



# LEAD SCORING CASE STUDY

## LOGISTIC REGRESSION

GROUP MEMBERS:

UNNATI GARG

SINGH NIDHI KUMARI SANJEEV

KUMARI NEHA

DS C-50

# CONTENTS

- Problem Statement
- Goals and Business Objective
- Python Libraries Used
- Steps to be followed for analysis
- How to implement the approach
- Analysis of Data
- Conclusion

# PROBLEM STATEMENT

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

# GOALS AND BUSINESS OBJECTIVE

- X education wants to know most promising leads. For that they want to build a Model which identifies the hot leads
- Build a logistic regression model to assign a lead score between 0 to 100 each of the leads which can be used the company to target potential leads
- Lead with higher score considered as Hot (most likely to get converted) while with lower score considered as Cold (less likely to get converted)
- Deployment of the model for the future use.

# STEPS TO BE FOLLOWED FOR ANALYSIS

- Import required python libraries and dataset
- Data Cleaning
- EDA
  - Inspecting the data frame
  - Checking for missing or null values
  - Analyze the columns and handling the missing values
  - Handle outliers
  - Create dummy variables
  - Convert to binary variables: Yes/No to 1/0
  - Group column features
- Building the Model
- Train-test Split
- Plot ROC Curve
- Find optimal cut off points
- Calculate Accuracy, Sensitivity and Specificity
- Calculate Precision and Recall
- Validating the model
- Deploying the model
- Prediction on test set
- Conclusion and Result

# PYTHON LIBRARIES USED

- Numpy
- Pandas
- Matplotlib
- Seaborn
- Sklearn.model\_selection import train\_test\_split
- Sklearn.preprocessing import StandardScaler
- Sklearn.linear\_model import LogisticRegression
- Statsmodel.api
- Statsmodel.stats.outliers\_influence import variance\_inflation\_factor
- Sklearn import metrics
- sklearn.metrics import confusion\_matrix
- sklearn.metrics import precision\_score, recall\_score
- sklearn.metrics import precision\_recall\_curve

# DATA CLEANING

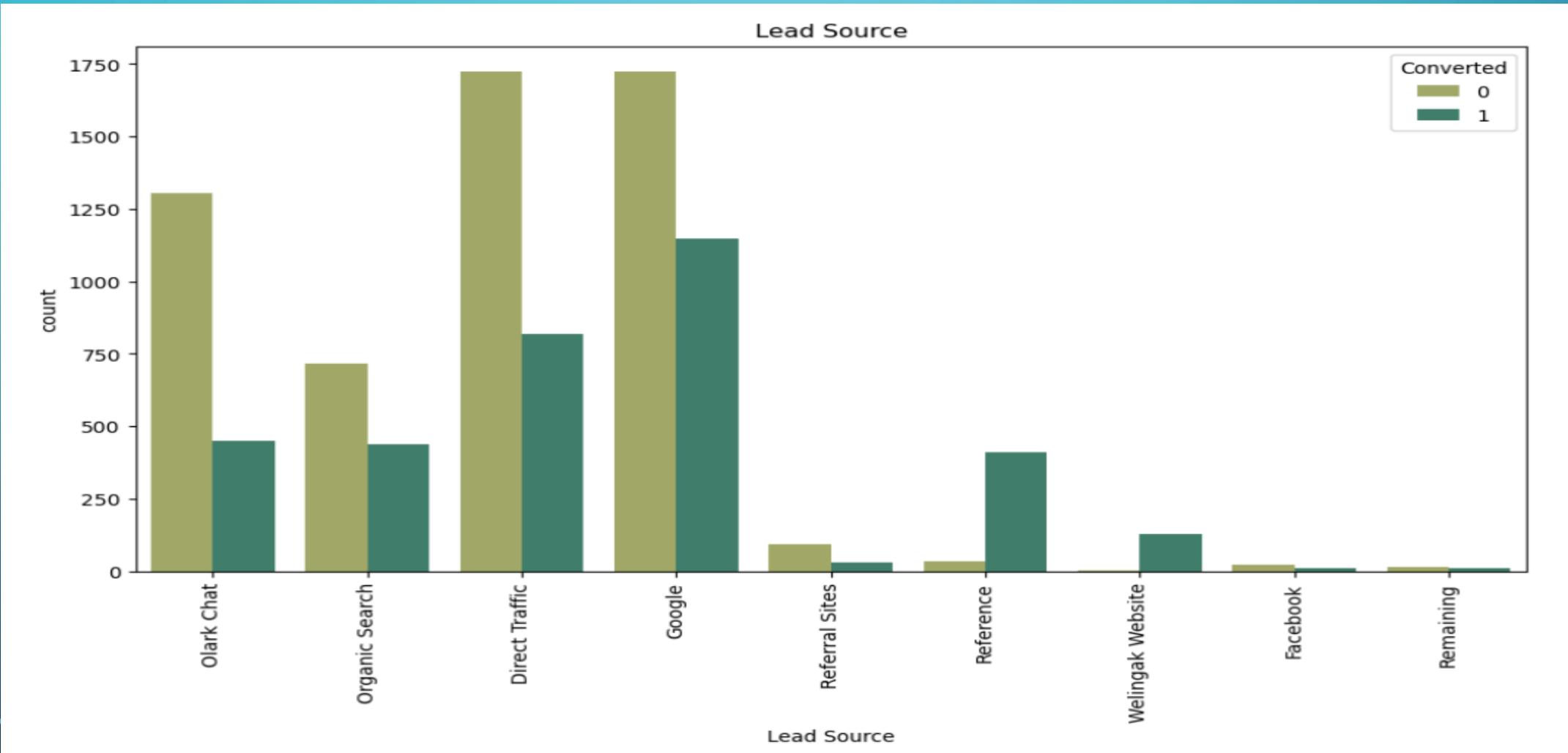
- Check number of rows and columns in the dataset
- Check for null or default values
- Replace default values with NAN
- Handling some columns with null values
- Drop null values



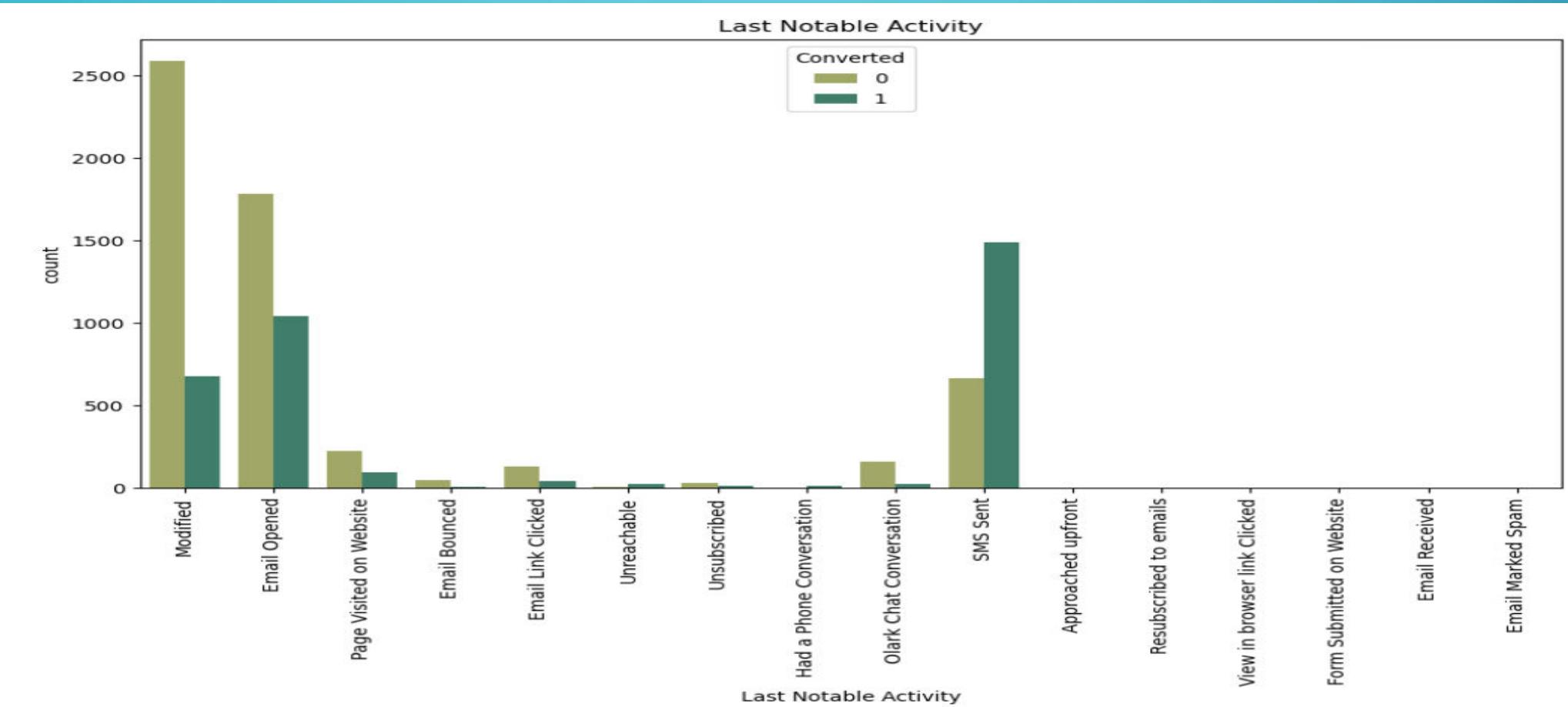
## EDA

- Visualization performed with attributes of dataset
- Most of the attributes were not found useful for analysis
- Useful attributes are “Lead Score”, “Last Activity”, “Lead Origin” and “Total time spent on website”

