

Predicting Appointment Cancellation and Isolating Contributing Factors (Medumo Group 6)

Jessica Sandler, Heather Johnson, Nicholas Pearce, Jacob Kozol
{sandlerj, heathdj, npearce, jkozol}@bu.edu

1. Project Task

When patients miss their appointments, it leaves empty spaces in hospital schedules that could otherwise be used to serve other patients. The goal of our project is to help hospitals predict which patients will not appear at their appointments so they can use that space to help other patients waiting for appointments. We also want to find which demographics are the most effective at predicting this information. We have decided to focus on the prediction model portion as determining which demographics does not involve machine learning techniques.

2. Related Work

[1] and [2] use, respectively, multinomial and binary logistic regression as the primary methods for solving the classification problem of predicting whether or not a patient will cancel/fail to appear at their appointment. [2] also uses L2-norm regularization to prevent overfitting, and 10-fold cross validation to assess the accuracy of their model.

3. Approach

Taking the factors of our data set as X , we will use binary logistic regression since we are only looking for the binary classification of took place (0) or missed appointment (1) with L2-norm regularization to counteract potential overfitting. We will also use the cross-entropy cost function in conjunction with gradient descent to optimize our solution. Because we are in search of the most accurate model, we have decided to expand our project to include multiple binary classification models. Aside from logistic regression, we will also create a support vector machine and a binary neural network classifier. We have begun implementing logistic regression and training it on some preliminary data.

4. Dataset and Metric

Our dataset has 1996 training examples and 222 test examples. Each example currently includes an index ID, Hospital ID, patient ID, registration date, procedure date, and whether or not the patient appeared at their

appointment. We were also given engagement data for each patient, as well as some basic demographic and enrollment information for each patient (gender, age, registration for email and SMS notifications). Before we begin to create our model, we should first compile all of this data so that all of the details are in one table, rather than two or three. However, over time we may need to modify our dataset, as one of our goals is to pinpoint which factors have the highest predicting ability. We have begun preprocessing our data by determining which features we want to include. The features we have so far include the time between registration and appointment date, the patient's age and gender, and if they are signed up for email and/or SMS notifications. We also plan to include if the appointment date is on a weekday or weekend and what season it falls in. We are working on determining the number of interactions a patient has with the application, how many modules they completed, and whether or not they added a family member to receive notifications.

The metric that we will use to evaluate success will be the receiver operating characteristic (ROC) curve, which plots the true positive rate (TPR) against the false positive rate (FPR). In this case, the TPR is defined as the number of correctly predicted cancellations/no shows divided by the total number of cancellations/no shows, and the FPR is defined as the number of successful appointments that have been predicted to be cancellations/no shows divided by the total number of successful appointments. The specific metric we will use is the area under the curve (AUC) of the ROC curve. The AUC is equal to the probability that a randomly chosen positive instance will be ranked higher than a randomly chosen negative instance by the classifier. Our goal is to achieve an AUC between .8 and .9.

5. Preliminary Results

We have created a basic logistic regression model using a subsection of the data with gender, SMS notifications, email notifications, and the difference

between the scheduling date and enrollment date as features. We got 79% accuracy with this model, however this was because our model was predicting that every patient would show up to their appointment, and in that dataset 79% of the patients showed up to their appointments. This says to us that we (obviously) need to add more features, which we plan to do as soon as we link our files together by patient ID and process the user interaction data for features that we can feed into our models.

6. Timeline and Roles

Task	Deadline	Lead
Data Preparation	11/18/18	Nick
Implement Logistic Regression	11/26/18	Jacob
Implement SVM	11/26/18	Heather
Implement Neural Network	11/26/18	Jessica
Validation/Analysis	throughout	Nick
Prepare Project Update	11/29/18	all
Prepare report and presentation	12/11/18	all

References

- 1) A. Alaeddini, K. Yang, P. Reeves, and C. Reddy. *A hybrid prediction model for no-shows and cancellations of outpatient appointments. IIE Transactions on Healthcare Systems Engineering*, 5(1):14-32, 2015.
- 2) H. Kurasawa, K. Hayashi, A. Fujino, K. Takasugi, T. Haga, K. Waki, T. Noguchi, and K. Ohe. *Machine-Learning-Based Prediction of a Missed Scheduled Clinical Appointment by Patients With Diabetes. Journal of Diabetes Science and Technology*, 10(3):730–736, 2016.

*Note: All italicized text was previously in the project proposal, anything else is new.