

Flu Shot Learning: Predicting H1N1 and Seasonal Flu Vaccination

Jorge Pais, José Baptista and Pedro Duarte

up201904841@edu.fe.up.pt; up201904814@edu.fe.up.pt; up201905050@edu.fe.up.pt

Abstract—

I. INTRODUCTION

Pandemics have rarely taken center stage in the way they have recently with COVID-19 in 2020. Vaccines are a key public health measure used to fight infectious diseases like COVID-19. They provide immunization for individuals, and enough immunization in a community can further reduce the spread of diseases through "herd immunity." In 2009, the H1N1 influenza virus, also known as swine flu, caused a global pandemic that was estimated to have resulted in 150,000 to 600,000 deaths worldwide. A vaccine for H1N1 became available in October of that year. In this study, a machine learning model was developed to help estimate the probability of a person receiving seasonal and H1N1 vaccines. For this, several classification methods were explored and compared between each other in order to figure out which one exhibited the best performance.

II. DATA RESOURCES

For the competition partaken in this project, the dataset was provided by DrivenData, and it comes from the National 2009 H1N1 Flu Survey (NHFS) which was collected through telephone interviews. The dataset consists of 36 attributes, varying from numerical with both ordinal and binary variables, and also categorical attributes. For the training data, there were 2 labels associated with each respondent, indicating whether or not each respondent had taken the H1N1 and Seasonal flu vaccines.

A. Data Cleaning and Pre-processing

The first step in the data analysis was to check if the data has any duplicate (i.e. duplicate respondent ids) and missing values, duplicates were not observed, but the dataset had many missing values in different attributes so, one of the first concerns was to clean the data. Some the attributes, namely employment concern and employment occupation had the highest missing data (13470 values). The *Pandas* library in Python identifies these missing values as NaN, and this can be identified and handled by different imputation classes from *SKLearn*. Another problem is that some of the data is categorical, which can be solved simply by encoding each possible category within a feature.

B. Feature Correlation

The next step involved checking of correlation between the attributes. This was done using the *Seaborn* library's correlation heatmap. It was noted that some attributes exhibited mild positive correlation, with two stand-out attributes named doctor recc H1N1 vaccine and doctor recc seasonal vaccine, which were positively correlated by 60%.

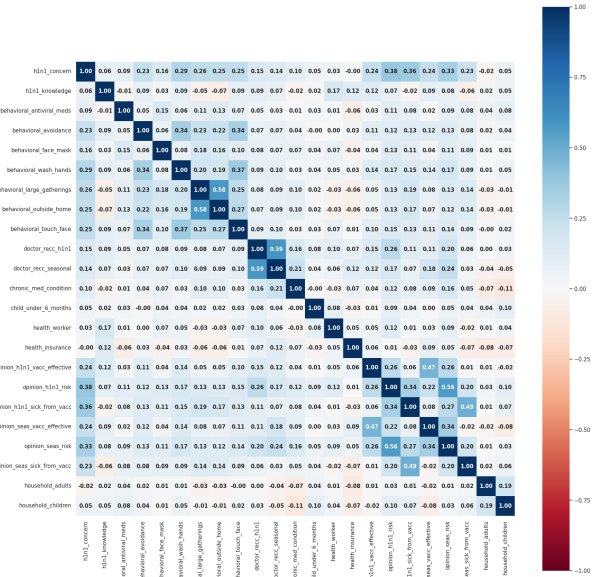


Figure 2.1 - Feature Correlation Heatmap

C. Class Balance and Label Correlation

Observing the distribution of the two target variables, shown in Figure 2.2, approximately half of individuals have been vaccinated for the seasonal flu, while only 20% have been vaccinated for H1N1. As for class balance, we can say that the distribution of individuals vaccinated for the seasonal flu is balanced, while the distribution of those vaccinated for H1N1 is imbalanced.

Taking the correlation (also known as the phi coefficient) between the two target variables a value of approximately $\phi = 0.377$ was obtained. This indicates a positive correlation between the two target variables, meaning that these aren't fully independent from each other.

