



# Online Shopping Products purchase prediction

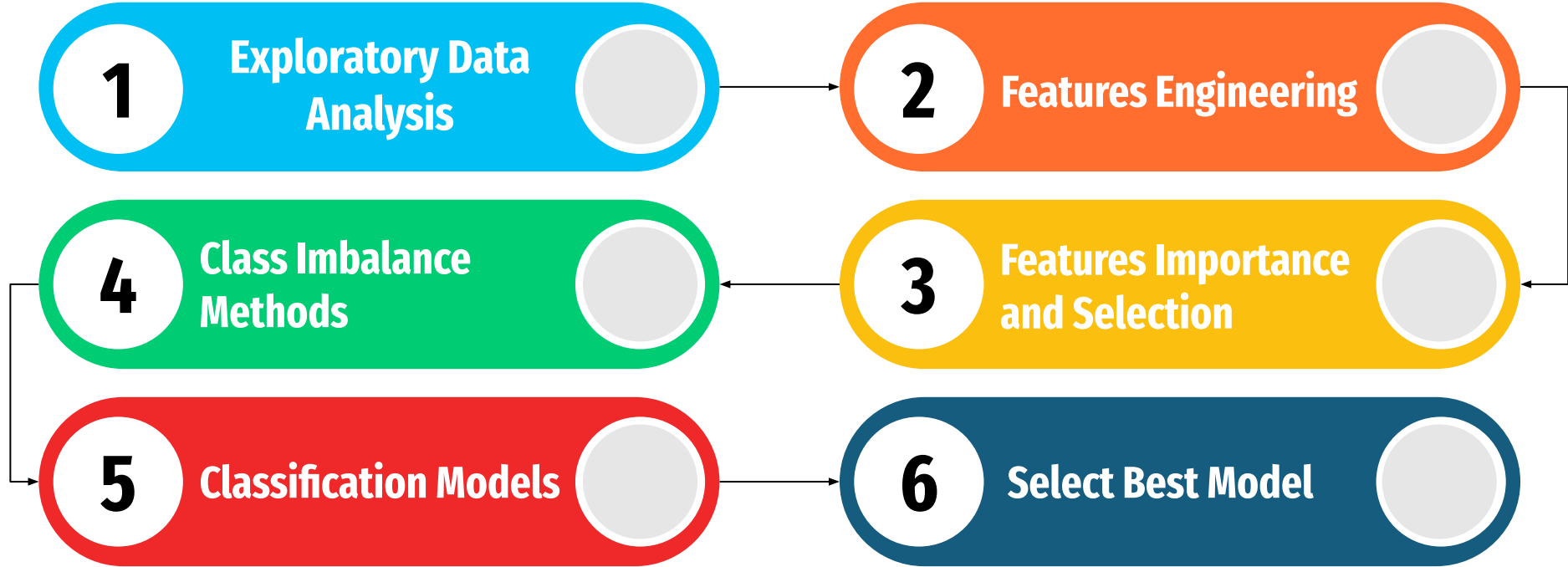
Raed Altuwaijri

# Problem statement

We are an online shopping store our goal is to increase sales and reduce logistics time and costs by recommending a products to the customer base on there prior orders



