

Abstract

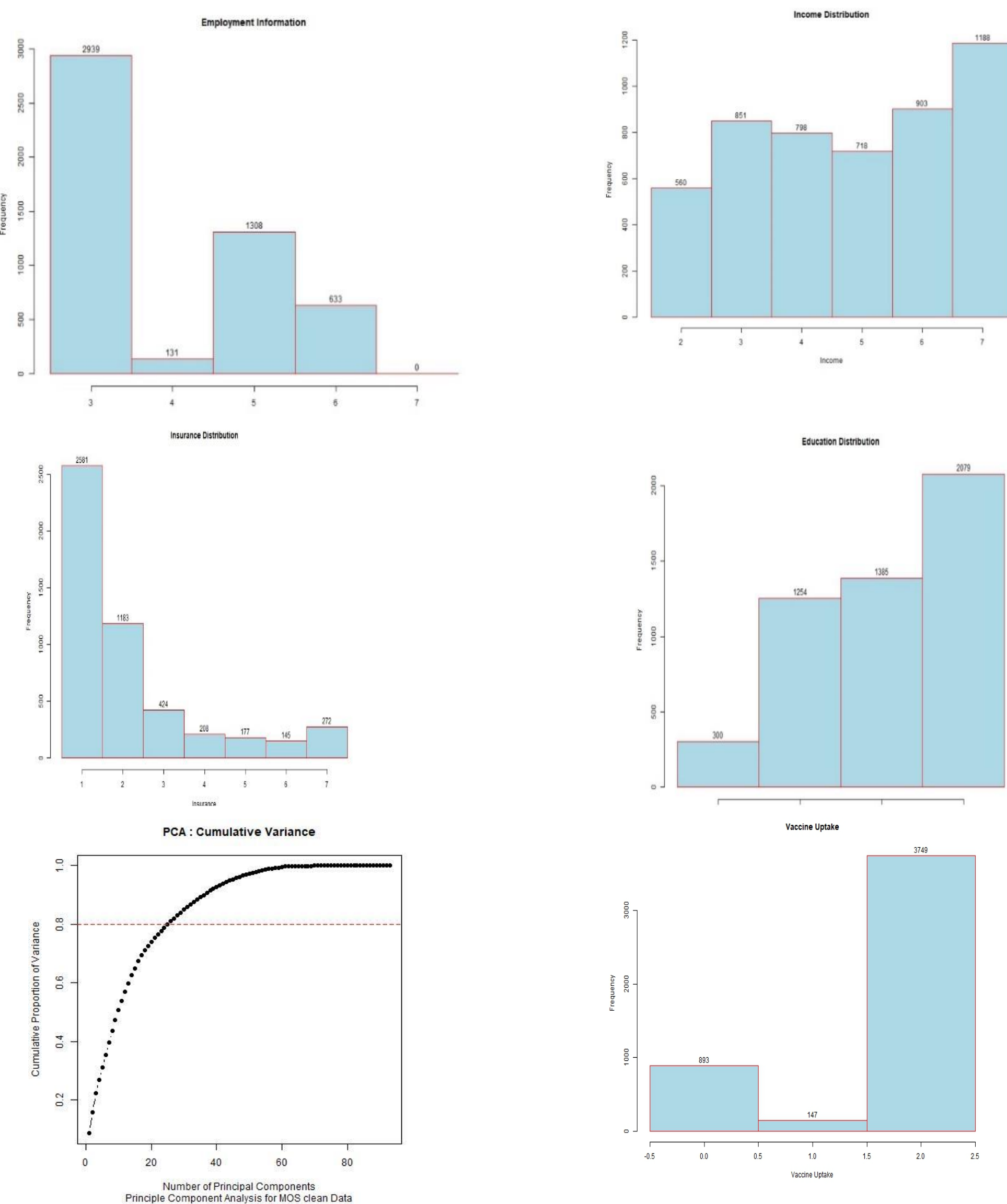
This project examines the factors that influence people's uptake of the COVID-19 vaccine and booster shots. The authors refer to the HSS COVID-19 Monthly Outcome Survey dataset found at HealthData.gov, which contains information on about 5,000 U.S. adults and covers beliefs, intentions, and behaviors related to the COVID-19 vaccine. The dataset contains information on approximately 5,000 U.S. adults covering beliefs, intentions, and behaviors related to the COVID-19 vaccine. The dataset has several features that make it suitable for analysis using classification models. The target variable is CAM5_VaccUptake, which indicates the number of times people have been vaccinated. The authors plan to identify key factors influencing vaccination intentions through exploratory data analysis (EDA) and the application of PCA, random forests, decision trees, SVM, and K-means clustering.

