

X Education: Hot Leads for a focussed sales follow-up.

Problem Statement:

X Education gets a lot of leads, but their lead conversion rate is poor.

Solution Statement:

Many leads in sales funnel are not actual potential students but some of them are very potential leads.

Solution is to identify the high potential prospective students and focus the sales follow-up activities on pursuing these hot leads.

Task:

Identify the potential leads from the response and sales activity data and score them accordingly.
The driving questions are:

- Which responses are significant in predicting whether the lead will convert to a sale or not?
- Generate a score for the convertible leads based on the probability of how potential a lead each of them could be.
- Find who are the Hot Leads based on the above score so that the sales team can focus their effort on those leads.

Solution

Step 1: Import and Inspect Data

```
In [1]: import warnings
warnings.filterwarnings('ignore')

In [2]: import pandas as pd
import numpy as np

In [3]: data=pd.read_csv('Leads.csv')
pd.set_option('max_columns',None)

In [4]: data.head()

Out [4]:
```

	ProspectID	Lead Number	Lead Origin	Lead Source	Do Not Email	Do Not Call	Converted	TotalVisits	Total Time Spent on Website	Page Views Per Visit	Last Activity	Country	Specialization
0	792762df-6bca-422b-b82c-b6e0bea6e20	660737	API	Clark Chat	No	No	0	0.0	0	0.0	Page Viewed on Website	NaN	Select
1	2a272436-4332-4326-a2f3-86fa-0cc8808942	660728	API	Organic Search	No	No	0	5.0	674	2.5	Email Opened	India	Select
2	0c02d445-7d44-4e39-9d49-19797f93b0cc	660727	Landing Page Submission	Direct Traffic	No	No	1	2.0	1532	2.0	Email Opened	India	Business Administration
3	0c02d445-7d44-4e39-9d49-19797f93b0cc	660719	Landing Page Submission	Direct Traffic	No	No	0	1.0	305	1.0	Unreachable	India	Media and Advertising
4	3250826b-6534-4026-9d53-4a8b08762352	660881	Landing Page Submission	Google	No	No	1	2.0	1428	1.0	Converted to Lead	India	Select

```
In [5]: data.columns

Out [5]: Index(['Prospect ID', 'Lead Number', 'Lead Origin', 'Lead Source', 'Do Not Email', 'Do Not Call', 'Converted', 'TotalVisits', 'Total Time Spent on Website', 'Page Views Per Visit', 'Last Activity', 'Country', 'Specialization',
      1, 'How did you hear about X Education', 'What is your current occupation', 'What matters most to you in choosing a course', 'Search', 'Magazine', 'Newspaper Article', 'X Education Forums', 'Newspaper', 'Digital Advertisement', 'Through Recommendations', 'Receive More Updates About Our Courses', 'Tags', 'Lead Quality', 'Update me on Supply Chain Content', 'Get updates on DM Content', 'Lead Profile', 'City', 'Asymmetrique Activity Index', 'Asymmetrique Profile Index', 'Asymmetrique Profile Score', 'I agree to pay the amount through cheque', 'A free copy of Mastering The Interview', 'Last Notable Activity', 'dtypes:object'])

In [6]: data.shape

Out [6]: (9240, 37)
```

