



# RETAILROCKET: E-COMMERCE USER PURCHASE PREDICTION

---

NURABIDAH JAMIL  
5 AUGUST 2021

# BACKGROUND

---

An e-commerce company wishes to understand visitors behaviour on its site and the factors that drive these site visitors to make a purchase. This is so that they can better target ads and improve future online sales.

# PROBLEM STATEMENT

Using Retailrocket e-commerce dataset from Kaggle, the goal is to build and compare different binary classification models that would best predict whether a user will make a purchase on the e-commerce site.



# TABLE OF CONTENTS

---

- 01 DATA  
Data Cleaning  
EDA
- 02 FEATURE ENGINEERING  
User Groups
- 03 MODELS  
Summary  
Model Selected
- 04 CONCLUSIONS /  
RECOMMENDATIONS



## EVENT

	timestamp	visitorid	event	itemid	transactionid
0	1433221332117	257597	view	355908	NaN
1	1433224214164	992329	view	248676	NaN



## CATEGORY TREE

	categoryid	parentid
0	1016	213.0
1	809	169.0

# DATA



## ITEM PROPERTIES

	timestamp	itemid	property	value
0	1435460400000	460429	categoryid	1338
1	1441508400000	206783	888	1116713 960601 n277.200



## ITEM PROPERTIES

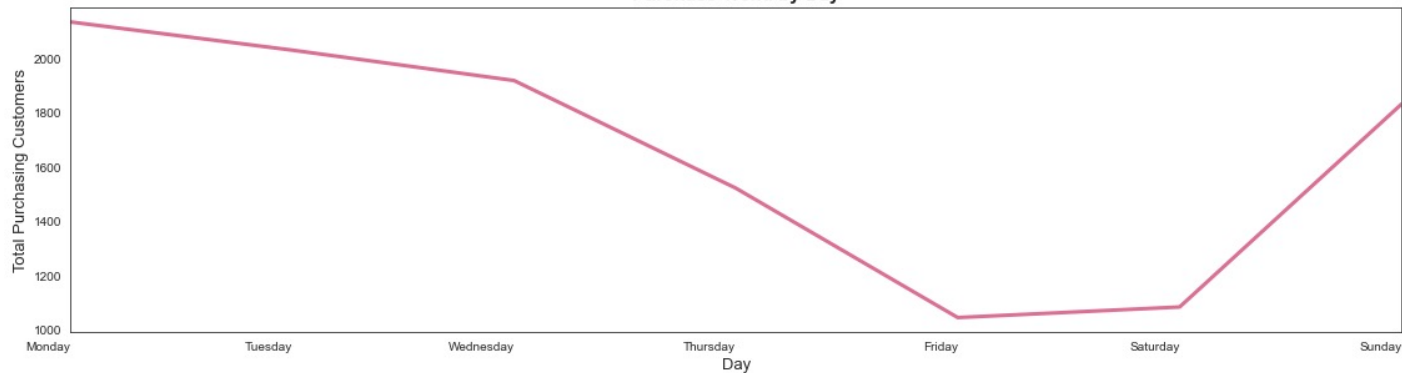
	timestamp	itemid	property	value
0	1433041200000	183478	561	769062
1	1439694000000	132256	976	n26.400 1135780

