## Abstract :

This project will be based on classification problem to predict customer online shopping intention. we aim to help companies to be successful in a highly competitive eCommerce environment by predict whether the customer, visiting web pages of an online shop, will end up with a purchase or not. The dataset contains 12330 observation and 18 features divided into 80% training and 20% testing.

I have used three classification models which are logistic regression, Random Forest and decision tree .The results show The random forest outperform the other model in accuracy , precision, recall and F1 Score with 88%,90%, 97%,93% respectively for not purchase . While for purchasing precision is 77%, recall is 53% ,F1 score is 62% .

## Design:

The dataset consists of 10 numerical and 8 categorical attributes. I used five attributes (Month , operating system, Browser , Visitor type , Weekend, Revenue) in Eda part and in model training. The 'Revenue' attribute is used as class label. The dataset is clean, there are no missing values, but the dataset is unbalanced. There is a risk of bias, so the analysis must take the unbalanced dataset into consideration.

## Algorithms:

I have used Random Forest, Logistic Regression and decision, the results showing in the figure below:

```
[          ]]
              precision    recall  f1-score   support

       False       0.90      0.97      0.93       985
        True       0.77      0.53      0.62       215

    accuracy                           0.89      1200
   macro avg       0.84      0.75      0.78      1200
weighted avg       0.88      0.89      0.88      1200
```

Figure 1 : Random forest performance

precision, recall and F1 Score with 90%, 97%,93% respectively for not purchase . While for purchasing the precision is 51%, recall is 62% ,F1 score is 62% .

```
              precision    recall  f1-score   support

       False       0.91      0.92      0.91       985
        True       0.61      0.56      0.58       215

    accuracy                           0.86      1200
   macro avg       0.76      0.74      0.75      1200
weighted avg       0.85      0.86      0.85      1200
```

Figure 2 : Decision Tree performance

precision, recall and F1 Score with 86%,91%, 92%,91% respectively for not purchase .
While for purchasing precision is 61%, recall  is 65% ,F1 score is 58% .

```
              precision    recall  f1-score   support

       False       0.85      0.99      0.92       985
        True       0.80      0.22      0.35       215

    accuracy                           0.85      1200
   macro avg       0.83      0.61      0.63      1200
weighted avg       0.84      0.85      0.81      1200
```

Figure 3 : Logistic Regression performance

## Tool :

- Numpy and Pandas for data manipulation

- Scikit-learn for modeling, confusion matrix and feature extraction

- Matplotlib and Seaborn for plotting