



# LEAD SCORING CASE STUDY

SHREYAS PATIL

KAJAL TOTAWAR

RAHUL SUPOLIA

# CONTENTS

- Problem Statement
- Business Objectives
- Data Set
- Solution Methodology
- Data preparation and EDA
- Model Building
- ROC/AUC Curve
- Conclusion

# PROBLEM STATEMENT

- X Education sells online courses to industry professionals. 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%.
- The company wants to identify the potential leads, also called Hot Leads.
- If company is able to identify this segment, then their sales team could call these leads , thus leading to better conversion rate.

# BUSINESS OBJECTIVES

- X education wants to select the most promising leads.
- X education wants to build a model to identify the hot leads.
- X education wants to deploy the model for future use as well in case requirement changes



# DATA SET

We have been provided with a leads dataset from the past with around 9000 data points:

- *This dataset consists of various attributes such as Lead Source, Total Time Spent on Website, Total Visits, Last Activity, etc.*
- *The target variable, in this case, is the column 'Converted' which tells whether a past lead was converted or not wherein 1 means it was converted and 0 means it wasn't converted.*
- *A data dictionary is also provided for more details on data.*

# SOLUTION METHODOLOGY

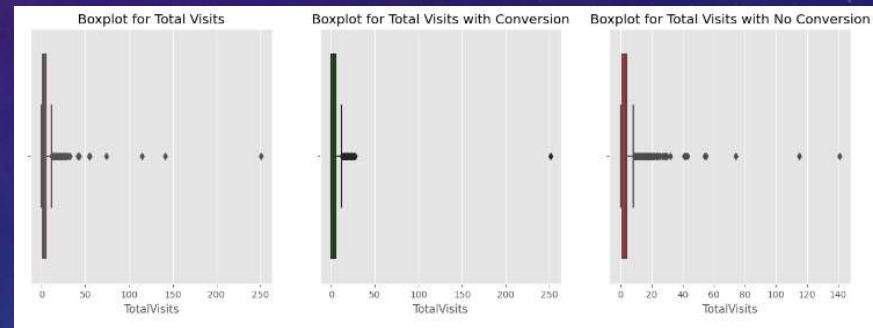
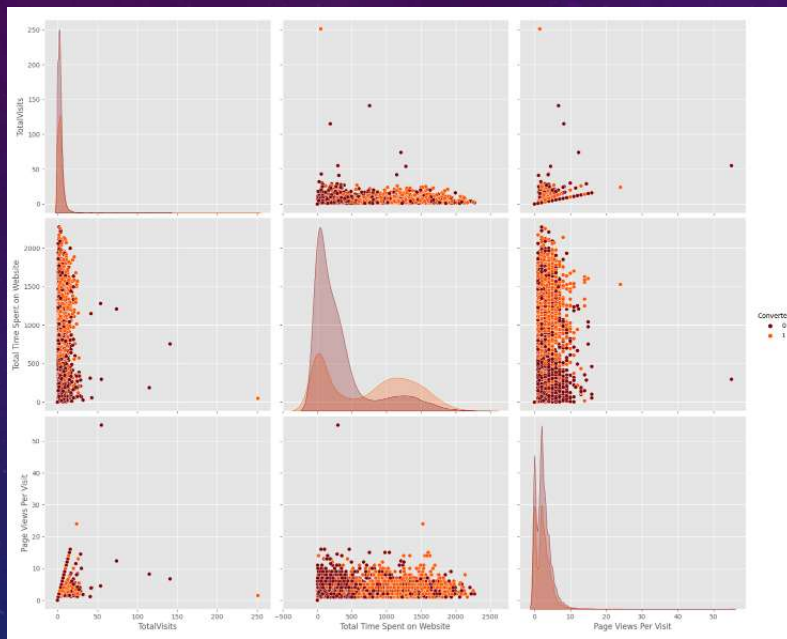
Below steps were used:

