
LEAD SCORING CASE STUDY

SUBMITTED BY:
ABHINAV CHOUDHARY & AKSHAY SAWANT

PROBLEM STATEMENT

An education company named X Education sells online courses to industry professionals.

The company gets a lot of leads, some of the leads get converted while most do not. The typical lead conversion rate of the company is around 30%.

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.

BUSINESS OBJECTIVES

- The business objective of this case study is **to identify the most promising leads**, that is the leads that are most likely to convert into paying customers for the X Education company.
- By successfully identifying the promising leads, the lead conversion rate would go up for the company as the sales team will focus more on communicating with the potential leads rather than making calls to everyone.
- We need **to build a model and assign a lead score to each of the leads** such that the customers with a higher lead score have a higher chance of conversion than the customers with a lower lead score.
- The ballpark target given by the company's CEO is **to get the lead conversion rate to be around 80%**.
- Also, the model should be able to adjust if the company's requirement changes in the near future.

ANALYSIS APPROACH

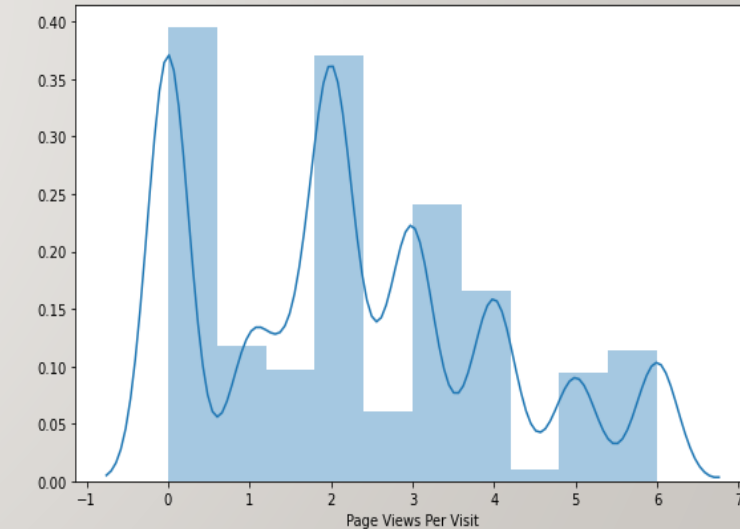
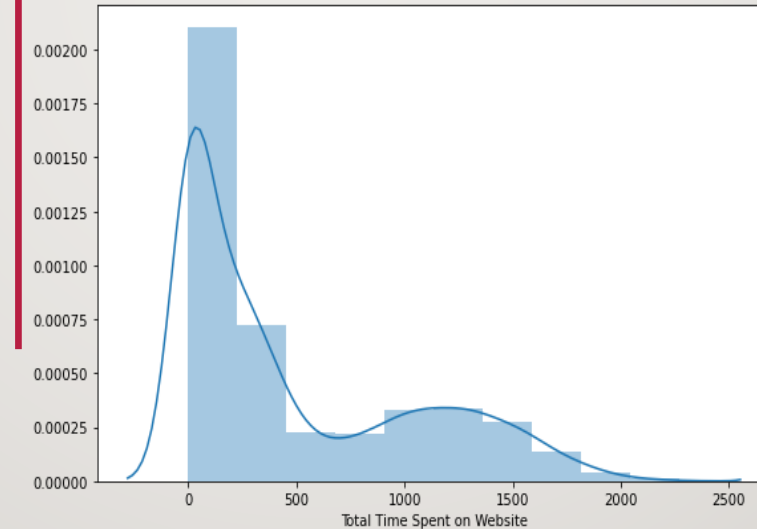
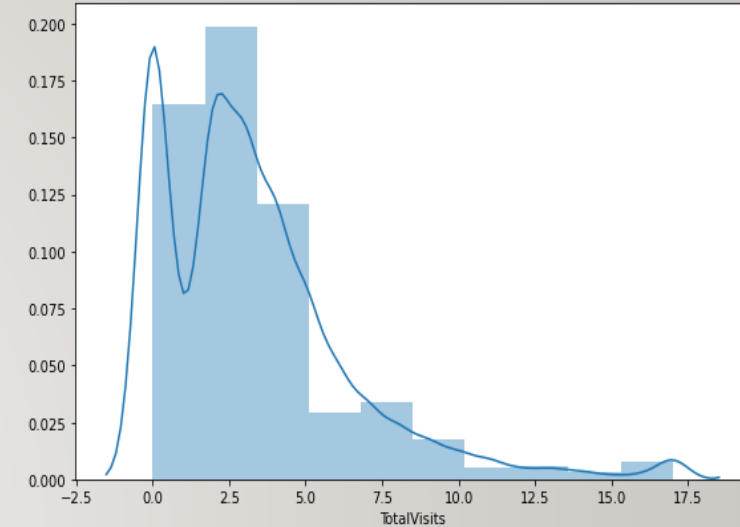
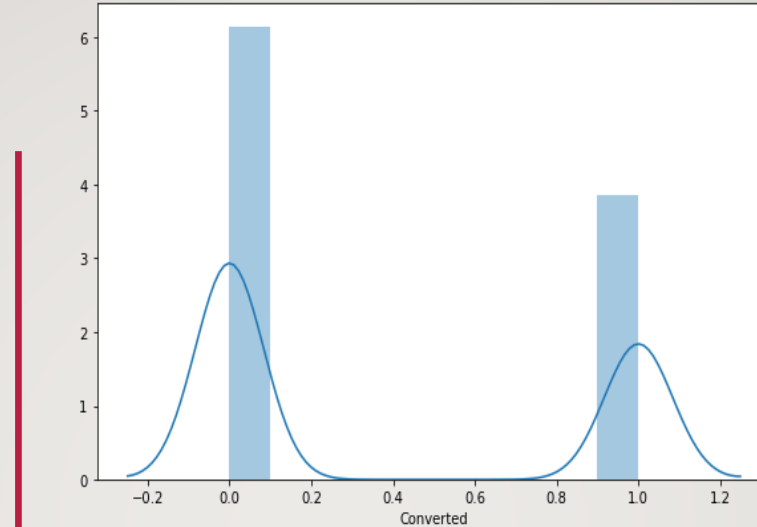
1. **Understanding & Cleaning the Data:** We have used the Jupyter Python Notebook to load and understand the Leads datasets. We have gone through the breadth & the depth of the features present in the datasets along with their definition to assess the data quality & its spread at a high level.
As part of the Data Cleaning process, we had found & treated all the irregularities in the dataset such as Missing Values; Outliers; Skewed Categorical features & Invalid data points.
2. **Exploratory Data Analysis:** After cleaning the data, we performed various types of Univariate, Bivariate & Multivariate analysis by plotting appropriate graphs with respect to the Target variable. This helped us to draw relevant insights & correlations present within the dataset.
3. **Data Preparation:** Here, we had firstly created the dummy variables of all the categorical features in the dataset. Then, we had split the dataset between training & test sets & after that, we performed the Standard Scaling of the independent features.
4. **Model Building:** Since the count of features were high, we first started with RFE to eliminate redundant features & then built the first Model. Over multiple iterations of refinement (through p-value & VIF), we concluded this step with a final model.
5. **Model Evaluation:** Here, with the final model, we first obtained the lead score and plotted the ROC Curve. After that, we determined the optimal cut-off to proceed further with the Evaluations Metrics.
6. **Making Predictions:** After evaluating the final model, we ran the model over the test dataset to make predictions & then, we reviewed the predicted results with the actual records. We finally concluded with our analysis findings & recommendations to the business.

