

# REAL-TIME CUSTOMER PURCHASING TENDENCY PREDICTION

SharkyData

Buyan Li, Yun Lin, Yu-Yuan Chang, Yafei Dong



[www.slidesgo.com](http://www.slidesgo.com)



# Table of contents

**01**

INTRODUCTION

**02**

EXPLORATION &  
VISUALIZATIONS

**03**

FEATURE  
ENGINEERING

**04**

MODEL BUILDING  
& EVALUATION

---

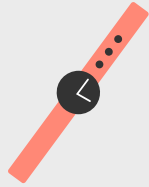


01

# Introduction



# Our Focus



Predicting the likelihood that customer will purchase during the shopping section in real-time



Analyzing the customer behaviors: page views, search history, and promo-code usage ect.



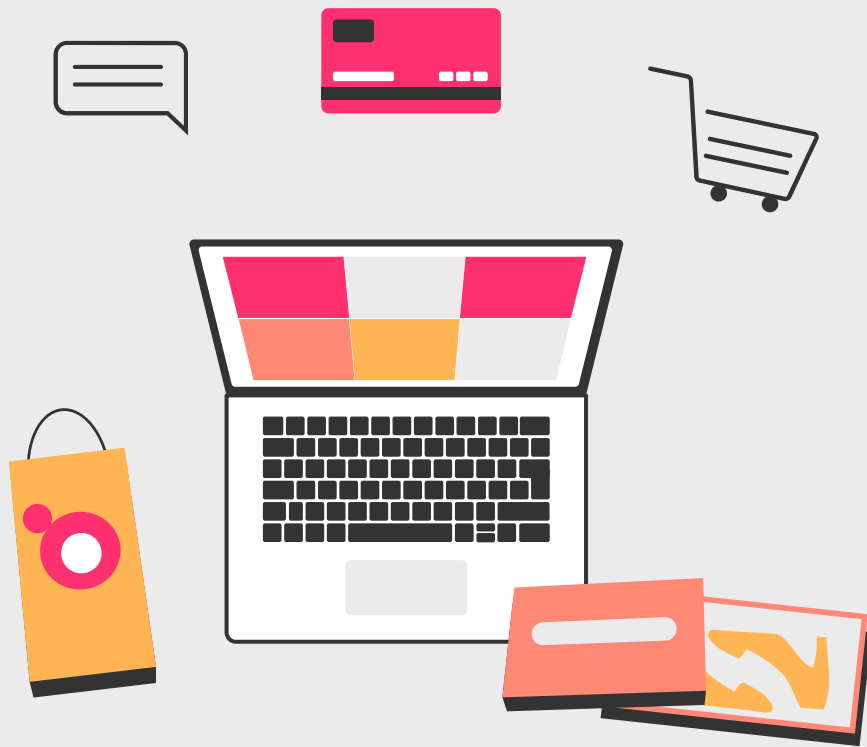
The implications for businesses seeking to optimize their e-commerce strategy

---

02

# EXPLORATION & VISUALIZATIONS

Yun Lin



# Initial Selection

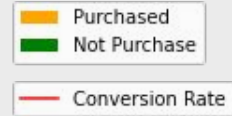
From the original dataset, we selected 17 features:

```
'visitid', 'pageurl', 'pagename', 'pageeventvar2',  
'pagetype', 'visitdatetime', 'hit_time_gmt',  
'productlist', 'searchterms', 'searchresults',  
'newvisit', 'post_evar27', 'evar28', 'evar83',  
'promocode', 'devicetype', and 'ordernumber'
```