1. Data Cleaning
2. EDA
3. Train- Test split
4. Model Building
5. Model Evaluation
6. Model Presentation
7. Conclusion

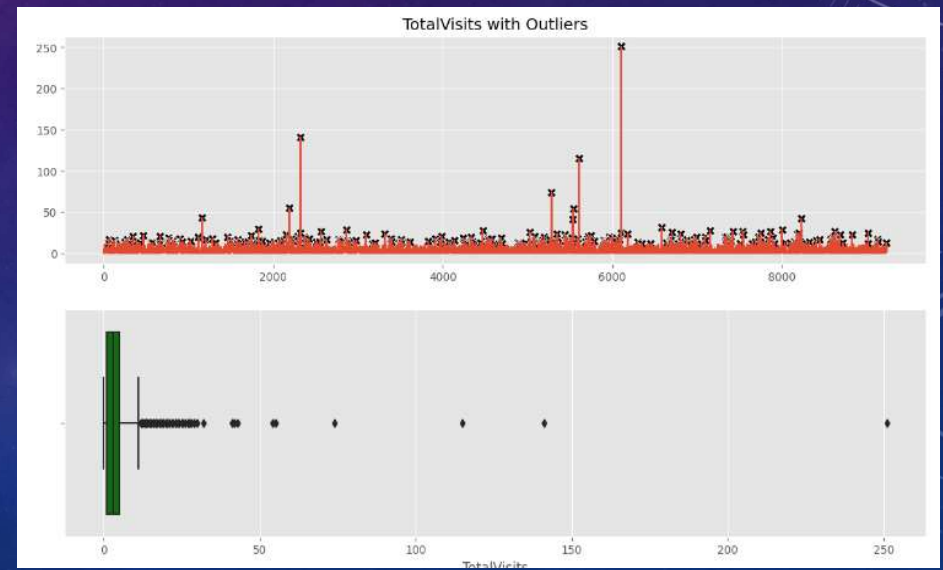
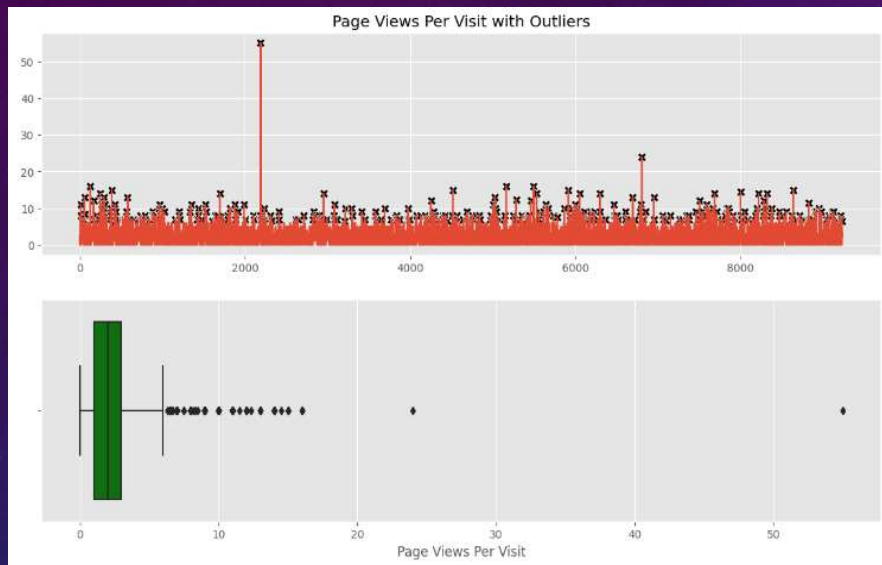
# DATA PREPARATION AND EDA

