

H1N1 Vaccine Prediction

Group 10

Allen Lee
Eric Chen
Hyoungmin (Stella) Lee
Marianna Carini
Smitha Kannanaikkal

11:52

USA TODAY

Getting a vaccine has been a huge undertaking. How all 50 states scramble to dole them out is the next massive challenge.

ELIZABETH WEISE | USA TODAY
6:15 am EST Dec. 8, 2020

The largest mass vaccination campaign ever attempted in the United States could begin as soon as this week, with the federal government turning over millions of doses to the states and territories.

Everything depends on them.



Like Dislike ← →

Introduction

- Using National 2009 H1N1 Flu Survey to predict vaccine demand in the United States
- Focusing on quantity of production

Business Problem

Why is this important?

- Pharmaceutical companies need an estimate as to not overproduce or underproduce
- Would benefit agencies like the World Health Organization and healthcare providers to plan distribution and receival of large quantities of vaccines

Methodology

Exploratory Data Analysis

Model Building

Model Evaluation

Data Cleaning

Preprocessing

Conclusion



Dataset Details

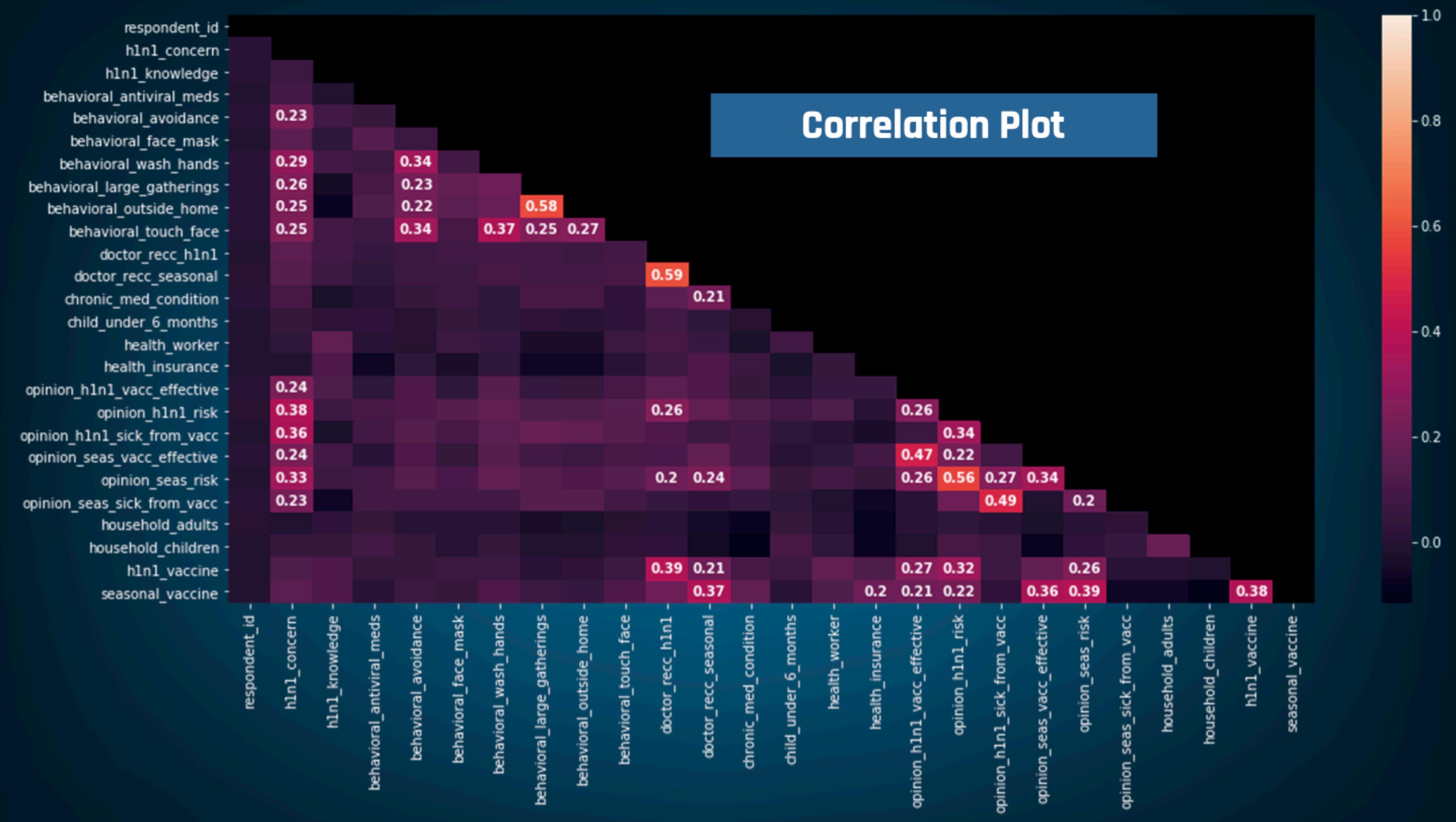
DRIVENDATA

Number of attributes: 38

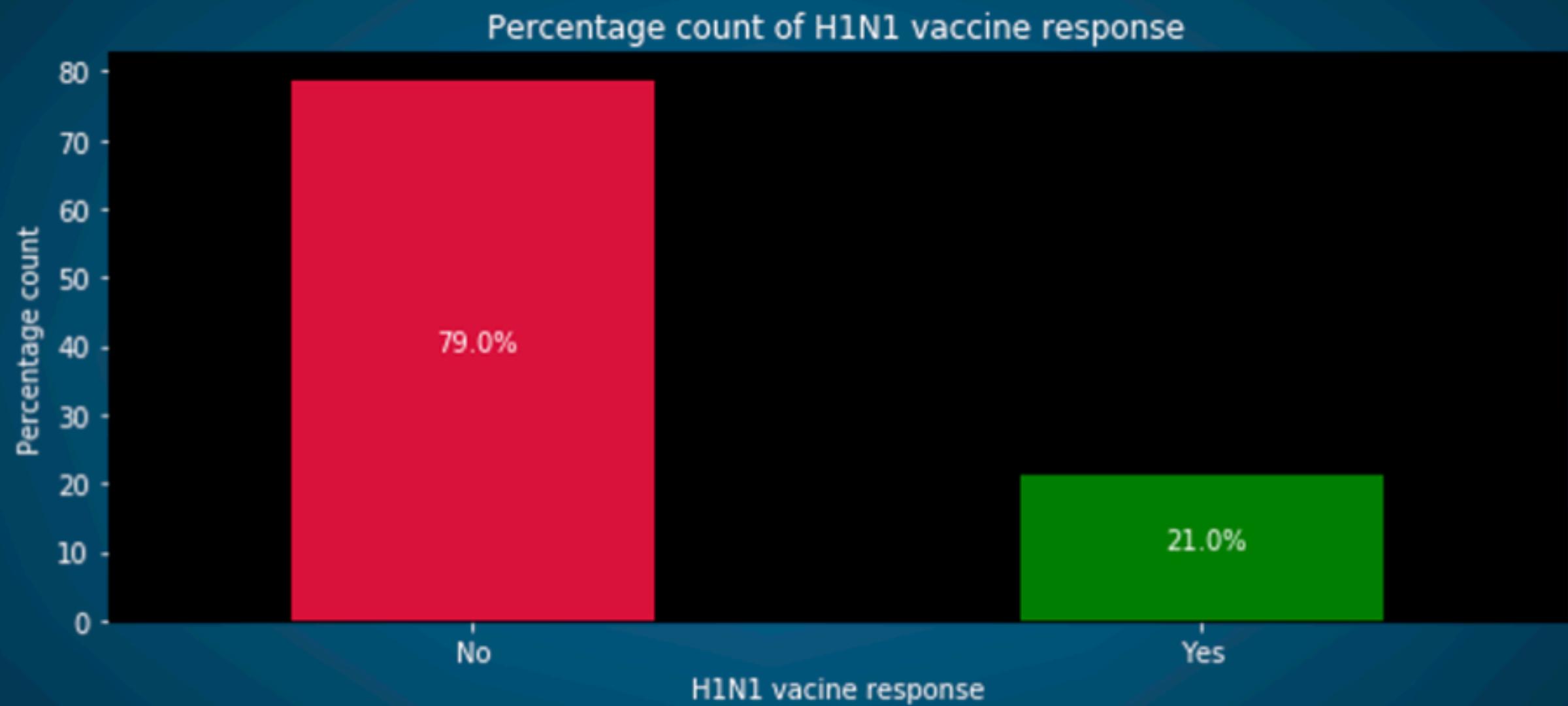
Number of Instances: 26706

| Behavioural | Opinion | Demographic | Others |
|--|--|---|--|
| <ul style="list-style-type: none">• antiviral_meds• avoidance• face_mask• wash_hands• large_gatherings• outside_home• touch_face | <ul style="list-style-type: none">• h1n1_vacc_effective• h1n1_risk• h1n1_sick_from_vacc• seas_vacc_effective• seas_risk• seas_sick_from_vacc• h1n1_concern | <ul style="list-style-type: none">• age_group• education• race• sex• income_poverty• marital_status• rent_or_own• household_adults• household_children• employment_status• census_msa | <ul style="list-style-type: none">• h1n1_knowledge• doctor_recc_h1n1• doctor_recc_seasonal• chronic_med_condition• child_under_6_months• health_worker• health_insurance• h1n1_vaccine• employment_industry• employment_occurrence• hhs_geo_region |

