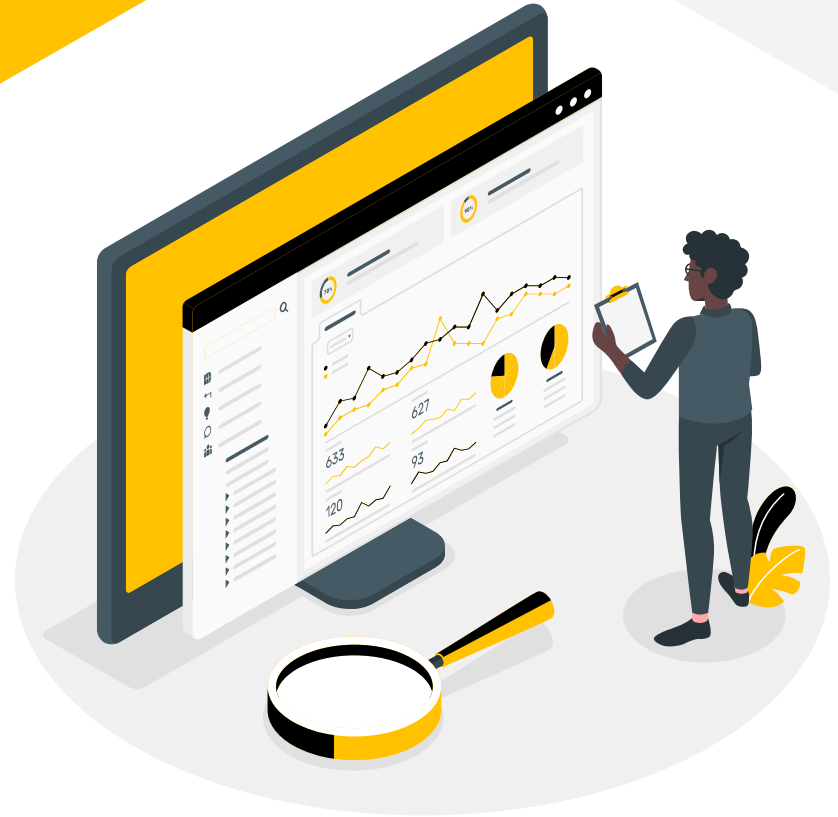# Online Shopper Intention

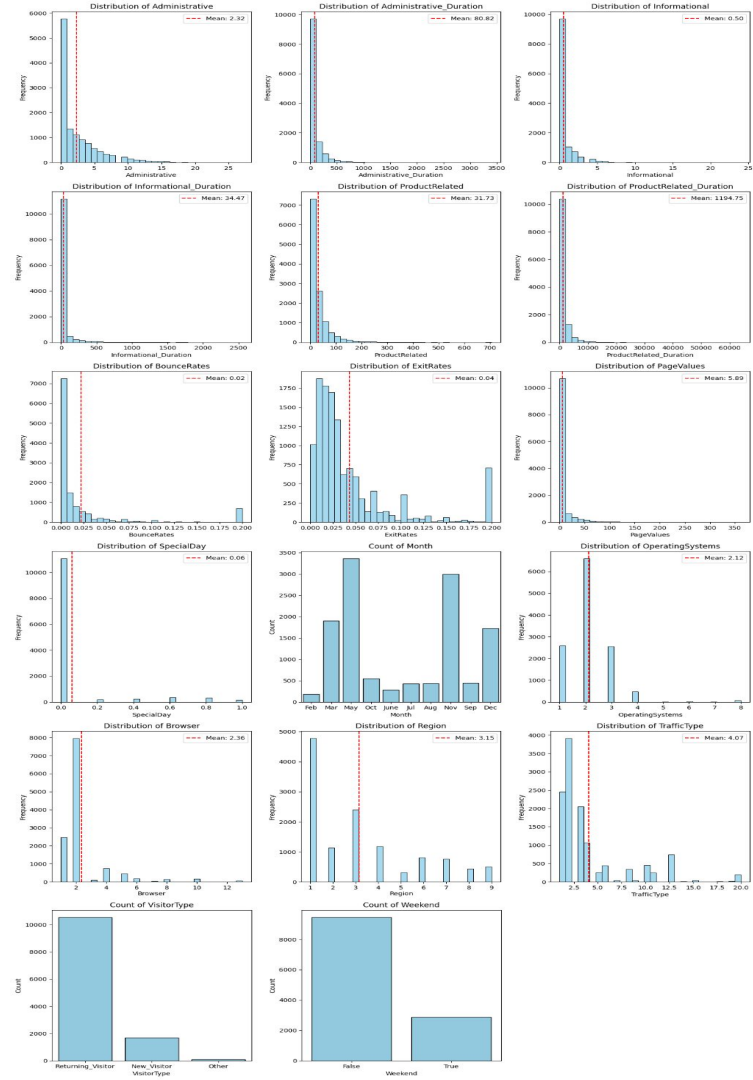By Thompson Pham

# Introduction & Background

01

- E-commerce - buying and selling of goods and services over the internet.
- Rapid growth of online shopping
- Understanding consumer behavior and predicting purchasing patterns importance

# Distribution

- Predominate left skew of most attributes
- Reflects the majority of classification of negative instance

```
Mean of continuous numerical features:
Administrative                 2.315166
Administrative_Duration       80.818611
Informational                  0.503569
Informational_Duration        34.472398
ProductRelated                31.731468
ProductRelated_Duration     1194.746220
BounceRates                    0.022191
ExitRates                      0.043073
PageValues                     5.889258
SpecialDay                     0.061427
