

Project Report on Customer Retention

For FLIPROBO TECHNOLOGIES



Shruti Raj

Batch no: 1842
Internship 29

Acknowledgement

I would like to express my special gratitude to the “Flip Robo” team, who has given me this opportunity to deal with a beautiful dataset and it has helped me to improve my analyzing skills.

A huge thanks to my academic team “Data trained” who has helped me grow from a Non Coder to what I am Now. Lastly, I would like to extend my Heartfelt thanks to my Husband and kids because without their support this project would not have been successful. And also thank you to many other persons who have helped me directly or indirectly to complete the project.

Contents

1.INTRODUCTION

2.Objective of the study

3.Data Analysis and Interpretation

3.1 Procedure

3.2 Understanding the dataset

3.3 Data Exploration

3.4 Data Visualization

4.Feature engineering

4.1 Encoding

4.2 Correlation

4.3 Scaling the Data

5.Model Building

5.1 HyperParameter Tuning

5.2 Saving the best Model

6.Feature Importance

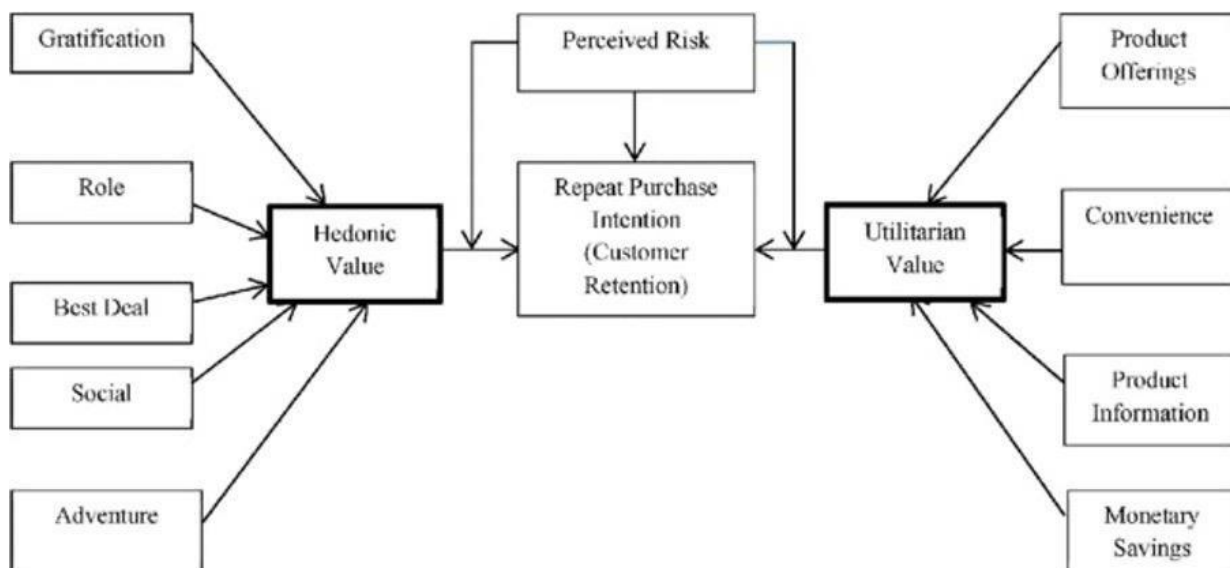
7.Conclusion

8.REFERENCES

1. INTRODUCTION

Customer satisfaction has emerged as one of the most important factors that guarantee the success of online stores; it has been posited as a key stimulant of purchase, repurchase intentions and customer loyalty. A comprehensive review of the literature, theories and models have been carried out to propose the models for customer activation and customer retention. Five major factors that contributed to the success of an e-commerce store have been identified as: service quality, system quality, information quality, trust and net benefit. The research furthermore investigated the factors that influence the online Customers' repeat purchase intention. The combination of both utilitarian value and hedonistic values are needed to affect the repeat purchase intention (loyalty) positively. The data is collected from the Indian online shoppers. Results indicate the e-retail success factors, which are very much critical for customer satisfaction.

E-commerce is expanding steadily in the country. Customers have an increasing choice of products at competitive rates. E-commerce is probably creating the biggest revolution in the retail industry, and this trend will continue in the years to come.



2. Objective of the study

- To apply Analytical skills to give findings and conclusions in detailed data analysis written in jupyter notebook.
 - To find out key factors influencing shoppers while choosing online retailers for buying products
 - To study the consumer perceptions towards online retailers in India
 - To compare the customer perceptions regarding selected online retailers

3. Data Analysis and Interpretation

3.1 Procedure

- The datasets were loaded to the Jupyter notebook and then checked for Missing values.
- The data is then filtered to remove the duplicate entries. The cleaned data will be further processed to get more details on the customers for better satisfaction from the e-commerce sellers.
- I have renamed the columns for better understanding of myself. In the dataset there were no numerical columns.
- The unique values of each feature is analyzed.
- All the categorical columns had been analyzed using univariate analysis with categorical plots to get better insight on the dataset.
- After visualizing the features it brought us a good insight into what actual customers are expecting from the e-commerce sellers. Observing the plots, we were able to get good measures to have customer retention.
- In this dataset we don't need to check outliers and skewness as all the columns are categorical.
- I then used Encoding techniques and encoded all datas.
- I have removed Pincode as it does not have any influence on my model.
- I have used Scaling before Building the model.
- I have used RandomForestClassifier to train and test my Model and I got the Accuracy of 100% and Cross Validation Score of 100%.
- Further I have created a feature Importance plot inferred from my Final Model.

3.2 Understanding the dataset

Features in the Dataset are:

- [illegible]

48. 'From the following, tick any (or all) of the online retailers you have shopped from;
49. Easy to use website or application',
50. Visual appealing web-page layout',
51. 'Wild variety of product on offer',
52. Complete, relevant description information of products',
53. Fast loading website speed of website and application',
54. Reliability of the website or application',
55. Quickness to complete purchase',
56. Availability of several payment options',
57. Speedy order delivery ',
58. Privacy of customers' information',
59. Security of customer financial information',
60. Perceived Trustworthiness',
61. Presence of online assistance through multi-channel',
62. Longer time to get logged in (promotion, sales period)',
63. Longer time in displaying graphics and photos (promotion, sales period)',
64. Late declaration of price (promotion, sales period)',
65. Longer page loading time (promotion, sales period)',
66. Limited mode of payment on most products (promotion, sales period)',
67. Longer delivery period',
68. Change in website/Application design',
69. Frequent disruption when moving from one page to another',
70. Website is as efficient as before',
71. Which of the Indian online retailer would you recommend to a friend?'

This Dataset contains data of Indian E-commerce websites that has the user data such as Gender, Age, Location and other features (Total=71) of data about the activity of the user in the website. The target is “Which website they would recommend to their friend”. We could analyze and find which factors affect the recommendation of the users. In the Jupyter notebook, I performed various Data Analysis techniques to find the factors for user's recommendation of websites.

The study was based on a survey conducted for online retailers like Amazon, Flipkart, Myntra, Paytm and Snapdeal in different cities. A sample of 269 customers was sent a questionnaire using convenience sampling technique. There were 181 Females and 88 Males in the sample. The responses were measured on a five-point Likert scale, ranging from 1 (strongly disagree) to 5 (strongly agree), to rate the extent of customers' perceptions towards services provided by online retailers.

Libraries Required

To run the program, we need `basic libraries to be imported like:

Importing Libraries and Loading the Dataset

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

import warnings
warnings.filterwarnings('ignore')
```

I have used Label Encoder, Ordinal Encoder and also MinMaxScaler.

```
cat=[i for i in train_df.columns if train_df[i].dtypes=='O']
from sklearn.preprocessing import OrdinalEncoder,LabelEncoder
from sklearn.preprocessing import MinMaxScaler
```

For building the model following Libraries were imported:

Building the Model

```
: from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split,cross_val_score
from sklearn.metrics import accuracy_score,confusion_matrix,classification_report
from sklearn.model_selection import RandomizedSearchCV
```

Software and Hardware Required

Hardware required: -

1. Processor — core i5 and above
2. RAM — 8 GB or above
3. SSD — 250GB or above

Software required: -

- 1.Anaconda

3.3 Data Exploration

I have imported the Data which was in excel format and other necessary libraries. All the Statistical analysis like Shape, Info, Null values, Nunique, Value Counts have been observed in Data Exploration. The value counts showed that there are duplicate entries in my dataset.

```
In [10]: for i in df.columns:
          print(df[i].value_counts())
          print('*****')
          *****
          Less than 10 times    114
          31-40 times          63
          41 times and above   47
          11-20 times          29
          21-30 times          10
          42 times and above    6
          Name: 6 How many times you have made an online purchase in the past 1 year?, dtype: int64
          *****
          Mobile internet      142
          Wi-Fi                76
          Mobile Internet       47
          Dial-up               4
          Name: 7 How do you access the internet while shopping on-line?, dtype: int64
```

- I have replaced the duplicates with the original entry.

There are several duplicate values/mistypes in the Dataset. We have to change it and replace them.

- 6 How many times you have made an online purchase in the past 1 year?
- 7 How do you access the internet while shopping on-line?

