



# LEAD SCORE CASE STUDY

P K VIVEK VIDHYA

# CONTENTS



- **PROBLEM STATEMENT**
- **BUSINESS OBJECTIVE**
- **SOLUTION METHODOLOGY**
- **DATA MANIPULATION**
- **EDA AND CATEGORICAL/NUMERICAL VARIABLE RELATIONSHIP**
- **MODEL BUILDING**
- **ROC CURVE AND EXACT CUT OFF**
- **CONCLUSION**



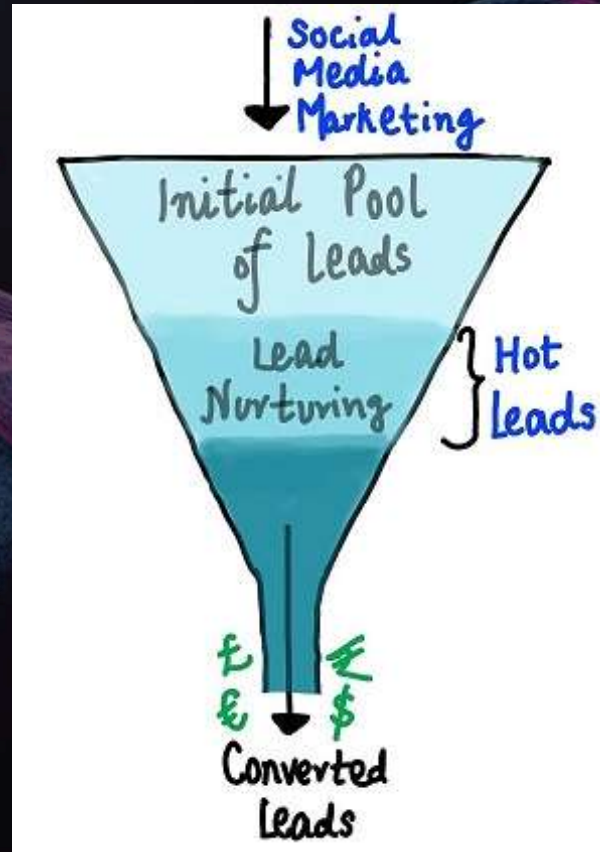
# PROBLEM STATEMENT

- A COMPANY NAMED X EDUCATION SELLS ONLINE COURSES TO ALL THE CUSTOMERS AND THEY GET LOT OF LEADS BUT THEIR RATE OF LEAD CONVERSION IS VERY LESS.
- FOR E.G. IF SAY THEY ACQUIRE 100 LEADS A DAY ONLY 30 ARE CONVERTED.

# PROBLEM STATEMENT

- IN ORDER TO MAKE THIS PROCESS MORE EFFICIENT OR FASTER THE COMPANY WISHES TO IDENTIFY POTENTIAL LEADS KNOWN AS **HOT LEADS**. THEY ALSO WANT TO FIND THE **COLD LEADS**.
- IF THEY SUCCESSFULLY IDENTIFY THOSE PEOPLE WHO CAN BE CONVERTED INTO SUCCESSFUL LEADS. THE RATE OF LEAD CONVERSION CAN GO HIGH AND THE SALES TEAM CAN FOCUS/COMMUNICATE WITH THEM AS THEY ARE IMPORTANT LEADS RATHER THAN MAKING CALLS TO EVERYONE.

# PROBLEM STATEMENT



THIS IS THE TYPICAL LEAD CONVERSION FUNNEL WHICH THE X EDUCATION COMPANY USES.



# BUSINESS OBJECTIVE

BUILD A LOGISTIC REGRESSION MODEL TO ASSIGN A LEAD SCORE BETWEEN 0 AND 100 TO EACH OF THE LEADS WHICH CAN BE USED BY THE COMPANY TO TARGET POTENTIAL LEADS. THEY WANT TO BUILD A MODEL TO IDENTIFY THE HOT LEADS AND TO DEPLOY THIS MODEL FOR THE FUTURE USE.

# SOLUTION METHODOLOGY



- UNDERSTANDING THE BUSINESS PROBLEM AND THE OBJECTIVE
- IMPORTING AND KNOWING THE INSIGHTS OF THE DATA
- DATA CLEANING AND MANIPULATION
  1. TREATING THE MISSING VALUES BY DROPPING THE COLUMNS WHICH HAVE HIGH RATE OF MISSING VALUES AND DROPPING THOSE COLUMNS/FEATURES WHICH AREN'T IMPORTANT FOR THE ANALYSIS
  2. CHECK THE OUTLIERS

# SOLUTION METHODOLOGY

- **EDA ANALYSIS ( EXPLORATORY DATA ANALYSIS)**
  - 1. FINDING THE DIFFERENT CLASSES OF CERTAIN FEATURES BY USING VALUE\_COUNTS()**
  - 2. VISUALIZATION AMONG THE NUMERICAL VARIABLES (BY USING PAIR PLOT)**
  - 3. VISUALIZATION AMONG THE CATEGORICAL VARIABLES WITH RESPECT TO DEPENDENT VARIABLE (BY USING BOX PLOT)**
  - 4. FIND THE FREQUENCY OF DISTRIBUTION AMONG THE VARIABLES (BY USING HISTOGRAM)**
  - 5. UNIVARIATE AND BIVARIATE ANALYSIS**