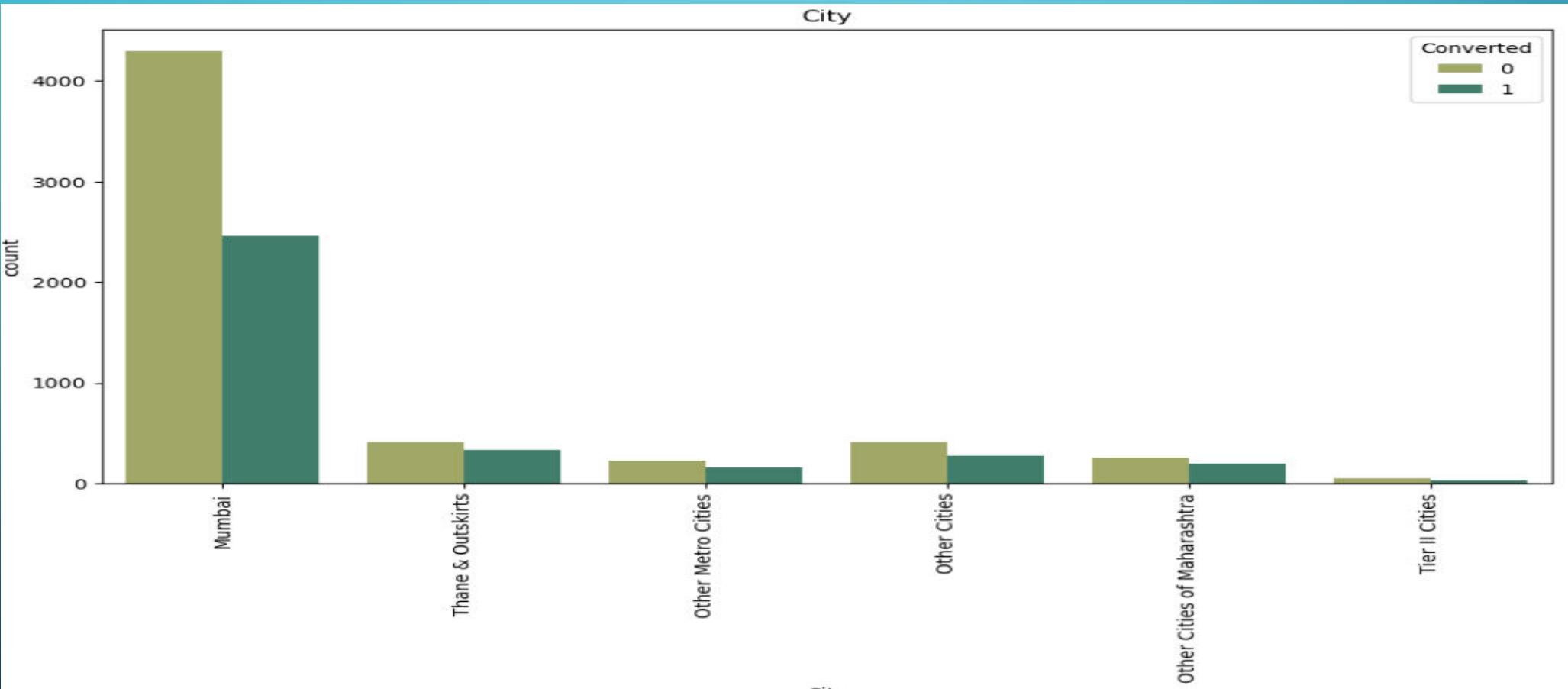
# LEAD SCORE



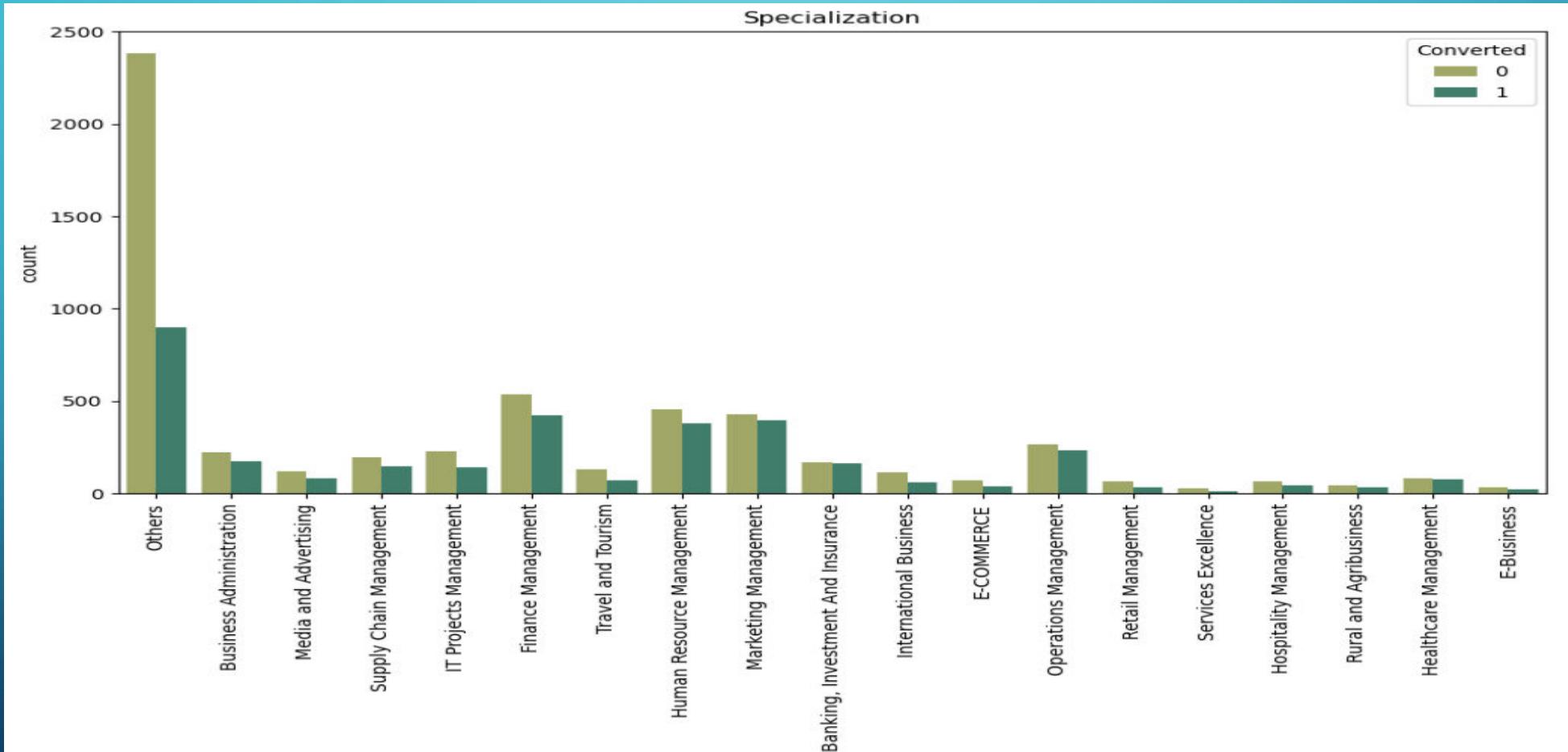
# LAST NOTABLE ACTIVITY



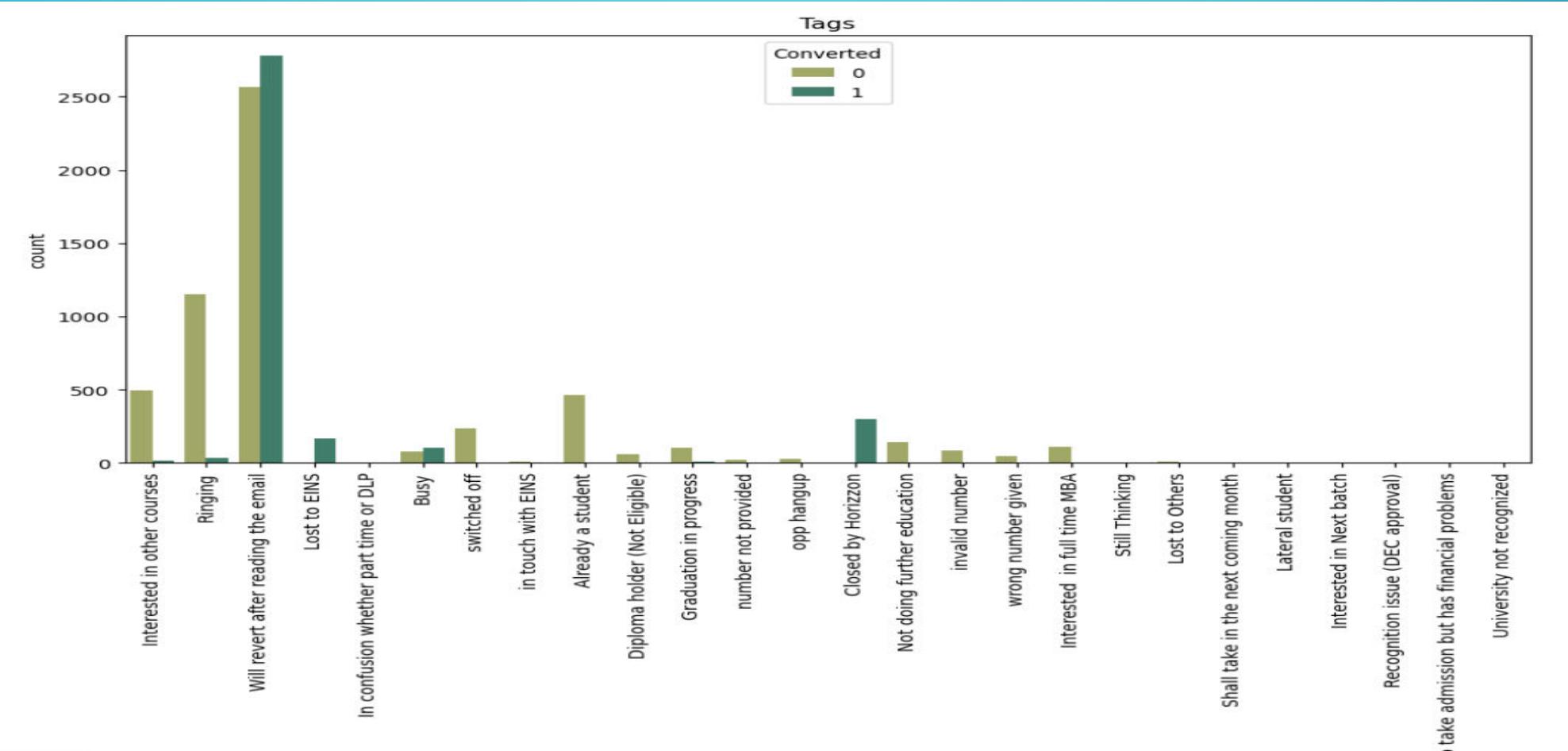
# CITY



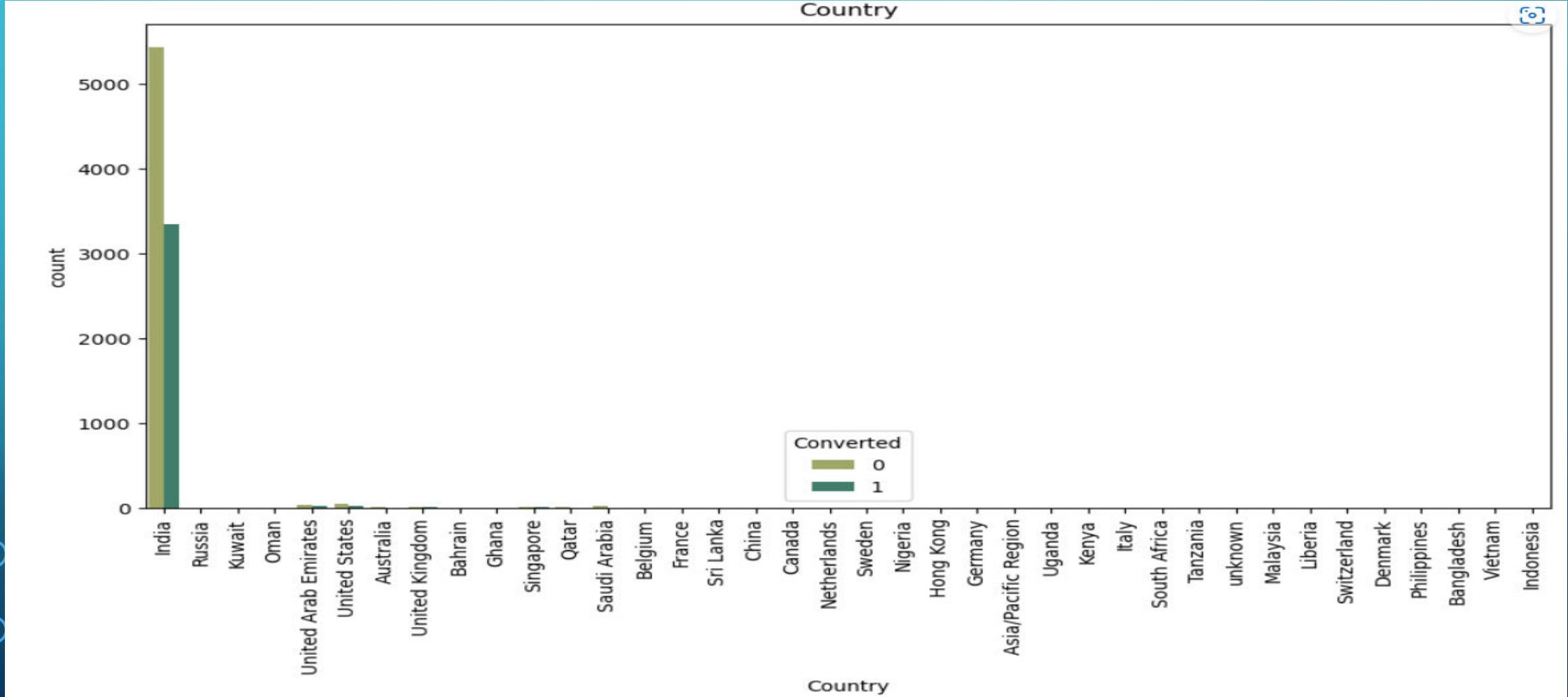
# SPECIALIZATION



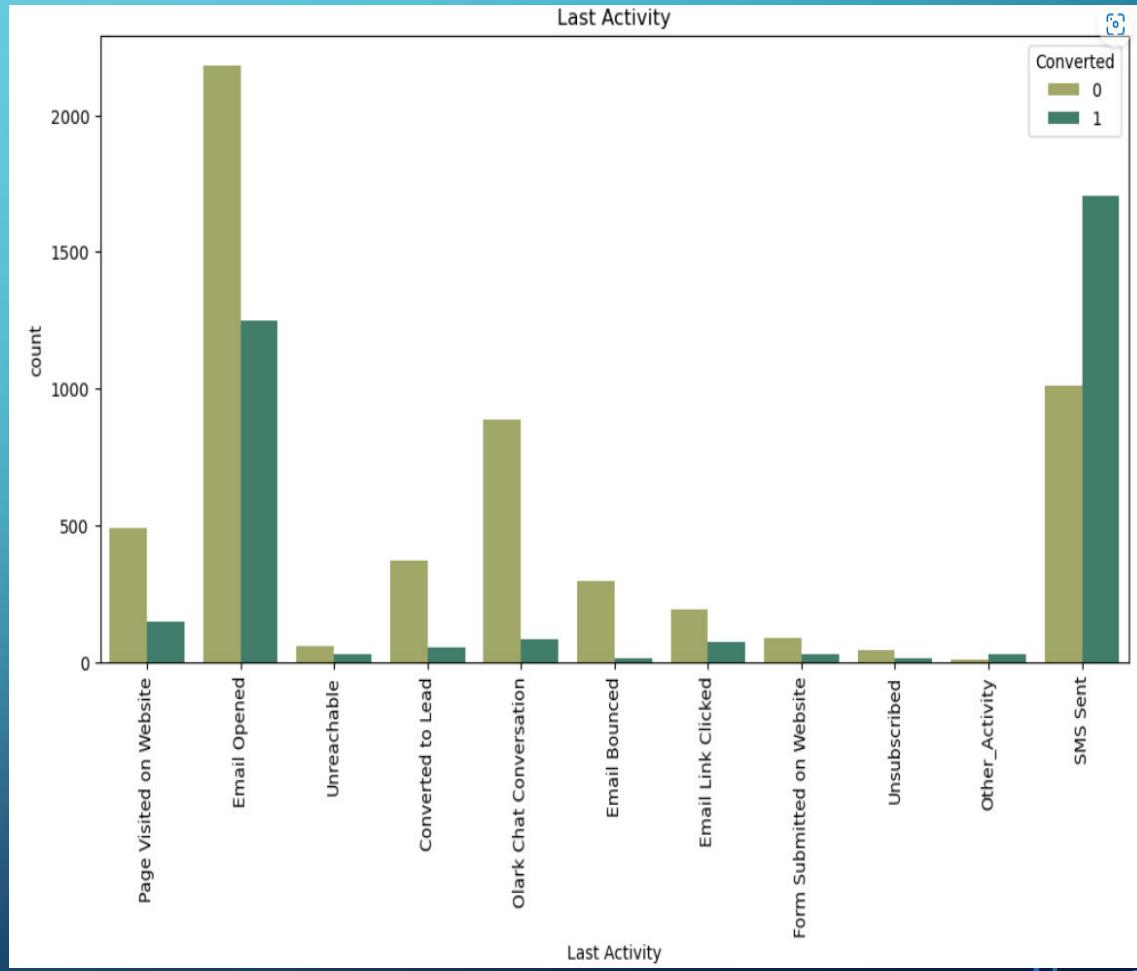
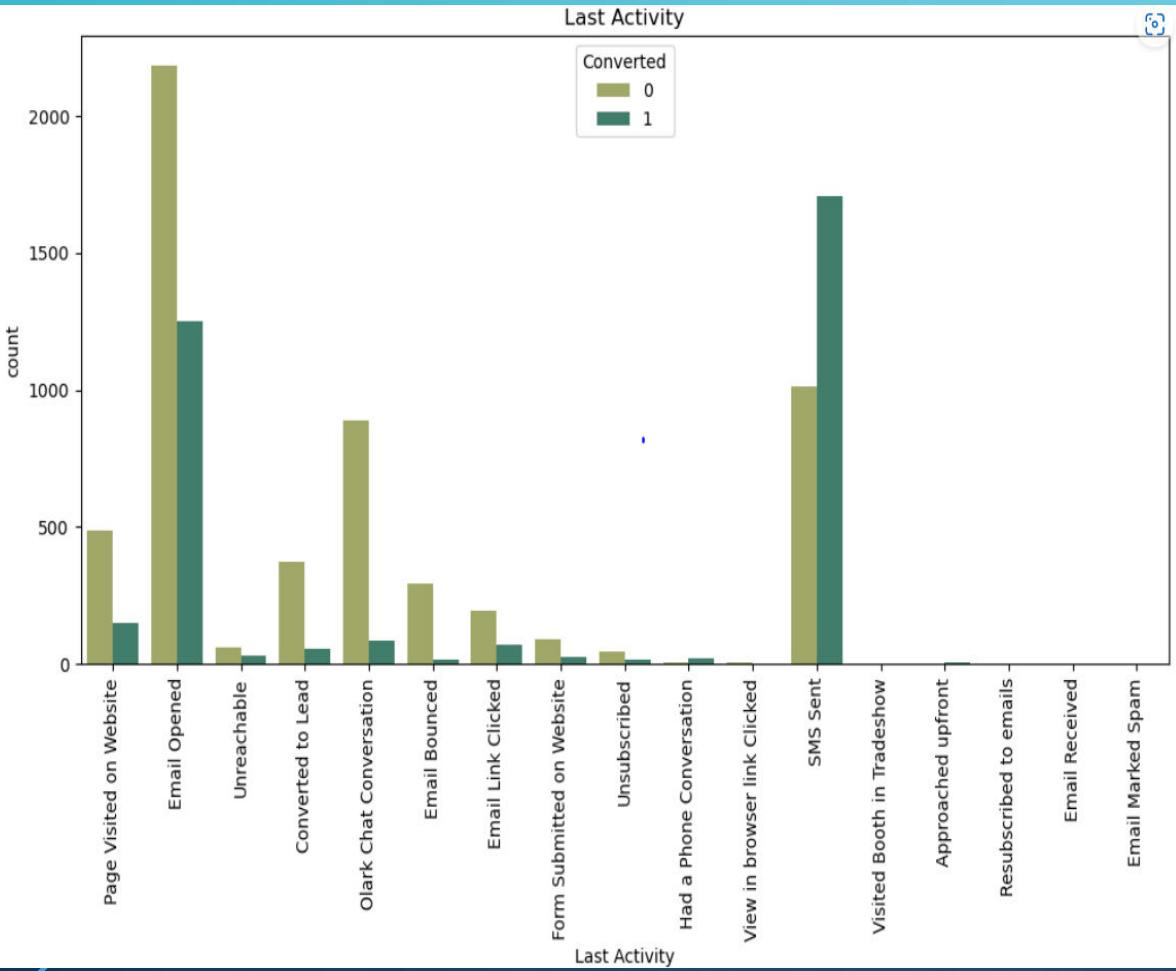
# TAGS



