

Machine Learning Capstone Project

Final Report

Definition

Project Overview

In recent years, technology has made it easier to bridge the gaps between siloed data sources, and provide a unified view of data across platforms, customer relationship management systems and marketing automation software. This enabled a new breed of products based on machine learning that aims to successfully predict whether a lead will convert into a customer.

Companies like Lattice Engines, Fliptop, 6sense and Infer are using machine learning to power sales and marketing teams, enabling them to quickly identify key characteristics of ideal target markets by uncovering purchasing patterns based on historical customer behavior and demographic data. In addition to historical data, some of them are using machine learning to analyze incoming data from external sources, like web searches and social media, to provide an even better understanding of who are the most promising customers.

As shown in the article "Probabilistic Modeling of a Sales Funnel to Prioritize Leads", by Brendan Duncan and Charles Elkan from the Department of Computer Science and Engineering at UCLA, from 2015, machine learning models can result in up to 307% increase in number of successful sales, as well as a dramatic increase in total revenue.

In fact, nearly two thirds of businesses have implemented predictive marketing analytics, including 42% that are expanding or upgrading their implementation within the next 12 months, according to Forresters From Insight to Action: How Predictive Analytics Improves B2B Marketing Outcomes, published in October 2015. Among businesses using predictive marketing analytics, 83% have experienced a positive business impact as a result of their implementation.

Selected predictive marketing analytics platform use cases

Lead prioritization	Net-new leads	Cross-sell/Upsell	Account-based marketing	Personas and segment building	Sales enablement
<ul style="list-style-type: none">• 6Sense• EverString• GrowthIntel• Infer• Lattice Engines• Leadspace• Mintigo• Radius• SalesPredict	<ul style="list-style-type: none">• 6Sense• EverString• GrowthIntel• Infer• Lattice Engines• Leadspace• Mintigo• Radius• SalesPredict	<ul style="list-style-type: none">• 6Sense• EverString• Lattice Engines• Mintigo• Radius	<ul style="list-style-type: none">• 6Sense• EverString• GrowthIntel• Infer• Lattice Engines• Leadspace• Mintigo• SalesPredict	<ul style="list-style-type: none">• Infer• Lattice Engines• Leadspace• Mintigo• Radius• SalesPredict	<ul style="list-style-type: none">• EverString• GrowthIntel• Infer• Lattice Engines• Leadspace• Mintigo• SalesPredict

Many predictive marketing analytics vendors are rooted in lead scoring but have broadened their capabilities to include predictive modeling, personalization, and product recommendations that push deeper into the purchase funnel. However, none of them are built with a specific market segment in mind. Despite the fact that there is no established market leader among these players, and that the crowded field

continues to attract new ventures, there is no specific solution tailored to the education market. This is our motivation.

Problem Statement

In June 2016 Udacity started to operate in Brazil. Since then, we have seen exponential growth in our Nanodegree students: average compounded weekly growth is 10-12% for the past 8 months. Our sales funnel is composed by the following stages:

1. Awareness: User discovered us through social network ad campaigns, PR, live events, or organically (SEO), and landed in our website;
2. Lead: User created a free account;
3. Marketing-qualified lead (MQL): User watched webinar(s), visited specific pages, opted-in our newsletter, followed Udacity in social network(s), and/or started a free course;
4. Sales-qualified lead (SQL): User visited our checkout more than once, made inquiries about how our Nanodegree programs work, and/or tried to pay (unsuccessfully);
5. Student in trial: User started Nanodegree trial;
6. Paying student: User successfully concluded the trial and/or (directly) enrolled in a Nanodegree.

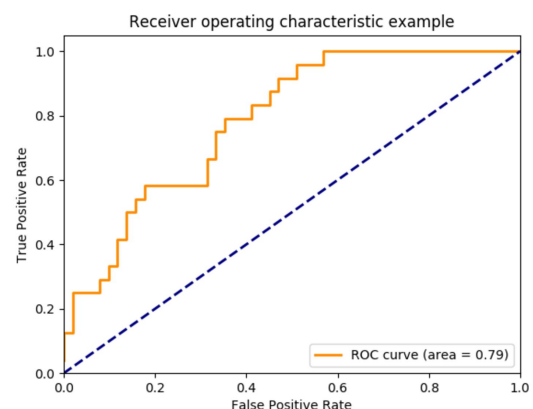
Although we currently use marketing automation, we are not using any user behaviour data to drive them toward the end of the funnel: we simply do promotional email campaigns to generate urgency to purchase. Usually, they are a one-size-fits-all campaigns targeted to MQLs. This raises the question:

Given all behaviour data we collect from our users through the sales funnel, can we train an algorithm to successfully predict whether or not they are potential Udacity Nanodegree students?