# FINAL DATASET

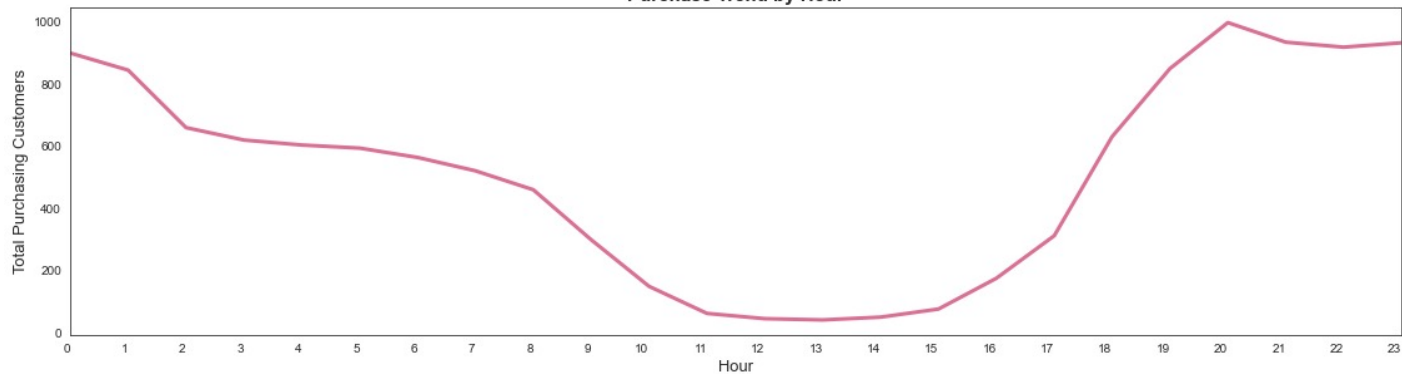
	visitorid	start_session	end_session	session_duration_m	cart	transaction	view	cart_cat	transaction_cat	view_cat
0	0	2015-09-11 23:49:49	2015-09-11 23:55:17	5.47	0	0	[285930, 357564, 67045]	0	0	[1188, 256, 333]
1	1	2015-08-13 20:46:06	2015-08-13 20:46:06	0.00	0	0	[72028]	0	0	[1192]
2	2	2015-08-07 20:51:44	2015-08-07 21:20:57	29.22	0	0	[325215, 325215, 259884, 216305, 342816, 34281...	0	0	[299, 299, 299, 299, 444, 444, 299, 299]
3	3	2015-08-01 10:10:35	2015-08-01 10:10:35	0.00	0	0	[385090]	0	0	[1171]
4	5	2015-07-17 04:45:56	2015-07-17 04:45:56	0.00	0	0	[61396]	0	0	[646]