OperatingSystems               2.124006
Browser                        2.357097
Region                         3.147364
TrafficType                    4.069586
dtype: float64

Variance of continuous numerical features
Administrative               1.103425e+01
Administrative_Duration      3.125085e+04
Informational                1.613297e+00
Informational_Duration       1.981036e+04
ProductRelated               1.978070e+03
ProductRelated_Duration      3.662130e+06
BounceRates                  2.351117e-03
ExitRates                    2.361624e-03
PageValues                   3.447868e+02
SpecialDay                   3.956808e-02
OperatingSystems             8.305129e-01
Browser                      2.949039e+00
Region                       5.767640e+00
TrafficType                  1.620199e+01
dtype: float64
```
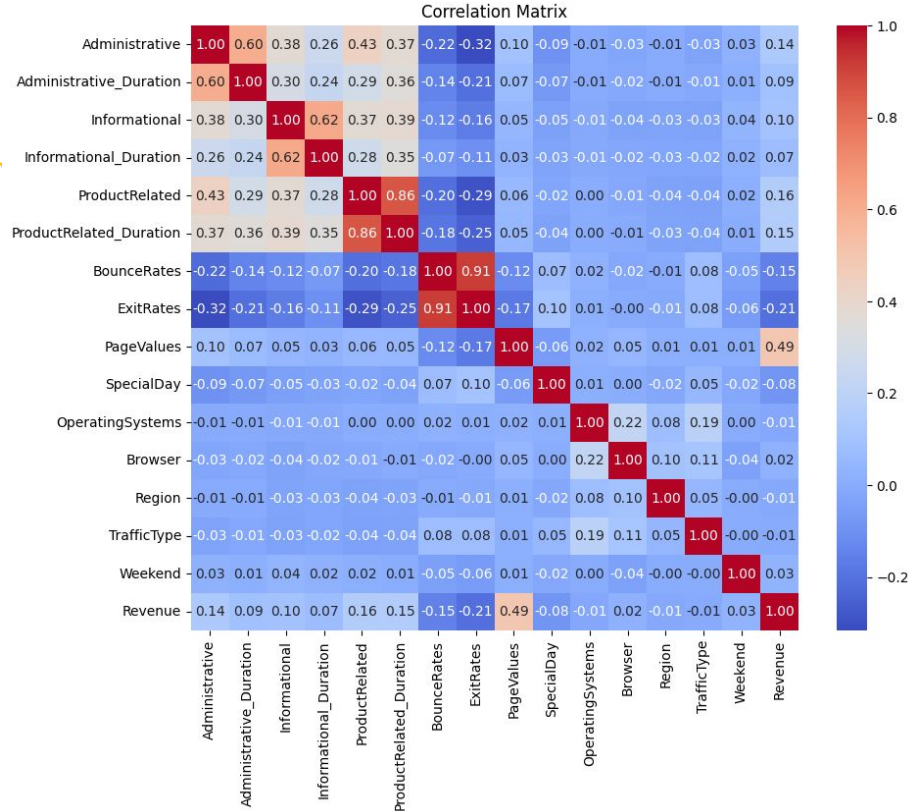
# Mean and Variance

