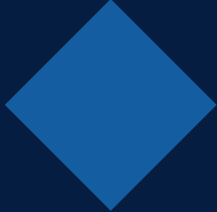# H1N1 AND SEASONAL FLU VACCINE
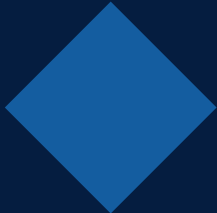
## ADMINISTRATION DATA

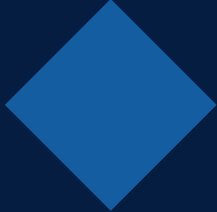*Presented by: Wallace Ouma*

# OVERVIEW

The National 2009 H1N1 Flu Survey (NHFS) was funded by the National Center for Immunization and Respiratory Diseases (NCIRD) and carried out collaboratively by NCIRD and the National Center for Health Statistics (NCHS), both part of the Centers for Disease Control and Prevention (CDC). The NHFS involved a telephone survey of households, utilizing a list-assisted random-digit-dialing method, aimed at tracking influenza immunization coverage during the 2009-10 season.
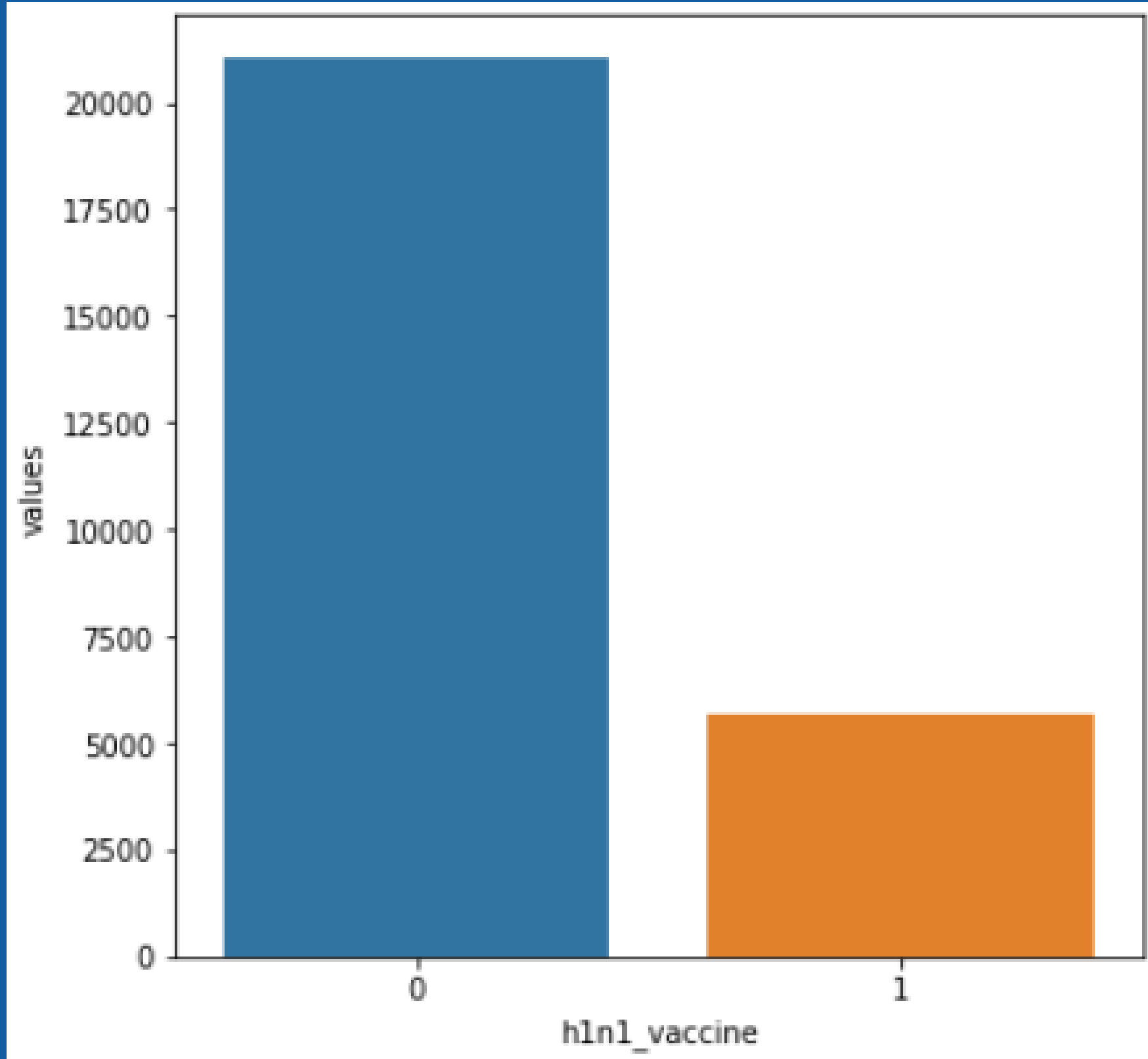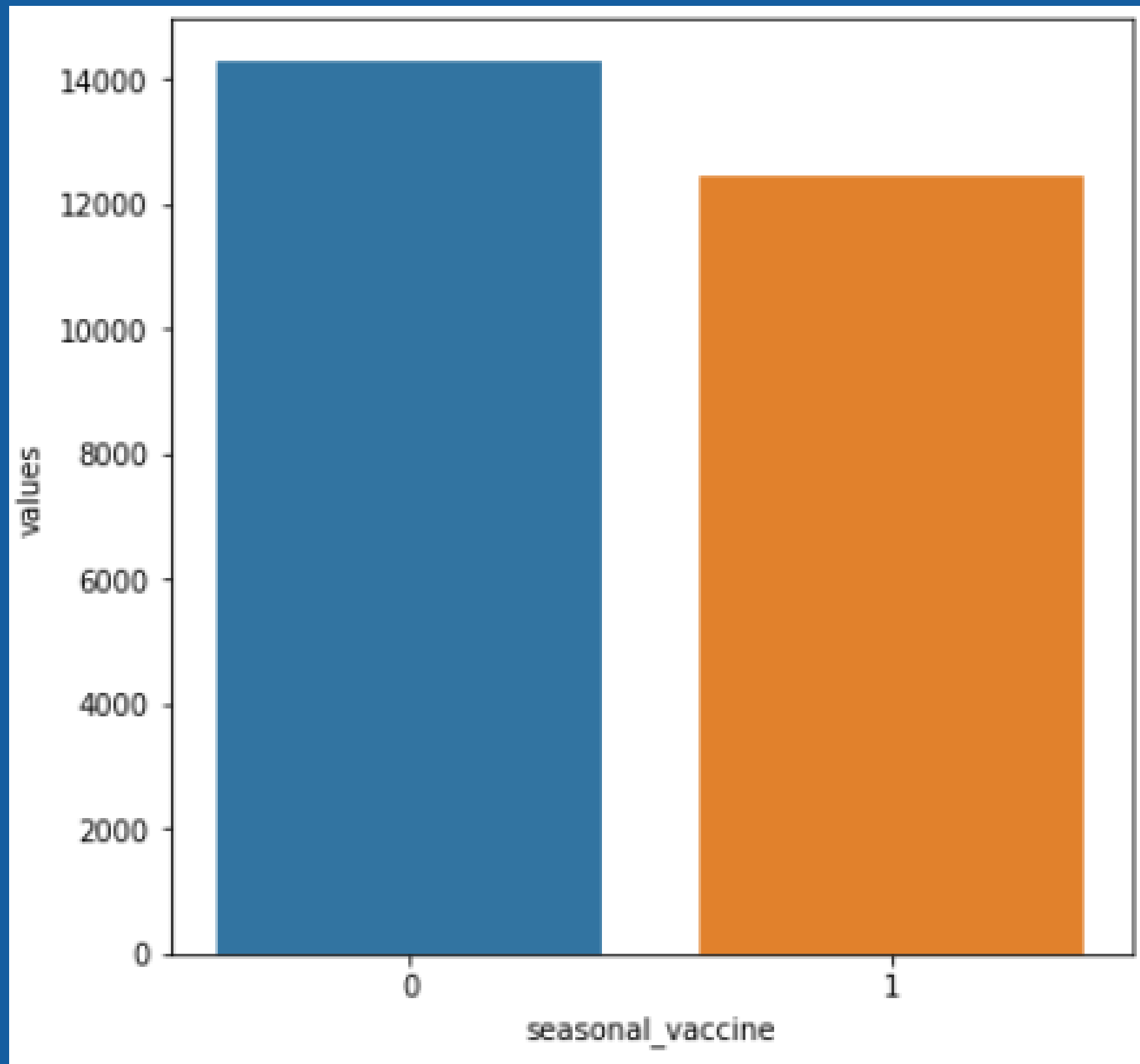
# PROJECT OBJECTIVE