# PURCHASE TRENDS

Purchase Trend by Day

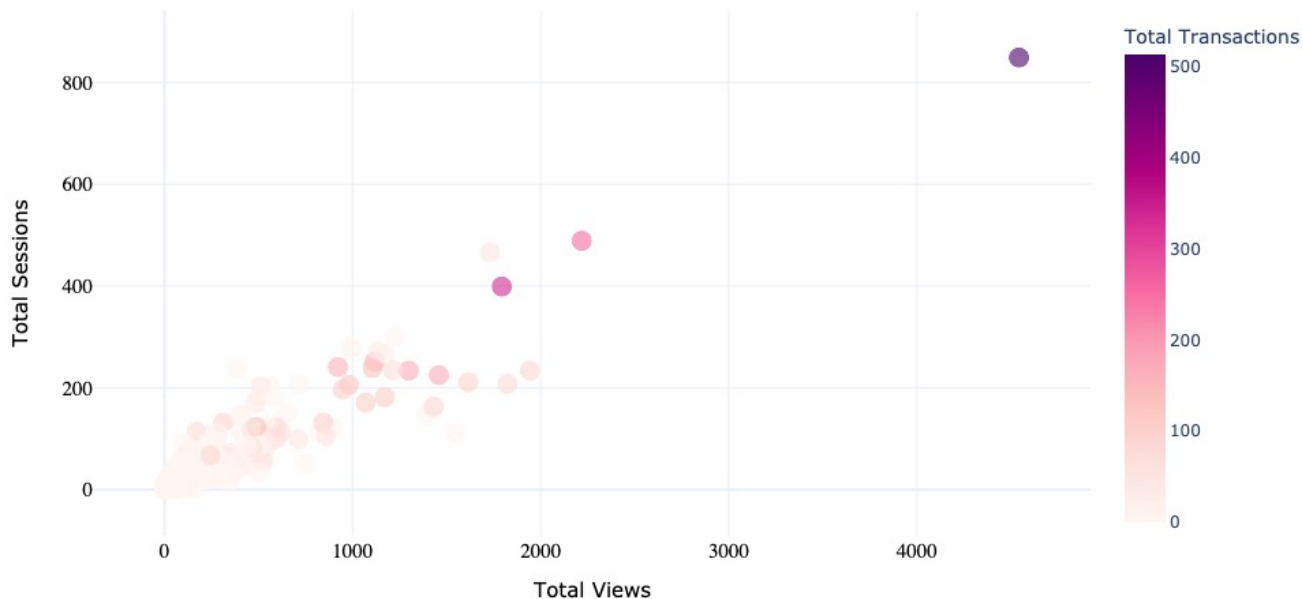


Purchase Trend by Hour



# ABNORMAL USERS

Users By Total Sessions, Views & Transactions



DATA  
REPORTING  
PERIOD

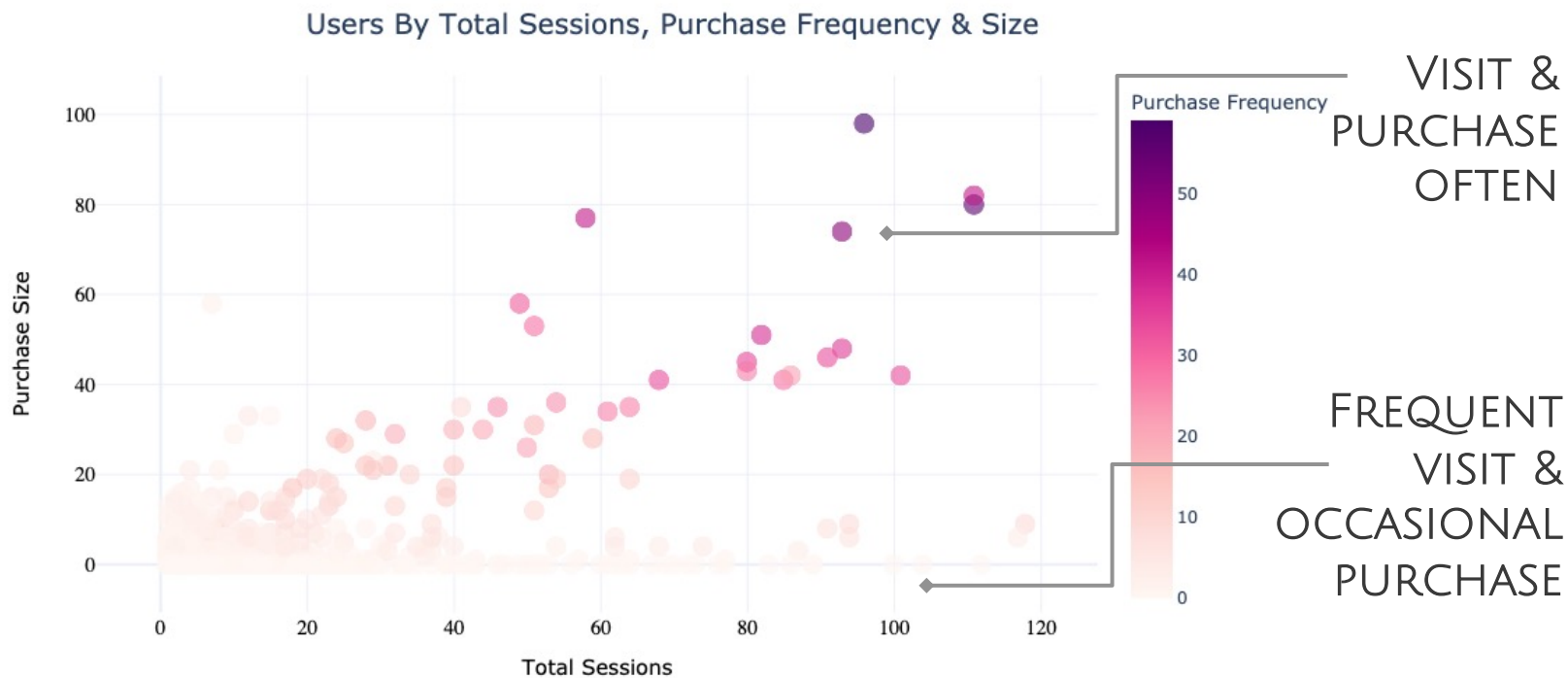
Less than 4  
months

ABNORMAL  
USERS

More than 150  
total sessions  
and/or 150 total  
transactions



# USERS






# USER GROUPS

---

## ▲ K-MEANS CLUSTERING

- Total clusters: 7

## ▲ CLUSTERS EXAMPLE:

GROUP NO.	 GROUP 3	 GROUP 1	 GROUP 0
VISIT FREQUENCY	Almost daily	Almost daily	Once
PURCHASE FREQUENCY	2-3 times a month	1-2 times in total	None yet