Metrics

Our dataset is very imbalanced: only approx. 10% of our samples are Udacity Nanodegree students. In order to avoid [accuracy paradox](#), we will use **ROC AUC score** to choose the best algorithm, and then to tune its parameters.

ROC AUC score means area under the ROC curve. ROC stands for Receiver Operating Characteristic, because it was created by Radar Engineers in World War-II. However, in machine learning the curve has nothing to do with radar signals or signal receivers. In an ROC curve, we plot 'True Positives' on Y-axis and 'True Negatives' on



X-axis: the greater the area under the curve, the better our model is.

Analysis

Data Exploration

In order to tackle the proposed problem, we used the datasets in the following table. Important to highlight: emails, names, or any other data that could identify users were not used.

Dataset name	Description	Rows	Total fields	Field examples
auth_user	Contains key data about users	60,859	11	id, last_login, username, first_name, last_name, email, date_joined
payment_app_product	Contains key data about Nanodegrees	26	12	id, name, code, price
payment_app_subscription	Contains all subscription data from paying students	5,760	39	id, access_until, cancel_requested, status, register_date, credit_card_retries, chosen_payment_type_code, product_id, user_id, first_payment_date, cohort_id, full_amount, instalment_amount
frontend_brazil.pages	Contains data from all website visits	4,707,359	47	anonymous_id, category, context_ip, name, path, referer, sent_at, timestamp, title, url, user_id, context_campaign_medium
frontend_brazil.identifies	Contains data that links visitors with user accounts	113,630	34	received_at, anonymous_id, user_id
frontend_brazil.tracks	Contains data from all events tracked in marketing website	1,092,373	31	received_at, event, event_text, timestamp, user_id, anonymous_id, context_ip
brazil_events.event_sign_up	Contains data from all webinars	26,931	54	email, enrollment_date, event, event_title, event_start_date, event_end_date, slug, user_id, user_interests, event_type
analytics_tables.course_enrollments	Contains data from all free course enrollments		15	user_id, course_key, join_time, leave_time, course_title, course_level

We gathered, cleaned and prepared all data in the **new_features.ipynb** Jupyter Notebook file. We have chosen to focus on **Digital Marketing students that enrolled in any moment after 2017-04-01 and before 2017-08-10**. We made it for one reasons: Digital Marketing Nanodegree is the largest program Udacity has in Brazil, with over 50% students - therefore, we have more data. We used as total universe users who are registered and visited Digital Marketing Nanodegree Overview Page (DMND NDOP) at least once in this period.

These are the general steps we followed in this phase:

1. First, we retrieved the data from **32,544 visitors** who came to DMND NDOP in the period.
2. We then retrieved who became paying student in this period - the target feature we wanted our model to predict. We stored that in the column **is_paying_student**, with **0 or 1**. We were able to collect data from **2,952** paying students.
3. Then, we started to gather data, that would be used to train our models. Our first hypothesis was that the likelihood of a student to enroll was influenced by how recent they had an account registered at Udacity. So, retrieved the joining date of each account and created **age_in_days** feature.
4. Doing webinars was one of our key strategies to nurture leads. So, we counted how many webinars each user watched, and saved this in **webinar_enrollments** feature.
5. Similarly, we counted how many times each user enrolled in a free course, and saved this in **course_enrollments** feature.
6. Then, we counted how many times each user visited each one of our **top 25 pages** in our website. We saved this data in 25 columns, each one with the prefix **is_**: the column **is_home** counted how many times a user visited our home page, while the column **is_checkout** counted how many times a user visited our checkout, for example.
7. We also counted how many visits were done in a mobile device, saving that in **is_mobile** column.
8. Then, we counted how many times each visited was generated from our **top 15 referrers**. We saved this data in 15 columns, each one with the prefix **is_referrer_**: the column **is_referrer_google** counted how many times a user came to our website from Google, while the column **is_referrer_facebook** counted how many times a user came from Facebook, for example.
9. One of the key strategies we used to nurture leads was email marketing. We counted how many times each user opened our emails, saving that in **opened_emails** column.
10. Finally, we reordered the 46 columns and saved the data in **data_prepared.csv** file.

