# Report

**OUR AGENDA FOR THE DAY**

5-Step Ordering Process

Introduction

Goals of
this project

Technical
Approach

Lets talk
business

Future work

1

2

3

4

5

Let me introduce myself :

About Me

**I'm Menna, aspiring AI engineer who is currently an AI intern at NTI. I was also an AI Intern at Samsung.**

Fresh Graduate From Faculty of Computer science Cairo university

## The goals of This project:

To help the the management team to :
- See the value from using your platform for advanced analytics and the data platform challenges that accompany it.
- See if using advanced analytics may increase sales in general.
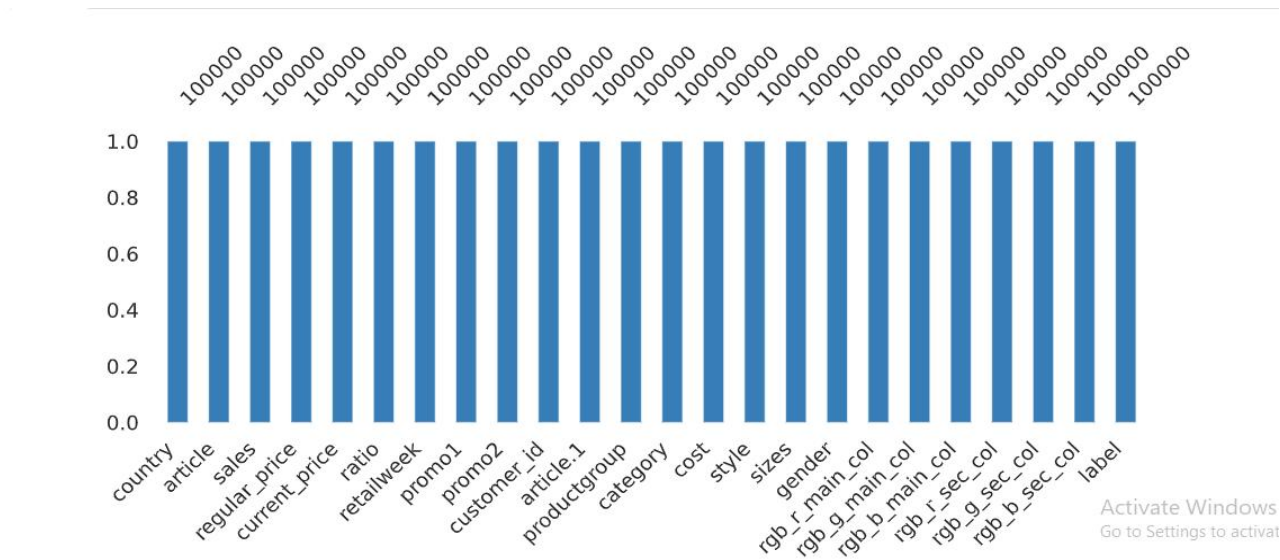
To help the the marketing team to :
- see if using advanced analytics you can help them increase their efficiency.
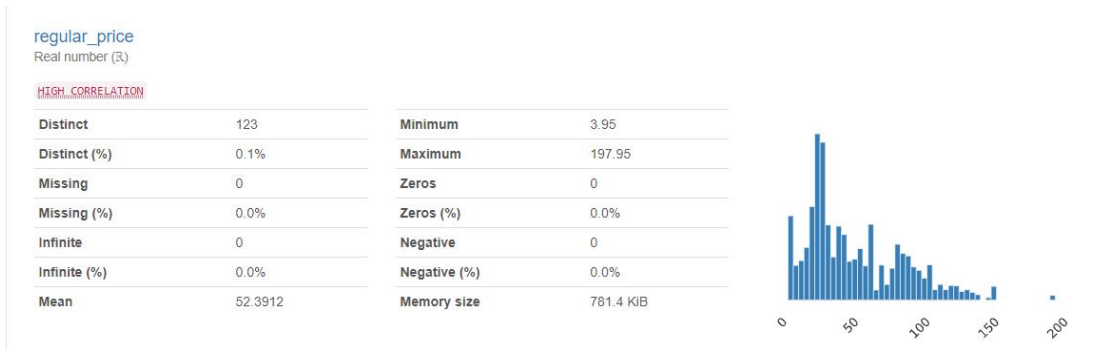
## Technical Approach :