Inspecting Missing Values

```
In [7]: data.isna().sum()

Out [7]:
```

	Prospect ID	Lead Number	Lead Origin	Lead Source	Do Not Email	Do Not Call	Converted	TotalVisits	Total Time Spent on Website	Page Views Per Visit	Last Activity	Country	Specialization
0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	0	0	0
5	0	0	0	0	0	0	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0	0	0	0	0	0	0
11	0	0	0	0	0	0	0	0	0	0	0	0	0
12	0	0	0	0	0	0	0	0	0	0	0	0	0
13	0	0	0	0	0	0	0	0	0	0	0	0	0
14	0	0	0	0	0	0	0	0	0	0	0	0	0
15	0	0	0	0	0	0	0	0	0	0	0	0	0
16	0	0	0	0	0	0	0	0	0	0	0	0	0
17	0	0	0	0	0	0	0	0	0	0	0	0	0
18	0	0	0	0	0	0	0	0	0	0	0	0	0
19	0	0	0	0	0	0	0	0	0	0	0	0	0
20	0	0	0	0	0	0	0	0	0	0	0	0	0
21	0	0	0	0	0	0	0	0	0	0	0	0	0
22	0	0	0	0	0	0	0	0	0	0	0	0	0
23	0	0	0	0	0	0	0	0	0	0	0	0	0
24	0	0	0	0	0	0	0	0	0	0	0	0	0
25	0	0	0	0	0	0	0	0	0	0	0	0	0
26	0	0	0	0	0	0	0	0	0	0	0	0	0
27	0	0	0	0	0	0	0	0	0	0	0	0	0
28	0	0	0	0	0	0	0	0	0	0	0	0	0
29	0	0	0	0	0	0	0	0	0	0	0	0	0
30	0	0	0	0	0	0	0	0	0	0	0	0	0
31	0	0	0	0	0	0	0	0	0	0	0	0	0
32	0	0	0	0	0	0	0	0	0	0	0	0	0
33	0	0	0	0	0	0	0	0	0	0	0	0	0
34	0	0	0	0	0	0	0	0	0	0	0	0	0
35	0	0	0	0	0	0	0	0	0	0	0	0	0
36	0	0	0	0	0	0	0	0	0	0	0	0	0
37	0	0	0	0	0	0	0	0	0	0	0	0	0

Checking percentage of missing values

```
In [8]: round(100*(data.isna().sum()/len(data.index)), 2)

Out [8]:
```

	Prospect ID	Lead Number	Lead Origin	Lead Source	Do Not Email	Do Not Call	Converted	TotalVisits	Total Time Spent on Website	Page Views Per Visit	Last Activity	Country	Specialization
0	0.00	0.00	0.00	0.39	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
1	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
2	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
3	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
4	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
5	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
6	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
7	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
8	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
9	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
10	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
11	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
12	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
13	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
14	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
15	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
16	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
17	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
18	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
19	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
20	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
21	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
22	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
23	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
24	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
25	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
26	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
27	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
28	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
29	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
30	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
31	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
32	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
33	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
34	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
35	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
36	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
37	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

Descriptive statistics

```
In [9]: data.describe()

Out [9]:
```

	Lead Number	Converted	TotalVisits	Total Time Spent on Website	Page Views Per Visit	Asymmetrique Activity Score	Asymmetrique Profile Score
count	9240.000000	9240.000000	9103.000000	9240.000000	9103.000000	5022.000000	5022.000000
mean	617188.435066	0.385390	3.445238	487.696268	2.362820	14.306252	16.344883
std	23405.905698	0.486714	4.854853	548.021466	2.161418	1.386694	1.811395
min	579533.000000	0.000000	0.000000	0.000000	0.000000	7.000000	11.000000
25%	598454.500000	0.000000	1.000000	12.000000	1.000000	14.000000	15.000000
50%	615479.000000	0.000000	3.000000	248.000000	2.000000	14.000000	16.000000
75%	637387.250000	1.000000	5.000000	936.000000	3.000000	15.000000	18.000000
max	660737.000000	1.000000	251.000000	2272.000000	55.000000	18.000000	20.000000

