

# X education leading case stud

# Problem Statement & Objective of the Study

- X Education gets a lot of leads, its lead conversion rate is very poor at around 30%
- ● X Education wants to make lead conversion process more efficient by identifying the most potential leads, also known as Hot Leads
- ● Their sales team want to know these potential set of leads, which they will be focusing more on communicating rather than making calls to everyone.

# Objective of the Study:

- To help X Education select the most promising leads, i.e., the leads that are most likely to convert into paying customers.
- The company requires us to build a model wherein we need to assign a lead score to each of the leads such that the customers with a higher lead score have a higher conversion chance and the customers with a lower lead score have a lower conversion chance.
- The CEO has given a ballpark of the target lead conversion rate to be around 80%.

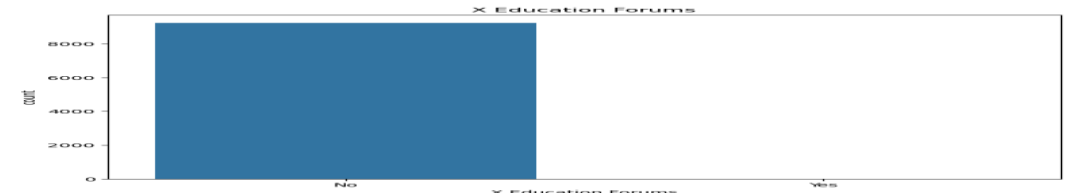
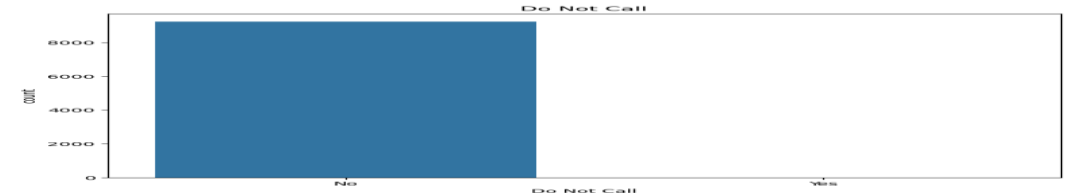
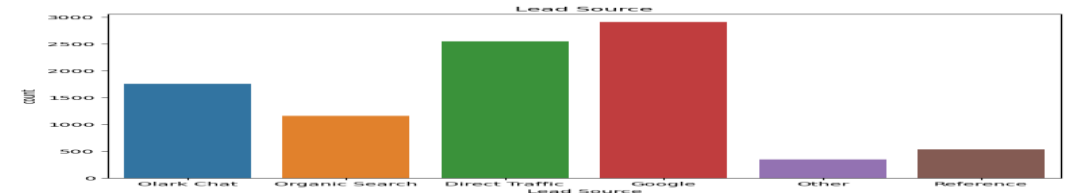
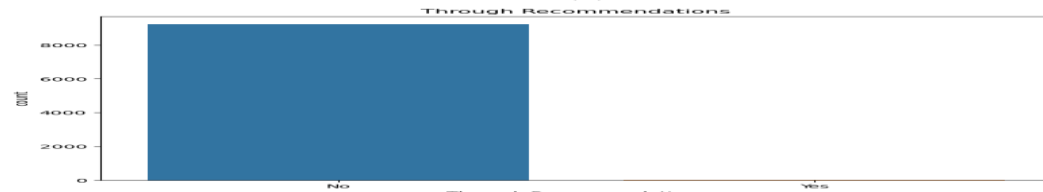
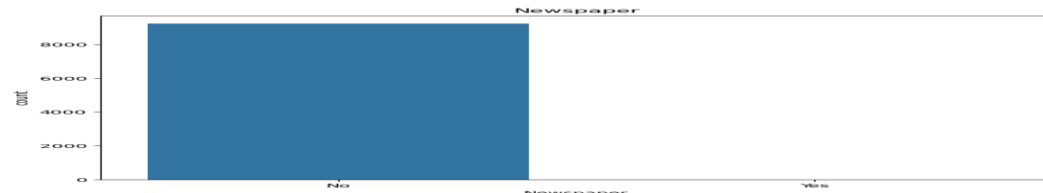
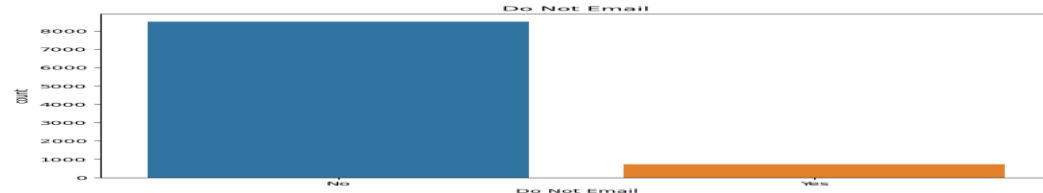
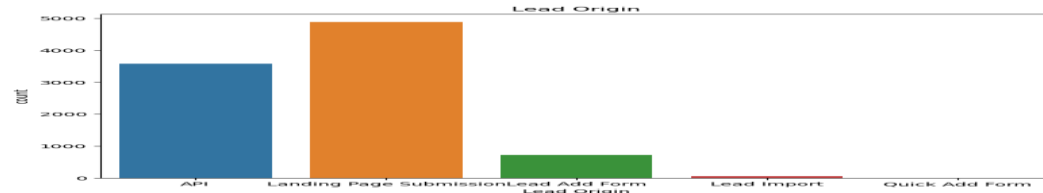
# Data Cleaning