Predict whether a participant will obtain the H1N1 and/or seasonal flu vaccinations based on the details they provided regarding their backgrounds, perspectives, and health practices.

# TARGET VARIABLES

# H1N1 (h1n1_vaccine)

Given that there are nearly four times more individuals who did not receive the H1N1 vaccine compared to those who did, it's necessary to utilize different metrics to evaluate the fit of the models. This imbalance in the dataset may require the use of specialized evaluation techniques to account for the disproportionate class distribution.

## Seasonal Flu Vaccine (seasonal_vaccine)

The number of individuals who received the flu shot is nearly as high as the number of those who did not. This balance contrasts with the difference in H1N1 vaccine recipients, which is more pronounced.
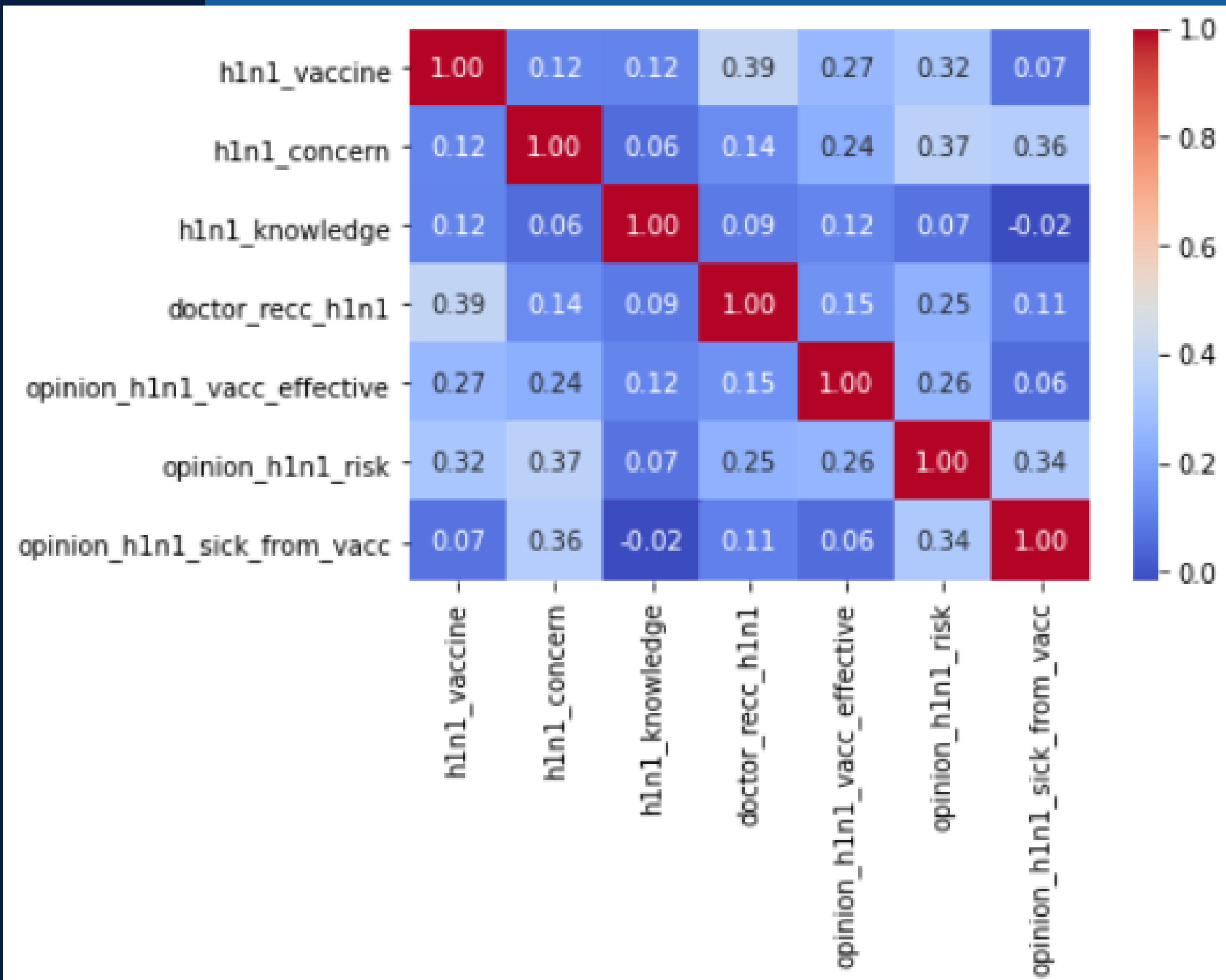
# EXPLORATORY DATA ANALYSIS