Inspect Data

```
In [10]: data.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
 1   Lead Number                         9240 non-null   int64
 2   Lead Origin                         9240 non-null   object
 3   Lead Source                         9204 non-null   object
 4   Do Not Email                       9240 non-null   object
 5   Do Not Call                         9240 non-null   object
 6   Converted                           9240 non-null   int64
 7   TotalVisits                         9103 non-null   float64
 8   Total Time Spent on Website         9240 non-null   object
 9   Page Views Per Visit                9103 non-null   float64
10   Last Activity                       9137 non-null   object
11   Country                             6779 non-null   object
12   Specialization                      7802 non-null   object
13   How did you hear about X Education  7033 non-null   object
14   What is your current occupation     6550 non-null   object
15   What matters most to you in choosing a course  6531 non-null   object
16   Search                             9240 non-null   object
17   Magazine                           9240 non-null   object
18   Newspaper Article                  9240 non-null   object
19   X Education Forums                 9240 non-null   object
20   Newspaper                          9240 non-null   object
21   Digital Advertisement              9240 non-null   object
22   Through Recommendations            9240 non-null   object
23   Receive More Updates About Our Courses  9240 non-null   object
24   Tags                               5887 non-null   object
25   Lead Quality                        4473 non-null   object
26   Update me on Supply Chain Content   9240 non-null   object
27   Get updates on DM Content           9240 non-null   object
28   Lead Profile                        7820 non-null   object
29   City                               5022 non-null   object
30   Asymmetrique Activity Index         5022 non-null   float64
31   Asymmetrique Profile Index          5022 non-null   object
32   Asymmetrique Profile Score          5022 non-null   float64
33   I agree to pay the amount through cheque  9240 non-null   object
34   A free copy of Mastering The Interview  9240 non-null   object
35   Last Notable Activity               9240 non-null   object
dtypes: float64(4), int64(3), object(30)
memory usage: 2.6+ MB
```

Step 2: Data Preparation (Missing Values, Insignificant Columns, Categorical Variables, Outliers etc)

Converting some binary variables (Yes/No) to 0/1

```
In [11]: # The following columns are Binary Variables:
# Do Not Email, Do Not call, Search, Magazine, Newspaper Article, X Education Forums, Newspaper,
# Digital Advertisement, Through Recommendations, Receive More Updates About Our Courses,
# Update me on Supply Chain Content, Get updates on DM Content, I agree to pay the amount through cheque,
# A free copy of Mastering The Interview

varlist=["Do Not Email","Do Not Call","Search","Magazine","Newspaper Article","X Education Forums",
"Newspaper","Digital Advertisement","Through Recommendations","Receive More
```



```
[68]: sns.boxplot(data['Page Views Per Visit'])
```

```
In [69]: data[['Page Views Per Visit']].fillna(data[['Page Views Per Visit']].mode()[0],inplace=True)
```

```
In [70]: data[['Last Activity']].value_counts()
```

```
Out[70]: Email Opened      3437
SMS Sent                2745
Olark Chat Conversation  973
Page Visited on Website 640
Converted to Lead       428
Email Bounced          326
Email Link Clicked     267
Form Submitted on Website 116
Unreachable            93
Unsubscribed           61
Had a Phone Conversation 30
Approached upfront      9
View in browser link Clicked 5
Email Received         2
Email Marked Spam      2
Visited Booth in Tradeshow 1
Resubscribed to emails  1
Name: Last Activity, dtype: int64
```

```
In [71]: data[['Last Activity']].fillna(data[['Last Activity']].mode()[0],inplace=True)
```

```
In [72]: data[['Asymmetrique Profile Index']].value_counts()
```

```
Out[72]: 02.Medium      2738
01.High             2131
03.Low              100
Name: Asymmetrique Profile Index, dtype: int64
```

```
In [73]: data[['Asymmetrique Profile Index']].fillna(data[['Asymmetrique Profile Index']].mode()[0],inplace=True)
```

```
In [74]: data.isna().sum()
```

```
Out[74]: Prospect ID      0
Lead Number        0
Lead Origin        0
Lead Source        0
Do Not Email      0
Do Not Call       0
Converted          0
TotalVisits       0
Page Views Per Visit 0
Time Spent on Website 0
Last Activity      0
Country           0
Specialization    0
What is your current occupation 0
Search           0
Magazine         0
Newspaper Article 0
X Education Forums 0
Newspaper        0
Digital Advertisement 0
Through Recommendations 0
Receive More Updates About Our Courses 0
Update me on Supply Chain Content 0
Get updates on DM Content 0
Lead Quality      0
Lead Profile      0
City              0
Asymmetrique Activity Index 0
Asymmetrique Profile Index 0
Asymmetrique Activity Score 0
Asymmetrique Profile Score 0
I agree to pay the amount through cheque 0
A free copy of Mastering The Interview 0
Last Notable Activity 0
Nigeria Career Prospects 0
Tags_DeadEnd     0
Tags_UnclearLeads 0
Tags_ReasonedLeads 0
dtype: int64
```

```
In [75]: data.head()
```

```
Out[75]:
```

