# Behavioral Risk Factor Surveillance

MOUNIKA VEMPALLI
*Master's in Data Analytics*
*San Jose State University*
San Jose, USA
mounika.vempalli@sjsu.edu

JAYA LAKSHMI GUNJI
*Master's in Data Analytics*
*San Jose State University*
San Jose, USA
jayalakshmi.gunji@sjsu.edu

VENU ANUPATI
*Master's in Data Analytics*
*San Jose State University*
San Jose, USA
venu.anupati@sjsu.edu

SAI SARAN PRASA
*Master's in Data Analytics*
*San Jose State University*
San Jose, USA
saisaran.parasa@sjsu.edu

*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.

## I. 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.

## II. LITERATURE REVIEW

1.The study, "Trends in Health Behaviors and Health Outcomes Among US Adults: The Behavioral Risk Factor Surveillance System 1984-2017," was conducted by Michael R. Kramer and colleagues (American Journal of Preventive Medicine, 2020). Using BRFSS data, this study looked at long-term changes in American adult health habits and outcomes. Data on physical activity, cigarette use, alcohol use, diet, and health outcomes like obesity, diabetes, and cardiovascular disease were all examined by the authors. They discovered that while some health behaviors (such as a decline in smoking rates) have been better over time, others have gotten worse (such as a decrease in physical activity). The study emphasizes the necessity of continuing to monitor health habits and results in order to pinpoint areas that require intervention.

2."State-Level Variations in the Prevalence of Chronic Obstructive Pulmonary Disease (COPD) in the United States: Findings from the Behavioral Risk Factor Surveillance System (BRFSS)" by Kavita Singh et al. (International Journal of Environmental Research and Public Health, 2018).The prevalence of COPD in the US varied at the state level, according to this study's analysis of BRFSS data. Data on self-reported diagnoses of COPD as well as risk variables like smoking and air pollution were examined by the authors. They discovered that there were significant regional variations in the prevalence of COPD and that risk variables including smoking and air pollution had a significant impact on this prevalence. The study emphasizes the requirement for focused initiatives to lower the prevalence of COPD, particularly in states with high-risk factor prevalence.

3."Prevalence of Sexual Orientation Identity and Health Disparities in Adults: The Behavioral Risk Factor Surveillance System, 2013" by Brian R. Hollenbeck et al. (American Journal of Public Health, 2018). The prevalence of sexual orientation identity and related health inequalities among US adults were investigated in this study using BRFSS data. Data on self-reported sexual orientation identity as well as health outcomes like mental health and substance use were examined by the authors. They discovered that compared to heterosexual adults, sexual minorities (including homosexual, lesbian, and bisexual people) had higher rates of mental health and drug use issues. This emphasizes the need for focused treatments to lessen health disparities in this community.

4.Bingrui Xie's cross-sectional study, "Association Between Sleep Duration and Cognitive Impairment: A Cross-Sectional Study Using the BRFSS Database," investigates the connection between sleep duration and cognitive impairment. The Behavioral Risk Factor Surveillance System (BRFSS) database, a sizable survey of the adult population in the United States, serves as the study's data source. For this investigation, a sample of almost 400,000 persons aged 45 and above was used. The study's key finding was that sleep duration and cognitive impairment had a U-shaped connection, meaning that both short and long sleep durations were linked to an increased risk of cognitive impairment. The study indicated that those who slept 7-8 hours per night were less likely to develop cognitive impairment than those who slept less than 5 hours or more than 9 hours per night. Overall, the research indicates that preserving cognitive function in middle-aged and older persons may depend on getting the recommended amount of sleep (7-8 hours each night).

## III. 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.

## IV. 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.
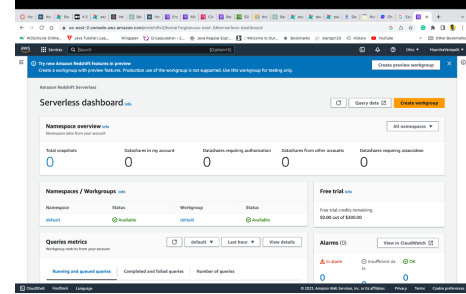
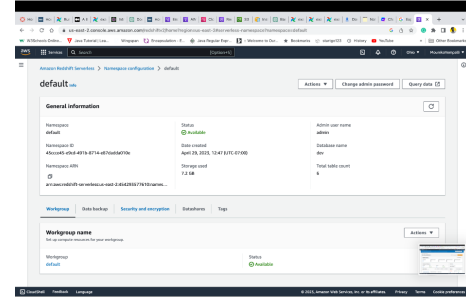

