# X Education Lead Scoring

By Vladimir Nikonov

# Solution Strategy

1. Business problem understanding.
2. Mapping a business problem to a data science / machine learning problem.
3. Data sourcing, EDA, pre-processing.
4. Model building.
5. Model evaluation.

# 1. Business Problem Understanding

- An education company called *X Education* provides online courses to industry professionals.

- The main source of revenue for the company is the customers buying their courses.

- The typical lead conversion rate at *X Education* is **30%**.

- *X Education* requires a machine learning model, that will calculate the lead score. The leads with a higher conversion probability will get higher scores, while the leads with lower conversion probability will get lower score.

# 2. Problem Mapping

- Customer:
  - *X Education*'s sales team. They communicate with the leads, the model can help them identify the leads that they need to focus on communicating with.

- Expected ML Solution:
  - A Logistic Regression Model.
  - The model computes the lead score, which is basically the probability of conversion.

- Data Source:
  - Leads dataset from the past, provided by X Education.

# 3. Data Sourcing, EDA, Pre-processing
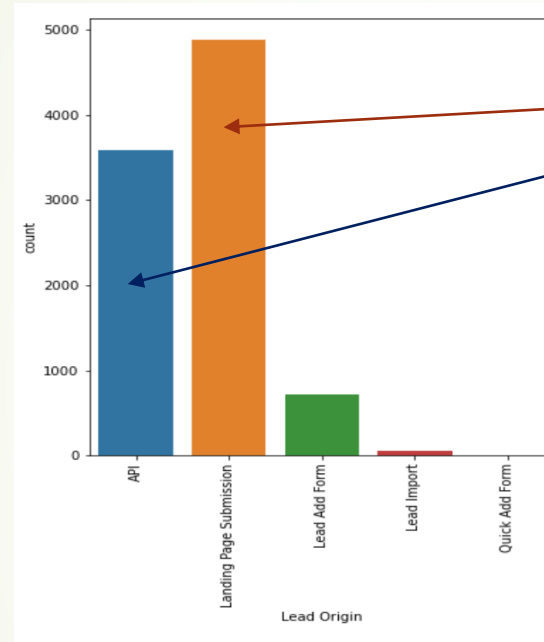
- 3.1 Data Sourcing
  - The leads data set with 9240 rows and 37 columns came from X Education. No duplicate rows were found in the data set.
- 3.2 Data Cleaning
  - After removing the columns with high percentage of missing values (>= 40%) and imbalanced data, only 10 features were left (3 numerical and 7 categorical).
  - The rest of the missing values was imputed with most frequent values for categorical features and median values for numerical features.
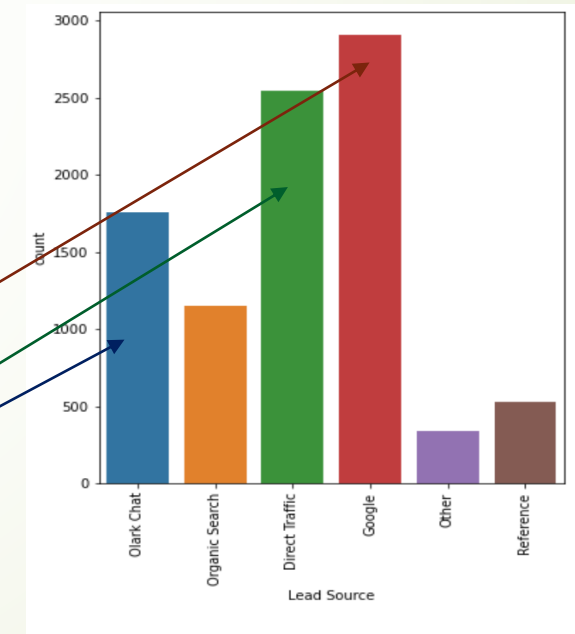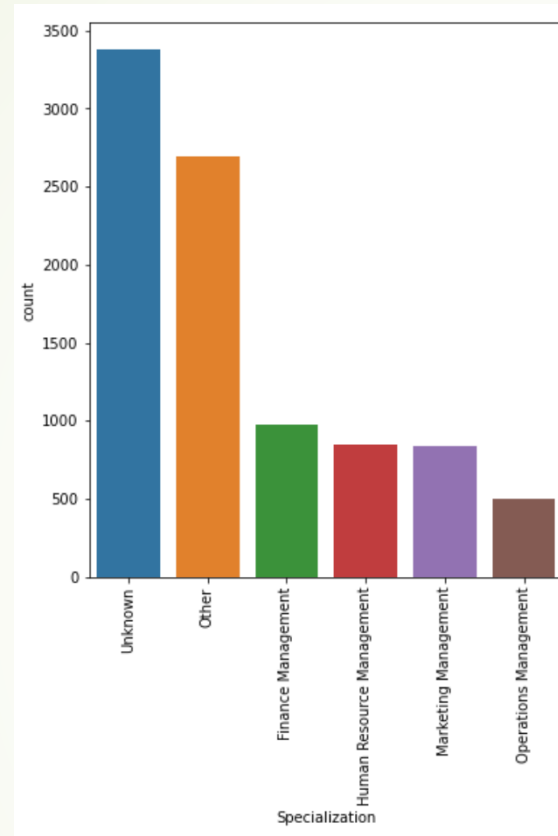
- 3.3 EDA