	Prospect ID	Lead Number	Lead Origin	Lead Source	Do Not Email	Do Not Call	Converted	TotalVisits	Total Time Spent on Website	Page Views Per Visit	Last Activity	Country	Specialization	What is your current occupation
0	79272cf-8ba4-429b-bf6c-beaf6e20	660737	API	Olark Chat	0	0	0	0.0	0	0.0	Page Visited on Website	India	Select	Unemployment
1	2a272436-5132-4136-bbf6-dcc8b88f4620	660728	API	Organic Search	0	0	0	5.0	674	2.5	Email Opened	India	Select	Unemployment
2	fcdcd111-4219-435d-a223-fdd2659dbda8	660727	Landing Page Submission	Direct Traffic	0	0	1	2.0	1532	2.0	Email Opened	India	Business Administration	Unemployment
3	0cc2f4b8-7fd4-4a39-8949-1979f70538cc	660719	Landing Page Submission	Direct Traffic	0	0	0	1.0	305	1.0	Unreachable	India	Media and Advertising	Unemployment
4	325e9628-e534-4526-9463-4a8b88782852	660681	Landing Page Submission	Google	0	0	1	2.0	1428	1.0	Converted to Lead	India	Select	Unemployment

```
In [76]: data[['Country']].value_counts()
```

```
Out[76]: India      9553
United States   69
United Arab Emirates 21
Singapore       24
Saudi Arabia    53
United Kingdom  15
Australia       13
Qatar           10
Bahrain         7
Hong Kong       7
France          6
Oman            5
unknown         5
Germany         4
South Africa    4
Canada          4
Kuwait          4
Nigeria         3
Sweden          3
Philippines     2
Uganda          2
Ghana           2
Asia/Pacific Region 2
China           2
Bangladesh      2
Italy            2
Malaysia         1
Liberia          1
Vietnam          1
Sri Lanka       1
Tanzania         1
Indonesia        1
Denmark          1
Malaysia         1
Switzerland      1
Russia           1
Kenya           1
Name: Country, dtype: int64
```

```
In [77]: data[['Country_India']] = data[['Country']].apply(lambda x: 1 if x=='India' else 0)
```

```
In [78]: data.drop('Country',inplace=True,axis=1)
```

```
In [79]: data.head()
```

```
Out[79]:
```

	Prospect ID	Lead Number	Lead Origin	Lead Source	Do Not Email	Do Not Call	Converted	TotalVisits	Total Time Spent on Website	Page Views Per Visit	Last Activity	Specialization	What is your current occupation
0	79272cf-8ba4-429b-bf6c-beaf6e20	660737	API	Olark Chat	0	0	0	0.0	0	0.0	Page Visited on Website	Select	Unemployment
1	2a272436-5132-4136-bbf6-dcc8b88f4620	660728	API	Organic Search	0	0	0	5.0	674	2.5	Email Opened	Select	Unemployment
2	fcdcd111-4219-435d-a223-fdd2659dbda8	660727	Landing Page Submission	Direct Traffic	0	0	1	2.0	1532	2.0	Email Opened	Business Administration	Unemployment
3	0cc2f4b8-7fd4-4a39-8949-1979f70538cc	660719	Landing Page Submission	Direct Traffic	0	0	0	1.0	305	1.0	Unreachable	Media and Advertising	Unemployment
4	325e9628-e534-4526-9463-4a8b88782852	660681	Landing Page Submission	Google	0	0	1	2.0	1428	1.0	Converted to Lead	Select	Unemployment

For categorical variables with multiple levels, create dummy features (one-hot encoded)

```
In [80]: dummy=pd.get_dummies(data[['Lead Origin','Lead Source','What is your current occupation','Better Career Prospects','Lead Quality','Asymmetrique Activity Index','Asymmetrique Profile Index','Lead Profile','City']],drop_first=True)
```

```
In [81]: #['Last Activity','','Lead Profile','Tags','Specialization'],
```

```
Out[81]: dummy.head()
```

```
Out[82]:
```