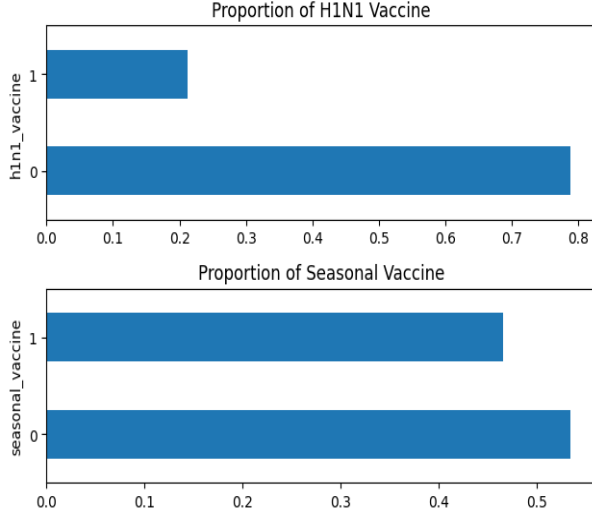


Figure 2.2 - Feature Correlation Heatmap

D. Feature Distributions

Firstly, it is observed that out of the people who received the seasonal vaccine, most of them were female. The same case was also observed with H1N1 vaccine through which one can conclude that women are more prone to get affected than men.

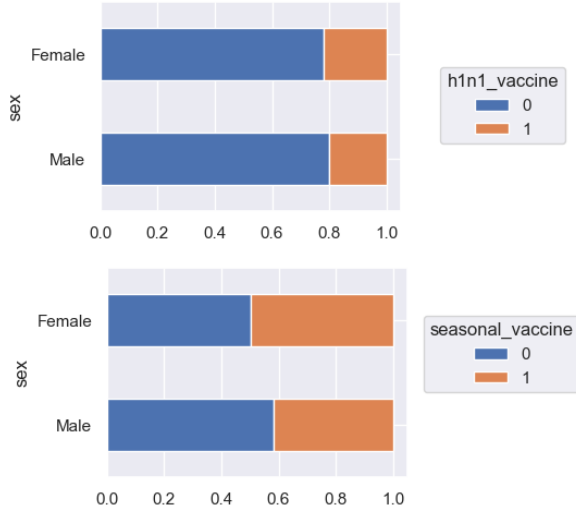


Figure 2 - Vaccination for male and female

The age group has a strong correlation with the seasonal flu vaccine but not with the H1N1 flu vaccine. It seems that people act appropriately when it comes to the seasonal flu as older individuals have a higher risk of complications. However, with H1N1 flu, even though older individuals have a higher risk of complications, they are less likely to get infected. This analysis does not provide information about causality, but it seems that the risk factors are reflected in vaccination rates. It appears that questions related to knowledge and opinions have a strong correlation with both target variables.

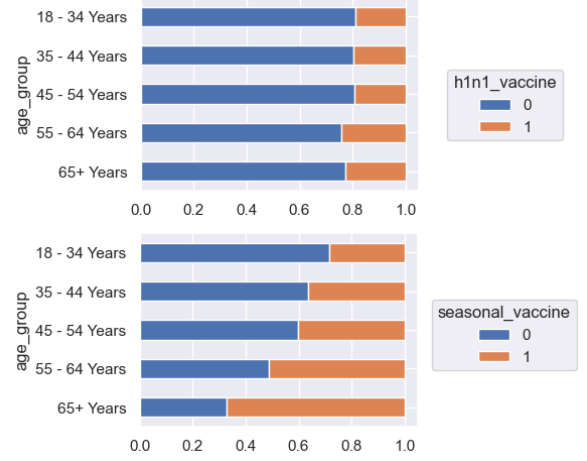
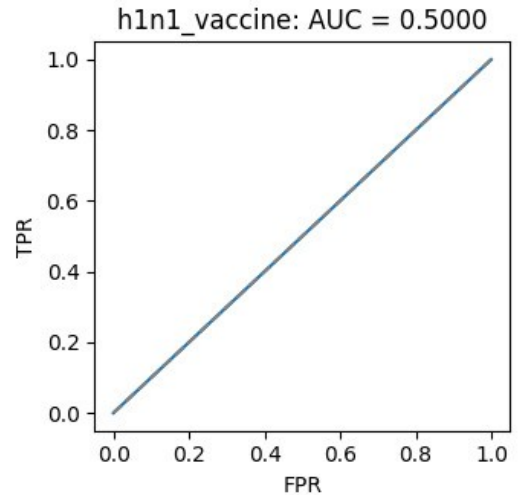