# Methodology



# Data Description

## Data sets

Data about each order and the user places that order

**Customer Orders**



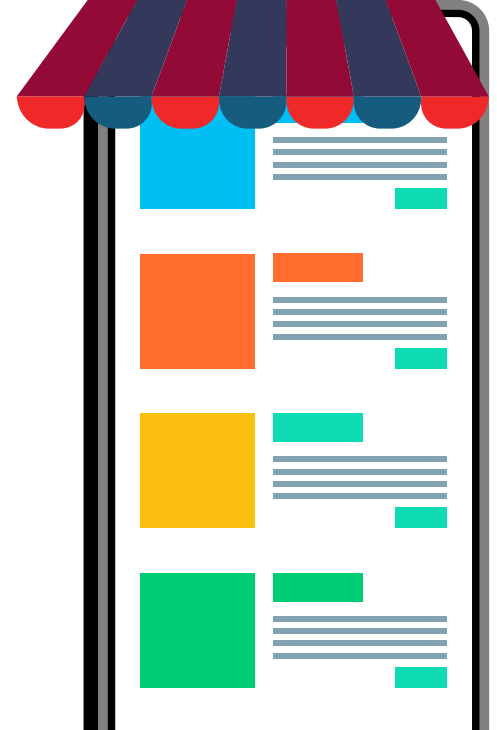
Data about each product

**Products**

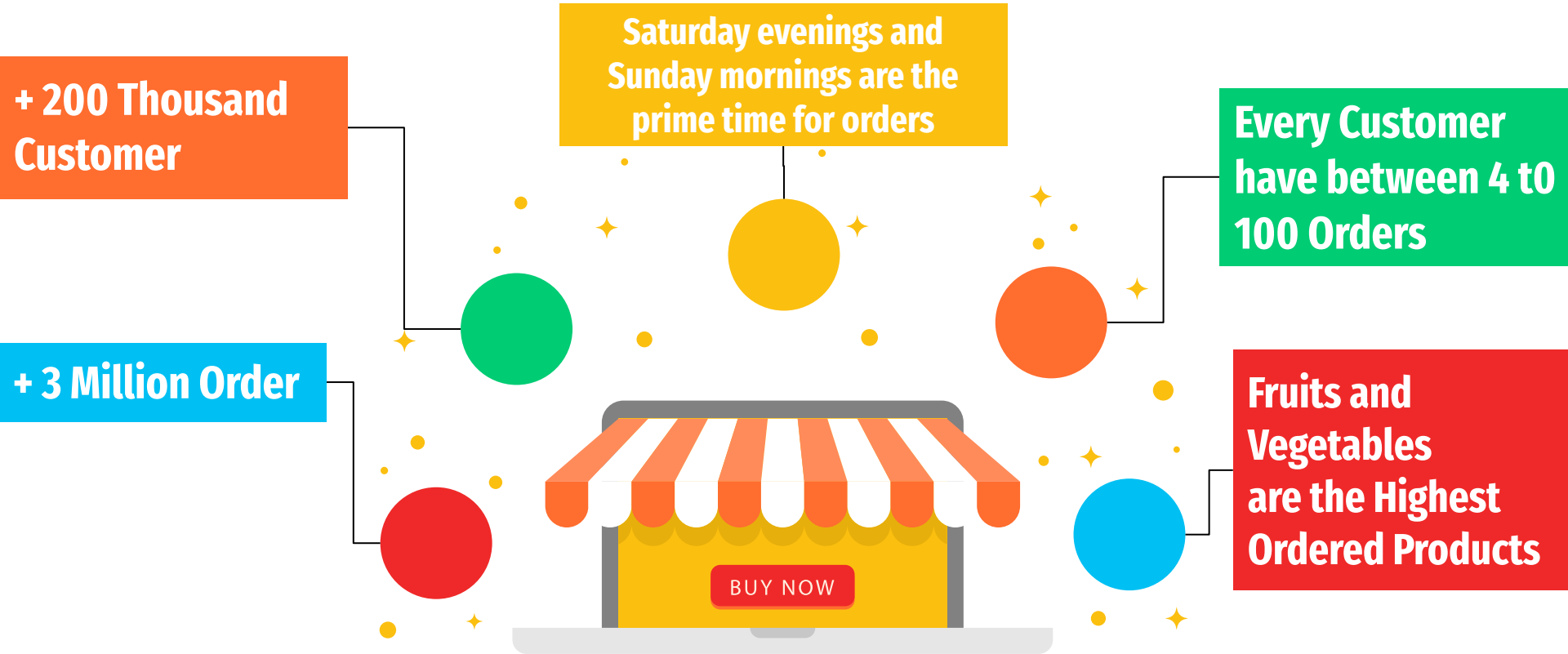


Data about products in each cart

**Cart Products**



# Exploratory Data Analysis



# Features

**In\_cart:** is product in cart? (**the target**)

**Order\_id:** unique identifier for each order

**Product\_id:** unique identifier for each product

**Add\_to\_cart\_order:** the sequence in which they have been added to the cart in that order.