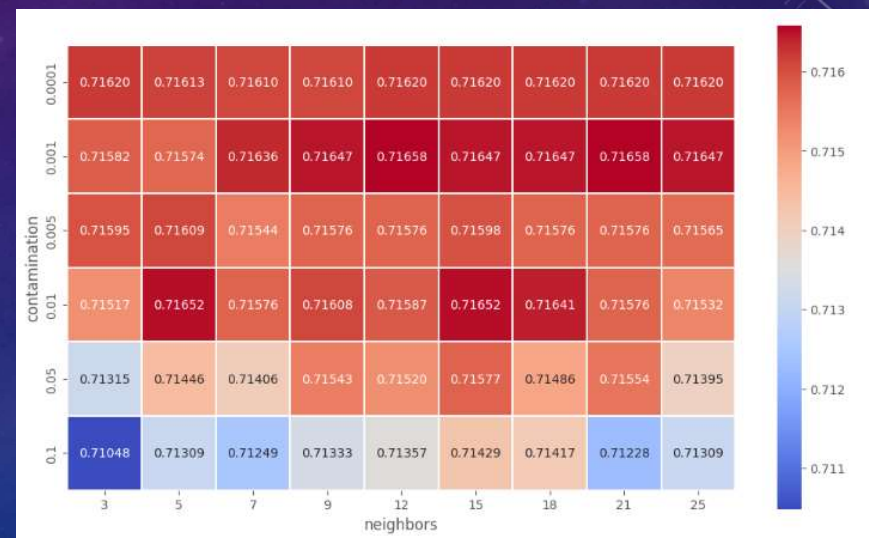
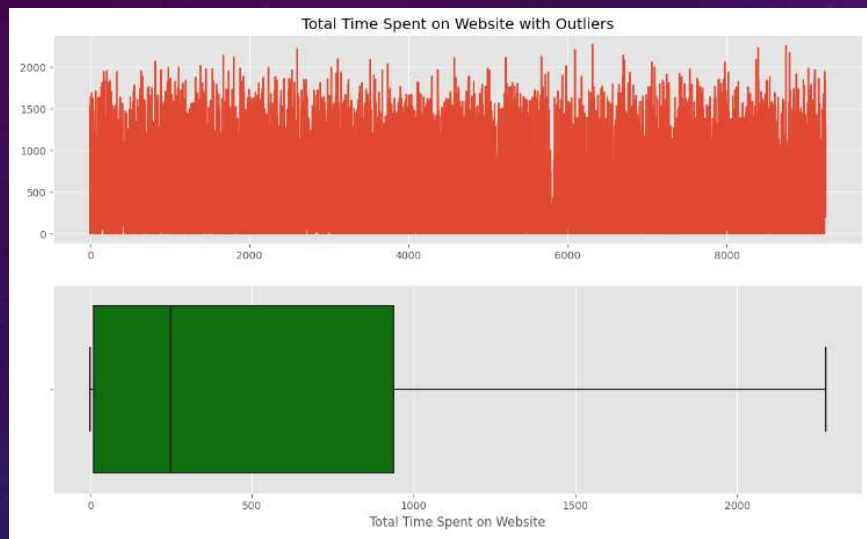
- Total Number of Rows =37, Total Number of Columns =9240.
- 4 columns with float values, 3 with integer value, and 30 with categorical value.
- Dropping 'Asymmetrique Activity Score', and 'Asymmetrique Profile Score' from numerical dataframe.
- Dropping from 'Lead Quality', 'Asymmetrique Profile Index', 'Asymmetrique Activity Index' categorical dataframe as null values are exceeding 40%.
- Created a new variable for those values which are either 'Select' or 'NaN' for these 4 columns.
- Imputed the frequent value in case of any null values present in our data.
- Addressed the outliers.