# EXPLORATORY DATA ANALYSIS (EDA)

Before proceeding with modeling data for prediction, it's essential to conduct exploratory data analysis (EDA). EDA involves observing and exploring variables along with their interactions to gain insights and understand the underlying patterns in the data.

# FEATURE INTERACTIONS

# H1N1

# SEASONAL FLU

# FEATURE ENGINEERING

# NULL VALUES

Null values were handled by replacing them with the mode (most frequent value) of each feature. This approach was particularly suitable given the prevalence of categorical variables in the dataset. However, due to a significant number of missing values in 'health_insurance', 'income_poverty', 'employment_industry', and 'employment_occupation', these variables were excluded from the analysis.

# ONE-HOT ENCODING

One-hot encoding is a method of representing categorical variables as binary vectors, where each category is represented by a binary value (0 or 1) in a separate dimension. This approach enables the inclusion of more features in the model, strengthening its predictive power compared to traditional methods.

# DATA PREPARATION

# TRAINING AND TEST DATA SPLIT

To accurately assess model performance, the dataset will be divided into training and test datasets. The models will be trained using the training data and then evaluated by comparing predictions against the actual results from the test set. In this scenario, the dataset was split into 80% for training and 20% for testing purposes.

# TRAINING AND TEST DATA SPLIT

Two distinct training and test datasets were divided. One set was designated for training and testing with the target variable labeled as 'h1n1_vaccine', while another set was created with the target variable labeled as 'seasonal_vaccine'.

# FEATURE SCALING

Feature scaling was implemented on both the training and test datasets. This process normalizes the data, making it easier to compute during modeling, and ensures standardization of the functionality range within the input dataset.

# FEATURE SELECTION

# LASSO (L1) REGULARIZATION

Lasso (L1-Regularization) was used for feature selection, which identifies the most significant variables for analysis while nullifying coefficients for less important features to prevent overfitting. This approach enhances model robustness and was applied separately to both the H1N1 and Seasonal Flu datasets.