**Reordered:** customer reorders count for this product

**User\_id:** unique identifier for each user

**Order\_number:** user order number

**Order\_dow:** order day of week

**Order\_hour\_of\_day:** order hour of day

**Days\_since\_prior\_order:** days between this order and the prior.

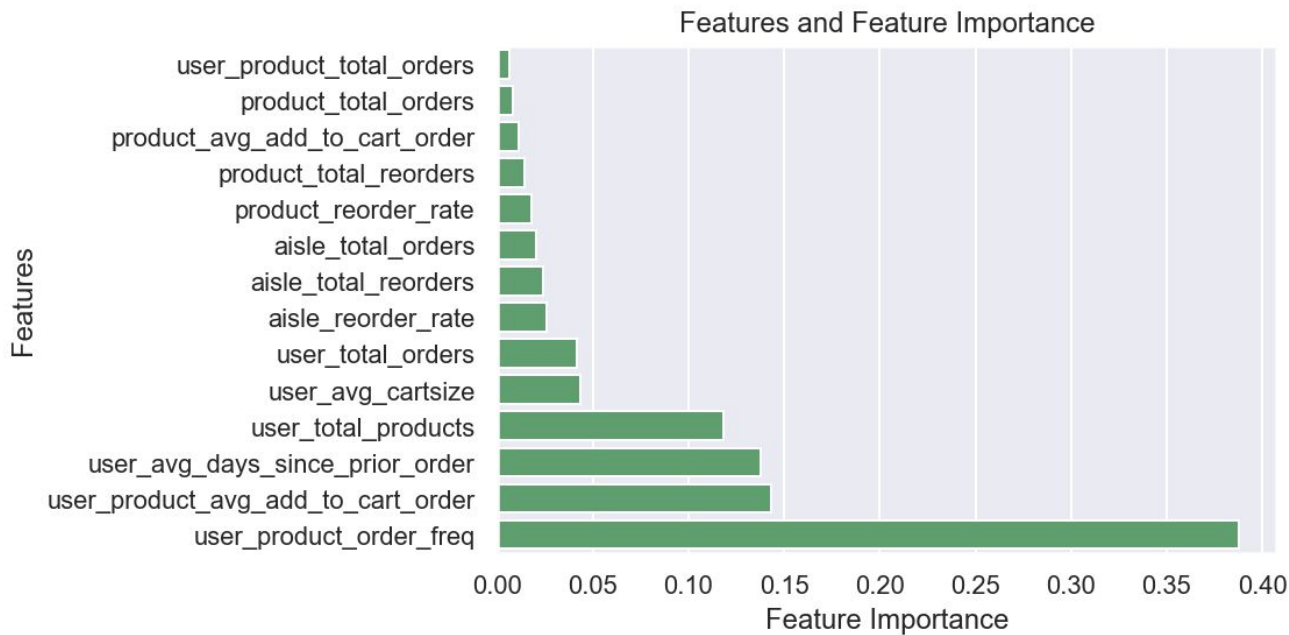


# Feature Engineering

User\_product\_order\_freq  
Product\_reorder\_rate  
User\_total\_products  
User\_avg\_days\_since\_prior\_order  
User\_avg\_cartsize  
User\_product\_total\_orders  
Product\_avg\_add\_to\_cart\_order  
User\_product\_avg\_add\_to\_cart\_order  
Product\_total\_orders  
Product\_total\_reorders  
User\_total\_orders  
Aisle\_reorder\_rate  
Aisle\_total\_reorders  
Aisle\_total\_orders

