



# IJRASET

International Journal For Research in  
Applied Science and Engineering Technology



---

# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

---

**Volume:** 10    **Issue:** VI    **Month of publication:** June 2022

**DOI:** <https://doi.org/10.22214/ijraset.2022.43693>

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

Call:  08813907089

E-mail ID: [ijraset@gmail.com](mailto:ijraset@gmail.com)

# Analyzing and Forecasting the Purchases Made on Black Friday

Dileep Kumar A S<sup>1</sup>, Ms. Sindhu D<sup>2</sup>, Dr. Ravikumar G K<sup>3</sup>

<sup>1</sup>Dept. of CSE, <sup>2</sup>Dept. of ISE, Professor & <sup>3</sup>Head(R&D)Dept. of CSE, BGS Institute of Technology Adichunchanagiri University, BG Nagar, Karnataka, India-571448.

**Abstract**— Things are sold at a substantial discount on the eve of Black Friday, resulting in sales that are 30 times larger than on normal flash sale days. Customers' data from purchases made on this day can be examined, resulting in a quick declaration of their preferences for specific products. We looked at data that contained packets of clients, and also the factors that influenced their purchases and the amounts they spent. This data is analysed and forecasted purely for the purpose of providing clients with customized discounts on goods depending on individual preferences and purchase budget. Four models were employed to forecast significant variations in training and test data (50:50, 70:30, 30:70), as well as a distinct sample training and testing dataset with two additional examples of prediction: xgboost, tfidftransform, both combination, and extra trees regressor. The two scenarios involve forecasting and analysing another dataset, as well as projected on the train data and testing data on a different testing data set. The dataset would be analyzed to learn about consumer behavior and trends of product the sale's popularity. For each of the five scenarios, the feature significance and benefit importance are displayed. All of the models' accuracy in various settings has been given in the manner of accuracy graphs and the accuracy findings have been displayed in the form of an RMSE score.

**Keywords**— Black Friday, XGBoost, Accuracy.

## I. INTRODUCTION

Various deals are held throughout the year to recognise the customer's enthusiasm in shopping. Black Friday are just a few of the notable sales. During these deals are sold both physically and internet at steep deals, and practically every e-commerce firms participate. Black Friday, commonly known as Thanksgiving Day, started in the US. Every year, on the fourth Thursday of November, is this sale held. This is the busiest purchasing day of the year, according to reports. This deal is taking place in the following countries: the USA, Canada, UK, India, and many others. With the exception of 2008, Black Friday is grown connectivity in every year. Then compared to the last year's trend, sales and attractiveness of this day grow by about 12% (about). The Monday following Thanksgiving is known as Cyber Monday (Black Friday). This sale was created solely to encourage people to shop online, and most things are sale at deep concession during this event. In 2017, Cyber Monday sales went up by \$6.59 billion, with 77 percent of online merchants admitting that the holiday had an impact on their sales. Flipkart and Amazon, respectively, organize the other two online deals Amazon Sales, Jio Mart. These two promotions are specifically targeted at Indian clients. These two promotions are specifically targeted at Indian clients. Since this is an article regarding digital shopping, the Black Friday bargain is included because it has the largest data which is the most recognized among the others. Another reason to celebrate Black Friday is that it should be acknowledged by web-based e-commerce websites and is observed globally.

Consumers had benefited from the e-commerce sector, according to research and has been a success in and of itself. It was made simple for a customer to purchase premium quality at a fair price. Vendor sales have also increased as a result of engagement with the e-commerce industry. However, in order to maintain popularity and reduce losses, Data science is now required across the board in this industry. Data science has demonstrated its efficacy in detecting fraud in sellers and customers, improving customer segmentation, and forecasting market pricing, among other things. As a result, contemporary data science methodologies required for the sector's continued to develop and to maintain its presence among competitors in the same industry as well as offline providers.