- ❖ "Select" level represents null values for some categorical variables, as customers did not choose any option from the list.
- ❖ Columns with over 40% null values were dropped.
- ❖ Missing values in categorical columns were handled based on value counts and certain considerations.
- ❖ Drop columns that don't add any insight or value to the study objective (tags, country)
- ❖ Imputation was used for some categorical variables.
- ❖ Additional categories were created for some variables.
- ❖ Columns with no use for modeling (Prospect ID, Lead Number) or only one category of response were dropped.
- ❖ Numerical data was imputed with mode after checking distribution.

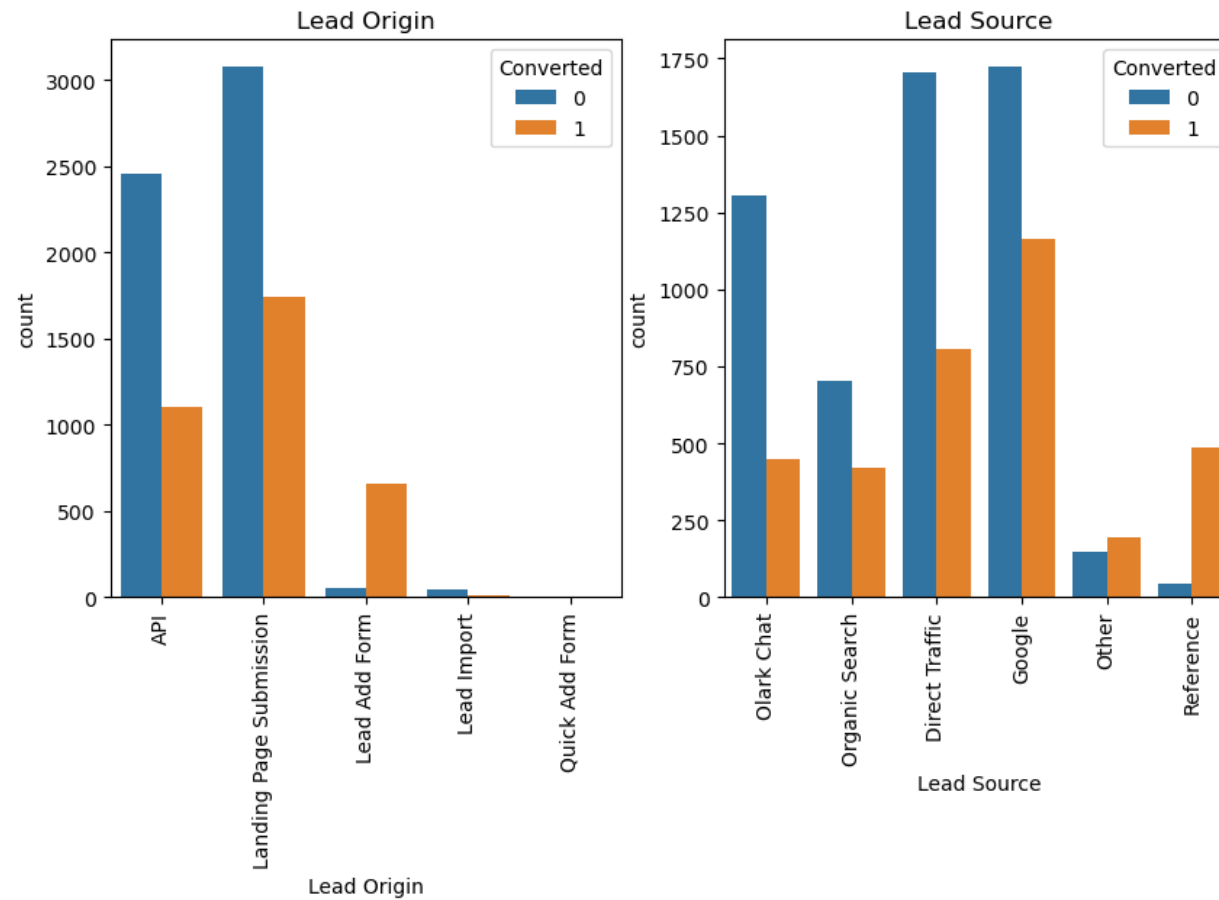
