



X EDUCATION - LEAD SCORING CASE STUDY

IDENTIFICATION OF HOT LEADS TO FOCUS MORE ON
THEM AND THUS ENHANCING THE CONVERSION
RATIO FOR X EDUCATION

Nguyen Thi Huong THAO (Ms.)
Master of Data Science Course – 5702
Nguyen.huong.thao87@gmail.com

BACKGROUND

X Education Company

- X Education , an education company named sells online courses to industry professionals
- Many interested professionals land on their website
- The company markets its courses on several websites like Google. Once these people land on the website, they might browse the courses or fill up a form for the course or watch some videos

BACKGROUND

X Education Company

- When these people fill up a form providing their email address or phone number, they are classified to be a lead
- Once these leads are acquired, employees from the sales team start making calls, writing emails, etc. Through this process, some of the leads get converted while most do not
- The typical lead conversion rate at X education is around 30%

PROBLEM STATEMENT

X Education Company's Problem

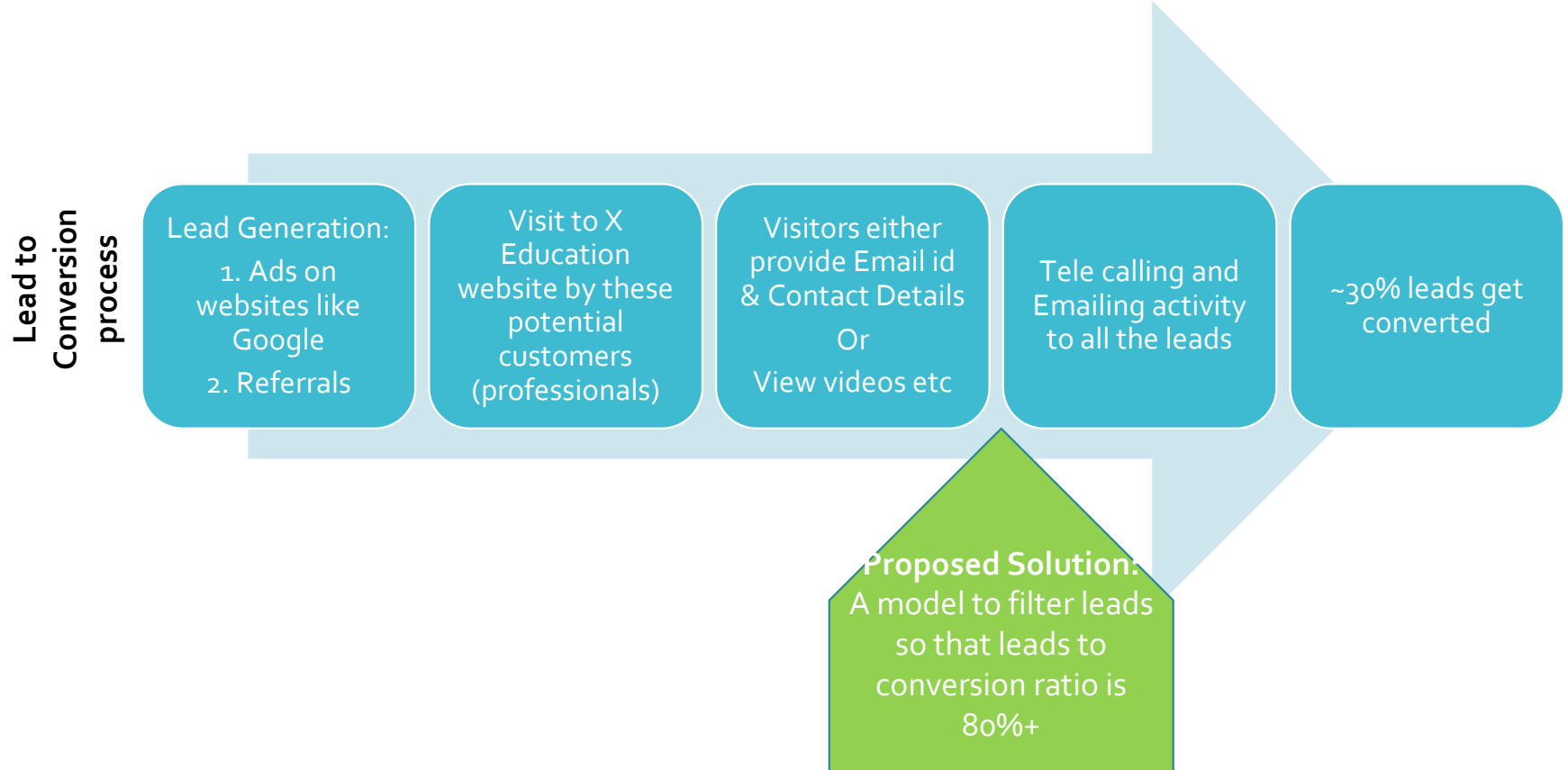
- X Education gets a lot of leads but its lead conversion rate is very poor
- To make this process more efficient, the company wishes to identify the most potential leads, also known as 'Hot Leads'
- If they successfully identify this set of leads, the lead conversion rate should go up as the sales team will now be focusing more on communicating with the potential leads rather than making calls to everyone

PROBLEM STATEMENT

X Education Company's Problem

- We will help them to select the most promising leads, i.e. the leads that are most likely to convert into paying customers.
- We are required to build a model wherein we need to assign a lead score to each of the leads such that the customers with higher lead score have a higher conversion chance
- The CEO, in particular, has given a ballpark of the target lead conversion rate to be 80%.

Lead – Conversion Process



PROPOSED SOLUTION

Selection of Hot Leads

Leads Clustering

We cluster the leads into certain categories based on their tendency or probability to convert, thus, getting a smaller section of hot leads to focus more on.