One of the most recent data science strategies is studying and projecting the history of a product purchased by a client and generating customized discounts for specific customers based on the prediction. Gender has an important impact in determining special offers. Online industries, such as Amazon and Snapdeal, Optimize their techniques on a daily basis to strengthen the customer-product interaction, and product purchasing event serves as a data that is analyzed and predicted in

order to increase the likelihood of customers purchasing the product while also provide them offers occasionally, thereby increasing sales and maintaining connectivity between them. The algorithm provides larger earnings on special days such as Daily deals, Holiday sales, Cyber Monday sales, Holi sales, Black Friday sales, Revolutionary Year sales, off-season sales, and Festival sales. We'll use the information of customers who went shopping on Black Friday's eve to accomplish this analysis. The dataset for the Black Friday sale was chosen because it takes place all over the world and is attended by the majority of e-commerce companies. As a result, the database will be vast, and the more data there is, the more accurate the forecast will be.

## II. RELATED WORK

Consumers' expectations throughout sales like Black Friday and Cyber Monday were explored by Esther et al. [1], who predicted why these sales had a large impact on gross profit. Fisher and his colleagues [2] presented many variables such as sex, gender, and consumer behaviour that influence marketing. Zhang [3] investigated & demonstrated different because of each individual's choice to shop online. They also discussed the high-quality products available for purchase through online shopping. Vijayasathy [4] has also contributed to the prediction of consumer intents to shop online. The Iris Flower and Titanic datasets were successfully studied and predicted before beginning the analyze and forecast of the consumer data on Black Friday. The 150 observations in the Iris dataset have the following attributes: length of sepal, width of sepal, length of petal, and width of petal. In the Titanic dataset, there are 150 observations with the following characteristics: people, survivor, pclass, name, age, fare, and cabin. Linda and her colleagues [5] have provided a comprehensive examination of Black Friday behaviour, highlighting factors such as career and marital status that influence sales on Black Friday. Bellizi et al. [6] have discussed how visualisation might influence customer psychology and marketing. They talked about how factors like the colour of the atmosphere might affect the customer's mood and likelihood of making a purchase. Donovan et al. [7] have carried out similar study. Margarete Sandelowski [8] focuses on research methods that integrate descriptive and analytical sampling, data collection, and analytic procedures in the analysis section. Burke [9] looked at what customers expect in a physical and online store. Thomas et al. [10] presented an exploratory study of black Friday consumption rituals in 2011. In this study about wolf habitat in the northern Great Lakes, Mladenoff [11] addresses research methodologies. Tashakkori et al. [12] In their research, they used a mixed methods approach, mixing qualitative and quantitative approaches. Ma X [13] has established a link between geoscience data and ogc standards in practise for data visualisation. In the field of political science, the International Encyclopedia of Political Science is a valuable resource. Brillinger [14] developed exploratory data analysis. Hammersley [15] explains how qualitative and quantitative data are related. For research methodologies in the psychology field, Hammersley [15] established a link between qualitative information.

## III. DESCRIPTION OF THE DATASET

Analytics Vidhya revealed the information throughout its Black Friday database in 2015. There were 550069 rows in the training dataset, 233600 rows in the testing dataset, and 233600 rows in the sample dataset. Although the test database simply provides standard purchases for each user with a quantity of 9000, it is supplied to ensure correctness. The training dataset includes parameters such as customer id, item id, gender, age, occupation, city category, stay in present location, relationship status, item category1, purchasing, item category2, and item category3. The test data has most of the similar attributes as the train dataset except for the 'purchase' attribute. The example dataset includes properties such as customer id, item id, and purchase. In the first three circumstances, just the training data will be used as an input dataset, which will be split into 90 percent 10 percent, 70 percent 30 percent, and 50 percent, with the former representing the updated training data and the latter representing the new test dataset [30]. To ensure correctness, a new dataset is prepared. Except that the 'purchase' feature will be eliminated for testing reasons, The sample dataset created from the present test data would have had the same characteristics as this dataset. Original training is the fourth scenario, where the algorithms are developed and implemented on a similar training sample. The fifth scenario uses Analytics Vidhya's original train, test, and sample dataset. Here, models developed with the training dataset are put to the test with the testing dataset.