EDA – UNIVARIATE & MULTIVARIATE ANALYSIS

NUMERIC FEATURES ANALYSIS:

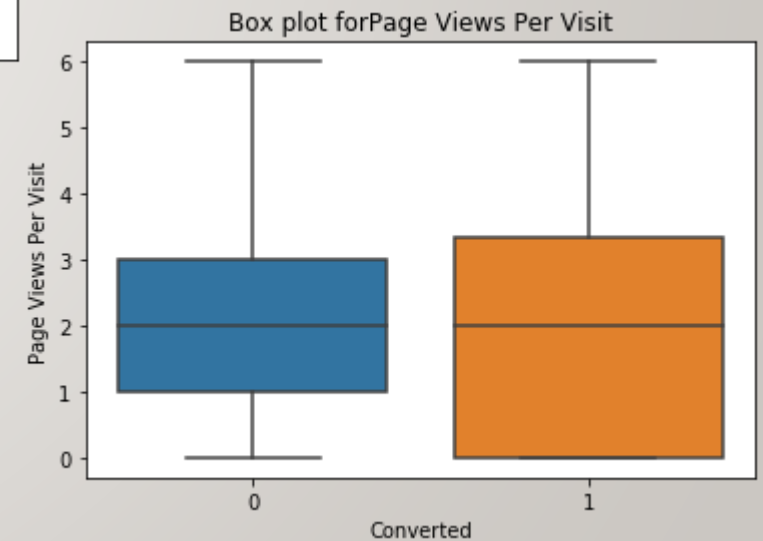
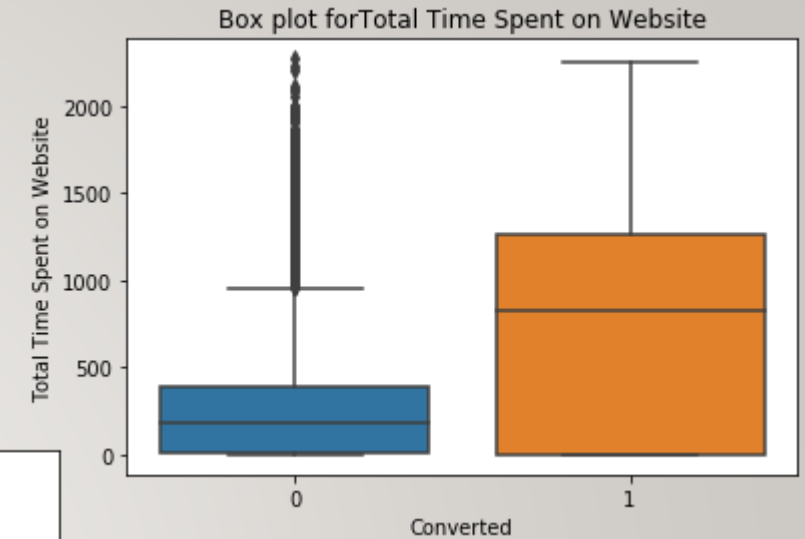
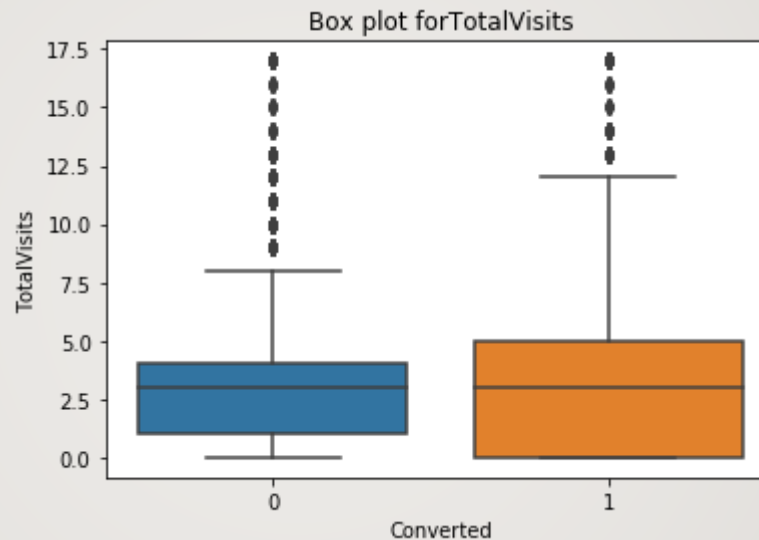
- After cleaning the data, we moved to the EDA for further insights.
- From the distribution plot(at the right), we can see that the spread of Total Visits & Total Time Spent per Visit both are right/positive skewed.
- It seems that most of the leads are not frequent visitors & their average engagement(time spent & page views per visit) on the website is slightly on a lower side.
- Here, the Class Imbalance of the Target Variable(Converted) is = **1.6**
This seems fine and we are good to proceed with further analysis.