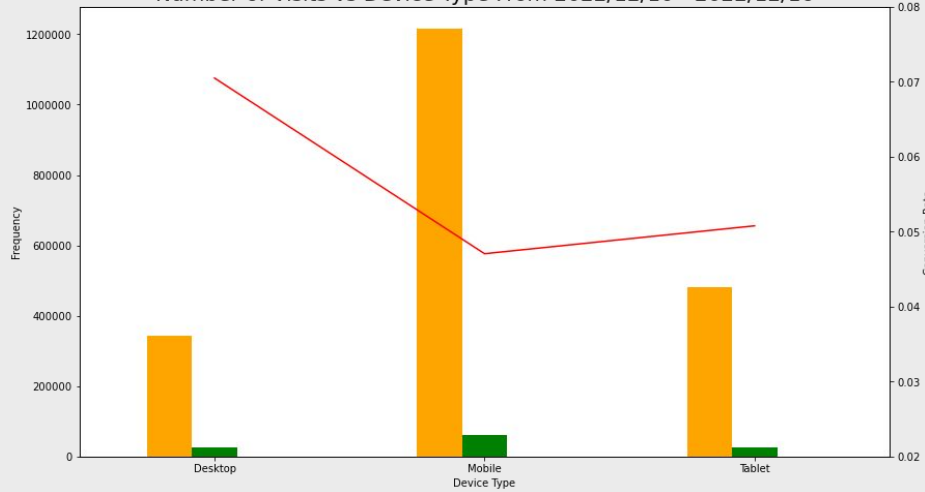




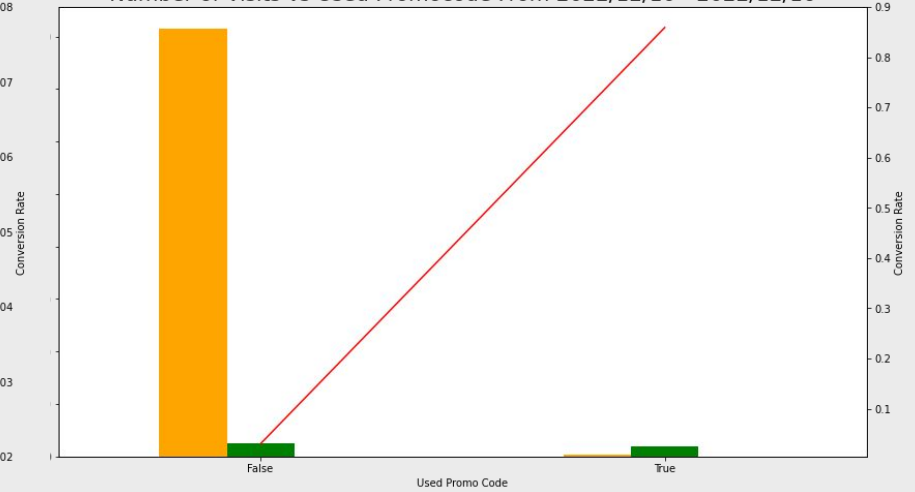
# Promo code & device Type groupby 'visitid'



Number of Visits vs Device Type From 2022/12/10 - 2022/12/16

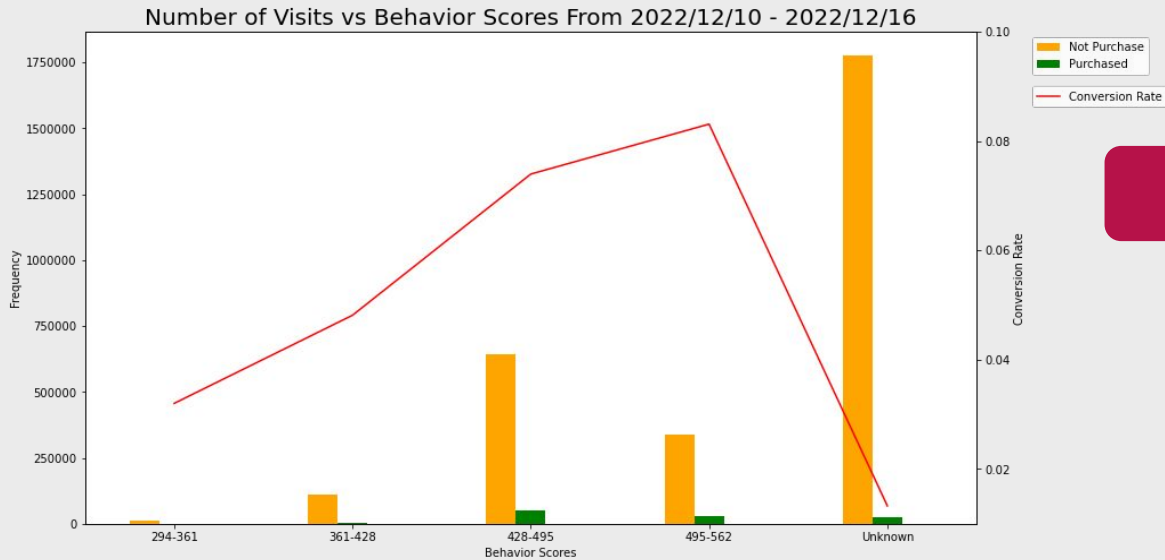


Number of Visits vs Used Promocode From 2022/12/10 - 2022/12/16





# Behavior Score (Evar83)



**Evar83**

**294-361**

**361-428**

**428-495**

**495-562**

**Unknown**



03

# FEATURE ENGINEERING





---

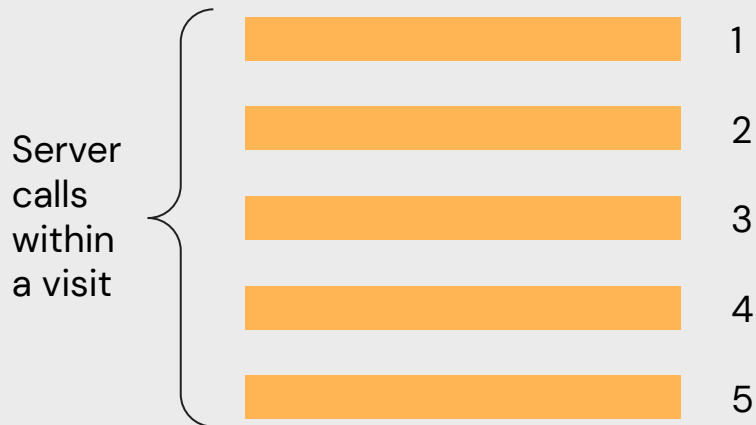
# Adding memory features

- What are “memory” Features and why do we need them?
  - Store information from previous server calls within the same visit
  - Needed for generating a prediction for each server call

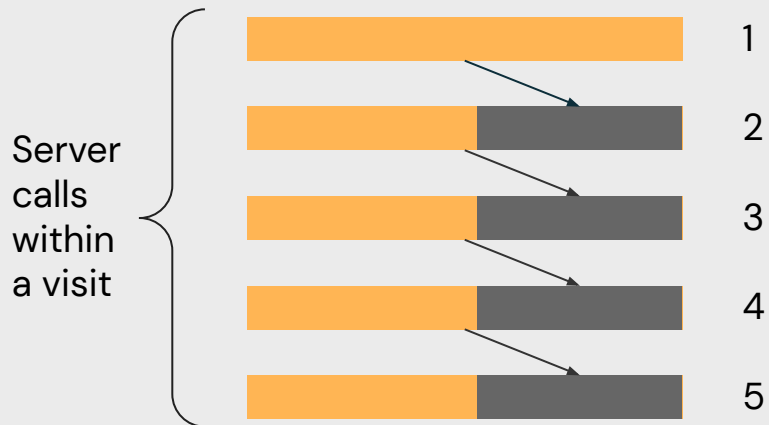


# Adding memory features

Before adding memory features



After adding memory features





---

# Adding memory features

Memory features added:

sum\_cart

sum\_prod

prod\_in\_cart

prod\_viewed

avg\_cart

avg\_prod

min\_cart

min\_prod

max\_cart

max\_prod

---



---

# Adding other features

Other features added:

`Timedelta`: Time between this server call and the previous one

`Pg_cat`: Category of the current page

`Session_time`: Total time spent since the start of this visit

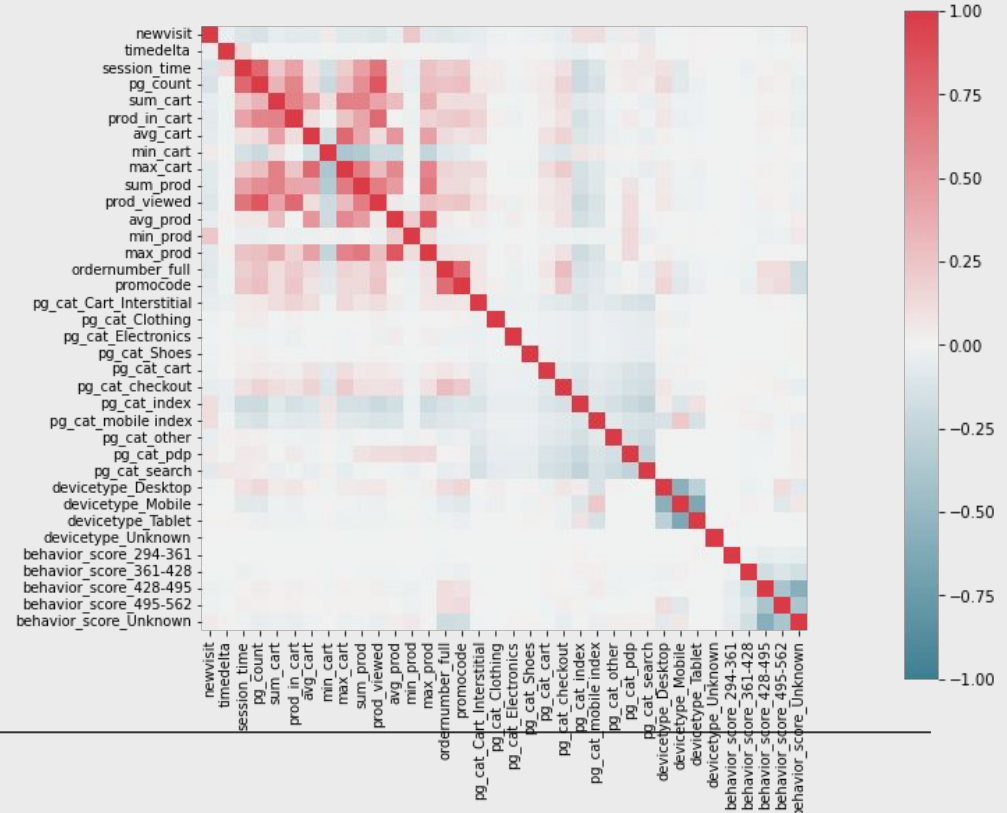
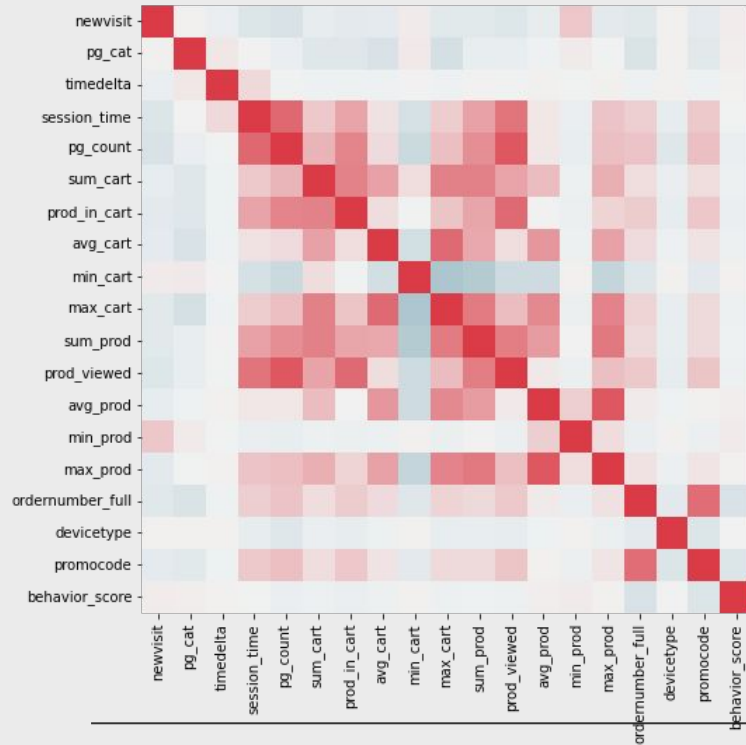
`Pg_count`: Total server calls generated since the start of this visit