## IV. RESEARCH METHODOLOGY

The quantitative approach in the data we're working with is subjected to exploratory analysis. Exploratory data analysis is a methodology for visually analysing data sets in order to highlight their important properties. Quantitative data includes numerical data such as iris datasets and scorecard data. Python, pandas, matplotlib, numpy array, seaborn, and python notebook were used to support the analysis. Python is a frequently utilised technology in the computing world (particularly in the analysis segment) due to its large number of libraries and ease of usage with various methods. As an

interpreter, it has become well-known for processing big files. Python is an excellent data munging language. Pandas is a data manipulation Data Frame object with integrated indexing that is fast and optimized. It accepts files in the csv files, text files, SQL database formats. It provides for high-speed dataset combining. Pandas is a performance-oriented programming language built on Cython.

The main goal is to analyse the data and anticipate what customers would buy based on different product IDs. We successfully deployed four distinct prediction models on databases with five different situations, including case 1: 90% training data, 10% testing data, and case 2: 70% training data, 10% testing data, 30 percent data from testing, 50 percent data from training, 50 percent data from testing. We also trained on the entire dataset rather than the split that was used to forecast a sample data set. Cases 4 and 5 fit into this category, with case 4 predicting the 'Purchase' percentage of training data using models developed from training data and case 5 predicting the 'Purchase' amount of testing data with models developed from training data. Thus, the study reports on the effectiveness of the models in various cases, as well as their accuracy when tested using data samples, in which each customer id has standard transactions equivalent to 9000 for case 5.

#### A. Analysis

For reading and modifying csv files, the Pandas library is utilised. To avoid making the data repetitive, absent cells are represented with the number 999. The data is shown using Seaborn and matplotlib in the forms of histograms, boxplots, bar graphs, and scatter plots, among other things. Because data has numerical values, the sort of analysis is quantitative. It was effectively discovered how variables such as gender, marital status, and occupation affect a customer's buy rate, and this was depicted using a graphical presentation. To fill in missing values, data preprocessing is performed. Each case's analysis was carried out separately, along with their unique fresh datasets. Case 1 employs a data from insights vidhya that would be split 90% train and 10% testing. Case 2 employs a dataset from insights vidhya that would be split 70/30 for training and testing. Case 3 uses a dataset from analytics vidhya that is split 50/50 between training and testing. Case 4 uses the analytics vidhya training dataset, which can also be utilized as the test data so that appropriate predictions may be given for the very same customers (user id). In example 5, the data for training and testing is separated from the analytics vidhya. Missing values are filled using data preprocessing. Every instance, as well as their new corresponding datasets, was subjected to separate analysis. For case 1, the analytics vidhya dataset is split among 90% train and 10% testing. Case 2 uses a dataset from analytics vidhya that is split 70/30 for training and testing. Case 3 uses a 50/50 training/testing split from the analytics vidhya dataset. For case 4, the train data is collected from vidhya, that is often used to the test dataset, that corrects the predictions for same customers (user id) may be provided. In example 5, the data for train set and test set is kept separated by the analytics vidhya.

#### B. Prediction

For each of the four scenarios, wide variety of classification algorithms are utilized to train the dataset. Xgboost, TfidfTransformer, and ExtraTreesRegressor are used for prediction and training. By integrating the input data with the received data from TfidfTransformer, a new model is developed (model3). The article displays the variable importance and gains importance of the qualities, as well as the performance of the prediction for multiple situations. The top three attributes that have demonstrated the greatest improvements during the prediction are product id (f8), user id (f10), and product category 1 (f5). Product categories 1 (f5), 2 (f6), and 3 (f7) are the top three qualities that have showed the greatest improvements during the prediction. The supervised learning method Xgboost, often described as 'Extreme Gradient Boosting,' is used for supervised learning problems (training based on data history). The procedure is a mathematical model that describes how to produce a  $y_i$  forecast for a given  $x_i$ . This is the most accurate way of forecasting, especially for events like the kaggle competition.

