# Case Study: Logistic Regression

Problem Statement:

An education company named X Education sells online courses to industry professionals. On any given day, many professionals who are interested in the courses land on their website and browse for courses.
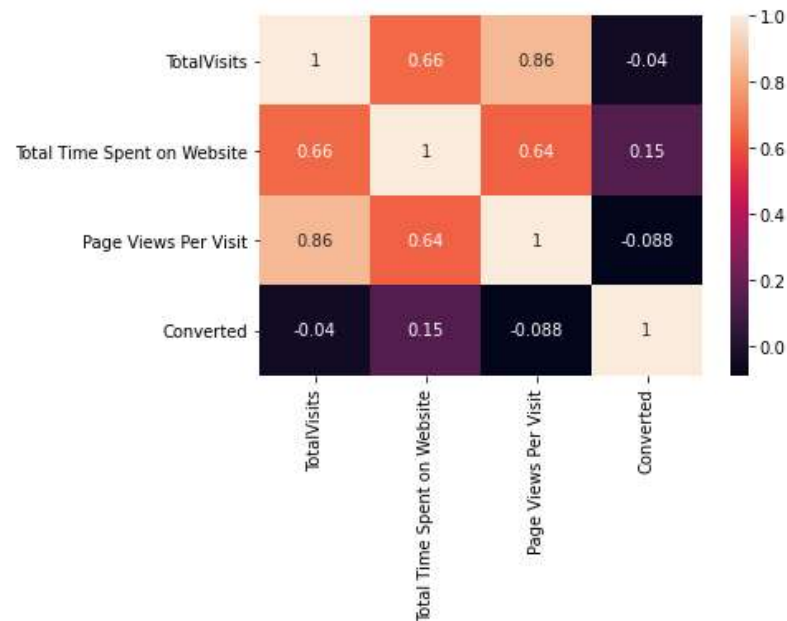
So based on various factors which are leads the company wants to predict the probability of leads converting into customers

# Data Cleaning:

- We checked for null values in columns and removed the columns which had more than 3000 null values.

- We removed the columns which were not usefull for our analysis

- Then we checked columns with very low variances and dropped them

- There were some categorical variables which had attribute called "select" which means that student has not selected for any option for that column,so this was as good as missing values and we dropped the columns which have more "select" values

# Transforming the data

- We analysed numerical colums by plotting paiplot and found that their result was bit skewed,so we transformed and we got better correlation and plots were more better
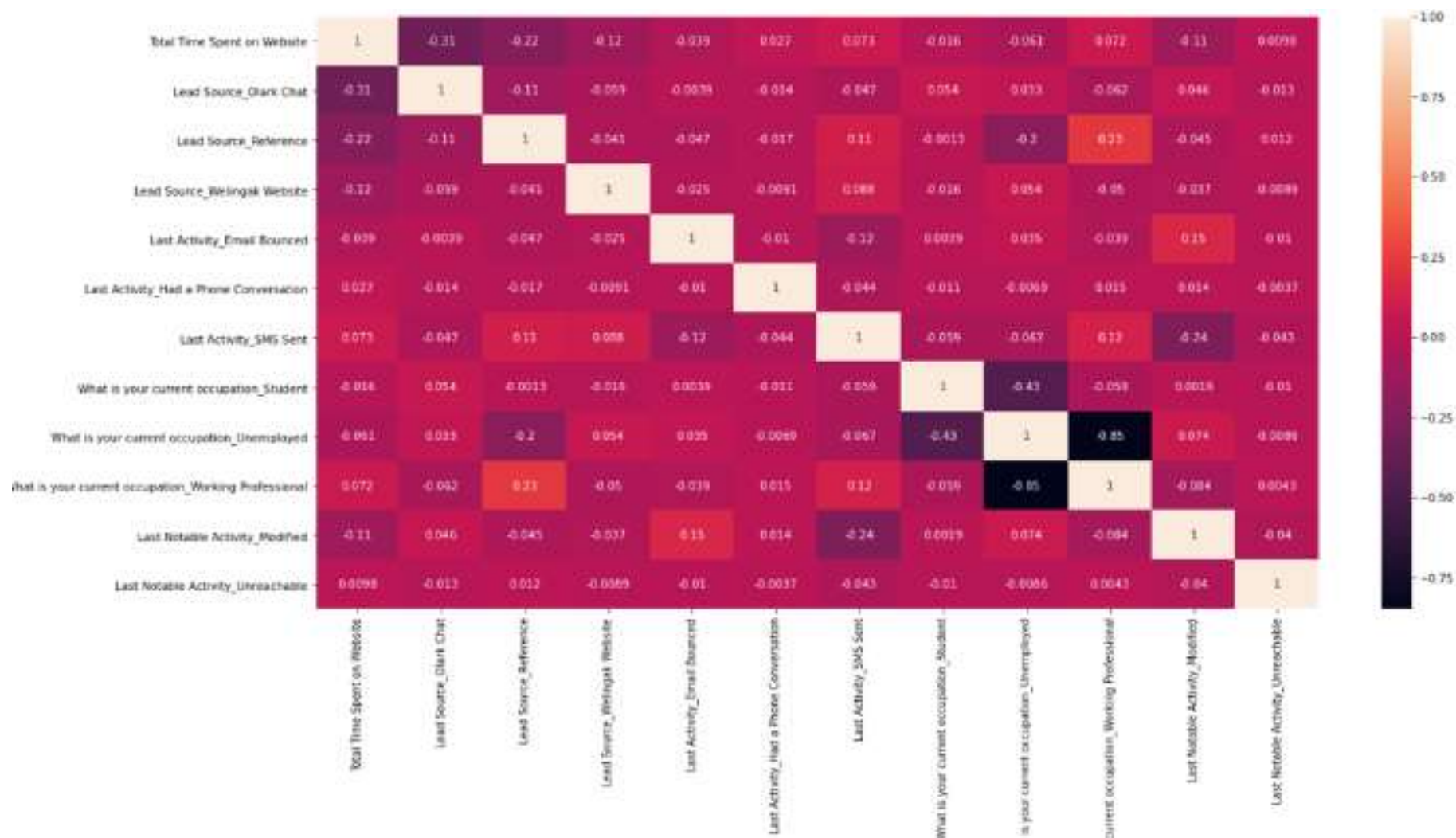
# Dummy variable creation

