Lead Scoring Case Study

Introduction

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%.

Problem Statement

- Lead conversion is low for X Education.
- Company wants to identify "Hot Leads" to increase the conversion.
- Initial lead generation is good but in the funnelling stage there are many leads which is not converted at the bottom.
- Create a model so that customer with high lead score have a higher conversion chance and vice versa.

Step 1: Import important libraries and warnings

```
import numpy as np
import pandas as pd

import warnings
warnings.filterwarnings('ignore')

import time, warnings
import datetime as dt

from IPython.display import display
pd.options.display.max_columns = None

import matplotlib.pyplot as plt
import seaborn as sns
```

Step 2: Read file and understand data type

```
lead= pd.read_csv('Leads.csv')
lead.head()
lead.shape
 (9240, 37)
lead.info()
 <class 'pandas.core.frame.DataFrame'>
 RangeIndex: 9240 entries, 0 to 9239
Data columns (total 37 columns):
 # Column
                                                   Non-Null Count Dtype
                                                   -----
 0 Prospect ID
                                                   9240 non-null
                                                                  object
     Lead Number
                                                   9240 non-null
                                                                  int64
     Lead Origin
                                                   9240 non-null
                                                                  object
     Lead Source
                                                   9204 non-null
                                                                  object
                                                   9240 non-null
 4 Do Not Email
                                                                  object
                                                   9240 non-null
 5 Do Not Call
                                                                  object
 6 Converted
                                                   9240 non-null
     TotalVisits
                                                  9103 non-null
                                                                  float64
     Total Time Spent on Website
                                                   9240 non-null
                                                                  int64
     Page Views Per Visit
                                                   9103 non-null
                                                                  float64
 10 Last Activity
                                                  9137 non-null
                                                                  object
                                                   6779 non-null
 11 Country
                                                                  object
 12 Specialization
                                                   7802 non-null
                                                                  object
```

Step 3: Check for Null value and remove null values

```
#checking for null
 lead.isnull().sum()
 Prospect ID
                                                     0
 Lead Number
                                                     0
 Lead Origin
                                                     36
 Lead Source
 Do Not Email
 Do Not Call
 Converted
                                                     0
 TotalVisits
                                                   137
 Total Time Spent on Website
 Page Views Per Visit
                                                   137
 Last Activity
                                                   103
 Country
                                                  2461
 Specialization
                                                  1438
 How did you hear about X Education
                                                  2207
 What is your current occupation
                                                  2690
 What matters most to you in choosing a course
                                                  2709
 Search
                                                     0
                                                     0
 Magazine
           . . . .
#Since there are not much distinct value in city and country hence dropping off
lead.drop(['City'], axis=1, inplace=True)
lead.drop(['Country'], axis=1, inplace=True)
#checking the % of null value
round(100*(lead.isnull().sum()/len(lead.index)),2)
Prospect ID
                                                  0.00
Lead Number
                                                  0.00
Lead Origin
                                                  0.00
```

Step 4: Variables which were considered for next set of process

<pre>lead.isnull().sum()</pre>								
Prospect ID	0							
Lead Number	0							
Lead Origin	0							
Lead Source	0							
Do Not Email	0							
Converted	0							
TotalVisits	0							
Total Time Spent on Website	0							
Page Views Per Visit	0							
Last Activity	0							
Specialization	0							
What is your current occupation	0							
A free copy of Mastering The Interview	0							
Last Notable Activity dtype: int64	0							

Step 5: Prepare data for modelling

Created dummy variable for the object data type

```
lead.head()
```

	Converted	TotalVisits	Total Time Spent on Website	Page Views Per Visit	Lead Origin_Landing Page Submission	Lead Origin_Lead Add Form	Lead Origin_Lead Import	Lead Source_Direct Traffic	Lead Source_Facebook	Lead Source_Google	Lead Source_Live Chat	Source_
0	0	0.0	0	0.0	0	0	0	0	0	0	0	
1	0	5.0	674	2.5	0	0	0	0	0	0	0	
2	1	2.0	1532	2.0	1	0	0	1	0	0	0	
3	0	1.0	305	1.0	1	0	0	1	0	0	0	
4	1	2.0	1428	1.0	1	0	0	0	0	1	0	