- Lets explore the data first:

| Dataset statistics | | Variable types | |
|---|---|---|---|
| Number of variables | 24 | Categorical | 13 |
| Number of observations | 100000 | Text | 1 |
| Missing cells | 0 | Numeric | 9 |
| Missing cells (%) | 0.0% | DateTime | 1 |
| Duplicate rows | 0 | | |
| Duplicate rows (%) | 0.0% | | |
| Total size in memory | 18.3 MiB | | |
| Average record size in memory | 192.0 B | | |

**The data has No Missing Values :**

**The Most Important Features out of 24 :**

**-Regular_Price**

regular_price
Real number (ℝ)

HIGH_CORRELATION

| | | | |
|---|---|---|---|
| Distinct | 123 | Minimum | 3.95 |
| Distinct (%) | 0.1% | Maximum | 197.95 |
| Missing | 0 | Zeros | 0 |
| Missing (%) | 0.0% | Zeros (%) | 0.0% |
| Infinite | 0 | Negative | 0 |
| Infinite (%) | 0.0% | Negative (%) | 0.0% |
| Mean | 52.3912 | Memory size | 781.4 KiB |

**- Current_Price**

current_price
Real number (ℝ)

HIGH_CORRELATION

| | | | |
|---|---|---|---|
| Distinct | 141 | Minimum | 1.95 |
| Distinct (%) | 0.1% | Maximum | 195.95 |
| Missing | 0 | Zeros | 0 |
| Missing (%) | 0.0% | Zeros (%) | 0.0% |
| Infinite | 0 | Negative | 0 |
| Infinite (%) | 0.0% | Negative (%) | 0.0% |
| Mean | 28.2908 | Memory size | 781.4 KiB |

**-Retailweek**

retailweek
Date

| | | | |
|---|---|---|---|
| Distinct | 123 | Minimum | 2014-12-28 00:00:00 |
| Distinct (%) | 0.1% | Maximum | 2017-04-30 00:00:00 |
| Missing | 0 | | |
| Missing (%) | 0.0% | | |
| Memory size | 781.4 KiB | | |

## -Sales

sales
Real number (ℝ)

| | | | |
|---|---|---|---|
| Distinct | 476 | Minimum | 1 |
| Distinct (%) | 0.5% | Maximum | 898 |
| Missing | 0 | Zeros | 0 |
| Missing (%) | 0.0% | Zeros (%) | 0.0% |
| Infinite | 0 | Negative | 0 |
| Infinite (%) | 0.0% | Negative (%) | 0.0% |
| Mean | 56.7818 | Memory size | 781.4 KiB |

## - Country

country
Categorical

HIGH CORRELATION

| | |
|---|---|
| Distinct | 3 |
| Distinct (%) | < 0.1% |
| Missing | 0 |
| Missing (%) | 0.0% |
| Memory size | 781.4 KiB |

Germany 49400
Austria 35140
France 15460

# Lets dive into Our Approach:

- **Data visualization**
- **Analytical insights with EDA**
- **Preprocessing**
- **Modeling**

# Data visualization :

**-How many people Buy our Products ?**
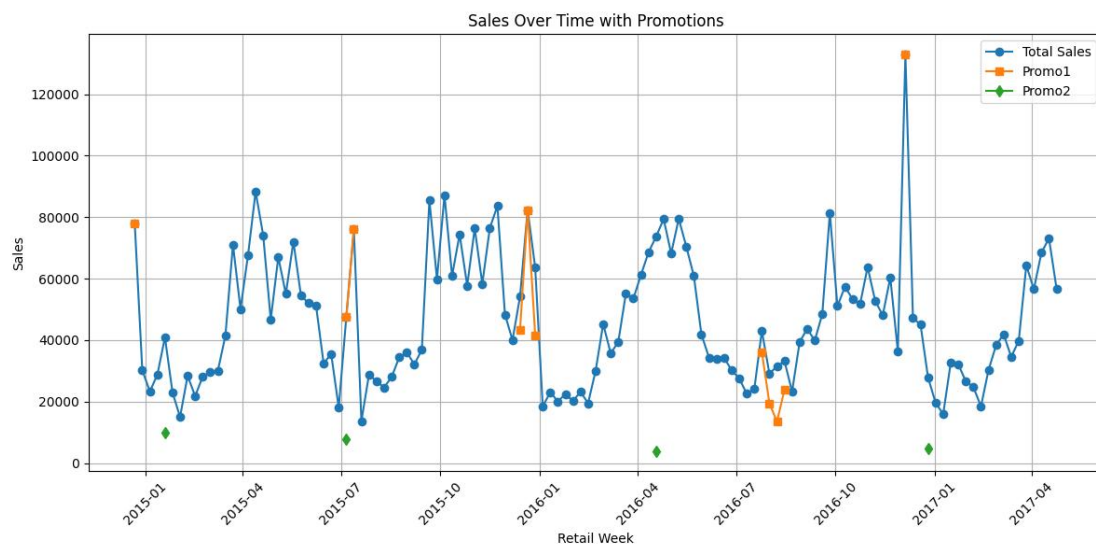
advertisement with customer buying And without