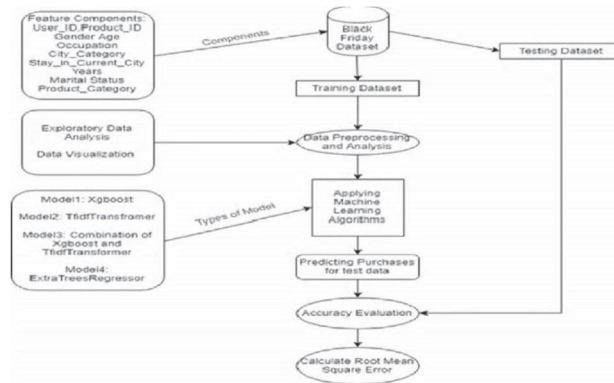


Fig1. Model for analysis and training

The Xgboost method library is offered in Python [34], and may be installed using pip or conda install. The data is converted into time-frequency appears in a files using TfidfTransform. It's a numerical metric that shows how significant a word is in a manuscript. It compensates for the reality that only a few words appear often in the dataset by increasing the ratio of instances a word occurs in the data and balancing the term's frequency in the same. This is done in order to reduce the impact of tokens that frequently appear in the data. TfidfTransformer may be found in sklearn.feature\_extraction.txt, It is pre-installed with the Anaconda framework and may be retrieved for use in other contexts. Extremely Randomized Trees' ExtraTreesRegressor provides a meta predictor employs averaging to increase projected system performance over-fitting by fitting a series of randomly selected decision trees on subsets of the dataset.

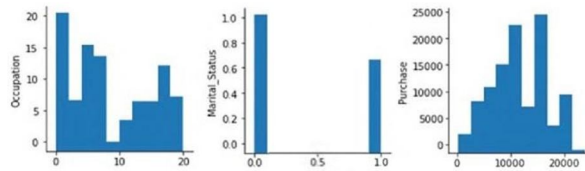


Fig2. Analysis and visualization

In a word, this strategy takes into account a number of different decision trees before averaging their precision for improved optimization. The algorithm is available in the sklearn. ensemble package, which is also available for the Anaconda framework and can be simply downloaded and deployed in other Python environments..

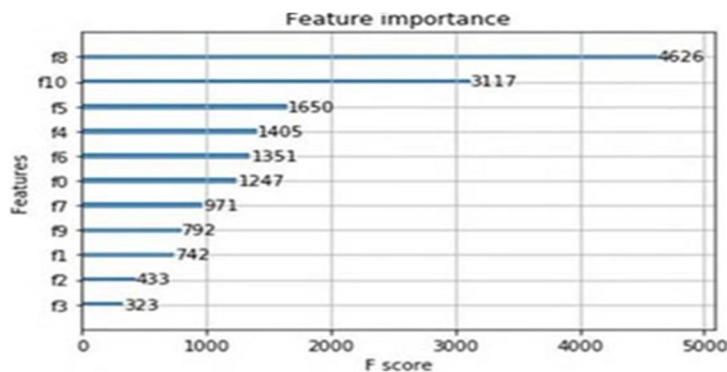


Fig3. Case1 Features

```
{'f0': 95783507.28227746,
'f1': 216730863.70215634,
'f10': 172154314.5323388,
'f2': 84781602.65357968,
'f3': 74766747.10216719,
'f4': 100222870.13209964,
'f5': 15486920770.45515,
'f6': 605541967.7616581,
'f7': 488772401.8939238,
'f8': 548647839.8479139,
'f9': 93567358.92462121}
```

Fig4. Gain of Case1

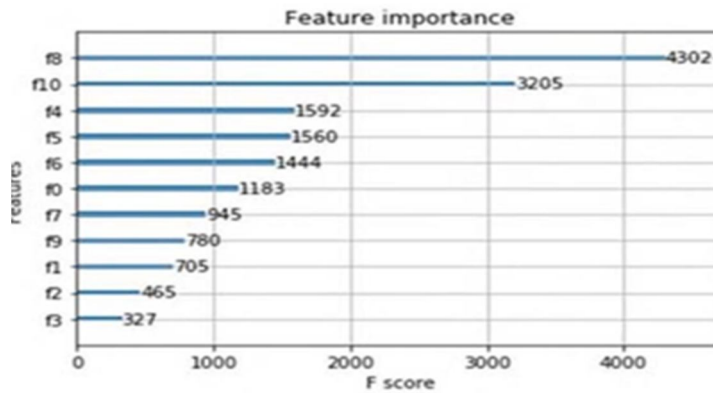


Fig5. Case2 Features

```
{'f0': 94880647.29152824,
'f1': 177820839.96666667,
'f10': 145338244.71123376,
'f2': 79826898.11328976,
'f3': 59943963.454301074,
'f4': 83115546.52774353,
'f5': 11950581016.86775,
'f6': 1045636457.0761495,
'f7': 2050991918.4542587,
'f8': 706737475.5845486,
'f9': 84552668.56644738}
```