# SOLUTION METHODOLOGY

- **CREATE A DUMMY VARIABLE FOR THOSE CATEGORICAL FEATURE WHICH IS HAVING MORE THAN TWO CLASSES.**
- **USE MIN MAX SCALER TO SCALE DOWN THE VALUES OF THE FEATURE WHICH IS NOT IN ZERO AND ONE.**

# SOLUTION METHODOLOGY

- **USE CLASSIFICATION TECHNIQUE AS OUR DEPENDENT/LABEL AS HAVING CATEGORICAL VALUE IN THE FORM OF ZERO AND ONE. SO, USE LOGISTIC REGRESSION MODEL.**
- **USE RFE FOR AUTO FEATURE SELECTION OF TOP FIFTEEN**
- **VALIDATION OF THE MODEL AND MODEL PRESENTATION**
- **GIVE A CONCLUSION AND RECOMMENDATION**



# DATA MANIPULATION

- INITIALLY OUR DATA SHAPE WAS 9240 ROWS AND 37 COLUMNS.
- WE DROPPED CITY AND COUNTRY AS IT WAS NOT REQUIRED FOR THE ANALYSIS
- ALSO FOUND FEW FEATURES HAVING A VALUE SELECT (THAT MEANS THE CUSTOMER HAS TO CHOSE). SO, DROPPED THE COLUMNS 'HOW DID YOU HEAR ABOUT X EDUCATION' AND 'LEAD PROFILE'.

# DATA MANIPULATION

- 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. THESE COLUMNS HAVE VALUES MAXIMUM NO .ITS NOT USEFUL FOR ANALYSIS.



# DATA MANIPULATION

- ALSO DROPPED FEW ROWS OF SOME COLUMNS WHICH HAVE NULL VALUES.
- FINALLY LEFT WITH **12 COLUMNS** AND **6373 ROWS**.

# DATA MANIPULATION

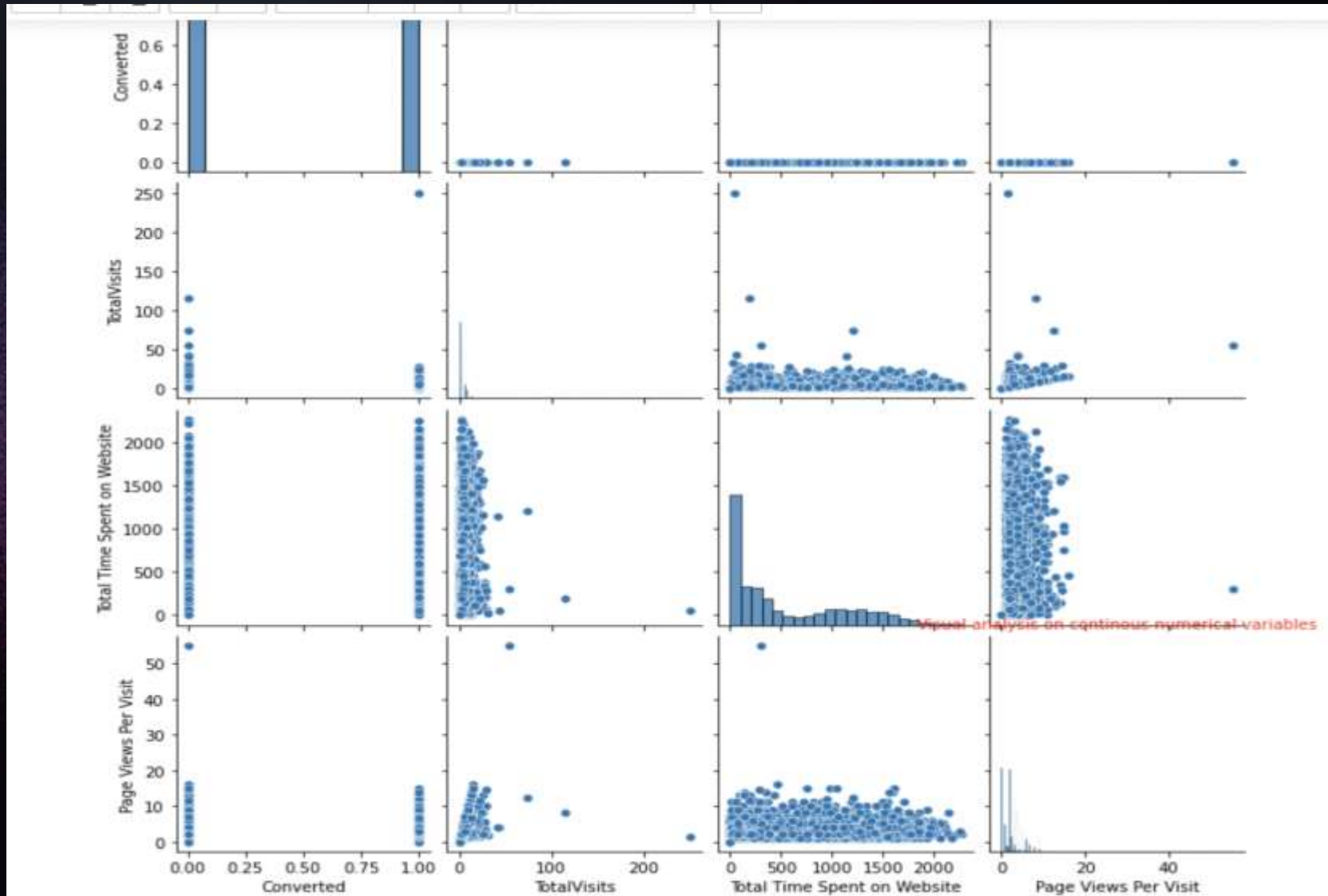
- ALSO CREATED DUMMY VARIABLE FOR FEW COLUMNS WHICH IS HAVING MORE THAN TWO CLASS AND IMPUTED INTO ONE DATAFRAME AND FINALLY CONCATENATED WITH OUR ORIGINAL DATAFRAME.
- ALSO DROPPED THE ORIGINAL COLUMNS WHOSE DUMMY VARIABLE IS CREATED. SO OUR FINAL DATAFRAME IS LEFT WITH **6373 ROWS** AND **75 COLUMNS**.





