# Lead Score Case Study

Submitted by :
Harshit Kumar
Gautami Pravin Wankhede

# Problem Statement

X Education sells online courses to industry professionals. The company markets its courses on several websites and search engines 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. When these people fill up a form providing their email address or phone number, they are classified to be a lead. Moreover, the company also gets leads through past referrals. 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%.
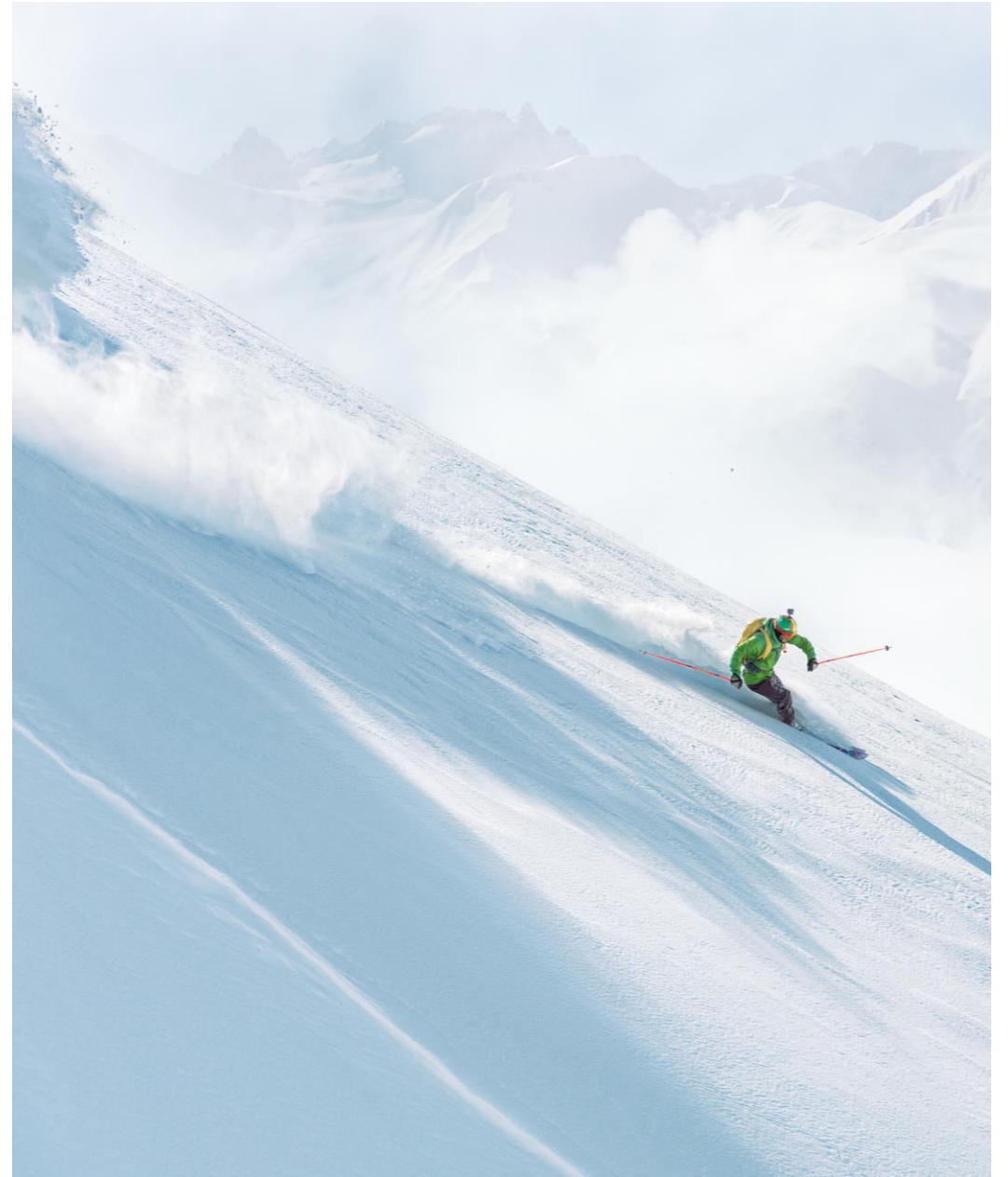
## Business Goal

