

# **E-retail factors for customer activation and retention: A case study from Indian e-commerce customers**

Submitted by:

Saranya.M

## **INTRODUCTION**

Customer satisfaction has emerged as one of the most important factors that guarantee the success of online store; 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.

### **Conceptual Background of the Domain Problem**

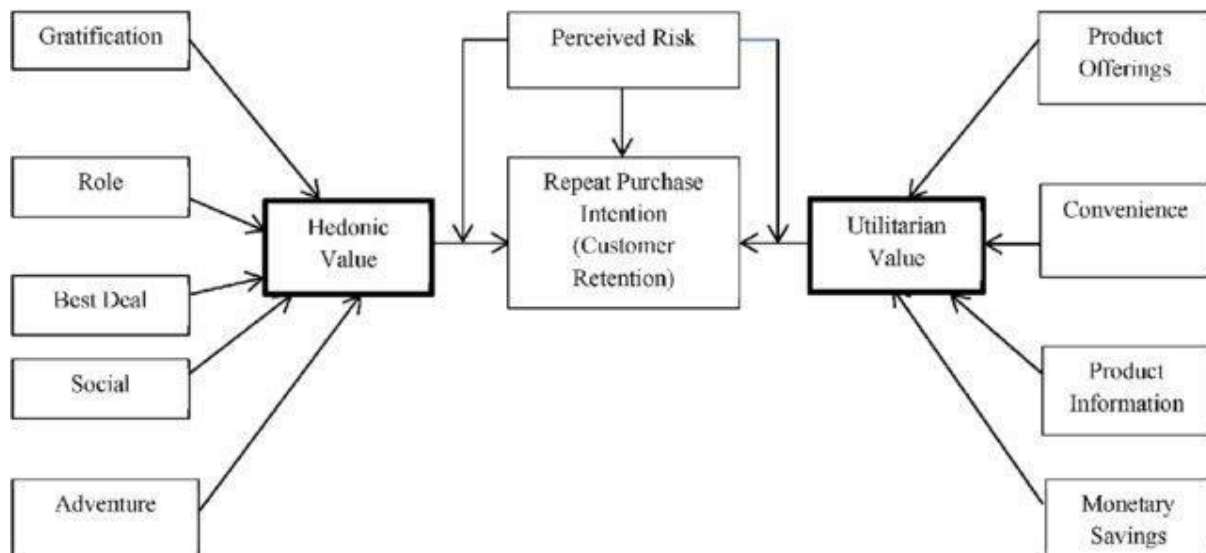
The data is collected from the Indian online shoppers. Results indicate the e-retail success factors, which are very much critical for customer satisfaction.

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.

## Motivation for the Problem Undertaken

Our main objective of doing this project is to analyse whether the users are shopping products from e-commerce websites, how did they give feedbacks to these websites on the basis of several positive and negative factors and also the details of the users on basis of factors like age, gender, etc.

## Diagrammatic Representation of Customer Retention



The Hedonic value consists of factors like Gratification, Role, Best Deal, Social and Adventure.

The Utilitarian value consists of factors like Product Offerings, Convenience, Product Information and Monetary Savings.

Customer Retention is based on 3 factors, according to the above diagram. They are:

Perceived Risk, Hedonic value and Utilitarian value

## Data Sources and their formats

The data is been given by a highly-confidential company and they gave it to us in an excel file. They also had provided the problem statement by explaining what they need from us and also the required criteria to be satisfied.

Let's check the data now. Below I have attached the snapshot below to give an overview.

```
In [18]: 1 df=pd.read_excel('F://flip robo intership31//Customer_retention_dataset -06-10-2022//customer_retention_dataset.xlsx')
2 df
```

Out[18]:

	1 Gender of respondent	2 How old are you?	3 Which city do you shop online from?	4 What is the Pin Code of where you shop online from?	5 Since How Long You are Shopping Online ?	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?	8 Which device do you use to access the online shopping?	9 What is the screen size of your mobile device?	10 What is the operating system (OS) of your device?	Longer time to get logged in (promotion, sales period)	Longer time in displaying graphics and photos (promotion, sales period)	Late declaration of price (promotion, sales period)
0	Male	31-40 years	Delhi	110009	Above 4 years	31-40 times	Dial-up	Desktop	Others	Window/windows Mobile	Amazon.in	Amazon.in	Flipkart.com
1	Female	21-30 years	Delhi	110030	Above 4 years	41 times and above	Wi-Fi	Smartphone	4.7 inches	IOS/Mac	Amazon.in, Flipkart.com	Myntra.com	snapdeal.com
2	Female	21-30 years	Greater Noida	201308	3-4 years	41 times and above	Mobile Internet	Smartphone	5.5 inches	Android	Myntra.com	Myntra.com	Myntra.com
3	Male	21-30 years	Karnal	132001	3-4 years	Less than 10 times	Mobile Internet	Smartphone	5.5 inches	IOS/Mac	Snapdeal.com	Myntra.com, Snapdeal.com	Myntra.com
4	Female	21-30 years	Bangalore	530068	2-3 years	11-20 times	Wi-Fi	Smartphone	4.7 inches	IOS/Mac	Flipkart.com, Paytm.com	Paytm.com	Paytm.com

-> There are totally 269 rows and 71 columns in this dataset

-> Our objective is to find the insights of the data and to do thorough data analysis.

# Hardware and Software Requirements and Tools Used

For doing this project, the hardware used is a laptop with high end specification and a stable internet connection. While coming to software part, I had used anaconda navigator and in that I have used **Jupyter notebook** to do my python programming and analysis.

For using an excel file, Microsoft excel is needed. In Jupyter notebook, I had used lots of python libraries to carry out this project and I have mentioned below with proper justification:

1. Pandas- a library which is used to read the data, visualisation and analysis of data.
2. NumPy- used for working with array and various mathematical techniques.
3. Seaborn- visualization tool for plotting different types of plot.
4. Matplotlib- It provides an object-oriented API for embedding plots into applications.