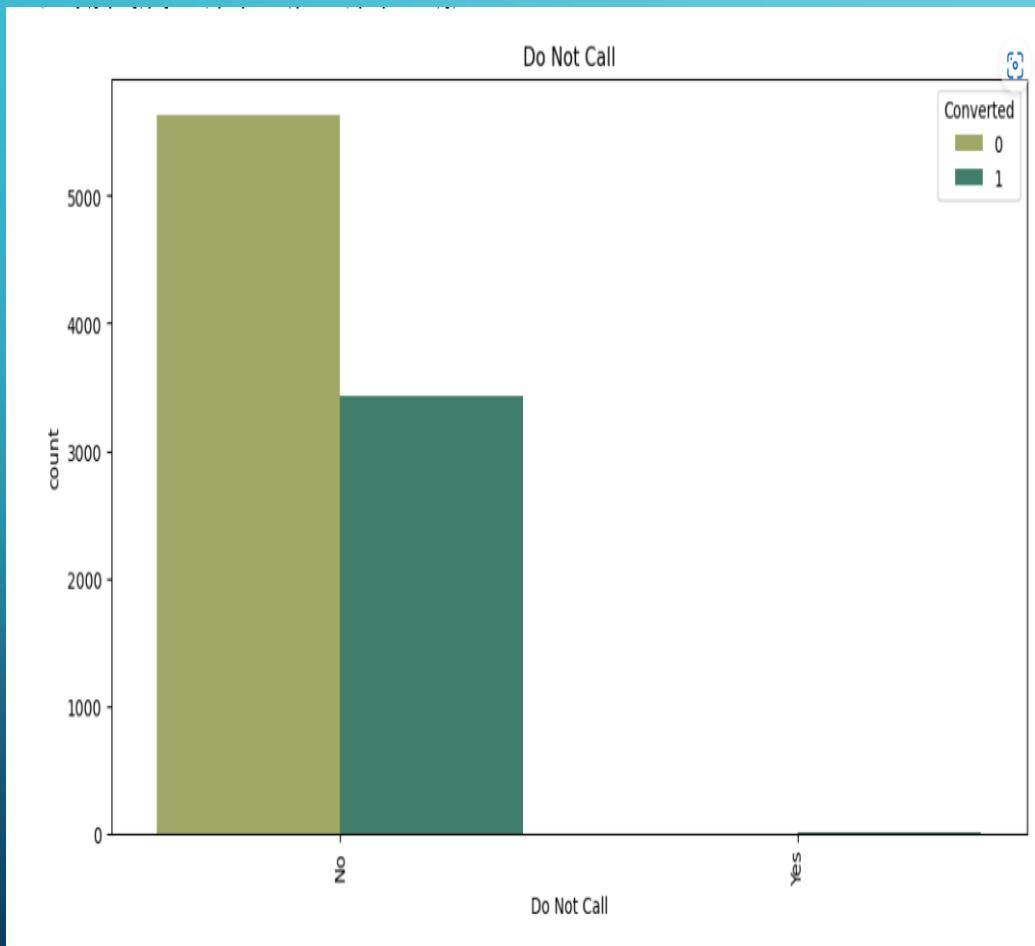
# COUNTRY



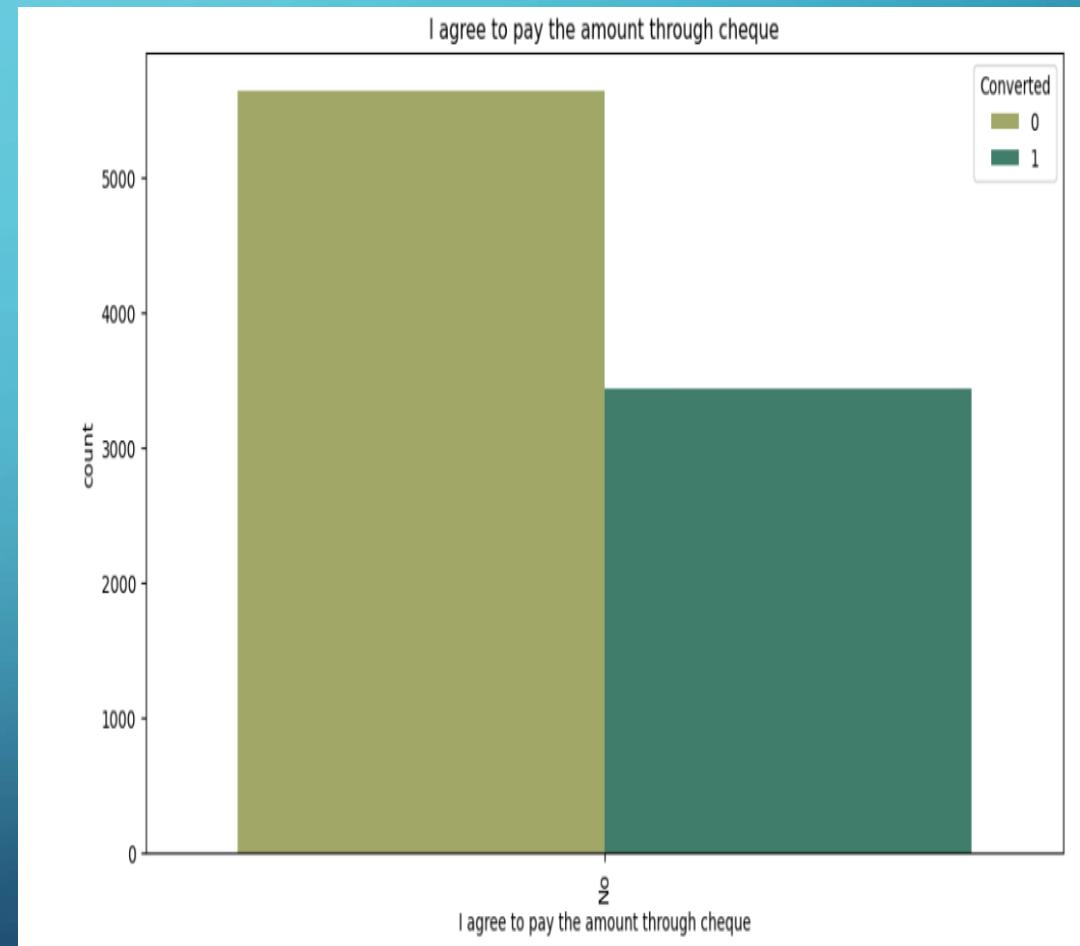
# LAST ACTIVITY - MERGE THE FIELDS WITH LESS COUNTS AS "OTHER\_ACTIVITY"