# EDA AND CATEGORICAL/NUMERICAL VARIABLE RELATIONSHIP

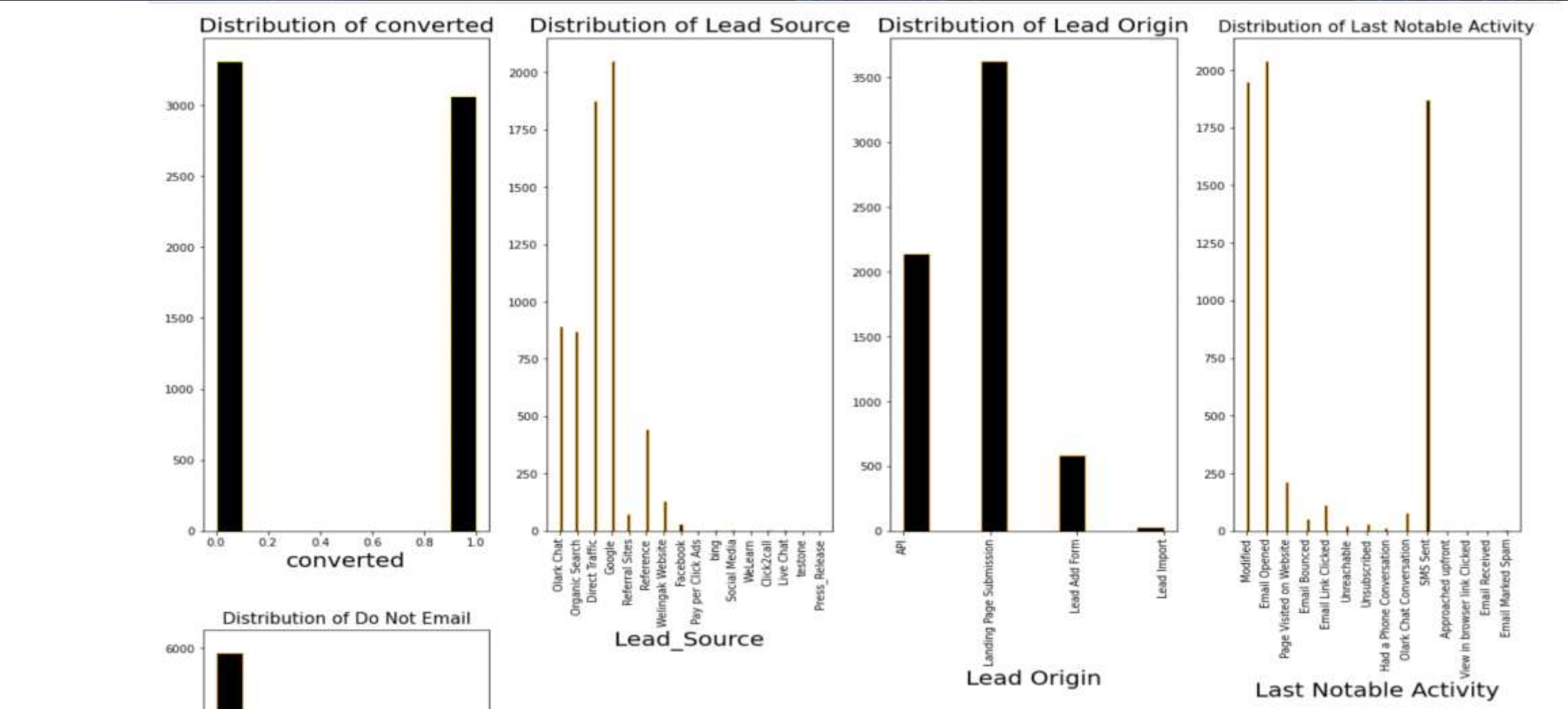
## Analysis and visualizing numerical variable by using pair plot.