	Better Career Prospects	Lead Origin_Landing Page Submission	Lead Source_Organic Search	Lead Source_Google	Lead Source_Olark Chat	Lead Source_Referal Sites	Lead Source_Welinkak Website	Lead Source_Referral Sites	Lead Source_Others	Sc
0	1	0	0	0	0	0	0	1	0	0
1	1	0	0	0	0	0	0	0	1	0
2	1	1	0	0	0	0	0	0	0	0
3	1	1	0	0	0	0	0	0	0	0
4	1	1	1	0	0	0	1	0	0	0

```
In [83]: data=pd.concat([data,dummy],axis=1)
```

```
In [84]: data.head()
```

```
Out[84]:
```

	Prospect ID	Lead Number	Lead Origin	Lead Source	Do Not Email	Do Not Call	Converted	TotalVisits	Total Time Spent on Website	Page Views Per Visit	Last Activity	Specialization	What is your current occupation
0	79272cf-8ba4-429b-bf6c-beaf6e20	660737	API	Olark Chat	0	0	0	0.0	0	0.0	Page Visited on Website	Select	Unemployment
1	2a272436-5132-4136-bbf6-dcc8b88f4620	660728	API	Organic Search	0	0	0	5.0	674	2.5	Email Opened	Select	Unemployment
2	fcdcd111-4219-435d-a223-fdd2659dbda8	660727	Landing Page Submission	Direct Traffic	0	0	1	2.0	1532	2.0	Email Opened	Business Administration	Unemployment
3	0cc2f4b8-7fd4-4a39-8949-1979f70538cc	660719	Landing Page Submission	Direct Traffic	0	0	0	1.0	305	1.0	Unreachable	Media and Advertising	Unemployment
4	325e9628-e534-4526-9463-4a8b88782852	660681	Landing Page Submission	Google	0	0	1	2.					


```
In [157]: numbers = [float(i)/10 for x in range(10)]
for i in numbers:
    y_train_pred_final[i] = y_train_pred_final.Converted_Prob.map(lambda x: 1 if x > 0.3 else 0)
y_train_pred_final.head()
```

```
Out[157]:
```

	Converted	Converted_Prob	predicted	0.0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9
3555	1	0.980599	1	1	1	1	1	1	1	1	1	1	1
8055	1	0.993261	1	1	1	1	1	1	1	1	1	1	1
1843	0	0.152428	0	1	1	0	0	0	0	0	0	0	0
1763	0	0.007727	0	1	0	0	0	0	0	0	0	0	0
3017	1	0.564242	1	1	1	1	1	1	1	0	0	0	0

Determining optimum probability cut-off

```
In [158]: cutoff_df = pd.DataFrame( columns = ['prob','accuracy','sensi','speci'])
from sklearn.metrics import confusion_matrix

# TP = confusion[1,1] # true positive
# TN = confusion[0,0] # true negatives
# FP = confusion[0,1] # false positives
# FN = confusion[1,0] # false negatives

num = [0,0,0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9]
for i in num:
    cml = metrics.confusion_matrix(y_train_pred_final.Converted, y_train_pred_final[i] )
    total=sum(sum(cml))
    accuracy = (cml[0,0]+cml[1,1])/total

    speci = cml[0,0]/(cml[0,0]+cml[0,1])
    sensi = cml[1,1]/(cml[1,0]+cml[1,1])
    cutoff_df.loc[i] = [i ,accuracy,sensi,speci]
print(cutoff_df)
```

prob	accuracy	sensi	speci	
0.0	0.0	0.389045	1.0000	0.000000
0.1	0.1	0.578120	0.9864	0.318136
0.2	0.2	0.751945	0.9856	0.666836
0.3	0.3	0.805322	0.8272	0.791391
0.4	0.4	0.836446	0.7868	0.868039
0.5	0.5	0.834734	0.7292	0.903936
0.6	0.6	0.824619	0.6488	0.936577
0.7	0.7	0.811080	0.5632	0.968925
0.8	0.8	0.783224	0.4696	0.982934
0.9	0.9	0.753929	0.3925	0.992613

```
In [159]: cutoff_df.plot.line(x='prob', y=['accuracy','sensi','speci'])
plt.show()
```