## Data Analysis

```
In [23]: 1 df.isnull().sum().any()
```

```
Out[23]: False
```

```
In [24]: 1 df.nunique()
```

```
Out[24]: Gender of respondent          2
How old are you?                     5
Which city do you shop online from?  11
What is the Pin Code of where you shop online from?  39
Since How Long You are Shopping Online ?  5
..
Longer delivery period                6
Change in website/Application design  7
Frequent disruption when moving from one page to another  8
Website is as efficient as before     8
Which of the Indian online retailer would you recommend to a friend?  8
Length: 71, dtype: int64
```

There are no null values in this dataset and 70 columns are of object datatype and only 1 column is of int data type.

```

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 (ipykernel)
In [24]: Out[24]: Gender of respondent 2
How old are you? 5
Which city do you shop online from? 11
What is the Pin Code of where you shop online from? 39
Since How Long You are Shopping Online ? 5
..
Longer delivery period 6
Change in website/Application design 7
Frequent disruption when moving from one page to another 8
Website is as efficient as before 8
Which of the Indian online retailer would you recommend to a friend? 8
Length: 71, dtype: int64

In [25]: 1 personal_info=['Gender of respondent','How old are you?','Which city do you shop online from?',
2 'What is the Pin Code of where you shop online from?','Since How Long You are Shopping Online ?',
3 'How many times you have made an online purchase in the past year?']

In [26]: 1 for i in personal_info:
2 if i!='What is the Pin Code of where you shop online from?':
3 plt.figure(figsize=(8,6))
4 df[i].value_counts().plot.pie(autopct='%1.1f%%')
5 centre=plt.Circle((0,0),0.7,fc='white')
6 fig=plt.gcf()
7 fig.gca().add_artist(centre)
8 plt.xlabel(i)
9 plt.ylabel('')
10 plt.figure()

```

We checked the value counts of all 71 columns above and we iterated using a for loop. We can see some value counts of the columns like gender, age, city, etc. Below I had attached the value counts of other columns.

```

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 (ipykernel)
In [22]: 1 df.dtypes
Out[22]: Gender of respondent object
How old are you? object
Which city do you shop online from? object
What is the Pin Code of where you shop online from? int64
Since How Long You are Shopping Online ? object
...
Longer delivery period object
Change in website/Application design object
Frequent disruption when moving from one page to another object
Website is as efficient as before object
Which of the Indian online retailer would you recommend to a friend? object
Length: 71, dtype: object

In [23]: 1 df.isnull().sum().any()
Out[23]: False

In [24]: 1 df.nunique()
Out[24]: Gender of respondent 2
How old are you? 5
Which city do you shop online from? 11
What is the Pin Code of where you shop online from? 39
Since How Long You are Shopping Online ? 5
..
Longer delivery period 6
Change in website/Application design 7

```

# Analysis of website feedbacks obtained

We can see that after column 47, there are both positive and negative feedbacks of the websites, which are given by the correspondents. We will analyse those data by using data analysis process.

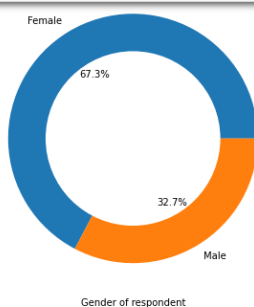
```
In [58]: 1 pca=PCA(n_components=29)
2 x=pca.fit_transform(x)
3 x=pd.DataFrame(x)
4 x.head()
```

```
Out[58]:
```

	0	1	2	3	4	5	6	7	8	9	...	19	20	21	22	
0	2.065419	-0.577759	-1.030081	-1.109784	0.652387	-1.137025	0.699876	-0.023177	-0.960103	-0.238855	...	0.598931	0.068875	-0.266070	-0.009322	-0.1
1	0.048667	-1.490547	1.081348	0.641617	0.066388	-0.820495	0.072214	-0.644870	0.087754	-0.296247	...	-0.176390	-0.008384	0.155024	0.313679	0.0
2	1.671684	-0.120022	0.775570	-1.481374	0.128287	0.836151	-0.793600	0.102789	0.448813	-0.515949	...	0.038239	-0.068419	0.008284	0.215976	-0.0
3	-0.009522	2.146296	0.753236	-0.363176	-1.348954	-0.176575	0.567430	-0.548924	-0.142604	-0.084665	...	0.025910	0.229481	-0.091051	0.190278	-0.0
4	0.051352	-0.187387	2.386865	0.914150	0.273219	-0.992250	-0.511792	0.701105	-0.225943	0.735107	...	-0.032298	0.130742	-0.195750	-0.163709	-0.0

5 rows x 29 columns

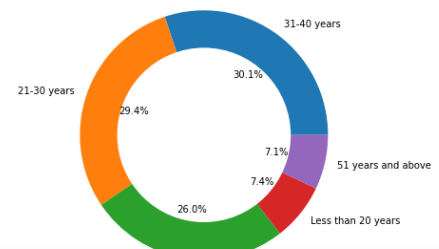
```
5 centre=plt.gca().add_artist(centre)
6 fig=plt.gcf()
7 fig.gca().add_artist(centre)
8 plt.xlabel(i)
9 plt.ylabel('')
10 plt.figure()
```



```
In [37]: 4 df['How many times you have made an online purchase in the
```

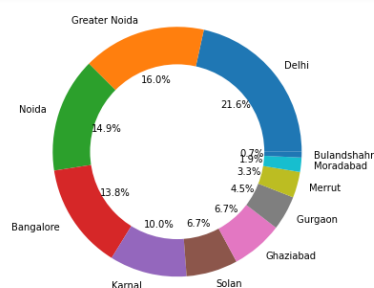
```
7 fig.gca().add_artist(centre)
8 plt.xlabel(i)
9 plt.ylabel('')
10 plt.figure()
```