Fig. 1. Amazon Redshift Serverless Cluster



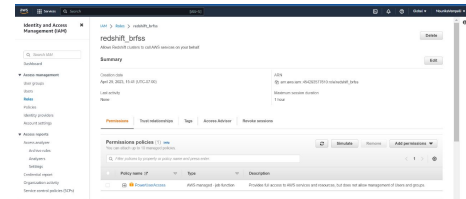Fig. 2. Amazon Redshift Serverless Cluster
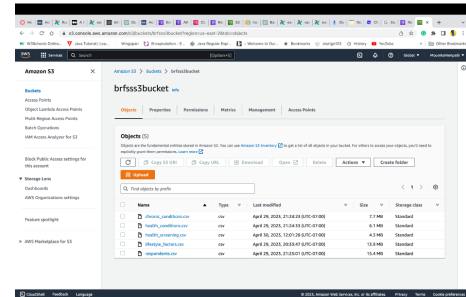


Fig. 3. Creating roles



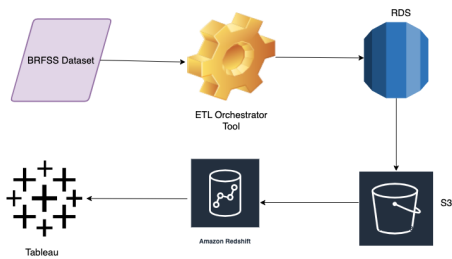Fig. 4. Uploaded Extarcted to S3 Bucket

Fig. 5. ARCHITECTURE DIAGRAM

## V. ER DIAGRAM

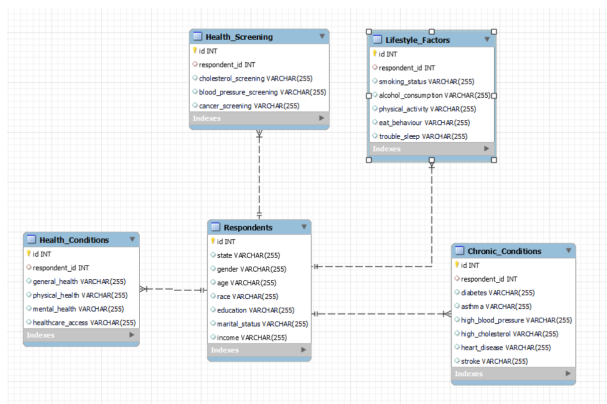We created the following ER diagram using the reverse engineering method in MySQL Workbench.
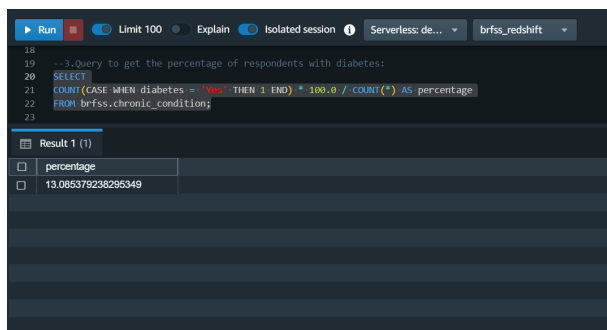


Fig. 6. ER Diagram

## VI. INSIGHTS



Fig. 7. Query to get the percentage of respondents with diabetes

We have done as many queries as possible to get more insights and to understand the data better .From all the querying and analysis that we have done, we found out the following observations.

• California, Colorado, and Kansas have the highest number of respondents, while Alaska, Arkansas, and the District of Columbia have the lowest number of respondents among the states.



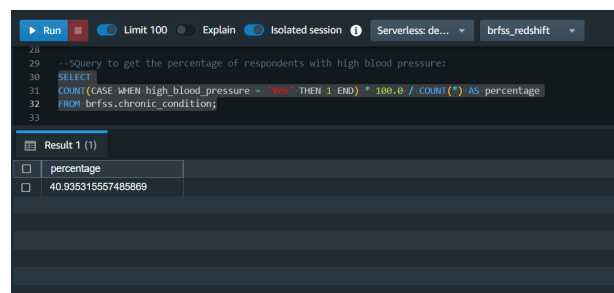Fig. 8. Query to get the percentage of respondents with asthma