---



# Correlation matrix

## before & after one-hot encoding



# 04

## MODEL BUILDING & EVALUATION



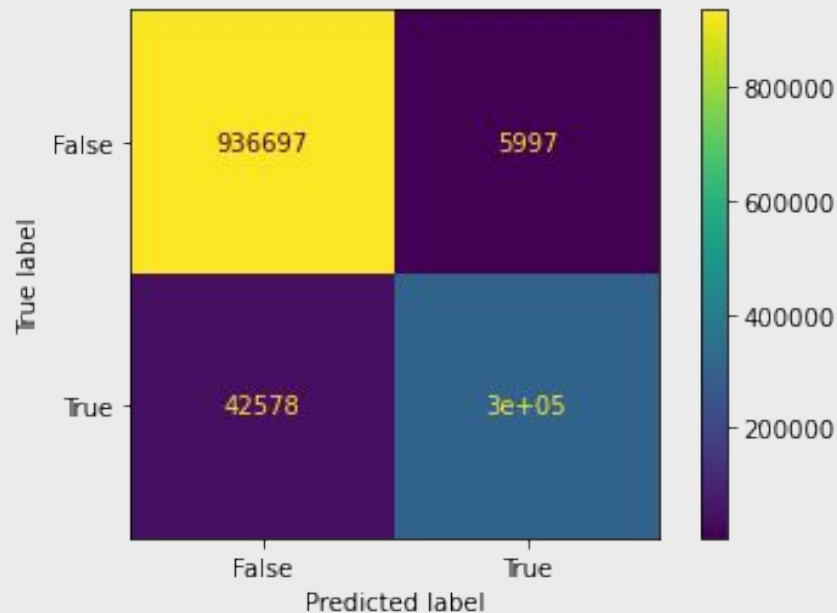
# Model Building

1. Random forest
  2. Gradient Boosting Classifier
  3. Logistic Regression
  4. SVM
  5. KNN
- 
-

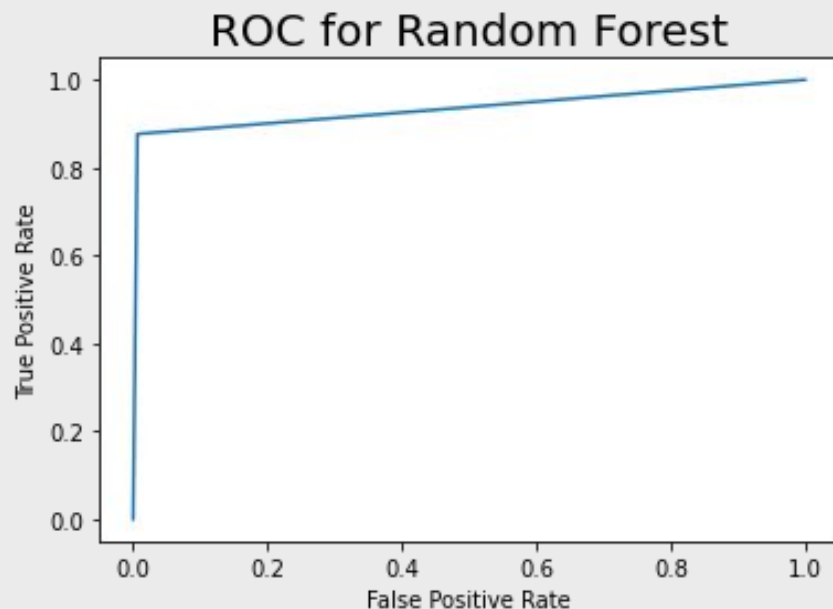


# Random Forest

Accuracy:	0.9624
Precision:	0.9806
Recall:	0.8765
F1-Score:	0.9256



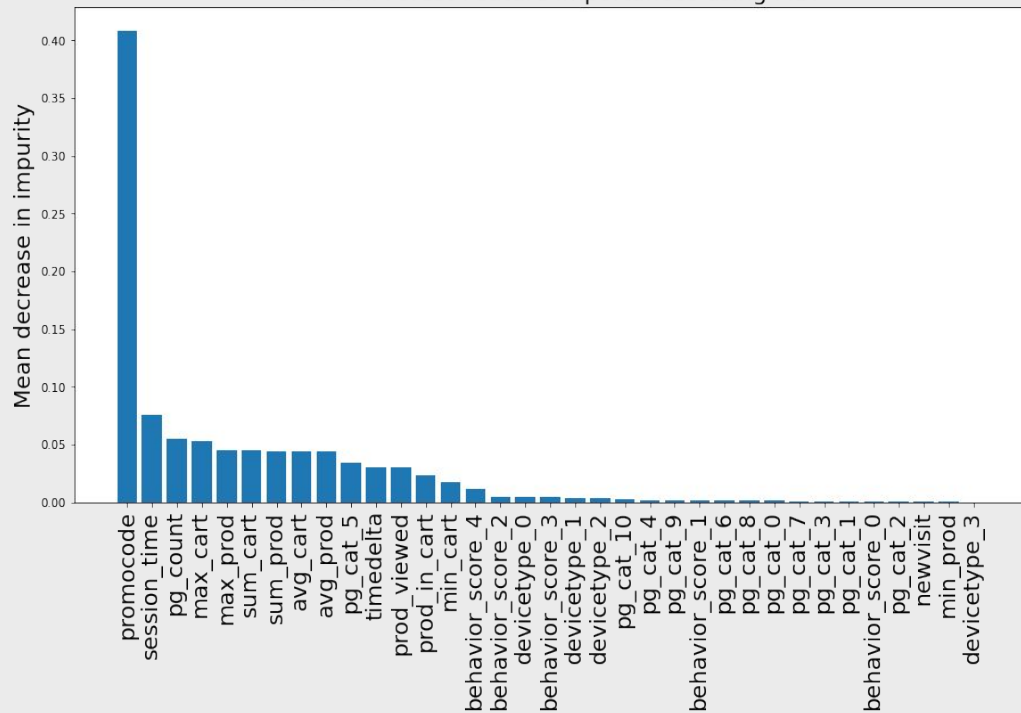
# Random Forest



1. **Parameter tuning:**  
criterion: "entropy" and "gini".  
max\_depth: [10, 20, 30]
  2. **Cross-Validation:**  
5 folds
  3. Performance the best!!
-