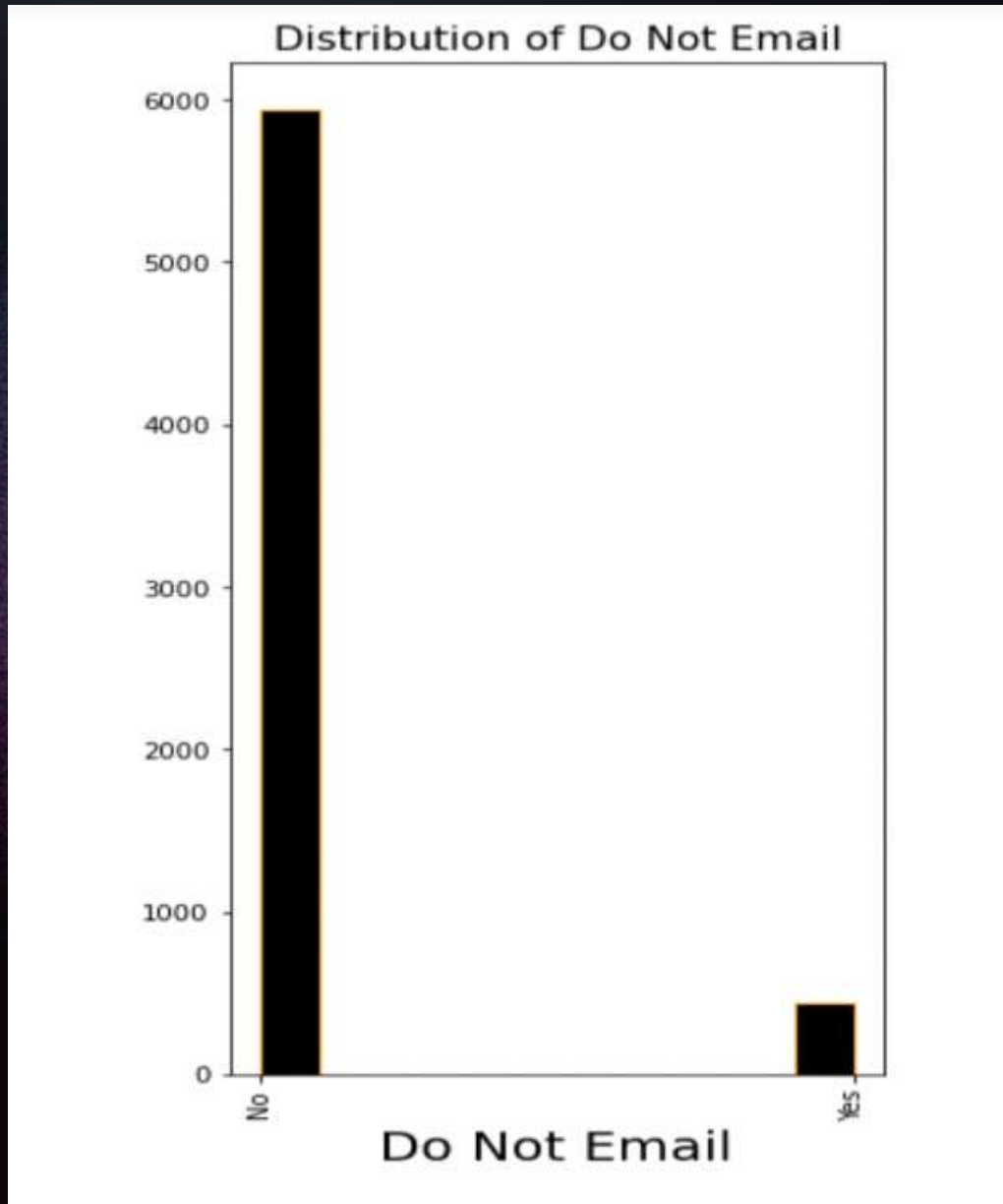
**OBSERVATION:** FOUND A KIND OF LINEAR RELATIONSHIP BETWEEN TOTAL VISITS AND PAGE VIEWS PER VISIT.



Frequency distribution of all categorical values by using histogram in sub plot.