# Data Cleaning

- Skewed category columns were checked and dropped to avoid bias in logistic regression models.
- Outliers in TotalVisits and Page Views Per Visit were treated and capped.
- Invalid values were fixed and data was standardized in some columns, such as lead source.
- Low frequency values were grouped together to “Others”.
- Binary categorical variables were mapped.
- Other cleaning activities were performed to ensure data quality and accuracy.
  - Fixed Invalid values & Standardizing Data in columns by checking casing styles, etc. (lead source has Google, google

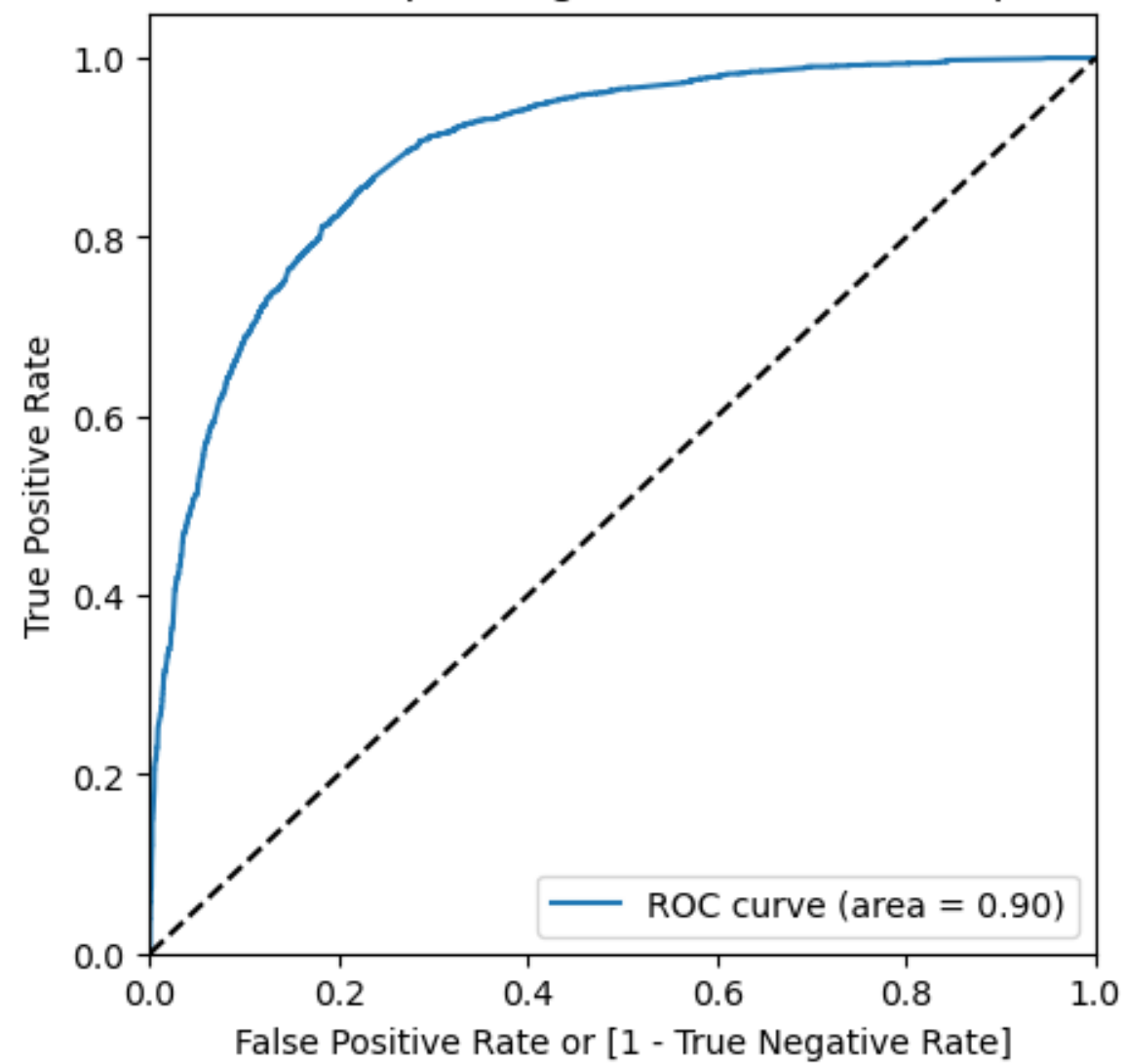
# Univariate Analysis



# Bivariate analysis



Receiver operating characteristic example





Accura,  
sensi,  
speci