Communicating with Hot Leads

Focus Communication

Since we would have a smaller set of leads to have communication with, we might make more impact with effective communication.

Conversion of Hot Leads

Increase conversion

Since we focussed on hot leads, which were more probable to convert, we would have a better conversion rate, and hence we can achieve the 80% target.

SOLUTION

Selection of Hot Leads

For our Problem Solution, the crucial part is to accurately identify hot leads.

The more accurate we obtain the hot lead, the more chance we get of higher conversion ratio.

Since we have a target of 80% conversion rate, we would want to obtain a high accuracy in obtaining hot leads.

The background is a dark gray gradient. In the corners, there are white, stylized circuit-like lines. These lines consist of straight segments and small circles, resembling a network or a circuit board layout. They are positioned in the top-left, top-right, bottom-left, and bottom-right corners, framing the central text.

IMPLEMENTATION

Loading & Observing the
past data provided by
the Company

Univariate, Bivariate, and
Heatmap for numerical and
categorical columns

Performing pre-requisites
for RFE and Logistic
Regression

Data
Gathering

Performing
EDA

Data
Cleaning

Data
Preparation

Model
Building

Duplicate removal, null value
treatment, unnecessary column
elimination, etc.

Outlier Treatment,
Feature-Standardization

Selection of top 25
features using RFE

Reduction of columns and
Model re-building

Verifying our Final Model
Accuracy etc. with model
built with PCA

Feature
Selection

Model
Building

Model
Improvement

Final Model

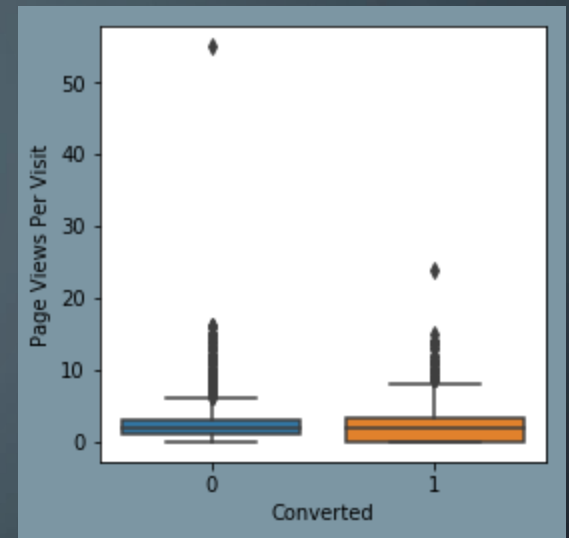
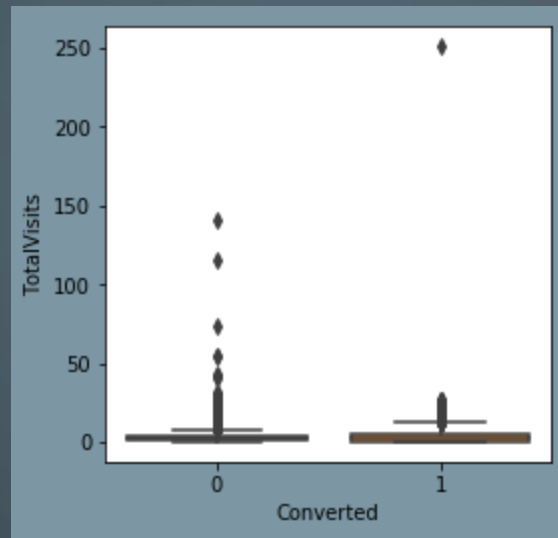
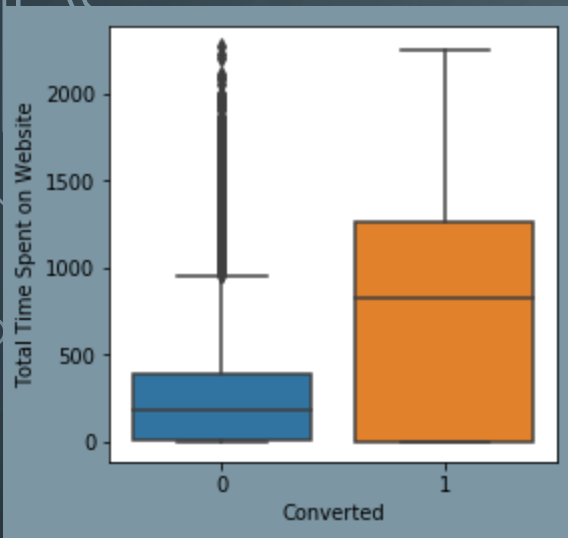
Verifying with
PCA