<Figure size 432x288 with 0 Axes>

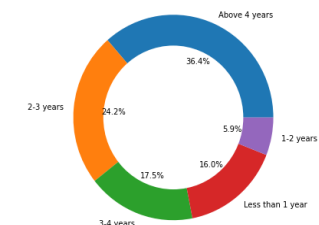


```
In [27]: 1 df['How many times you have made an online purchase in the
```

```
5 centre=plt.gca().add_artist(centre)
6 fig=plt.gcf()
7 fig.gca().add_artist(centre)
8 plt.xlabel(i)
9 plt.ylabel('')
10 plt.figure()
```



```
9 plt.ylabel('')
10 plt.figure()
```



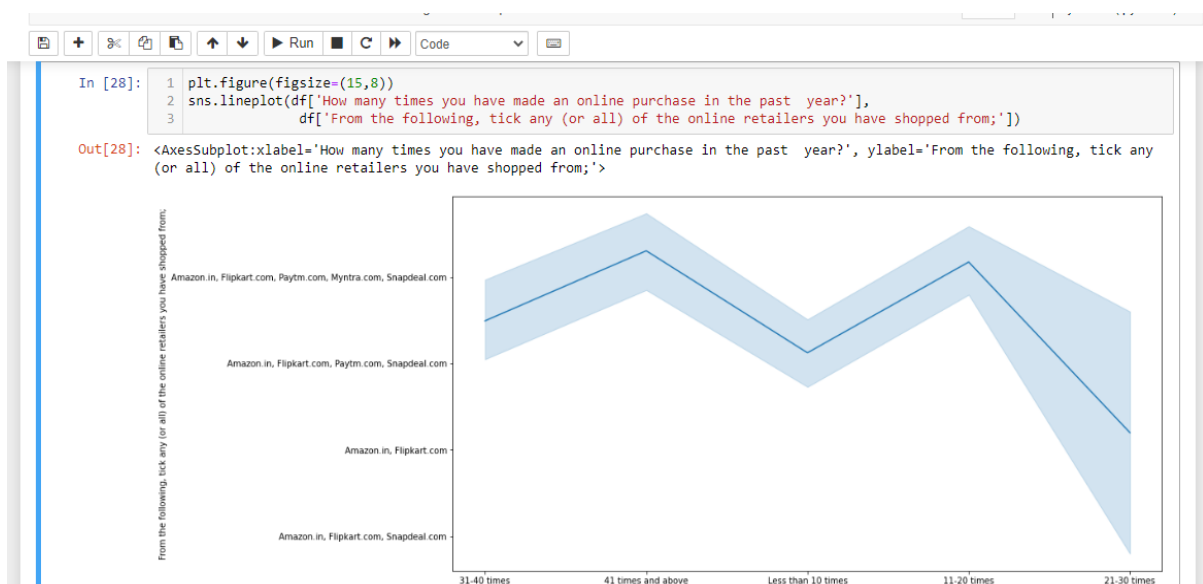
```
In [27]: 1 df['How many times you have made an online purchase in the past year?']
2
```

ipe here to search

There is double the number of women than men who have taken this survey. -Most of the people are in their 30's followed by 20's, teenagers and senior citizen are the least in number. -Most of the people belong from delhi, noida and banglore, ambiguity can also be seen as noida has two categories (noida and grater noida) which need to be handled -Most of the people shopping online have been shopping from a long time. -Majority of people shop online 10 times a year, amiguity can also be seen for range 42 times and above which needs to be handled

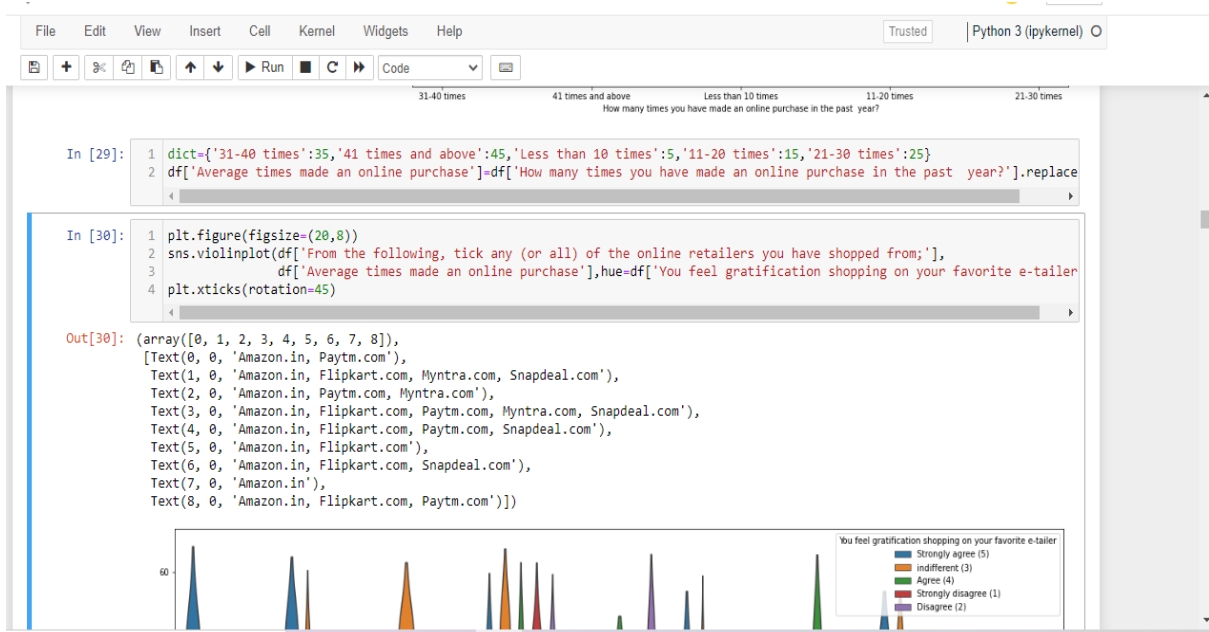
## #Analysis on the basis of Various following factors