Fig6. Gain of Case2

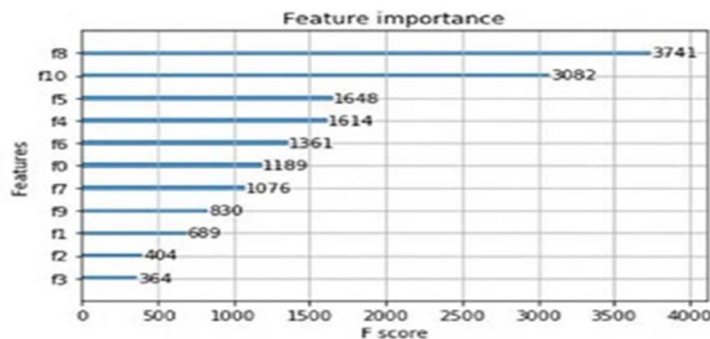


Fig7. Case3 Features

```
{'f0': 70060371.20117746,
'f1': 150959442.24687955,
'f10': 113336469.909695,
'f2': 62917851.95544554,
'f3': 54799929.557692304,
'f4': 69418926.51908302,
'f5': 8162556395.881068,
'f6': 770935033.9008082,
'f7': 897951689.1148698,
'f8': 547157523.5263566,
'f9': 66450518.57493976}
```

Fig8. Gain of Case3

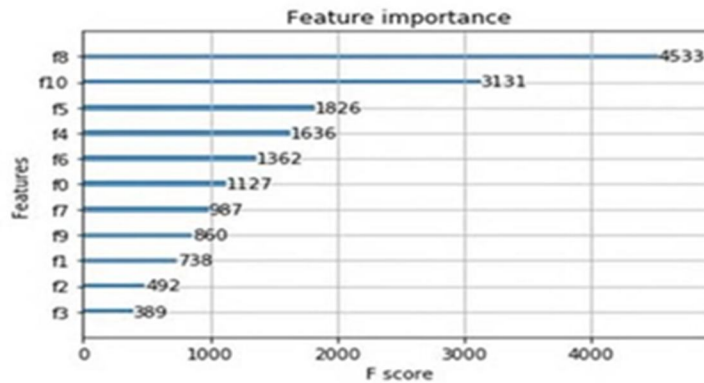


Fig9. Case4&5 Features

```
{'f0': 126083841.39129546,
'f1': 238378383.7134553,
'f10': 193910198.75650272,
'f2': 92954207.15447155,
'f3': 97840719.41645244,
'f4': 109823364.47848412,
'f5': 16769464181.98379,
'f6': 1144273232.8295887,
'f7': 865824058.9723403,
'f8': 684836281.6213766,
'f9': 102460940.94430234}
```

Fig10: Gain of Case4&5

## V. RESULTS AND DISCUSSIONS

When compared to their corresponding sample dataset, the results are provided in four sections and five row, with the numbers including the RMSE. The root of MSE, which can be derived from sklearn.metrics, was used to find the RMSE (Root mean square error). Sklearn metrics comes pre- installed in the framework and may be installed by using pip. From the first three examples, the Analytics Vidhya training dataset is used as an input data, Further train and test data are produced based on the case situation from the input dataset. In these type of instances, the example dataset will constructed by introducing the these parameters userid, product id, and buy, after this they eliminating the 'purchase' feature from the test data. The RMSE is use to verify the precision of the sample datasets when compared to the projected 'buy' quantity of the test. In the first scenario (case 1), the input data is divided into 90 percent and 10%, with 90% going to the training dataset and 10% going to the testing dataset. For each model, the RMSE score is generated

and displayed in the table. Figure also includes a graph depicting the similarity between anticipated and real 'buy' data. Figure 11 shows a graphical depiction of the values compared with each other.