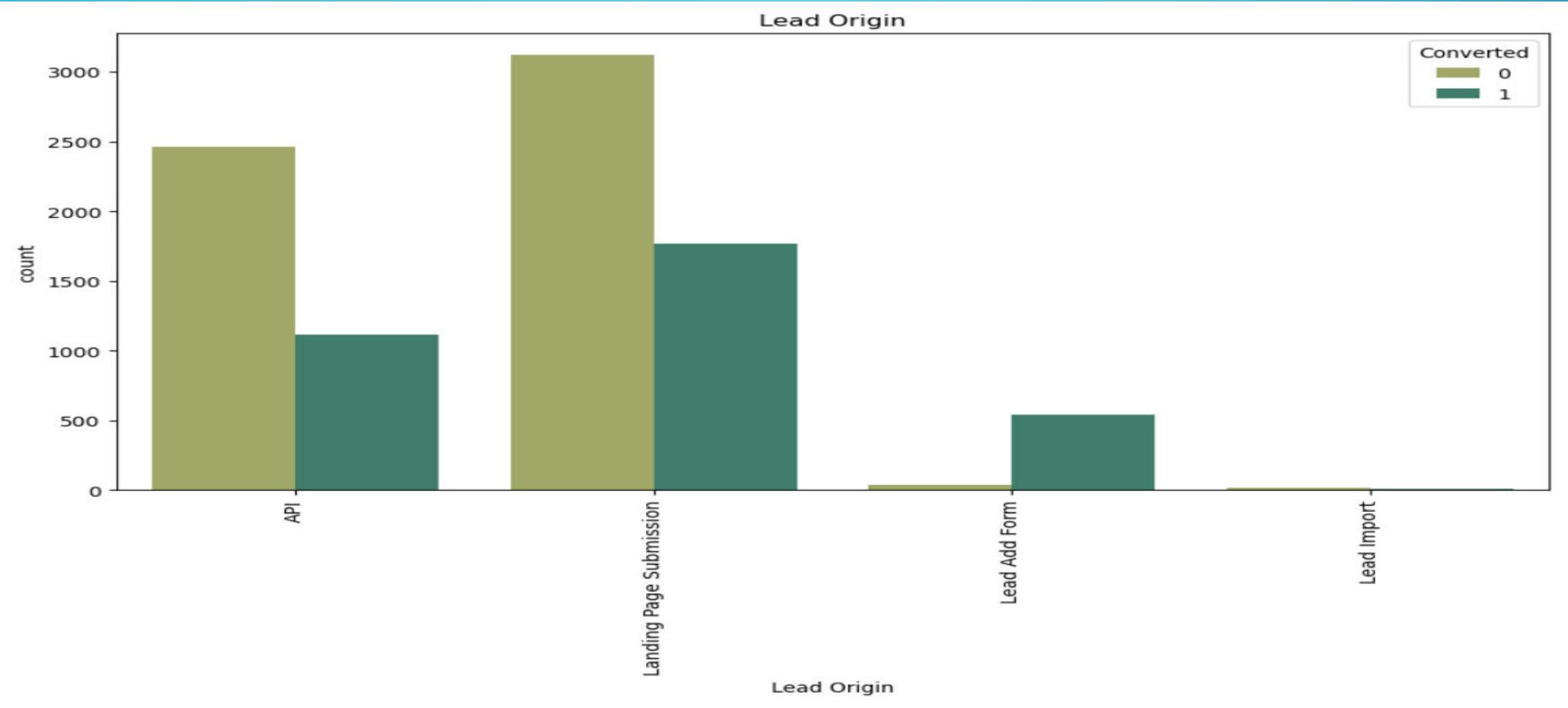
# DO NOT EMAIL



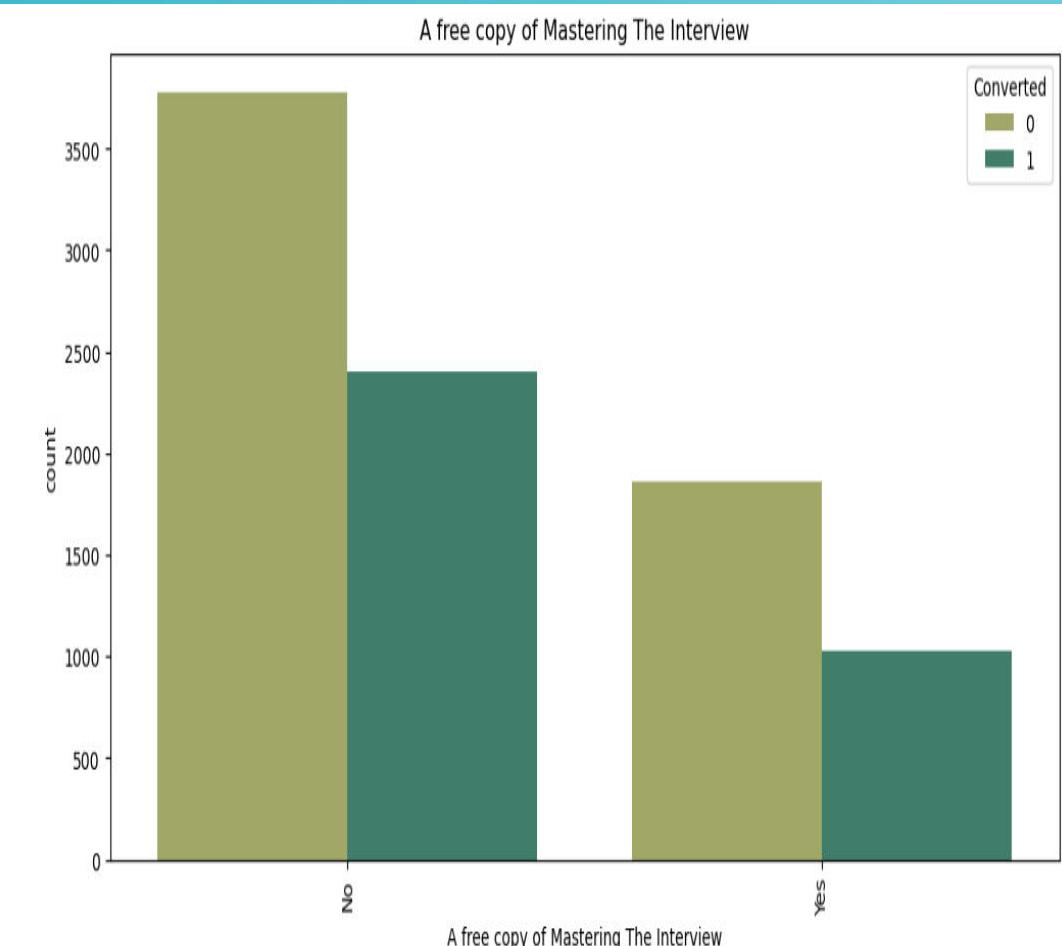
# I AGREE TO PAY THE AMOUNT THROUGH CHEQUE



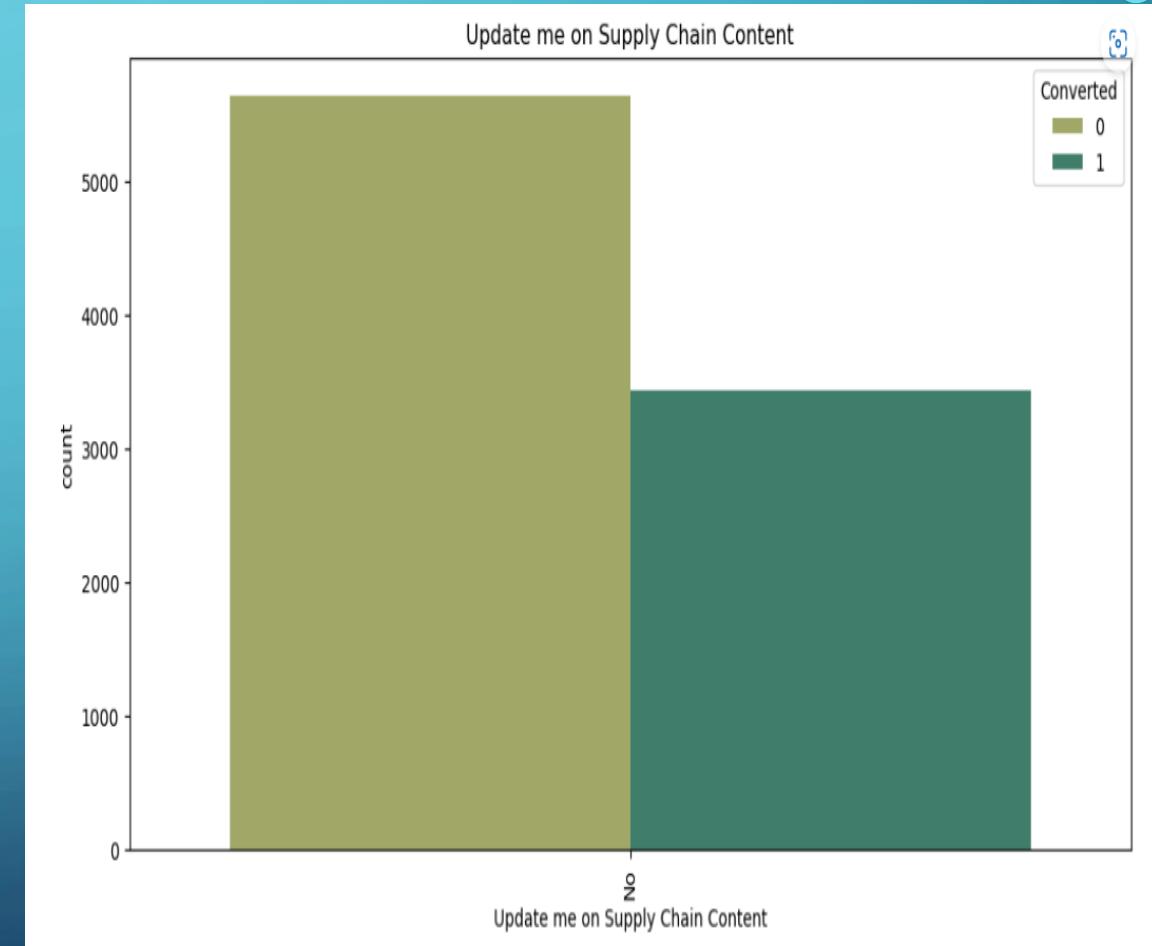
# LEAD ORIGIN



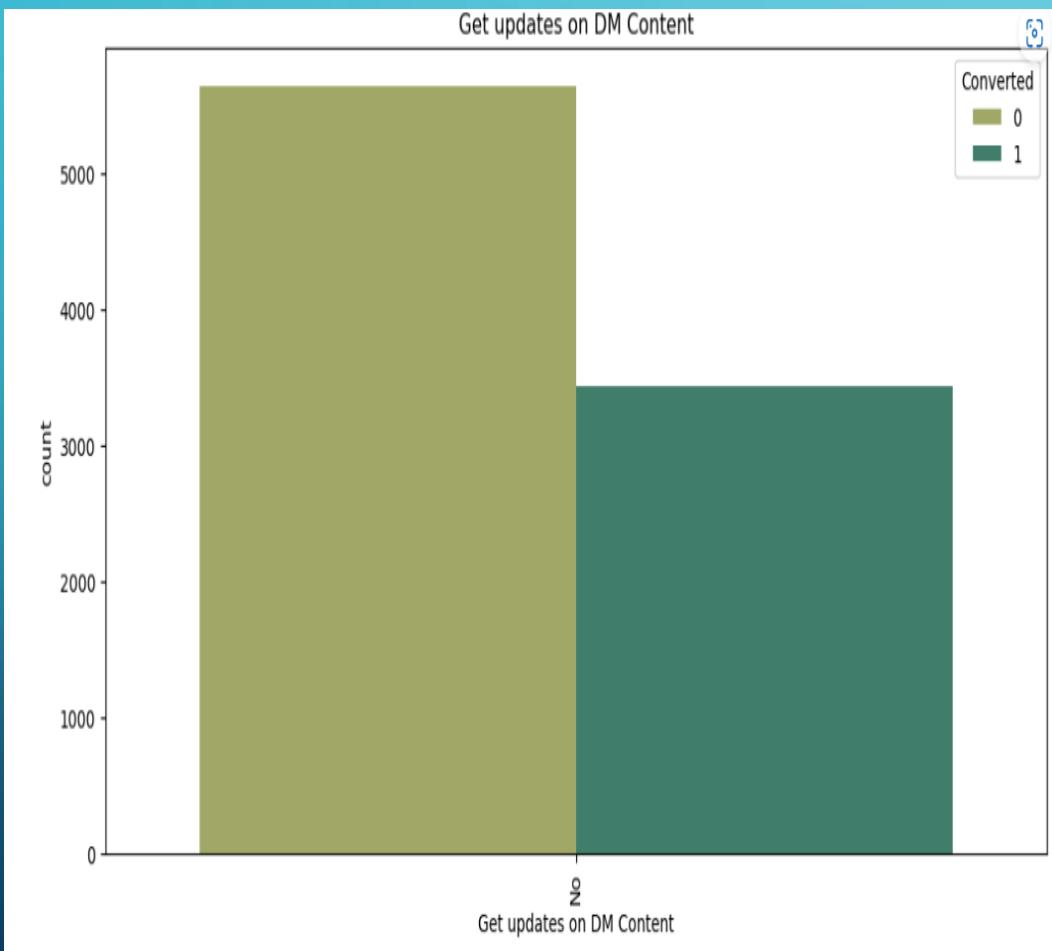
# A FREE COPY OF MASTERING THE INTERVIEW



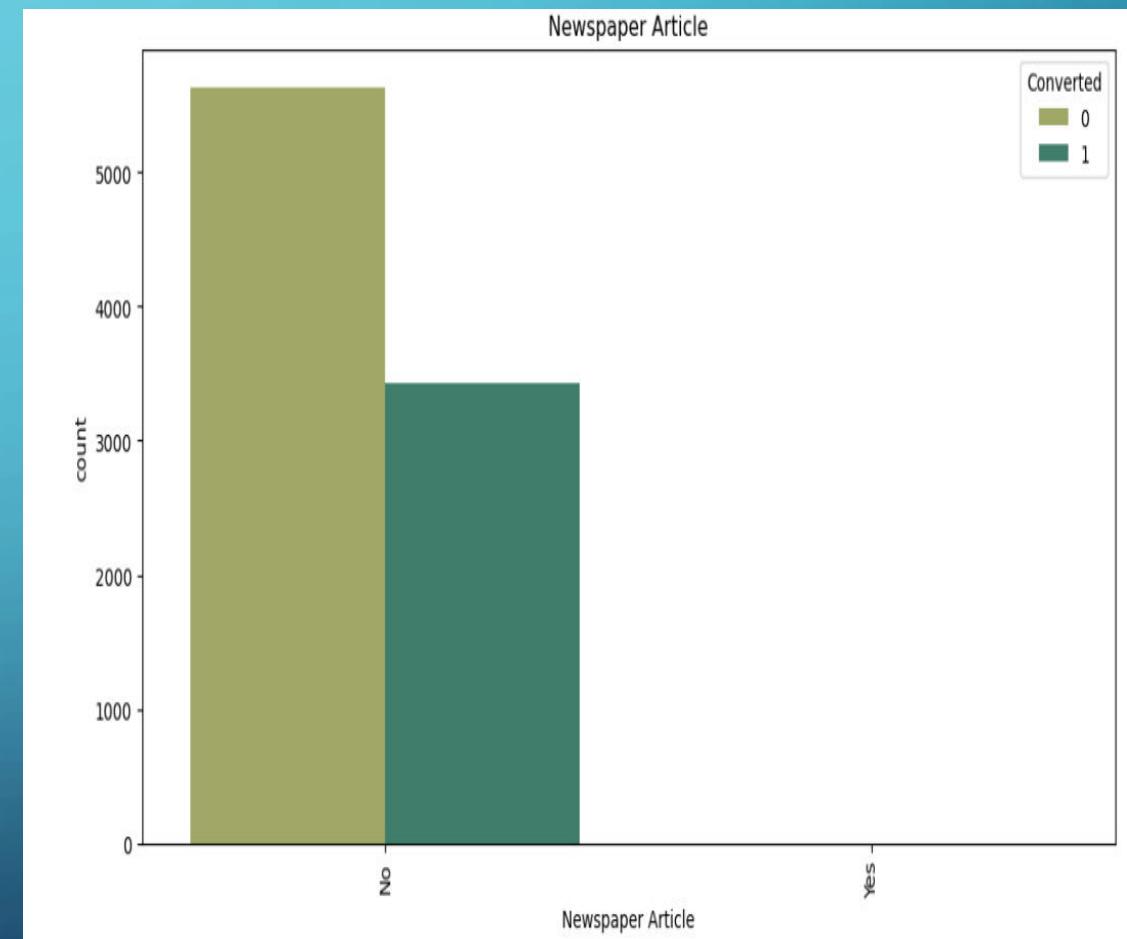
# UPDATE ME ON SUPPLY CHAIN CONTENT



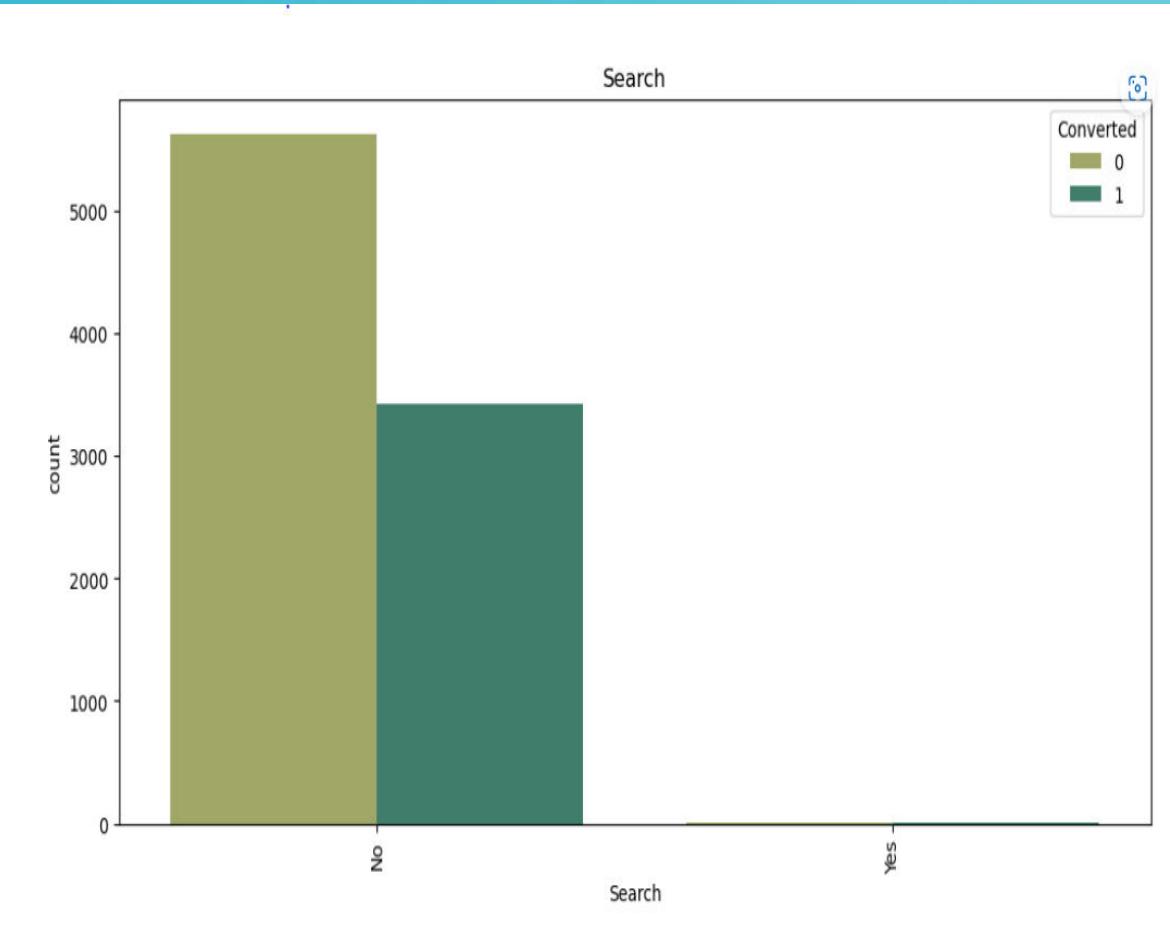
# GET UPDATES ON DM CONTENT



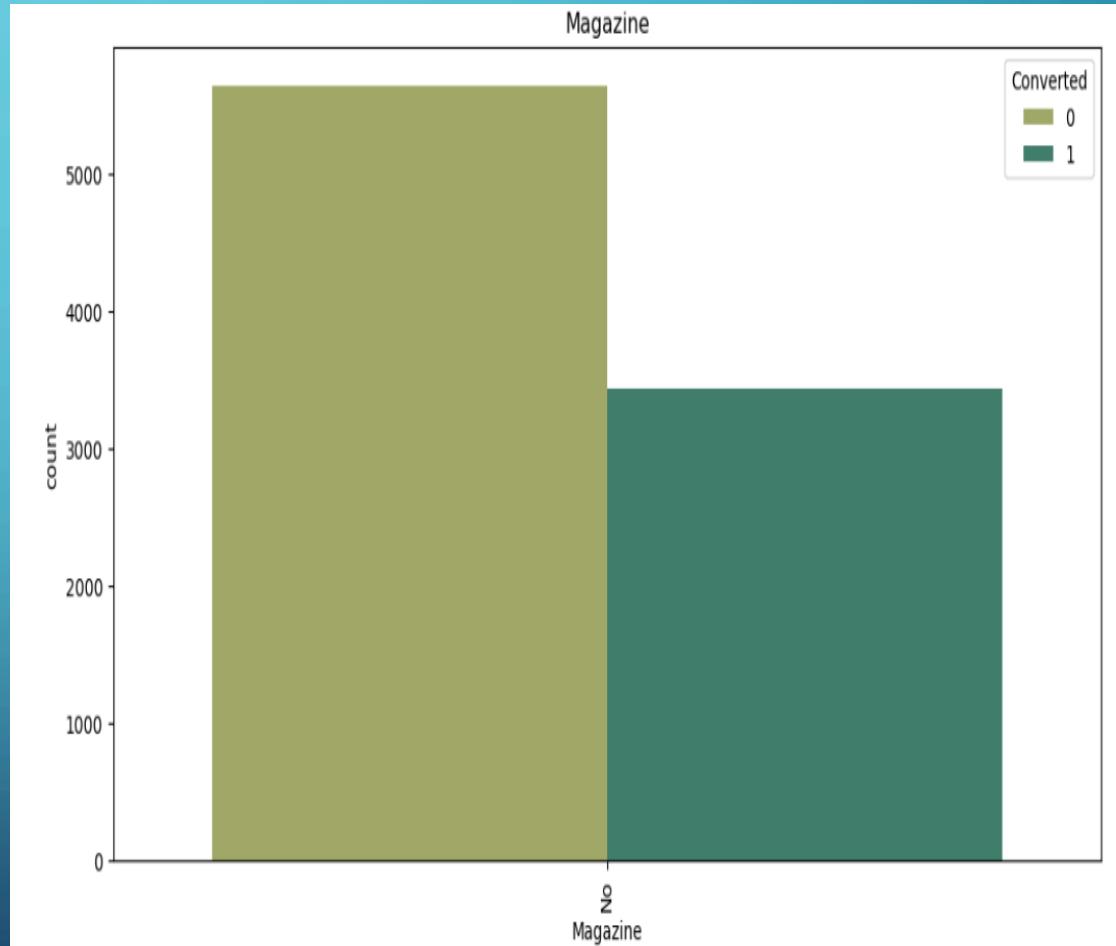
# NEWSPAPER ARTICLE



# SEARCH

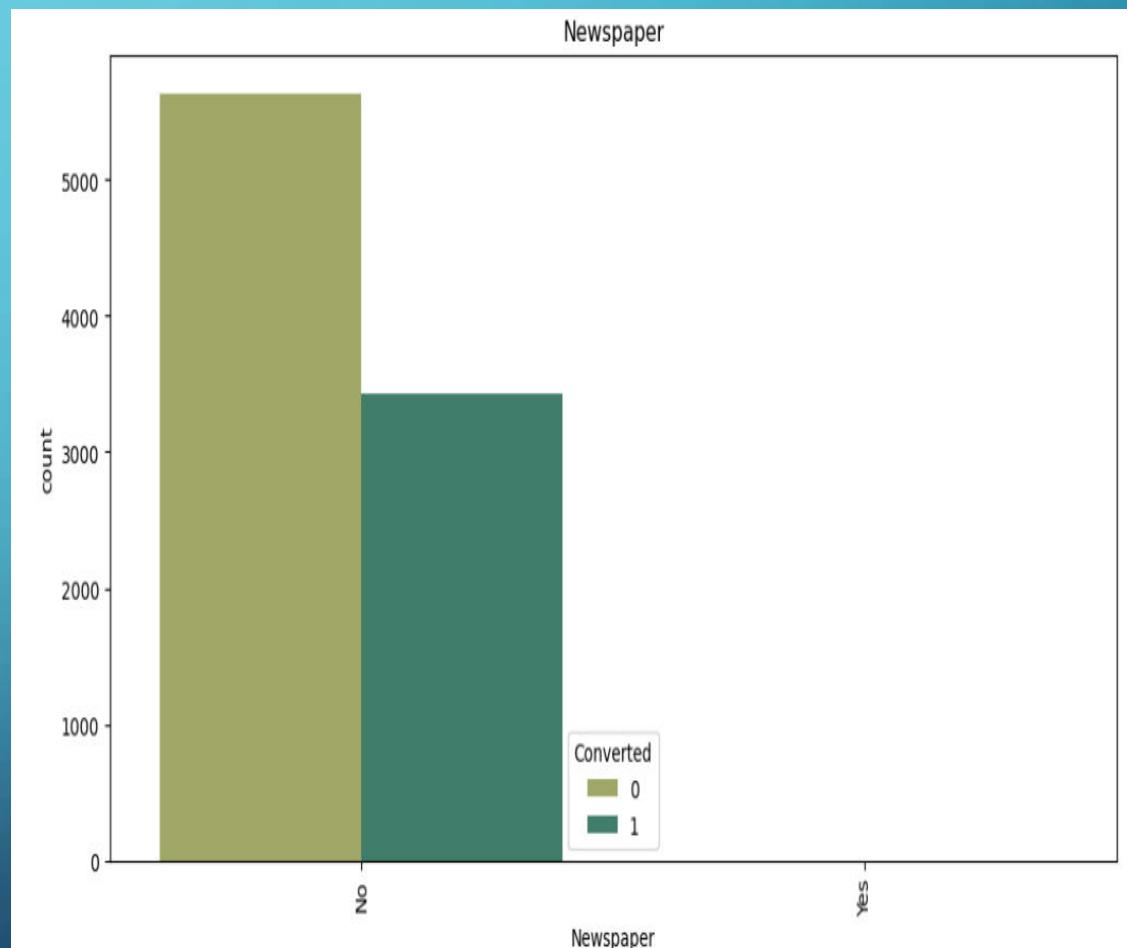
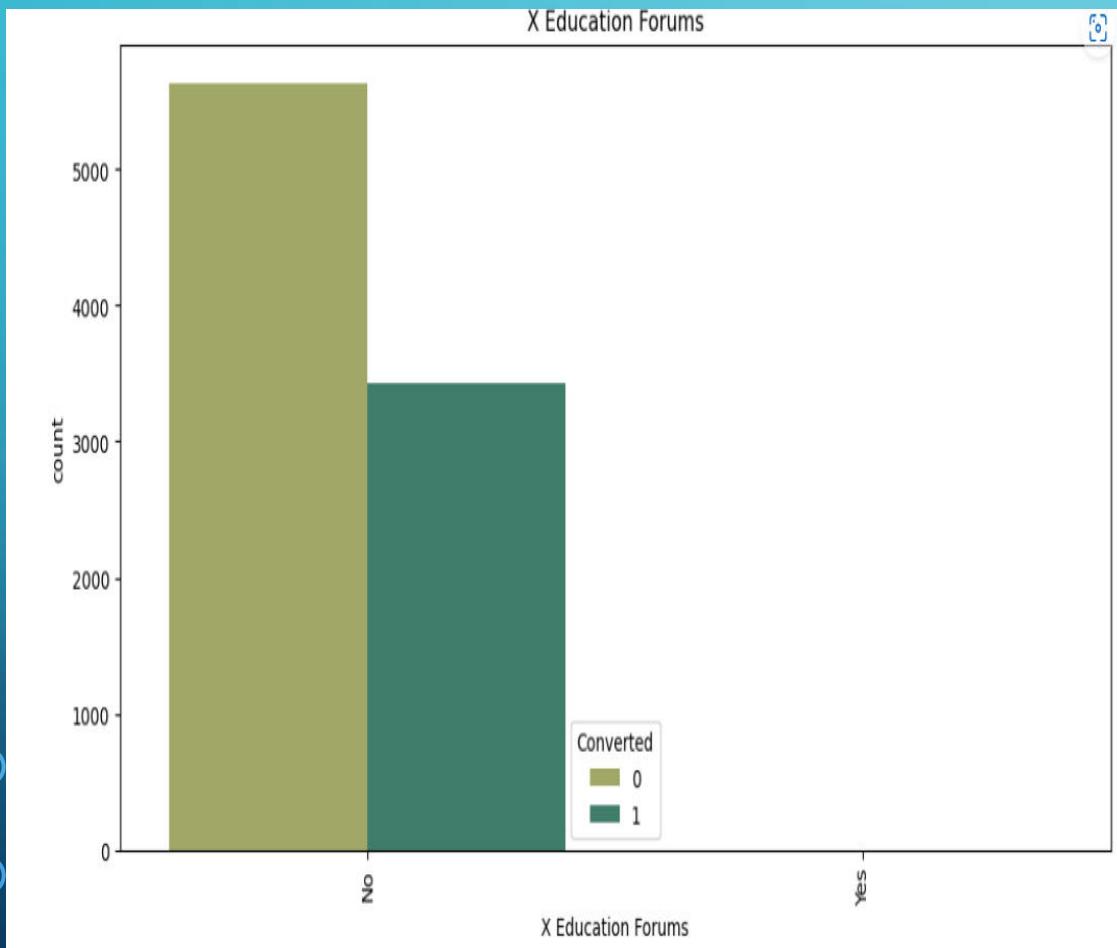