### \*Intention of Repeat purchase:



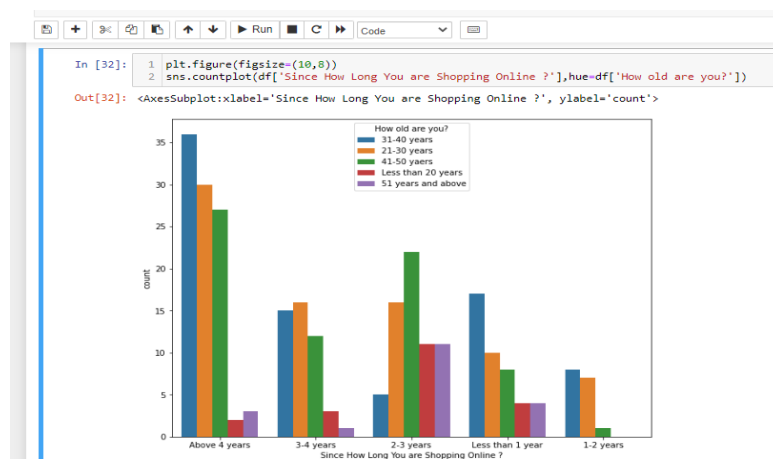
Now, we Heavy shoppers who shop more than 41 times a year shop from all the online brands, some of the people who shop for 32-40 and less than 10 times a year seem to exclude myntra. People shop from Amazon and flipkart whatever be the case.



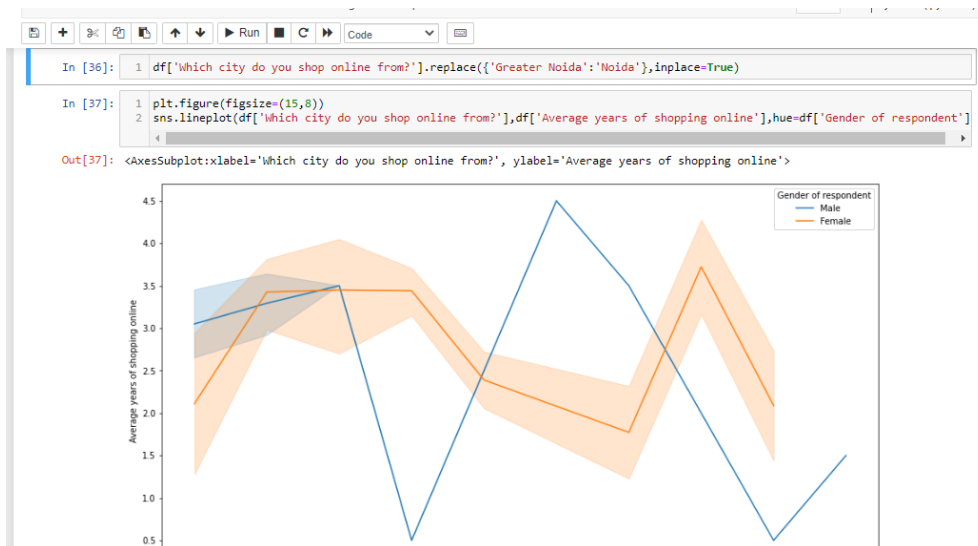


Almost all the people who have shopped from amazon, flipkart and paytm are satisfied. People who shop from a more number of online brands dosent seem to be satisfied.

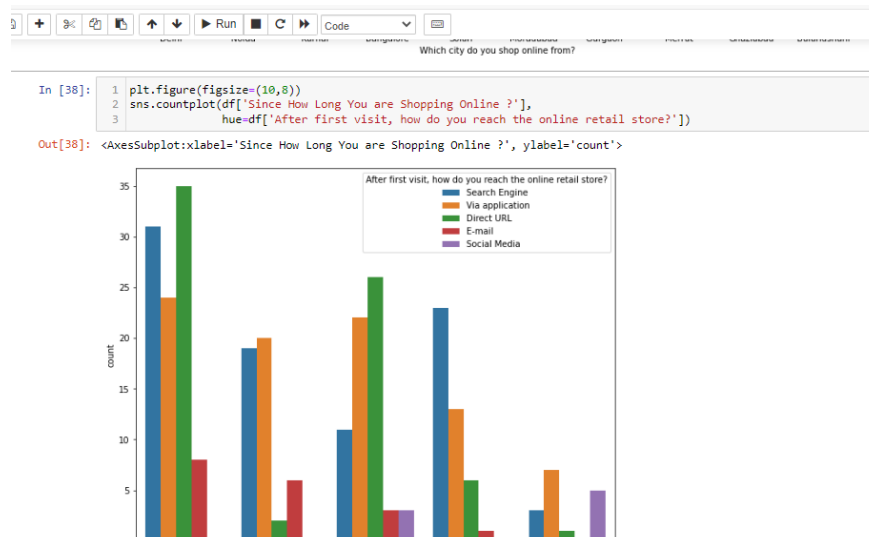
### \*Online Retailing:



Highest number of people have been shopping online for above 4 years except for the age group below 20 years and above 50 years. People who are shopping online for 1-2 years does not include teenagers and elder people.