- 3.3.1 Univariate Analysis



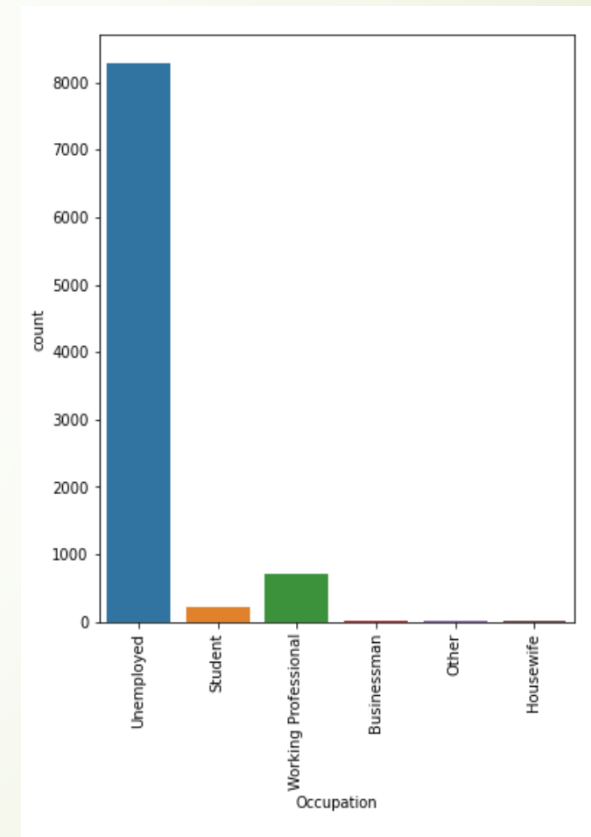Leads mostly originate from landing a page submission or through the API.

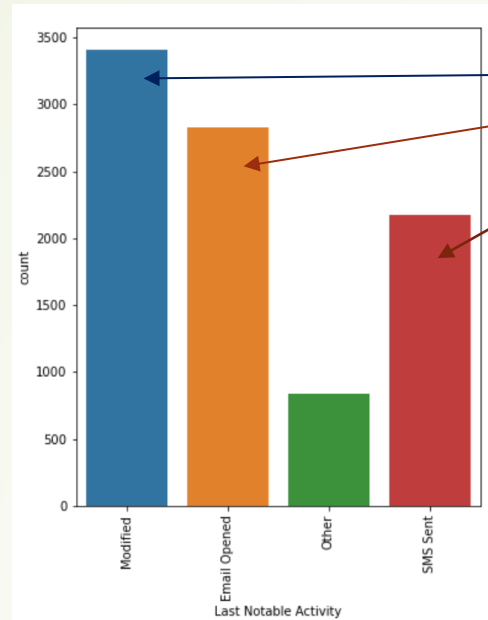The three main sources of leads are:
Google
Direct Traffic
Olark Chat