```
df["6 How many times you have made an online purchase in the past 1 year?"].replace("42 times and above", "41 times and above", inplace=True)
df["7 How do you access the internet while shopping on-line?"].replace("Mobile internet", "Mobile Internet", inplace=True)
```

There were no Missing values in my Dataset.

```
In [7]: df.isnull().sum().any()
Out[7]: False
```

- I have renamed the column names for better understanding as they had many spaces, mistypes.

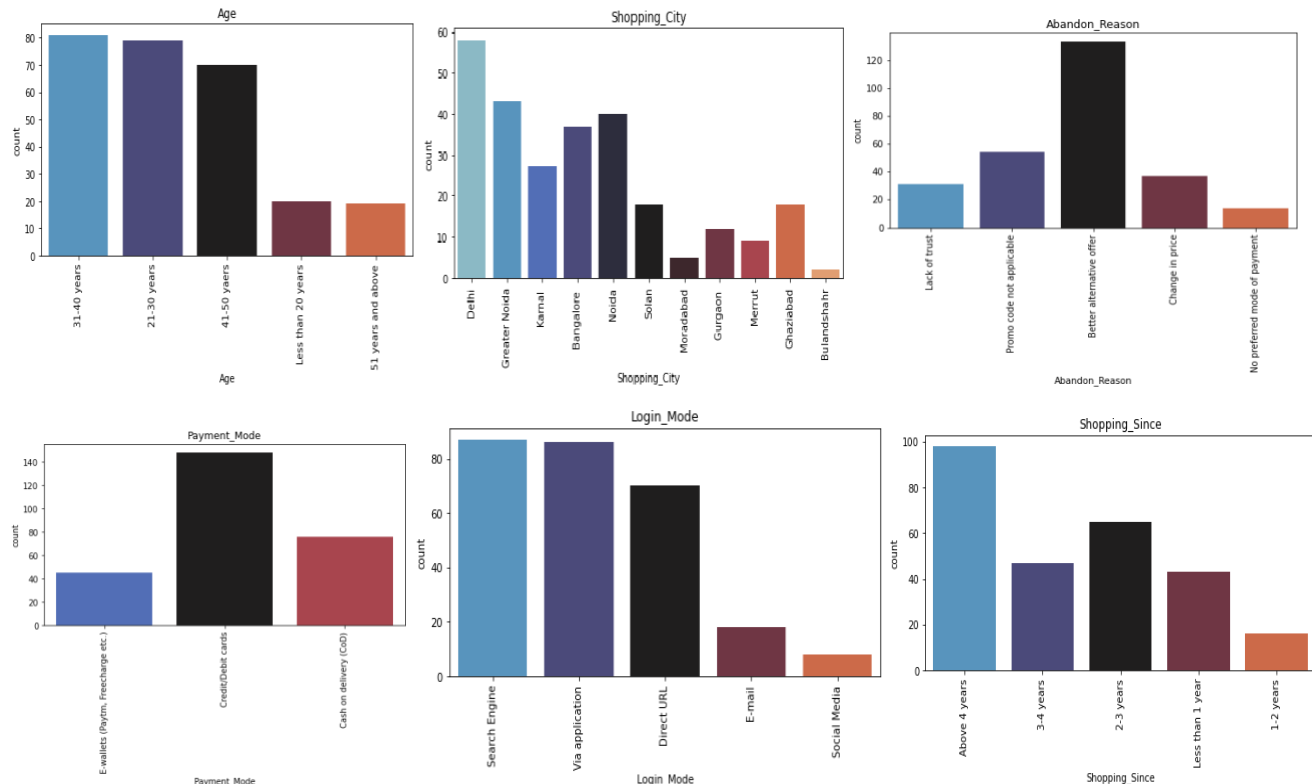
```
#Renaming the column name
Rename=[ 'Gender', 'Age', 'Shopping_City', 'PinCode', 'Shopping_Since', 'Shopping_Frequency', 'Internet_Access', 'Device_Used',
         'Screen_Size', 'Operating_System', 'Browser_Used', 'Channel_FirstUsed', 'Login_Mode', 'TimeSpent_ForPurchase',
         'Payment_Mode',
         'Abandon_Frequency',
         'Abandon_Reason', 'Content_Readability', 'Similar_ProductInfo', 'Seller_ProductInfo', 'ProductInfo_Clarity',
         'Ease_Navigation',
         'Loading_ProcessingSpeed', 'UserFriendly_Interface', 'Conveninet_PaymentMode', 'TimelyFulfilment_Trust',
         'Customer_Empathy',
         'CustPrivacy_Guarantee', 'VariousChannel_Responses', 'Benefit_Discount', 'Enjoy_OnlineShopping',
         'Convenience_Flexibility',
         'Returns_ReplacementPolicy',
         'Loyalty_ProgramAccess',
         'QualityInfo_Satisfaction', 'WebsiteQuality_Satisfaction', 'NetBenefit_Satisfaction', 'User_Trust',
         'Product_SevealCategory', 'Relevant_ProductInfo', 'Monetary_Savings',
         'Patronising_Convenience', 'Adventure_Sense', 'Enhances_SocialStatus', 'Gratification_Shopping',
         'Role_Fulfilment', 'Money_Worthy', 'Shopped_From', 'Easy_WebApp',
         'Visually_AppealingWebApp', 'Product_Variety', 'Complete_ProductInfo', 'Fast_WebApp', 'Reliable_WebApp',
         'Quick_Purchase', 'PaymentOptions_Availability',
         'Fast_Delivery', 'CustInfo_Privacy', 'FinancialInfo_Security', 'Perceived_Trustworthiness', 'MultiChannel_Assist',
         'Long_LoginTime', 'LongPhoto_DisplayTime',
         'LatePrice_Declare', 'Long>LoadingTime', 'Limited_PaymentMode', 'Late_Delivery', 'ChangeWebApp_Design',
         'Page_Disruption', 'WebApp_Efficiency',
         'Recommendation']
```

I have also changed the datatype of Pin Code for better analysis.

3.4 Data Visualization

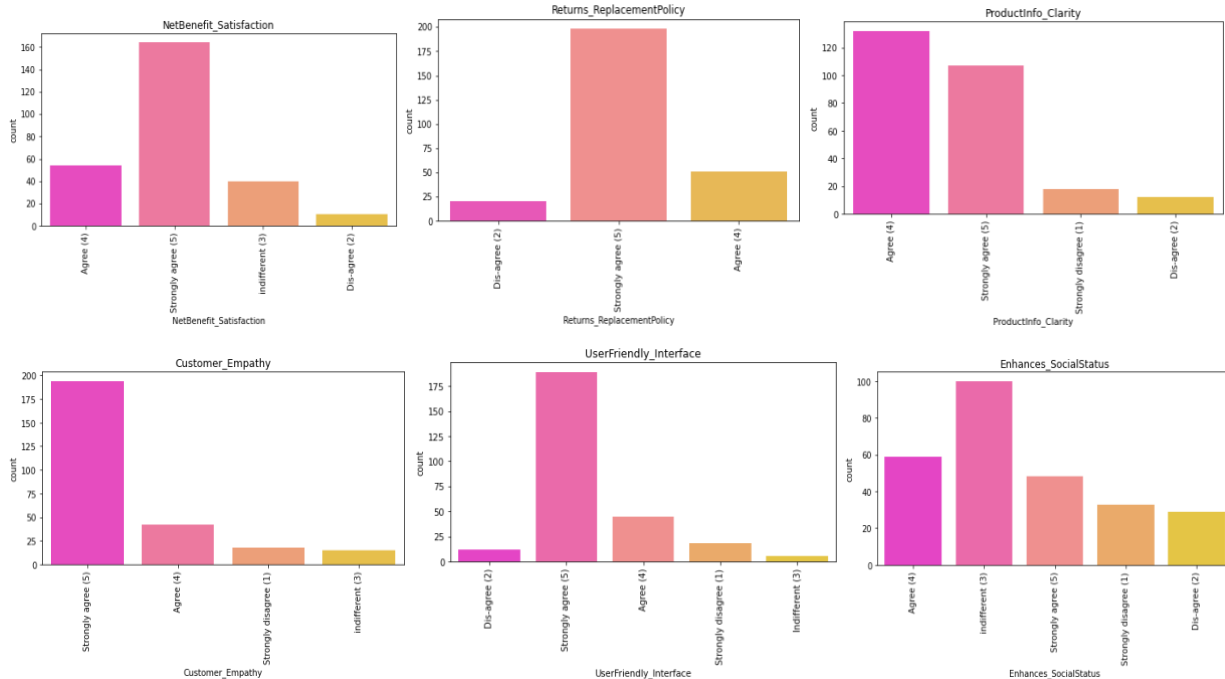
Since all are categorical Datatypes, we are using categorical plotting for analysis.

