

## **Here is a summary of the key steps and insights:**

### **Data Preparation:**

- The code begins by importing necessary libraries and loading a dataset containing information about leads.
- Exploratory data analysis (EDA) is conducted to understand the dataset, including data types, null values, and unique values in columns.
- Columns with irrelevant or no-value information are identified and dropped.
- Null values and 'Select' values are handled appropriately.
- Binary variables are converted to numeric form, and dummy features are created for categorical variables with multiple levels.

### **Exploratory Data Analysis (EDA):**

- The dataset is explored further, and correlations between variables are visualized using a heatmap.
- Outliers in continuous variables are checked, and object data types are converted to numeric where necessary.

### **Data Splitting and Feature Scaling:**

- The data is split into training and testing sets for model evaluation.
- Feature scaling is applied to relevant features using StandardScaler.

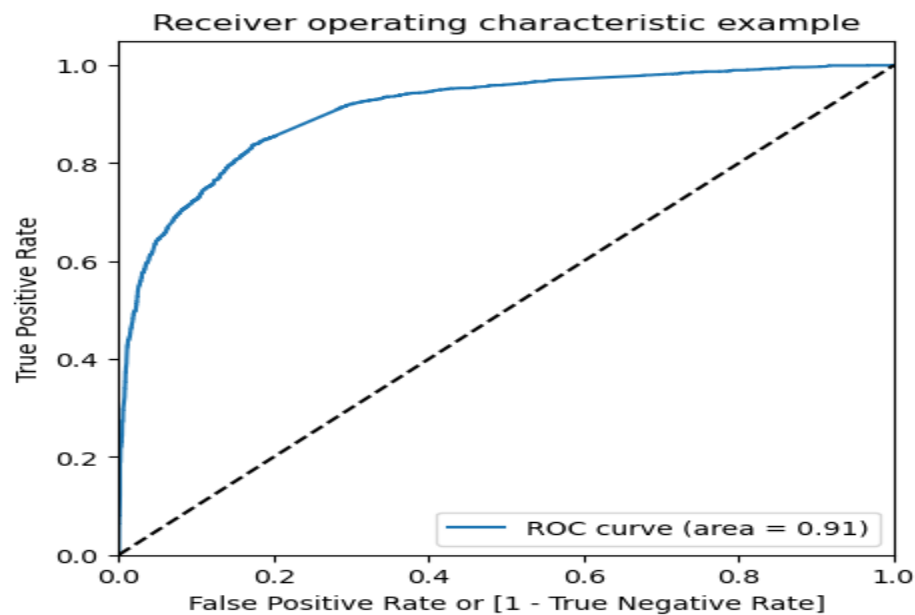
### **Logistic Regression Model Building:**

- An initial logistic regression model is built using the StatsModels library.
- Recursive Feature Elimination (RFE) is employed for feature selection.
- Multicollinearity is checked using Variance Inflation Factor (VIF).

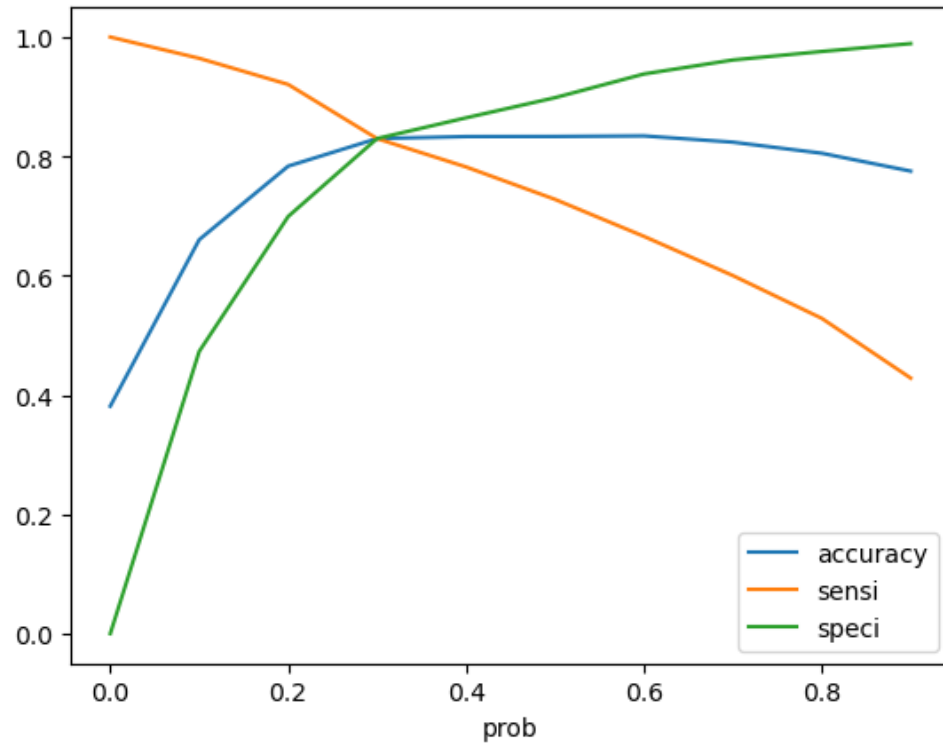
	Features	VIF
1	LeadOrigin_Lead Add Form	1.74
11	LeadQuality_Might be	1.53
13	LastNotableActivity_Modified	1.44
4	LeadSource_Google	1.40
8	What is your current occupation_Working Profes...	1.35
2	LeadSource_Direct Traffic	1.31
9	LeadQuality_High in Relevance	1.29
7	LeadSource_Welingak Website	1.28
10	LeadQuality_Low in Relevance	1.19
0	Total Time Spent on Website	1.17
5	LeadSource_Organic Search	1.13
12	LeadQuality_Worst	1.12
6	LeadSource_Referral Sites	1.02
3	LeadSource_Facebook	1.01