- We created dummies for categorical variables
- Created dummies for specialization column separately because it had some attributes of value "select".We dropped the dummies column which had "select"

# Test Train Split

- We split the data into training and testing for our model
- Put all features except target to variable x and target to y
- We split the data into 70% training and 30% test data
- We then scaled the data for better results.We used standardscaler.

# Model Building¶

- We used Recursive feature elimination to select features for our model.We selected 15 models for this

- We used confusion matrix to analyse our model.Our model was found to be 78.8% accurate

- We then used  variance inflation factor (VIF) to check multicollinearity and the result we got after few iteration was found to be 78.66

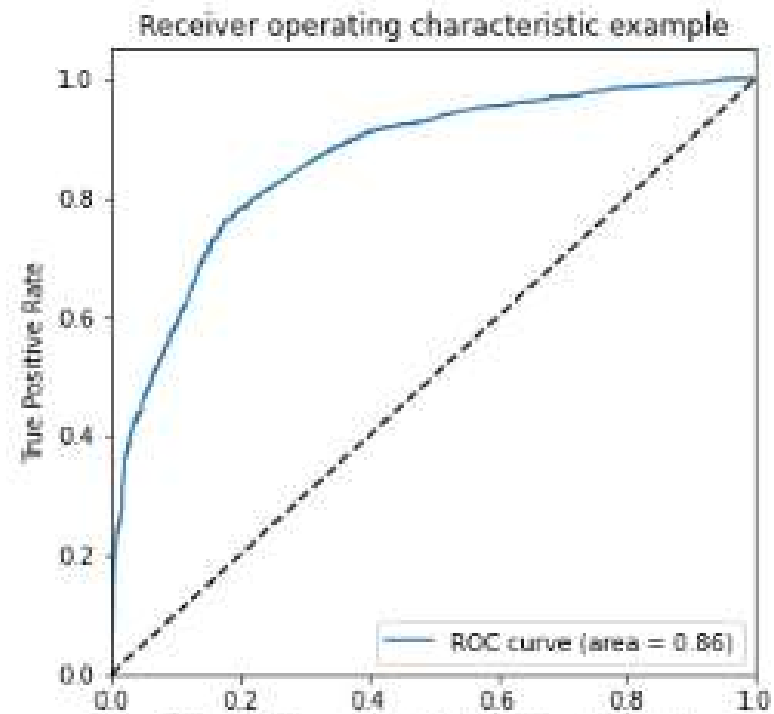- We checked validated the model and collinearity using heatmap

Correlation heatmap of lead scoring features including Total Time Spent on Website, Lead Source_Olark Chat, Lead Source_Reference, Lead Source_Welingak Website, Last Activity_Email Bounced, Last Activity_Had a Phone Conversation, Last Activity_SMS Sent, What is your current occupation_Student, What is your current occupation_Unemployed, What is your current occupation_Working Professional, Last Notable Activity_Modified, and Last Notable Activity_Unreachable.