X Education needs help in selecting the most promising leads, i.e. the leads that are most likely to convert into paying customers. The company needs a model wherein you a lead score is assigned to each of the leads such that the customers with higher lead score have a higher conversion chance and the customers with lower lead score have a lower conversion chance. The CEO, in particular, has given a ballpark of the target lead conversion rate to be around 80%.
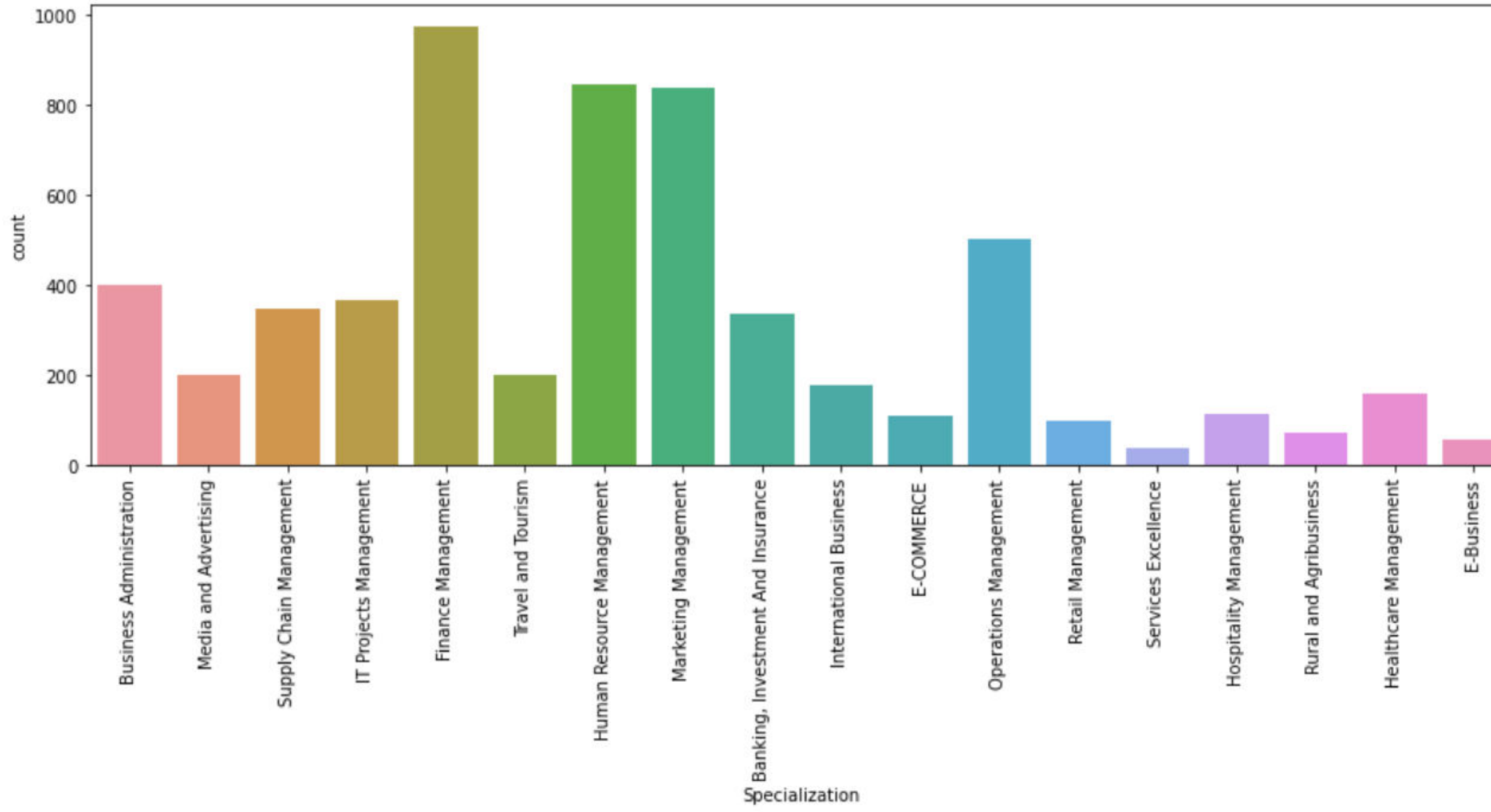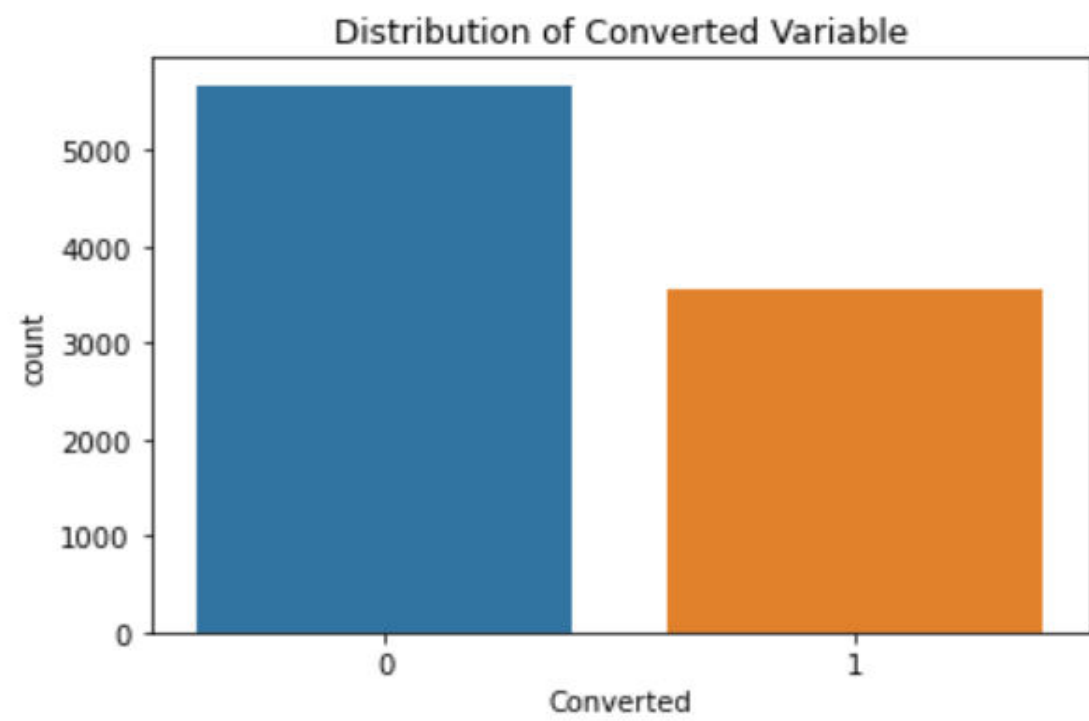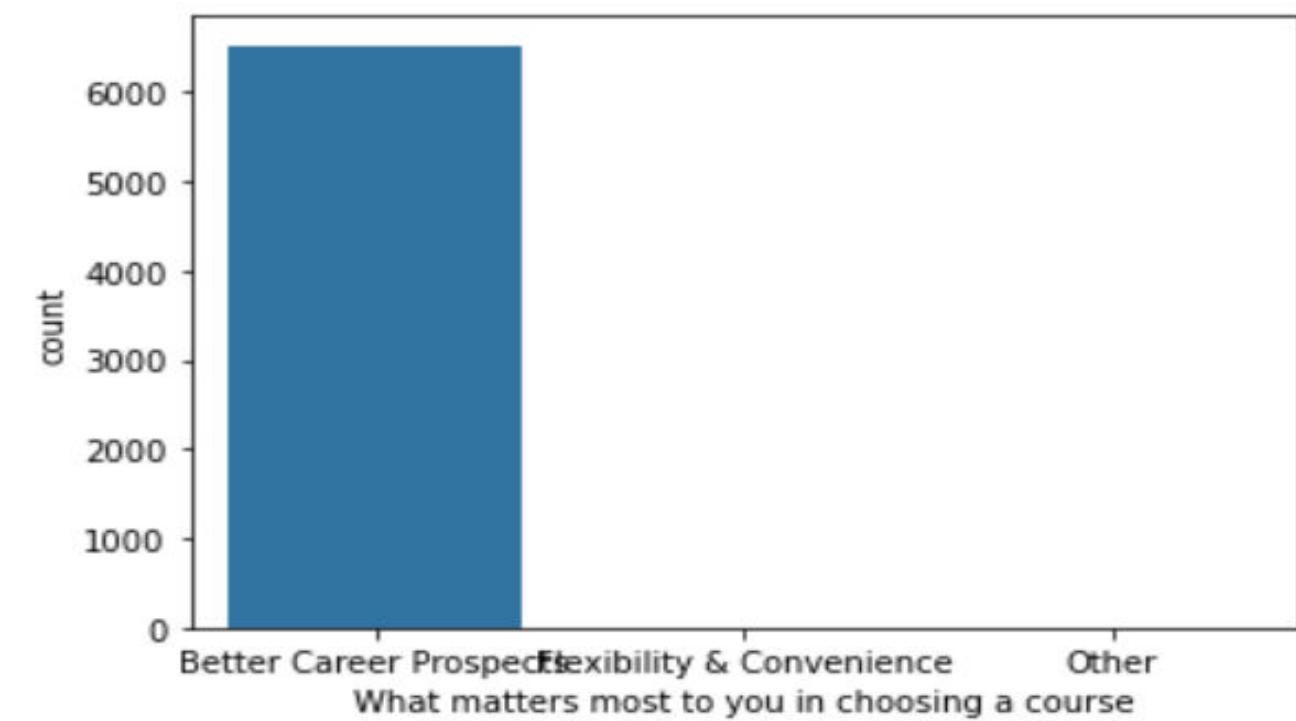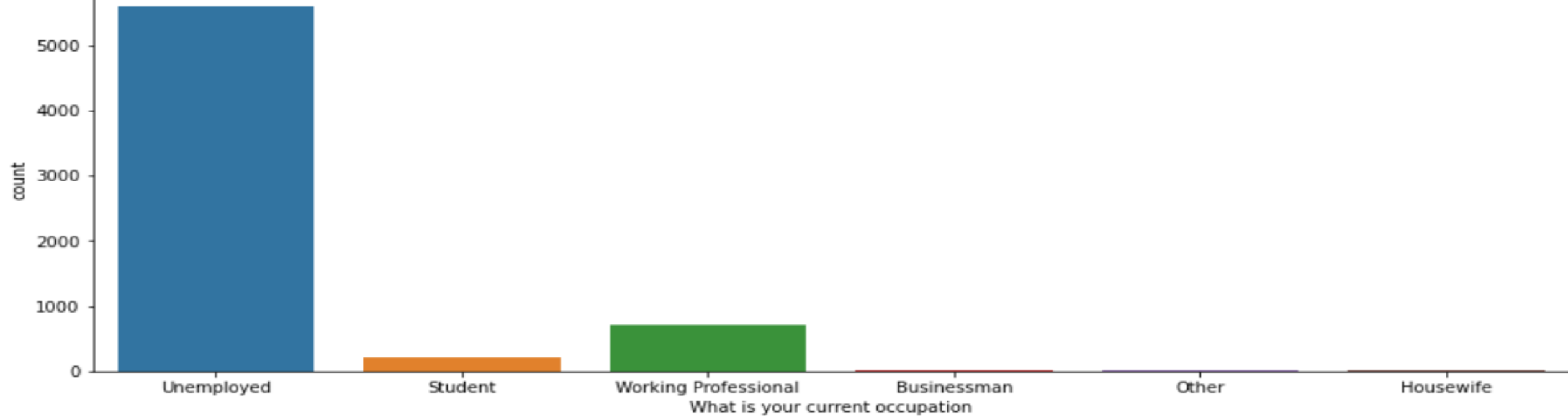
# Strategy