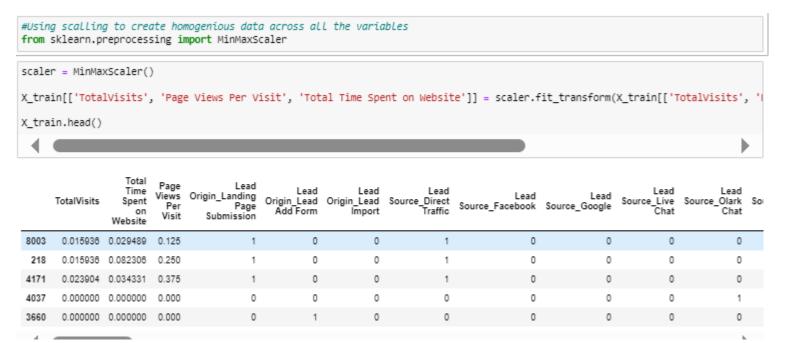
Step 6: Train and Split data

Data was split into 70% train and 30% test

```
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
import statsmodels.api as sm
from statsmodels.stats.outliers_influence import variance_inflation_factor
#X = Lead.drop(['Converted'], axis=1, inplace =True)
X=lead
X=X.drop(['Converted'],axis=1)
X.head()
                Total
                      Page
                                     Lead
                Time
                                                Lead
                                                            Lead
                                                                         Lead
                                                                                                                   Lead
                                                                                                                                Lead
                      Views
                            Origin_Landing
                                                                                         Lead
                                                                                                       Lead
   TotalVisits
                                                                                                             Source_Live Source_Olark Source
                Spent
                                          Origin_Lead Origin_Lead Source_Direct
                                                                               Source_Facebook Source_Google
                       Per
                                     Page
                                                                                                                   Chat
                                            Add Form
                                                           Import
                                                                                                                                Chat
                  on
                       Visit
                               Submission
             Website
                                                                                                                      0
0
                                       0
                                                   0
                                                               0
                                                                                            0
                                                                                                          0
         0.0
                        0.0
                 674
                        2.5
                                       0
                                                   0
                                                               0
                                                                                            0
                                                                                                                      0
                                                                                                                                   0
1
         5.0
2
         2.0
                1532
                        2.0
                                                               0
                                                                                            0
                                                                                                                      0
3
                                                   0
                                                               0
                                                                                            0
                                                                                                          0
                                                                                                                      0
                                                                                                                                   0
         1.0
                 305
                        1.0
                                                               0
                                                                                                                      0
         2.0
                1428
                        1.0
                                                                                                                                   0
y=lead['Converted']
y.head()
     0
0
3
4
Name: Converted, dtype: int64
#SpLit the data into 70% train and 30% test data
X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=0.7, test_size=0.3, random_state=100)
```

Step 7: Scaling the data

Using scaling to create homogeneous data across all the variables



Step 8: Model Building

Used RFE to identify variable required to create the model

```
# Import 'LogisticRegression' and create a LogisticRegression object
from sklearn.linear model import LogisticRegression
logreg = LogisticRegression()
# Import RFE and select 15 variables
from sklearn.feature_selection import RFE
rfe = RFE(logreg,n_features_to_select=15)
rfe = rfe.fit(X_train, y_train)
#Features have been selected by RFE
list(zip(X_train.columns, rfe.support_, rfe.ranking_))
[('Totalvisits', True, 1),
 ('Total Time Spent on Website', True, 1),
 ('Page Views Per Visit', False, 11),
 ('Lead Origin_Landing Page Submission', False, 2),
  ('Lead Origin_Lead Add Form', True, 1),
 ('Lead Origin_Lead Import', False, 49),
  ('Lead Source_Direct Traffic', False, 16),
 ('Lead Source Facebook', False, 42),
  ('Lead Source_Google', False, 35),
  ('Lead Source_Live Chat', False, 39),
  ('Lead Source_Olark Chat', True, 1),
  ('Lead Source_Organic Search', False, 34),
  ('Lead Source_Pay per Click Ads', False, 33),
 ('Lead Source_Press_Release', False, 53),
 ('Lead Source_Reference', False, 4),
  ('Lead Source_Referral Sites', False, 36),
 ('Lead Source_Social Media', False, 58),
  ('Lead Source_WeLearn', False, 25),
  ('Lead Source_Welingak Website', True, 1),
 ('Lead Source_bing', False, 21),
 ('Lead Source testone', False, 32),
```