# MODELING

# MODELS TESTED

- Logistic Regression
- K-Nearest Neighbors
- Gradient Boosting Classifier
- Decision Tree
- Decision Tree with Tuning
- Random Forest
- Random Forest with Tuning
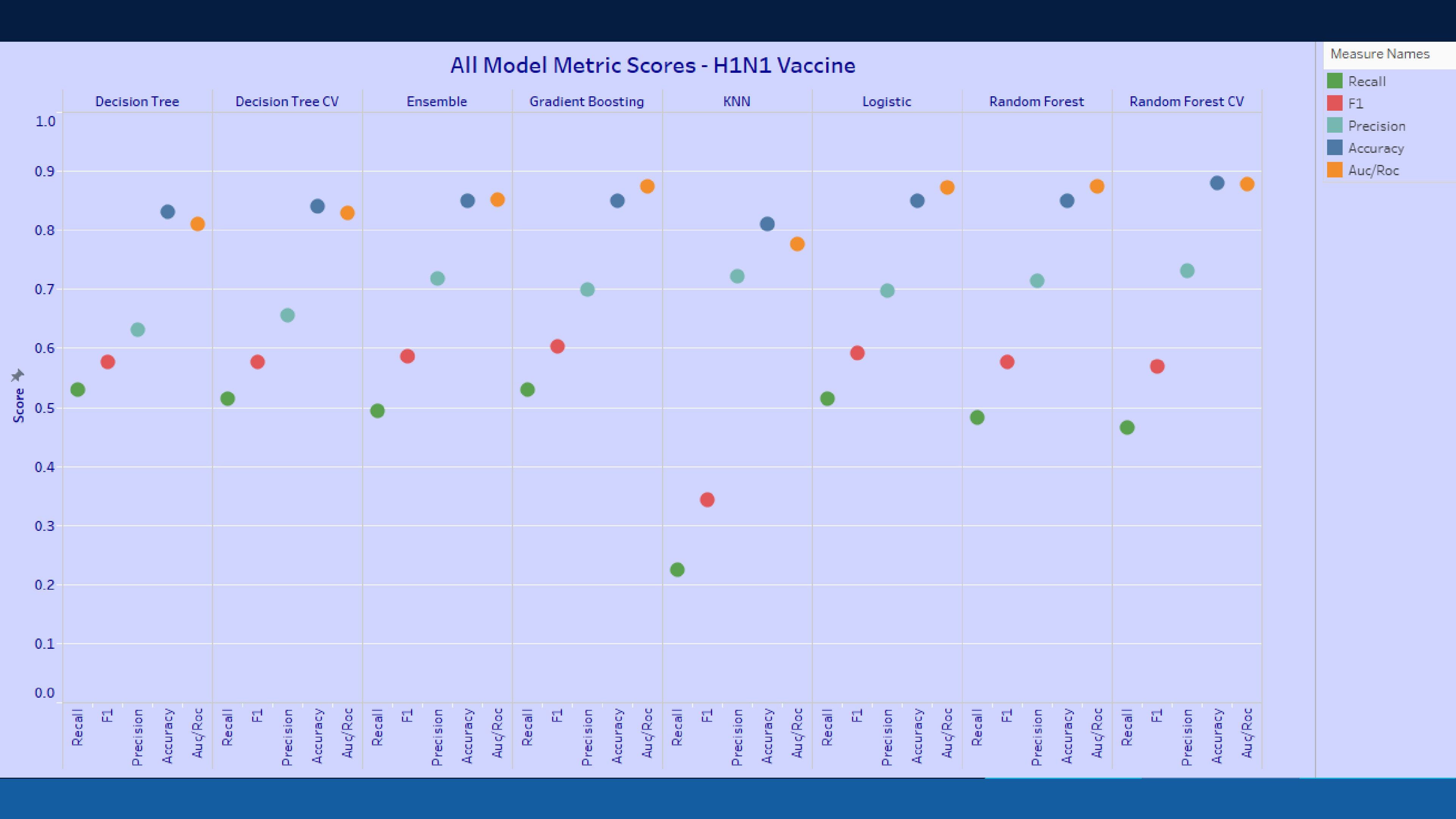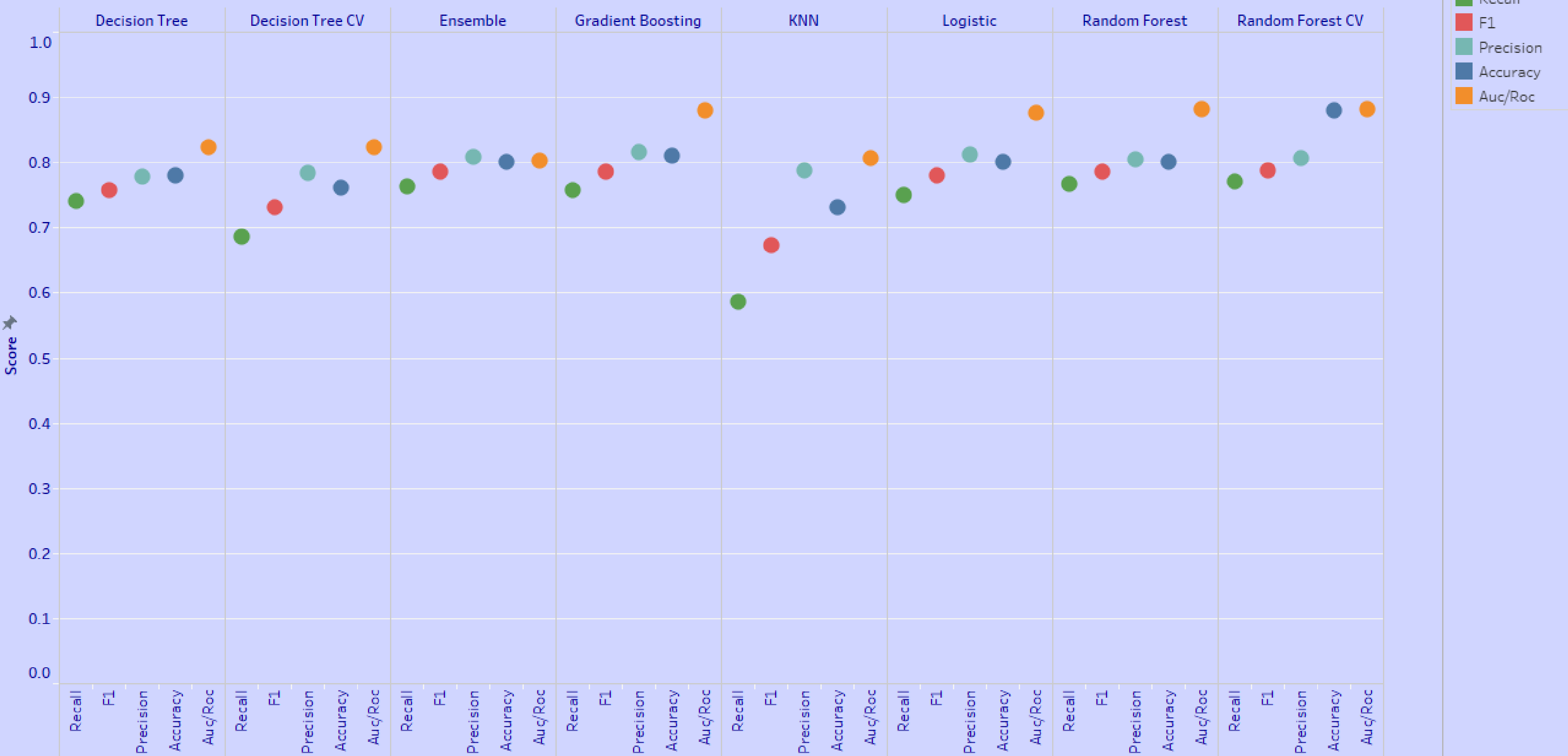- Voting Classifier Ensemble