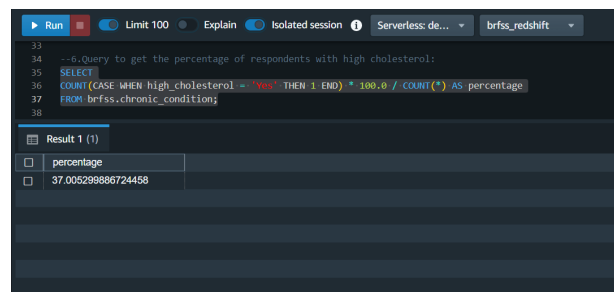Fig. 9. Query to get the percentage of respondents with high blood pressure



Fig. 10. Query to get the percentage of respondents with high cholesterol



Fig. 11. Query to get the percentage of respondents with heart disease

Fig. 12. Query to get the percentage of respondents who have had a stroke



Fig. 13. Query to get the average score for physical health



Fig. 14. Query to get the average score for mental health



Fig. 15. Query to find the percentage of respondents with diabetes by gender



Fig. 16. Query to find the percentage of respondents with asthma by gender



Fig. 17. Query to find the percentage of respondents with high blood pressure by age group



Fig. 18. Query to find the percentage of respondents who consume alcohol by age group

• The majority of respondents fall into the 65-74 age group, followed by the 45-54 age group. The 18-24 age group has the lowest representation.

• Diabetes prevalence is slightly higher in males (13.59percentage) than in females (12.72percentage).

• Asthma prevalence is significantly higher in females (15.60percentage) than in males (11.40percentage).

• High blood pressure and high cholesterol are prevalent among respondents, affecting 40.94percentage and 37.00 percentage of them, respectively.

• Heart disease and stroke have a lower prevalence, affecting 5.78percentage and 4.29percentage of respondents, respectively.

• General health scores are mostly rated as "Good" (3 out of 5), with average physical health and mental health scores at 11 and 10 (out of 30), respectively.

• A large percentage of respondents (86.75) have had cholesterol screenings, while 10.11percentage have had blood pressure screenings, and only 5.13percentage have had cancer screenings.

• Current smokers make up 9.63percent of respondents, and 4.99percent consume alcohol.

• Regular physical activity is reported by 46.25 percent of respondents, and 84.07 percent have good eating behavior.

• Higher education levels correlate with higher average income, and married respondents have a higher average income compared to other marital statuses.

• The percentage of respondents with healthcare access varies by state, with Alabama having the highest percentage (14.03) and the District of Columbia having the lowest (5.86).

• The prevalence of high blood pressure increases with age, and alcohol consumption is slightly higher among the 25-34 and 45-54 age groups.

• Respondents with multiple chronic conditions are relatively rare.

• A strong correlation exists between having a stroke and having high blood pressure, with 73.32 percent of respondents who had a stroke also having high blood pressure.

• A smaller percentage of respondents (9.42) who smoke also have heart disease.

• Cancer screening rates vary by age group and gender, with the highest percentages found in the 65-74 and 55-64 age groups for both males and females.

## VII. 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.

## VIII. VISUALIZATIONS

Data visualization is the act of transforming raw data into meaningful graphical representations such as charts, graphs, pictures, and films. It is advantageous to gain insights from it in this manner since it simplifies the explanation of the digits and numbers. The relevant images are generated using the previously stated energy consumption algorithm.

## IX. TECHNICAL DIFFICULTY

In order to do analysis and get insights on risk factors producing various health conditions, we have used a variety of approaches throughout this project. This project has be implemented using AWS cloud infrastructure and services like AWS Redshift. These technologies has improved the project's



Fig. 19. Query to get the percentage of respondents who smoke and have heart disease



Fig. 20. Query to get the percentage of respondents who smoke and have heart disease
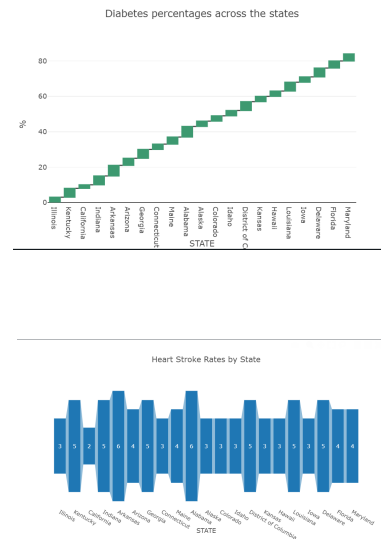