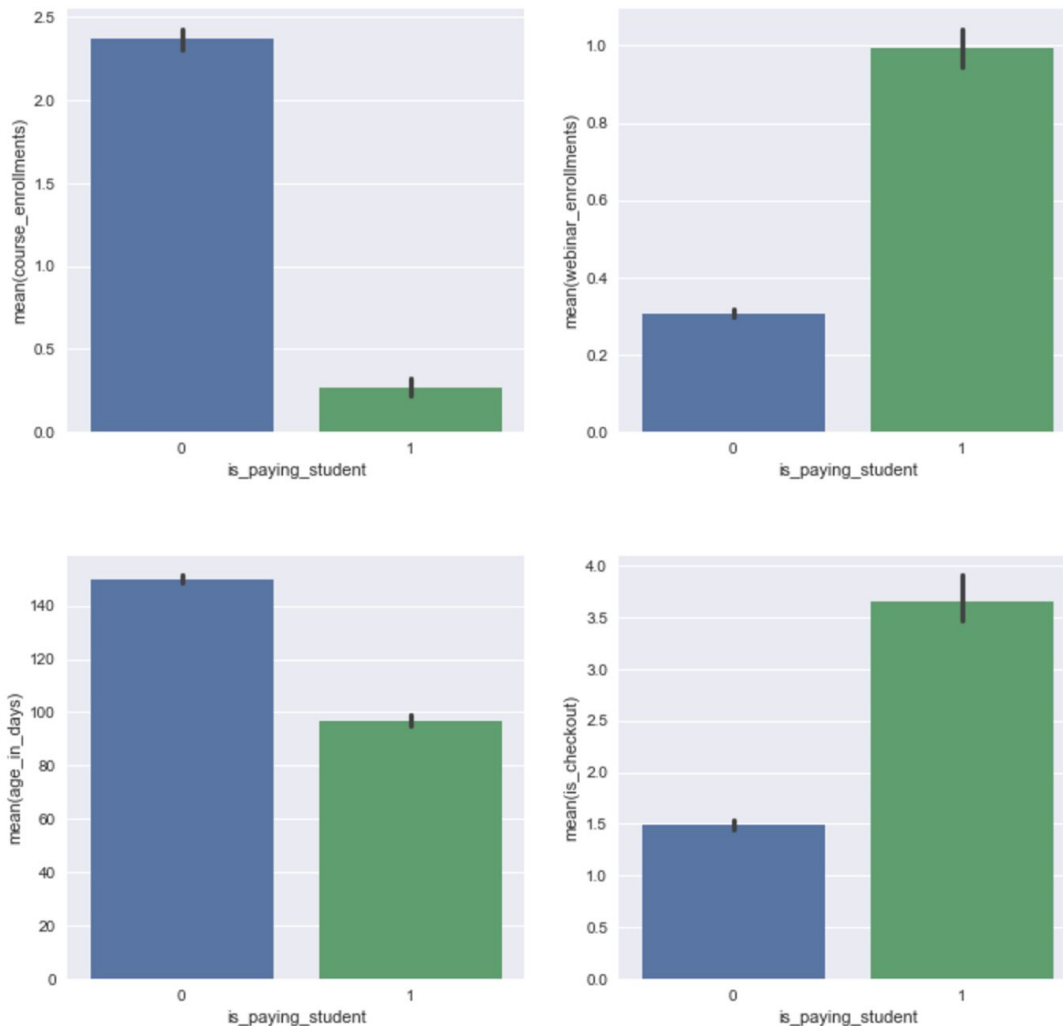
Before start training models, we checked statistics of our dataset main features. In the table below you can see a sample of that. For a complete description of all feature, please check **capstone_final_project.ipynb**.

	count	mean	std	min	25%	50%	75%	max
opened_emails	31991	29	49	0	3	11	35	1211
age_in_days	31991	145	116	1	60	112	196	417
webinar_enrollments	31991	0.37	0.93	0	0	0	0	21
course_enrollments	31991	2.2	5.2	0	0	0	2	128
is_checkout	31991	1.7	4.5	0	0	0	2	306
is_mobile	31991	1.0	5.9	0	0	0	0	210
count_visits	31991	50	88	1	10	23	55	4130
is_ndop	31991	8.8	20.4	0	1	4	10	1583
is_catalog_all	31991	4.6	8.8	0	0	2	5	427
is_home	31991	10.9	30.3	0	1	3	10	2948
is_referrer_infomoney	31991	0.16	0.73	0	0	0	0	49
is_signin	31991	3.3	7.0	0	0	1	4	393

Exploratory Visualization

We drew some interesting conclusions from exploratory visualization, as seen in the charts below. We saw, as expected, that leads who became paying students ended up accessing over 2x more the checkout page vs. who didn't convert. We also saw, as expected, that leads who became paying students ended up watching close to 3x more webinars than who didn't convert.