Model building using RFE
for selected columns

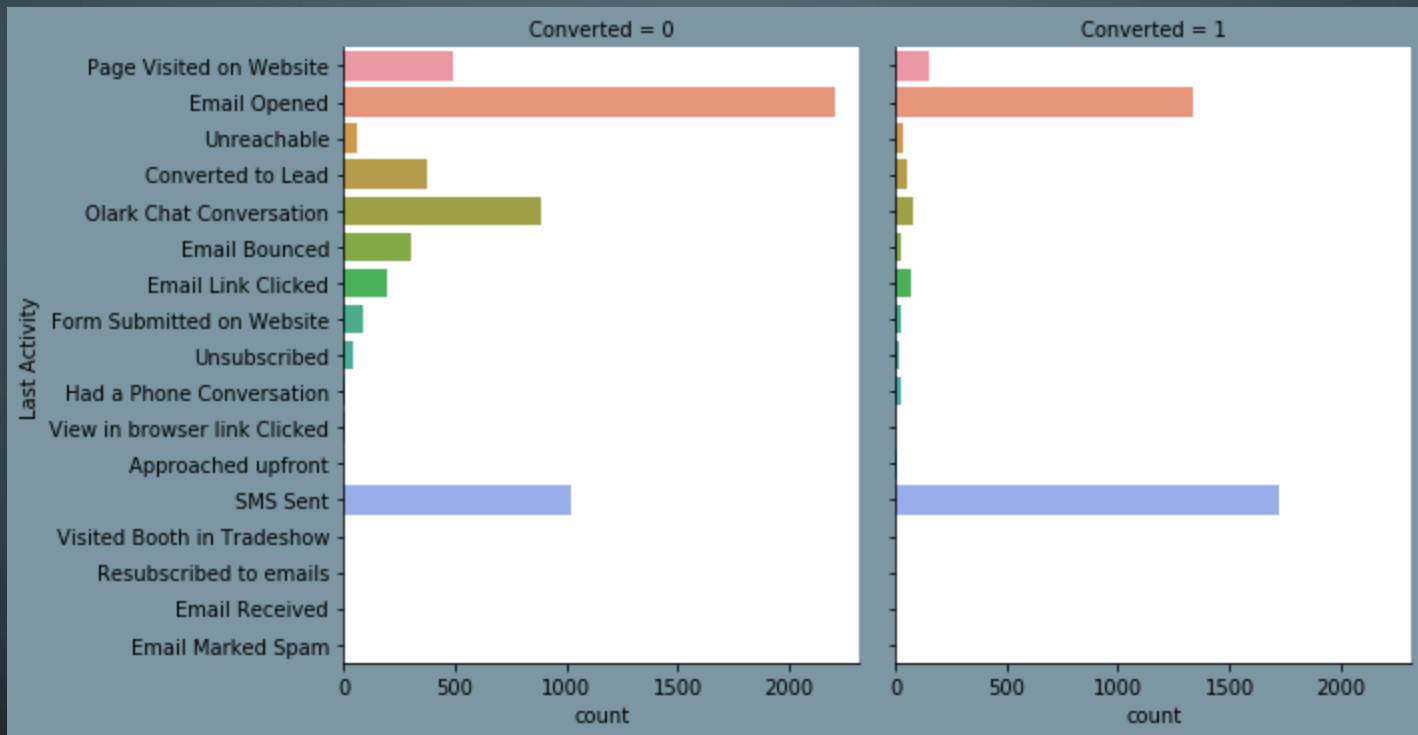
Final Model Analysis and
performance on Test Data

The image features a dark gray background with a subtle radial gradient. In the four corners, there are decorative white line art elements resembling circuit board traces or neural network connections, each ending in a small circle.

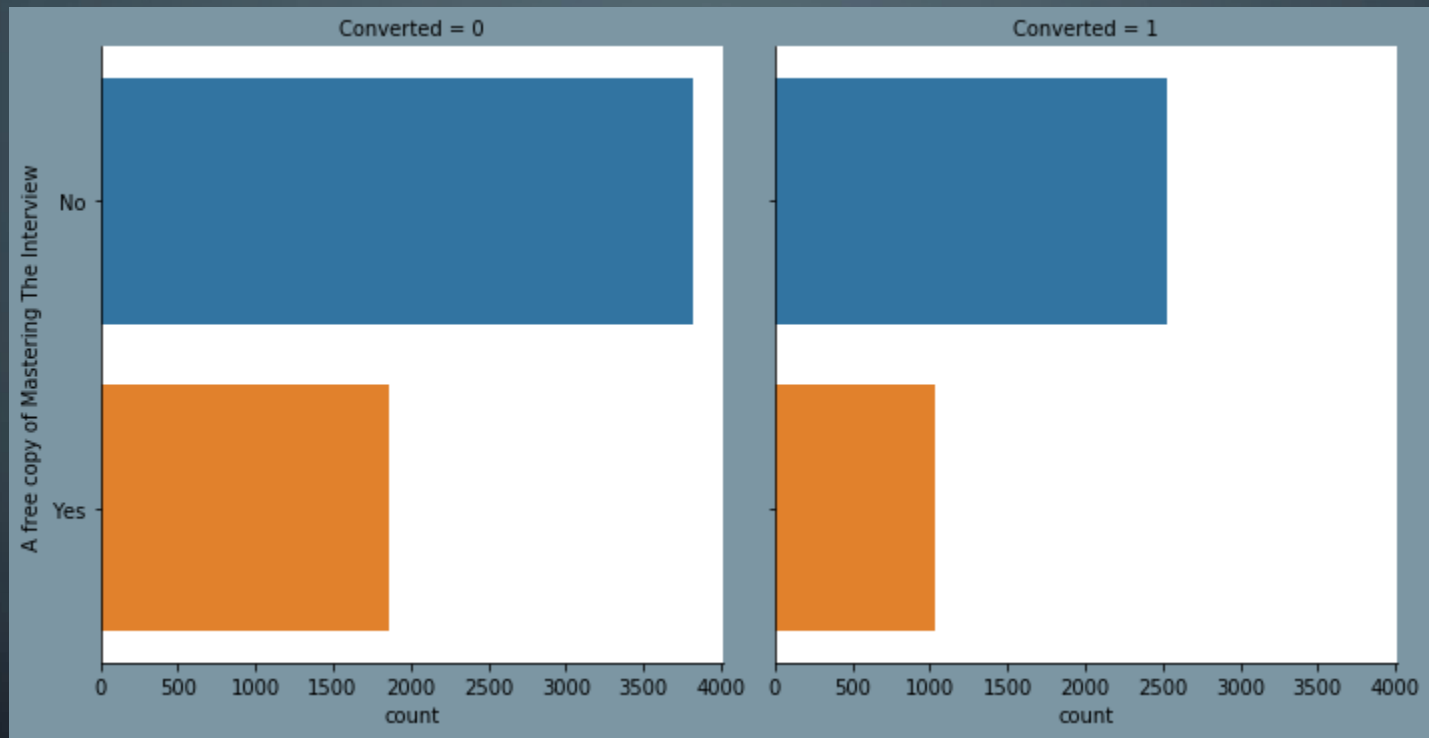
PLOTS (VISUALIZATION)