Note:

At the Probability of approx. 0.3, the sensitivity and specificity curves cross each other. This gives us the optimum cut-off point. We can proceed calculating the metrics at this point.

```
In [160]: y_train_pred_final['final_predicted'] = y_train_pred_final.Converted_Prob.map(lambda x: 1 if x > 0.3 else 0)
y_train_pred_final.head()
```

```
Out[160]:
```

	Converted	Converted_Prob	predicted	0.0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	final_predicted
3555	1	0.980599	1	1	1	1	1	1	1	1	1	1	1	1
8055	1	0.993261	1	1	1	1	1	1	1	1	1	1	1	1
1843	0	0.152428	0	1	1	0	0	0	0	0	0	0	0	0
1763	0	0.007727	0	1	0	0	0	0	0	0	0	0	0	0
3017	1	0.564242	1	1	1	1	1	1	1	0	0	0	0	1

```
In [161]: metrics.accuracy_score(y_train_pred_final.Converted, y_train_pred_final.final_predicted)

Out[161]: 0.805322188515406

In [162]: confusion2 = metrics.confusion_matrix(y_train_pred_final.Converted, y_train_pred_final.final_predicted)
confusion2

Out[162]: array([[13107,   819],
                [ 432, 2068]], dtype=int64)

In [163]: TP = confusion2[1,1] # true positive
TN = confusion2[0,0] # true negatives
FP = confusion2[0,1] # false positives
FN = confusion2[1,0] # false negatives

In [164]: # Sensitivity
TP / float(TP+FN)

Out[164]: 0.8272

In [165]: # Specificity
TN / float(TN+FP)

Out[165]: 0.7913907284768212

In [166]: print(FP/ float(TN+FP))

0.20860927152317882

In [167]: print (TP / float(TP+FP))

0.7163145133356426

In [168]: print (TN / float(TN+ FN))

0.8779316191014411

In [169]: # Precision
confusion2[1,1]/(confusion2[0,1]+confusion2[1,1])

Out[169]: 0.7163145133356426

In [170]: # Recall
confusion2[1,1]/(confusion2[1,0]+confusion2[1,1])

Out[170]: 0.8272

In [171]: from sklearn.metrics import precision_score, recall_score

In [172]: precision_score(y_train_pred_final.Converted, y_train_pred_final.final_predicted)

Out[172]: 0.7163145133356426

In [173]: recall_score(y_train_pred_final.Converted, y_train_pred_final.final_predicted)

Out[173]: 0.8272

In [174]: X_test[['Total Time Spent on Website','Page Views Per Visit','TotalVisits','Asymmetrique Activity Score',
'Asymmetrique Profile Score']]=scaler.transform(X_test[['Total Time Spent on Website','Page Views Per Visit',
'TotalVisits','Asymmetrique Activity Score','Asymmetrique Profile Score']])

In [175]: X_test=X_test[X_train_4.columns]

In [176]: X_test_sm=sm.add_constant(X_test)

In [177]: y_test_pred=real.predict(X_test_sm)
```

Task 2: Probability Score for the potential leads

```
In [178]: y_test_pred_final=pd.DataFrame({'Converted':y_test,'Converted_prob':y_test_pred})
y_test_pred_final

Out[178]:
```

	Converted	Converted_prob
6187	0	0.569190
7610	0	0.178430
7084	0	0.254505
491	0	0.143002
4222	0	0.728445
...
7246	0	0.094872
2727	1	0.995274
501	0	0.129145
8330	1	0.999239
506	0	0.071211

2755 rows × 2 columns

Task 3: Hot Leads among the potential leads

```
In [179]: y_test_pred_final['final_predicted'] = y_test_pred_final.Converted_prob.map(lambda x: 1 if x > 0.3 else 0)

In [180]: y_test_pred_final

Out[180]:
```