Motivation

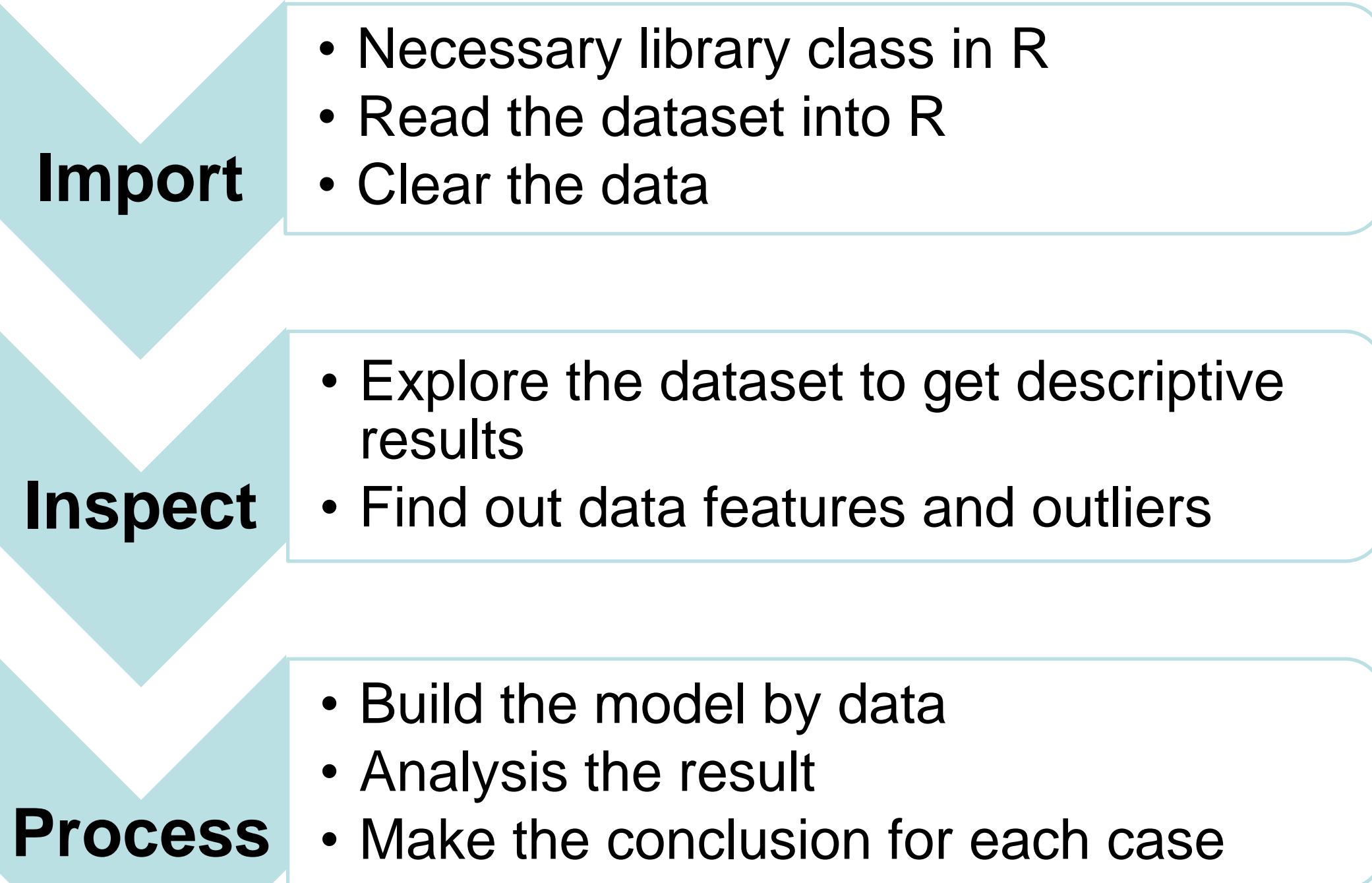
The motivation of this project indicate the High vaccination rates are essential to protect susceptible populations and suggests that the government should develop policies to increase vaccination rates, such as mandatory vaccination for federal employees or workers in the health care industry.COVID-19 has caused more deaths than World War II and emphasizes the importance of vaccines and booster shots to reduce deaths.