Let's have a look at the spread w.r.t. converted vs. non-converted leads in the next slide...



NUMERIC FEATURES ANALYSIS:

- Here are the box plots w.r.t. converted vs. non-converted leads.
- It is evident that leads who got converted have visited & engaged more on the website & therefore have higher no. of visits & time spent on website compared to the non-converted leads
- For 'Page Views per Visit' metric, we can say that it is slightly better for converted leads. But otherwise, there doesn't seem to be a major difference we can notice for this metric between converted & non-converted leads.



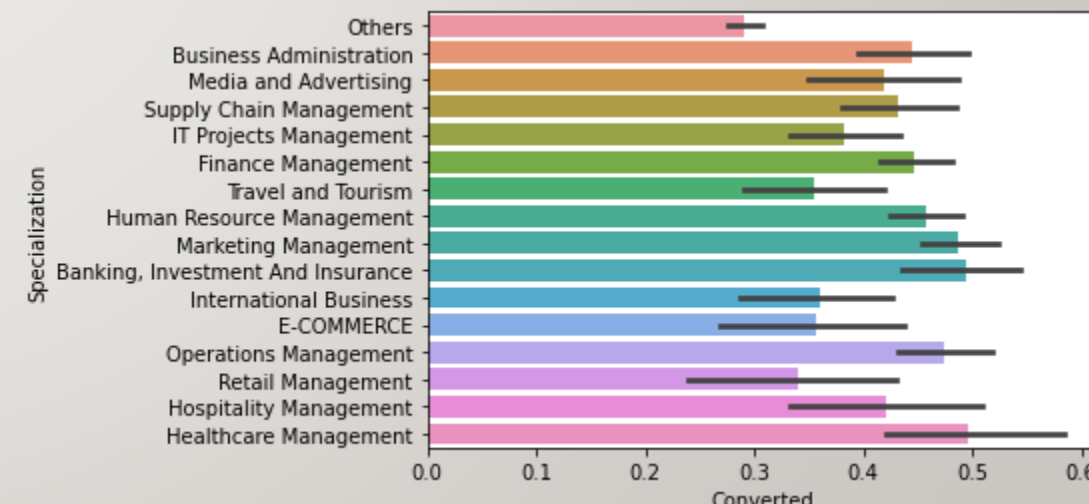
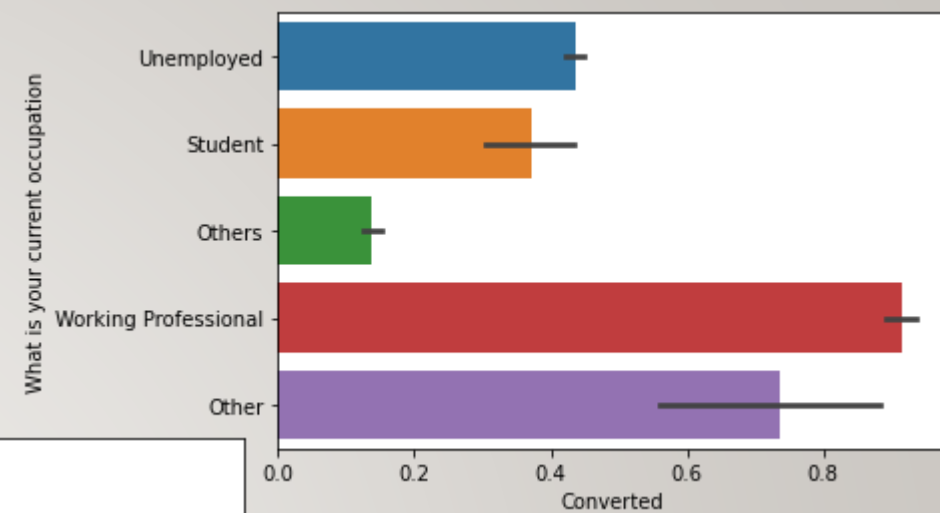
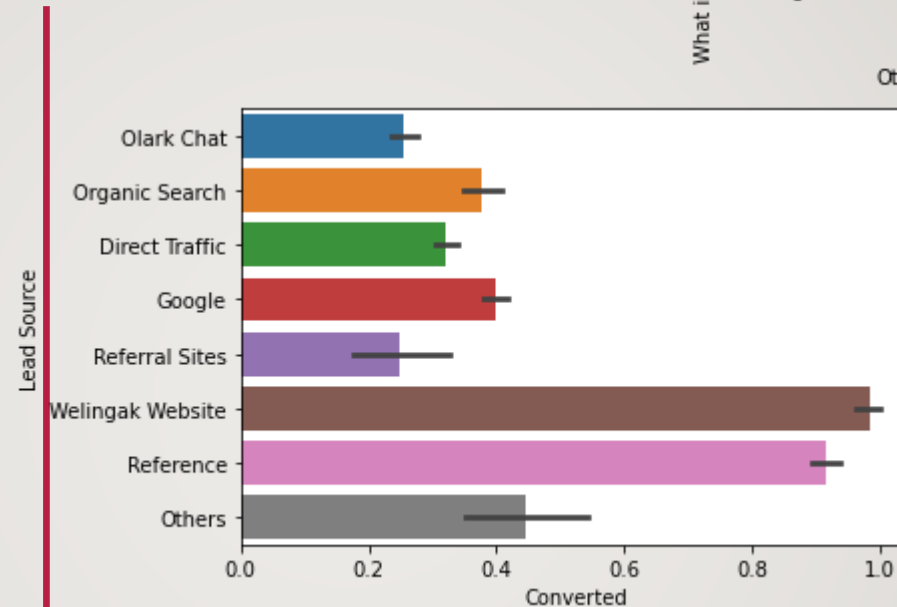
NUMERIC FEATURES ANALYSIS:

- In terms of correlation, from the correlation matrix, we have noticed a decent correlation of 0.7 between 'Total Visits' & 'Page Views per Visit' features. Otherwise, there isn't any major correlations we could see between the numeric variables.



CATEGORICAL FEATURES ANALYSIS:

- Here, we plotted each categorical feature with respect to their respective conversion rate.
- It has been observed that 'Working Professions' leads have higher conversion rate than anyone else.
- Leads who have come from 'Welingak website' or through some 'Reference' had relatively higher conversion rate
- Leads who have their specialization in 'Healthcare Management' or 'Banking Investments & Insurance' have a higher chances of conversion.



A rectangular sign with a thick black border and a white inner border. The sign has a solid red background. Centered on the red background is the text "MODEL BUILDING" in white, uppercase, sans-serif font.