- Mean collaborate with distribution, excluding ProductRealated_Duration
- High variance → users completing transaction
- Concern: Noise, overfitting, complexity

# Correlation



Correlation Matrix

## Variance

Interesting: Attributes which had high variance has high correlation

## Highlight

Good:
Pagevalues, ExitRates, and BounceRates

Bad:
Multicollinearity & less interpretable attributes

# Models



**Logistic Regression**
Accuracy: 87.31%

**Random Forest**
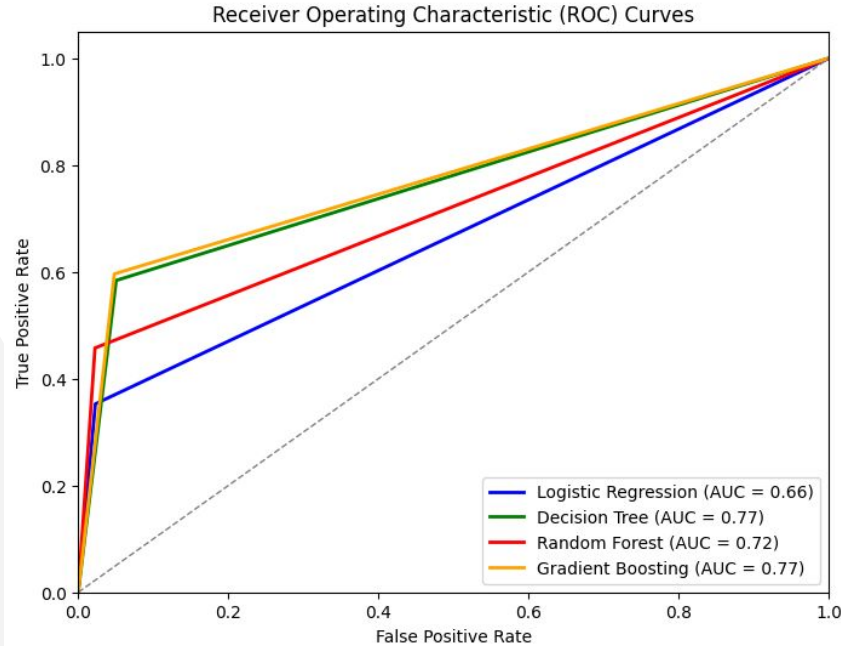Accuracy: 88.93%

**Decision Tree**
Accuracy: 88.85%

**Gradient Boost**
Accuracy: 89.13%

**Receiver Operating Characteristic (ROC) Curves**

Logistic Regression (AUC = 0.66)
Decision Tree (AUC = 0.77)
Random Forest (AUC = 0.72)
Gradient Boosting (AUC = 0.77)