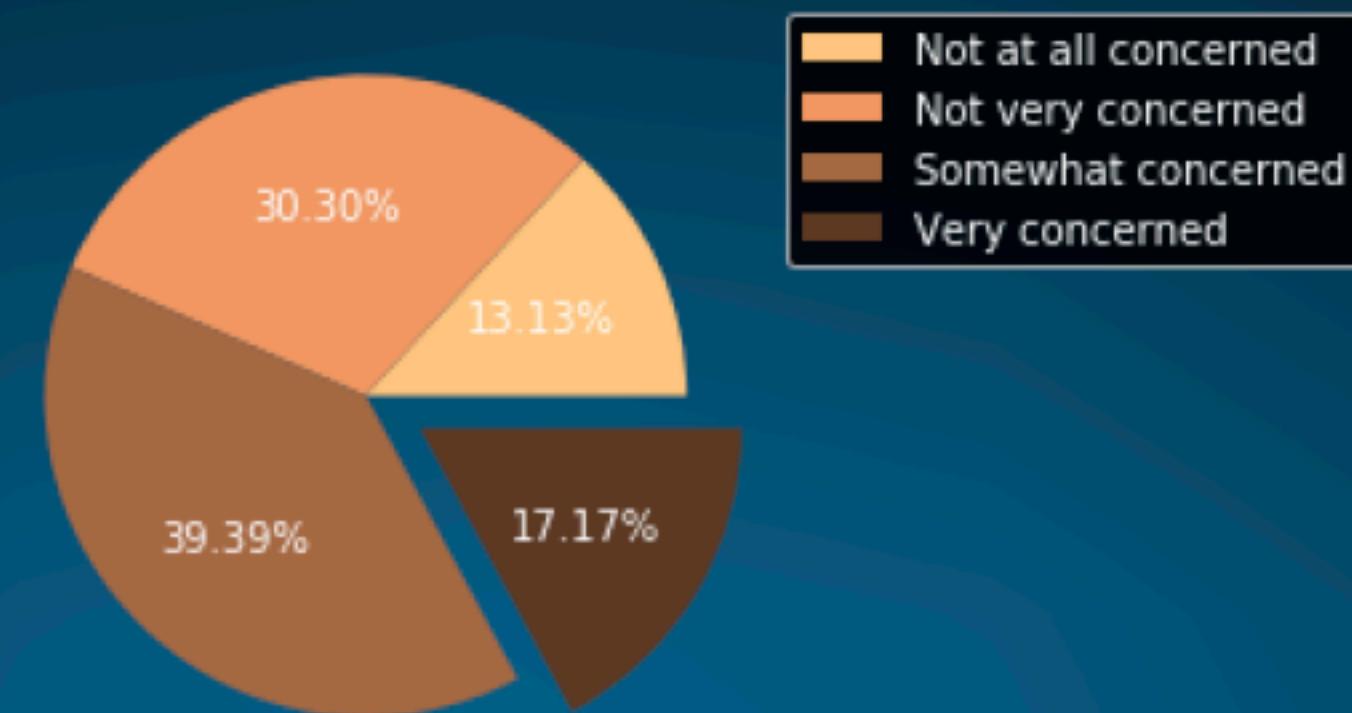
Exploratory Data Analysis



Class Imbalance



Percentage count of H1N1 concern level

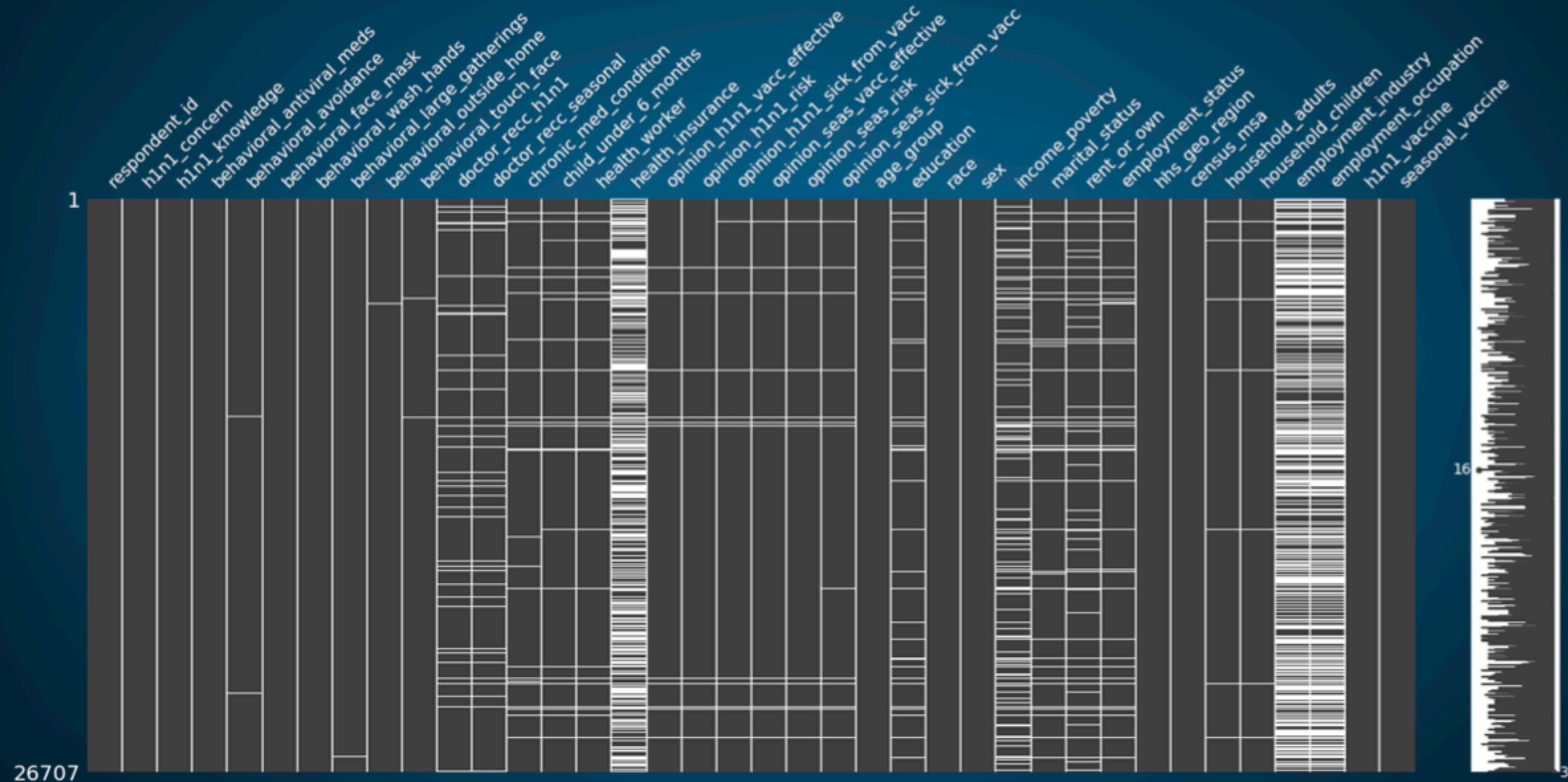


Percentage of H1N1 vaccine response across each level



Data Cleaning

Missing Value in Data



Missing Value Imputation

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Naïve Bayes
Imputed attribute
by prediction