a) Shoppers Information:



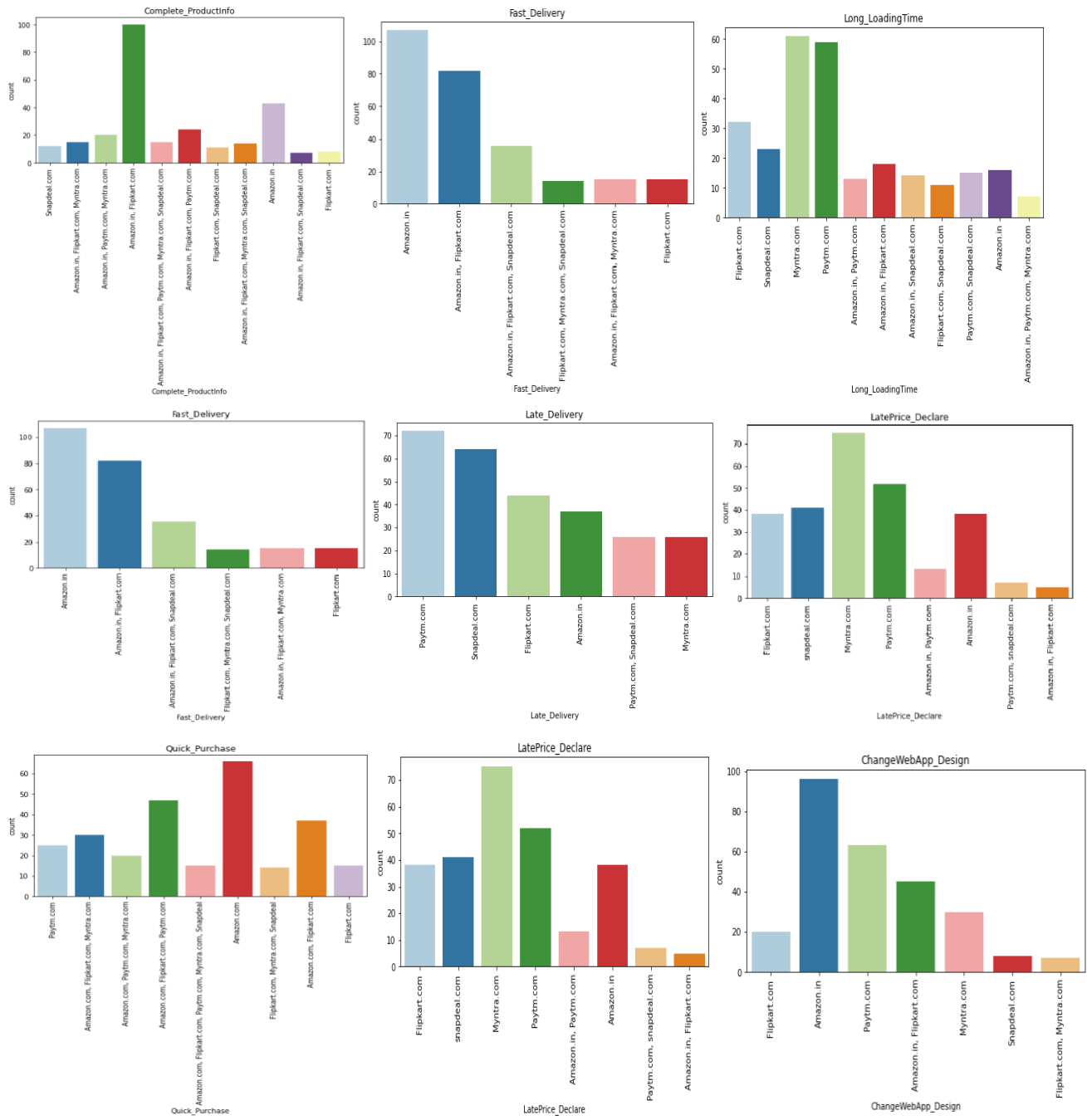
- Number of females is more than males.
- Most of the shoppers are 31-40 years old followed by 21-30 years old. Least number of shoppers are 51 years old.
- According to the data, most shoppers shop from Delhi and greater Noida.
- Most shoppers are shopping online for more than 4 years,
- The shopping frequency for the past 1 year is Less than 10 times for most of the shoppers.
- Mobile Internet is used by most of the shoppers.
- Many use Smartphones followed by laptops.
- Operating systems used by most shoppers are windows and followed by Android.
- Most of the shoppers use Google chrome as their browser.
- Most of the shoppers came to the website for the first time through search engines.
- Most of the shopper's login to the website through Search Engines and Mobile Applications.
- Most shoppers spent more than 15 minutes before purchasing.
- A few of the shoppers notified that they used less than a minute for the purchase.
- E-retail customers mostly use Debit/credit card for purchasing.
- Most of the Shoppers abandon the cart sometimes and very few abandon very frequently.
- Most shoppers abandon the cart because they get a Better alternative Offer.

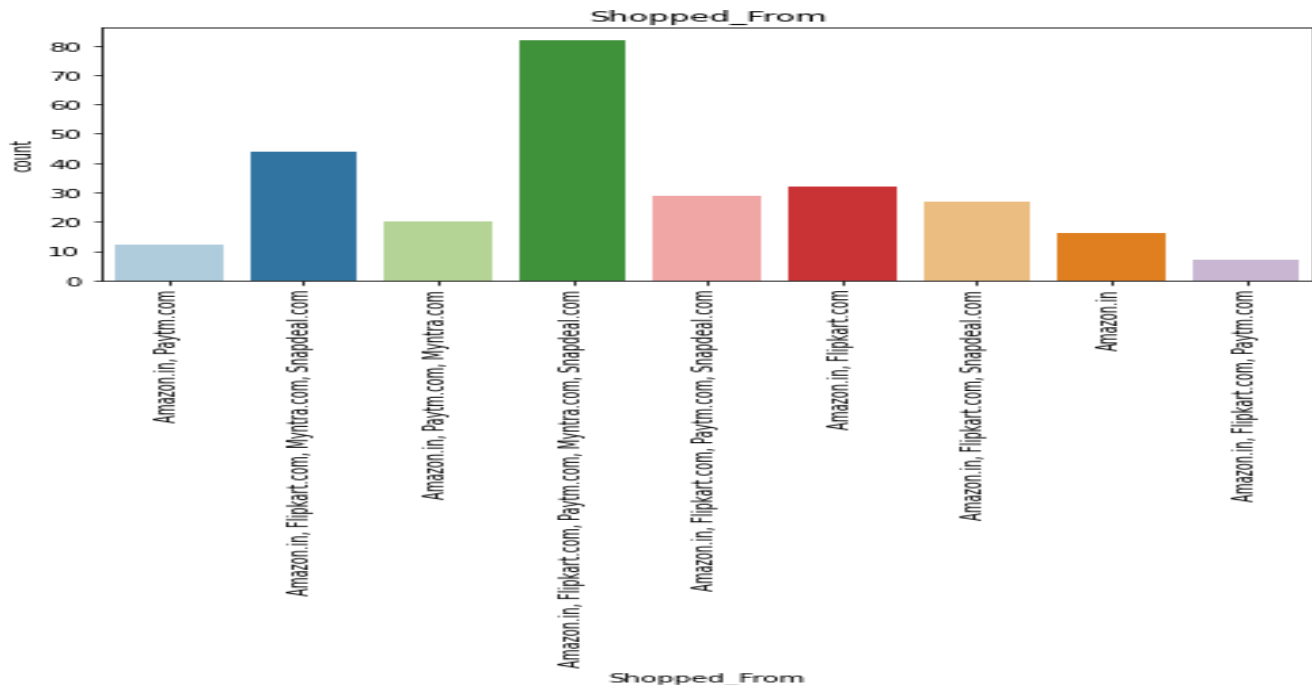
b) Likert Scale:



- Most shoppers **strongly agree** that contents are easy to read and understand.
- Most shoppers **strongly agree** that similar products to the content highlighted is important for product comparison.
- Most shoppers **strongly agree** that seller and product complete information is important for purchase decisions.
- Maximum shoppers **strongly agree** to have convenient payment methods.
- Most shoppers Trust that the online retail store will fulfill its part of the transaction at stipulated time.
- Most Shoppers **strongly agree to have Empathy** (readiness to assist with queries).
- Most Shoppers **strongly agree to have the ability to guarantee** the privacy of the customer.
- Most Shoppers **strongly agree to have Responsiveness**, availability of several communication channels (email, online rep, twitter, phone etc.).
- Most Shoppers **strongly agree that Online shopping gives monetary benefits and discounts.**
- Most Shoppers **strongly agree that Enjoyment** is derived from shopping online.
- Most Shoppers **strongly agree that Shopping online is convenient and flexible.**
- Most Shoppers **strongly agree** that Return and replacement policy of the e-tailer is important for purchase decisions.