Fig. 21. Query to get the percentage of respondents who smoke and have heart disease
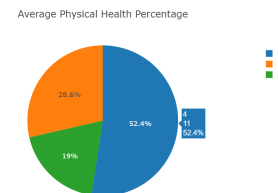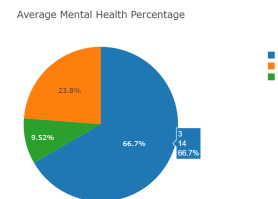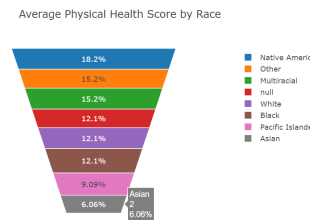


Fig. 22. Query to get the percentage of respondents who smoke and have heart disease
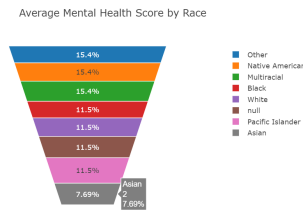
Fig. 23. Query to get the percentage of respondents who smoke and have heart disease

productivity while also given us a chance to practice utilizing these services and gain knowledge of ideas like data modeling, exploratory data analysis, and warehousing.

## X. IMPACT

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.

## XI. RELEVANCE TO THE COURSE

This project provided us with an opportunity to work with more complex queries and analyses. We have covered the concepts such as normalization and indexing which are important considerations when designing and managing a database. It has also helped us to cover the concepts such as designing a relational schema, defining the relationships between them, querying the database to retrieve specific subsets of data, joining tables together to combine related data or aggregating data to calculate statistics, maintaining data integrity by performing ETL operations and visualizing the data.

## XII. UNIQUENESS

The dataset that we took for this project is vast with over 120,000 entries from individuals across different states in the US, providing a large and diverse sample from both national and regional levels. Data related to the health status of the people along with their exact geographical location of each data entry, their low confidence and high confidence limits enables us to examine and identify potential discrepancies in terms of quality of life between gender, age, geographic areas, or other groups making it a comprehensive resource for understanding health outcomes.

## XIII. MoM(PAIR PROGRAMMING)

Agenda:
• Review progress on Data model and ETL tool
• Discuss potential issues or concerns
• Make decisions and assign action items
Discussion:
• Venu presented their proposed solution to the data modeling and suggested using a dataset which is available on Kaggle.

• Mounika challenged the proposed solution, expressing concerns about its feasibility and meaningful entities, and further suggesting to use the dataset from the original BRFSS archives.
• A lively discussion ensued, with both parties presenting their perspectives and providing evidence to support their arguments.
• After careful consideration of all viewpoints, the group decided to complete the data model first and ensure that the dataset has enough scope to perform data modeling to draw meaningful entities.

## XIV. SCRUM MEETING MOM

Minutes of Meeting
Project Name: BRFSS Dataset Project
Date: April 17 2023
Attendees: Venu, Saran, Mounika, Jaya
Agenda:
• Review progress on data warehousing, ETL orchestration, data cleaning and data preprocessing, and visualizations
• Identify any potential roadblocks and challenges
• Assign action items
Discussion:
• Saran and Jaya shared updates on data cleaning and data preprocessing, highlighting that the team had made significant progress in cleaning and preprocessing the data, and that they would soon start visualizing the data.
• Mounika discussed her progress in developing the ETL orchestration and mentioned that the team would soon have a functional ETL process.
• Venu shared an update on data warehousing using Redshift and mentioned that progress was on track and that the team could soon start loading data into the warehouse.
• The team discussed potential roadblocks and identified that the integration of the various components of the project could pose a challenge.
• Venu proposed scheduling a meeting to discuss integration and identify any potential issues that may arise.
• The team agreed to schedule a meeting to discuss integration and identified the need for frequent communication among the team members.
Action Items:
• Saran and Jaya will continue working on data cleaning and preprocessing and start visualizing the data
• Mounika will continue developing the ETL orchestration and work on integrating it with the data warehousing.
• Venu will continue working on data warehousing using Redshift.
• The team will schedule a meeting to discuss integration and identify any potential issues that may arise.

## XV. VERSION CONTROL

We used GitHub to publish and make our project public. https://github.com/monikavempalli/brfss-db

## XVI. CONCLUSION

• This project analyzed a diverse sample of respondents across various states, age groups, and backgrounds. The data reveals several key trends and insights that can inform public health policies and interventions.