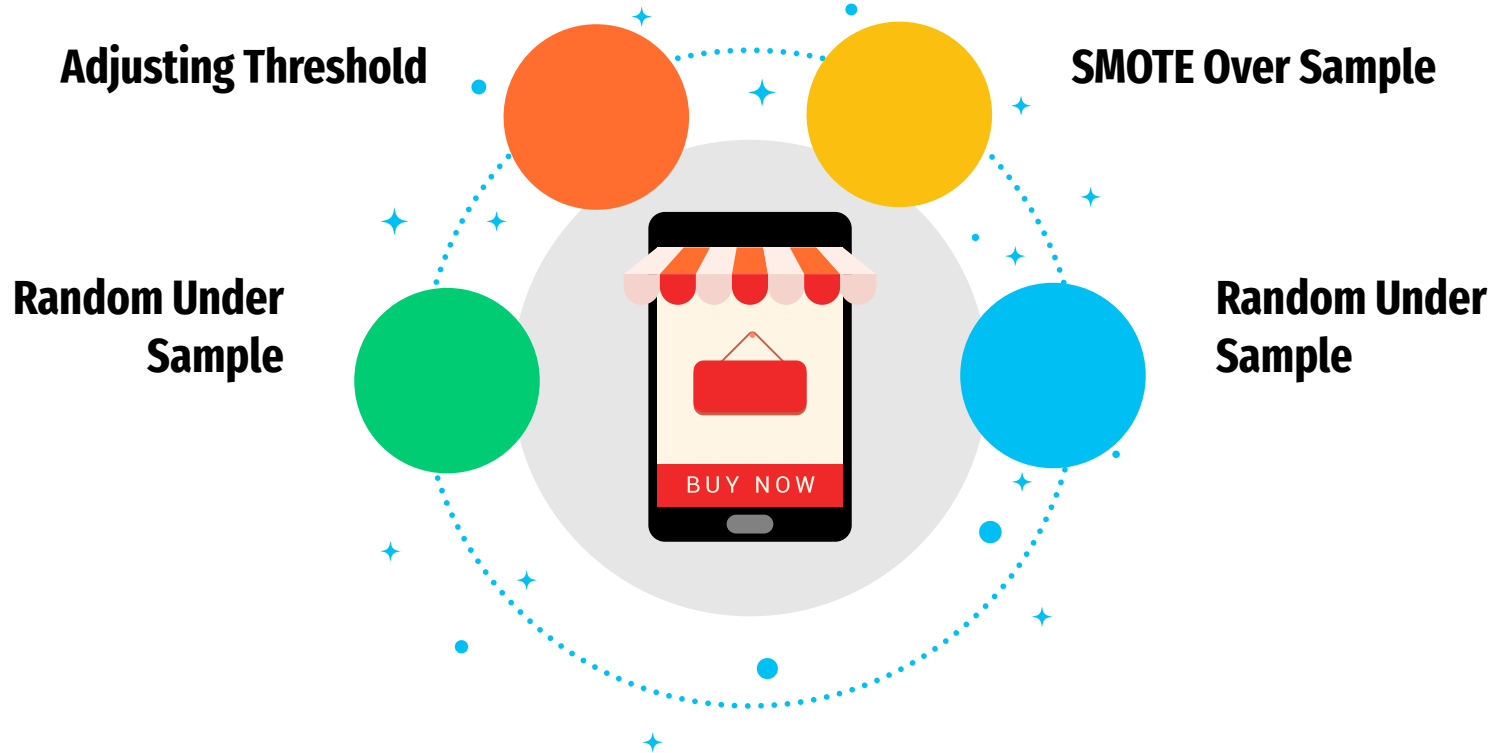


# Feature Importance



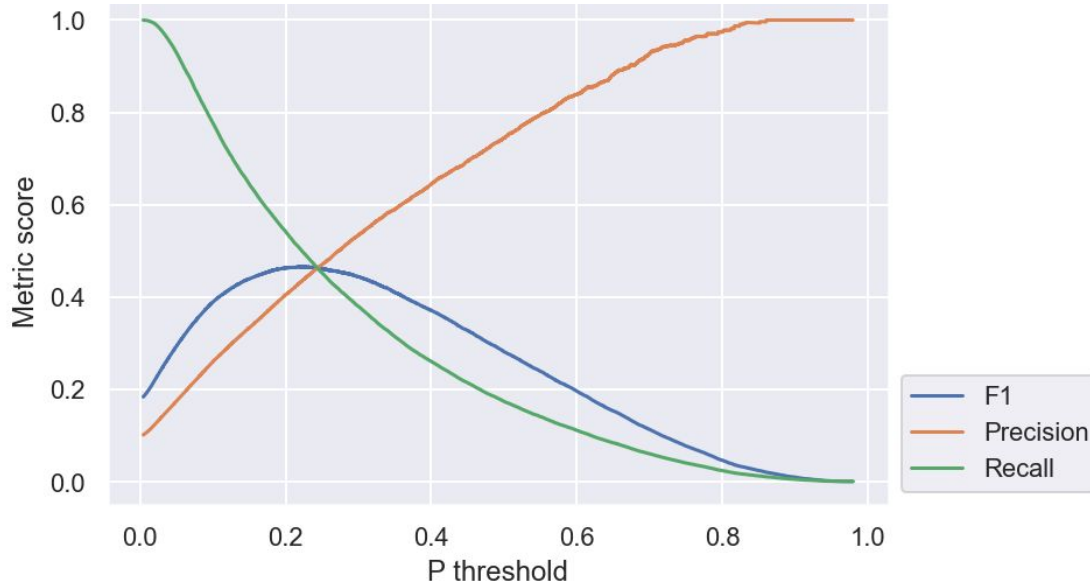


# Class Imbalance Methods



# Class Imbalance - Adjusting Threshold

## XGB Metric Scores vs Probability Threshold

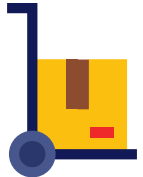
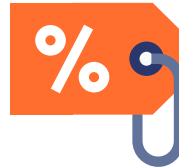


**Best Threshold : 0.22**

**Best Train F1 Score : 0.465**



# Classification Models



# Model Benchmarking

## Scores Based on Competition Leaderboards



Rank	F1 Score
1ST	0.40914
2ND	0.40820
3RD	0.40810
4TH	0.40744



Included in The Prize



Not Included in The Prize



## F1 Scores (Testing)



Model	Adjust threshold	SMOTE	Random Over Sample	Random Under Sample
LR	0.327740	0.317308	0.318361	0.294702
DT	0.223622	0.223012	0.221520	0.219156
RF	0.239094	0.345359	0.256328	0.350634
KNN	0.085911	0.217676	0.214491	0.218282
XGB	0.370872	0.334560	0.337023	0.329331

## Best F1 Scores (Testing)



Model	Adjust threshold	SMOTE	Random Over Sample	Random Under Sample
LR	0.327740	0.317308	0.318361	0.294702
DT	0.223622	0.223012	0.221520	0.219156
RF	0.239094	0.345359	0.256328	0.350634
KNN	0.085911	0.217676	0.214491	0.218282
XGB	0.370872	0.334560	0.337023	0.329331

# Future Work

**Use Multiprocessing  
to Handle Big Data**

**Cross Validation  
Based on Time**

**Hyperparameter  
tuning**