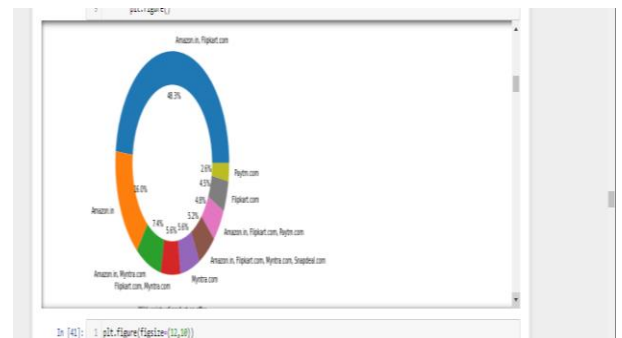
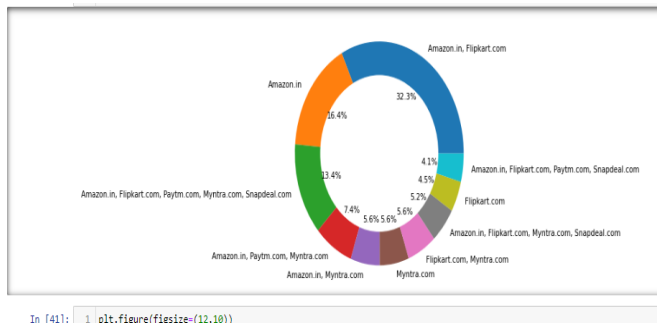
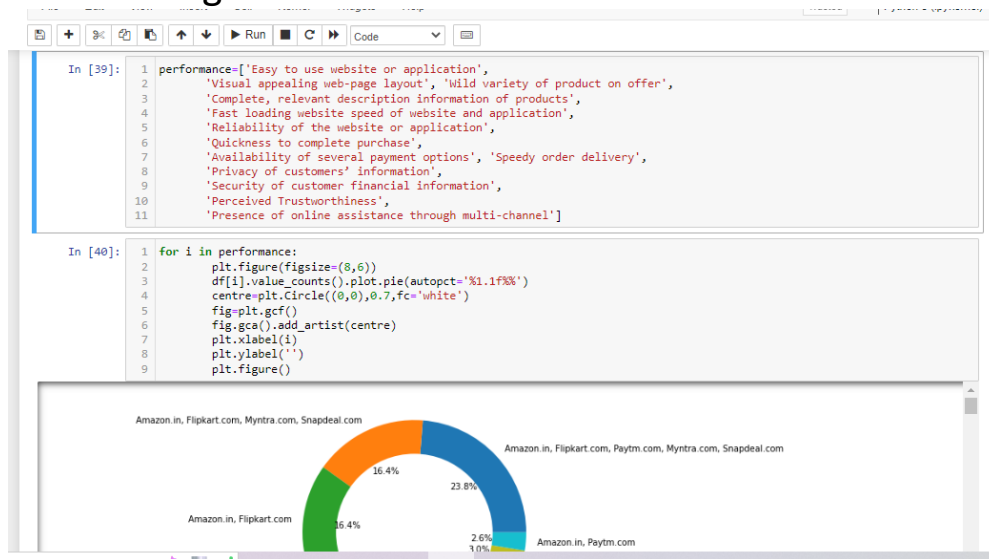


In lines, we can see that density of female customers is more than male. Men living in banglore and ghaziabad shop have shopped online for less than 1 year. Highest number of men shopping online belong from delhi and noida, while men from moradabad have been shopping online for the longest. Women from meerut and noida have shopped the longest.

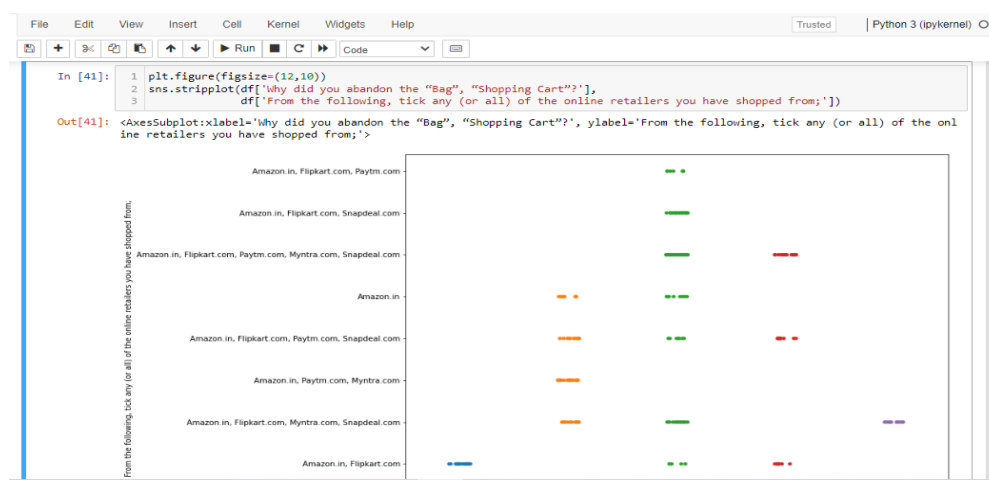


Even though people who are shopping online for more than 3 years donot use the application rather use search engine and direct url's in large number which indicates that online brands should update all their platforms rather than just application.