**From a Technical perspective The label is imbalanced as it has**

- 13.9% of people dont do advertising
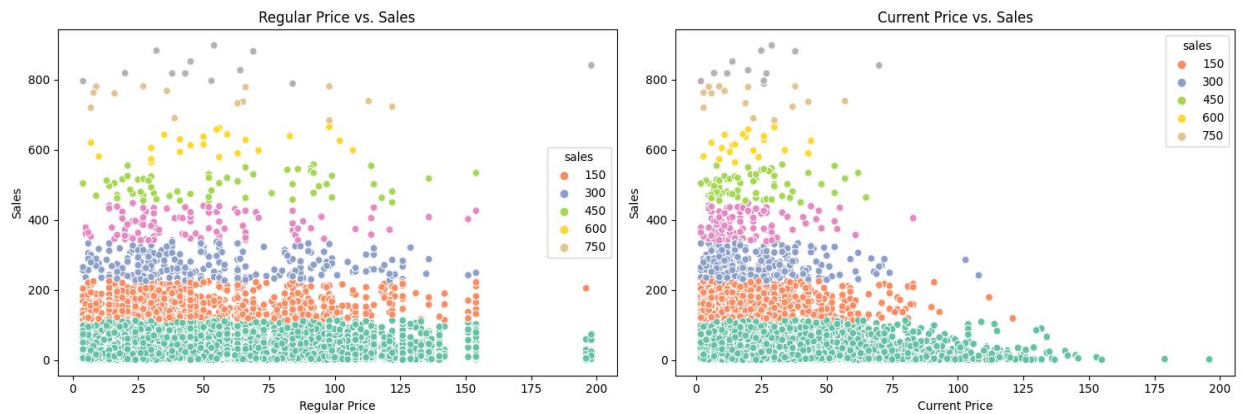- 86.1 do advertising

**From Sales And marketing perspective its good that the highest percentage of people do advertising**

**-What do you think the Impact of sales Over Time with Promotions ?**



**- Sales increase during promotional periods so it's a positive indicator of the effectiveness of our promotional strategies**
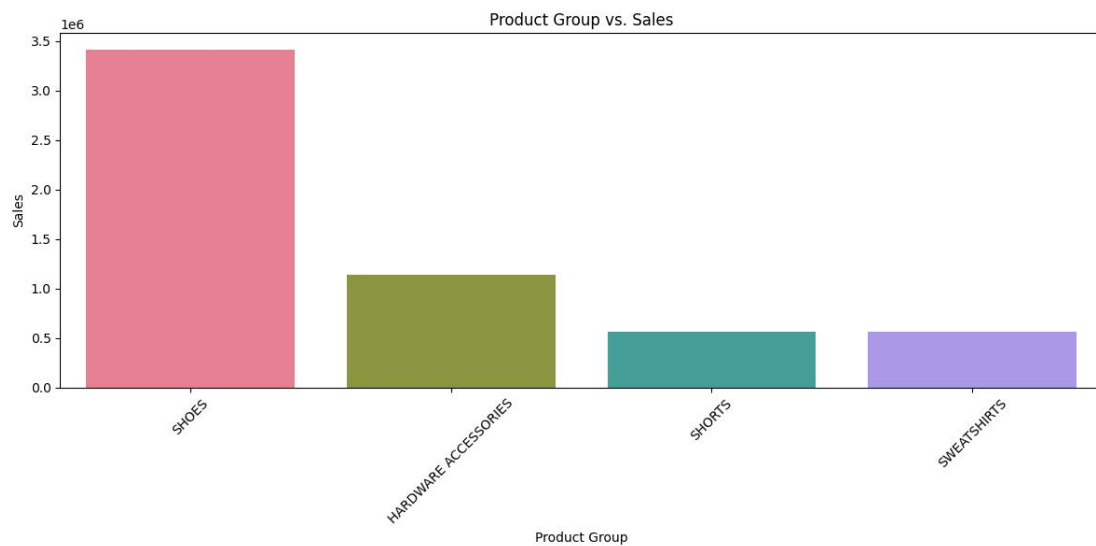
## How changes in regular and current prices affect sales ?