- Most Shoppers **strongly agree that Gaining access to loyalty programs** is a benefit of shopping online.
- Most Shoppers **strongly agree that Displaying quality Information** on the website improves customer satisfaction.
- Most Shoppers **agree that users derive satisfaction** while shopping on a good quality website.
- Most Shoppers are **indifferent** that Net Benefit derived from shopping online can lead to users satisfaction.
- Most Shoppers **agree to have Monetary savings.**
- Most Shoppers are **indifferent** towards that Shopping on your preferred e-tailer enhances your social status.

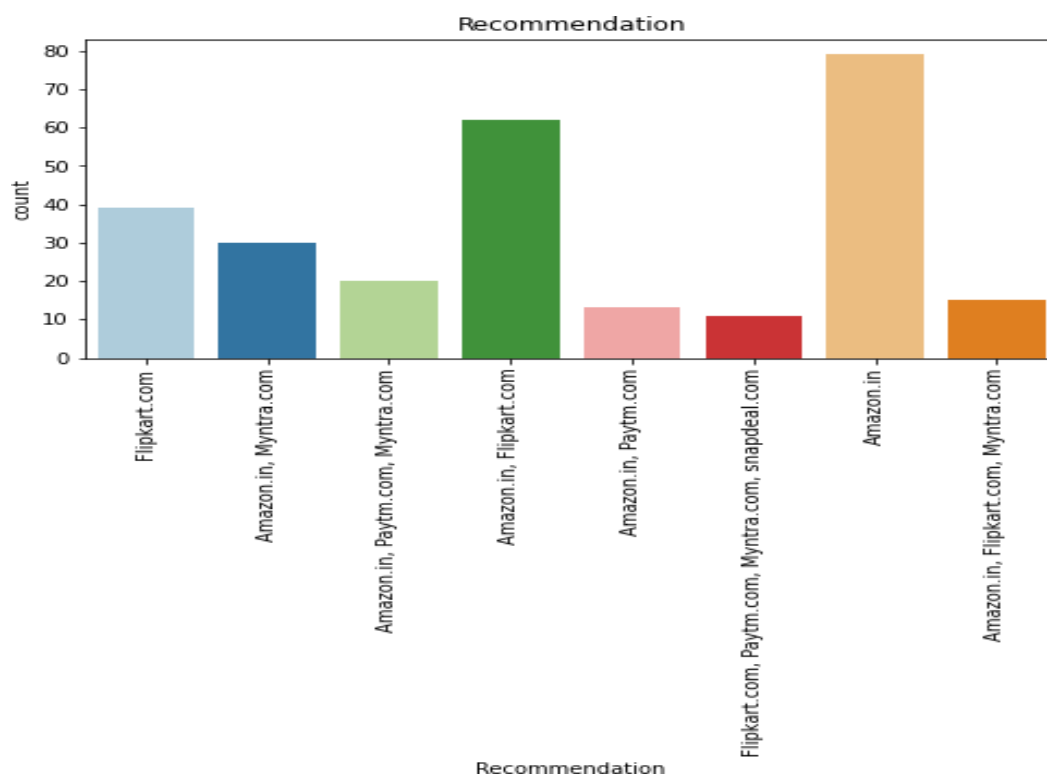
c) Website Preferences:





- Most shoppers shop through **all online ecommerce portals**.
- Most shoppers say all online ecommerce portals are easy to use.
- Most shoppers agree that **Amazon and flipkart are the most visually** appealing sites.
- **Product variety is more in amazon and flipkart** by almost 75% of the shoppers.
- Most shoppers agree that **Amazon and Flipkart have the most product Info**.
- According to most shoppers **Amazon is the most reliable** website for online purchasing.
- Most shoppers agree that **Amazon is the quickest to complete purchases**.
- **Payment option availability is more for Amazon and Flipkart** according to shoppers.
- **Amazon delivers faster** according to most shoppers.
- Most shoppers agree that **Amazon and Flipkart have privacy of customer information**.
- **Financial information is most secure with Amazon** according to most shoppers.
- **Amazon is the most trusted** by most shoppers.
- **All ecommerce portals have multi channels** according to most shoppers.
- **Login time is very high for Amazon** during sales and promotions time.
- **Image display time is more for Amazon and Flipkart** during sales and promotions time.
- **Price declaration is late for Myntra** during sales and promotions time according to most shoppers.
- **Long page loading time for Myntra and Paytm** according to most shoppers.
- **Snapdeal has limited payment options** according to most shoppers.
- **Delivery is late using Paytm**.
- **Almost all have page disruptions and least page disruption is for Snapdeal**.

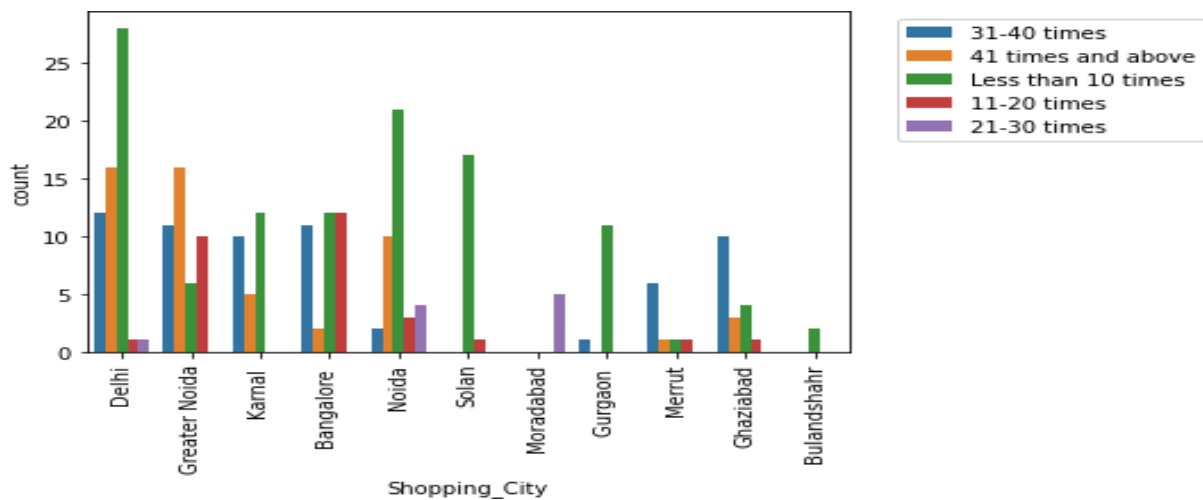
Target Variable/Recommendation



- In this dataset, Amazon.in is the most recommended Ecommerce website followed by Flipkart.com.