# MAGAZINE



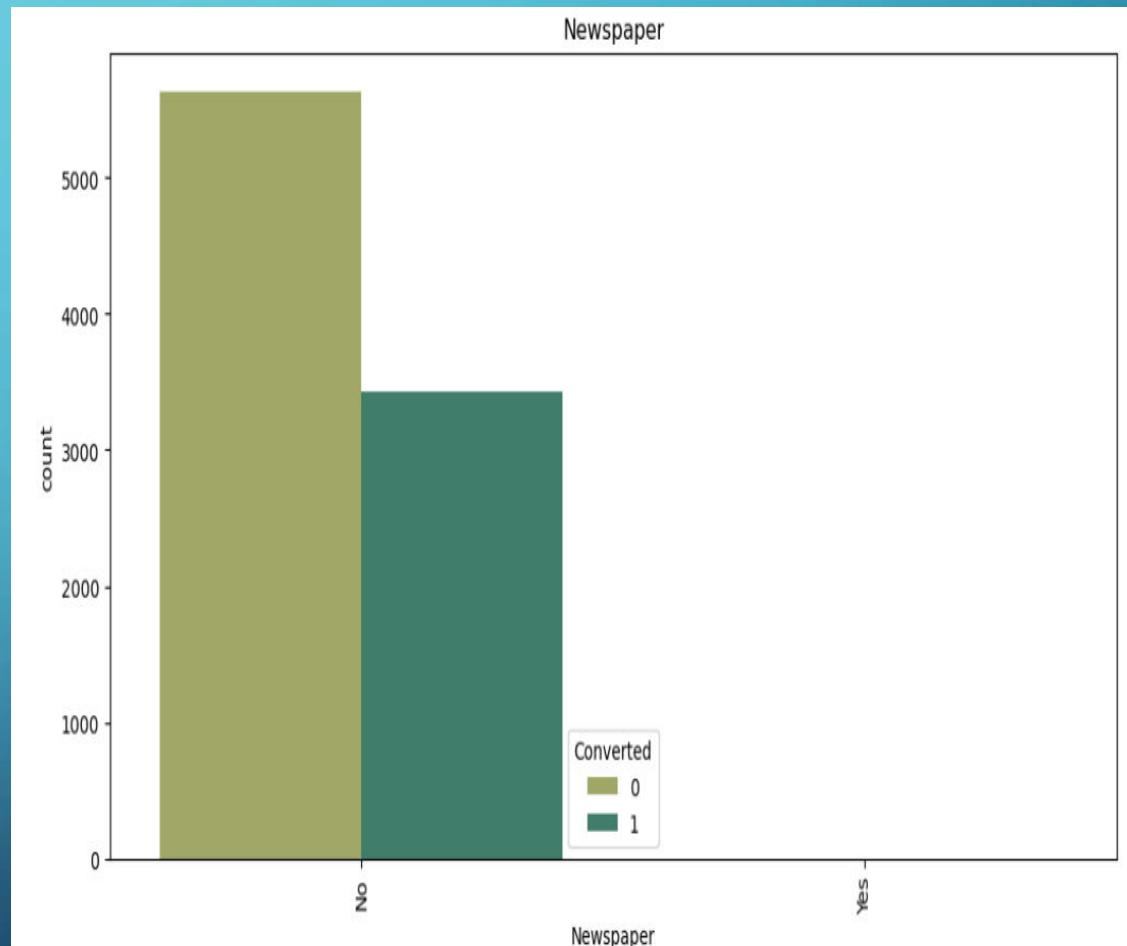
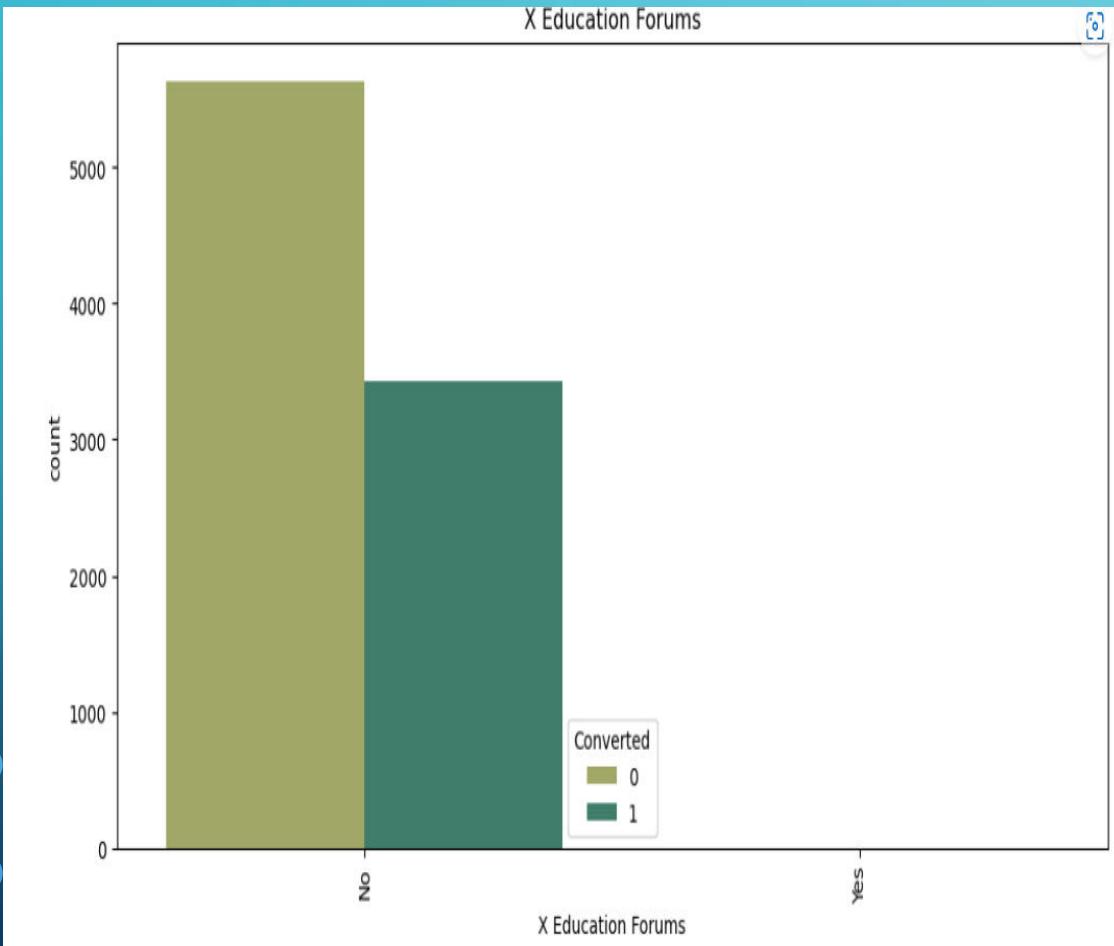
# X EDUCATION FORUMS