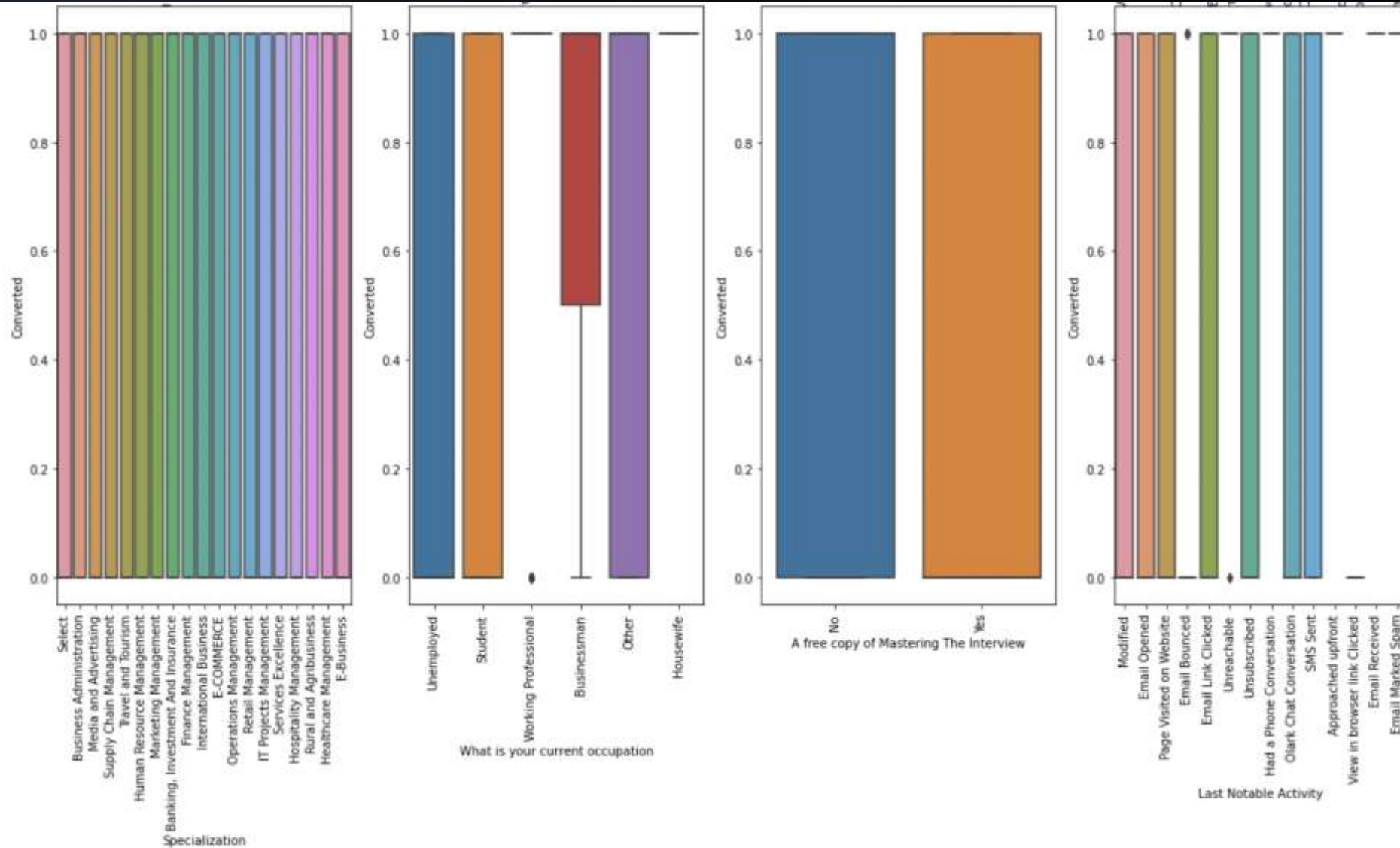


## Frequency distribution of all categorical values by using histogram in sub plot.



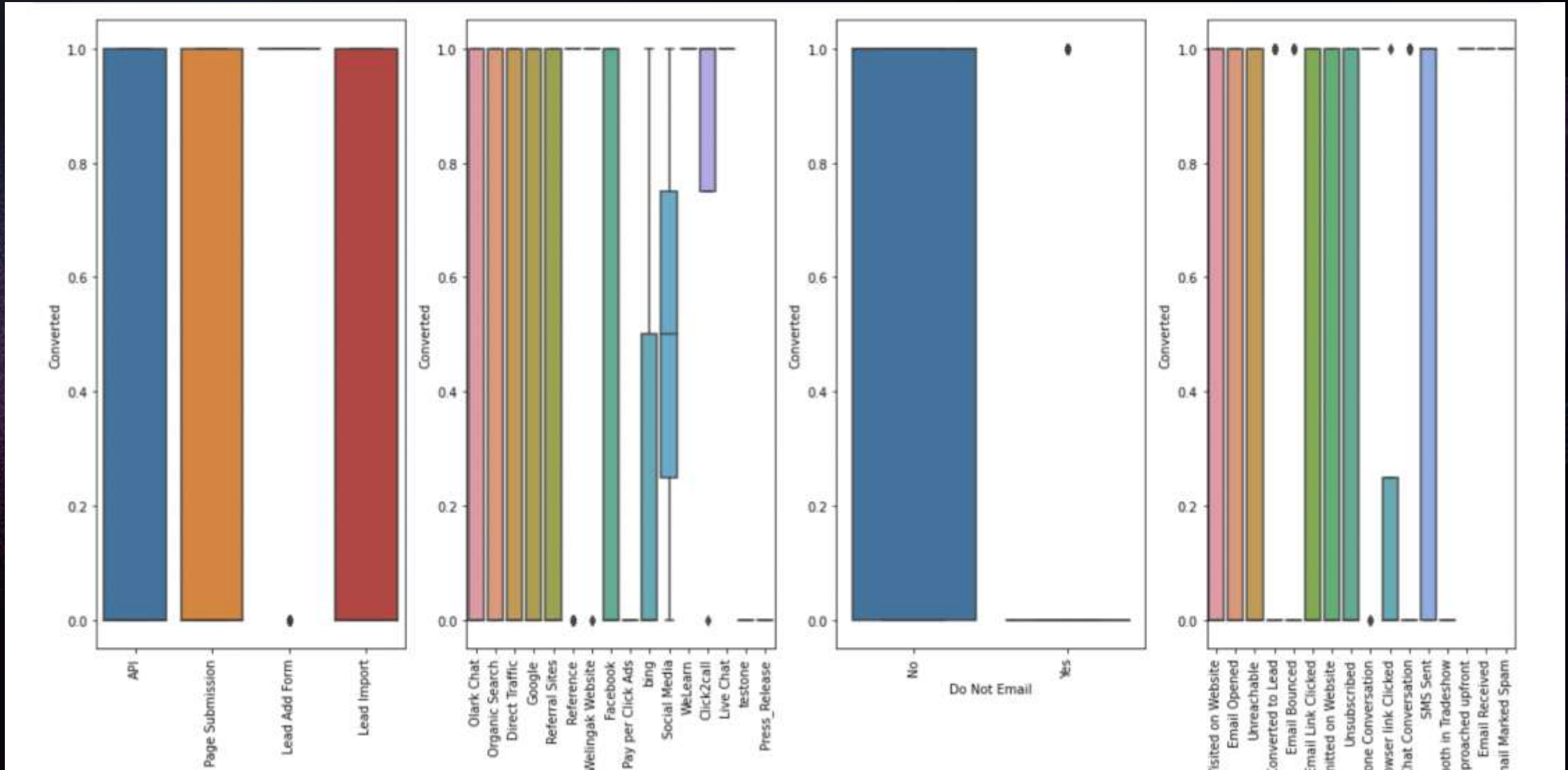
- **OBSERVATION:-** FOUND DISTRIBUTION OF CONVERTED RATE IS LESS THAN NO LEAD ,AND IN LEAD SOURCE GOOGLE HAVE MORE DISTRIBUTION ,LEAD ORIGIN = LANDING PAGE HAVE MORE DISTRIBUTION ,LN LAST ACTIVITY EMAIL OPENED AND SMS SENT HAVE MORE DISTRIBUTION AND DO NOT EMAIL SENT **NO** HAVE MORE DISTRIBUTION IN OUR DATASET