We also observed, not so obviously, that the older an account is, less likely the lead is to convert. Also, the more a lead enrolls in a free course, the less likely it is for him/her to become a Nanodegree student.



Algorithms and Techniques

We used all key supervised learning key methods, not only base estimators (e.g. Stochastic Gradient Descent, Naive Bayes, Nearest Neighbors) but also ensemble methods (e.g. Gradient Tree Boosting), and compare all results - using scikit-learn implementation.

Also, we used 75% of the data for training and 25% of the data for testing.

Benchmark

Udacity US created a model to predict if a student would start a trial. We used that as benchmark, as it is a comparable approach used to the same purpose.

After extracting similar data from US data sources, the team from headquarters treated the data, divided into test/train set (% of elements go to the train set) and ran the training set over the following algorithms: various support vectors machines with different parameters and kernels, stochastic gradient descent, linear and logistic regression.

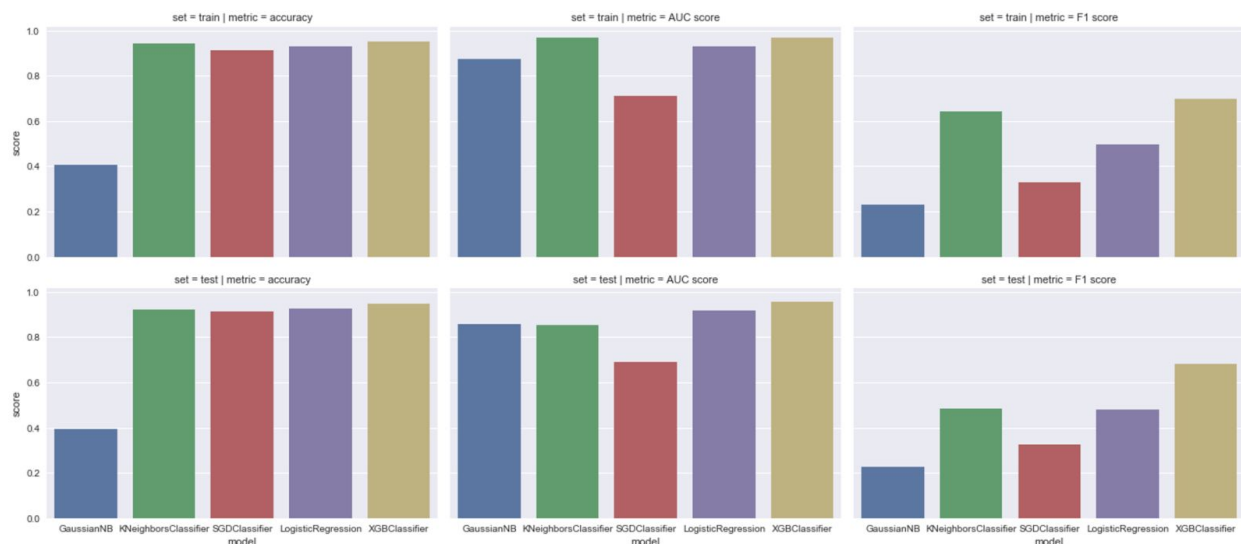
They compared the results against accuracy (over the test dataset), where SVM and logistic regression gave approximately 72% and other ~64%. The highest value was our benchmark.

Methodology

All preprocessing steps were documented in either `new_features.ipynb` or `capstone_final_project.ipynb`.

Implementation

After i) preparing the data for modeling, training and testing, and ii) initializing helper functions used for training and testing the supervised learning models, we runned the models and selected the best ones based on their metrics.