Figure 2.2 - Vaccination for different age groups

III. PERFORMANCE METRIC

To measure the performance of the classifications performed between different classifiers, the ROC (Receiver Operating Characteristic) metric was utilized. The ROC is a type of plot used in binary classifiers, which measures the true positive rate (TPR) against the false positive rate (FPR) for different classifier thresholds. To obtain a quantitative measurement of the performance obtained, it is possible to take the area under the curve (AUC). One way of interpreting AUC is as the probability that the model ranks a random positive example more highly than a random negative example. Note that this area can be between 0 and 1, the latter characterizing a perfect classifier. Firstly, the AUC was calculated for a *Dummy Classifier*. It is a classifier model that makes predictions without trying to find patterns in the data, serving as a simple baseline to compare with other more complex classifiers. The strategy used to generate predictions was "uniform", in order to generate predictions uniformly at random. The value obtained for the AUC of the H1N1 vaccine can be seen in the graph in figure 3.1.

Figure 3.1 - ROC curve for *Dummy Classifier* (H1N1 vaccine)

After that, using the same method, the classification performance of different experimented classifiers was measured. One

Figure 3.2 - ROC curve for *CatBoost* (seasonal vaccine)

Figure 4.1 - Model Structure used

A. Logistic Regression

B. Naive Bayes

C. K-Nearest Neighbors

case of classification, the algorithm computes the k closest neighbors to the observation, and classifies the observation based on what class held the majority among those neighbors. This method can achieve very good results depending on what value of k is utilized. A larger k will generally suppress the effect of noise on the data, but makes the decision boundary less clear which might affect the ROC score. To use the model, *SKLearn* includes the `KNeighborsClassifier` class, which was used in combination with `GridSearchCV` to find the value k that produces the best results.

D. Decision Trees

Decision Trees are models that attempt to match the target variable by learning decisions from the training data and applying these decisions to the features of an observation's features. These models have the advantage of being conceptually easy to understand, but are quite prone to overfitting the training data.

Similarly to the previous models, the *SKLearn* implementation of the decisions trees, `DecisionTreeClassifier`, was utilized.

E. XGBoost

XGBoost (which stands for eXtreme Gradient Boosting), is a library that provides many gradient boosting algorithms. Essentially, gradient boosting gives a result based on an ensemble of many weaker learning models, as for example decision trees, combining these sequentially and having each stage correct the errors of the previous one.

This library provides many different models and several interfaces for different programming environments. For the purposes of this work, the `XGBClassifier` class was used, as it easily integrates with the *SKLearn*.

F. CatBoost

The final model that was utilized in this project was *CatBoost*. Similarly to *XGBoost*, and as the name implies, *CatBoost* is also a gradient boosting library. The advantage of *CatBoost* is that it natively supports categorical features without any need for encoding these. This library provides a classification model, also compatible with *SKLearn*, through the class `CatBoostClassifier`.

This model was utilized in two ways, firstly by integrating it into the pipeline developed previously (in the notebook `Project_mainModel.ipynb`, including categorical feature encoding) and in a separate preprocessing pipeline in the file `Project_catBoostedModel.ipynb`. This was done since the first the pipeline utilized *SKLearn*'s `ColumnTransformers`, which made it challenging to specify what columns were categorical in *CatBoost*.