- The least recommended websites for online shopping are PayTm, Myntra, and Snapdeal.

Bivariate Analysis

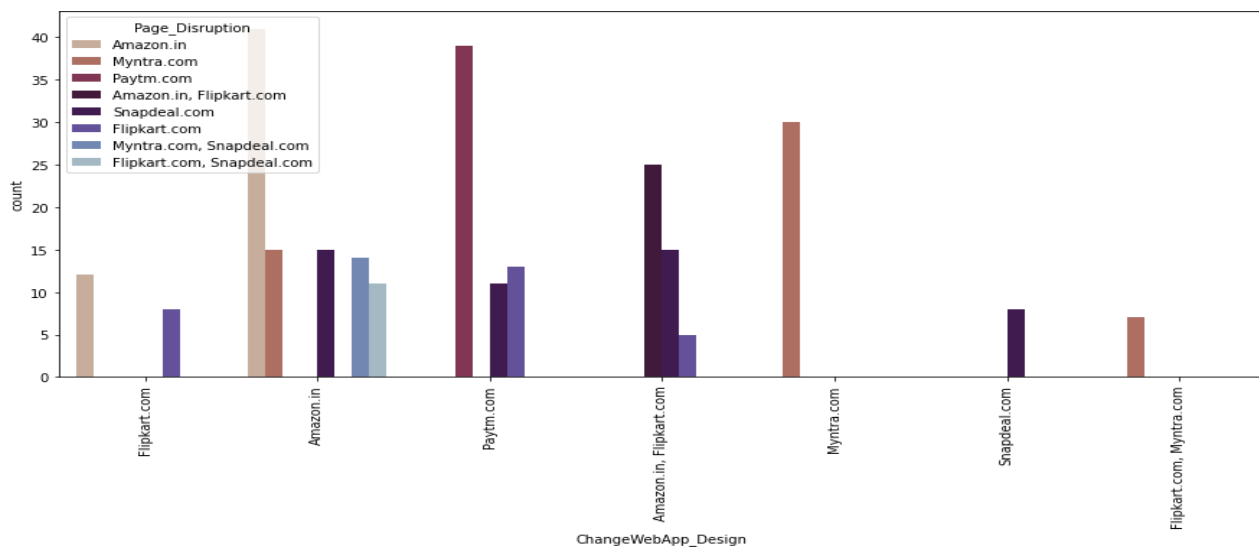
a) Shopping City and Frequency of shopping for the past 1 year:



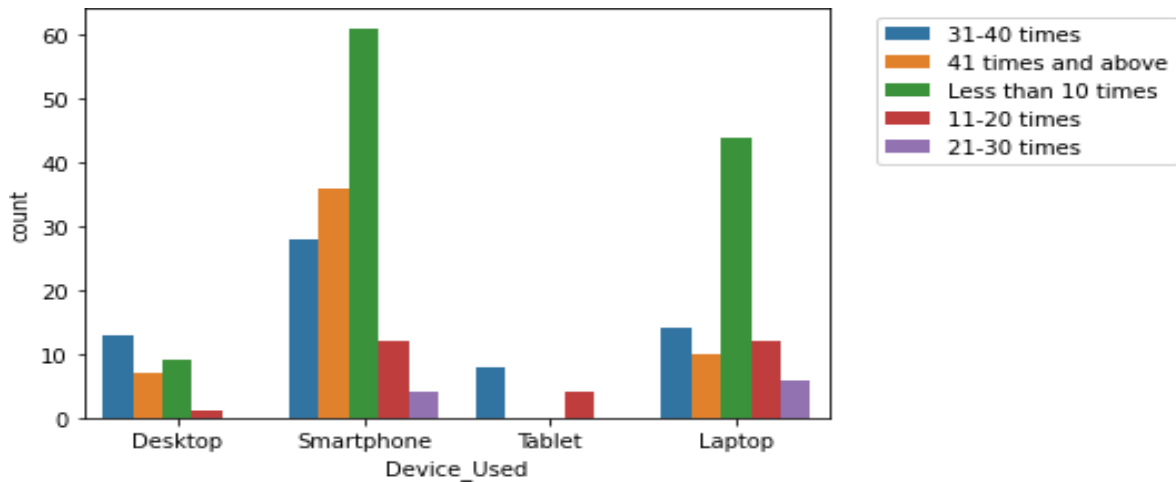
- More people have shopped from Delhi followed by Greater Noida and Noida.
- Least number from Bulandshahr
- Most of the shoppers have shopped for less than 10 times in the past year

b) Change web/App design Vs Page Disruption

Most shoppers say Amazon has the best design for the website and also a major drawback is that it is also the Most page disrupted.

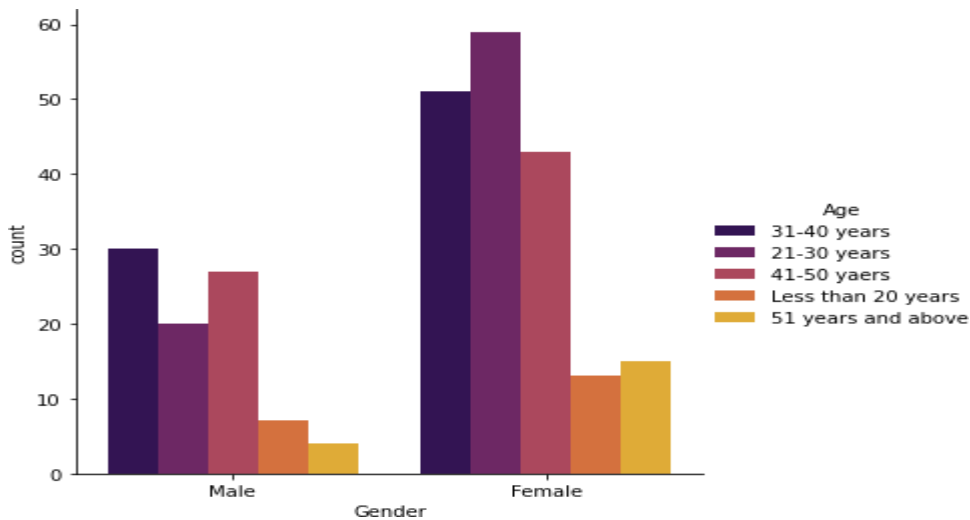


c) Device Used Vs Frequency of shopping



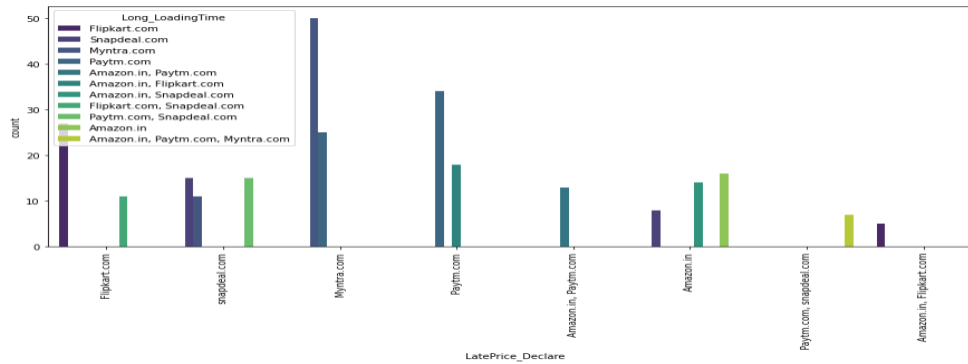
- Most people used smartphones and shopping frequency was less than 10. Second most shopping frequency is above 40 times for smartphone users.
- Tablets users' shopping frequency is 10-20 and 30-40.
- Laptop users shopping frequency is more at 0-10

d) Gender and Age



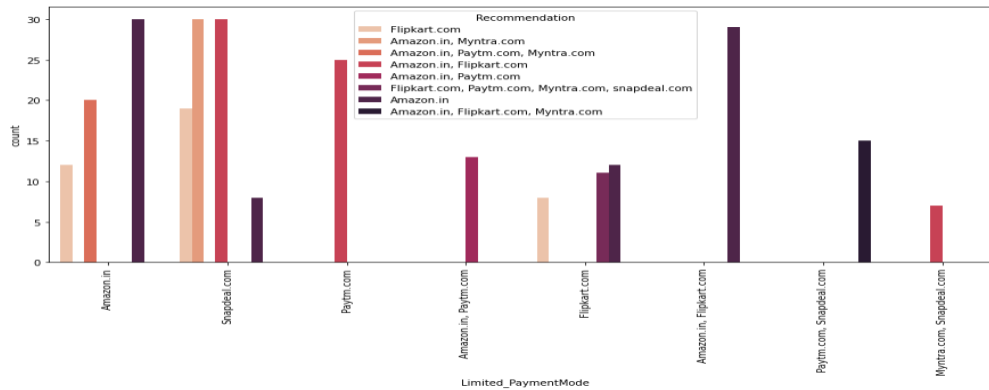
- Accordingly Females are the most online shoppers.
- Males below 20 years and above 51 years use less ecommerce websites.
- Females between the age of 21-30 shop more according to the data.
- Males aged between 31-40 shop more among the males category.

e) Late Price declaration Vs Long Loading Time



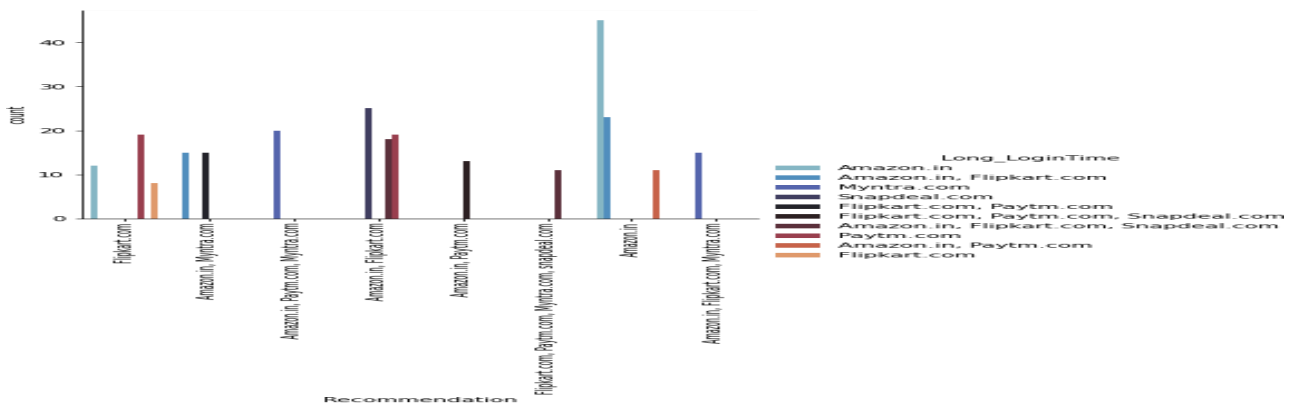
- Major drawback for Myntra.com is Late declaration of price and Long Loading time during sale and Promotion time.

f) Limited Payment mode Vs Recommendation:



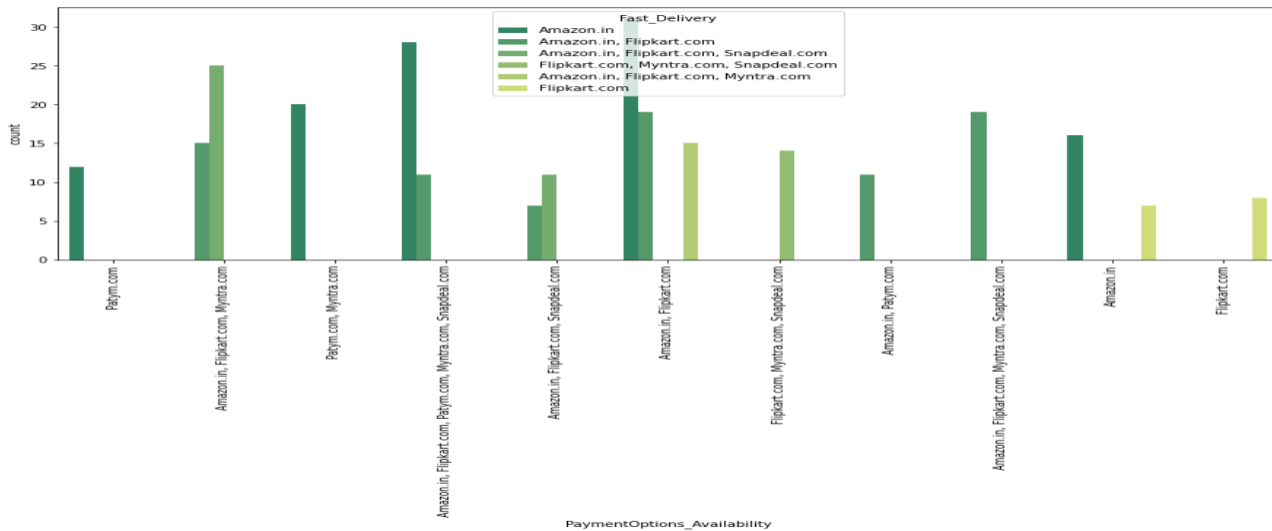
- Most people say that Snapdeal has the least payment options followed by Amazon. Even then, Amazon is the most recommended but Snapdeal comes in least recommended.

g) Recommendation Vs Long Login Time



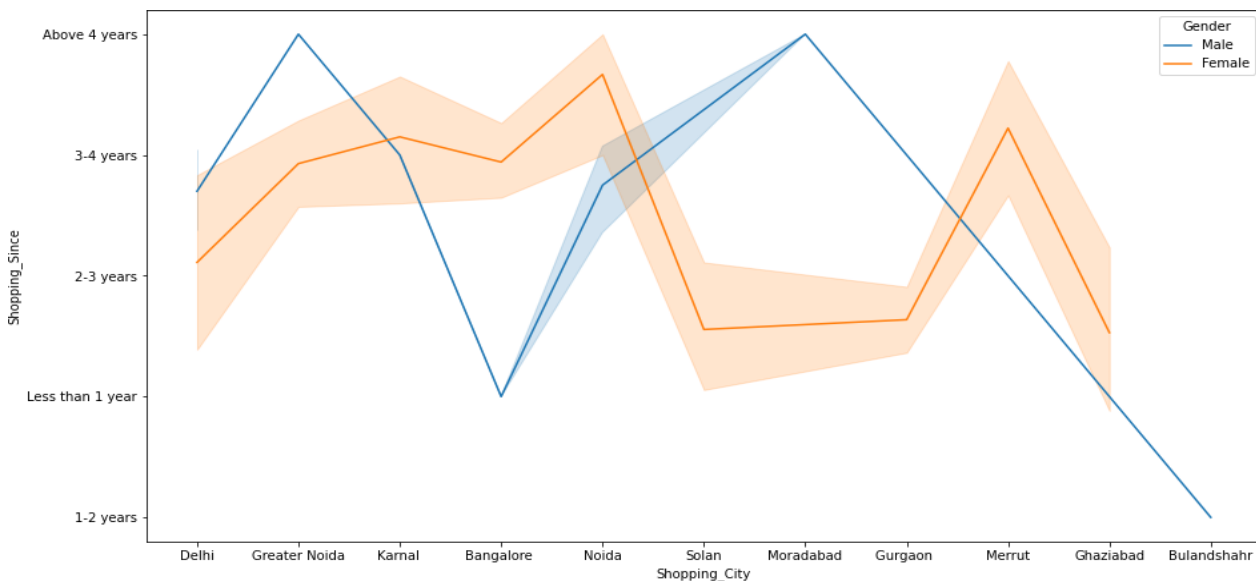
Even though Login time is more for Amazon, It is the most recommended.

h) Payment Option Availability Vs Fast Delivery



- Best Payment option and faster delivery is Amazon and Flipkart by most of the shoppers.

i) Shopping City Vs Gender

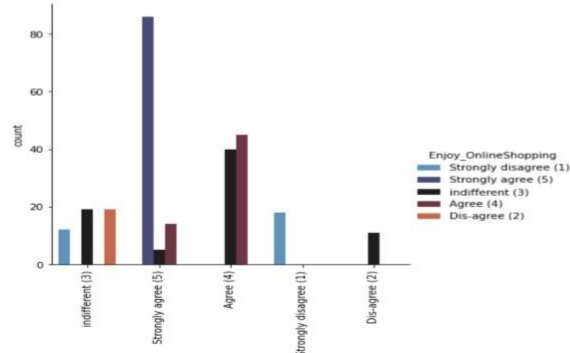


- We can see that the density of female customers is more than male.
- Men living in Bangalore and Ghaziabad have shopped online for less than 1 year.
- More Number of men shopping online from greater Noida and Moradabad are using online portals for more than 4 years.
- Women from Meerut and Noida have shopped the longest.

j) Benefit Discount Vs Enjoy Online shopping