# Visualization of categorical variable vs dependent variable using box plot in sub plot





# Visualization of categorical variable vs dependent variable using box plot in sub plot



## Observation on Visualization of categorical variable vs dependent variable using box plot in sub plot

- Free copy of mastering the interview both values yes and no are converted
- Customer who choose to respond by emailing do not email (No) are converted
- Also found customer of different specialization are converted

## Observation on Visualization of categorical variable vs dependent variable using box plot in sub plot

- maximum from google search ,organic ,direct traffic etc. are converted
- In the last labeled activity customers who opened email, visited the website, chat with olark and link clicked were converted
- Customers who opened email, SMS, visited the website and clicked the link in their last activity were converted.



# MODEL BUILDING

- SPLIT THE DATA INTO 'X' AND 'Y' (INDEPENDENT AND DEPENDENT VARIABLE). AGAIN SPLIT 'X' AND 'Y' INTO TRAIN AND TEST DATASET.
- I HAVE CHOSE 70-30 RATIO FOR SPLITTING, 70 PERCENTAGE IS TRAIN PART AND 30 PERCENTAGE IS TEST PART.
- USED MIN MAX SCALER TO SCALE DOWN THE VALUES FOR THOSE VARIABLE WHOSE VALUES ARE NOT IN THE FORM OF ZERO AND ONE.

# MODEL BUILDING

- **USED RFE FEATURE WITH AN OUTPUT OF TOP 15 SELECTED FEATURE**
- **MADE PREDICTION ON TRAIN DATA FIRST AND ULTIMATELY PREDICTED ON TEST DATA.**
- **I REPEATEDLY MADE MODELS TO IMPROVE THE ACCURACY AND THE FIFTH MODEL WAS THE FINAL MODEL WITH DESIRED ACCURACY AS THERE IS A MULTICOLLINEARITY BETWEEN THE DEPENDENT VARIABLE AND ALSO WITH A HIGH 'P' VALUE.**