V. RESULTS

A. Result Analysis

Examining the experimental results for all the non gradient boosting models, the following results were obtained:

Model	ROC AUC score
Logistic regression (C = 0.01)	0.84555
Logistic regression (C = 0.1)	0.84655
Logistic regression (C = 1)	0.84640
Logistic regression (C = 10)	0.84635
Gaussian Naive Bayes	0.72947
Multinomial Naive Bayes	0.79467
Decision Tree	0.66738
KNearestNeighbors (k=169)	0.82478

Table 5.1 - Scores for some of the models tested

It is possible to see that all the Logistic Regression models exhibited the best classification performance, having pretty much the same score for all the values of C . This parameter is essentially the inverse of the regularization strength. The worst performing model out of all, was clearly the Decision tree classifier. Observing the ROC plot for the Decision Tree, shown in figure 5.1, it is possible to see that it is composed of two straight lines, due to this classifier not being able to estimate the probability of each label, but instead assigning a hard label to each observation.

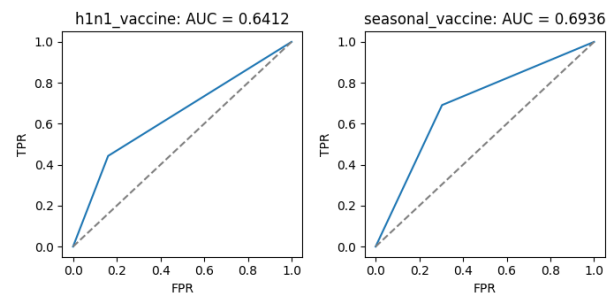


Figure 5.1 - ROC plot for the Decision Tree classifier

In terms of the Naive Bayes classifiers, it is possible to see that the multinomial distribution approximation presents much promising results than the Gaussian approximation.

Taking a look at the results K-Nearest Neighbors, first the score plot of Figure 5.2 was obtained by varying the number of neighbors used and averaging the score across 10 folds of the dataset. It is possible to observe that above a certain threshold of k , that the score seems to stagnate.

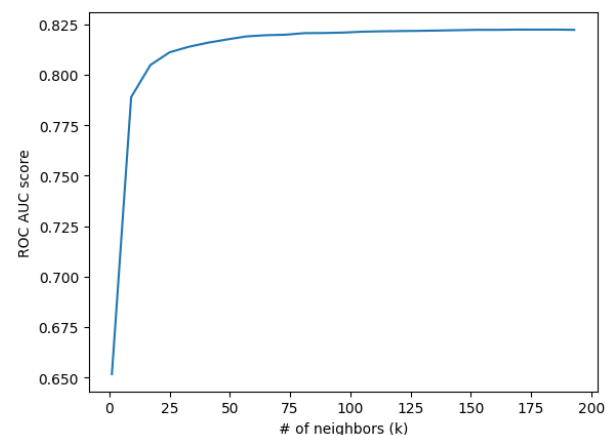


Figure 5.2 - K-Nearest Neighbors ROC AUC Score for varying number of neighbors

As for the gradient boosting models, the results were the following:

Model	ROC AUC score
XGBoost (Default parameters)	0.83732
CatBoost with One-Hot encoding	0.85069
CatBoost with specified cat_features	0.86907

Table 5.2 - Scores for the gradient boosting models

It is possible to see that *CatBoost* performs much better than *XGBoost*, which itself performed worse than the Logistic Regression models used initially. These results also make evident that to realize the full potential of *CatBoost*, there isn't much preprocessing needed, only simple imputation and categorical column identification is needed.

B. Competition Results

Applying the unseen competition test data to some of the best performing models, the following results were obtained:

Model	ROC AUC score
Logistic Regression (C=0.1)	0.8356
K-Nearest Neighbors (k=169)	0.8122
CatBoost with One-Hot encoding	0.8430
CatBoost with specified cat_features	0.8616

Table 5.2 - Scores for the gradient boosting models

With these results, it was very clear that *CatBoost* demonstrated exceptional performance when compared to every other model tested. The final score of 0.8616 granted our team, at the time of writing, the 285th position on this competition's leaderboard. The top submission at the time of writing is a 0.8658, meaning that our classifier's performance is within 0.5% of the top 1 result.

VI. CONCLUSION