### Model Assessment and Refinement:

- The model's statistical performance is evaluated, including p-values and VIF.
- Variables with high p-values are dropped, and the model is re-evaluated.
- Sensitivity, specificity, and other performance metrics are calculated.
- The ROC curve is plotted to visualize the tradeoff between sensitivity and specificity.

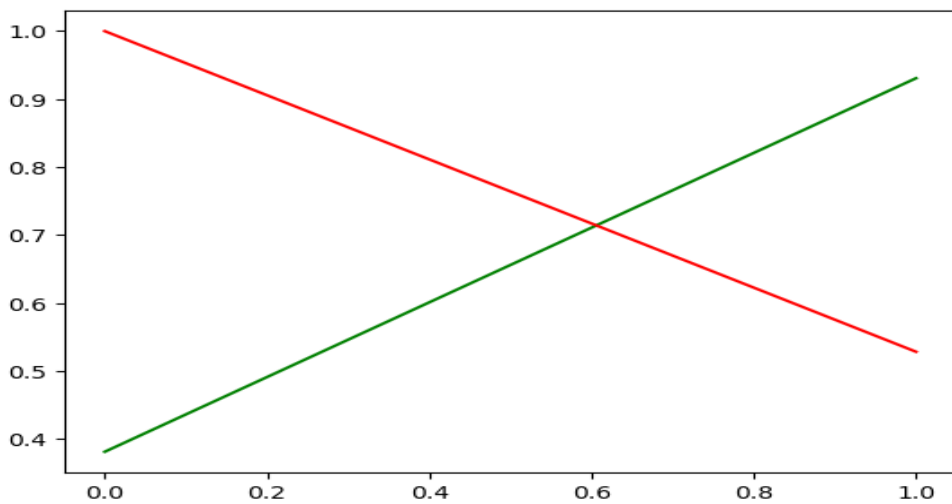


- The optimal cutoff point for predictions is determined.



### Precision and Recall:

- Precision and recall are calculated, and the precision-recall tradeoff is explored.



### Testing on the Test Set:

- The trained model is applied to a separate test set to assess its generalization performance.
- The final model's accuracy metrics-79.44%,

confusion matrix- array([[1642, 35],  
[ 535, 560]], dtype=int64),

sensitivity-51.14%, and

specificity-97.91% .

### Conclusion:

- The important variables for potential buyers are as follows:
  1. The total time spent on website
  2. Total number of visits
  3. When the lead source was- olark chat, wellingak website
  4. When the last activity was- sms, olark chat conversation
  5. When the lead origin is lead add form
  6. When the current occupation was- working professional, student, unemployed, other

X Education can increase all the potential buyers to change their mind and buy their courses.

The Model seems to predict the Conversion Rate very well and we should be able to give the company confidence in making good calls based on this model

Overall, the code serves as a valuable guide for implementing a logistic regression-based lead scoring model and offers insights into various stages of the data science pipeline, from data cleaning to model testing.