```
In [35]: df.groupby("Benefit_Discount")["Enjoy_OnlineShopping"].value_counts()
Out[35]: Benefit_Discount    Enjoy_OnlineShopping
Agree (4)                  Agree (4)          45
                        Indifferent (3)       40
Dis-agree (2)              Indifferent (3)     11
Strongly agree (5)         Strongly agree (5)  86
                        Agree (4)           14
                        Indifferent (3)      5
Strongly disagree (1)      Strongly disagree (1) 18
Indifferent (3)            Dis-agree (2)      19
                        Indifferent (3)      19
                        Strongly disagree (1) 12
Name: Enjoy_OnlineShopping, dtype: int64

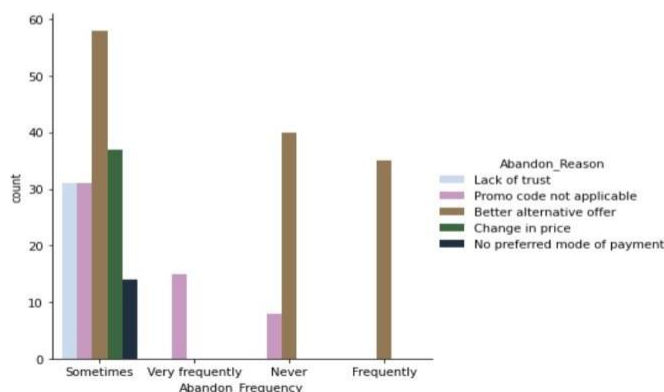
In [36]: sns.factorplot('Benefit_Discount', kind='count', data=df, hue='Enjoy_OnlineShopping', palette="icefire")
plt.xticks(rotation=90);
```



k) Abandon Frequency Vs Abandon Reason