**They have the same effect on sales and it increase with both of them and that means :**
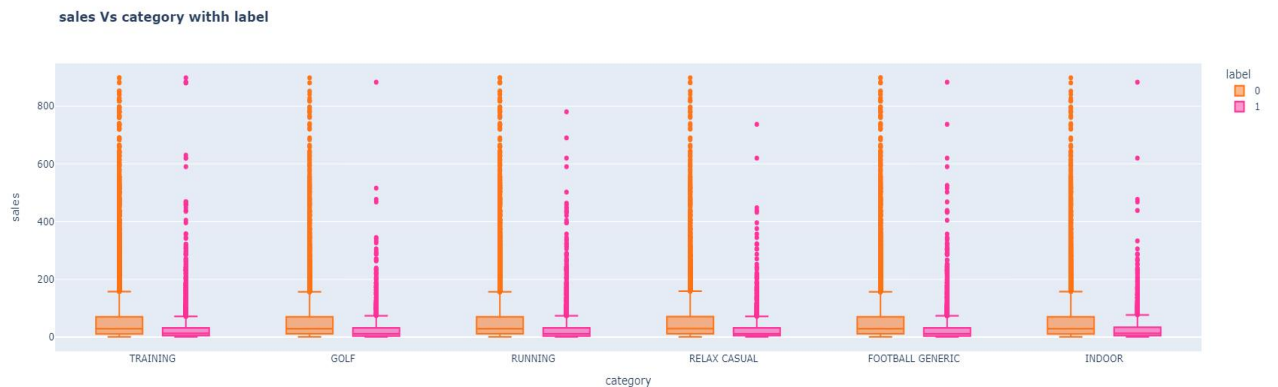
- there is a strong demand for the product
- Customers might be willing to pay more for the product due to its popularity and unique features.

## The Most Popular Product :



- **As we see the most popular Product that customers are interested in purchasing is Shoes**
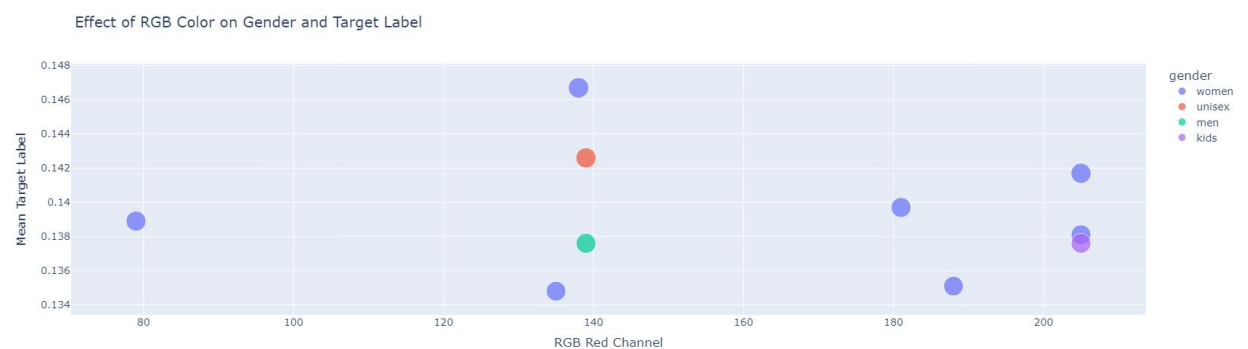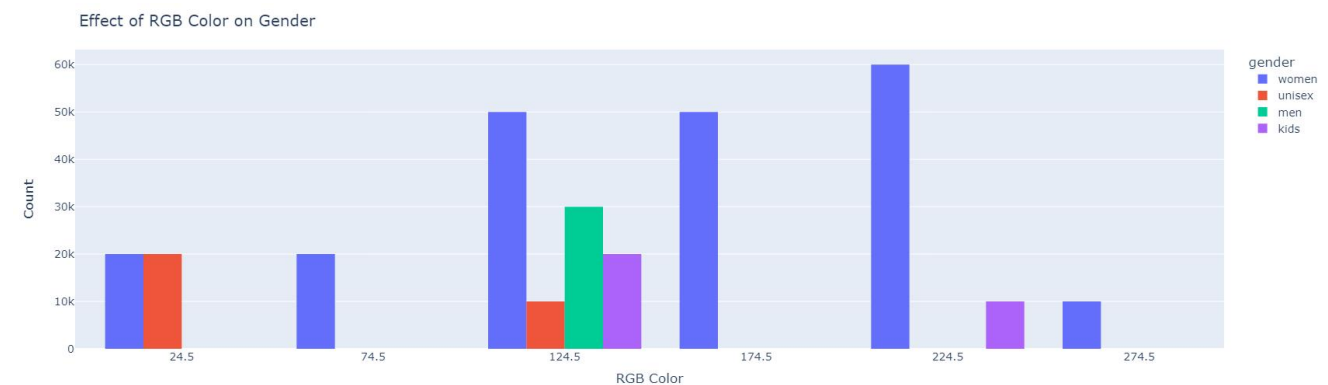
## How sales vary across different product categories ?

