certain conditions, and these differences should be taken into account in designing targeted interventions and healthcare services.
4. Age is a significant factor influencing the prevalence of chronic conditions, highlighting the need for age-specific interventions and healthcare services.
5. Education and marital status have an impact on respondents' average income, emphasizing the importance of addressing socioeconomic disparities in health outcomes and access to healthcare services.
6. The distribution of respondents across states and the differences in healthcare access by state indicate that state-specific policies and interventions may be necessary.
7. Lifestyle factors such as smoking, alcohol consumption, and physical activity play a significant role in respondents' health.
8. The findings reveal a complex interplay of factors influencing health, and public health policies and interventions should be designed to address these factors comprehensively and in a targeted manner.

The background features a light blue-to-white gradient on the left and a light green-to-white gradient on the right. These gradients are separated by a vertical line. On the left side, there are several overlapping, wavy, light blue shapes that curve upwards and to the right. On the right side, there are several overlapping, wavy, light green shapes that curve upwards and to the left.

THANK YOU