```
In [37]: df.groupby("Abandon_Frequency")["Abandon_Reason"].value_counts()
Out[37]: Abandon_Frequency    Abandon_Reason
Frequently                Better alternative offer  35
Never                    Better alternative offer  40
                        Promo code not applicable  8
Sometimes                Better alternative offer  58
                        Change in price           37
                        Lack of trust             31
                        Promo code not applicable  31
                        No preferred mode of payment 14
Very frequently          Promo code not applicable 15
Name: Abandon_Reason, dtype: int64

In [38]: #Factor plot for Abandon_Frequency
sns.catplot('Abandon_Frequency', kind='count', data=df, hue='Abandon_Reason', palette="cubehelix_r")
Out[38]: <seaborn.axisgrid.FacetGrid at 0x7f873d752880>
```



- Most people abandon cart due to better alternates.
- Very frequently has been marked by few shoppers and the reason is Promo code not applicable.

4. Feature engineering

4.1 Encoding

After visualization, I have done the encoding of Data using Ordinal and Label Encoder.

```
encode=OrdinalEncoder()  
lencode=LabelEncoder()  
  
for i in cat:  
    train_df[i]=encode.fit_transform(train_df[i].values.reshape(-1,1))  
  
target_df=lencode.fit_transform(target_df)
```

```
52]: train_df.head()
```

```
52]:
```

	Gender	Age	Shopping_City	Shopping_Since	Shopping_Frequency	Internet_Access	Device_Used	Screen_Size	Operating_System	Browser_Used	Channel_F
0	1.0	1.0	2.0	3.0	2.0	0.0	0.0	3.0	2.0	0.0	
1	0.0	0.0	2.0	3.0	3.0	2.0	2.0	0.0	1.0	0.0	
2	0.0	0.0	4.0	2.0	3.0	1.0	2.0	2.0	0.0	0.0	
3	1.0	0.0	6.0	2.0	4.0	1.0	2.0	2.0	1.0	3.0	
4	0.0	0.0	0.0	1.0	0.0	2.0	2.0	0.0	1.0	3.0	

4.2 Correlation

- There is multicollinearity in the dataset
- Correlation is done for Training variables only.

4.3 Scaling the Data

Using MinMaxScaler I have scaled the data.

```
: xd=scaler.fit_transform(train_df)  
x=pd.DataFrame(xd,columns=train_df.columns)
```

```
: x.head()
```

```
:
```

	Gender	Age	Shopping_City	Shopping_Since	Shopping_Frequency	Internet_Access	Device_Used	Screen_Size	Operating_System	Browser_Used	Channel_F
0	1.0	0.25	0.2	0.75	0.50	0.0	0.000000	1.000000	1.0	0.0	
1	0.0	0.00	0.2	0.75	0.75	1.0	0.666667	0.000000	0.5	0.0	
2	0.0	0.00	0.4	0.50	0.75	0.5	0.666667	0.666667	0.0	0.0	
3	1.0	0.00	0.6	0.50	1.00	0.5	0.666667	0.666667	0.5	1.0	
4	0.0	0.00	0.0	0.25	0.00	1.0	0.666667	0.000000	0.5	1.0	

5. Model Building

5.1 HyperParameter Tuning

```
params={'n_estimators':[100, 300, 500, 700],
        'min_samples_split':[1,2,3,4],
        'min_samples_leaf':[1,2,3,4],
        'max_depth':[None,1,2,3,4,5,6,7,8,9,10,15,20,25,30,35,40]}
g=RandomizedSearchCV(RandomForestClassifier(),params,cv=10)
g.fit(xtrain,ytrain)
print("Best Estimator:",g.best_estimator_)
print("Best Parameter",g.best_params_)
print("Best Score",g.best_score_)
```

```
Best Estimator: RandomForestClassifier(max_depth=20, min_samples_leaf=3, min_samples_split=4)
Best Parameter {'n_estimators': 100, 'min_samples_split': 4, 'min_samples_leaf': 3, 'max_depth': 20}
Best Score 1.0
```

Using the best Parameters, I have built the Final Model, Found the Accuracy Score, Cross validation Score and Built the confusion Matrix.

```
Final_Model=RandomForestClassifier(max_depth=20, min_samples_leaf=3,min_samples_split=4,n_estimators=100)
Final_Model.fit(xtrain,ytrain)
pred=Final_Model.predict(xtest)
acc= accuracy_score(pred,ytest)
score=cross_val_score(Final_Model,x,y,cv=10)
print('The Accuracy score is:', acc*100)
print('The cross validation score', (score.mean()*100))
```

```
The Accuracy score is: 100.0
The cross validation score 100.0
```

```
print('Confusion Matrix')
print(confusion_matrix(pred,ytest))
```

```
Confusion Matrix
[[25  0  0  0  0  0  0  0  0]
 [ 0 21  0  0  0  0  0  0  0]
 [ 0  0  4  0  0  0  0  0  0]
 [ 0  0  0  6  0  0  0  0  0]
 [ 0  0  0  0  4  0  0  0  0]
 [ 0  0  0  0  0  5  0  0  0]
 [ 0  0  0  0  0  0 12  0  0]
 [ 0  0  0  0  0  0  0  4  0]]
```

- Got the Accuracy score of 100% and CV Score of 100% which is the best Score.

5.2 Saving the best Model

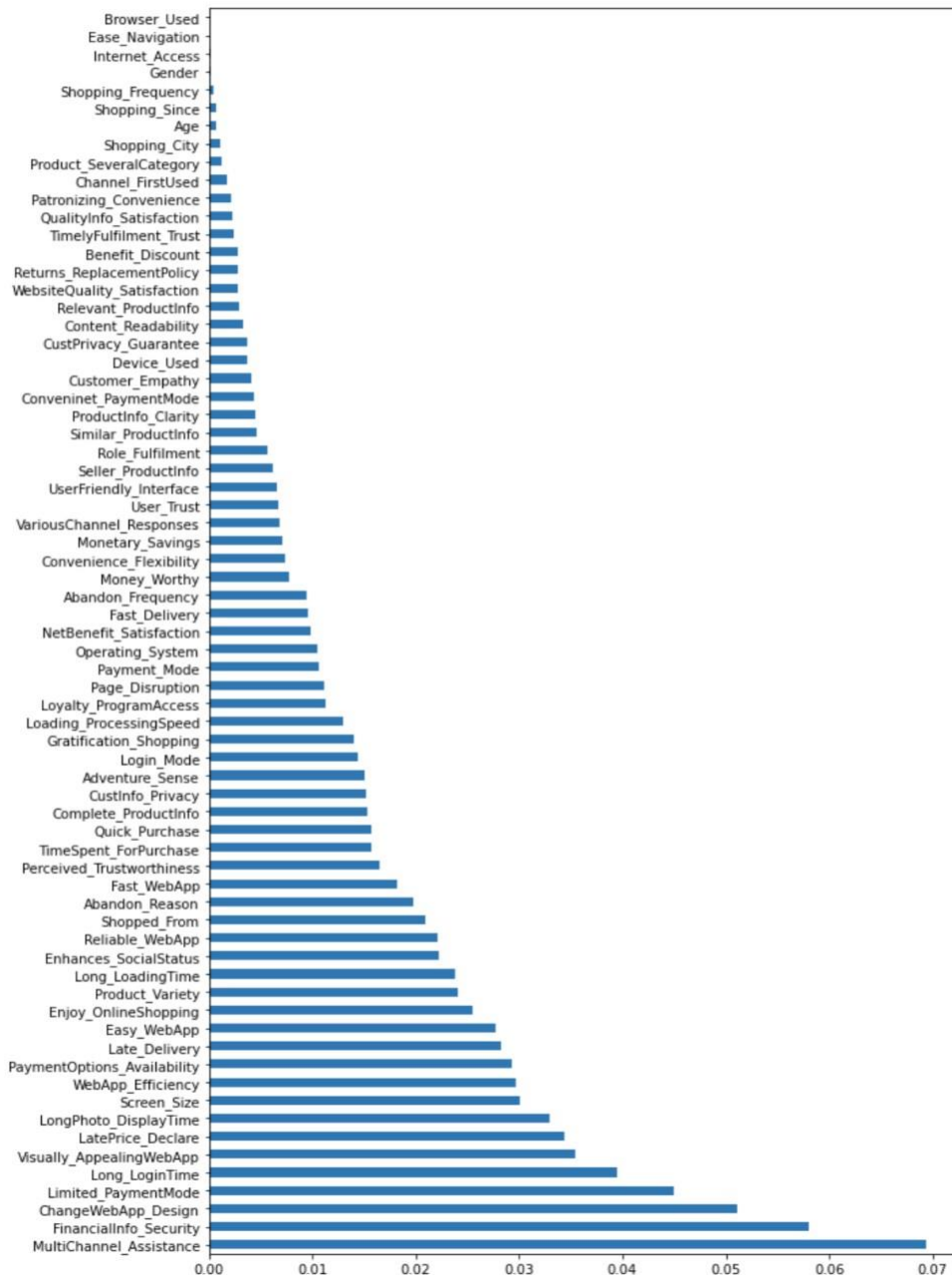
Using Joblib, I have saved the Final Model.

6. Feature Importance

Using the Final Model I have plotted Feature importance barplot and below is what it look like:

```
feat_importances = pd.Series(Final_Model.feature_importances_, index=a.columns)
feat_importances.nlargest(70).plot(kind='barh', figsize=(10,18))
```

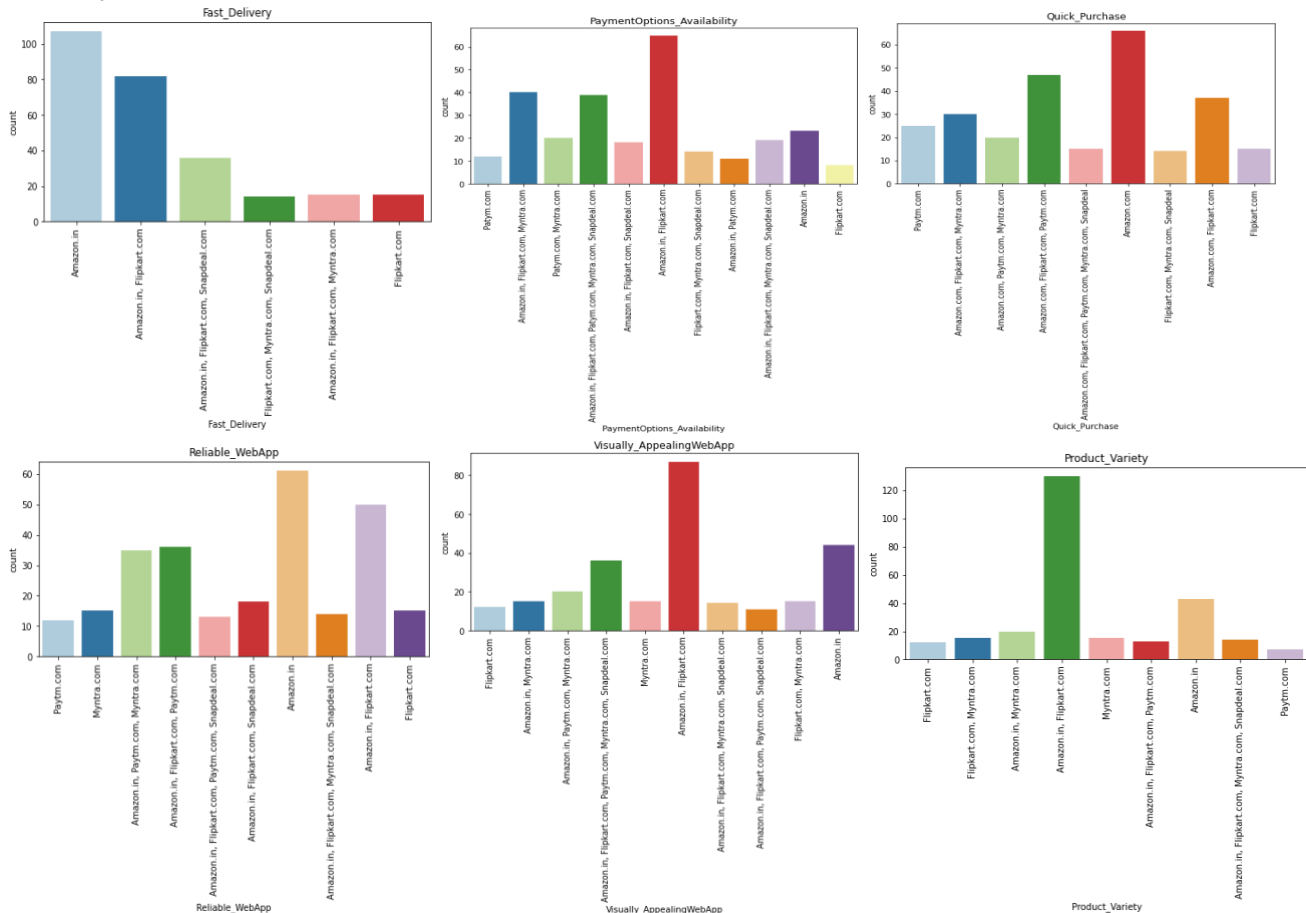
Out[82]: <AxesSubplot:>



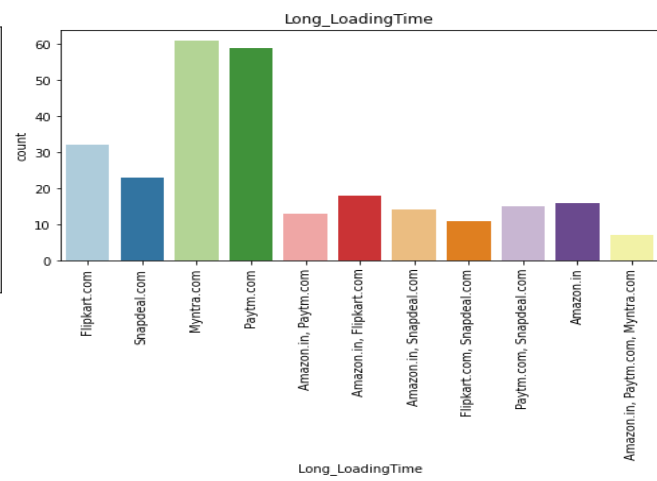
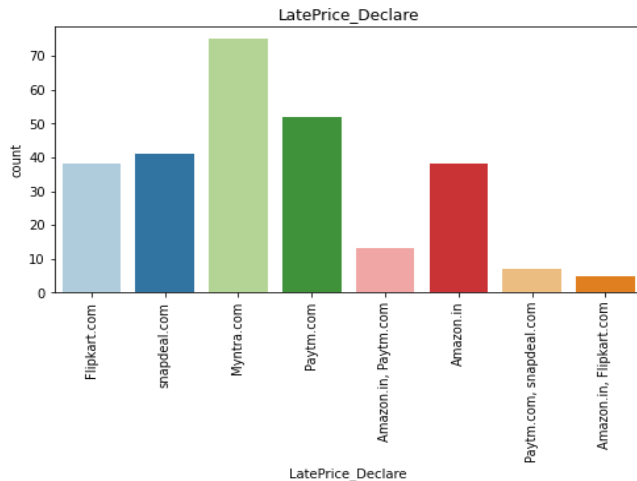
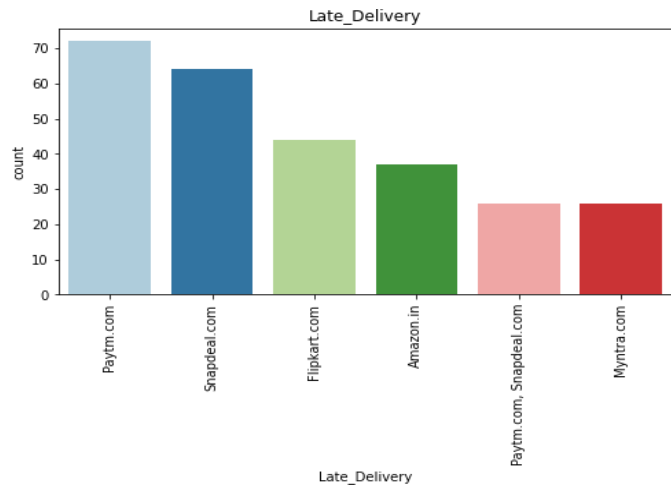
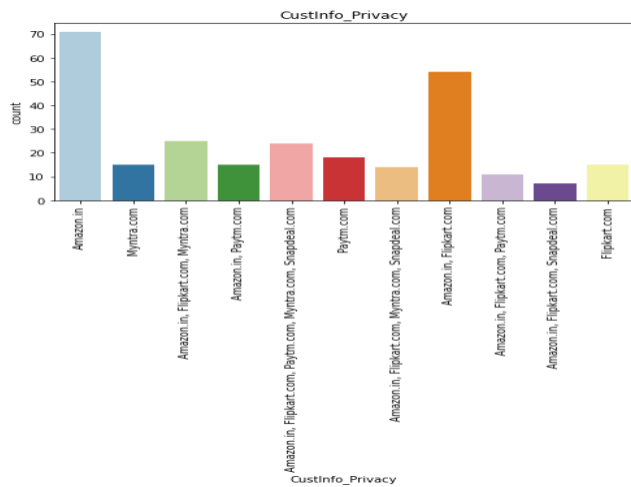
7. Conclusion

The results of this study suggest following outputs which might be useful for E-commerce websites to extend their business:

- **Females are more prone to shopping** and so more feminine related products attract more females and hence improves the Customer retention.
- The Financial Security, Perceived trustworthiness, The reliability of the Ecommerce website, all play an equally important role in deciding the buying behavior of online customers.
- The shoppers want to be sure that it will be possible to return the product if he does not like it in real life.
- The logistics factor, which includes Cash on delivery option, One day delivery and Descriptive factors like the Product information and Loading factors like Long Loading Time, Price declaration etc plays a secondary role in this process though these are Must-be-quality.
- All the websites were not equally preferred by online customers.
- **Amazon was the most preferred followed by Flipkart.** This can be explained easily by the previous result that we got. These two companies are most trusted in the industry and hence, have a huge reliability. These websites have the most lenient return policies as compared to others and also the time required to process a return is low for these.



- There is high risk of customer churn with:
 - Myntra.com
 - Snapdeal.com
 - Paytm.com



- The reasons why Customer retention is low for Myntra, Snapdeal and Paytm can be clearly inferred from above plots.
 - Low Customer Information Privacy
 - Late Delivery time
 - Long Loading time than other websites.
 - Price declaration is very late during Sale and promotions.
 - Low Payment Option
- It is crucial for E-commerce to consider their customer satisfaction because this will retain customer loyalty as well as attract potential customers.
- To conclude, Customer's Trust on the Company is the most important factor for Customer Retention. Factors like, Return Policy, Refund Policy, Fast Delivery, Wide Payment Channels, etc. helps the high-risk E Commerce websites to improve their Customer Retention Score.

1.