**Logistic Regression**
Lowest AUC (0.66) More prone to misclassify class 1 as class 0
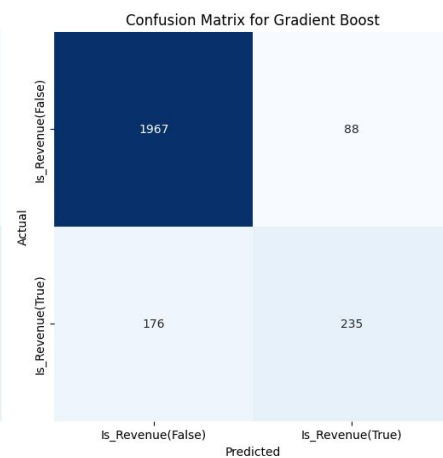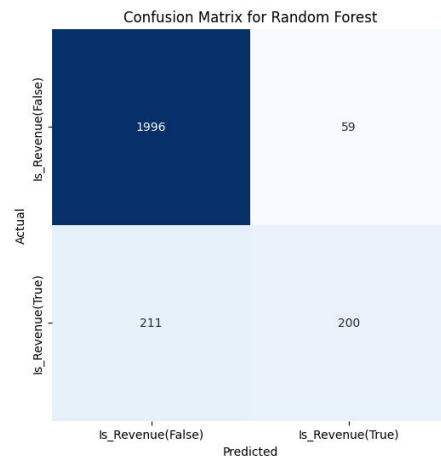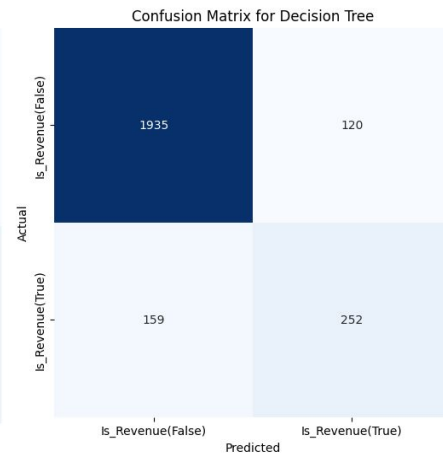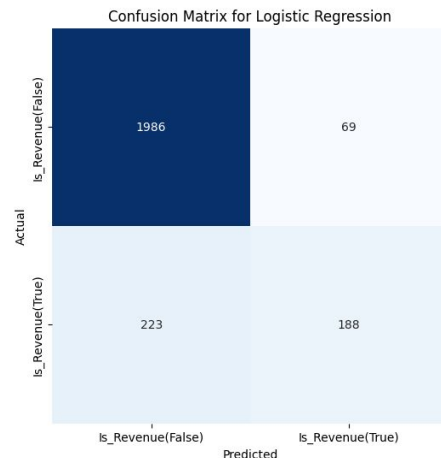
**DT & GB**
Highest AUC (0.77) Moderately well at seperating class 1 and class 2

# Confusion Matrix

- DT edges over RF(misclassification less pronounced)
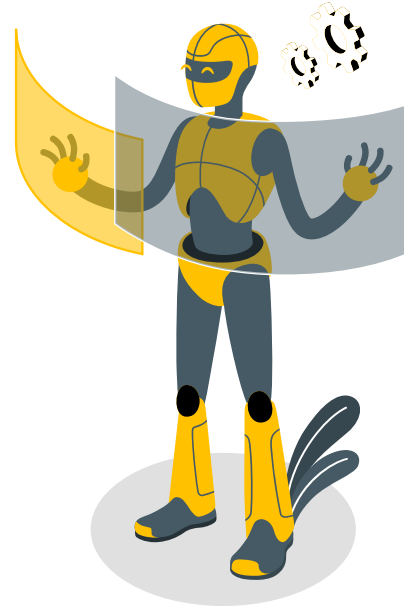- Number of misclassifications (top right & bottom left), reflect accuracy



**Confusion Matrix for Logistic Regression**

|  | Is_Revenue(False) | Is_Revenue(True) |
|---|---|---|
| Is_Revenue(False) | 1986 | 69 |
| Is_Revenue(True) | 223 | 188 |

**Confusion Matrix for Decision Tree**

|  | Is_Revenue(False) | Is_Revenue(True) |
|---|---|---|
| Is_Revenue(False) | 1935 | 120 |
| Is_Revenue(True) | 159 | 252 |

**Confusion Matrix for Random Forest**

|  | Is_Revenue(False) | Is_Revenue(True) |
|---|---|---|
| Is_Revenue(False) | 1996 | 59 |
| Is_Revenue(True) | 211 | 200 |

**Confusion Matrix for Gradient Boost**

|  | Is_Revenue(False) | Is_Revenue(True) |
|---|---|---|
| Is_Revenue(False) | 1967 | 88 |
| Is_Revenue(True) | 176 | 235 |

# Conclusion

**Best:**
Gradient Boosting - Balance of accuracy and predictive performance

**Alternative:**
Decision Tree - Marginally less accurate, but scalable

# Citations

Slide Template:
Slidesgo

**Sources**

- "Online Shoppers Purchasing Intention Dataset." *UCI Machine Learning Repository*, archive.ics.uci.edu/dataset/468/online+shoppers+purchasing+intention+dataset. Accessed 8 May 2024.