**sales Vs category withh label**
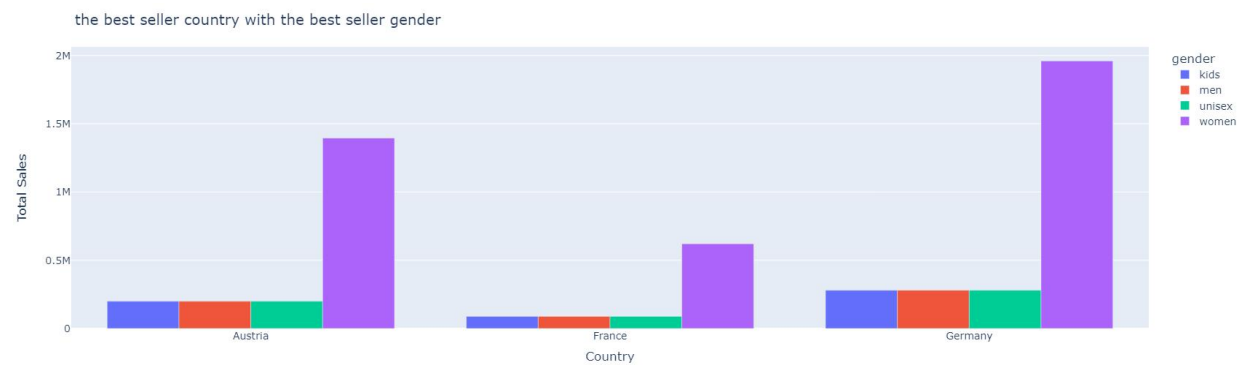


**Sales Don't Vary Much that mean :**

- Customers are purchasing products from various categories, indicating a balanced demand but the customers still dont buy that much as we said the data is imbalanced in the advertising process

## Which Gender interested to buy the product depending on colors ?



Effect of RGB Color on Gender



Effect of RGB Color on Gender and Target Label

**Women has the most interest in Colors and the product has more popularity with women**

## Which country is the best Buyer country with the best Buyer gender:



the best seller country with the best seller gender

**The best buyer country is Germany and women is the best buyer Gender**

# Preprocessing :

**- Discovering Outliers :**

```
country           2.000000
article         241.000000
sales            54.000000
regular_price    54.000000
current_price    26.000000
ratio             0.344409
retailweek       62.000000
promo1            0.000000
promo2            0.000000
customer_id    3553.250000
article.1         5.000000
productgroup      0.000000
category          4.000000
cost              7.310000
style             2.000000
sizes             0.000000
gender            1.000000
rgb_r_main_col   67.000000
rgb_g_main_col   77.000000
rgb_b_main_col  148.000000
rgb_r_sec_col    91.000000
rgb_g_sec_col    56.000000
rgb_b_sec_col   100.000000
label             0.000000
dtype: float64
```

**- Dropping 2 columns : ['article','customer_id']**

**-Discovering Important Features:**

```
        Weight      Feature
0.1273 ± 0.0001    ratio
0.0956 ± 0.0010    sales
0.0774 ± 0.0006    current_price
0.0768 ± 0.0014    retailweek
0.0715 ± 0.0009    regular_price
0.0531 ± 0.0012    country
0.0144 ± 0.0002    promo1
0.0065 ± 0.0003    cost
0.0048 ± 0.0002    rgb_g_main_col
0.0047 ± 0.0003    article.1
0.0044 ± 0.0002    rgb_b_main_col
0.0044 ± 0.0004    category
0.0030 ± 0.0002    rgb_r_main_col
0.0014 ± 0.0001    gender
0.0012 ± 0.0001    style
0.0010 ± 0.0001    rgb_g_sec_col
0.0009 ± 0.0001    rgb_r_sec_col
0.0009 ± 0.0002    productgroup
0.0008 ± 0.0001    promo2
0.0007 ± 0.0001    rgb_b_sec_col
```