Step 9: Fitting logistic regression on the train data set Used RFE to identify variable required to create the model

Specialization Select -0.3400

```
: # Fitting Logistic Regression model on X train post adding constant
  X_train_sm = sm.add_constant(X_train)
  logm2 = sm.GLM(y_train, X_train_sm, family = sm.families.Binomial())
  res = logm2.fit()
  res.summary()
  Generalized Linear Model Regression Results
     Dep. Variable:
                                                           4461
                          Converted
                                      No. Observations:
            Model:
                              GLM
                                          Df Residuals:
                                                           4445
     Model Family:
                           Binomial
                                             Df Model:
                                                             15
     Link Function:
                              Logit
                                                Scale:
                                                          1.0000
           Method:
                              IRLS
                                        Log-Likelihood:
                                                         -2087.2
             Date:
                   Tue. 20 Feb 2024
                                             Deviance:
                                                          4134.4
             Time:
                           11:48:44
                                          Pearson chi2: 4.83e+03
                                22 Pseudo R-squ. (CS):
                                                         0.3676
     No. Iterations:
   Covariance Type:
                          nonrobust
                                                               std err
                                                                           z P>|z|
                                                                                       [0.025
                                                                                                0.975]
                                                        coef
                                                     -0.9490
                                                                                                 0.233
                                              const
                                                                0.603
                                                                       -1.573 0.116
                                                                                       -2.131
                                          TotalVisits
                                                     10.2343
                                                                2.636 3.882 0.000
                                                                                        5.068
                                                                                               15.401
                          Total Time Spent on Website
                                                                                                 4.768
                                                      4.4045
                                                                0.188 23.735 0.000
                                                                                        4.041
                          Lead Origin_Lead Add Form
                                                                0.259 16.363 0.000
                                                                                        3.729
                                                                                                 4.744
                                                                                        1.372
                             Lead Source_Olark Chat
                                                     1.6324
                                                                0.133 12.267 0.000
                                                                                                 1.893
                       Lead Source_Welingak Website
                                                     2.3444
                                                                1.038 2.258 0.024
                                                                                        0.310
                                                                                                 4.379
                                   Do Not Email_Yes
                                                                                       -1.895
                                                                                                -1.141
                                                     -1.5177
                                                                0.192 -7.892 0.000
                Last Activity_Had a Phone Conversation
                                                                0.987 1.186 0.235
                                                                                        -0.764
                                                                                                3.108
                                                     1.1713
                              Last Activity_SMS Sent
                                                                0.082 14.305 0.000
                                                                                        1.017
                                                                                                 1.340
            What is your current occupation_Housewife 22.6104 2.45e+04 0.001 0.999
                                                                                      -4.8e+04 4.8e+04
               What is your current occupation_Student
                                                                0.634 -1.776 0.078
                                                                                        -2.369
                                                                                                 0.117
          What is your current occupation_Unemployed
                                                     -1.2968
                                                                0.598 -2.169 0.030
                                                                                        -2.468
                                                                                               -0.125
   What is your current occupation_Working Professional
                                                                0.627
                                                                       1.992 0.046
                                                                                        0.020
                                                                                                 2.476
                                                                       0.001 0.999 -4.09e+04 4.1e+04
        Last Notable Activity_Had a Phone Conversation 23.0108 2.09e+04
                    Last Notable Activity_Unreachable
                                                                0.807 3.429 0.001
                                                                                        1.186
```