Out of known specializations, the ones connected to management are more popular with the leads.
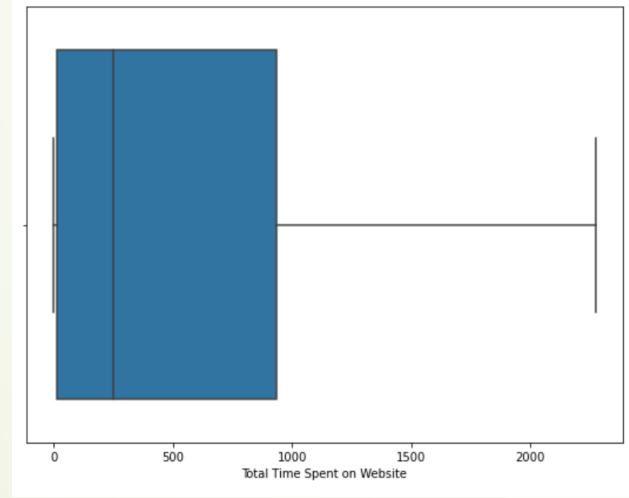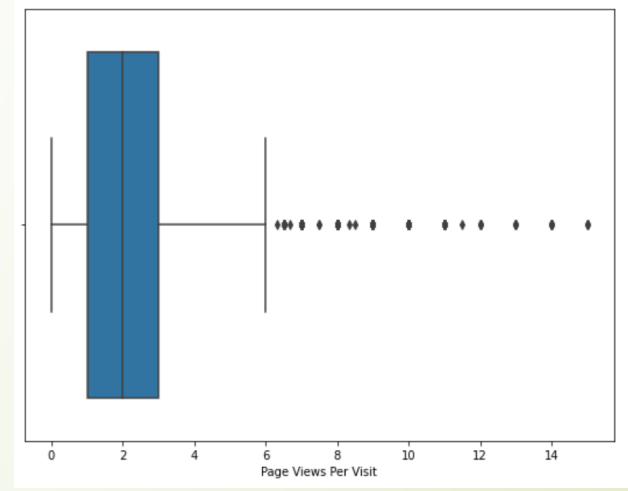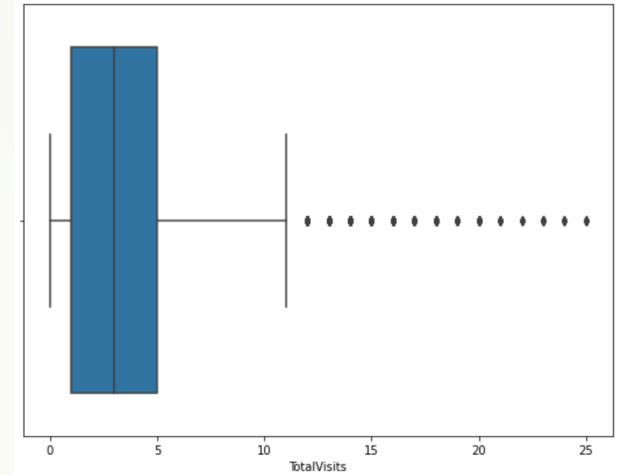


Most of the leads are unemployed.

The most occurring last notable activity is Modified with
Email Opened and
SMS Sent being second and third respectively.
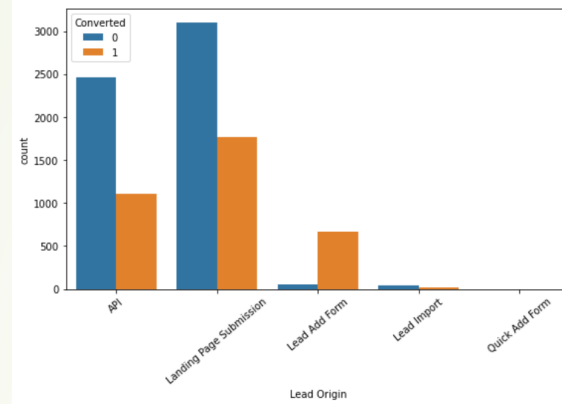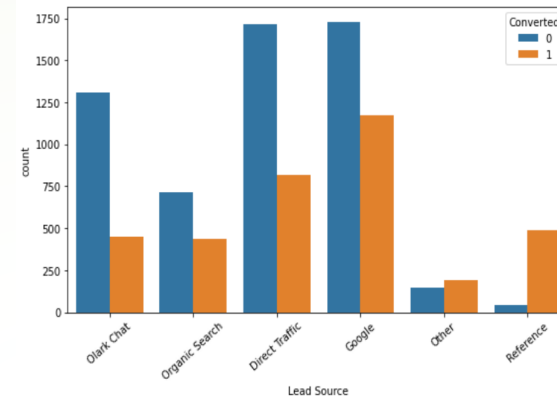