# Metrics beyond simply accuracy¶

- We checked various other factors to measure our model

- Sensitivity : TP / (TP+FN) = 73.5 %

- Specificity : TN /(TN+FP) = 83.4 %

- False positive rate : FP/ (TN+FP)=16 %
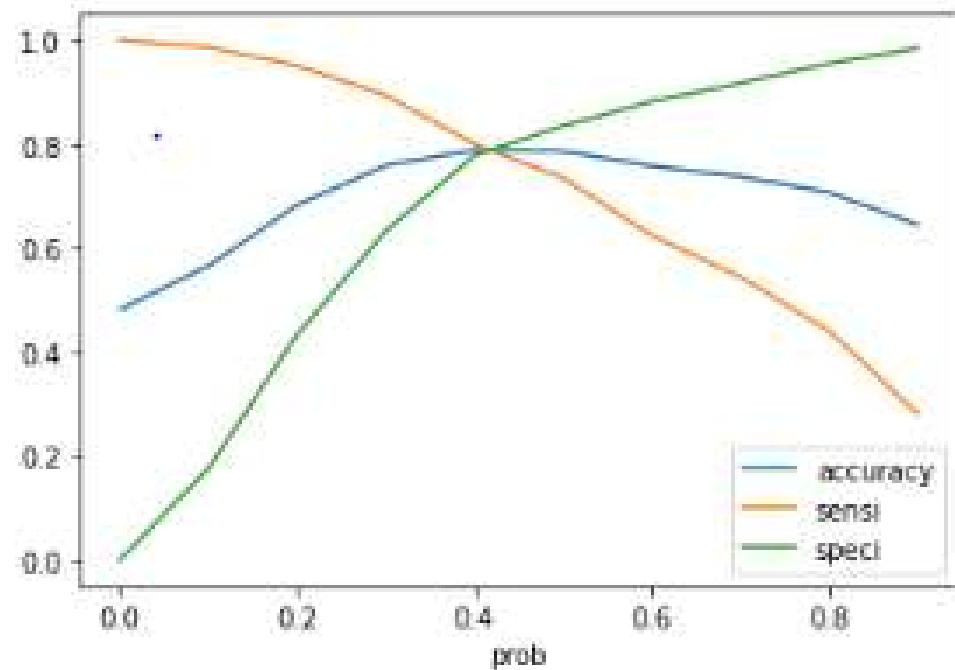
- Positive predictive value: TP / (TP+FP) = 80%

# ROC Curve

- An ROC curve demonstrates several things: •It shows the tradeoff between sensitivity and specificity (any increase in sensitivity will be accompanied by a decrease in specificity). •The closer the curve follows the left-hand border and then the top border of the ROC space, the more accurate the test. •The closer the curve comes to the 45-degree diagonal of the ROC space, the less accurate the test.



Receiver operating characteristic example

ROC curve (area = 0.86)