# Result: Feature Importances

Random Forest Feature importances using MDI

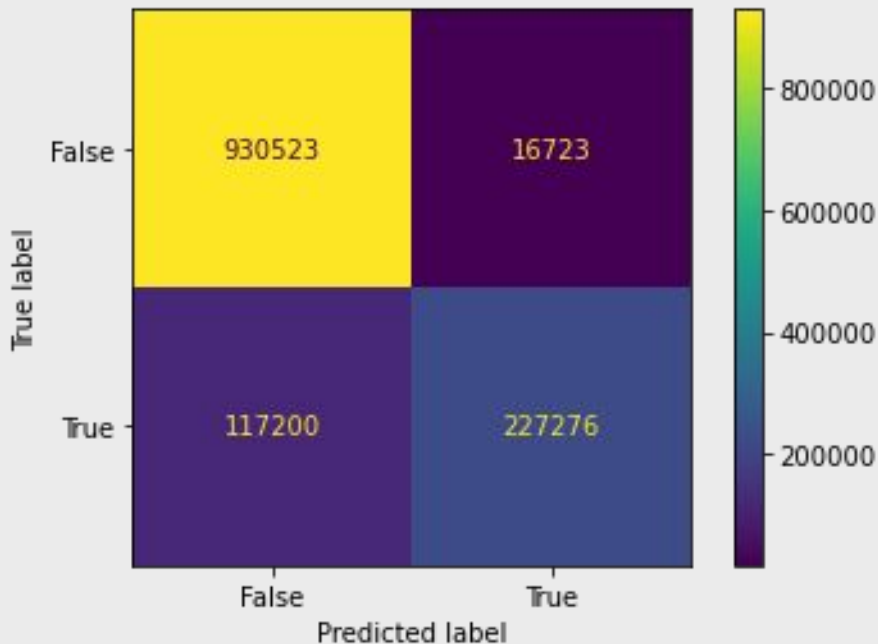


Feature Importance:

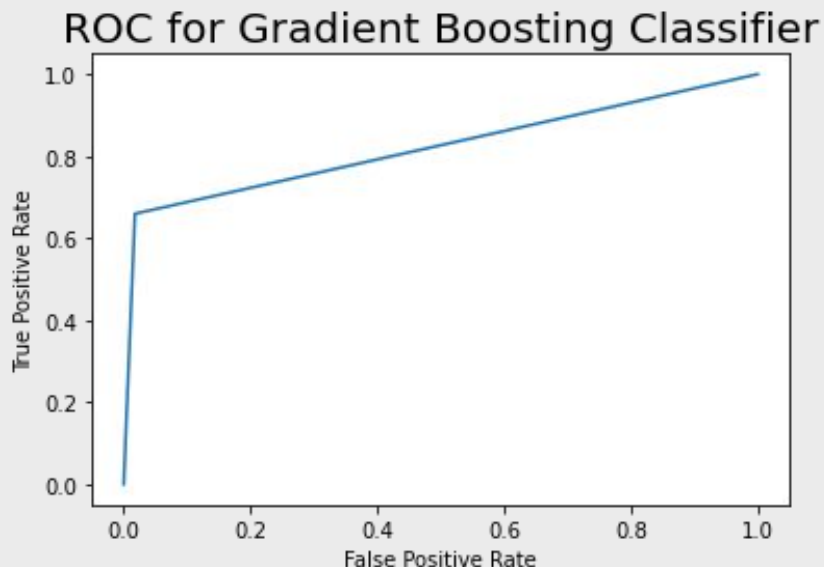
1. Promo code
2. Session time
3. Page count (Number of action)

# Gradient Boosting Classifier

Accuracy	0.8963
Precision	0.9314
Recall	0.6597
F1-Score	0.7724



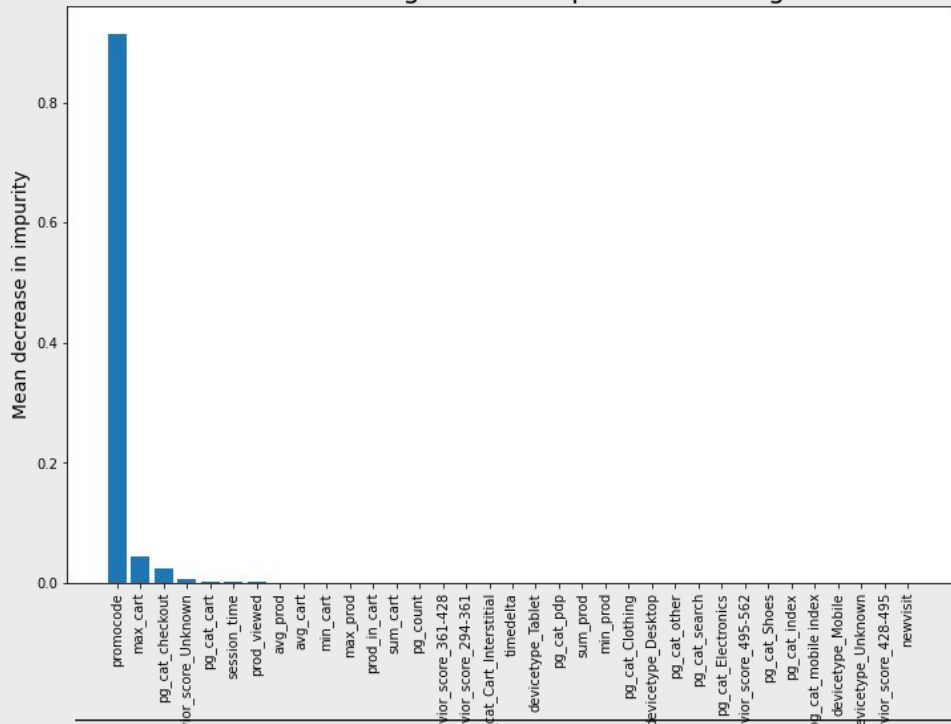
# Gradient Boosting Classifier



1. **Parameter tuning:**  
Learning rate: [0.01, 0.1, 0.5, 1]  
n\_estimator: [10, 50, 100, 150, 200]
2. **Cross-Validation:**  
5 folds