MODEL BUILDING

MODEL BUILDING :

- We had built a model with all the features included and found that there were many insignificant variables present in our model.
- With the help of RFE, we reduced the no. of features to 15 most significant features & built the model again.
- In the next model, from the summary stats, we decided to drop the insignificant feature with high p-value.
- In the 3rd iterations of model building & feature selection, we concluded with a statistically significant & most stable model with 14 statistically significant features. This model has all the p-values & VIFs within the permissible range.

Generalized Linear Model Regression Results

Dep. Variable:	Converted	No. Observations:	6468
Model:	GLM	Df Residuals:	6453
Model Family:	Binomial	Df Model:	14
Link Function:	logit	Scale:	1.0000
Method:	IRLS	Log-Likelihood:	-2682.5
Date:	Sun, 07 Mar 2021	Deviance:	5365.1
Time:	19:28:40	Pearson chi2:	8.50e+03
No. Iterations:	7		
Covariance Type:	nonrobust		

	Features	VIF
0	const	5.04
2	Lead Source_Olark Chat	1.45
14	What is your current occupation_Working Profes...	1.34
13	What is your current occupation_Unemployed	1.32
8	Last Activity_Olark Chat Conversation	1.30
1	Total Time Spent on Website	1.25
3	Lead Source_Reference	1.14
9	Last Activity_SMS Sent	1.12
5	Do Not Email_Yes	1.09
10	Last Activity_Unsubscribed	1.07
6	Last Activity_Converted to Lead	1.06
12	What is your current occupation_Student	1.05
4	Lead Source_Welingak Website	1.03
7	Last Activity_Had a Phone Conversation	1.01
11	Specialization_Hospitality Management	1.01

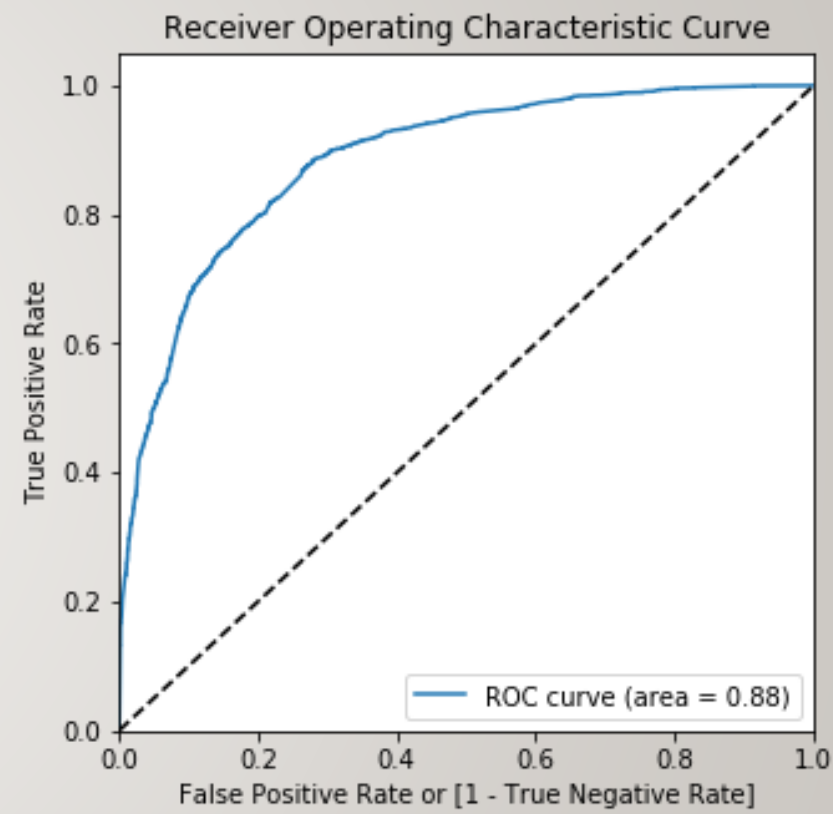
	coef	std err	z	P> z	[0.025	0.975]
const	-2.2147	0.087	-25.354	0.000	-2.386	-2.044
Total Time Spent on Website	1.0594	0.039	27.116	0.000	0.983	1.136
Lead Source_Olark Chat	1.2232	0.103	11.872	0.000	1.021	1.425
Lead Source_Reference	3.4195	0.203	16.832	0.000	3.021	3.818
Lead Source_Welingak Website	5.2346	0.724	7.226	0.000	3.815	6.654
Do Not Email_Yes	-1.4200	0.168	-8.428	0.000	-1.750	-1.090
Last Activity_Converted to Lead	-1.2426	0.219	-5.674	0.000	-1.672	-0.813
Last Activity_Had a Phone Conversation	2.1693	0.676	3.208	0.001	0.844	3.495
Last Activity_Olark Chat Conversation	-1.2312	0.165	-7.445	0.000	-1.555	-0.907
Last Activity_SMS Sent	1.2195	0.074	16.481	0.000	1.075	1.365
Last Activity_Unsubscribed	1.1960	0.461	2.593	0.010	0.292	2.100
Specialization_Hospitality Management	-0.8886	0.319	-2.786	0.005	-1.514	-0.264
What is your current occupation_Student	1.2243	0.235	5.210	0.000	0.764	1.685
What is your current occupation_Unemployed	1.1323	0.085	13.352	0.000	0.966	1.298
What is your current occupation_Working Professional	3.6524	0.197	18.540	0.000	3.266	4.038