We compared **accuracy**, **AUC score** and **F1 score** for 5 models:

1. Gaussian Naive Bayes
2. K-Nearest Neighbors
3. Stochastic Gradient Descent
4. Logistic Regression
5. Extreme Gradient Boosting ([tutorial](#))

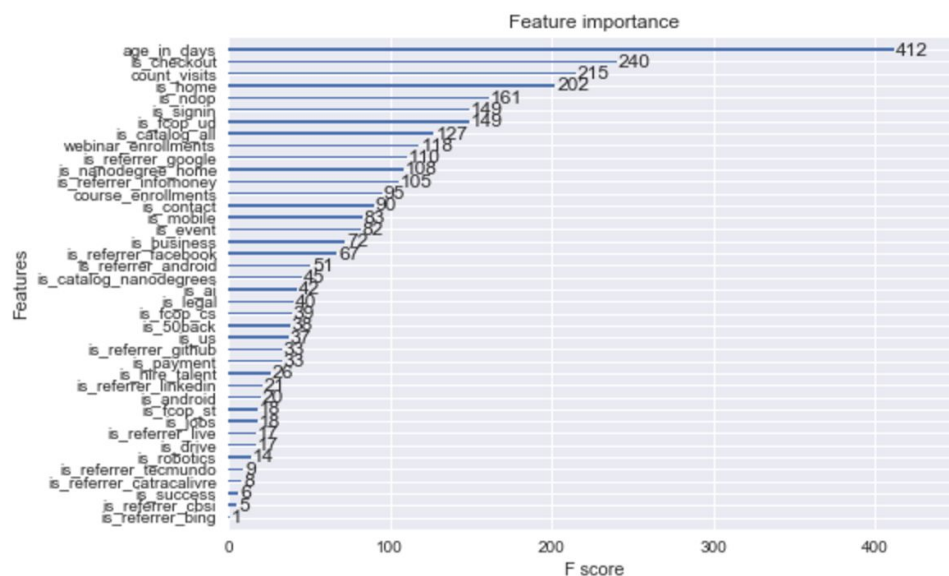
Extreme Gradient Boost won all metrics in both training and test datasets. Therefore, we have chosen this algorithm to move forward to parameter tuning/refinement.

Refinement

To do parameter tuning and model improvement, we used **grid search** (GridSearchCV) library. Also, we followed the step-by-step approach described in [Complete Guide to Parameter Tuning in XGBoost](#) article.

Before getting started, we implemented a helper function to enable us to see which were the best parameters and perform **cross-validation** (using the built-in function provided by XGBoost). Also, we **discarded accuracy, focusing on AUC score** only, as our **dataset was very imbalanced**.

AUC score mean of cross-validation sets, prior to tuning, was **0.9651**. The optimal number of trees was 502, and we were able to get the importance of each feature plotted in the picture below:



Then, we tuned the parameters in the following order:

1. **max_depth** and **min_child_weight**
2. **gamma**
3. **subsample** and **colsample_bytree**
4. **reg_alpha** and **reg_lambda**

As most of tuning was done varying pairs of parameters, we used several heat maps to visually see the best parameters. Below you can find a sample of the best charts. Please refer to **capstone_final_project.ipynb** for all step-by-step details.



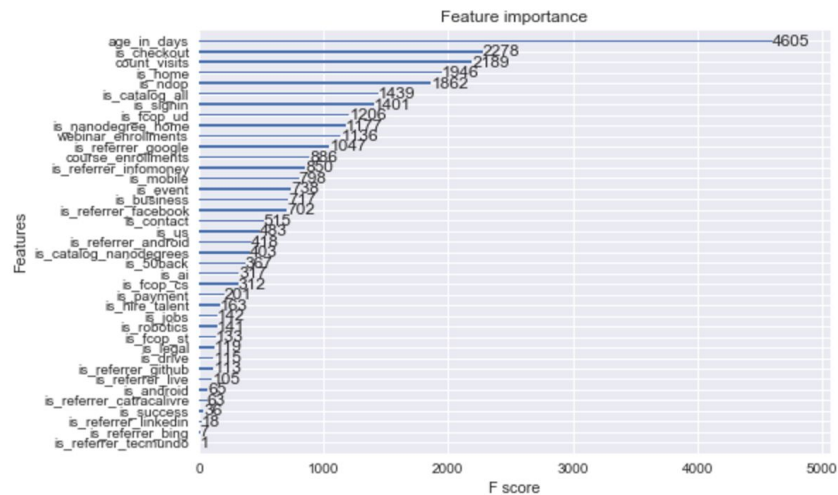
Results

Model Evaluation and Validation

After all parameter tuning, we were able to reach an **AUC score of 0.96691**. Lastly, we lowered the learning rate, applying all tuned parameters and added significantly more trees, we reached the final results:

1. AUC score in train dataset: **0.9798**
2. AUC score mean of cross-validation sets: **0.9664**
3. AUC score in test dataset: **0.9631**

Parameter tuning slightly improved ROC AUC score. The high scores in train dataset, CV datasets (mean) and test datasets indicates that we were able to generate a very strong prediction model, without overfitting. Also, we were able to read the most important features - the chart below summarizes the main ones:



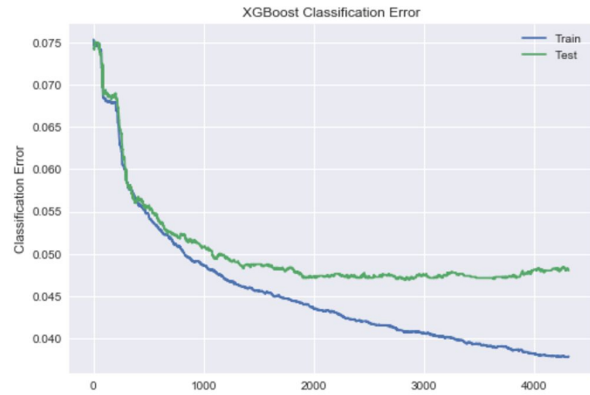
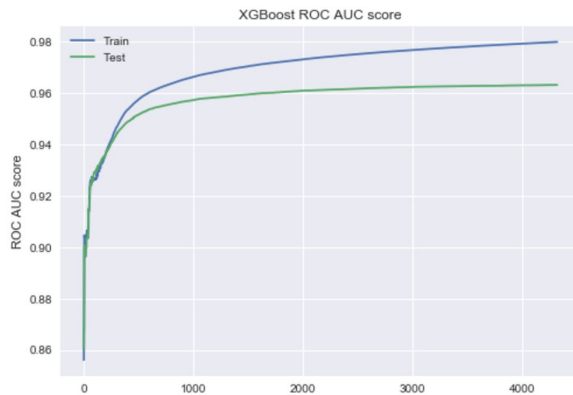
Justification

Accuracy of our model was significantly higher than the benchmark used, **0.95** vs. **0.71**. As our dataset was pretty skewed, with only 10% of samples with `is_paying_student == 1`, we re-did all the exercise undersampling the population of `is_paying_student == 0` to reach 1:1 ratio instead of 1:10. The result was still significantly better than the benchmark: we achieved **0.87** vs. **0.71**.

Conclusion

Free-Form Visualization

We retrieved the performance of our model on the evaluation dataset and plotted it to get insight into how learning happened during the training. As we can see in the Learning Curves below, the performance of our model stops improving around 4,000 iterations.



Reflection

In this capstone project, we leveraged everything we learned throughout the Nanodegree program to solve a problem of lead conversion prediction applying machine learning algorithms and techniques. We defined the problem: train an algorithm to successfully predict whether or not a Udacity website visitor becomes a Nanodegree student. We analyzed the problem, visualizing and exploring the data available. Then, we implemented different algorithms and metrics, choosing XGBoost. After this, we documented preprocessed, refined, and postprocessed. Afterwards, we collected the results about the performance of our model and validated/justified these values. Finally, we drew conclusions about our results.

During this whole process, we have chosen to highlight 3 aspects:

1. Gathering high quality data and making sure we build a strong dataset to be used later requires significant amount of work, and is a very important part of the process. Different than the other projects, when the dataset was provided, we had to create the capstone dataset, engineering every feature, one-by-one.
2. The end-to-end process was extremely iterative: we runned the end-to-end process dozens of times, tweaking/sampling the dataset, adding features, removing features, changing the metrics, etc. We only started to really better understand how the process really works after several iterations.
3. Finally, the parameter tuning phase required a lot of patience, as searching the parameters grid after optimal values took a lot of time. However, it paid out: we achieved a high performance model, which can (and will) be tested at Udacity Brazil.

Improvement

We believe the model can be significantly improved by gathering more data in 2 fronts. One would be web scraping visitors' cookies and trying to retrieve LinkedIn data. The other would be collecting and using time spent in each page and each visit.