

# LEAD SCORING CASE STUDY

## 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.

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%.

Now, although X Education gets a lot of leads, its lead conversion rate is very poor. For example, if, say, they acquire 100 leads in a day, only about 30 of them are converted. 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.

## Goals of the Case Study

There are quite a few goals for this case study.

Build a logistic regression model to assign a lead score between 0 and 100 to each of the leads which can be used by the company to target potential leads. A higher score would mean that the lead is hot, i.e. is most likely to convert whereas a lower score would mean that the lead is cold and will mostly not get converted.

Steps of model building:

1. First, we have loaded and cleaned the data.
2. Then we have performed univariate and bivariate analysis of the features under EDA.
3. Based on the univariate analysis we have seen that many columns are not adding any information to the model, hence we dropped them.
4. Then we performed dummy variable creation for the categorical features.
5. Once data was ready, we did train test split of the data.

6. Then scaling was performed using fit\_transform function of StandardScaler.
7. Automatic feature selection was done using RFE function.
8. Then we built model and done some manual feature elimination using VIF and PValue.
9. After building final model we have done model evaluation on test data where we have calculated metrics such as accuracy, sensitivity and specificity.

#Accuracy : 81.5 %

#Sensitivity : 68.5 %

#Specificity : 88.8 %

10. Below are the finally selected features of the model.

|  |         |
|--|---------|
| Lead Source_Welingak Website<br>99                               | 0.6872  |
| Lead Source_Reference<br>57                                      | 0.4493  |
| What is your current occupation_Housewife<br>14                  | 0.3965  |
| Last Activity_Had a Phone Conversation<br>47                     | 0.3660  |
| Last Notable Activity_Unreachable<br>36                          | 0.2913  |
| What is your current occupation_Working Professional<br>66       | 0.2808  |
| Last Activity_SMS Sent<br>32                                     | 0.1968  |
| Total Time Spent on Website<br>90                                | 0.1833  |
| Lead Source_Olark Chat<br>25                                     | 0.1744  |
| Last Activity_Converted to Lead<br>21                            | -0.0733 |
| Last Notable Activity_Modified<br>83                             | -0.0932 |
| Last Activity_Olark Chat Conversation<br>45                      | -0.1066 |
| Specialization_not provided<br>99                                | -0.1147 |
| Specialization_Select<br>07                                      | -0.1204 |
| Lead Origin_Landing Page Submission<br>99                        | -0.1349 |
| Do Not Email<br>88   | -0.1472 |
| What matters most to you in choosing a course_not provided<br>38 | -0.1618 |

11. Recommendations of the model:

1. The company should make calls to the leads coming from the lead sources "Welingak Websites" and "Reference" as these are more likely to get converted.
2. The company should make calls to the leads who are the "House Wives" as they are more likely to get converted.
3. The company should make calls to the leads whose last activity was SMS Sent as they are more likely to get converted.
4. The company should make calls to the leads who spent "more time on the websites" as these are more likely to get converted.
5. The company should not make calls to the leads whose last activity was "Olark Chat Conversation" as they are not likely to get converted.
6. The company should not make calls to the leads who chose the option of "Do not Email" as "yes" as they are not likely to get converted.