# METRICS

Accuracy
Precision
Recall
F1 Score
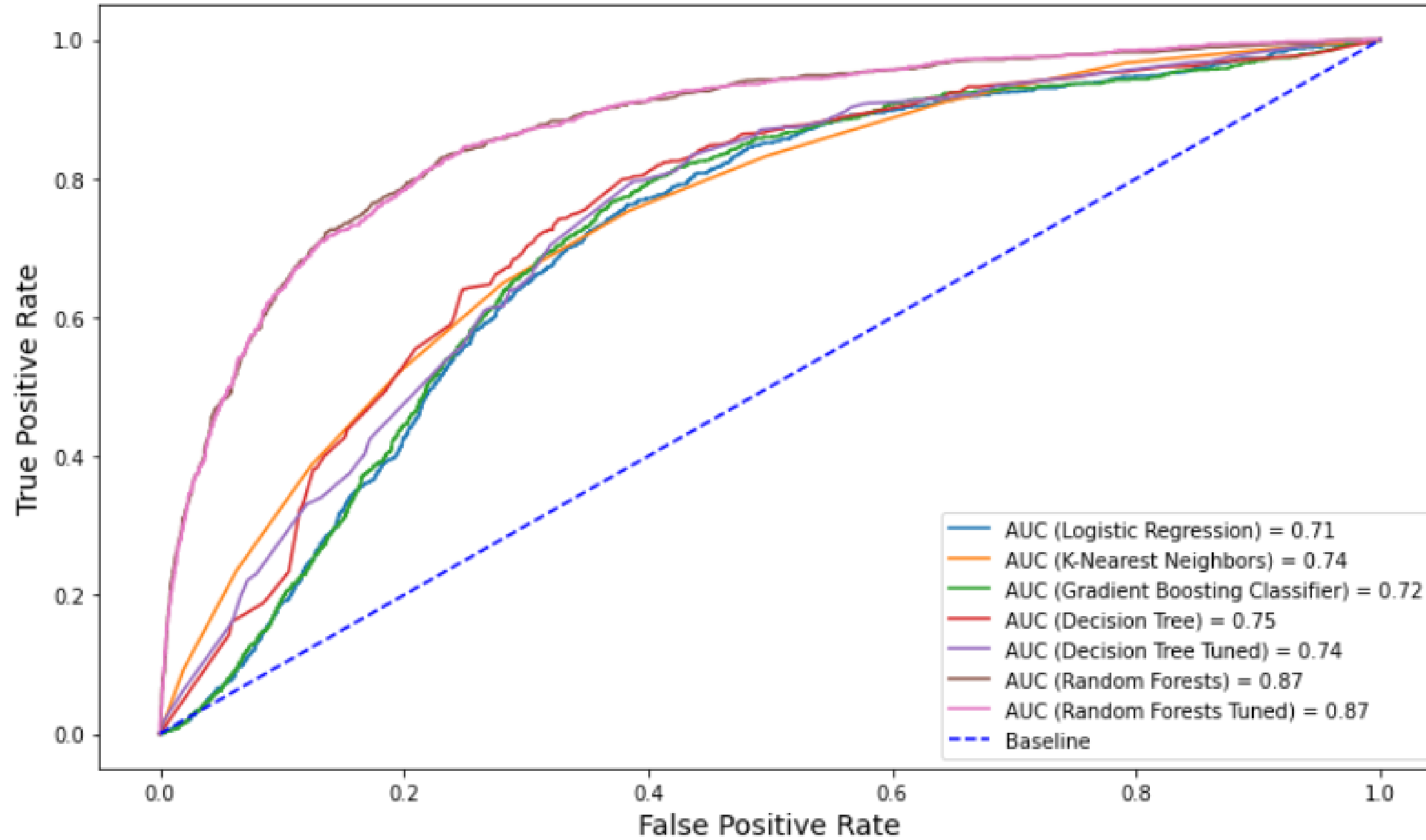Misclassified Samples
AUC/ROC Score