# EDA





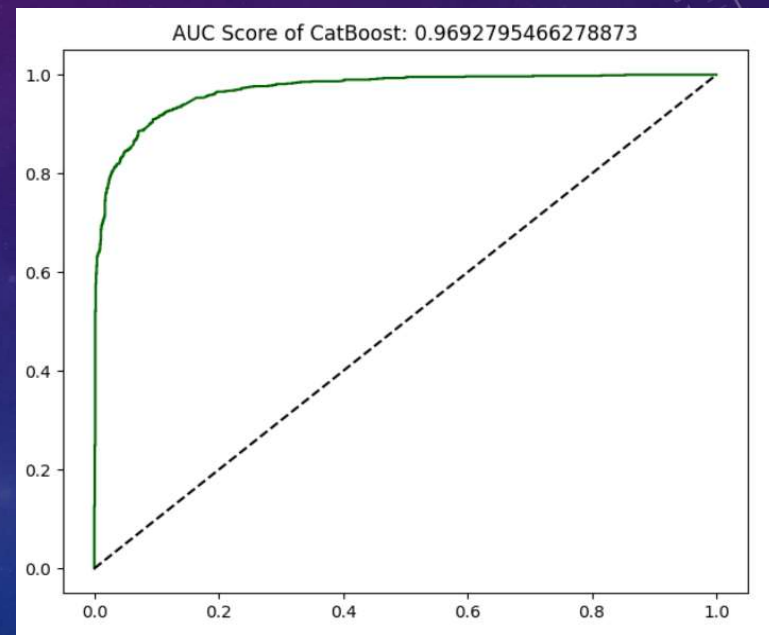
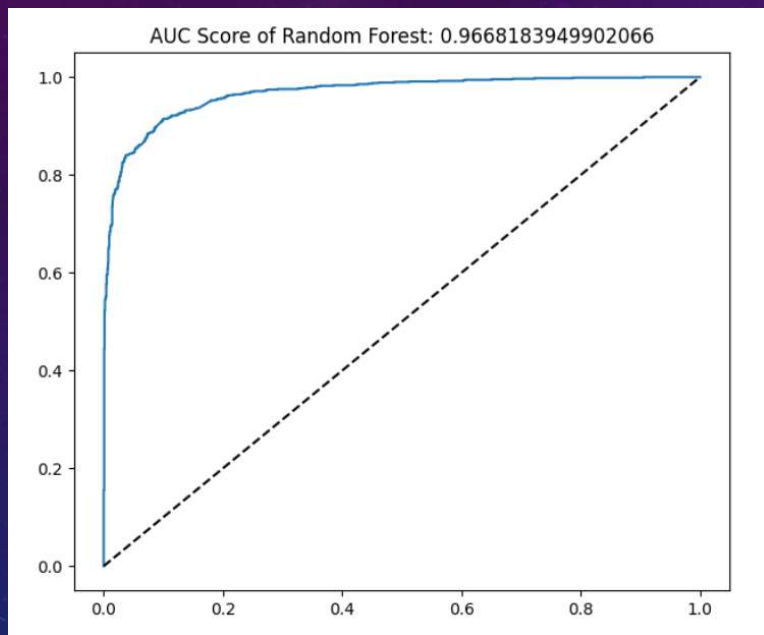




# MODEL BUILDING

- Splitting dataset into training and testing sets.
- The split was done at 75% and 25% for train and test data respectively.
- Based on the results, we chose the RandomForest Classifier, GradientBoosting, LightGBM & Catboost on which we tested the other metrics to see in depth performance of these 4 models based on several different metrics to choose the best model for our analysis.
- Predictions on test dataset.
- Hyperparameter tuning done with a focus on Random Forest and CatBoost models for enhanced results
- Overall accuracy comes to 91%

# PLOTTING ROC/AUC CURVE





# CONCLUSION

- It was found that the variables that mattered the most in the potential buyers:
  1. The total time spend on the Website.
  2. Total number of visits.
  3. When the lead source was:
    - a. Google
    - b. Direct traffic
    - c. Organic search
    - d. Welingak website
  4. When the last activity was:
    - a. SMS
    - b. Olark chat conversation
  5. When the lead origin is Lead add format.
  6. When their current occupation is as a working professional.