- Source the data for analysis.

- Clean and prepare the data.

- Exploratory Data Analysis.

- Feature Scaling.

- Splitting the data into Test and Train dataset.

- Building a logistic Regression model and calculate Lead Score.

- Evaluating the model by using different metrics - Specificity and Sensitivity or Precision and Recall.

- Applying the best model in Test data based on the Sensitivity and Specificity Metrics. And Validation of the model.

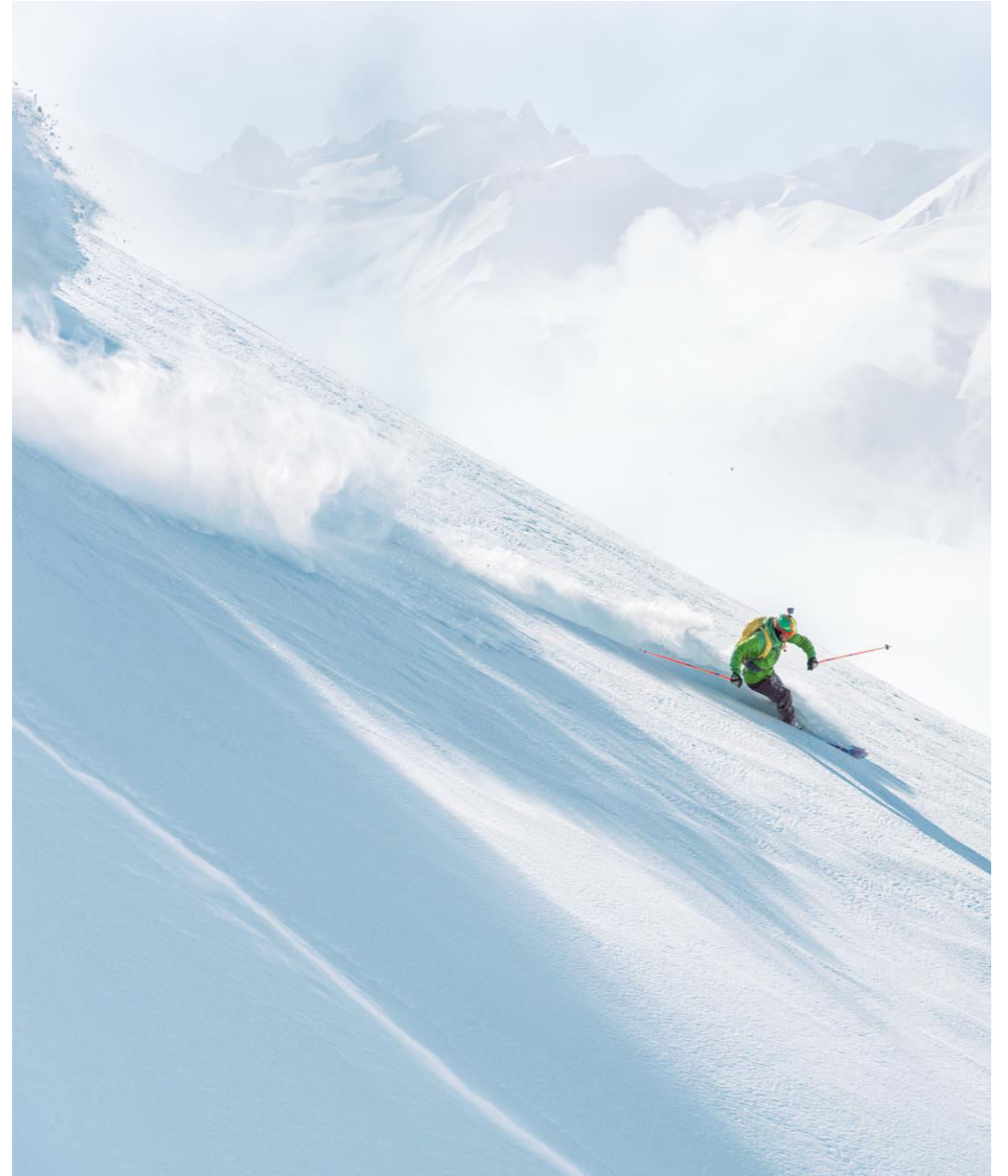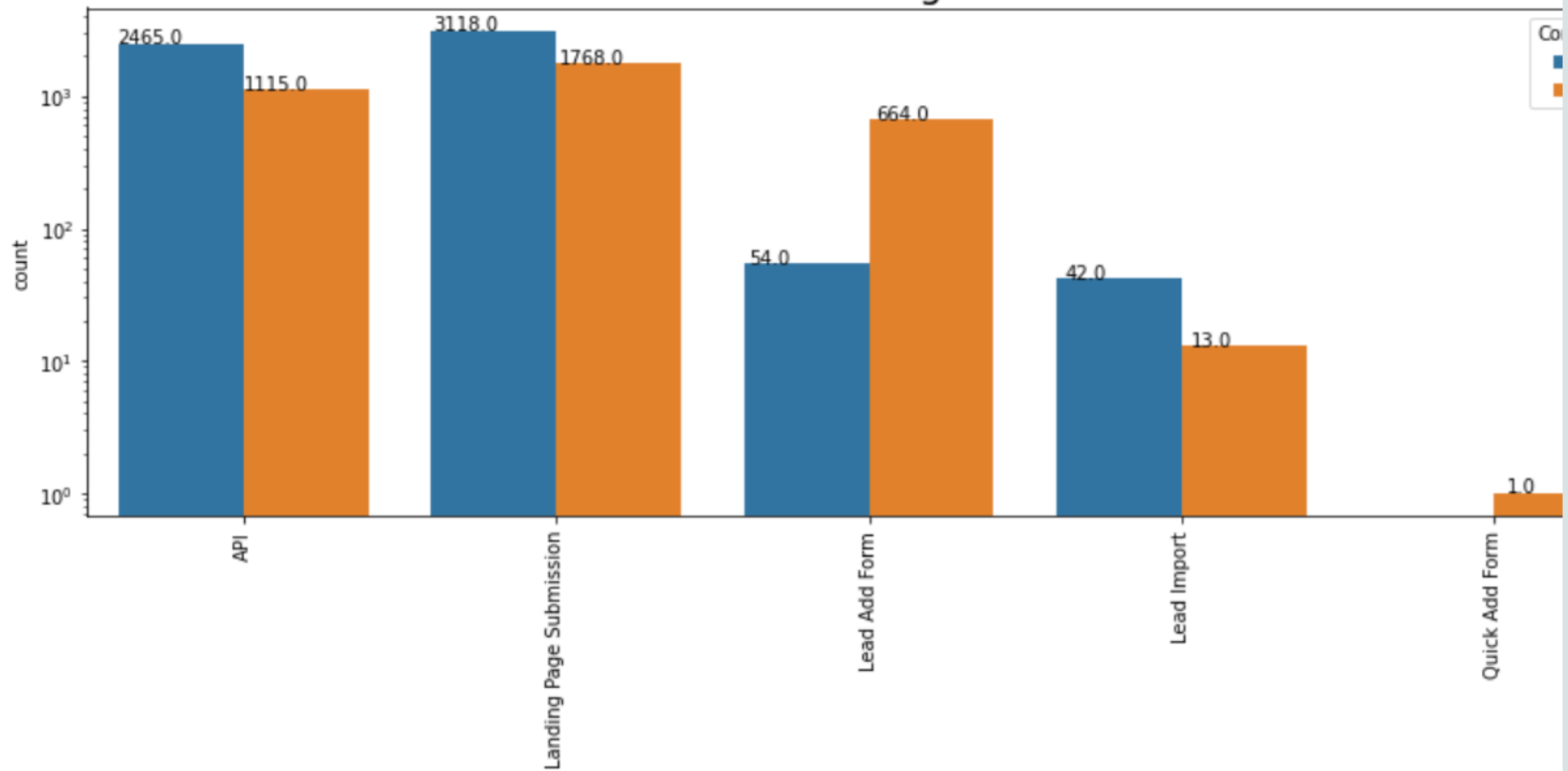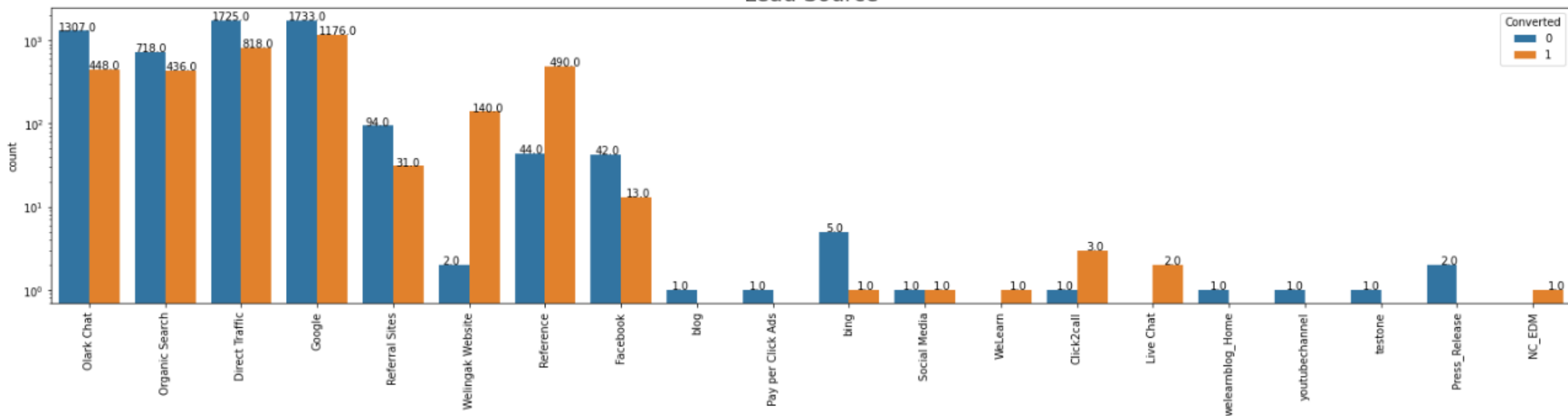- Conclusions and recommendations.
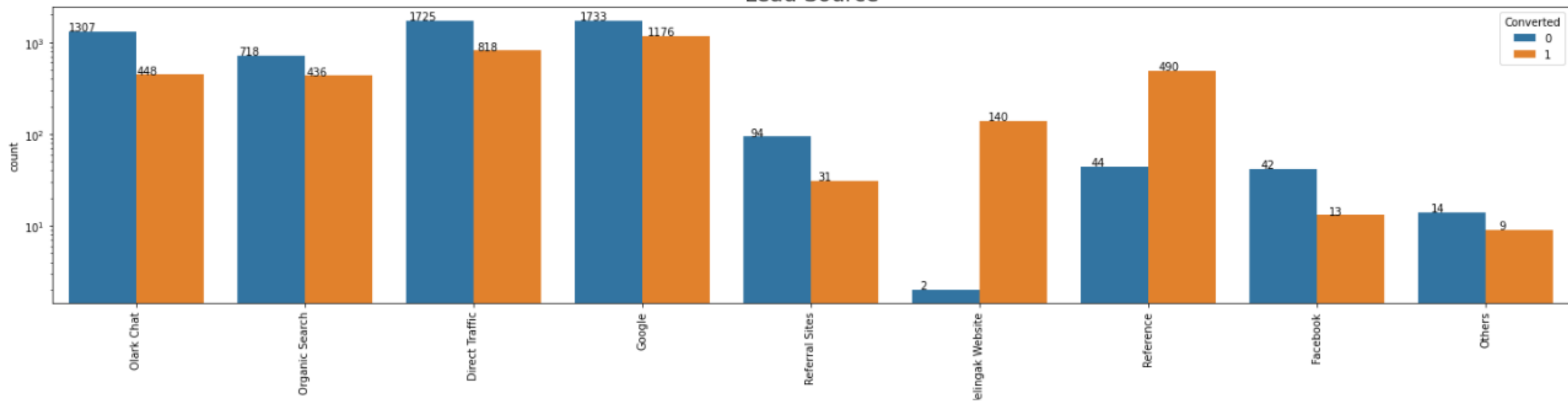
# EDA