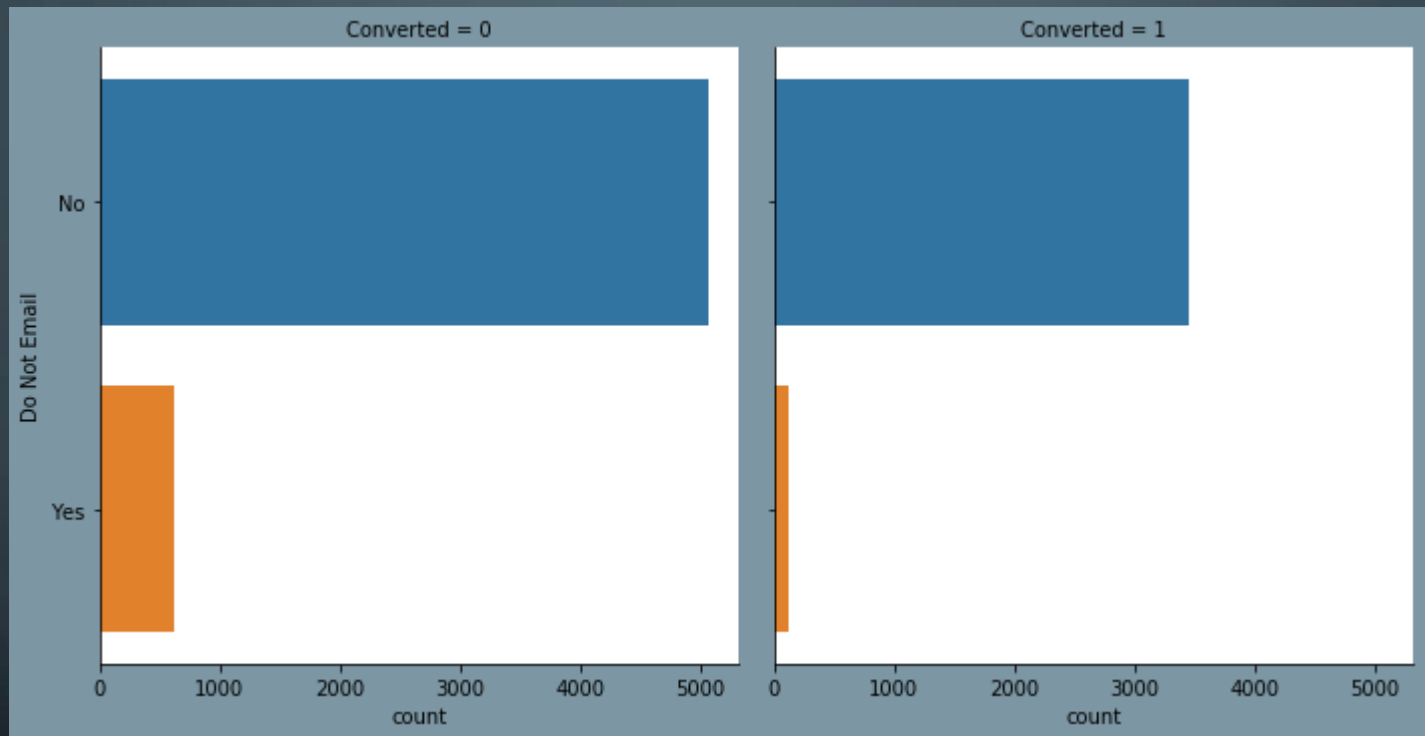
EDA plots depicting variation in numerical columns for those who Converted and those who didn't.