0.098 -3.464 0.001

-0.532

-0.148

Step 10: Checking VIF

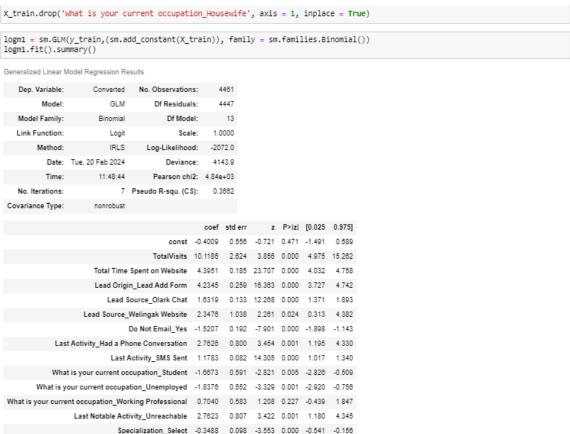
Checking Variation inflation factor to manually remove any variable where VIF is >5

```
vif = pd.DataFrame()
vif['Features'] = X_train.columns
vif['VIF'] = [variance_inflation_factor(X_train.values, i) for i in range(X_train.shape[1])]
vif['VIF'] = round(vif['VIF'], 2)
vif = vif.sort_values(by = "VIF", ascending = False)
vif
                                      Features VIF
        What is your current occupation_Unemployed 4.13
             Last Activity_Had a Phone Conversation 2.44
 12 Last Notable Activity_Had a Phone Conversation 2.43
                      Total Time Spent on Website 2.39
                            Specialization_Select 1.90
  2
                      Lead Origin Lead Add Form 1.71
                         Lead Source Olark Chat 1.66
  0
                                      TotalVisits 1.63
                          Last Activity_SMS Sent 1.59
 11 What is your current occupation_Working Profes... 1.58
                   Lead Source_Welingak Website 1.37
  9
            What is your current occupation_Student 1.10
  5
                               Do Not Email_Yes 1.09
  8
          What is your current occupation_Housewife 1.01
                  Last Notable Activity_Unreachable 1.01
```

Since VIF for the all the variable was within 5 hence it good to go and now our focus is on p value >0.05

Step 11: Checking p-value

For all the variable where p-value is greater than 0.05 remove it



Removed variables- 'What is your current occupation_Working Professional, What is your current occupation_Housewife'

y_train_pred_final.head()

Step 12: Model Evaluation

```
# 'predict' is used to predict the probabilities on the train set
y_train_pred = res.predict(X_train_sm)
y_train_pred[:10]
       0.315577
218
        0.151844
4171
      0.135876
      0.278192
        0.959650
        0.156043
      0.143676
      0.952580
       0.079814
       0.981919
dtype: float64
# Reshaping it into an array
y_train_pred = y_train_pred.values.reshape(-1)
y_train_pred[:10]
array([0.3155766 , 0.1518439 , 0.13587609, 0.27819235, 0.95965009,
       0.1560432 , 0.14367596, 0.95258003, 0.07981364, 0.98191931])
# Create a new dataframe containing the actual conversion flag and the probabilities predicted by the model
y_train_pred_final = pd.DataFrame({'Converted':y_train.values, 'Conversion_Prob':y_train_pred})
y_train_pred_final.head()
   Converted Conversion_Prob
                   0.315577
                   0.151844
                   0.135876
                   0.278192
                   0.959650
#Creating new column 'Predicted' with 1 if Paid_Prob > 0.5 else 0
y_train_pred_final['Predicted'] = y_train_pred_final.Conversion_Prob.map(lambda x: 1 if x > 0.5 else 0)
```

Step 12: Model Evaluation

Sensitivity = 79%, Specificity -78%, Accuracy = 79% - hence this is a good model.

```
: # Make predictions on the test set using 0.45 as the cutoff
 y_pred_final['final_predicted'] = y_pred_final.Conversion_Prob.map(lambda x: 1 if x > 0.42 else 0)
: # Check y_pred_final
  y_pred_final.head()
     Converted Conversion_Prob final_predicted
                     0.998712
                     0.137945
                     0.717442
                     0.311448
                     0.731041
 metrics.accuracy_score(y_pred_final['Converted'], y_pred_final.final_predicted)
0.7855648535564853
confusion2 = metrics.confusion_matrix(y_pred_final['Converted'], y_pred_final.final_predicted )
  confusion2
: array([[779, 217],
[193, 723]], dtype=int64)
: TP = confusion2[1,1] # true positive
  TN = confusion2[0,0] # true negatives
  FP = confusion2[0,1] # false positives
  FN = confusion2[1,0] # false negatives
: # Calculate sensitivity
  TP / float(TP+FN)
0.7893013100436681
: # Calculate specificity
  TN / float(TN+FP)
0.7821285140562249
```