## Modeling :
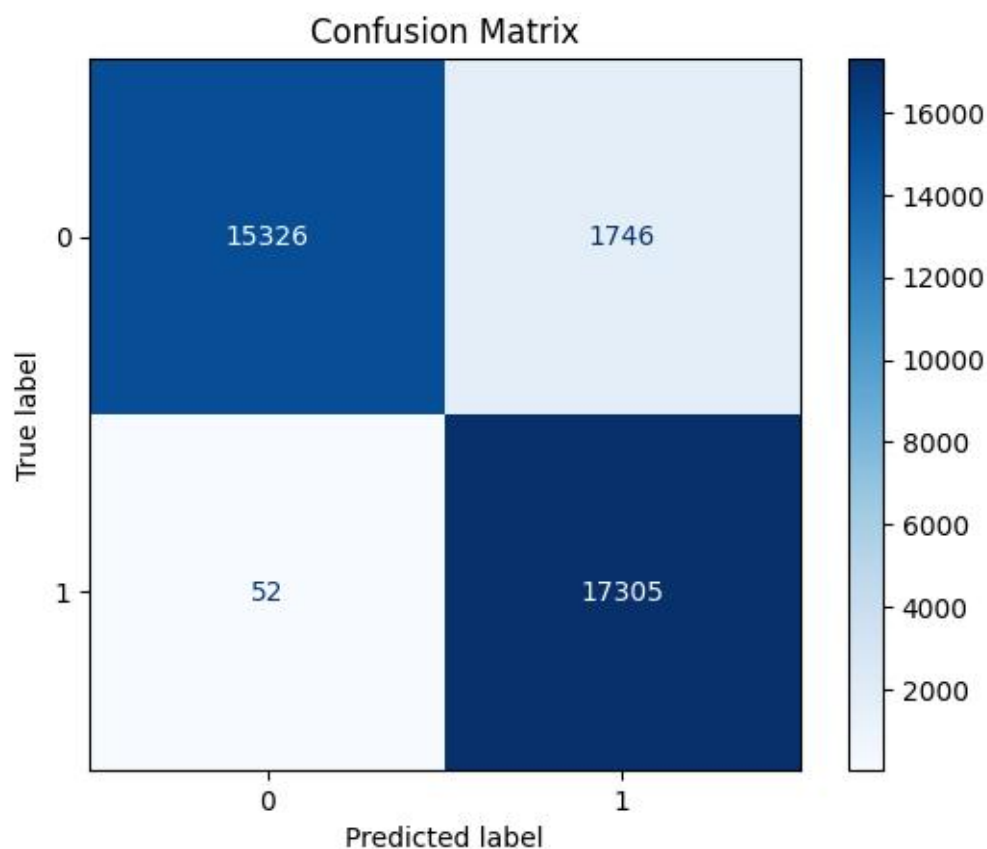I tried Logistic Regression,Decision Tree,KNN,XGBOOST
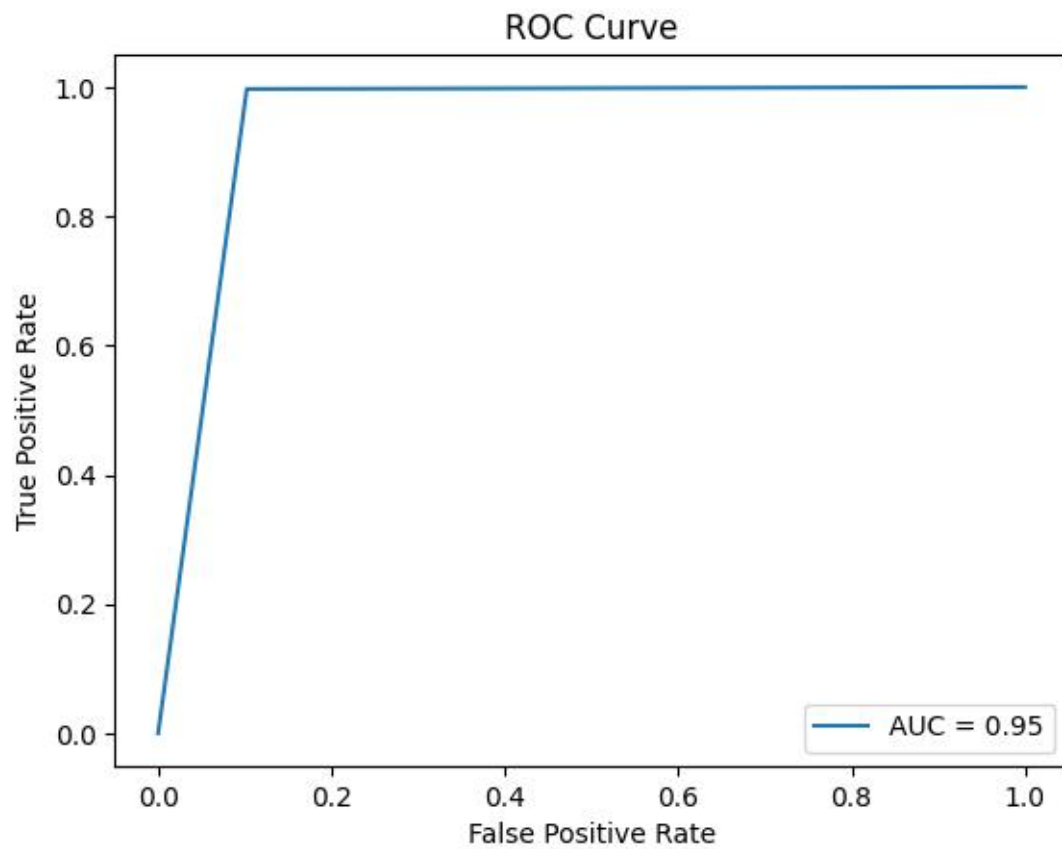Random Forest and xGBoost were the best Models so we would go with any of them
I chose these Models because they are suitable for Classification Problem like our Problem