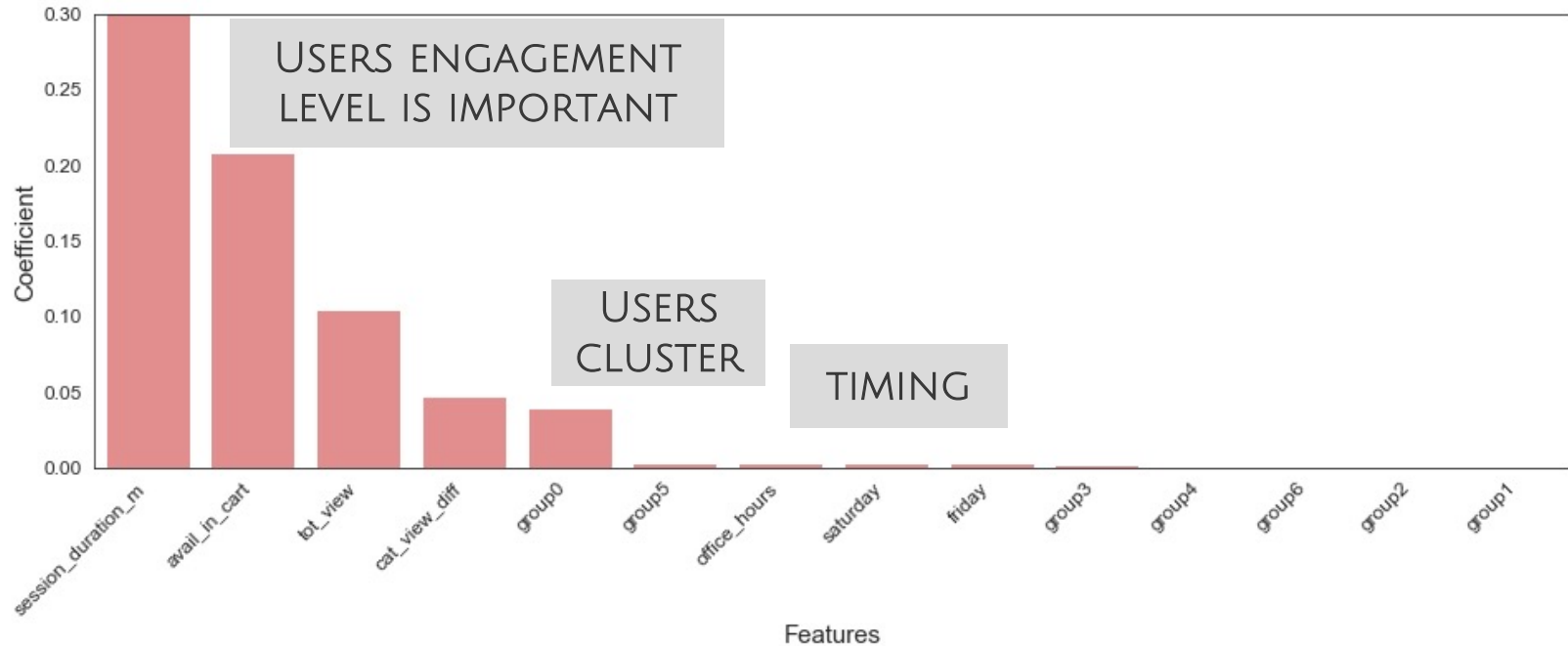
# MODEL SUMMARY

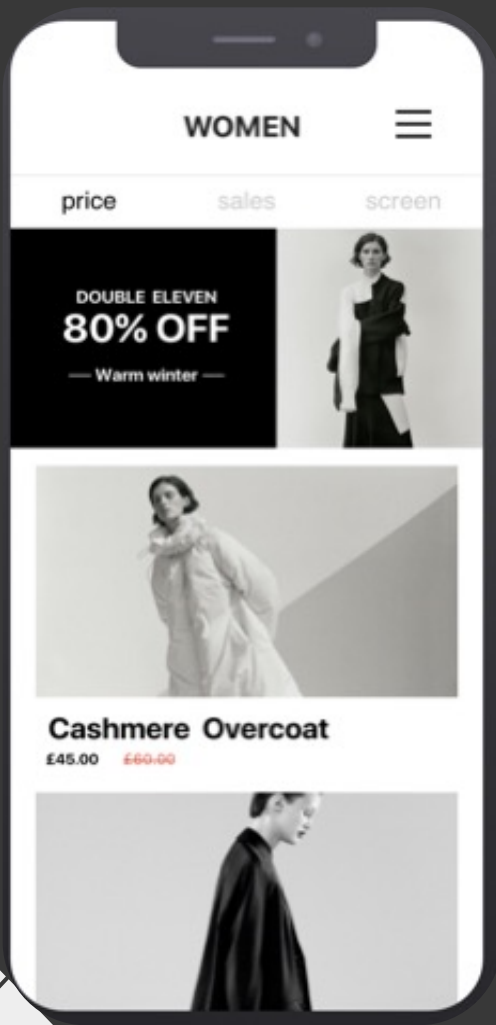
MODEL	TRAIN			TEST		
	PRECISION	RECALL	F1	PRECISION	RECALL	F1
Random Forest	0.27	0.90	0.41	0.14	0.78	0.23
Logistic Regression	0.11	0.87	0.19	0.11	0.82	0.19
XG Boost	0.13	0.93	0.24	0.07	0.97	0.13
ADA Boost	0.11	0.94	0.19	0.05	0.98	0.10