# MODEL BUILDING

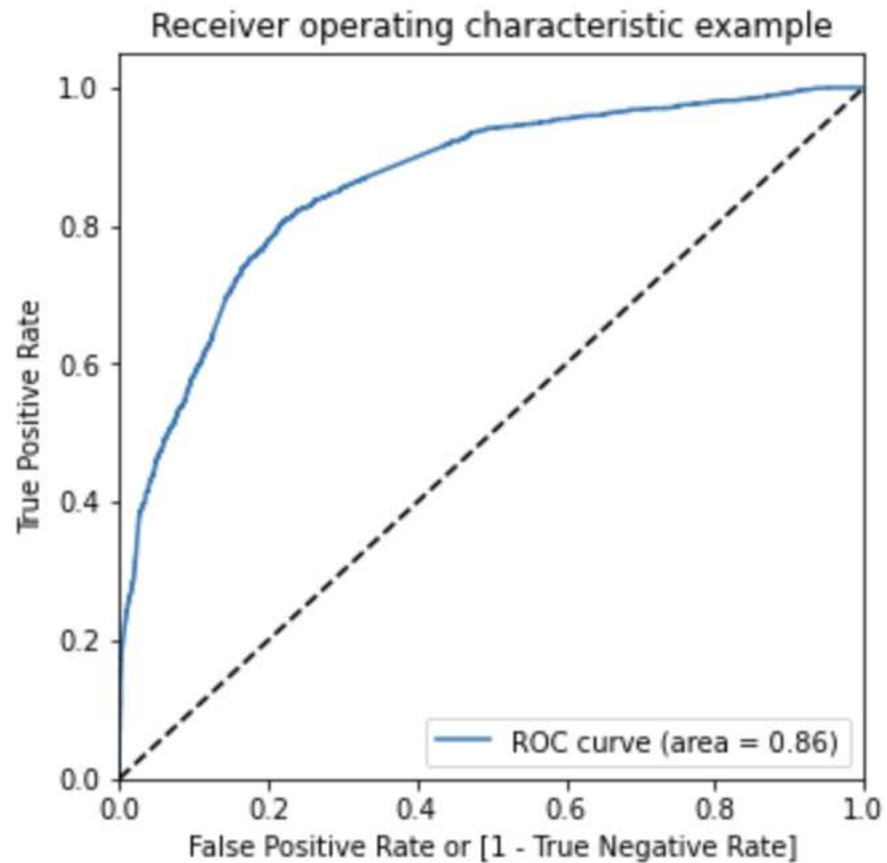
- I HAVE USED A CLASSIC LIST COMPREHENSION METHOD TO CREATE A DATAFRAME WITH VARIABLES NAMED ACTUAL, PREDICTED AND CONVERTED PREDICTED.
- I HAVE FOUND CUT-OFF BEING 0.4 AND DECIDED TO CUT THE PREDICTION VALUE.
- MY FINAL MODEL ACCURACY SCORE ON TEST DATA IS 0.78



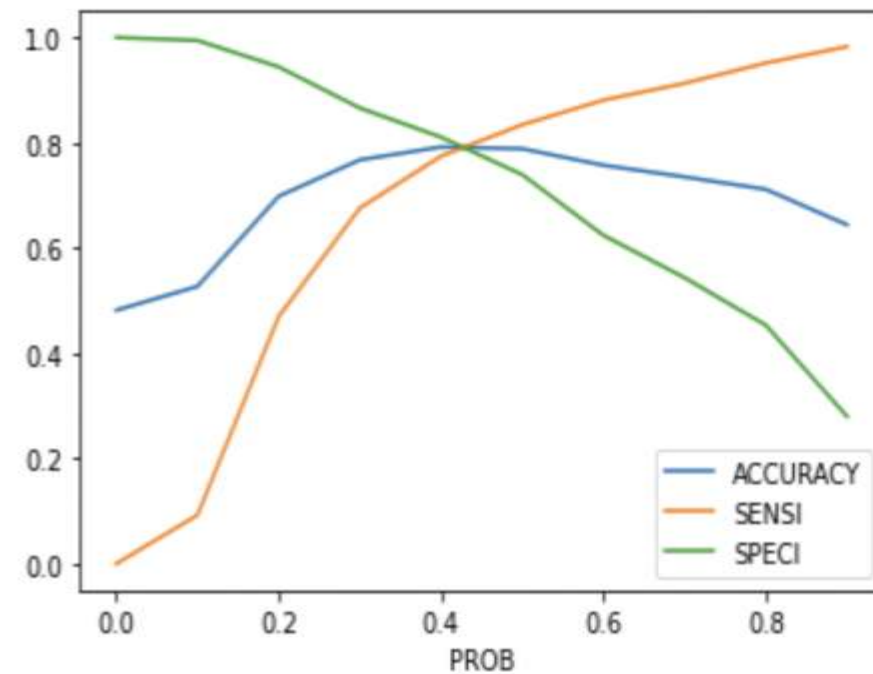
# MODEL BUILDING

- SENSITIVITY(TRUE POSITIVE RATE) = 0.78
- SPECIFICITY(TRUE NEGATIVE RATE)=0.79
- RECALL=0.78
- PRECISION=0.77
- F1 SCORE IS 0.78

# ROC CURVE AND EXACT CUT-OFF



```
# Plotting for exact cutoff  
new_cut.plot.line(x='PROB', y=['ACCURACY', 'SENSI', 'SPECI'])  
plt.show()
```



# ROC CURVE AND EXACT CUT-OFF

## OBSERVATION

- FOUND REGION UNDER CURVE IS 0.86 .IT LOOKS LIKE GOOD MODEL .BUT NEED TO CHECK THE IN CUT OFF RANGE
- AS WE FOUND AROUND 0.42 IS OPTIMAL VALUE WHERE ALL THE THREE METRICS INTERSECT BETWEEN EACH OTHER ( ACCURACY, SENSITIVITY AND SPECIFICITY)