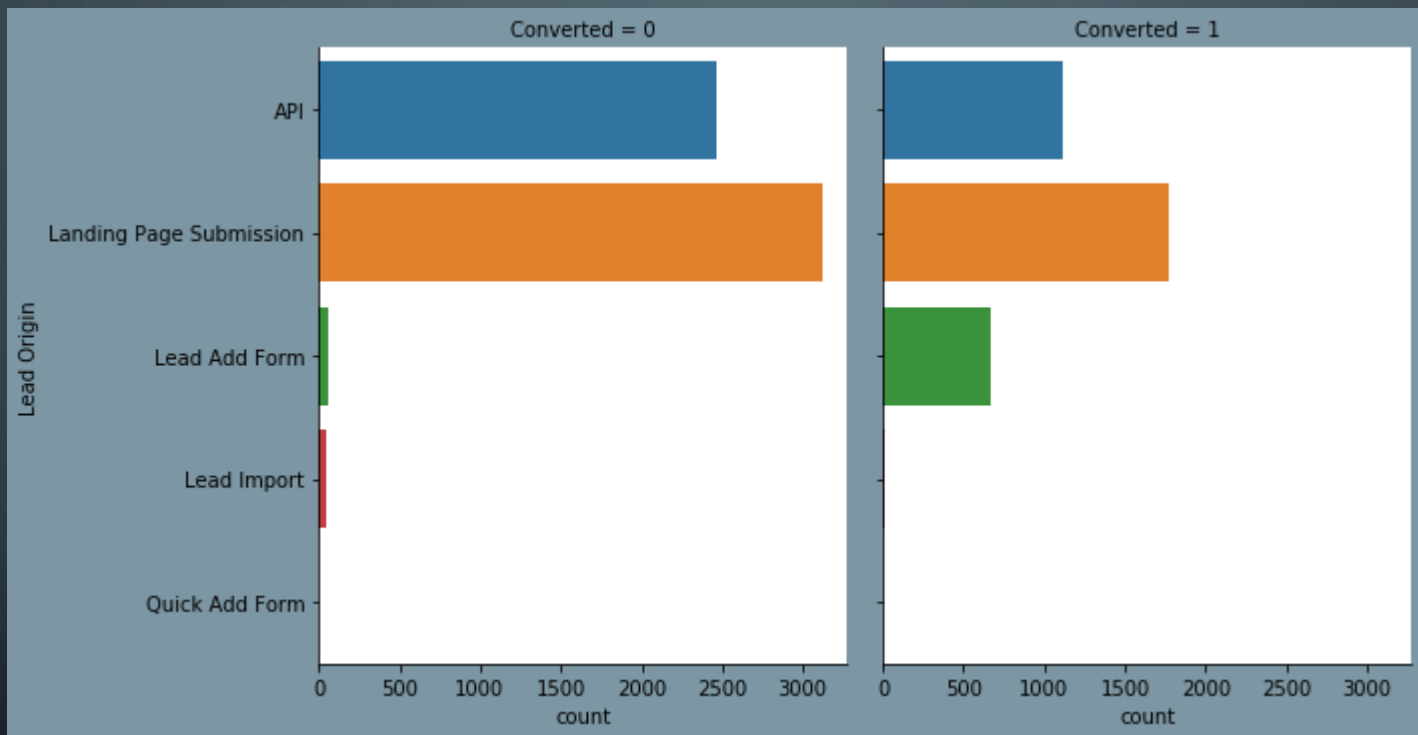
EDA plots depicting variation in categorical column (Last Activity) for those who Converted and those who didn't.



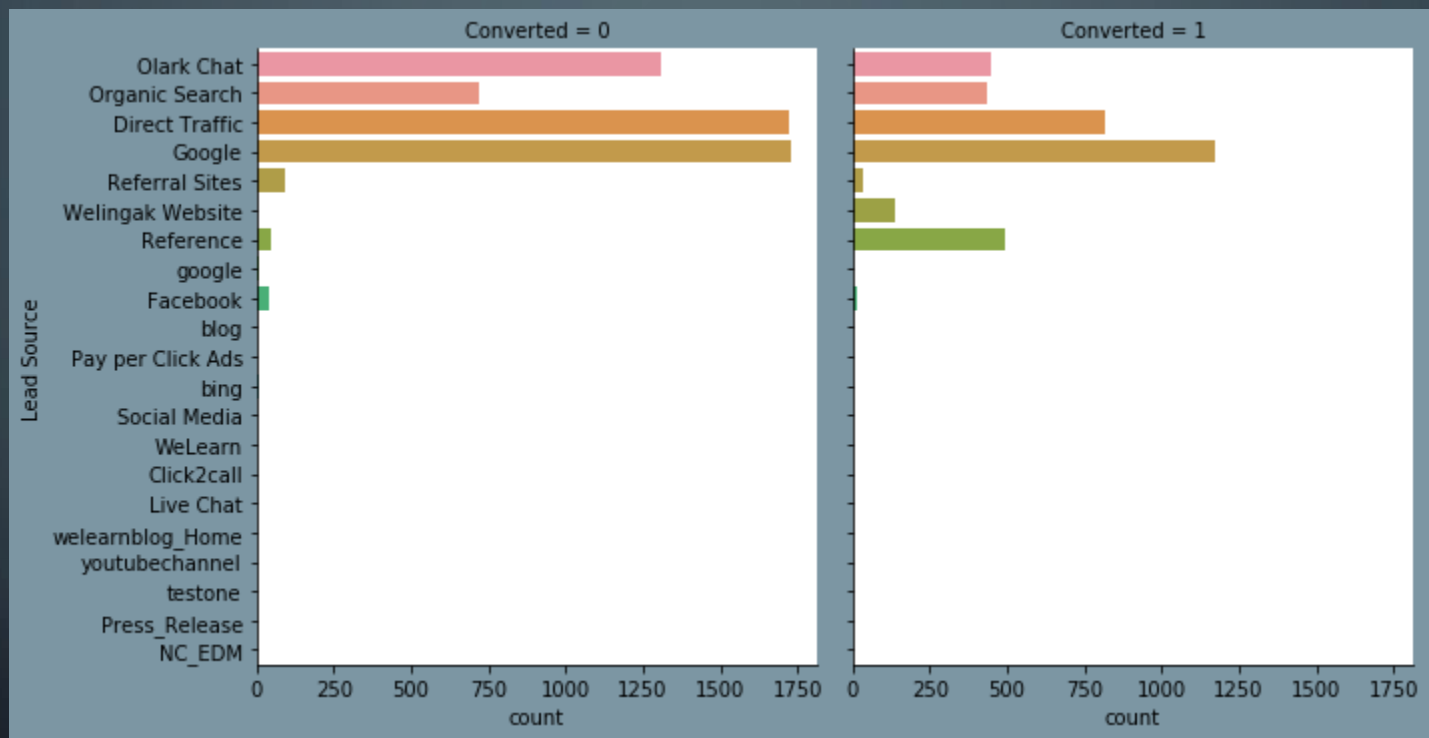
EDA plots depicting variation in categorical column (A free copy of Mastering The Interview) for those who Converted and those who didn't.