All Model Metric Scores - H1N1 Vaccine

All Model Metric Scores - Seasonal Flu Vaccine
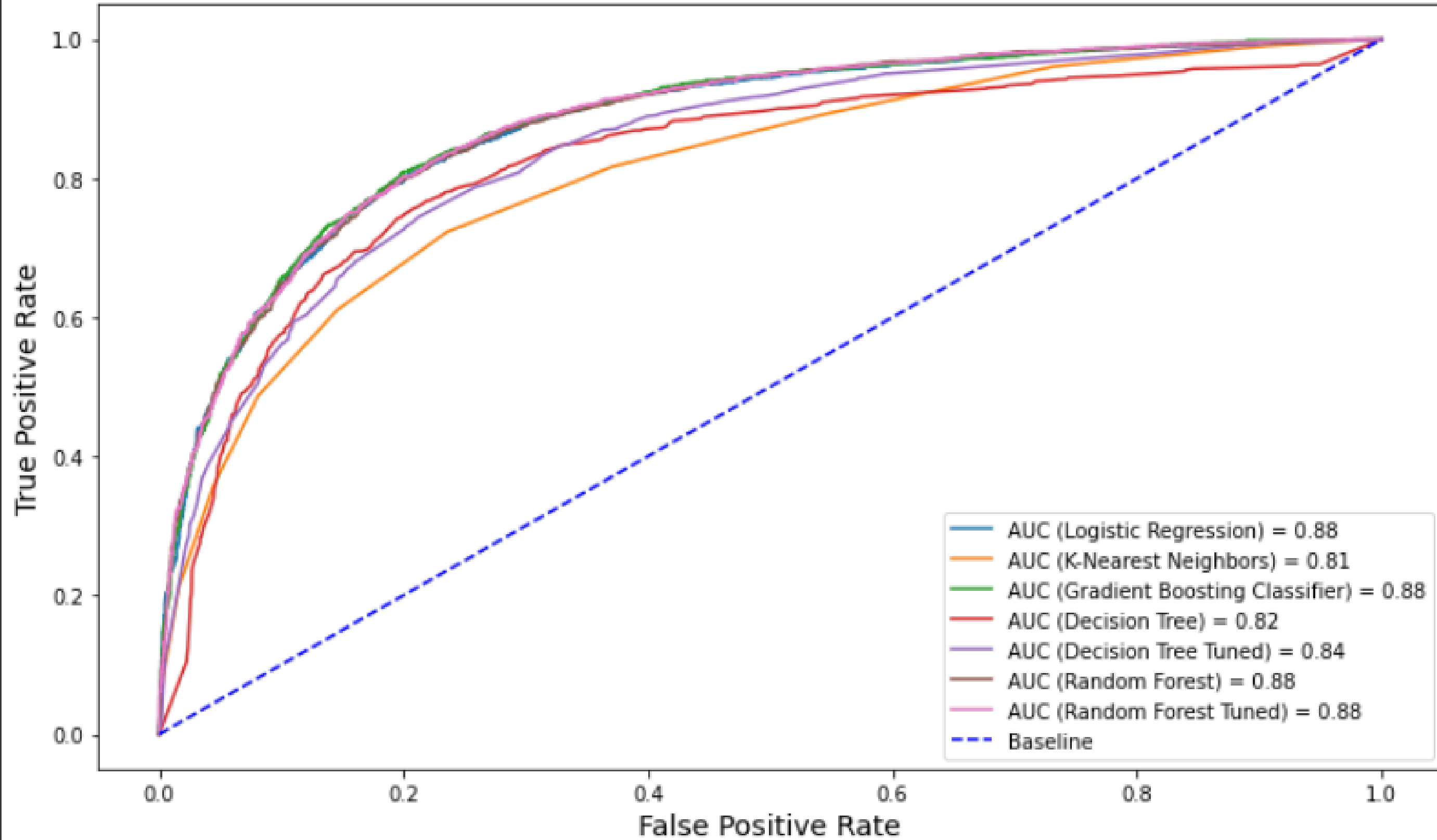
# WHY USE ROC/AUC TO DETERMINE THE BEST MODEL?

The ROC is plotted between true positive (1.0) and false positive (0.0) rates.
The area under the curve (AUC) is the summary of this curve that tells generally how well a model is performing.
The AUC measures the quality of the model's predictions irrespective of what classification threshold is chosen – ideal for an unbalanced problem like ours.

ROC Curve - H1N1

AUC (Logistic Regression) = 0.71
AUC (K-Nearest Neighbors) = 0.74
AUC (Gradient Boosting Classifier) = 0.72
AUC (Decision Tree) = 0.75
AUC (Decision Tree Tuned) = 0.74
AUC (Random Forests) = 0.87
AUC (Random Forests Tuned) = 0.87
Baseline

ROC Curve - Seasonal Flu

AUC (Logistic Regression) = 0.88
AUC (K-Nearest Neighbors) = 0.81
AUC (Gradient Boosting Classifier) = 0.88
AUC (Decision Tree) = 0.82
AUC (Decision Tree Tuned) = 0.84
AUC (Random Forest) = 0.88
AUC (Random Forest Tuned) = 0.88
Baseline

# Proposed Models

| H1N1 | Random Forest with Cross Validation |
|------|-------------------------------------|

| Seasonal Flu | Random Forest with Cross Validation |
|--------------|-------------------------------------|

# QUESTIONS?