# CONCLUSION

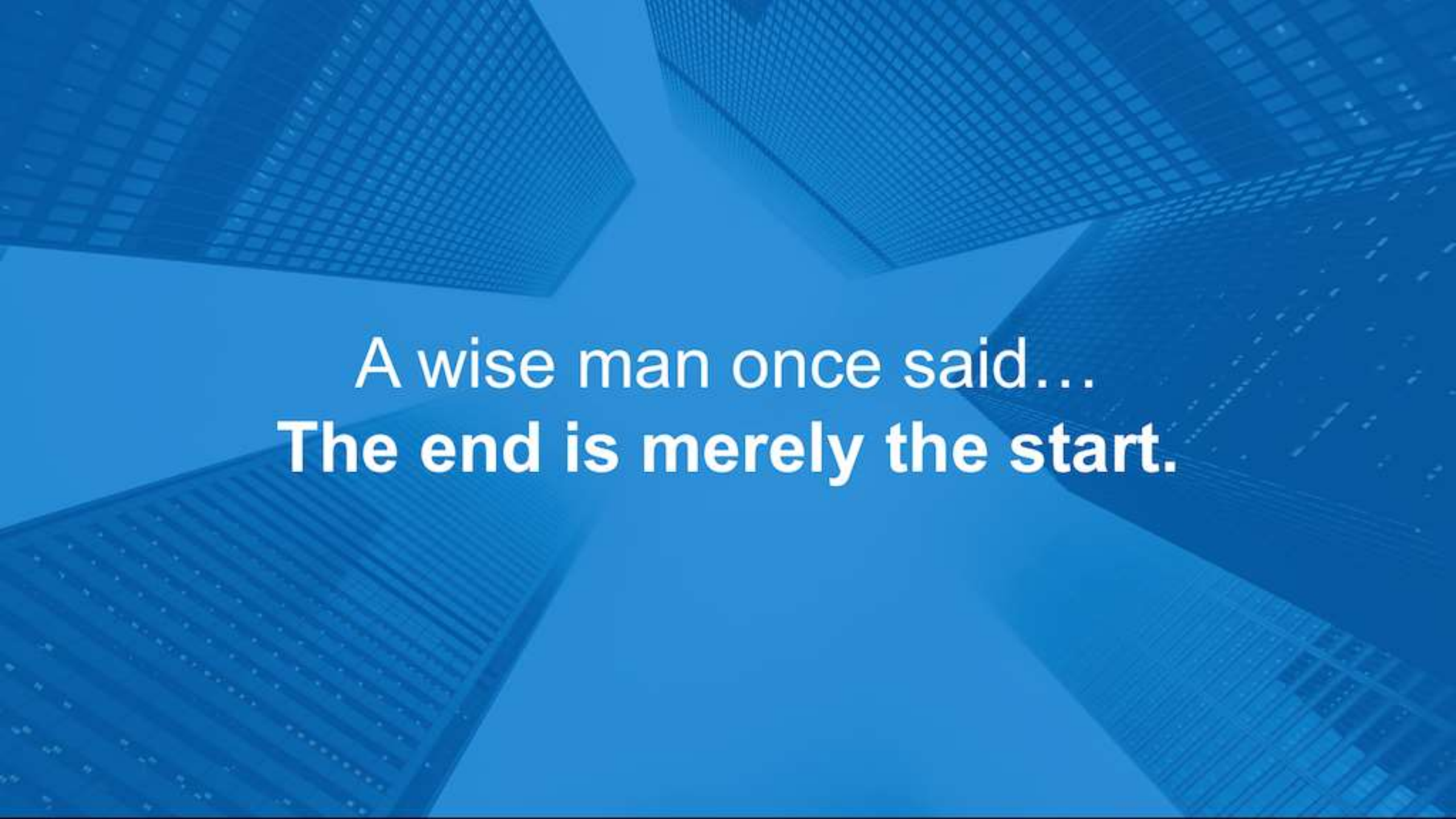
- FOUND THAT THE VARIABLES WHICH ARE HAVING NEGATIVE AND POSITIVE IMPACT ON LEAD CONVERSION
- **POSITIVE VARIABLES** USED FOR LEAD CONVERSION
  - 1.THE TOTAL TIME SPEND ON THE WEBSITE.
  - 2.TOTAL NUMBER OF VISITS.
  - 3.WHEN THE LEAD SOURCE WAS:
    - a. WELINGAK WEBSITE
    - b. OLARK CHAT
  - 4.WHEN THE LAST ACTIVITY WAS:
    - a. SMS
    - b. PHONE CONVERSATION
  - 5.WHEN THE LEAD ORIGIN IS LEAD ADD FORM.

# CONCLUSION

- **NEGATIVE VARIABLES** WHICH HAVE A NEGATIVE IMPACT ON LEAD CONVERSATION
  - A. IF THE CURRENT OCCUPATION IS STUDENT AND UNEMPLOYED

Above aspect we should Keep in mind, The X Education can improve by seeing there negative and positive variables .By using Above features in mind X education can achieve their goal.



The background is a solid blue color with several large, semi-transparent geometric shapes overlaid. These shapes are composed of fine, parallel lines that create a grid-like or woven texture. The shapes are arranged in a way that they appear to be receding into the distance, creating a sense of depth and perspective. The overall effect is modern and architectural.

A wise man once said...  
**The end is merely the start.**



The  
End