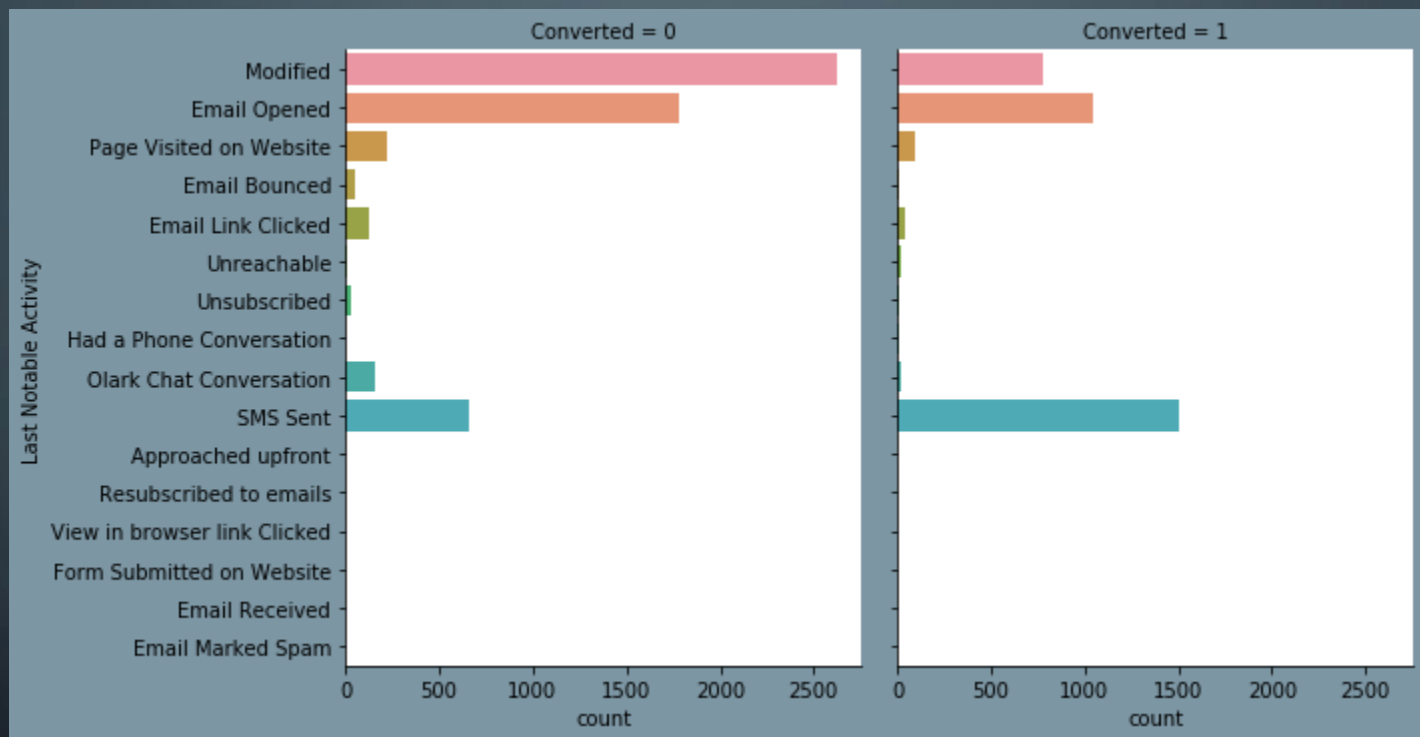
EDA plots depicting variation in categorical column (Do Not Email) for those who Converted and those who didn't.



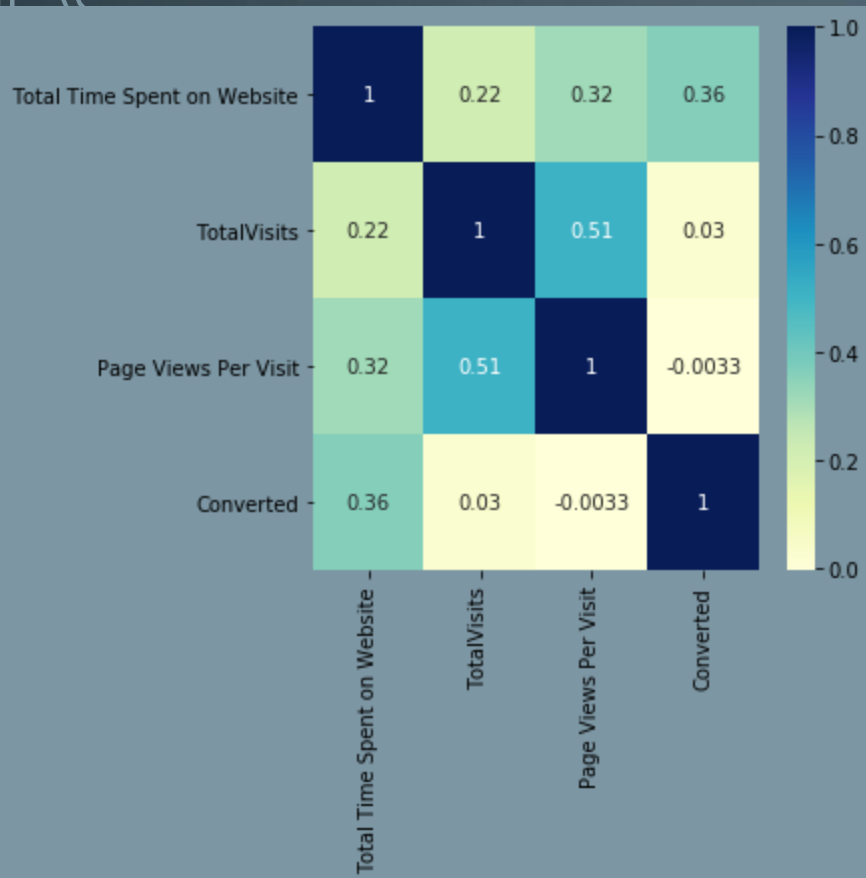
EDA plots depicting variation in categorical column (Lead Origin) for those who Converted and those who didn't.