In the second scenario, the input data is divided into two parts: 70% and 30%. 70% of the training dataset is maintained, and 30% of the testing dataset is kept. The RMSE score for each model is calculated by comparing the predicted 'buy' value of the testing and sample datasets. Figure shows a graph depicting the similarity between anticipated and real 'buy' data. Figure 12 shows a graphical depiction of the four models' values compared to one another and sample dataset. In the third scenario (case 3), the input dataset is split evenly, 50 percent to 50 percent. Both the training and testing datasets have the same amount of rows. In the table, the RMSE value for this case is shown. Figure shows a graph depicting the similarity between anticipated and real 'buy' data. Figure 13 shows a graphical depiction compared to each other.

In the fourth scenario, the train data is obtained from analytic vidhya, which is also used as a test dataset (case 4). This was done to see if the 'Purchase' value forecast was accurate for the same clients. With the same dataset, the approaches are used to identify the 'Purchase' value, which is then evaluated to the actual 'Purchase' value. The graph depicts the RMSE value for this case. Figure 14 shows the four models' values in relation with each other and with the real cost.

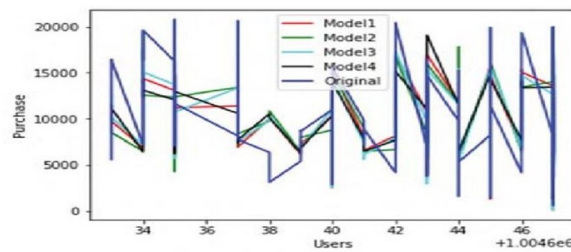


Fig 11. Case1 graph

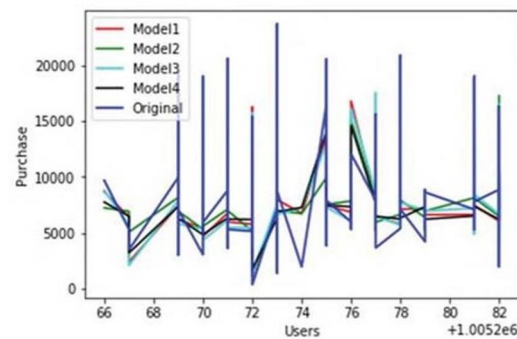


Fig 12. Case2 graph

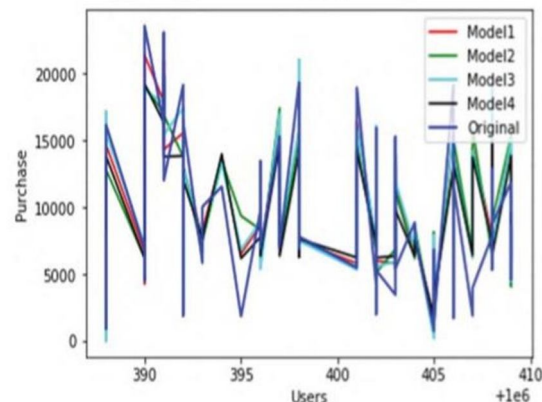


Fig 13. Case3 graph



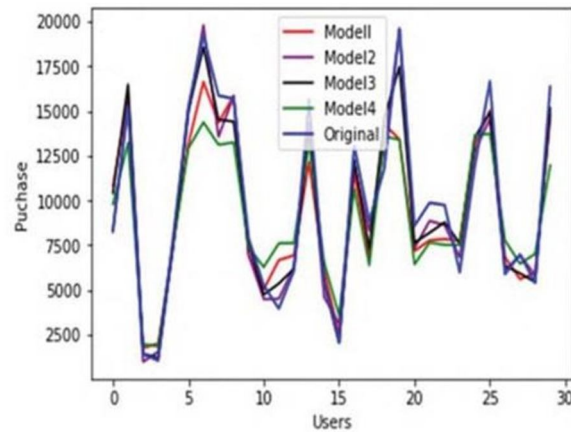


Fig14. Case14 Graph

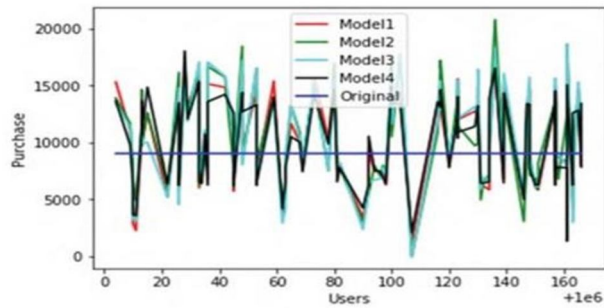


Fig15. Case15 Graph

## VI. CONCLUSION

ML techniques are trained on to use the Black Friday data to predict the 'purchase' values for specific product and user. Taking a look at the model's anticipated values and comparing them to the actual ones 90% of the data in Case 1 is for training and 10% is for testing. In this instance, overfitting may occur, in which the noise is misinterpreted as data to be trained, resulting in incorrect predictions with lower accuracy. In Case 2, 70% of the data is for training and 30% is for testing. This can be used for training and testing because, when compared to other methods, it has the highest likelihood of overfitting and underfitting data, and the sample data contains field data. Case 3 has 50 percent for train data and 50 percent for test data, and appears to be the most accurate. It is a viable option, but it is not suggested as the first choice. This is due to the fact that the training data is not very large, which can result in underfitting for some tests. Case 4 is made of a single database that is used