# Finding Optimal Cutoff Point

- Optimum cutoff was fount to be around .42
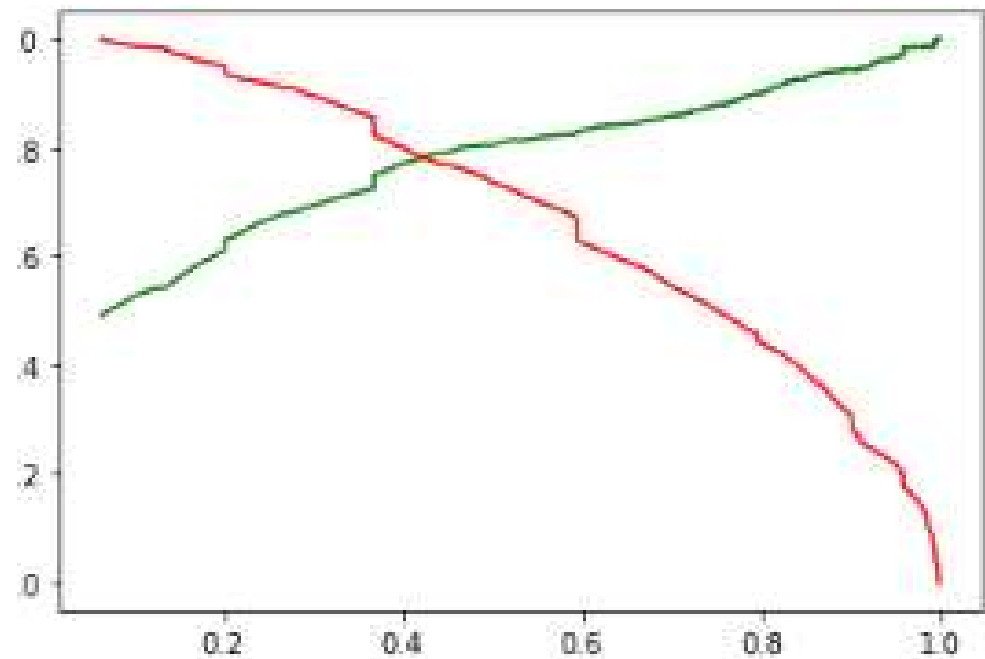
- overall accuracy:78.9%

# Precision and Recall

- Precision = 80.4
- Recall = 73.5

# Precision and Recall Tradeoff

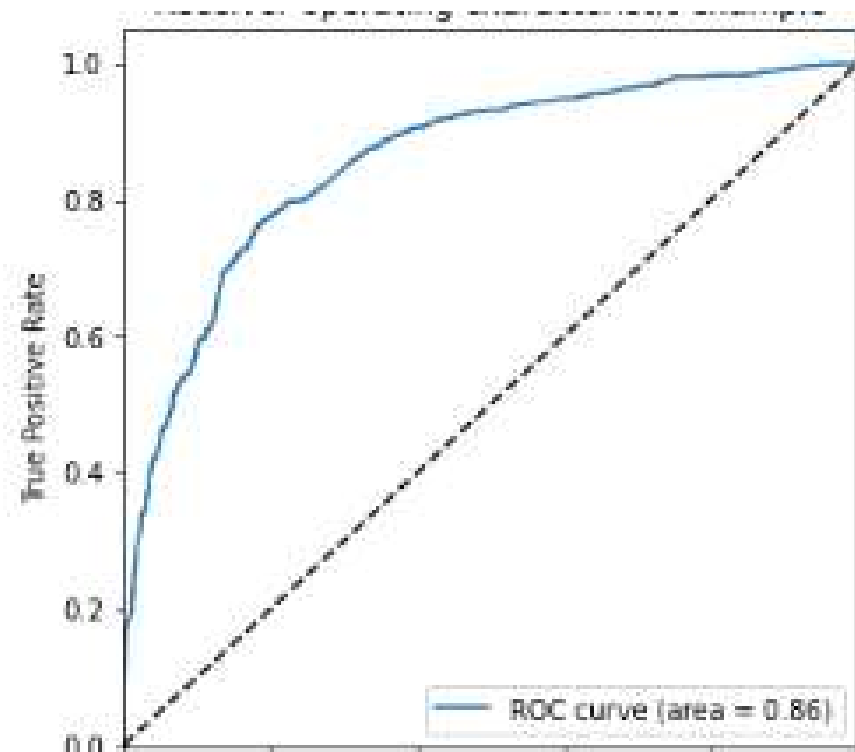- Cutoff between precision is around .42 which is same

# Making predictions on the test set

- Precision : 77.48
- Recall : 78.16
- Sesitivity: 78.16
- Specificity : 79.11
- False positive rate: 77.48

# ROC curve on test set

- Its similar to train set ROC curve

# Conclusion

1. Our Logistic Regression Model is decent and accurate enough, when compared to the model derived using PCA.

2. X Education Company needs to focus on following key aspects to improve the overall conversion rate:

a. Increase user engagement on their website since this helps in higher conversion

b. Increase on sending SMS notifications since this helps in higher conversion

c. Get TotalVisits increased by advertising etc. since this helps in higher conversion

d. Improve the Olark Chat service since this is affecting the conversion negatively