	Converted	Converted_prob	final_predicted
6187	0	0.569190	1
7610	0	0.178430	0
7084	0	0.254505	0
491	0	0.143002	0
4222	0	0.728445	1
...
7246	0	0.094872	0
2727	1	0.995274	1
501	0	0.129145	0
8330	1	0.999239	1
506	0	0.071211	0

2755 rows × 3 columns

```
In [181]: ## Accuracy
metrics.accuracy_score(y_test_pred_final.Converted, y_test_pred_final.final_predicted)

Out[181]: 0.8010889292196007

In [182]: confusion2 = metrics.confusion_matrix(y_test_pred_final.Converted, y_test_pred_final.final_predicted )
confusion2

Out[182]: array([[1162,   356],
                [ 192,  845]], dtype=int64)

In [183]: TP = confusion2[1,1] # true positive
TN = confusion2[0,0] # true negatives
FP = confusion2[0,1] # false positives
FN = confusion2[1,0] # false negatives

In [184]: # Sensitivity
TP / float(TP+FN)

Out[184]: 0.8148505303760849

In [185]: # Specificity
TN / float(TN+FP)

Out[185]: 0.7927823050058207

In [186]: precision_score(y_test_pred_final.Converted, y_test_pred_final.final_predicted)

Out[186]: 0.7035803497085762

In [187]: recall_score(y_train_pred_final.Converted, y_train_pred_final.final_predicted)

Out[187]: 0.8272
```

Assignment Subjective Questions:

1. Which are the top three variables in your model which contribute most towards the probability of a lead getting converted?

Lead Source_Weingak Website
Lead Origin_Lead Add Form
What is your current occupation_Working Professional

1. What are the top 3 categorical/dummy variables in the model which should be focused the most on in order to increase the probability of lead conversion?

Lead Source_Weingak Website
Lead Origin_Lead Add Form
What is your current occupation_Working Professional

1. X Education has a period of 2 months every year during which they hire some interns. The sales team, in particular, has around 10 interns allotted to them. So during this phase, they wish to make the lead conversion more aggressive. So they want almost all of the potential leads to them. So during this phase, they wish to make the lead conversion more aggressive. So they want almost all of the potential leads to them. (i.e. the customers who have been predicted as 1 by the model) to be converted and hence, want to make phone calls to as much of those people as possible. Suggest a good strategy they should employ at this stage.

This is essentially a peak calling season for X Education and thus to enhance as much mileage as possible and maximize on the caller availability, all the people have been predicted by the model (probability > 0.3) can be called. To increase probability of getting conversions early-on in the process, the list can be prioritized based on the scores and can be called in the descending order of scores.

1. Similarly, at times, the company reaches its target for a quarter before the deadline. During this time, the company wants the sales team to focus on some new work as well. So during this time, the company's aim is to not make phone calls unless it's extremely necessary, i.e. they want to minimize the rate of useless phone calls. Suggest a strategy they should employ at this stage.

This is a sensitive calling season for X Education and thus to be as efficient as possible by maximizing the conversion rate, the callers can be asked to focus on people with most contributions. Based on domain knowledge, we can focus on top 3 or top 5 people. As a matter of fact, The model has identified 3 specific types of Leads:

1. Leads who have come from the Weingak website, **Lead Source_Weingak Website**
2. Leads who have been identified as a lead by triggering the 'Lead Add Form' user action, **Lead Origin_Lead Add Form**
3. Leads who are working professionals and possibly have a specific reason in their minds for the course and perhaps lesser monetary challenges in taking up a course. **What is your current occupation_Working Professional**

The actual focus group of leads for the callers can be decided based on a detailed business discussion backed by the model statistics and business accumen.

Note:

The below three variables are the most significant influencers for the model.

Lead Quality_Worst
Tags_ReasonedLeads (These are leads whom the sales people have tagged as have been reasoned by the lead during a sales call)
Lead Source_Weingak Website

We have to note that the first two are negatively influencing the conversion. Essentially, these are people who are more likely will not convert.