

Machine Learning Engineer Nanodegree

Capstone Proposal

Edgar Valente
January 13, 2019

Domain Background and Problem Statement

In the city of Vitória (state of Espírito Santo, Brazil) there are numerous public health facilities. In these, the number of patients that don't attend to their medical appointments has reached nearly 30% from 2014 to 2015. This pattern is similar all over the country, so it isn't a local issue. This indice represents nearly R\$20 million of public cash being wasted by appointments that result in no-shows, taking into account operational costs, SMS sent, confirmation call and employee's time. That's at least R\$1 million monthly waste.

Datasets and Inputs

The used dataset is available on [Kaggle as version 3](#). I'll rewrite their names during the analysis for further clarity, with their updates registered in the dictionary between brackets ('[]'). It contains 300,000 observations and 15 variables (300000, 15), them being:

- **Age**: integer, age of person to make the appointment [age]
- **Gender**: categorical, gender of person to make the appointment, allowing either male or female [gender]
- **AppointmentRegistration**: integer, the day the person asked for an appointment [app_registration]
- **AppointmentData**: date, the day of the scheduled appointment [app_date]
- **DayOfTheWeek**: categorical, the day of the week of the scheduled appointment [week_day]

- **Status**: categorical label, indicating if the person showed up (show-up) or not (no-show) on the appointment [show-up]
- **Diabetes**: categorical binary, whether the person has diabetes [diabetes]
- **Alcoholism**: categorical binary, whether the person is an alcoholic [alcoholism]
- **HiperTension**: categorical binary, whether the person has hypertension [hypertension]
- **Handicap**: categorical binary, whether the person has a handicap [handicap]
- **Smokes**: categorical binary, whether the person is a smoker [smokes]
- **Scholarship**: categorical binary, whether the person receives monetary help from the government, called 'Bolsa família' [monetary_help]
- **Tuberculosis**: categorical binary, whether the person has tuberculosis [tuberculosis]
- **SMS_Reminder**: categorical binary, whether the person has received an SMS reminder for the appointment one day before it [sms_reminder]
- **AwaitingTime**: integer, days from appointment registration to scheduled appointment [days_to_appointment]

Solution Statement

The idea is to identify patterns between people that miss appointments and the ones that show-up and, possibly, influence those patterns. Since at the registration of the appointment it is possible to choose the date, I'm hoping to be able to recommend different dates for different people to optimize show-ups. Even if the recommendation is not possible, I intend to develop a prediction model to at least identify possible no-shows, since it would be possible to act on them with reminders, if applicable.

Benchmark Model

Since there is no way to make an AB Test during the modeling, I intend to create up to two prediction models, depending on the evaluation metrics results: one focusing on f1-score and another focusing on 0.5-score, detailed on the Evaluation Metrics section. I intend to use Logistic Regression (so I have a basis for benchmarking) and some ensemble algorithm - possibly Random Forests or Gradient Boosting. Then I intend to create a clustering model - possibly a Gaussian Mixture Model - to check if there are patterns between no-shows and show-up. If there are patterns, I intend to recommend a day of the week for different test individuals and check on the prediction model if the individual will show-up. I'll be able to tell if the recommendation is good based on the

amounts of show-ups taking into consideration the prediction model's accuracy or f1-score.

Evaluation Metrics

For the main prediction model, I'm going to use confusion matrix with f1 and fbeta scores. I'm going to prioritize precision over recall (using $\beta = 0.5$), since I want the predicted values to be true positives, so I can try and recommend only on highly assertive no-shows. For validating the recommendation, I'm gonna use f1-score, since I wanna be as sure as possible if the recommendation is correct. For the clustering model, I'm gonna use silhouette score.

Project Design

The input data will be cleaned for clarification. Preprocessing for categorical variables (week_day and gender) will be applied so that I can apply the algorithms. I'll have either two or three models:

- Prediction model to identify no-shows;
- Prediction model to benchmark recommendations;
- Clustering model to identify patterns between show-ups and no-shows and try to recommend a day of the week for the appointment;

The models will be created using python's scikit-learn library, using its Pipelines where applicable.

The dataset, both original and edited will be found in a 'data/' directory, while any necessary Jupyter Notebook in a 'notebooks' directory.