After getting rid of most outliers using the IQR method, we can see that 75% of the leads visit the website less than 5 times and view less than 6 pages per visit.

75% of the leads spent less than 1000 seconds on X Education website

# 3.3.2 Bivariate Analysis







In the column 'Lead Origin', much more leads that originated from the lead add forms got converted than not. Lead add form can be a good indicator of a lead becoming converted.

The percentage of converted leads is much higher than non-converted among the leads that were referred to X Education. Being referred can be a good indicator of a future conversion.

Most leads among working professionals got converted, which can be explained by their motivation to get better career prospects at the place of their current employment. Lead being a working professional can be a good indicator of a lead becoming converted.

## 3.4 Pre-processing

- Train-test split with 0.7 ratio was applied.

- The numeric variables were scaled using the Standard Scaler.

- The correlation matrix was drawn. Even though some variables had high correlation, it was decided not to remove them and let the EDA deal with them.

# 4. Model Building

- RFE was performed to select 15 most important features.

- Logistic regression model was built using statsmodels.api.

- Multicollinearity and significance of the variables were checked using the VIF scores and p-values. If a variable was insignificant or had high multicollinearity, it was removed and the process was repeated.

- In the end, 9 features were left in the final model.

## Generalized Linear Model Regression Results

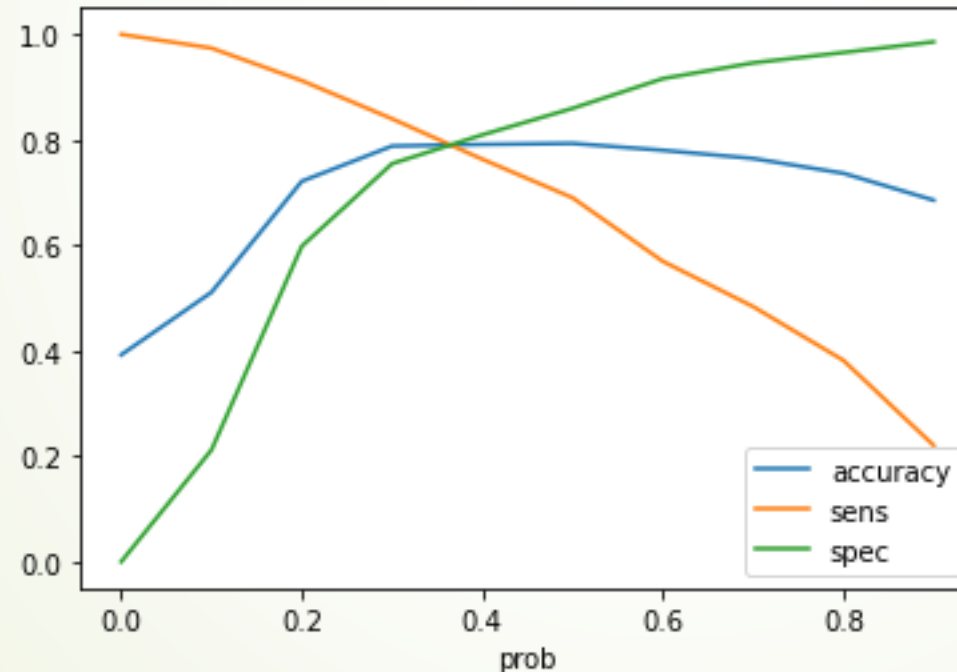| Dep. Variable: | Converted | No. Observations: | 6449 |
|---|---|---|---|
| Model: | GLM | Df Residuals: | 6439 |
| Model Family: | Binomial | Df Model: | 9 |
| Link Function: | Logit | Scale: | 1.0000 |
| Method: | IRLS | Log-Likelihood: | -2912.7 |
| Date: | Sun, 13 Nov 2022 | Deviance: | 5825.5 |
| Time: | 23:12:43 | Pearson chi2: | 6.72e+03 |
| No. Iterations: | 6 | Pseudo R-squ. (CS): | 0.3535 |
| Covariance Type: | nonrobust | | |