• Firstly, chronic conditions such as high blood pressure, high cholesterol, and diabetes are highly prevalent among the respondents. Specifically, high blood pressure is the most common condition, affecting over 40percent of respondents. This underscores the need for effective public health measures to prevent and manage these conditions, including promoting lifestyle changes and ensuring access to appropriate healthcare services.

• Gender differences are also evident in the prevalence of certain conditions. For instance, females have a higher prevalence of asthma, while males have a higher prevalence of diabetes. These differences should be taken into account when designing targeted interventions and healthcare services.

• Age is another significant factor influencing the prevalence of chronic conditions. Older respondents have higher rates of high blood pressure, high cholesterol, and stroke. This highlights the importance of age-specific interventions and healthcare services, as well as the need to address the unique health challenges faced by older populations.

• Education and marital status both have an impact on respondents' average income. College graduates earn the highest average income, while those with less than a high school education earn the lowest. Similarly, married respondents earn the highest average income, while separated respondents earn the lowest. These findings underscore the importance of addressing socioeconomic disparities in health outcomes and access to healthcare services.

• The wide variation in the distribution of respondents across states and the differences in healthcare access by state indicate that state-specific policies and interventions may be necessary to address the unique health challenges faced by each state's population.

• Lifestyle factors such as smoking, alcohol consumption, and physical activity also play a significant role in respondents' health. Smoking is linked to heart disease, and alcohol consumption varies by age group, with the highest consumption observed in the 25-34 age group. Regular physical activity is reported by just over 46percent of respondents, suggesting that efforts to promote physical activity could have a significant impact on overall health outcomes.

• In conclusion, the findings from this project reveal a complex interplay of factors influencing the health of respondents, including chronic conditions, age, gender, socioeconomic status, and lifestyle choices. Public health policies and interventions should be designed to address these diverse factors in a comprehensive and targeted manner, with the goal of improving overall health outcomes and reducing disparities in access to healthcare services.

## REFERENCES

[1] "Behavioral Risk Factor Surveillance System" 1993-2010 Health-Related Quality of Life Data Dataset: https://www.kaggle.com/datasets/thedevastator/behavioral-risk-factor-surveillance

[2] "Behavioral Risk Factor Surveillance System-Centers of Disease Control and Prevention" Basic Information about Dataset: https://www.cdc.gov/brfss/index.html

[3] B. Xie, "Association Between Sleep Duration and Cognitive Impairment: A Cross-Sectional Study Using the BRFSS Database," 2020 International Conference on Public Health and Data Science (ICPHDS), Guangzhou, China, 2020, pp. 108-115, doi: 10.1109/ICPHDS51617.2020.00030.

[4] "Trends in Health Behaviors and Health Outcomes among US Adults: The Behavioral Risk Factor Surveillance System 1984-2017" by Michael R. Kramer et al. (American Journal of Preventive Medicine, 2020)

[5] "State-Level Variations in the Prevalence of Chronic Obstructive Pulmonary Disease (COPD) in the United States: Findings from the Behavioral Risk Factor Surveillance System (BRFSS)" by Kavita Singh et al. (International Journal of Environmental Research and Public Health, 2018)

[6] "Prevalence of Sexual Orientation Identity and Health Disparities in Adults: The Behavioral Risk Factor Surveillance System, 2013" by Brian R. Hollenbeck et al. (American Journal of BRFSS Public Health, 2018)

[7] Remington PL, Smith MY, Williamson DF, Anda RF, Gentry EM, Hogelin GC. Design, characteristics, and usefulness of state-based behavioral risk factor surveillance: 1981-87. Public Health Rep. 1988 Jul-Aug;103(4):366-75. PMID: 2841712; PMCID: PMC1478092.

[8] Macera CA, Ham SA, Yore MM, Jones DA, Ainsworth BE, Kimsey CD, Kohl HW 3rd. Prevalence of physical activity in the United States: Behavioral Risk Factor Surveillance System, 2001. Prev Chronic Dis. 2005 Apr;2(2):A17. Epub 2005 Mar 15. PMID: 15888228; PMCID: PMC1327711.

[9] Centers for Disease Control and Prevention (CDC). Methodologic changes in the Behavioral Risk Factor Surveillance System in 2011 and potential effects on prevalence estimates. MMWR Morb Mortal Wkly Rep. 2012 Jun 8;61(22):410-3. PMID: 22672976.