MODEL EVALUATION

MODEL EVALUATION : LEAD SCORE & ROC CURVE

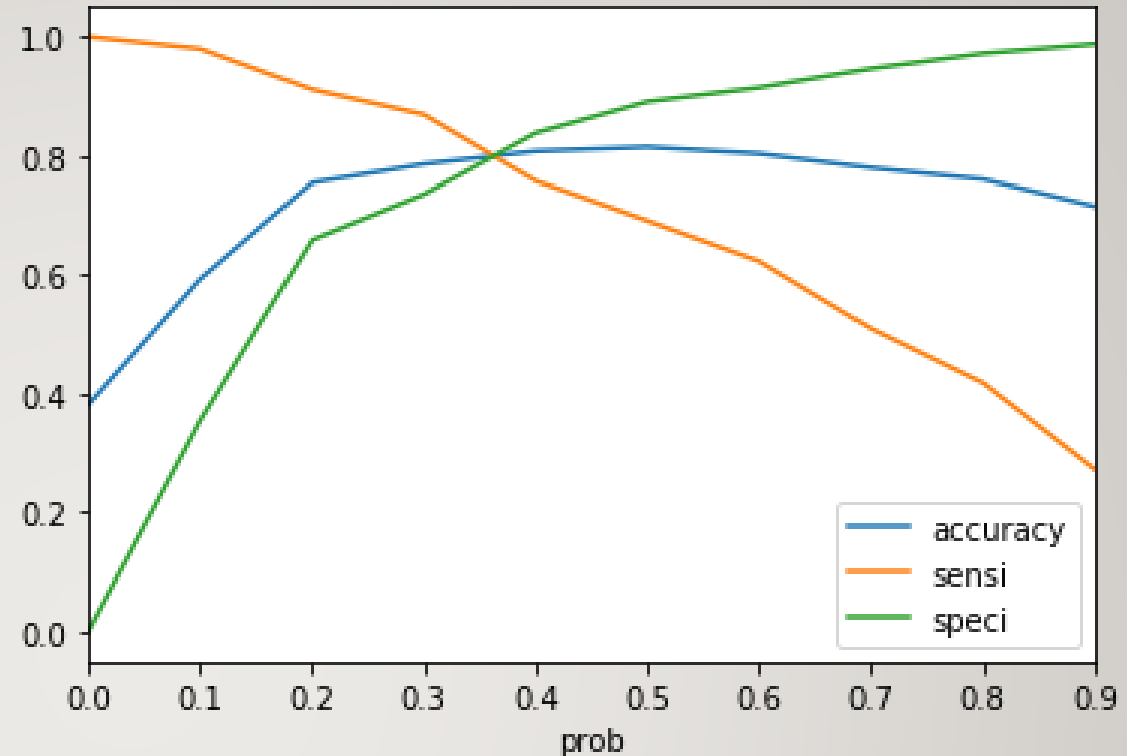
- After building the final model, we predicted the probability of getting converted on our train dataset, which is termed as the “**Lead_Score**”(varies between 0 to 1). With this, we later determined the optimal probability cut-off to determine if the leads were converted or not.
- From the final model, after making predictions on it(on train dataset), we created the **ROC Curve** to determine the model stability with **AUC Score** (Area Under the Curve). As we can see from the graph plotted on the right side, this score is 0.88 which is a great score & represents quite stable & reliable model.
- In other words, as we see the ROC Curve is leaned towards the left side of the border, this means that the performance of our final model would be great.