| | coef | std err | z | P>|z| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| const | -0.3172 | 0.108 | -2.925 | 0.003 | -0.530 | -0.105 |
| Total Time Spent on Website | 1.1068 | 0.038 | 29.025 | 0.000 | 1.032 | 1.182 |
| LeadOrigin_Lead Add Form | 3.9960 | 0.408 | 9.805 | 0.000 | 3.197 | 4.795 |
| LeadSource_Direct Traffic | -1.5322 | 0.122 | -12.520 | 0.000 | -1.772 | -1.292 |
| LeadSource_Google | -1.0607 | 0.111 | -9.593 | 0.000 | -1.277 | -0.844 |
| LeadSource_Organic Search | -1.1337 | 0.130 | -8.739 | 0.000 | -1.388 | -0.879 |
| LeadSource_Reference | -1.0591 | 0.451 | -2.348 | 0.019 | -1.943 | -0.175 |
| Specialization_Unknown | -0.7688 | 0.087 | -8.829 | 0.000 | -0.939 | -0.598 |
| NotableActivity_Email Opened | 0.7235 | 0.077 | 9.395 | 0.000 | 0.573 | 0.874 |
| NotableActivity_SMS Sent | 2.0948 | 0.085 | 24.761 | 0.000 | 1.929 | 2.261 |

| | Features | VIF |
|---|---|---|
| 1 | LeadOrigin_Lead Add Form | 4.35 |
| 5 | LeadSource_Reference | 4.16 |
| 7 | NotableActivity_Email Opened | 1.54 |
| 3 | LeadSource_Google | 1.49 |
| 8 | NotableActivity_SMS Sent | 1.48 |
| 2 | LeadSource_Direct Traffic | 1.36 |
| 6 | Specialization_Unknown | 1.36 |
| 0 | Total Time Spent on Website | 1.20 |
| 4 | LeadSource_Organic Search | 1.20 |

# 5. Model Evaluation

- In order to evaluate the model's performance, an optimal cutoff for probability had to be chosen.

- In order to do that, a plot of accuracy, sensitivity and specificity for all cutoffs was plotted.



- Judging by the plot, the optimal cutoff was somewhere between 0.35 and 0.38.
- In order to decide the exact cutoff, accuracy, sensitivity and specificity were calculated for the probabilities of 0.35, 0.36, 0.37, 0.38.

```
Pred_0.35
Accuracy: 0.7942316638238487
Sensitivity: 0.8090909090909091
Specificity: 0.7846389385047206
-------------------------------

Pred_0.36
Accuracy: 0.7948519150255854
Sensitivity: 0.8023715415019763
Specificity: 0.7899974483286553
-------------------------------

Pred_0.37
Accuracy: 0.791905721817336
Sensitivity: 0.78300395256917
Specificity: 0.7976524623628477
-------------------------------

Pred_0.38
Accuracy: 0.7915955962164677
Sensitivity: 0.7758893280632411
Specificity: 0.801735136514417
-------------------------------
```

0.35 was chosen as the optimal cutoff with the accuracy of ~0.79, sensitivity ~0.81 and specificity ~0.78

- Finally, the predictions were made on the test set.
- The model performed well on the test set with the following scores:
  - Accuracy – 0.78
  - Sensitivity – 0.81
  - Specificity – 0.76

- The most important variables for identifying the hot leads are:

  - LeadOrigin_Lead Add Form
  - NotableActivity_SMS Sent
  - LeadSource_Direct Traffic
  - LeadSource_Organic Search
  - Total Time Spent on Website

  - LeadSource_Google
  - LeadSource_Reference
  - Specialization_Unknown
  - NotableActivity_Email Opened