# Categorical Variable Relation

Lead Source

## Lead Source

Converted
- 0
- 1

| | 1307 | | | 1725 | | 1733 | | | | | | | | | | | | |
| 448 | | 718 | 436 | | 818 | | 1176 | | | | | | 490 | | | | | |
| | | | | | | | | 94 | 31 | | 140 | | | 44 | | 42 | 13 | 14 | 9 |
| | | | | | | | | | | 2 | | | | | | | | | |

Olark Chat · Organic Search · Direct Traffic · Google · Referral Sites · lelingak Website · Reference · Facebook · Others

## Do Not Email

Converted
- 0
- 1

| 5063 | 3443 | | |
| | | 616 | 118 |

No · Yes

Last Activity of Lead
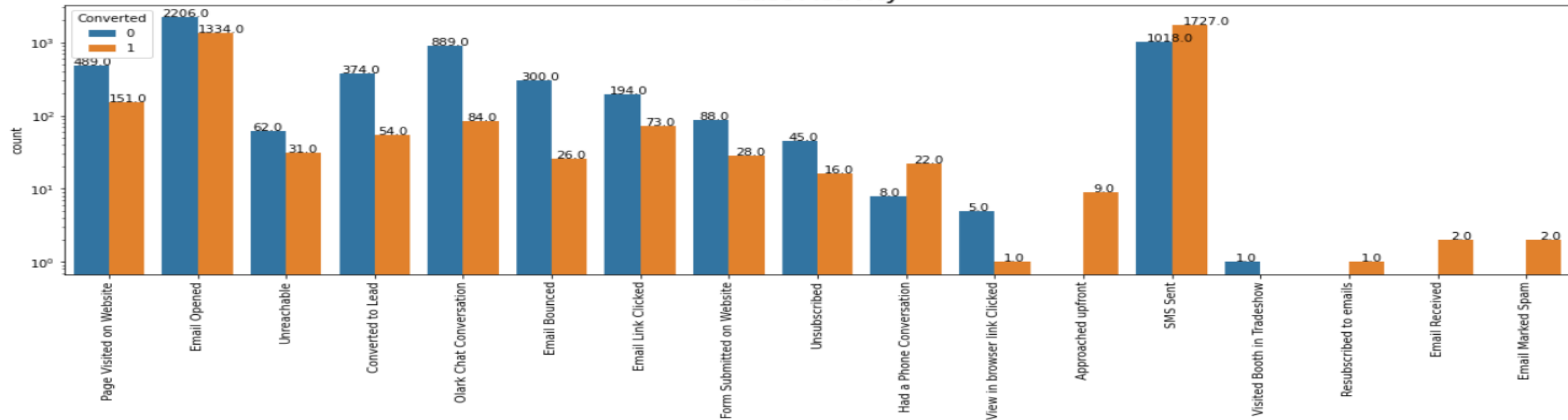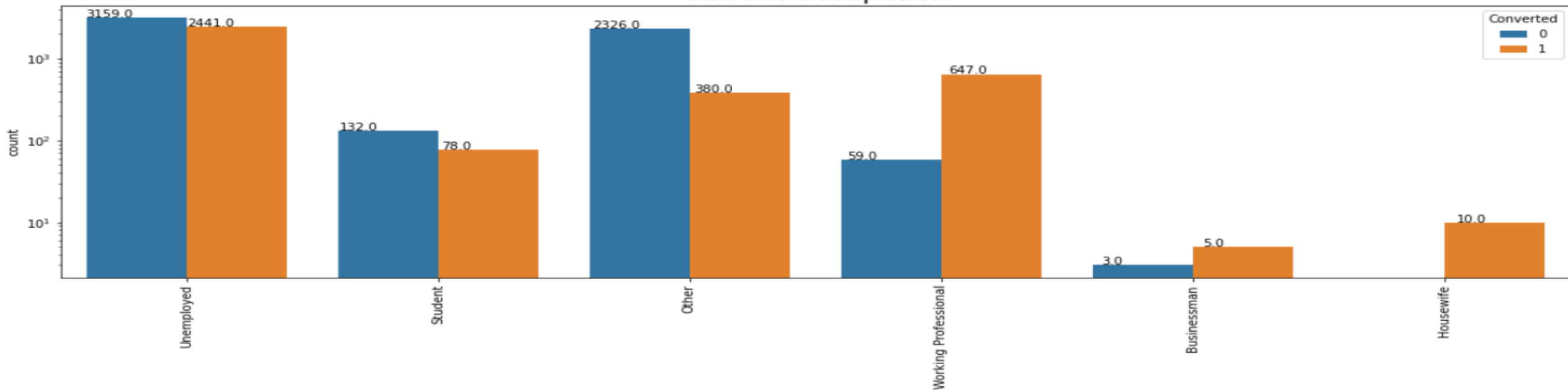
Last Activity