Null Pattern
Imputed attribute
based on the Null
value pattern



Median/Mode
Imputed attribute
based on central
measure

Model Building

Machine Learning Models

Model 1

**Naïve
Bayes**

Model 2

**Decision
Tree**

Model 3

**Random
Forest**

- Several versions of each model were built
- Forward feature selection
 - Oversampling
 - 10-fold cross validation

Performance measure

Precision Recall Fscore

Prediction

| Actual | Prediction | |
|----------|------------|----------|
| | Positive | Negative |
| Positive | TP | FN |
| Negative | FP | TN |

$$\uparrow \text{Precision} = \frac{TP}{TP+FP}$$

$$\uparrow \text{Recall} = \frac{TP}{TP+FN}$$

$$\uparrow \text{Fscore} = \frac{2 * \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

Model Performance

| Naïve Bayes | Accuracy | Class 1 Precision | Class 1 Recall | Class 1 F1-Score | Class 1 ROC_AUC |
|-----------------------------|----------|-------------------|----------------|------------------|-----------------|
| Benchmark | 79.99% | 0.526 | 0.633 | 0.574 | 0.828 |
| Forward Selection | 84.35% | 0.692 | 0.474 | 0.563 | 0.849 |
| Oversampling | 73.97% | 0.740 | 0.740 | 0.740 | 0.822 |
| Oversampling Forward Select | 78.25% | 0.793 | 0.769 | 0.781 | 0.856 |

| Decision Tree | Accuracy | Class 1 Precision | Class 1 Recall | Class 1 F1-Score | Class 1 ROC_AUC |
|-----------------------------|----------|-------------------|----------------|------------------|-----------------|
| Benchmark | 84.31% | 0.670 | 0.521 | 0.586 | 0.798 |
| Forward Selection | 85.01% | 0.693 | 0.533 | 0.602 | 0.839 |
| Oversampling | 83.83% | 0.827 | 0.855 | 0.841 | 0.863 |
| Oversampling Forward Select | 83.91% | 0.830 | 0.825 | 0.861 | 0.860 |

| Random Forest | Accuracy | Class 1 Precision | Class 1 Recall | Class 1 F1-Score | Class 1 ROC_AUC |
|-----------------------------|----------|-------------------|----------------|------------------|-----------------|
| Benchmark | 84.69% | 0.721 | 0.460 | 0.562 | 0.868 |
| Forward Selection | 84.78% | 0.671 | 0.554 | 0.607 | 0.637 |
| Oversampling | 91.23% | 0.893 | 0.937 | 0.914 | 0.974 |
| Oversampling Forward Select | 91.48% | 0.896 | 0.939 | 0.917 | 0.974 |

Model Evaluation

Training Set Performance

| Model | Process |
|---------------|--------------------------------|
| Naïve Bayes | Oversampling Forward Select |
| Decision Tree | Oversampling |
| Random Forest | Oversampling Forward Select |

| Accuracy | Class 1 Precision | Class 1 Recall | Class 1 F1-Score | Class 1 ROC_AUC |
|----------|-------------------|----------------|------------------|-----------------|
| 79.25% | 0.800 | 0.780 | 0.790 | 0.865 |
| 85.00% | 0.836 | 0.871 | 0.853 | 0.870 |
| 91.74% | 0.903 | 0.935 | 0.919 | 0.974 |

Testing Set Performance

| Accuracy | Class 1 Precision | Class 1 Recall | Class 1 F1-Score | Class 1 ROC_AUC |
|----------|-------------------|----------------|------------------|-----------------|
| 79.56% | 0.519 | 0.777 | 0.612 | 0.868 |
| 77.70% | 0.491 | 0.747 | 0.592 | 0.758 |
| 81.20% | 0.552 | 0.713 | 0.622 | 0.864 |

Conclusion

- Naïve Bayes was the best performing model in terms of minimizing potential loss in sales/revenue
- Using the confusion matrix we estimate 107.1 million doses will be needed



Thank you

Do you have any questions?