EDA and data cleaning

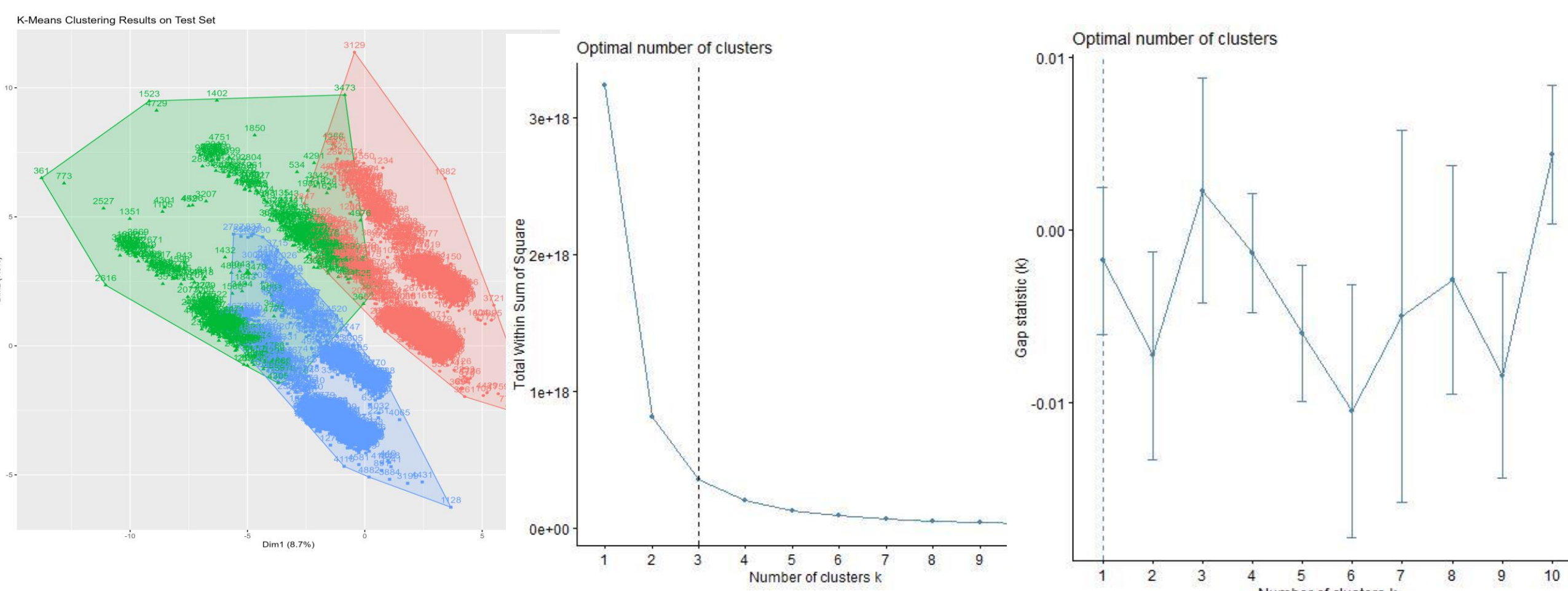
For the data cleaning , remove people refuse to answer and invalid data column like data and ID for each participant. Then use the histogram to see income, education, insurance, employment factor distribution. The next is do the PCA to see how each feature contribution to my target value CAM5_VaccUptake. As we can see around 31 feature contribution over 80% of Cumulative Variance. This indicate my hypothesis is in the correct direct which some have strong effect on the target value.