## #Brand image



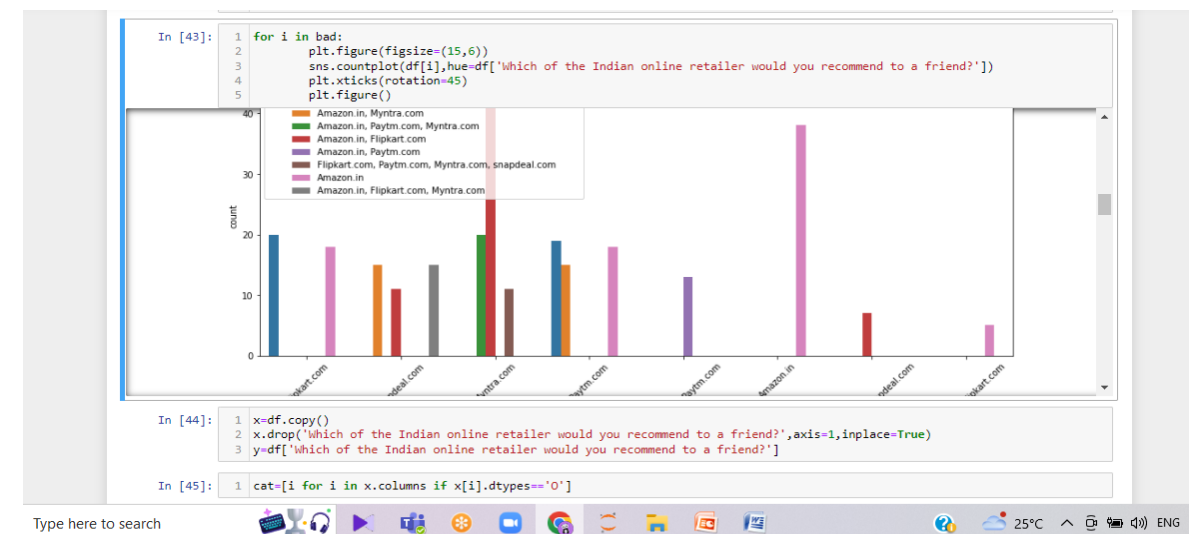
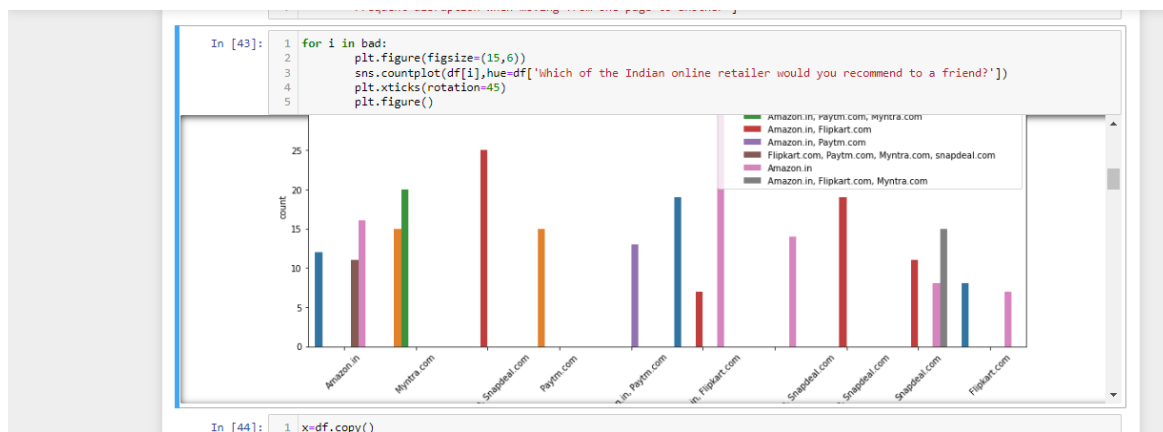
Amazon, Flipkart have been had the highest votes for having all the positive points and have maintained a very good brand image followed by paytm and the myntra.



We can clearly see that most of the time people abandon the bag is because they get a better alternative offer or promo code not applicable. There is also lack of trust seen in amazon, flipkart and paytm by some people.

## #Loyalty

Loyal customers are those who keep using the same brand even if it is not good as other brands



Customers seem to be more loyal to amazon, flipkart and paytm as even though many of them have given negative remarks about them still they would recommend these platforms to their friend

## #Processing the dataframe

```
In [44]: 1 x=df.copy()
2 x.drop('Which of the Indian online retailer would you recommend to a friend?',axis=1,inplace=True)
3 y=df['Which of the Indian online retailer would you recommend to a friend?']

In [45]: 1 cat=[i for i in x.columns if x[i].dtypes=='O']

In [46]: 1 from sklearn.preprocessing import OrdinalEncoder,LabelEncoder
2 encode=OrdinalEncoder()
3 labe=LabelEncoder()

In [47]: 1 for i in cat:
2     x[i]=encode.fit_transform(x[i].values.reshape(-1,1))
3     y=labe.fit_transform(y)

In [48]: 1 from sklearn.preprocessing import MinMaxScaler
2 scaler=MinMaxScaler()

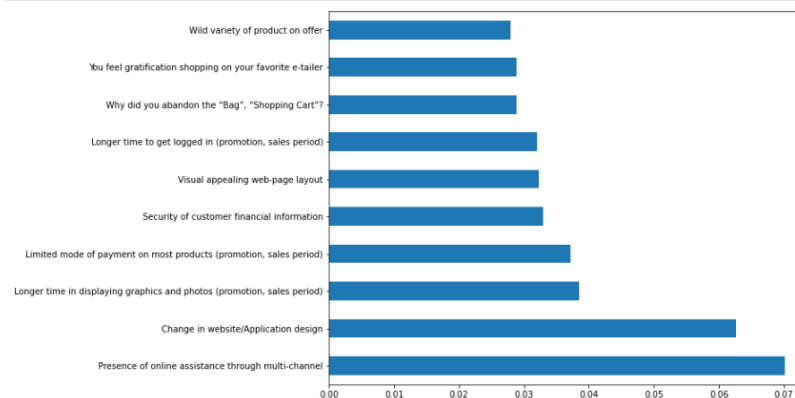
In [49]: 1 xd=scaler.fit_transform(x)
2 x=pd.DataFrame(xd,columns=x.columns)

In [50]: 1 from sklearn.ensemble import RandomForestClassifier
2 m=RandomForestClassifier()
3 m.fit(x,y)

Out[50]: RandomForestClassifier()

In [51]: 1 feat_importances = pd.Series(m.feature_importances_, index=x.columns)
2 plt.figure(figsize=(10,8))
3 feat_importances.nlargest(10).plot(kind='barh')
```

```
In [51]: 1 feat_importances = pd.Series(m.feature_importances_, index=x.columns)
2 plt.figure(figsize=(10,8))
3 feat_importances.nlargest(10).plot(kind='barh')
4 plt.show()
```



In the above chart we can see that above features are of most importance in determining which platform will a customer recommend to his friend.