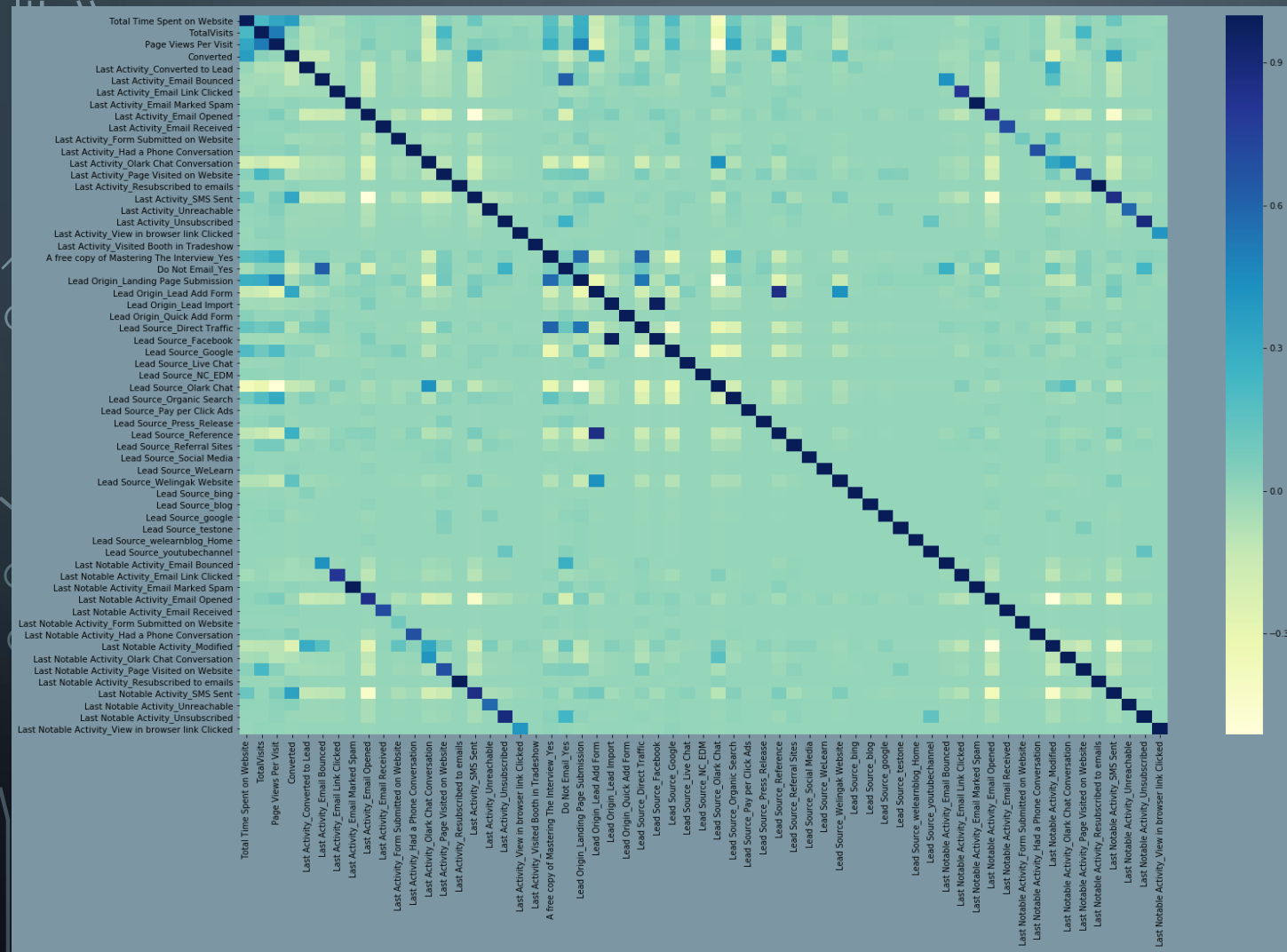
EDA plots depicting variation in categorical column (Lead Source) for those who Converted and those who didn't.



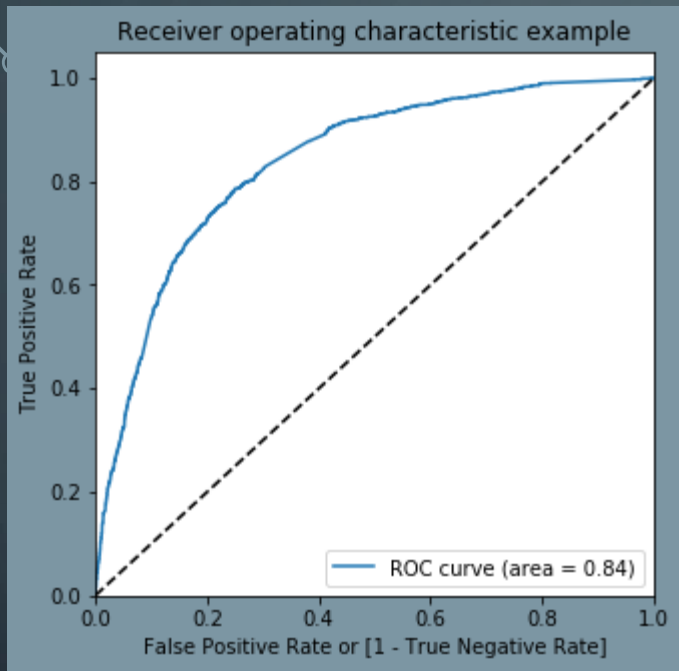
EDA plots depicting variation in categorical column (Last Notable Activity) for those who Converted and those who didn't.