# GB: Feature Importances

Gradient Boosting Feature importances using MDI

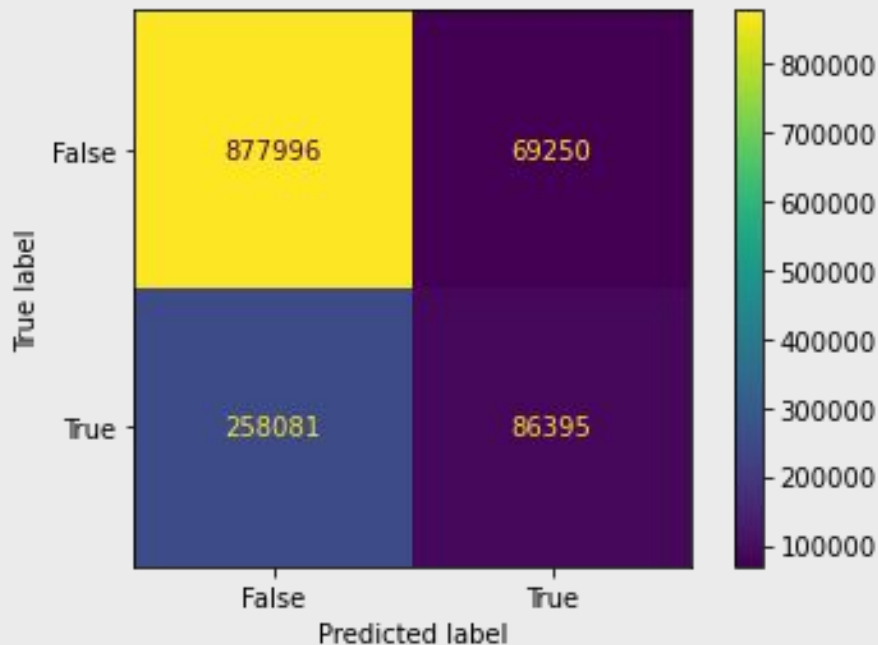


Feature Importance:

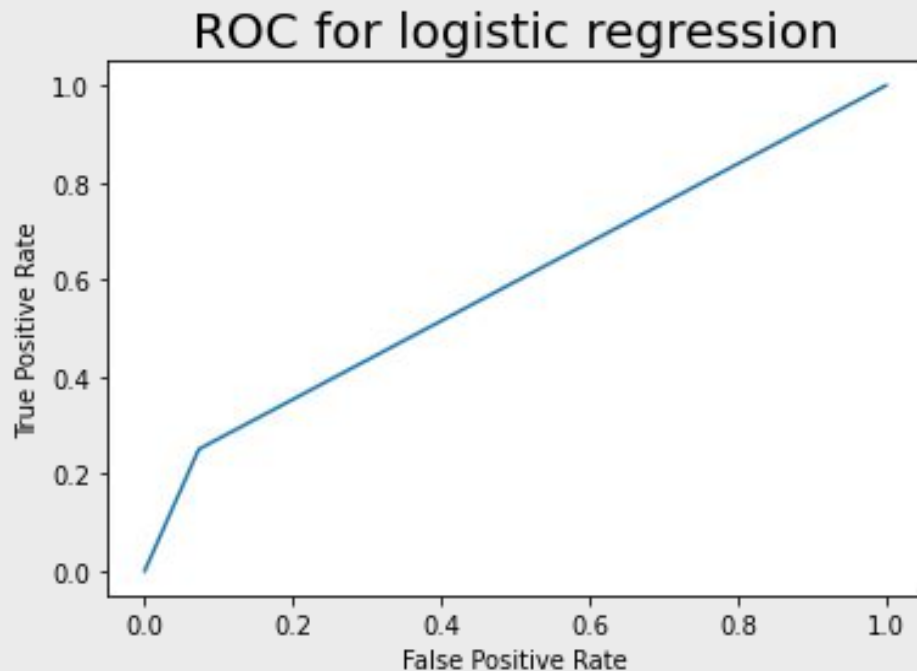
1. Promo code
2. max\_cart
3. pg\_cat\_checkout (Number of action)