### Random Forest :

```
Accuracy: 0.9477765836939789

              precision    recall  f1-score   support

           0       1.00      0.90      0.94     17072
           1       0.91      1.00      0.95     17357

    accuracy                           0.95     34429
   macro avg       0.95      0.95      0.95     34429
weighted avg       0.95      0.95      0.95     34429
```



Confusion Matrix

ROC Curve

**-XGBoost:**

```
0.9372912370385431
Accuracy: 0.9372912370385431

              precision    recall  f1-score   support

           0       1.00      0.88      0.93     17072
           1       0.89      1.00      0.94     17357

    accuracy                           0.94     34429
   macro avg       0.94      0.94      0.94     34429
weighted avg       0.94      0.94      0.94     34429
```
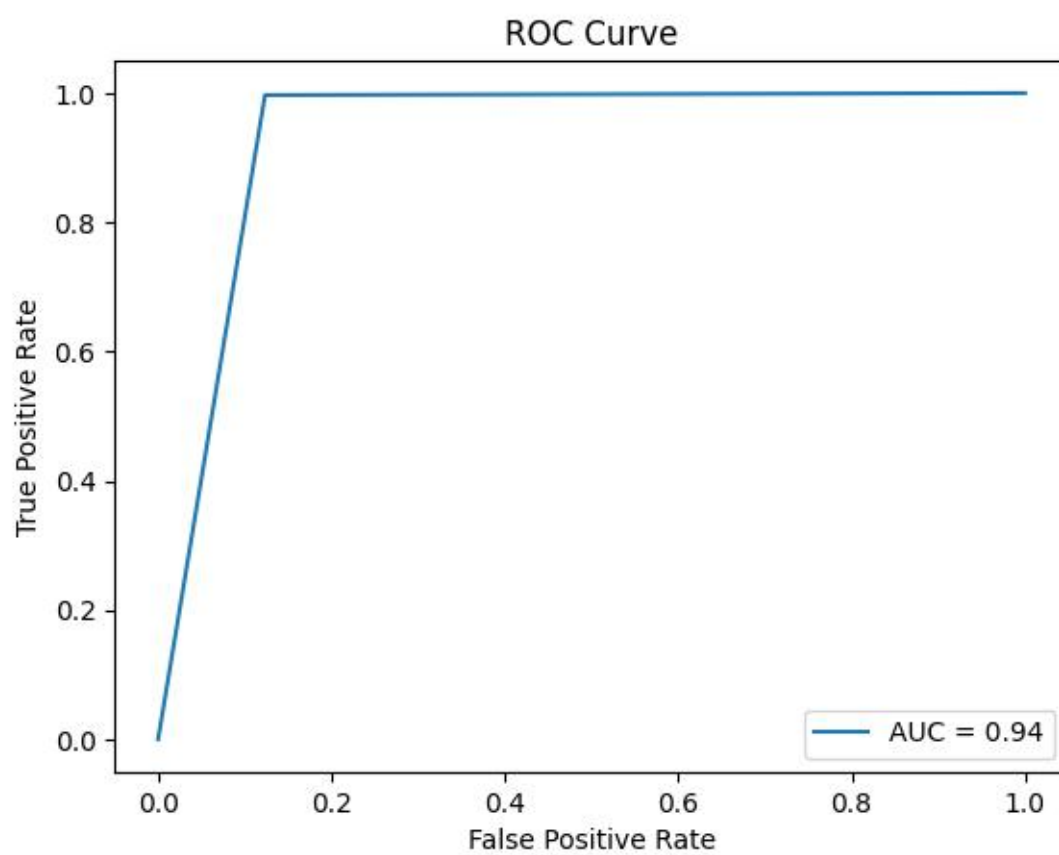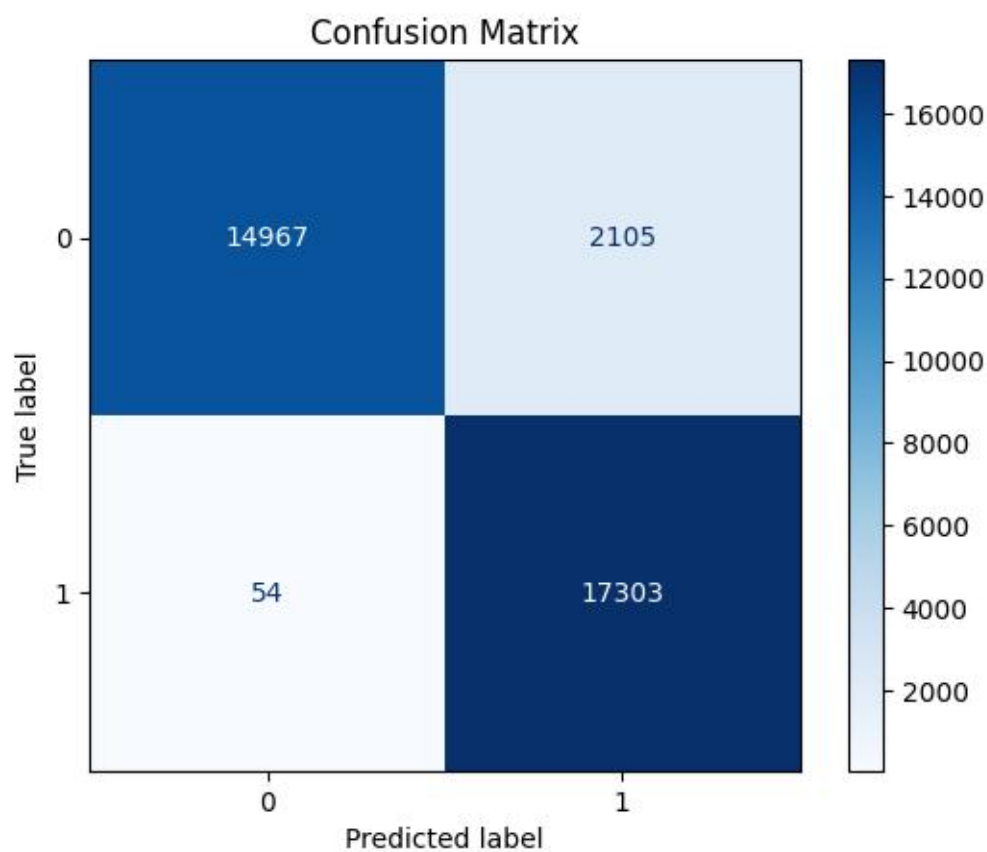
Confusion Matrix



ROC Curve

**-The accuracy of the different Models we Tried :**

| Models | Percentage |
|--------|-----------|
| knn | 0.92 |
| DecisionTree | 0.93 |
| RandomForest | 0.94 |
| LogisticRegression | 0.78 |
| xgboost | 0.93 |

**Future work :**
**We would like to work more on the products to be more good and popoular among different people**