as both a trained and a test dataset. Example 4 will likewise be ignored because the train and prediction were performed on same dataset in this case. Though the model's accuracy and forecasts have improved, the algorithm may still fail to predict proper values for fresh data. Case 5 is made up of training images tested data and example dataset derived from data analytics vidhya, However, since the data isn't connected to an actual circumstance or data, it's not taken into account.

## REFERENCES

- [1] Swilley, Esther, and Ronald E. Goldsmith. "Black Friday and Cyber Monday: Understanding consumer intentions on two major shopping days." *Journal of retailing and consumerservices*, vol. 20,1,2013, pp.43-50.
- [2] Fischer, Eileen, and Stephen J. Arnold. "Sex, gender identity, gender role attitudes, and consumer behavior." *Psychology & Marketing*, vol.11, 2, 1994, pp.163-182.
- [3] Song, Ji Hee, and Jason Q. Zhang. "Why do people shop online?: Exploring the quality of online shopping experience." *American Marketing Association. Conference Proceedings*. 2004.



- [4] Vijayarathy, Leo R. "Predicting consumer intentions to use on-line shopping: the case for an augmented technology acceptance model." *Information & management*, vol. 41,6, 2004, pp.747-762.
- [5] Simpson, Linda, et al. "An analysis of consumer behavior on Black Friday." *American International Journal of Contemporary Research*, 2011.
- [6] Bellizzi, Joseph A., and Robert E. Hite. "Environmental color, consumer feelings, and purchase likelihood." *Psychology & marketing*, vol.9, 5, 1992, pp.347-363.
- [7] Donovan, Robert J., et al. "Store atmosphere and purchasing behavior." *Journal of retailing*, vol.70, 3, 1994, pp. 283-294.
- [8] Sandelowski, Margarete. "Focus on research methods combining qualitative and quantitative sampling, data collection, and analysis techniques." *Research in nursing & health*, vol. 23, 3, 2000, pp. 246-255.
- [9] Burke, Raymond R. "Technology and the customer interface: what consumers want in the physical and virtual store." *Journal of the academy of Marketing Science*, vol.30,4 ,2002, pp. 411-432.
- [10] Boyd Thomas, Jane, and Cara Peters. "An exploratory investigation of Black Friday consumption rituals." *International Journal of Retail & Distribution Management*, vol. 39, 7 ,2011, pp. 522-537.
- [11] Mladenoff, David J., et al. "A regional landscape analysis and prediction of favorable gray wolf habitat in the northern Great Lakes region." *Conservation Biology* , vol. 9,2 ,1995, pp. 279-294.





10.22214/IJRASET



45.98



IMPACT FACTOR:  
7.129



IMPACT FACTOR:  
7.429



# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24\*7 Support on Whatsapp)