Module Development And Analysis

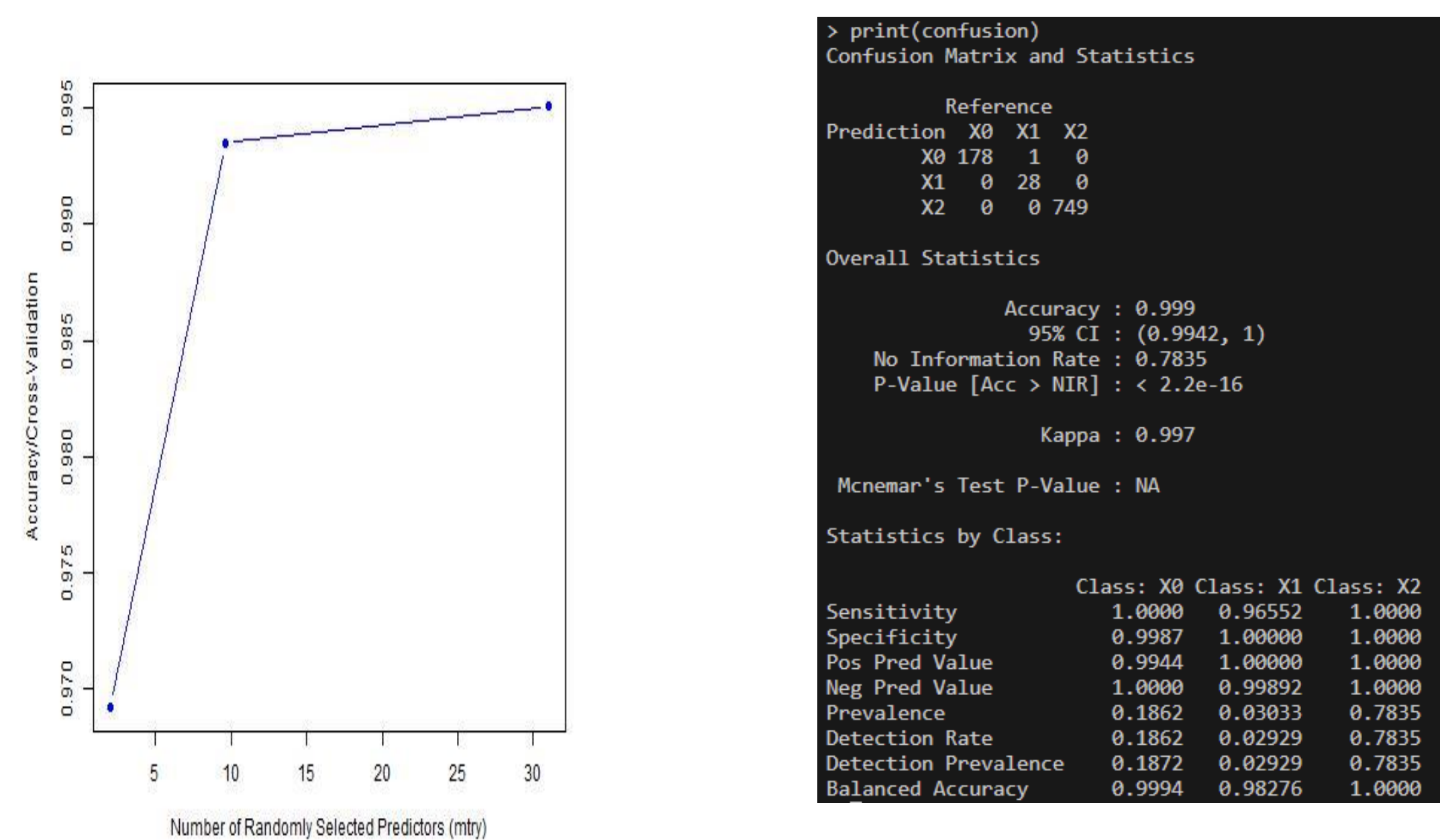


K-mean



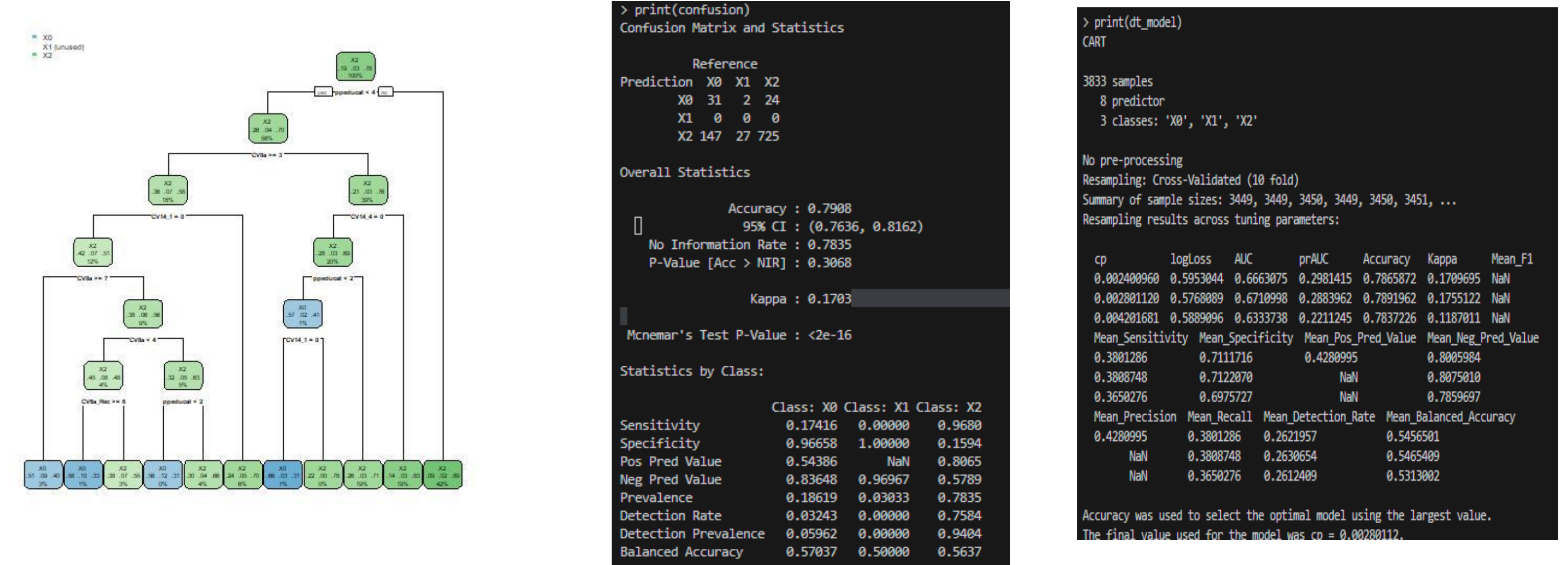
The fviz_nbclust (elbow rule) and clusGap methods were used to determine the optimal number of clusters, k. Both methods indicated that k=3 was the appropriate number of clusters. However, model performance analysis showed that there was overlap between clusters and that the sum of squares between groups was only 13% of the total sum of squares, suggesting that the model failed to effectively differentiate between different groups. Therefore, the application of k-mean clustering in this project is not satisfactory.

Random Forest



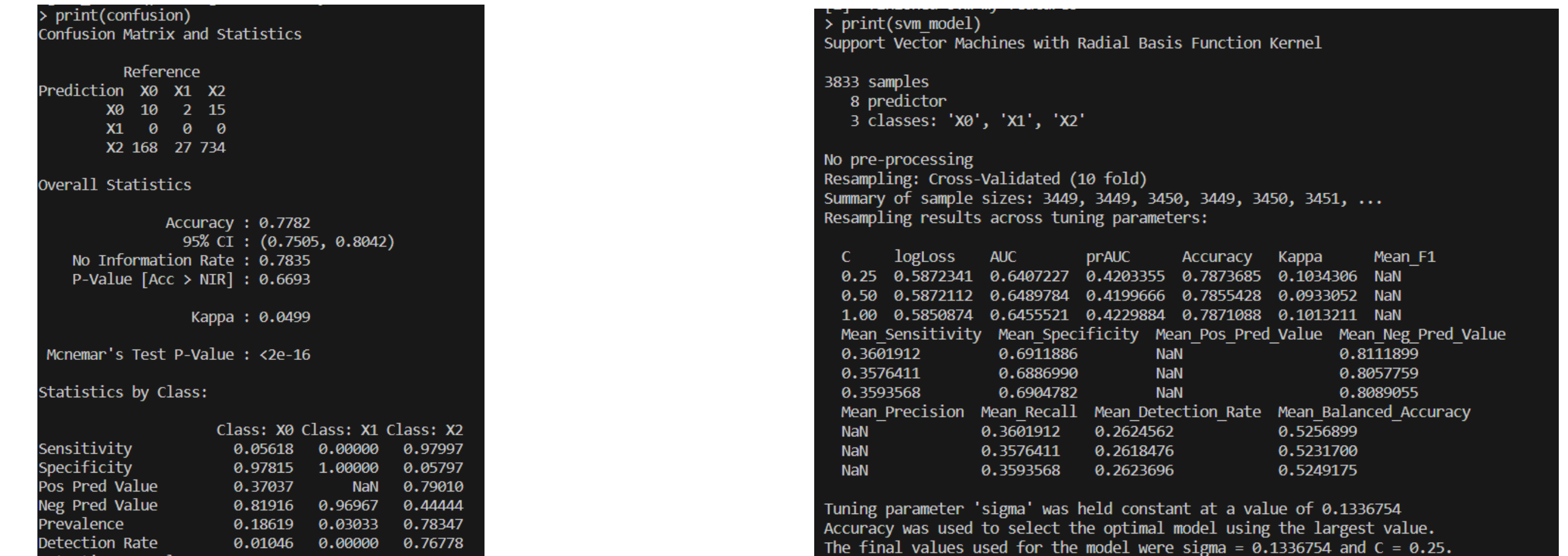
This model can handle complex relationships and does not require data cleaning. By splitting 80% of the training set and 20% of the test set, as well as k-fold cross-validation, the authors determined the optimal mtry parameter to be 31, indicating that choosing 31 features for prediction is optimal. The model performs well on the test set with 99.9% accuracy and sensitivities for each group close to 1. The results also show that the accuracy is close to 99% even with only 5-10 features, indicating that only a few key features are needed for effective prediction. This suggests that vaccination rates can be significantly improved by controlling for a few key variables. The authors plan to further analyze this in depth using decision tree modeling.