	Converted	Lead_Score
0	0	0.310625
1	0	0.254190
2	0	0.356241
3	0	0.829693
4	0	0.310625



MODEL EVALUATION : FINDING THE OPTIMAL CUT-OFF

- In order to find the most optimal cut-off, we had plotted the graph between 'Accuracy', 'Sensitivity', and 'Specificity' at different probability/lead_score values.
- From this plot(on the right), we looked at the intersection point of accuracy, sensitivity and specificity which came out **to be at 0.35**, where all the score are in a close range which is the ideal point to select and hence it was selected.
- Therefore, the final probability/lead_score cut-off value that we decided was **0.35**
- With this cut-off value, we had created the **Confusion Matrix**(see table at the right) to check the Accuracy, Sensitivity, Specificity, Precision & Recall of the model



Confusion Matrix on Training Dataset

(Cut-off = 0.35) ACTUAL	PREDICTED	
	0 (Non-Converted)	1 (Converted)
0 (Non-Converted)	3228	774
1 (Converted)	520	1946

MODEL EVALUATION : EVALUATION METRICS

- With the probability cut-off value of 0.35 & the Confusion Matrix, we calculated the Evaluation Metrics here. >>
- **Precision and Recall** have a very important role in model evaluation & business decision making as it tells how our model will behave on unknown datasets.
- Our one of the Business Objectives was **to achieve a Recall of 80%** which means that the business wanted most of the hot leads to be identified so that the sales team can take appropriate actions to convert those hot leads.
- Therefore, our final model is now apt enough to identify all such hot leads for the sales team.

EVALUATION METRICS	SCORE <small>(rounded)</small>
Accuracy	80%
Sensitivity	79%
Specificity	81%
Precision	72%
Recall	79%

MAKING PREDICTIONS

MAKING PREDICTIONS :

- After evaluating the final model, we ran the model over the test dataset to make predictions & then, we reviewed the predicted results with the actual records.
- Here, again we evaluated the model on the test dataset with the help of the Evaluation Metrics.
- The table at the right has these metrics scores.
- **The results shows that our model is very much stable even on unknown datasets.**

EVALUATION METRICS	SCORE <small>(rounded)</small>
Accuracy	80%
Precision	73%
Recall	80%

CONCLUSION & KEY RECOMMENDATIONS

CONCLUSION & KEY RECOMMENDATIONS

1. As we had seen that our model performed equally well on the test dataset as it had on the training dataset. This shows that the model is quite stable and has a very good Accuracy & Recall.
2. Also, by changing the probability cut-off, the model has the ability to adjust with the change in company's requirements in the near future.
3. Below are the top 3 Important Features that the company should focus on to further increase the conversion rate of the leads:
 - i. Lead Source as "Welingak Website"
 - ii. Lead Source as "Reference"
 - iii. Current Occupation as "Working Professional"

THANK YOU 😊