# NEWSPAPER

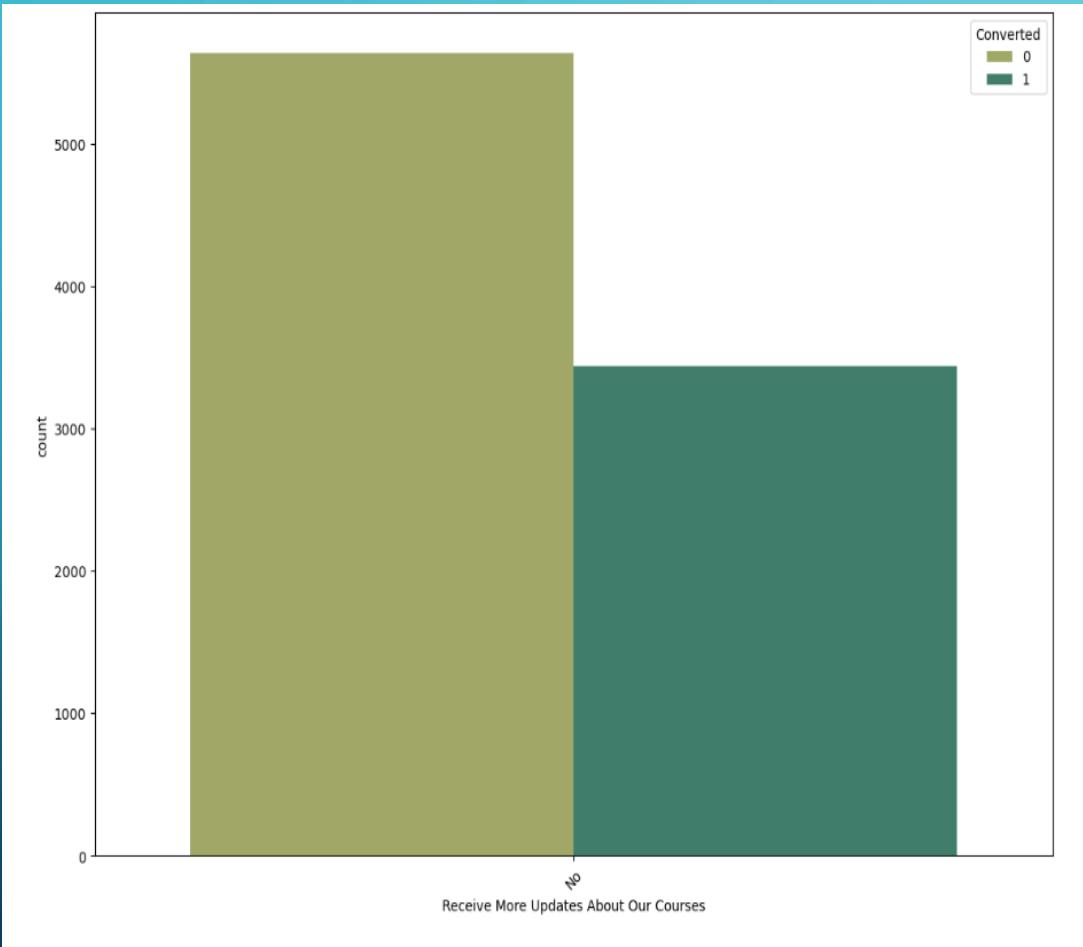


# X EDUCATION FORUMS

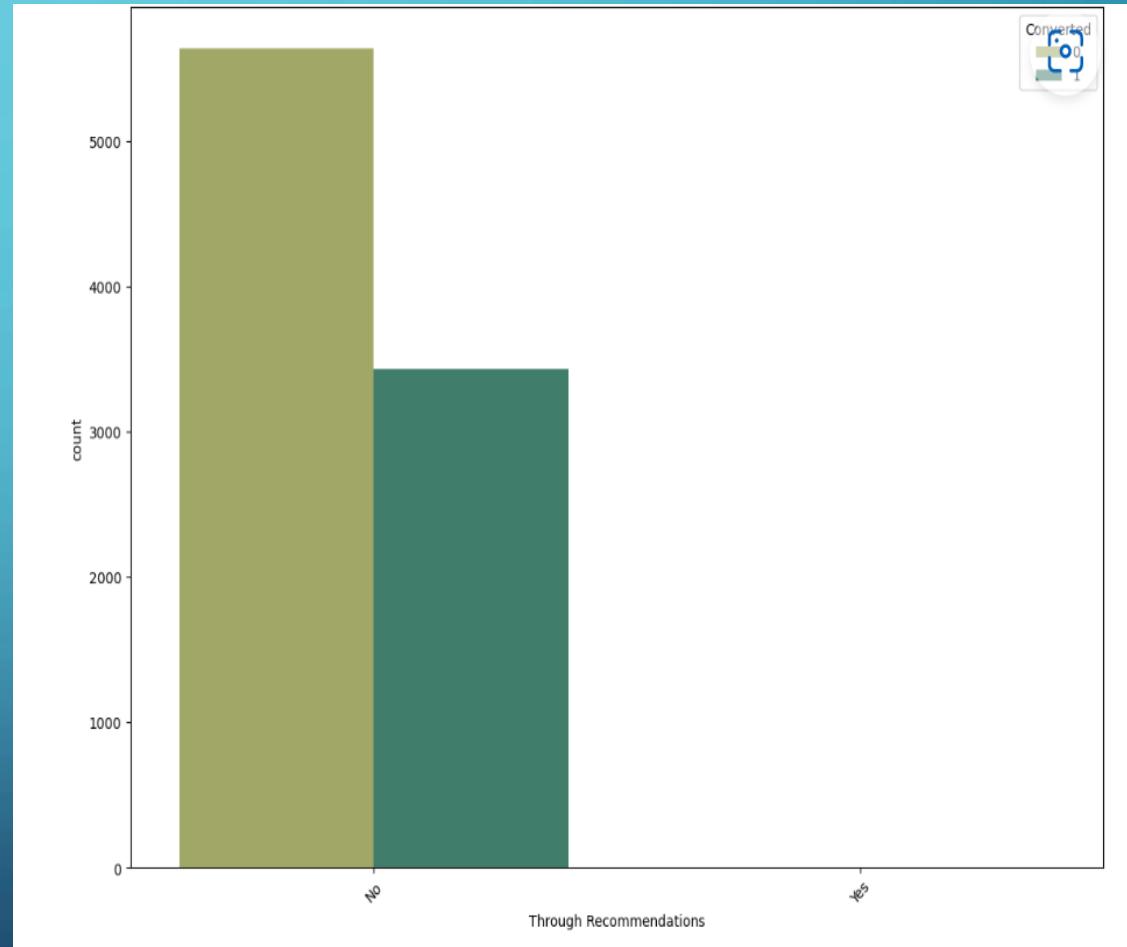
# NEWSPAPER



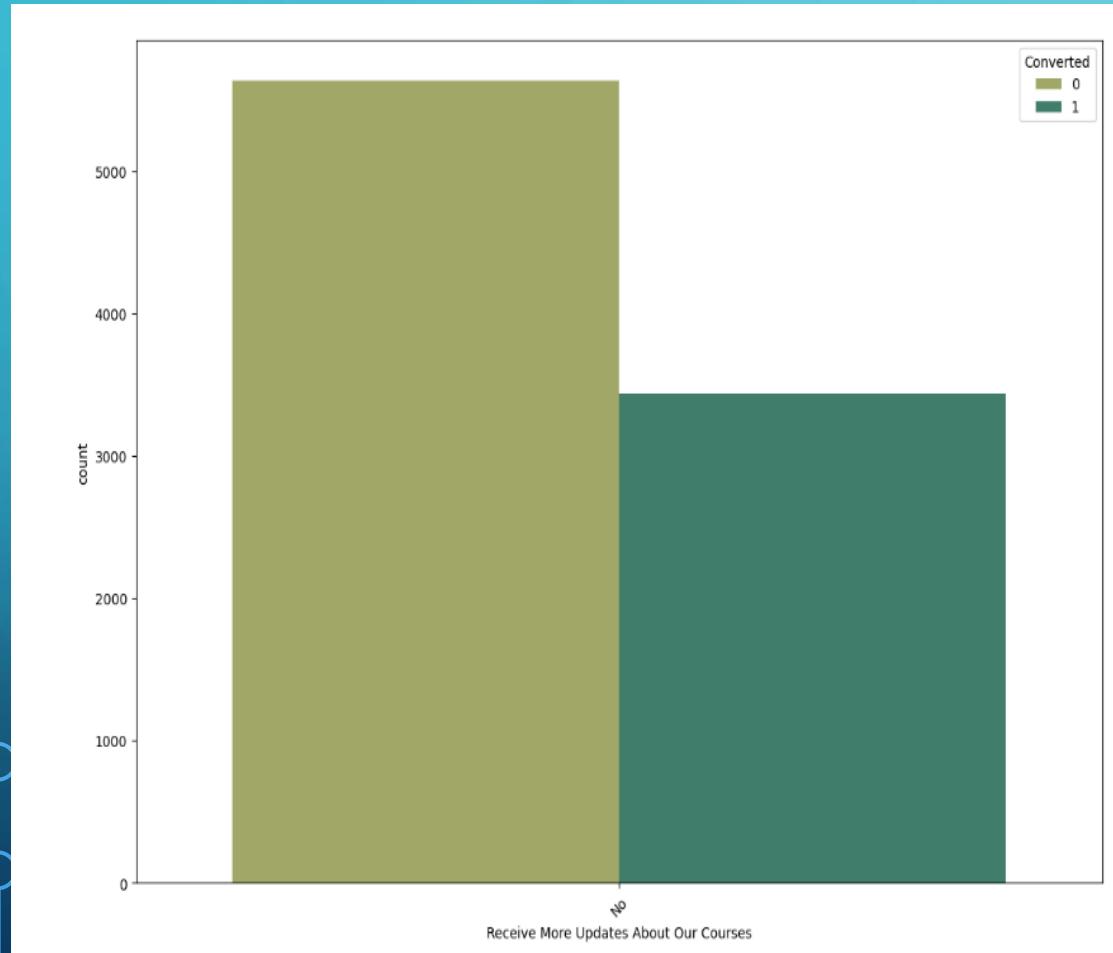
RECEIVE MORE UPDATES  
ABOUT OUR COURSES



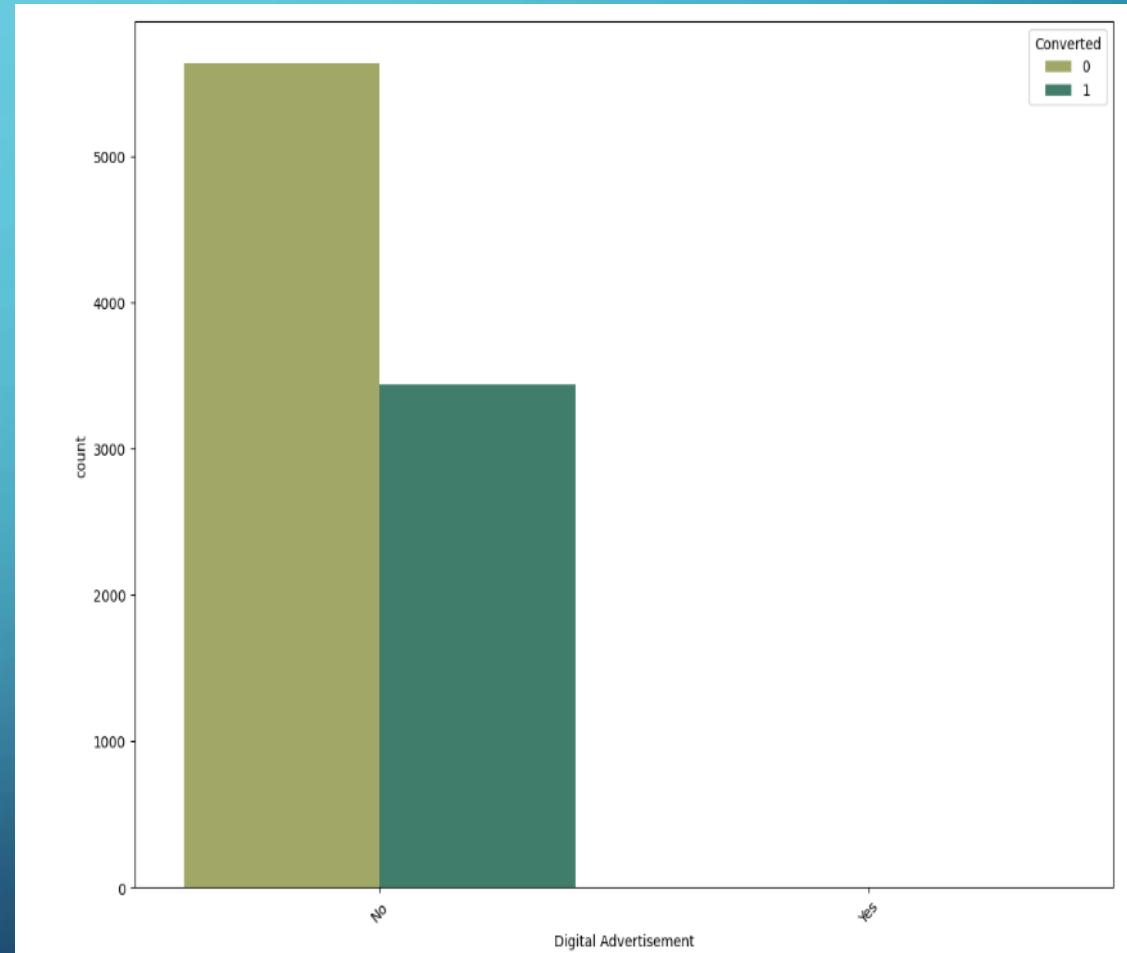
THROUGH  
RECOMMENDATIONS



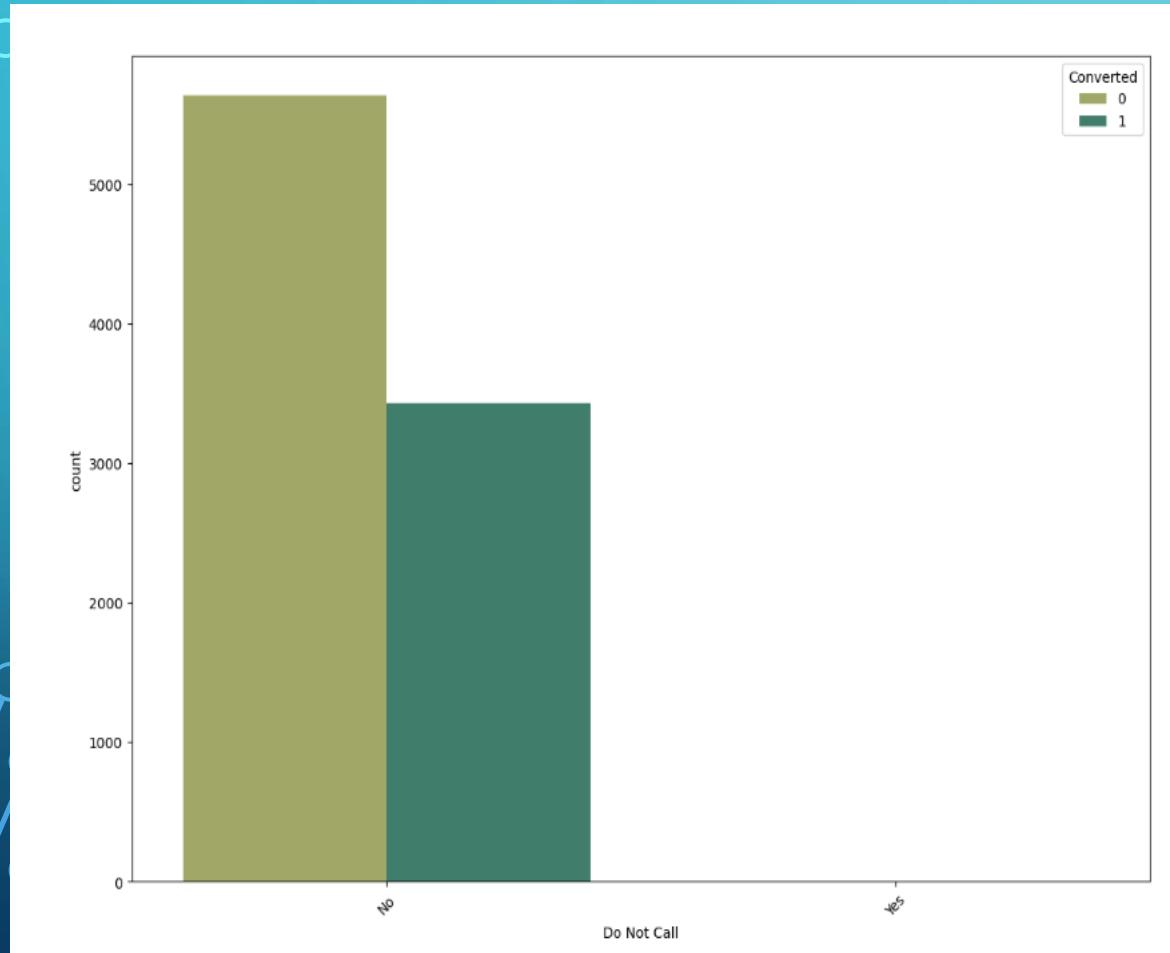
RECEIVE MORE UPDATES  
ABOUT OUR COURSES



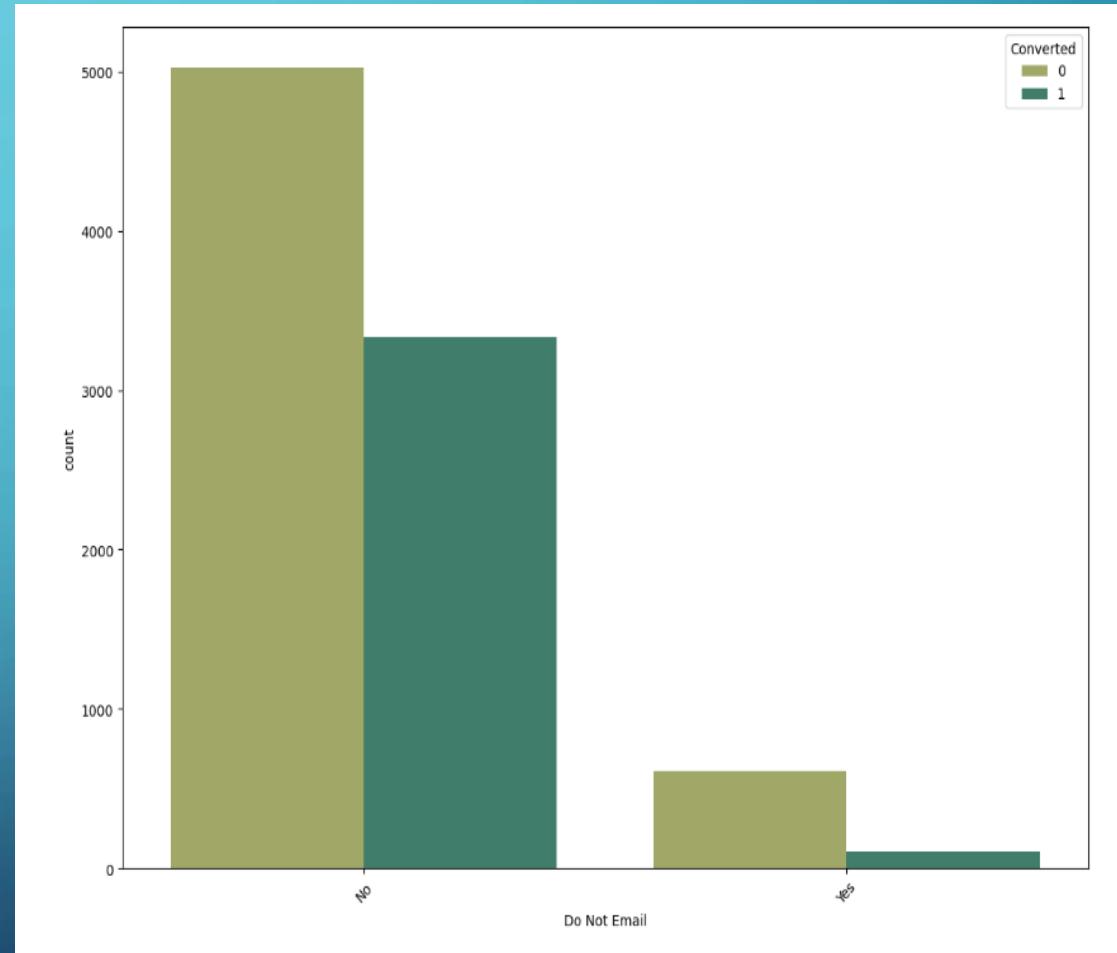
DIGITAL ADVERTISEMENT



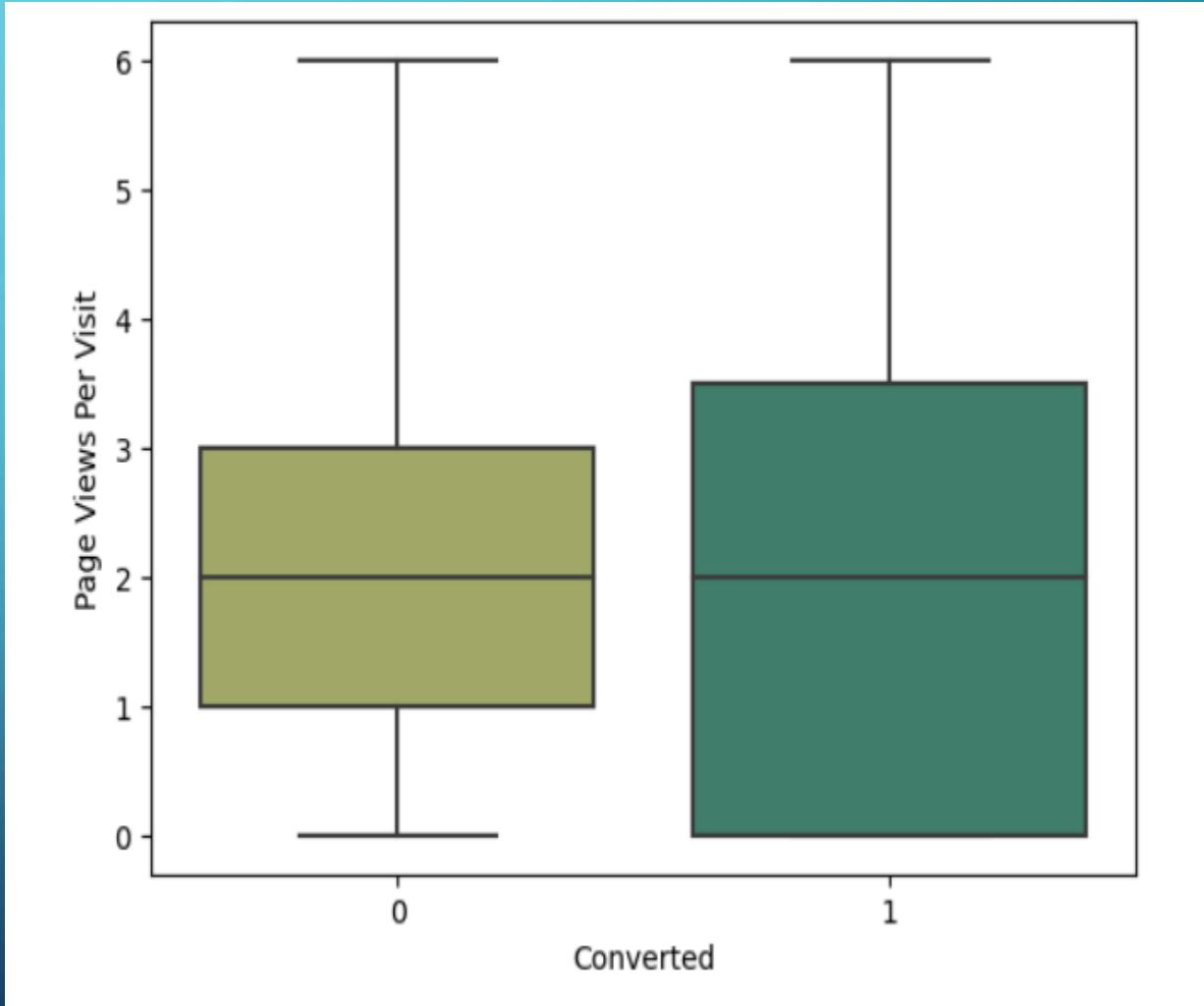
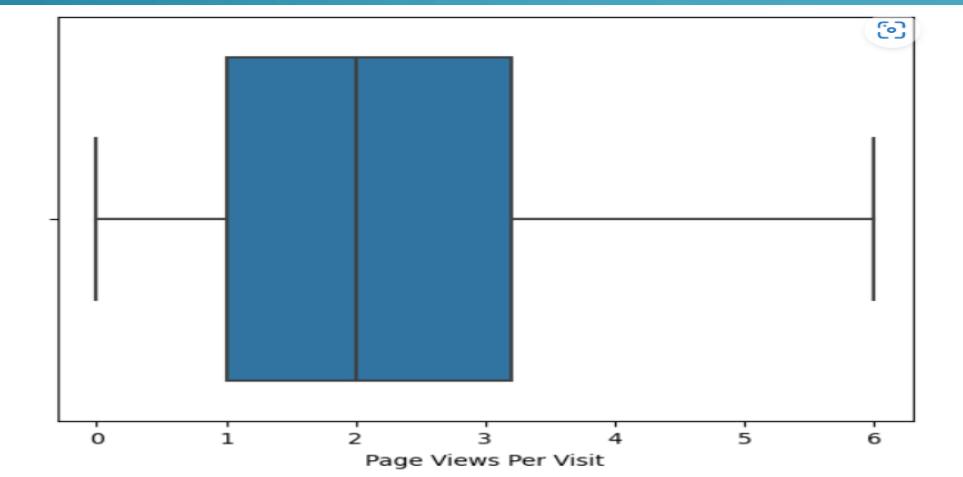
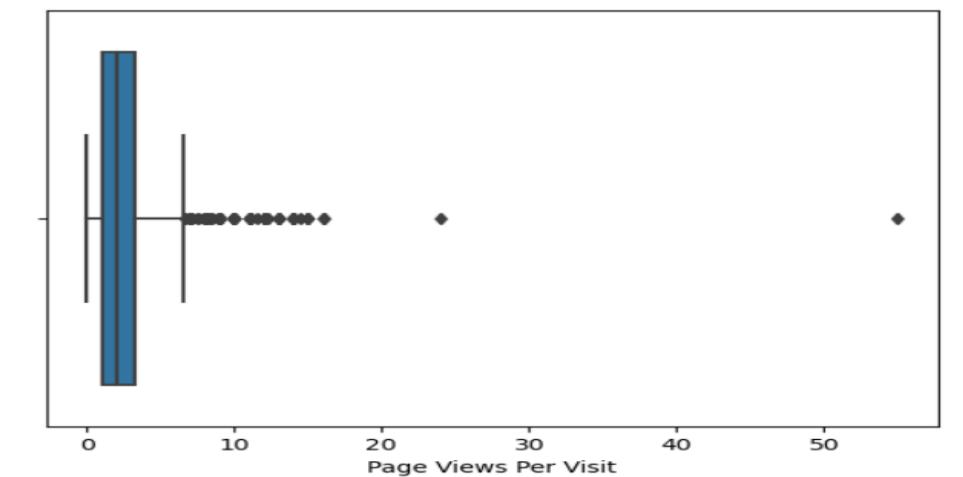
# DO NOT CALL