## Current Occupation

Converted
- 0
- 1

| Unemployed | Student | Other | Working Professional | Businessman | Housewife |
|---|---|---|---|---|---|
| 3159.0 / 2441.0 | 132.0 / 78.0 | 2326.0 / 380.0 | 59.0 / 647.0 | 3.0 / 5.0 | 10.0 |

## Specialization

Converted
- 0
- 1

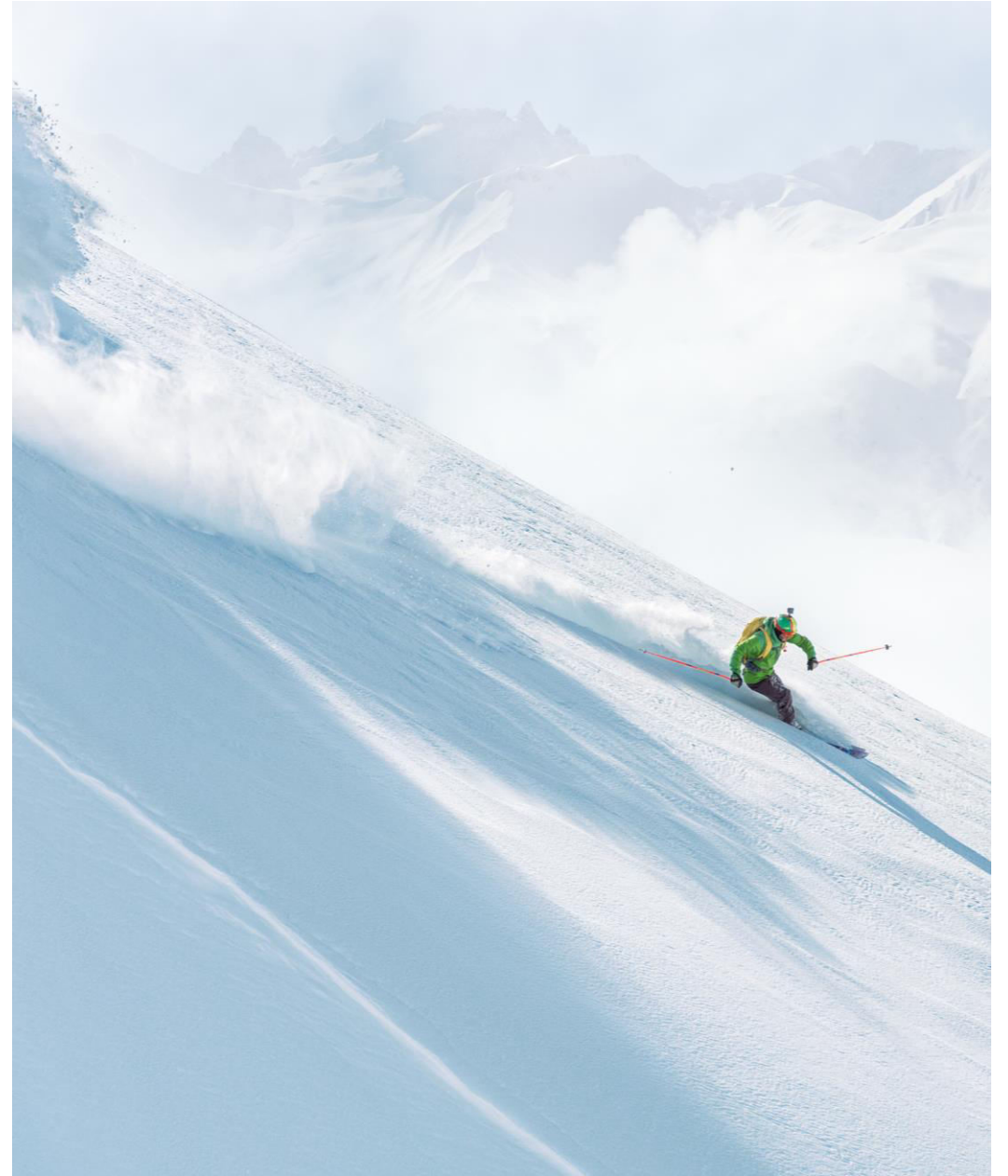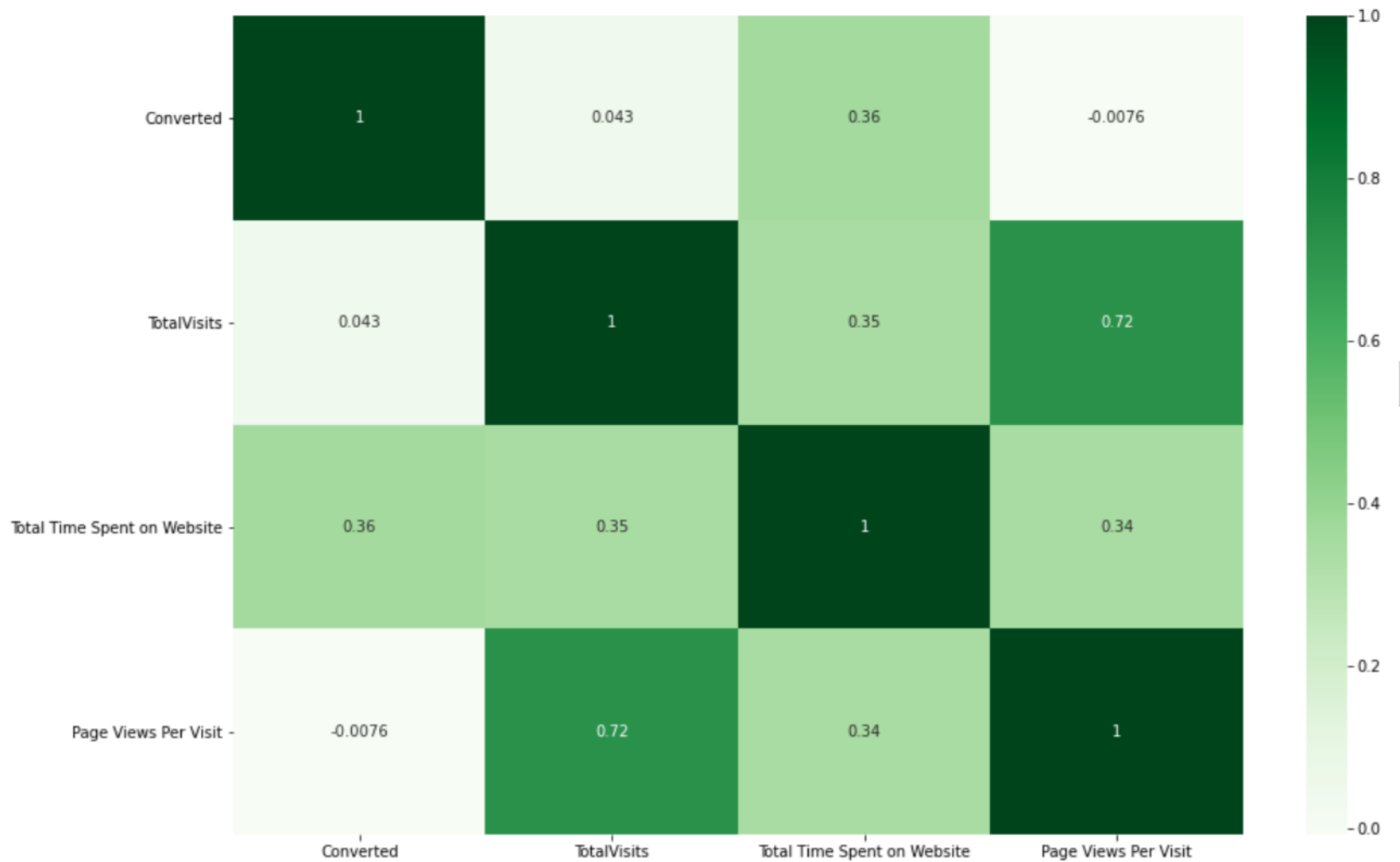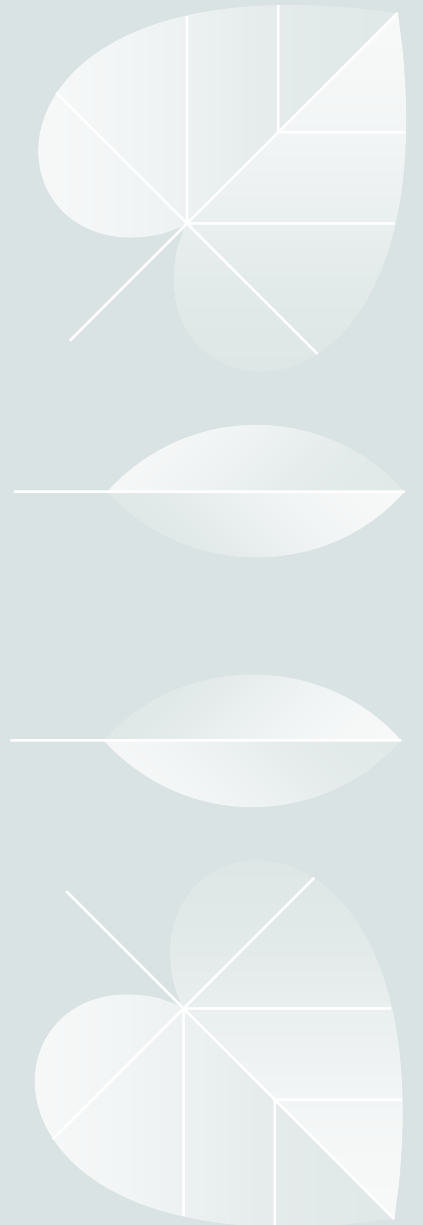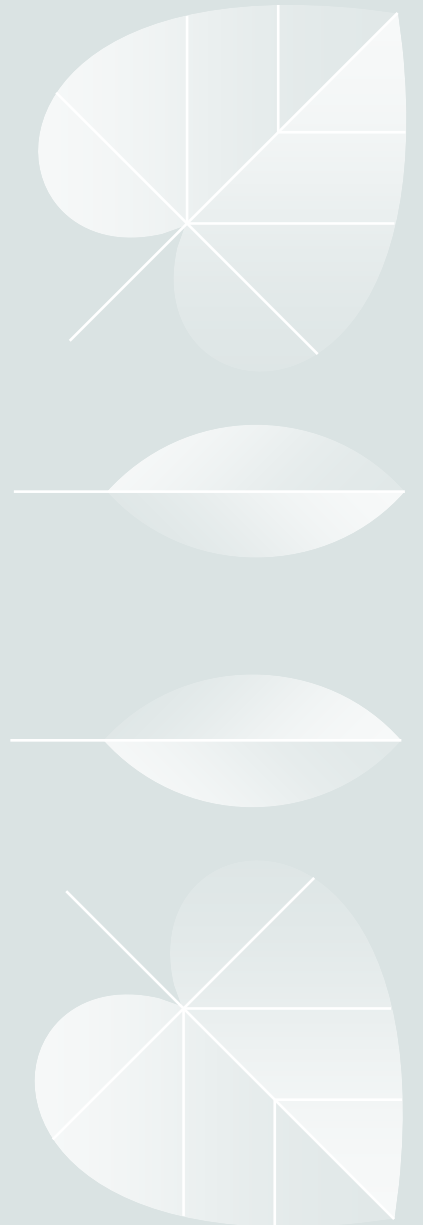| Other | Business Administration | Media and Advertising | Supply Chain Management | IT Projects Management | Finance Management | Travel and Tourism | Human Resource Management | Marketing Management | Banking, Investment And Insurance | International Business | E-COMMERCE | Operations Management | Retail Management | Services Excellence | Hospitality Management | Rural and Agribusiness | Healthcare Management | E-Business |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 2411 / 969 | 224 / 179 | 118 / 85 | 198 / 151 | 226 / 140 | 540 / 436 | 131 / 72 | 460 / 388 | 430 / 408 | 171 / 167 | 114 / 64 | 72 / 40 | 265 / 238 | 66 / 34 | 29 / 11 | 66 / 48 | 42 / 31 | 80 / 79 | 36 / 21 |