[10] Siegel, Paul Z., et al. "Behavioral Risk Factor Surveillance, 1986–1990." Morbidity and Mortality Weekly Report: Surveillance Summaries, vol. 40, no. SS-4, 1991, pp. 1–23. JSTOR, http://www.jstor.org/stable/24675437. Accessed 9 Mar. 2023.

**Appendix**

| Criteria | Pts | Comments |
|---|---|---|
| Presentation Skills Includes time management | 5 pts | |
| Code Walkthrough | 3 pts | The Code walkthrough during the presentation would cover the architecture and orchestration of the project and other related tools. All the code and scripts used in the project are stored on GitHub, and the project report includes screenshots of each outcome. |
| Discussion / Q&A | 4 pts | |
| Demo | 5 pts | |
| Version Control Use of Git / GitHub or equivalent; must be publicly accessible | 3 pts | https://github.com/monikavempalli/brfss-db<br><br>Used git for our project. |
| Significance to the real world | 5 pts | The dataset from the Behavioral Risk Factor Surveillance System (BRFSS) is significant in the real world as it provides crucial information on the health behaviors and risk factors of adults in the United States. This information can be used to develop and implement public health programs and policies aimed at improving the overall health of the population. |
| Lessons learned Included in the report and presentation? How substantial and unique are they? | 5 pts | Yes, we have covered all the relevant sections. |
| Innovation | 5 pts | Using the BRFSS dataset, we were able to draw insights of health patterns across several states in the United States. The innovation in this project is the key challenge to data model a dataset of huge attributes and clean it. |
| Teamwork | 5 pts | As a team, everyone has actively taken part in each and every phase of the project and we practiced agile methodology during the course of the project. |
| Technical difficulty | 4 pts | We have learnt so many technologies like AWS Redshift, RDS, S3 and techniques like warehousing. |

| | | |
|---|---|---|
| This criterion is linked to a learning outcome Practiced pair programming? See: https://en.wikipedia.org/wiki/Pair_programming Li nksLinksLinksLinksLinksLinksLinksLinksLinksL in ksLinksLinks Links to an external site. to an external site. to an external site. to an external site. to an external site. to an external site. to an external site. to an external site. to an external site. to an external site. to an external site. to an external site. to an external site | 2 pts | Yes, we have practiced pair programming and attached the minutes of meeting as a proof. |
| Practiced agile / scrum (1-week sprints)? Submit evidence on Canvas - meeting minutes, other artifacts | 3 pts | Attached MOM of one of our scrum meetings as a proof. |
| Used Grammarly / other tools for language? Grammarly free version is sufficient; can use other tools as well. Submit report screenshot on Canvas. | 2 pts | Grammarly and Google docs were used |
| Slides | 5 pts | Have attached it in Canvas. |
| Report Format, completeness, language, plagiarism, whether turnItIn could process it (no unnecessary screenshots), etc | 5 pts | We utilized the IEEE format for the document created in overleaf and selectively incorporated features in Microsoft PowerPoint to design the slides for our presentations. |
| Used unique tools E.g.: LaTeX for writing reports (submit .tex that is not generated from another format such as .docx; generating from .lyx and similar LaTeX editor outputs is fine. Also checkout https://www.overleaf.com/LinksLinksLinks Links to an external site. to an external site.) Unique features of Prezi or powerpoint, etc | 5 pts | We have used LaTex using overleaf website to generate the project report |
| Performed substantial analysis using database techniques Project must include an analytics component | 3 pts | We used Tableau and Redshift to visualize our findings. |
| Used a new database or data warehouse tool not covered in the HW or class | 3 pts | AWS Redshift and AWS RDS |
| Used appropriate data modeling techniques | 5 pts | Out of 330 columns in the dataset, selected (27) columns and normalized the dataset into 3NF separating entities. For visual reference, we used draw.io to capture relationships. |
| Used ETL tool | 1 pts | Developed a custom ETL orchestration in python to perform data loading. |
| Demonstrated how Analytics support business | 3 pts | We were able to draw conclusions of |

| decisions | | health patterns which may help any business that is involved across health sector |
|---|---|---|
| Used RDBMS Idea is to exercise as many topics from the course as possible | 1 pts | We have used RDBMS to run various queries on MySQL Workbench |
| Used Datawarehouse Idea is to exercise as many topics from the course as possible | 1 pts | Used AWS redshift for querying |
| Includes DB Connectivity / API calls Possibly using Python | 1 pts | Used MySQL Connector to connect to an RDS instance. |
| Used NOSQL | 1 pts | |