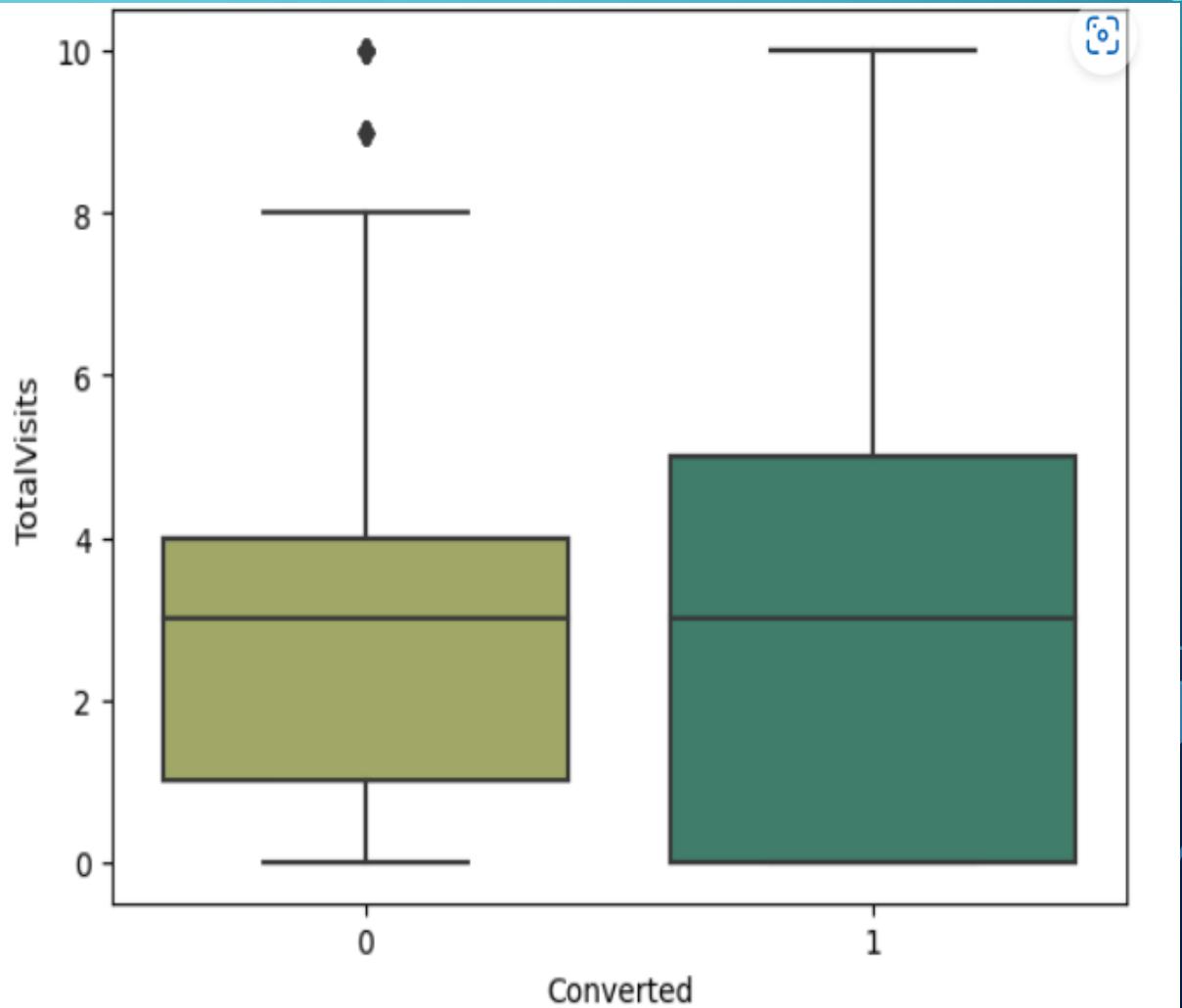
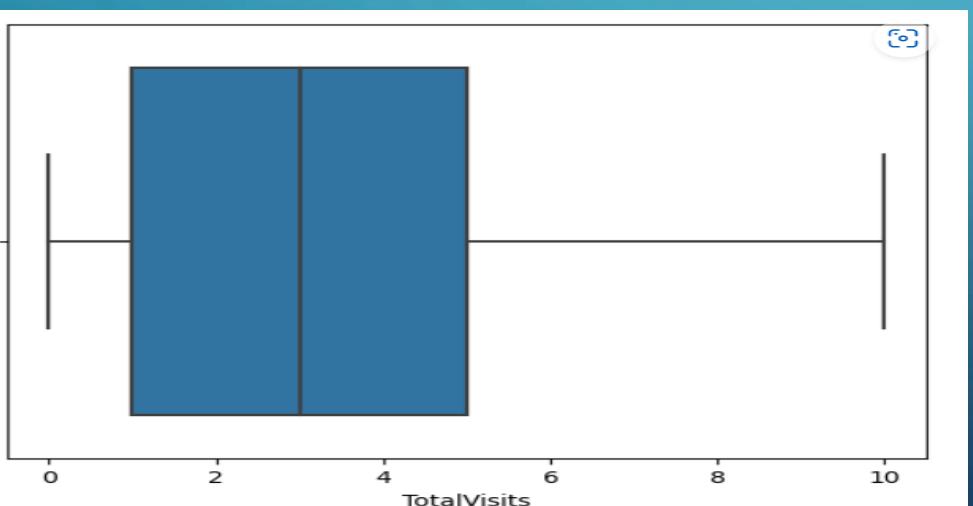
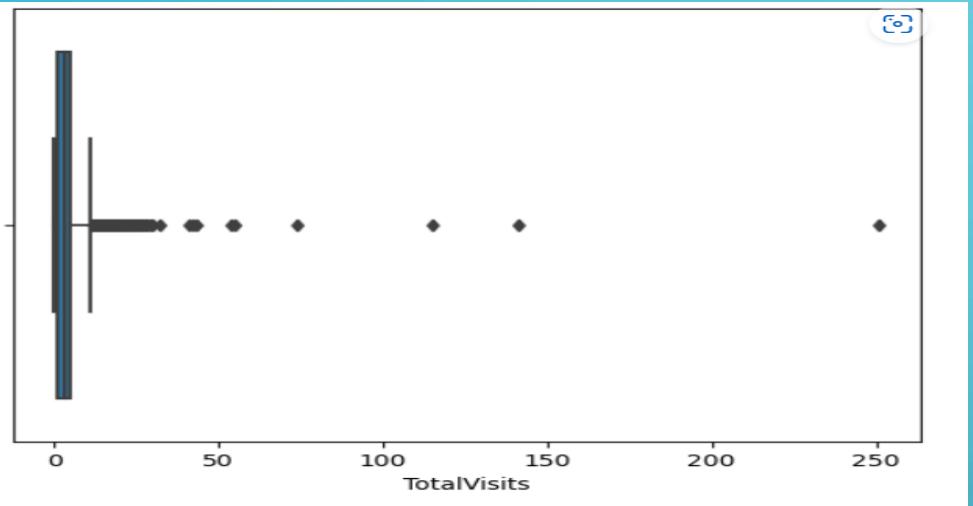
# DO NOT EMAIL



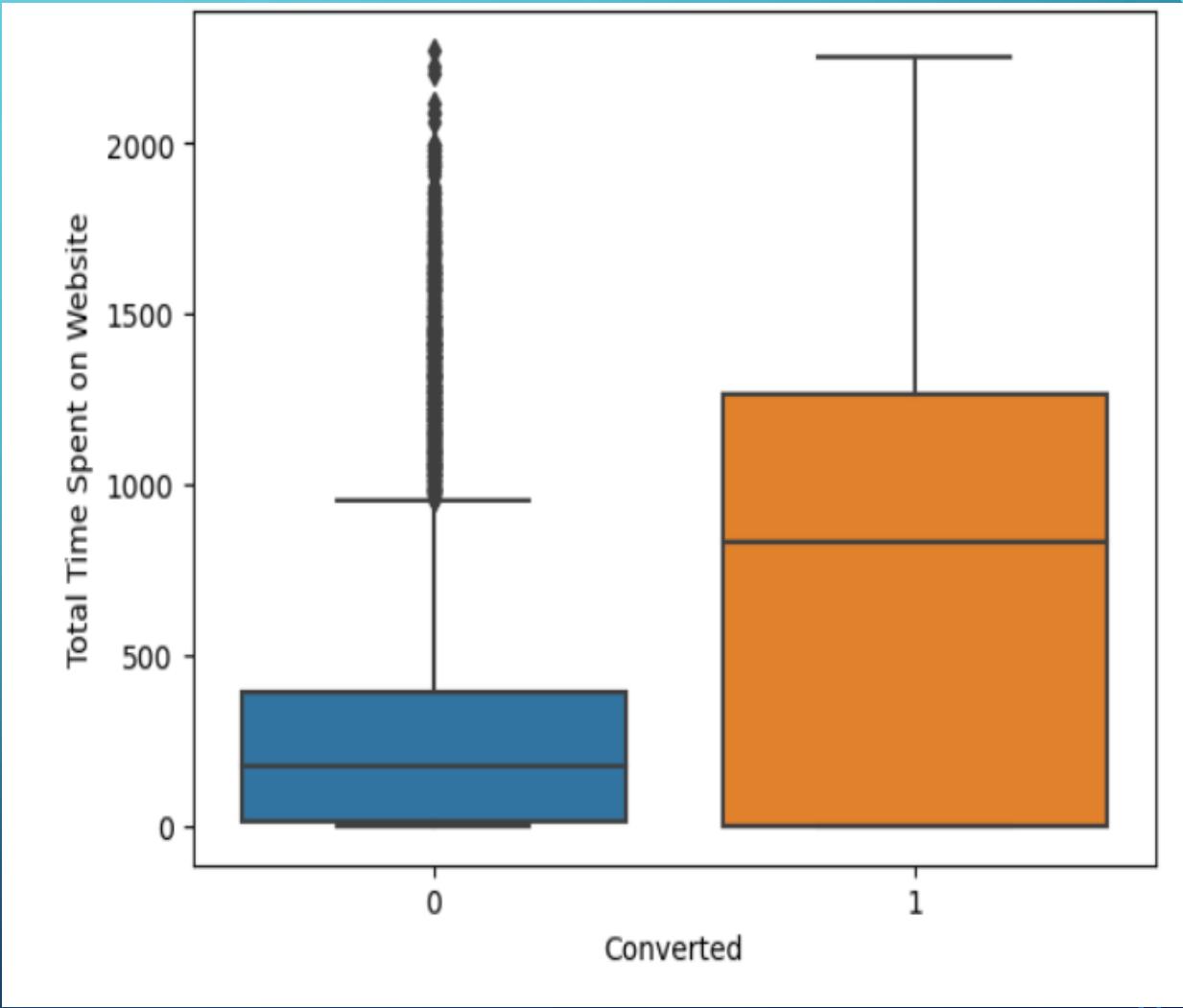
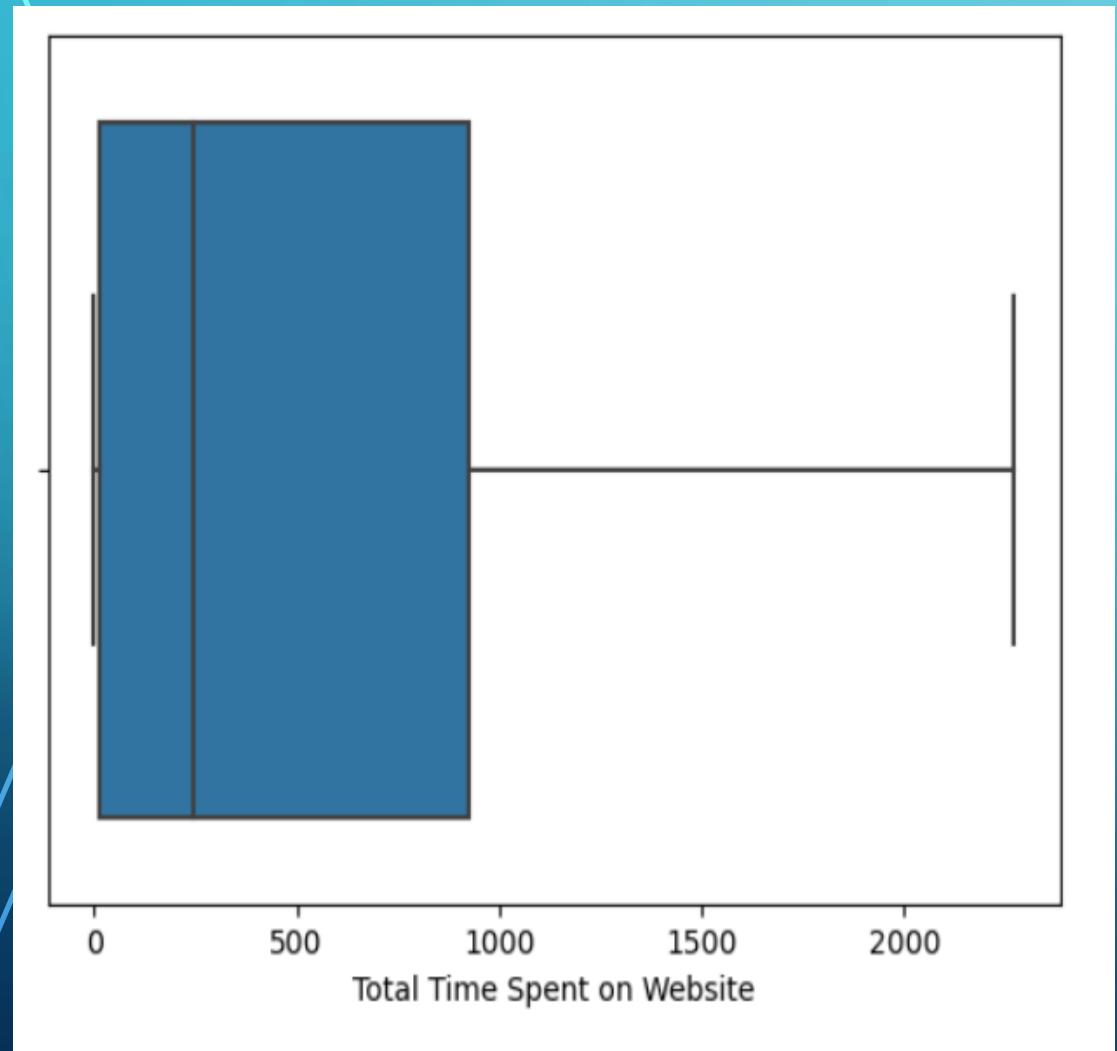
# PAGE VIEWS PER VISIT – REMOVE OUTLIERS



## TOTAL VISITS – REMOVE OUTLIERS



# TOTAL TIME SPENT ON WEBSITE



# DATA PREPARATION FOR MODELLING

# DATA PREPARATION FOR MODELING

- Binary conversion to 1 and 0

```
: 1 vars = ['Do Not Email', 'Do Not Call']
 2
 3 def binary_map(x):
 4     return x.map({'Yes': 1, "No": 0})
 5
 6 df[vars] = df[vars].apply(binary_map)
```

- Dummy Variables creation for categorical variables

```
: 1 # Dummy Variables creation categorical variables
 2 dummy_vari= pd.get_dummies(df[['Lead Origin', 'Lead Source', 'Last Activity', 'Specialization','What is your current occupat
   'City','Last Notable Activity']]], drop_first=True)
 3 dummy_vari.head()
 4
```

- Dropping variables for which dummies created

```
1 df = df.drop(['Lead Origin', 'Lead Source', 'Last Activity', 'Specialization','What is your current occupation',
2                           'City','Last Notable Activity'], axis = 1)
```

# MODEL BUILDING

# MODEL BUILDING

- Splitting the data in to training and testing dataset
- Perform regression on train – test split. Chosen 70:30 ratio
- Use RFE for feature selection
- Building model by removing the variables whose p-value and VIF is greater than 0.05 and 5 respectively
- For creation for model steps were followed by removing the attributes with high values to almost zero
- With vif ranging below 3 total 14 models were made and 14th model was considered as the final model. Also, with 0 p-value and VIF value below 3
- Prediction on test data set
- Calculate Accuracy, Sensitivity and Specificity

- Split the data for training and testing and start model training

```
1 X = df.drop(['Prospect ID','Converted'], axis=1)  
2 y = df['Converted']
```

```
1 X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=0.7, test_size=0.3, random_state=100)
```

```
1 scaler = StandardScaler()  
2 X_train[['TotalVisits','Total Time Spent on Website','Page Views Per Visit']] = scaler.fit_transform(X_train[['TotalVisits',  
3 X_train.head()
```

## RFE Modeling

```
: 1 from sklearn.feature_selection import RFE  
: 1 logreg = LogisticRegression()  
: 2 rfe = RFE(logreg, step=20)  
: 3 rfe = rfe.fit(X_train, y_train)  
: 1 rfe.support_  
: 1 cols = X_train.columns[rfe.support_]  
2 cols
```

## Building the Model

```
: 1 X_train_sm = sm.add_constant(X_train[cols])  
2 logm1 = sm.GLM(y_train, X_train_sm, family = sm.families.Binomial())  
3 result = logm1.fit()  
4 result.summary()  
: 1 col1 = cols.drop('What is your current occupation_Housewife')
```

## Rebuilding the model

```
1 X_train_sm = sm.add_constant(X_train[col1])
2 logm2 = sm.GLM(y_train,X_train_sm, family = sm.families.Binomial())
3 res = logm2.fit()
4 res.summary()
```

```
: 1 col1 = col1.drop('Last Notable Activity_Had a Phone Conversation')
```

## VIF modeling

```
1 vif = pd.DataFrame()
2 vif['Features'] = X_train[col1].columns
3 vif['VIF'] = [variance_inflation_factor(X_train[col1].values, i) for i in range(X_train[col1].shape[1])]
4 vif['VIF'] = round(vif['VIF'], 2)
5 vif = vif.sort_values(by = "VIF", ascending = False)
6 vif
```