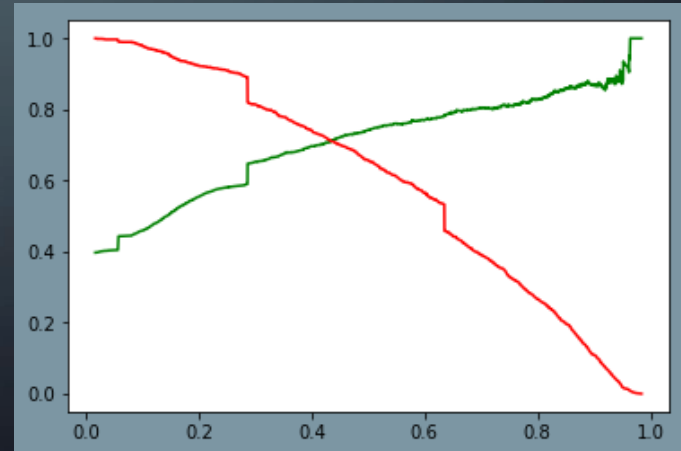
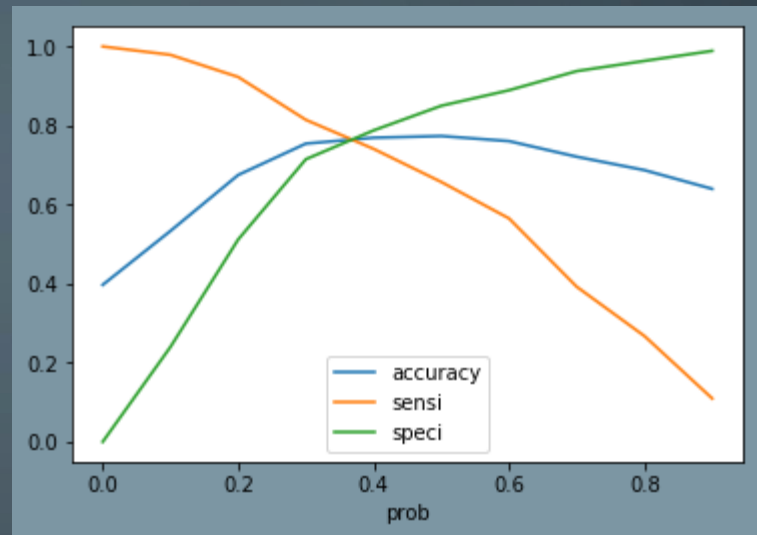
EDA plots depicting correlation (Heat Map) of all selected numerical columns.



EDA plots depicting correlation (Heat Map) of all selected columns (numerical columns and dummy columns).



**Linear Regression Final
Model Parameters
Area under ROC = 0.84
Intermediate cut-off = 0.35
Final cut-off = 0.42**





EDA plots depicting correlation (Heat Map) of all selected columns (numerical columns and dummy columns) in our final Model.

The background is a dark gray gradient. In the four corners, there are white, stylized circuit-like lines. These lines consist of straight segments and small circles, resembling a network or a circuit board layout. The lines are more dense in the bottom-left and top-left corners and more sparse in the top-right and bottom-right corners.

INFERENCE / CONCLUSION

MODEL ANALYSIS

Performance of our Final Model

Overall accuracy on Test set: 0.786

Sensitivity of our logistic regression
model: 0.733

Specificity of our logistic regression
model: 0.823

INFERENCES FROM MODEL

Business Insights Derived from our
Model

Top 3 variables in model, that
contribute towards lead conversion are:

- Total Time Spent on Website
- Last Notable Activity_SMS Sent
- TotalVisits

INFERENCES FROM MODEL

Business Insights Derived from our
Model

Top 3 variables in my model, that
should be focused are:

- Last Activity_SMS Sent (positively impacting)
- Last Activity_Olark Chat Conversation (negatively impacting)
- Lead Source_Olark Chat (negatively impacting)

CONCLUSION 1 (LR MODEL)

OUR LOGISTIC REGRESSION MODEL IS DECENT AND ACCURATE ENOUGH, WHEN COMPARED TO THE MODEL DERIVED USING PCA, WITH 78.6 % ACCURACY ON TEST SET, 73.3 % SENSITIVITY AND 82.3 % SPECIFICITY.

WE CAN VARY THESE PARAMETERS BY VARYING THE CUT-OFF VALUE AND THUS PREDICT HOT LEADS BASED ON SCENARIOS LIKE AVAILABILITY OF EXTRA RESOURCES AND VICE-VERSA.

CONCLUSION 2 (RECOMMENDATION)

X EDUCATION COMPANY NEEDS TO FOCUS ON FOLLOWING KEY ASPECTS TO IMPROVE THE OVERALL CONVERSION RATE:

- INCREASE USER ENGAGEMENT ON THEIR WEBSITE SINCE THIS HELPS IN HIGHER CONVERSION
- INCREASE ON SENDING SMS NOTIFICATIONS SINCE THIS HELPS IN HIGHER CONVERSION
- GET TOTAL VISITS INCREASED BY ADVERTISING ETC. SINCE THIS HELPS IN HIGHER CONVERSION
- IMPROVE THE OLARK CHAT SERVICE SINCE THIS IS AFFECTING THE CONVERSION NEGATIVELY



THANK YOU!

Nguyen.huong.thao87@gmail.com