# Logistic Regression

Accuracy	0.7465
Precision	0.5550
Recall	0.2508
F1-Score	0.3454



# Logistic Regression

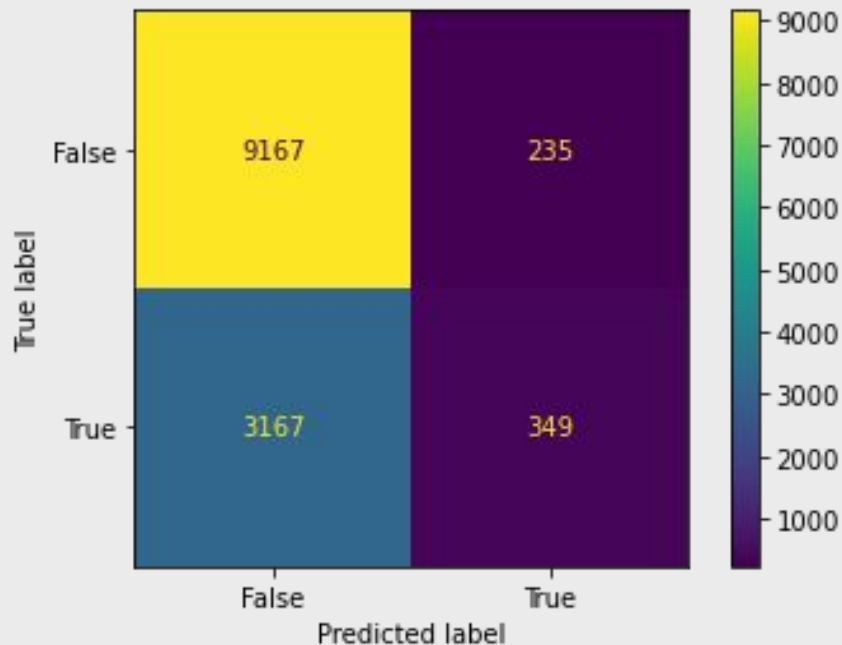


1. **Parameter tuning:**  
penality: 'l1', 'l2', 'elasticnet', None
2. **Cross-Validation:**  
5 folds
3. Did not performance well



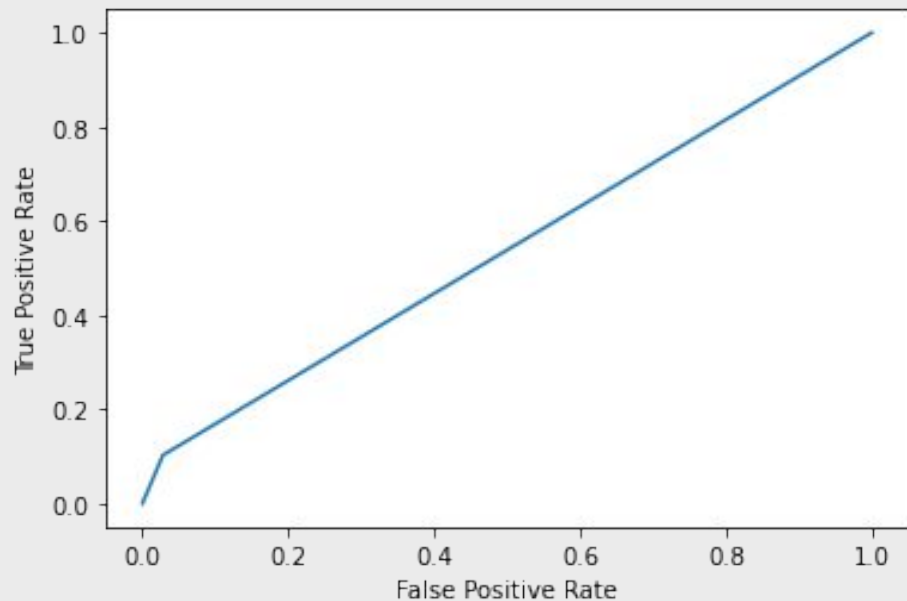
# SVM

Accuracy	0.7366
Precision	0.5976
Recall	0.0992
F1-Score	0.1702



# SVM

ROC for SVM



1. **Parameter tuning:**  
kernel: "linear", "poly" and "rbf"
2. **Cross-Validation:**  
5 folds
3. Only use 1% of the sample.
4. Did not performance well

# Summary

1. Analyzing the frequency of pageviews for a particular product.
2. Design new feature and construct model.
3. Being able to provide insights into customer preferences and forecast which products are likely to be purchased in the future.

# Future Work

1. Factors external: economic conditions, weather, promo code and social.
2. Recommendation system
3. We can limited our model into first few click/pages

# Limitation

1. Features Mapping were not consistent
  2. Not sure when 'Promocode' applied
  3. The data is during Holiday season
  4. The data contains only one week timeframe
-

**Thanks for listening**

---

# **Backup slides**



# NewVisit → NewVisitor

	visitid	newvisit
0	001	0
2	001	0
3	001	0
4	001	0
5	001	0
...	...	...
6458604	999838842448347828201007282656244896618	1
6458605	99997363407775132790292379433452812751	1
6458606	99997363407775132790292379433452812751	0
6458607	99997363407775132790292379433452812751	0

	visitid	newvisit
0	001	0
2	001	0
3	001	0
4	001	0
5	001	0
...	...	...
6458604	999838842448347828201007282656244896618	1
6458605	99997363407775132790292379433452812751	1
6458606	99997363407775132790292379433452812751	1
6458607	99997363407775132790292379433452812751	1



# Promocode → checkout?

So we assume promo code applied in cart

s.eVar82

Promo Code Errors in Checkout (eVar 82)

eVar82 is empty

2041286 rows x 4 columns

```
[184]: p_t = test[test['promocode']==1]
       p_t
```