## Drop high VIF values

```
1 col1 = col1.drop('What is your current occupation_Unemployed')
```

## Rebuilding the model

```
1 X_train_sm = sm.add_constant(X_train[col1])
2 logm9 = sm.GLM(y_train,X_train_sm, family = sm.families.Binomial())
3 res = logm9.fit()
4 res.summary()
```

## Rebuilding the model

```
1 X_train_sm = sm.add_constant(X_train[col1])
2 logm14 = sm.GLM(y_train,X_train_sm, family = sm.families.Binomial())
3 res = logm14.fit()
4 res.summary()
```

## Check VIF again

```
1 vif = pd.DataFrame()
2 vif['Features'] = X_train[col1].columns
3 vif['VIF'] = [variance_inflation_factor(X_train[col1].values, i) for i in range(X_train[col1].shape[1])]
4 vif['VIF'] = round(vif['VIF'], 2)
5 vif = vif.sort_values(by = "VIF", ascending = False)
6 vif
```

	Features	VIF
10	Specialization_Others	2.16
3	Lead Source_Olark Chat	2.03
12	Last Notable Activity_Modified	1.78
2	Lead Origin_Landing Page Submission	1.70
6	Last Activity_Olark Chat Conversation	1.59
8	Last Activity_SMS Sent	1.57
1	Total Time Spent on Website	1.29
4	Lead Source_Reference	1.24
0	Do Not Email	1.21
11	What is your current occupation_Working Profes...	1.19
5	Lead Source_Welingak Website	1.09
9	Last Activity_Unsubscribed	1.08
7	Last Activity_Other_Activity	1.01

From the above, can observed that the p values for columns are now 0 and VIF values are not too high. Hence model number 14 can be considered as the final model for evaluation

# PREDICTION ON TRAINING SET

```
: 1 y_train_pred = res.predict(x_train_sm)
 2 y_train_pred[:10]
```

```
: 1 y_train_pred = y_train_pred.values.reshape(-1)
 2 y_train_pred[:10]
```

```
: 1 y_train_pred_final = pd.DataFrame({'Converted':y_train.values, 'Converted_prob':y_train_pred})
 2 y_train_pred_final['Prospect ID'] = y_train.index
 3 y_train_pred_final.head()
```

Consider the cut off point to be 0.5

```
: 1 y_train_pred_final['predicted'] = y_train_pred_final.Converted_prob.map(lambda x: 1 if x > 0.5 else 0)
 2 y_train_pred_final.head()
```

Check the Confusion Matrix

```
: 1 con = metrics.confusion_matrix(y_train_pred_final.Converted, y_train_pred_final.predicted )
 2 print(con)
```

## Check the Accuracy

```
1 print('Accuracy :',metrics.accuracy_score(y_train_pred_final.Converted, y_train_pred_final.predicted))  
Accuracy : 0.8170366871358841
```

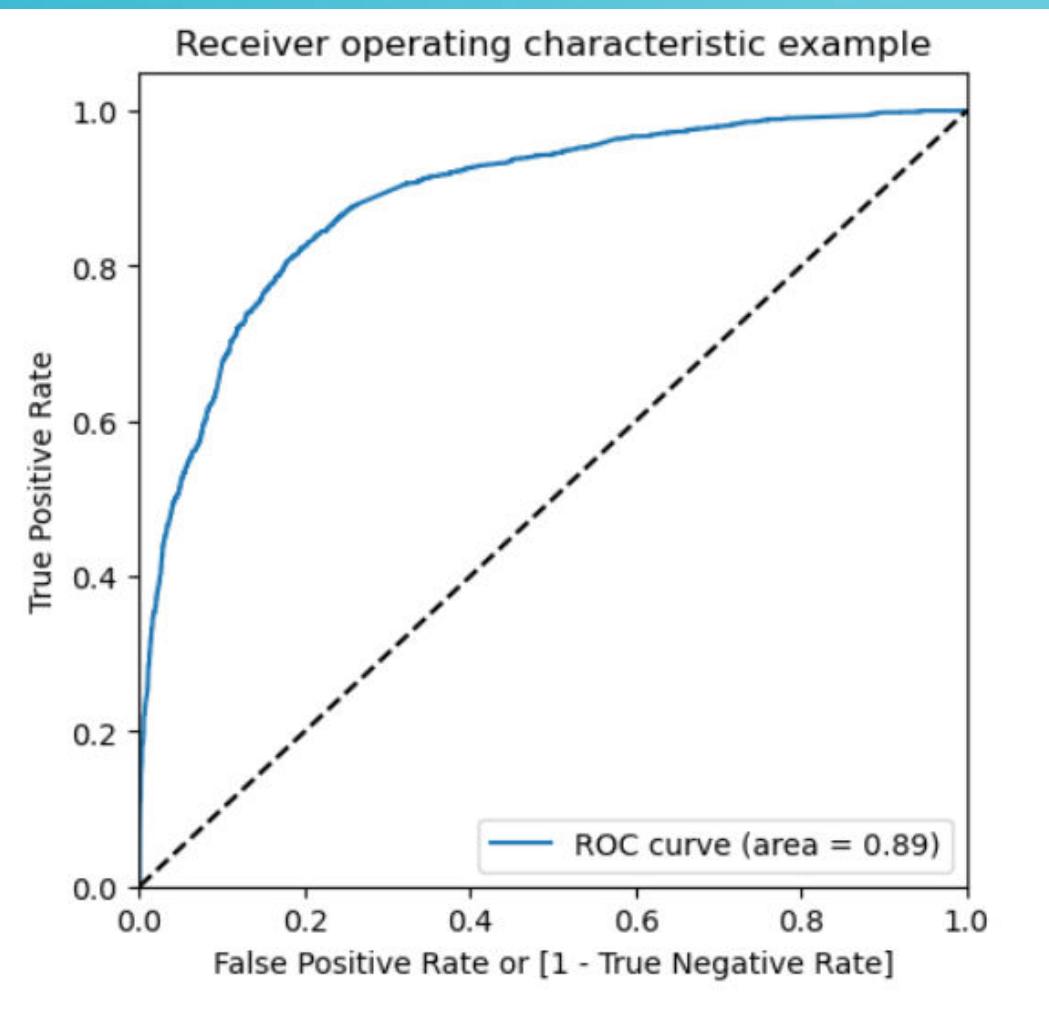
## Check TP TN FP FN and calculate Sensitivity, Specificity, False Positive Rate, Positive Predictive Value, Negative Predictive Value

```
1 # true positive  
2 TP = con[1,1]  
3 # true negatives  
4 TN = con[0,0]  
5 # false positives  
6 FP = con[0,1]  
7 # false negatives  
8 FN = con[1,0]
```

```
1 print("Sensitivity : ",TP / float(TP+FN)*100)  
2 print("Specificity : ",TN / float(TN+FP)*100)  
3 print("False Positive Rate : ",FP / float(TN+FP)*100)  
4 print("Positive Predictive Value : ",TP / float(TP+FP)*100)  
5 print ("Negative predictive value : ",TN / float(TN+ FN)*100)
```

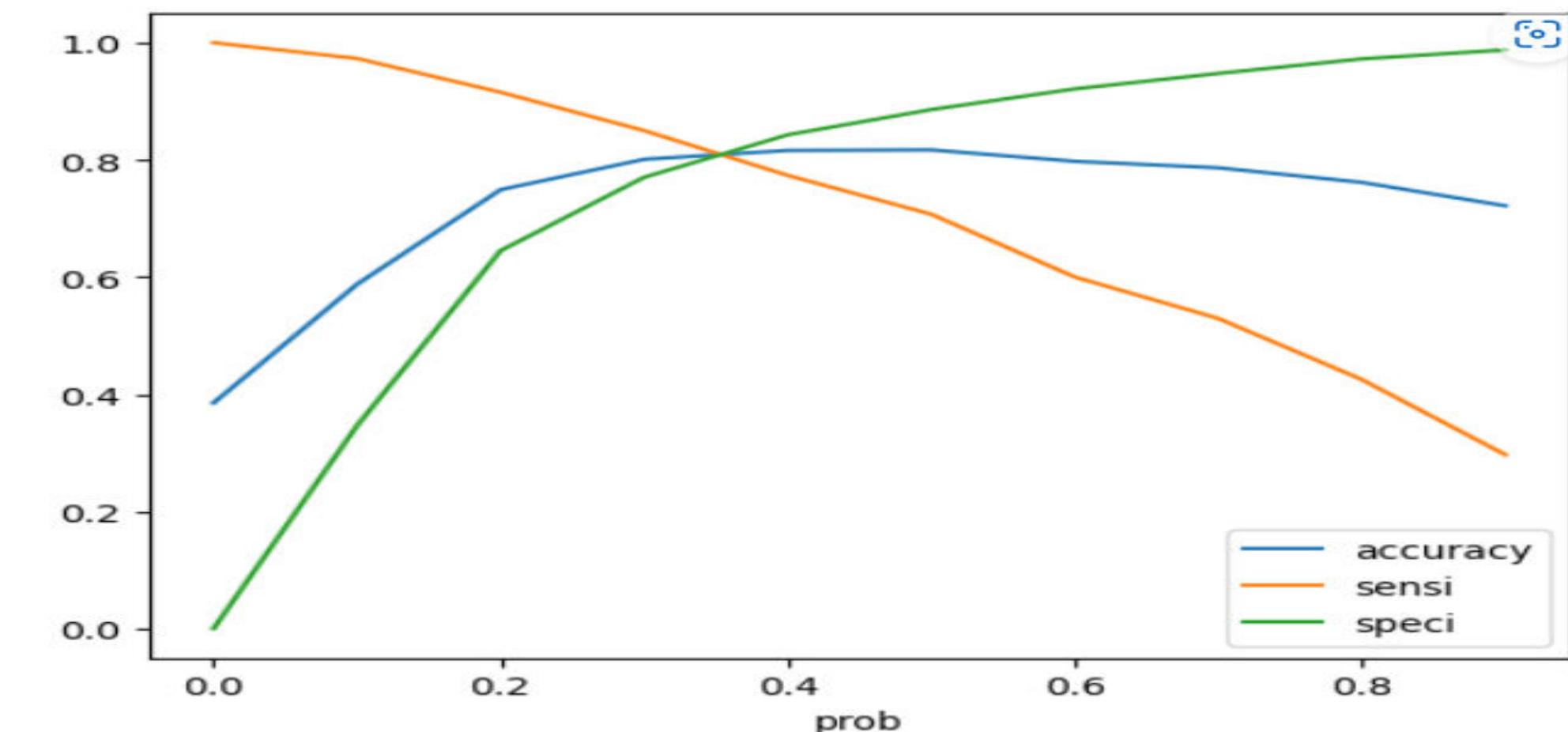
1. Our model specificity is determined to be good (88.6%) but our sensitivity is just 70.7%. As a result, this will need to be resolved.
2. Our model sensitivity was 70.7%, mostly as a result of the arbitrary cut-off point of 0.5 that we used. The ROC curve will be used to optimise this cut-off point in order to obtain a respectable value for sensitivity.

# ROC CURVE



```
1 def roc_c( actual, probs ):
2     fpr, tpr, thresholds = metrics.roc_curve( actual, probs,
3                                             drop_intermediate = False )
4
5     auc_score = metrics.roc_auc_score( actual, probs )
6
7     plt.figure(figsize=(5, 5))
8     plt.plot( fpr, tpr, label='ROC curve (area = %0.2f)' % auc_score )
9     plt.plot([0, 1], [0, 1], 'k--')
10    plt.xlim([0.0, 1.0])
11    plt.ylim([0.0, 1.05])
12    plt.xlabel('False Positive Rate or [1 - True Negative Rate]')
13    plt.ylabel('True Positive Rate')
14    plt.title('Receiver operating characteristic example')
15    plt.legend(loc="lower right")
16    plt.show()
17
18
19    return None
```

```
[1]: numbers = [float(x)/10 for x in range(10)]
for i in numbers:
    y_train_pred_final[i] = y_train_pred_final.Converted_prob.map(lambda x: 1 if x > i else 0)
y_train_pred_final.head()
```



## Now, Lets Evaluate the model again

```
: 1 print("Accuracy :",metrics.accuracy_score(y_train_pred_final.Converted, y_train_pred_final.final_predicted)*100)
Accuracy : 81.10533774208785

: 1 confusion_mat = metrics.confusion_matrix(y_train_pred_final.Converted, y_train_pred_final.final_predicted )
2 confusion_mat
: array([[3148,  757],
       [ 443, 2003]], dtype=int64)

: 1 # true positive
2 TP = confusion_mat[1,1]
3 # true negatives
4 TN = confusion_mat[0,0]
5 # false positives
6 FP = confusion_mat[0,1]
7 # false negatives
8 FN = confusion_mat[1,0]

: 1 print("Sensitivity : ",TP / float(TP+FN)*100)
2 print("Specificity : ",TN / float(TN+FP)*100)
3 print("False Positive Rate : ",FP/ float(TN+FP)*100)
4 print("Positive Predictive Value : ",TP / float(TP+FP)*100)
5 print ("Negative predictive value : ",TN / float(TN+ FN)*100)

Sensitivity : 81.88879803761243
Specificity : 80.6145966709347
False Positive Rate : 19.3854033290653
Positive Predictive Value : 72.57246376811595
Negative predictive value : 87.66360345307713
```

From the above Sensitivity and Specificity are 81.9% and 80.6% respectively.

# PREDICTION ON TEST SET

```
1 X_test[['TotalVisits','Total Time Spent on Website','Page Views Per Visit']] = scaler.transform(X_test[['TotalVisits',
2                                         'Total Time Spent on Website',
3                                         'Page Views Per Visit']])
4
5 X_test = X_test[col1]
6 X_test.head()
```

```
1 y_test_df = pd.DataFrame(y_test)
2
3 y_test_df[ 'ID' ] = y_test_df.index
4 y_pred1.reset_index(drop=True, inplace=True)
5 y_test_df.reset_index(drop=True, inplace=True)
```

*Appending y\_test\_df and y\_pred1*

```
1 y_predict = pd.concat([y_test_df, y_pred1],axis=1)
2 y_predict
```

```
1 y_predict=y_predict.rename(columns={ 0 : 'Converted_probality'})
2 y_predict = y_predict.reindex(columns=['ID','Converted','Converted_probality'])
```

```
1 y_predict['final_predicted'] = y_predict.Converted_probality.map(lambda x: 1 if x > 0.34 else 0)
```

```
1 y_predict
```

```

1 print("Accuracy :",metrics.accuracy_score(y_predict.Converted, y_predict.final_predicted)*100)
Accuracy : 80.42600073448402

1 confusion = metrics.confusion_matrix(y_predict.Converted, y_predict.final_predicted )
2 confusion
array([[1394, 340],
       [193, 796]], dtype=int64)

1 # true positive
2 TP = confusion[1,1]
3 # true negatives
4 TN = confusion[0,0]
5 # false positives
6 FP = confusion[0,1]
7 # false negatives
8 FN = confusion[1,0]

1 print("Sensitivity : ",TP / float(TP+FN)*100)
2 print("Specificity : ",TN / float(TN+FP)*100)
3 print("False Positive Rate : ",FP/ float(TN+FP)*100)
4 print("Positive Predictive Value : ",TP / float(TP+FP)*100)
5 print ("Negative predictive value : ",TN / float(TN+ FN)*100)

Sensitivity : 80.48533872598584
Specificity : 80.3921568627451
False Positive Rate : 19.607843137254903
Positive Predictive Value : 70.07042253521126
Negative predictive value : 87.83868935097668

```

## Assinging the score to final Test data

```

1 y_predict['Lead_Score'] = y_predict.Converted_probability.map( lambda x: round(x*100))
2 y_predict.head()

```

	ID	Converted	Converted_probability	final_predicted	Lead_Score
0	3271	0	0.129393	0	13
1	1490	1	0.968762	1	97
2	7936	0	0.111751	0	11
3	4216	1	0.804505	1	80
4	3830	0	0.132210	0	13

# MODEL EVALUATION

# MODEL EVALUATION

## TRAIN DATA

- ACCURACY : 81.1 %
- SENSITIVITY : 81.8 %
- SPECIFICITY : 80.6 %

## TEST DATA

- ACCURACY : 80.4 %
- SENSITIVITY : 80.5 %
- SPECIFICITY : 80.4 %

# CONCLUSION

# CONCLUSION

- It was found that the variables that mattered the most for the potential buyers are following:
  - What matters most to you in choosing a course
  - Tags will revert after reading the email
  - Last Notable Activity modified
  - Tags: Ringing and Other
  - Last Activity: SMS sent and Olark chat conversation
  - Occupation: Working Professional
  - Lead Origin lead import

# CONCLUSION

- Keeping the above in mind, the X Education can flourish as they have very high chance to get almost all the potential buyers to change their mind and buy their courses
- Thus, we have succeeded in achieving our objective of estimating the desired lead conversion rate to be somewhere around 80%. We should be able to provide the CEO confidence to make wise decisions based on this model in order to achieve a higher lead conversion rate of 80% because the model appears to anticipate the conversion rate quite effectively



THANK YOU