# Bivariate Analysis

# Data Conversion

- Numerical Variables are Normalised

- Dummy Variables are created for object type variables

- Total Rows for Analysis: 8792
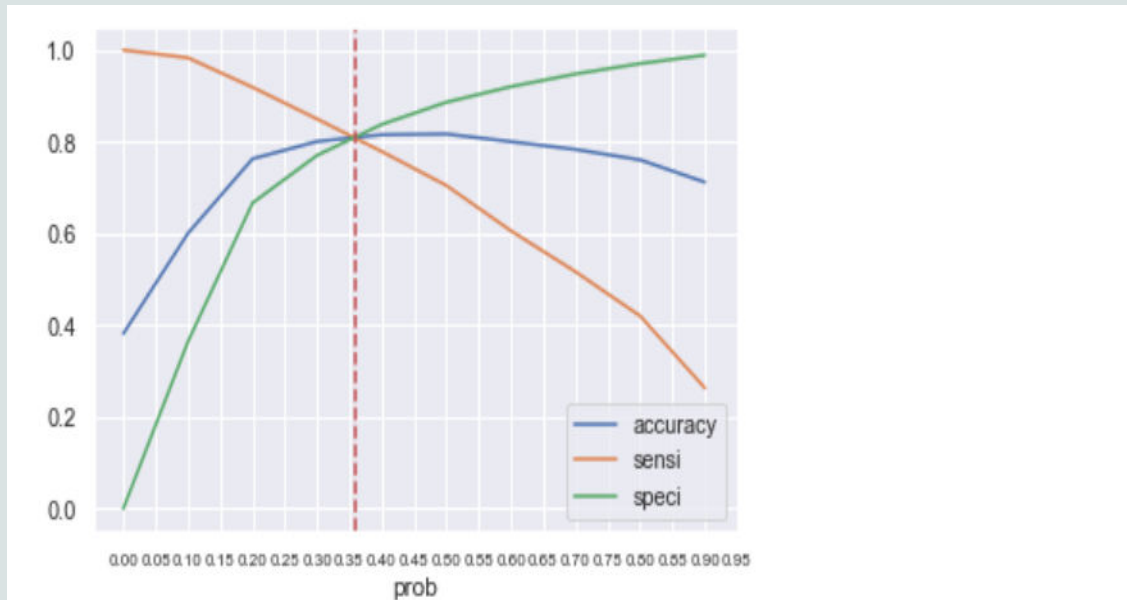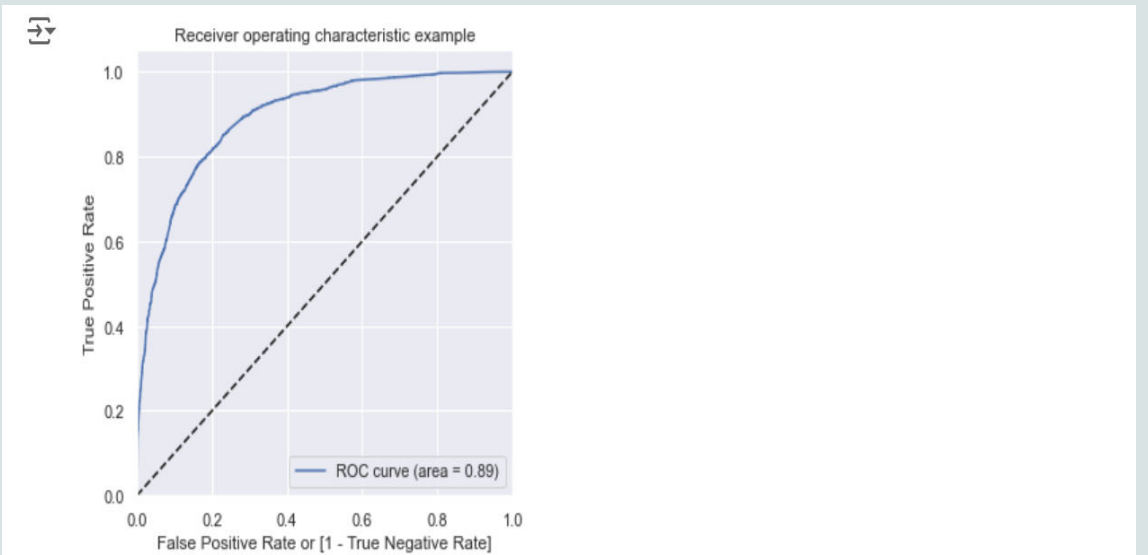
- Total Columns for Analysis: 43

# Model Building

- Splitting the Data into Training and Testing Sets.

- The first basic step for regression is performing a train-test split, we have chosen 70:30 ratio.

- Use RFE for Feature Selection.

- Running RFE with 15 variables as output.

- Building Model by removing the variable whose p- value is greater than 0.05 and vif value is greater than 5.

- Predictions on test data set ⯈ Overall accuracy 81%.

# ROC Curve



From the curve above, it seems that 0.358 is optimal cutoff point to take .



Observation

We are getting a good value of 0.89 indicating a good predictive model.As ROC Curve should be a value close to 1.

# Conclusion

➢ It was found that the variables that mattered the most in the potential buyers are (In descending order) :

➢ The total time spend on the Website.

➢ Total number of visits.

➢ When the lead source was:
   a. *Google*
   b. *Direct traffic*
   c. *Organic search*
   d. *Welingak website*

➢ When the last activity was:
   a. *SMS*
   b. *Olark chat conversation*

➢  When the lead origin is Lead add format.

➢ When their current occupation is as a working professional. Keeping these in mind the X Education can flourish as they have a very high chance to get almost all the potential buyers to change their mind and buy their courses.

# Thank you

Harshit Kumar

Gautami Pravin Wankhede