## #Using chi2 test

```
In [52]: 1 from sklearn.feature_selection import SelectKBest
2 from sklearn.feature_selection import chi2

In [53]: 1 selection = SelectKBest(score_func=chi2)
2 fit = selection.fit(x,y)

In [54]: 1 dfscores = pd.DataFrame(fit.scores_)
2 dfcolumns = pd.DataFrame(x.columns)
3 featureScores = pd.concat([dfcolumns,dfscores],axis=1)
4 featureScores.columns = ['Features','Score']

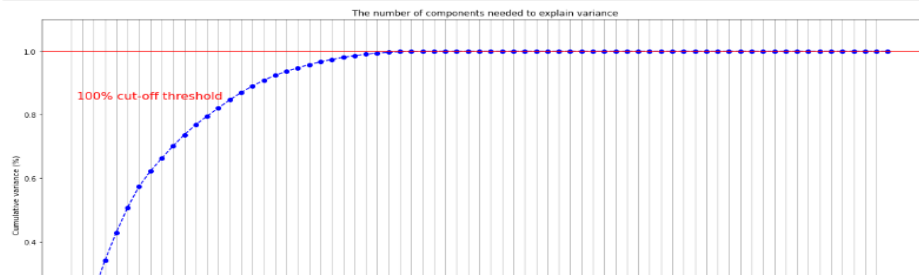
In [55]: 1 print(featureScores.nlargest(10,'Score'))
2 feat=list(featureScores.nlargest(10,'Score')['Features'])

16 Why did you abandon the "Bag", "Shopping Cart"? 75.754028
22 Loading and processing speed 59.810983
42 Shopping on the website gives you the sense of... 59.253509
10 What browser do you run on your device to acce... 57.171099
67 Change in website/Application design 55.301526
49 Visual appealing web-page layout 54.245760
65 Limited mode of payment on most products (prom... 53.269266
61 Longer time to get logged in (promotion, sales... 48.222655
62 Longer time in displaying graphics and photos ... 48.130643
50 Wild variety of product on offer 47.605973

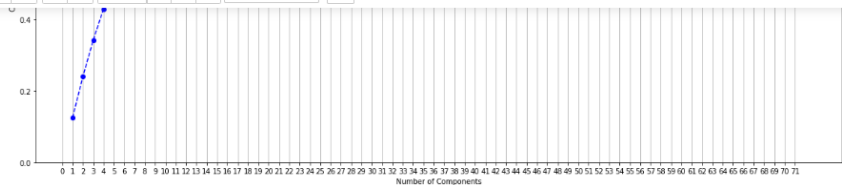
In [56]: 1 from sklearn.decomposition import PCA
2 pca = PCA().fit(x)
```

## #PCA

```
In [57]: 1 fig, ax = plt.subplots(figsize=(20,10))
2         xi = np.arange(1, 73, step=1)
3         yi = np.cumsum(pca.explained_variance_ratio_)
4
5         plt.ylim(0.0,1.1)
6         plt.plot(xi, yi, marker='o', linestyle='--', color='b')
7
8         plt.xlabel('Number of Components')
9         plt.xticks(np.arange(0, 72, step=1)) #change from 0-based array index to 1-based human-readable Label
10        plt.ylabel('Cumulative variance (%)')
11        plt.title('The number of components needed to explain variance')
12
13        plt.axhline(y=1, color='r', linestyle='-')
14        plt.text(0.5, 0.85, '100% cut-off threshold', color = 'red', fontsize=16)
15
16        ax.grid(axis='x')
17        plt.show()
```



```
File Edit View Insert Cell Format Windows Help
+ - * < > Run Code
```



```
In [58]: 1 pca=PCA(n_components=29)
2         x=pca.fit_transform(x)
3         x=pd.DataFrame(x)
4         x.head()
```

```
Out[58]:
```

	0	1	2	3	4	5	6	7	8	9	...	19	20	21	22	
0	2.065419	-0.577759	-1.030081	-1.109784	0.652387	-1.137025	0.699876	-0.023177	-0.960103	-0.238855	...	0.598931	0.068875	-0.266070	-0.009322	-0.1
1	0.048667	-1.490547	1.081348	0.641617	0.066388	-0.820495	0.072214	-0.644870	0.087754	-0.296247	...	-0.176390	-0.008384	0.155024	0.313679	0.0
2	1.671684	-0.120022	0.775570	-1.481374	0.128287	0.836151	-0.793600	0.102789	0.448813	-0.515949	...	0.038239	-0.068419	0.008284	0.215976	-0.0
3	-0.009522	2.146296	0.753236	-0.363176	-1.348954	-0.176575	0.567430	-0.548924	-0.142604	-0.084665	...	0.025910	0.229481	-0.091051	0.190278	-0.0
4	0.051352	-0.187387	2.386865	0.914150	0.273219	-0.982250	-0.511792	0.701105	-0.225943	0.735107	...	-0.032298	0.130742	-0.195750	-0.163709	-0.0