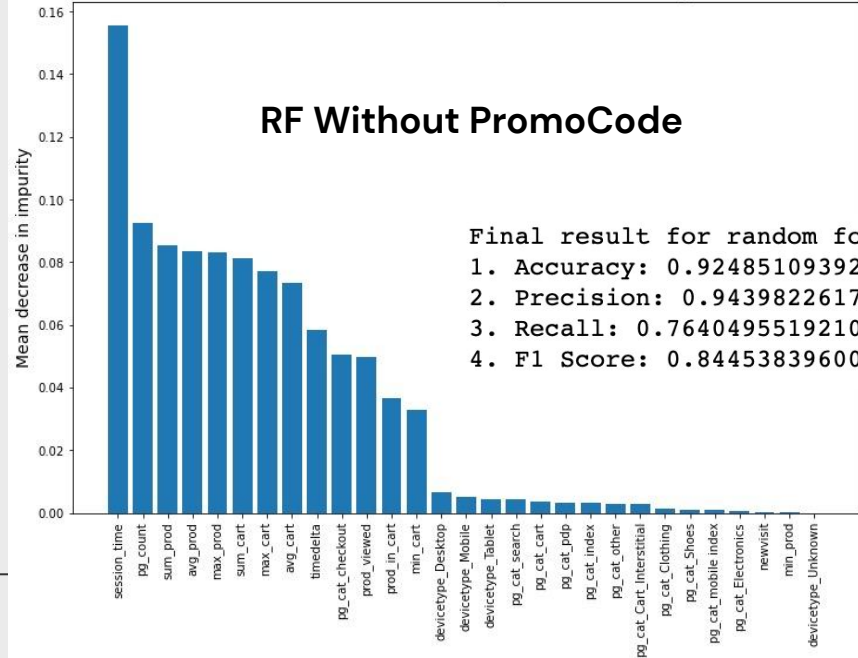
t[184]:

	promocode	eVar82	order
visitid			
100010825755995377356241128984439462981	True	False	True
100034118541581008486617653558704628597	True	False	True
100068898285134821927141456802434564309	True	False	True
10008266340795712901655509953153309891	True	False	True
10008266340795712901655509953153309892	True	False	True
...	...	...	...
99898932752452050033614126798257750096	True	False	True
99898932752452050033614126798257750097	True	False	True
99966830539311880448190947143702206652	True	False	True
99981909984864820644946664237114942271	True	False	True
99986567744921905273434973551979382621	True	False	False

55284 rows x 3 columns

Errors submitting a promo code in checkout

Random Forest Feature importances using MDI



# Back Up Slides

## Variable Importance for RF

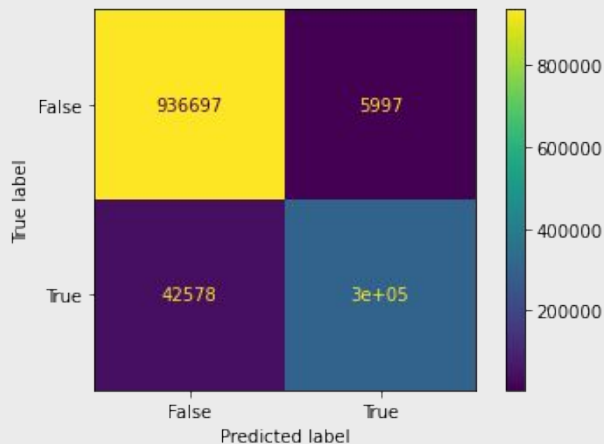
### 1. Mean Decrease in Impurity (MDI)

- Same as Bagging.
- Calculate the total amount that the MSE (for regression) or Gini index (for classification) is decreased due to splits over a given predictor in different nodes (weighted by the probability of reaching that node (which is approximated by the proportion of samples reaching that node)).
- The weighted decrease in purity as a result of the splits over a given predictor is averaged over all trees, and is used as a measure of the importance of variable  $j$  in the random forest.
- The default in Scikit-learn `feature_importances_`



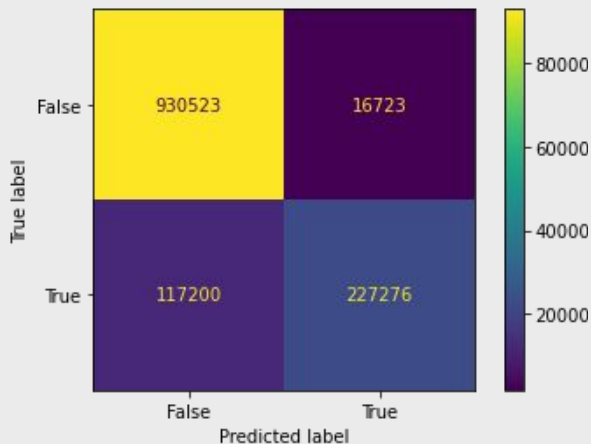
# Random Forest

F1-Score:	0.9256
-----------	--------



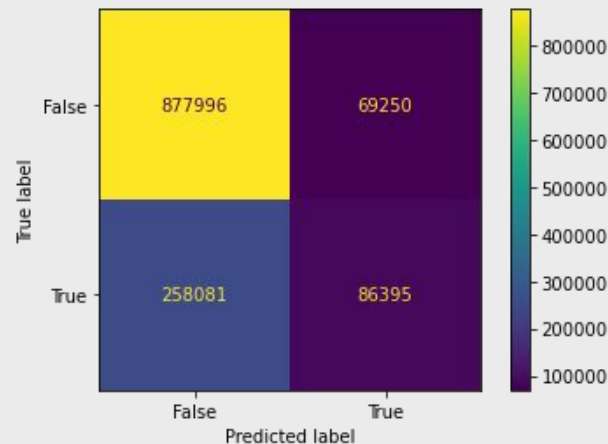
# Gradient Boosting

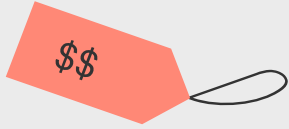
F1-Score:	0.7724
-----------	--------



# Logistic Regression

F1-Score:	0.3454
-----------	--------





# Resources



Slide template: [www.slidesgo.com](http://www.slidesgo.com)