\* Best Parameters: {'rf\_\_max\_depth': 30, 'rf\_\_max\_features': 0.6, 'rf\_\_min\_samples\_leaf': 2, 'rf\_\_min\_samples\_split': 3, 'rf\_\_n\_estimators': 200}

# MODEL SELECTED

## Feature Importance





# CONCLUSIONS & RECOMMENDATIONS

---

- ▲ HIGH FALSE POSITIVE CASES AND MODEL IS NOT ROBUST
- ▲ TIMING, USER ENGAGEMENT AND ITEM AVAILABILITY ARE IMPORTANT
- ▲ AVOID LAUNCHING SHORT-TERM CAMPAIGNS AT CERTAIN TIME
- ▲ IMPROVE USERS ENGAGEMENT ON THE SITE

# CONCLUSIONS & RECOMMENDATIONS

---

## ▲ IMPLEMENT DIFFERENT MARKETING STRATEGIES FOR EACH USER GROUP:



GROUP 3

- Loyalty programs
- Market new products



GROUP 1

- Target them for short-term campaigns like flash sales

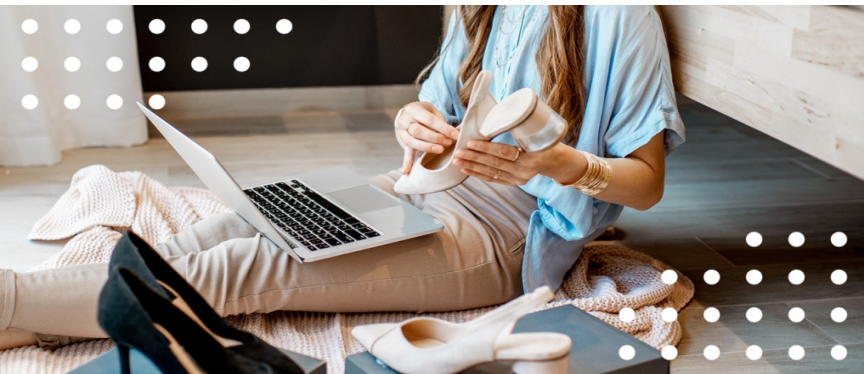


GROUP 0

- Aggressive price incentives

# NEXT STEPS

---



## INCORPORATE OTHER DATA

- The marketing channel the users come from
- Purchase value

## REFRESH CUSTOMER SEGMENTATION

- Incorporate missing key information

## OTHER MODELS

- Neural Network

THANK YOU

---



# DATA DICTIONARY

---

VARIABLE	DESCRIPTION
group0	A dummy variable to indicate whether the users belong to group 0.
group1	A dummy variable to indicate whether the users belong to group 1.
group2	A dummy variable to indicate whether the users belong to group 2.
group3	A dummy variable to indicate whether the users belong to group 3.
group4	A dummy variable to indicate whether the users belong to group 4.
group5	A dummy variable to indicate whether the users belong to group 5.
group6	A dummy variable to indicate whether the users belong to group 6.
friday	A dummy variable to indicate whether the session occurs on a Friday.
saturday	A dummy variable to indicate whether the session occurs on a Saturday.
office_hours	A dummy variable to indicate whether the session occurs within 8 am to 5 pm.
avail_in_cart	The percentage of items in the cart that are available.
tot_view	The total number of items viewed during a particular session.
cat_view_diff	A variable to show how different are the items being viewed (scaled between 0.005 and 1 , with 1 indicating that all items viewed are different).
session_duration_m	The duration during which there are regular active interactions coming from a user on the website.