5 rows x 29 columns

We can clearly see that with 29 features all the information can be retained

## #Modeling Phase

```
5 rows x 29 columns
```

```
In [61]: 1 from sklearn.model_selection import train_test_split, cross_val_score
2
3         from sklearn.ensemble import RandomForestClassifier
4
5         from sklearn.metrics import accuracy_score, confusion_matrix, classification_report, roc_auc_score, roc_curve

In [62]: 1 xtrain, xtest, ytrain, ytest = train_test_split(x, y, test_size=0.3, random_state=7)

In [63]: 1 model = RandomForestClassifier()
2         model.fit(xtrain, ytrain)
3         p = model.predict(xtest)
4         s = cross_val_score(model, x, y, cv=10)

In [64]: 1 print('Accuracy', np.round(accuracy_score(p, ytest), 4))
2         print('-----')
3         print('Mean of Cross Validation Score', np.round(s.mean(), 4))
4         print('-----')
5         print('Confusion Matrix')
6         print(confusion_matrix(p, ytest))
7         print('-----')
8         print('Classification Report')
9         print(classification_report(p, ytest))

Accuracy 1.0
```

```

7 print('-----')
8 print('Classification Report')
9 print(classification_report(p,ytest))

Accuracy 1.0
-----
Mean of Cross Validation Score 0.9963
-----
Confusion Matrix
[[26  0  0  0  0  0  0  0]
 [ 0 22  0  0  0  0  0  0]
 [ 0  0  4  0  0  0  0  0]
 [ 0  0  0  4  0  0  0  0]
 [ 0  0  0  0  5  0  0  0]
 [ 0  0  0  0  0  7  0  0]
 [ 0  0  0  0  0  0 11  0]
 [ 0  0  0  0  0  0  0 2]]
-----
Classification Report
              precision    recall  f1-score   support

0               1.00        1.00        1.00        26
1               1.00        1.00        1.00        22
2               1.00        1.00        1.00         4
3               1.00        1.00        1.00         4
4               1.00        1.00        1.00         5
5               1.00        1.00        1.00         7
6               1.00        1.00        1.00        11
7               1.00        1.00        1.00         2

 accuracy
macro avg          1.00        1.00        1.00        81
weighted avg        1.00        1.00        1.00        81

In [65]: 1 print('Accuracy',np.round(accuracy_score(p,ytest),4))
          2 print('-----')
          3 print('Mean of Cross Validation Score',np.round(s.mean(),4))
          4 print('-----')
```

## Hyperparameter Tuning

```

In [66]: 1 from sklearn.model_selection import RandomizedSearchCV

In [67]: 1 params={'n_estimators':[100, 300, 500, 700],
              2 'min_samples_split':[1,2,3,4],
              3 'min_samples_leaf':[1,2,3,4],
              4 'max_depth':[None,1,2,3,4,5,6,7,8,9,10,15,20,25,30,35,40]}

In [68]: 1 g=RandomizedSearchCV(RandomForestClassifier(),params,cv=10)

In [69]: 1 g.fit(xtrain,ytrain)
Out[69]: RandomizedSearchCV(cv=10, estimator=RandomForestClassifier(),
                           param_distributions={'max_depth': [None, 1, 2, 3, 4, 5, 6, 7,
                           8, 9, 10, 15, 20, 25, 30,
                           35, 40],
                           'min_samples_leaf': [1, 2, 3, 4],
                           'min_samples_split': [1, 2, 3, 4],
                           'n_estimators': [100, 300, 500, 700]})

In [70]: 1 print(g.best_estimator_)
          2 print(g.best_params_)
          3 print(g.best_score_)

RandomForestClassifier(max_depth=40, min_samples_leaf=2, min_samples_split=4,
                       n_estimators=700)
{'n_estimators': 700, 'min_samples_split': 4, 'min_samples_leaf': 2, 'max_depth': 40}
0.9947368421052631

In [71]: 1 m=RandomForestClassifier(max_depth=20, min_samples_leaf=4, min_samples_split=4,n_estimators=700)
          2 m.fit(xtrain,ytrain)
          3 p=m.predict(xtest)
          4 score=cross_val_score(m,x,y,cv=10)
```

## #Finalizing the best Model

### \*Evaluation Metrics

```

2 print('-----')
3 print('Mean of Cross Validation Score',np.round(s.mean(),4))
4 print('-----')
5 print('Confusion Matrix')
6 print(confusion_matrix(p,ytest))
7 print('-----')
8 print('Classification Report')
9 print(classification_report(p,ytest))

Accuracy 1.0
-----
Mean of Cross Validation Score 0.9963
-----
Confusion Matrix
[[26  0  0  0  0  0  0  0]
 [ 0 22  0  0  0  0  0  0]
 [ 0  0  4  0  0  0  0  0]
 [ 0  0  0  4  0  0  0  0]
 [ 0  0  0  0  5  0  0  0]
 [ 0  0  0  0  0  7  0  0]
 [ 0  0  0  0  0  0 11  0]
 [ 0  0  0  0  0  0  0 2]]
-----
Classification Report
              precision    recall  f1-score   support

0               1.00        1.00        1.00        26
1               1.00        1.00        1.00        22
2               1.00        1.00        1.00         4
3               1.00        1.00        1.00         4
4               1.00        1.00        1.00         5
5               1.00        1.00        1.00         7
6               1.00        1.00        1.00        11
7               1.00        1.00        1.00         2

 accuracy
macro avg          1.00        1.00        1.00        81
weighted avg        1.00        1.00        1.00        81
```

## Saving the Model

```
In [76]: 1 import joblib
          2 joblib.dump(model, 'Retention.obj')

Out[76]: ['Retention.obj']
```

## Conclusion:

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

The cost of the product, the reliability of the E-commerce company and the return policies all play an equally important role in deciding the buying behaviour of online customers. The cost is an important factor as it was the basic criteria used by online retailers to attract customers. The reliability of the E-commerce company is also important, as it is even required in offline retail. It is important because customers are paying online, so they need to be sure of security of the online transaction. The return policies are important because in online retail customer does not get to feel the product. Thus, he wants to be sure that it will be possible to return the product if he does not like it in real. Whereas, the logistics factor, which included Cash on delivery option, One day delivery and the quality of packaging plays a secondary role in this process though these are Must-be-quality. This is so because these all does not interfere with the real product and people believe that this is the basic value that E-commerce websites provide.

All the websites were not equally preferred by online customers. Amazon was the most preferred followed by Flipkart. This can be explained easily by previous result that we got. These two companies are most trusted in the industry and hence, have a huge reliability. Also, the sellers listed on these websites are generally from Tier 1 cities as compared to Snapdeal and PayTM which have more sellers from tier 2 and 3 cities. Also, 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.