Decision Tree



In the project, the authors used a decision tree model for classification and the best model parameters were determined after segmentation and K-fold cross-validation on 80% training set and 20% test set. The model had an accuracy of 79.08% on the test set, but performed poorly in predicting those who were vaccinated only once, with a sensitivity of 0. This may be due to the small number of people in the sample who were vaccinated only once. In addition, since the time point selected for the dataset was near the end of the epidemic, this may have affected the performance of the model, suggesting that temporal factors and changes in vaccine policy have an impact on vaccination behavior.

Support vector machine(SVM)



The dataset was first divided into 80% training set and 20% test set and the optimal C and sigma parameters were determined by K-fold cross validation. Eight features were selected to train the SVM model using Gaussian kernel. The test results show an accuracy of 78% when sigma is 0.1336754 and C is 0.25. Application to the test set showed that the model was inaccurate in predicting the category of only one vaccination with an accuracy of 77.82% and a 95% confidence interval ranging from 0.7505 to 0.8042. The p-value of McNemar's Test was very low indicating significant differences in prediction accuracy between categories. The SVM model was biased in predicting some of the categories, which could be due to training imbalance in the categories of the data or improper selection of model parameters.

Summary

Factors Influencing COVID-19 Vaccination Behavior A variety of classification models were used for this, including random forests, decision trees, support vector machines (SVMs), and K-means clustering, and factors such as income, education, insurance, employment, and good behavior were found to have a significant effect on vaccination. Adjustments in data preprocessing and model selection during the project had a significant impact on the final results. In future studies, there are plans to explore additional data balancing techniques and the application of more advanced models, such as deep learning networks, as well as the implementation of more comprehensive feature engineering. Despite the positive preliminary results, more research and analysis are needed to more fully understand the factors that influence vaccination behavior.

Resources:

- [1] U.S. Department of Health & Human Services.. HHS COVID-19 Monthly Outcome Survey Wave 28. HealthData.gov. Retrieved from https://healthdata.gov/Health/HHS-COVID-19-Monthly-Outcome-Survey-Wave-09/6itx-cwvh/about_data
 - [2] U.S. Census Bureau. (2022, March 17). Educational Attainment in the United States: 2021. Retrieved from <https://www.census.gov/newsroom/press-releases/2022/educational-attainment.html>
 - [3] Peter G. Peterson Foundation. (2023, November). The Share of Americans Without Health Insurance in 2022 Matched a Record Low. Retrieved from <https://www.pgpf.org/blog/2023/11/the-share-of-americans-without-health-insurance-in-2022-matched-a-record-low>
 - [4] Centers for Disease Control and Prevention. (n.d.). Weekly Updates by Select Demographic and Geographic Characteristics. Retrieved from https://www.cdc.gov/nchs/nvss/vsrr/covid_weekly/index.htm
- Link to GitHub repository:
https://github.com/wei-jun7/DataAnalytics2023_WeiJunLi