



sciendo

BALTIC JOURNAL OF LAW & POLITICS

A Journal of Vytautas Magnus University
VOLUME 15, NUMBER 4 (2022)
ISSN 2029-0454

Cite: Baltic Journal of Law & Politics 15:4 (2022): 277-286
DOI: 10.2478/bjlp-2022-004029

Higher Accuracy of Spam Email Prediction using K-Nearest Neighbor Algorithm Comparing with Multinomial Naive Bayes Algorithm

Putta Charan

Research scholar, Department of Computer Science and Engineering, Saveetha School of Engineering, Saveetha Institute of Medical and Technical Sciences, Saveetha University, Chennai, Tamil Nadu, India, Pincode: 602105.

P. Sriramya

Project Guide, Corresponding Author, Department of Data Science, Saveetha School of Engineering, Saveetha Institute of Medical and Technical Sciences, Saveetha University, Chennai, Tamil Nadu, India, Pincode: 602105.

Received: August 8, 2022; reviews: 2; accepted: November 29, 2022.

Abstract

Aim: To make an Innovative Spam Prediction of spam emails using Machine learning modeling techniques and to evaluate their performance. **Materials and Methods:** The initiative's main goal is to collect samples from two different groups. The K-Nearest Neighbor Algorithm is responsible for Group-1, whereas the Multinomial Naive Bayes Algorithm is responsible for Group-2. For both Algorithms, the same sample sizes were used. 80% of the G-Power will be used in the test set. **Result:** Data is trained in the given model so that Machine learning can function effectively. Emails are used as inputs for the Multinomial Naive Bayes algorithm, which gives us a probabilistic index of the email and determines if it is spam or not. The K-Nearest Neighbor Algorithm outperforms the Multinomial Naive Bayes Algorithm, and our hypothesis is significant with a significance value of 0.011. **Conclusion:** These results were achieved through machine learning models such as Multinomial Naive Bayes, and K-Nearest Neighbors. In this paper, have demonstrated that for the spam filtering method the most efficient algorithms are KNN and MNB given as they have the highest level of accuracy.

Keywords

Classifier, Filtering, Innovative Spam Prediction, K-nearest neighbor, Machine learning, Multinomial Naive Bayes.

INTRODUCTION

The purpose of Spam Email classification is to automatically classify new emails as spam or ham based on their contents. There has been a significant growth in the number of emails received, necessitating effective approaches such as Text Mining and Natural Language Processing to automatically categorize emails as spam or ham nearly 4.1 billion Email accounts are created throughout the world and More than 196 billion Emails will be sent day by day. Spam-Emails are one of the main threats to Email Users (Kontsewaya, Antonov, and Artamonov 2021). In this paper, we compare the performance of two machine learning techniques for spam detection including K-Nearest Neighbor classifier

Comparing with Multinomial Naive Bayes classifier. Multinomial Naive Bayes Classifier takes more time during the training period but its classification speed is better than other classifiers. An unwanted Email sent in bulk to an unknown recipient is referred to as a spam Email (Akinyelu 2021). It refers to the use of an email system to send unsolicited emails, particularly marketing emails to a large number of people. These accounts perform all email traffic worldwide. Unsolicited emails indicate that the receiver has not been permitted to receive them. Spam emails have grown in popularity over the last decade and are a problem that most email users confront for filtering methods. The applications of the research are Users and emails (Hossain, Uddin, and Halder 2021), (Kumar, Sonowal, and Nishant 2020). Botnets or networks of infected computers may send massive amounts of spam emails.

Innovative Spam Prediction using K-Nearest Neighbor Algorithm comparing with Multinomial Naive Bayes Algorithm. In GoogleScholar this article is published 772 times, and in ScienceDirect, this article is published 72 times in the past 5 years. In these 2 databases, the most cited articles and their findings are, that the preliminary discussion in the research background looks at how machine learning methods are used in the email spam filtering processes of the top internet service providers (Dada et al. 2019). That they suggested a new spam detection approach that combines an artificial bee colony algorithm with a logistic regression classification model ("Spam Filtering Using a Logistic Regression Model Trained by an Artificial Bee Colony Algorithm" 2020). This suggested method identifies e-mail spam in both textual and speech-enabled e-mails. In terms of text extraction speed, performance, cost efficiency, and accuracy, the suggested GDTPNLP technique gives a greater spam detection rate (Ismail et al. 2022). Here The technology recognizes the required features for categorizing spam emails automatically. The suggested system is based on the Genetic Algorithm and the Random Weight Network (Faris et al. 2019). From the above literature analysis and study, the paper (Kontsewaya, Antonov, and Artamonov 2021) is most relevant to our study and done most of the analysis.

Previously our team has a rich experience in working on various research projects across multiple disciplines (Venu and Appavu 2021; Gudipaneni et al. 2020; Sivasamy, Venugopal, and Espinoza-González 2020; Sathish et al. 2020; Reddy et al. 2020; Sathish and Karthick 2020; Benin et al. 2020; Nalini, Selvaraj, and Kumar 2020). On Daily basis, Spam Email is continuously increasing day by day. The rapidly increasing Spam Emails are responsible for over 77% of the whole global email traffic, these motivated me to do the Research on Spam mail Prediction. The team in the department has much experience in research on Machine learning models, so it's helpful to come up with innovative ideas in machine learning approaches for developing efficient algorithms with higher accuracy in the spam email prediction and this shows experience in our lab for research of spam email prediction. The aim is to increase the accuracy value of the email spam prediction using Machine Learning techniques and predict if the email is spam or not and make an Innovative Spam Prediction of spam emails using Machine learning modeling techniques and evaluate their performance (Gaurav et al. 2019).

MATERIALS AND METHODS

This Research paper for Spam Email Prediction research is done in the Software Engineering Lab, Saveetha School of Engineering, SIMATS. The Dataset has been taken from Kaggle and this has an open-source license to download and use the data for the research. In this project, there will be mainly two groups of samples taken in the project. That Group-1 belongs to the K-Nearest Neighbor Algorithm and Group-2 Belongs to the Multinomial Naive Bayes Algorithm. The sample sizes of both groups are 30% and 70% total sample sizes taken are 100% of the data. The Same set of Sample sizes will have for both algorithms. Iteration-1 for the Training set and Iteration-2 for the Test set will have

80% of the G-Power (Maguluri et al. 2019). This helps to create a more Accurate Prediction for the Spam Mail using Machine Learning models.

Data Collection

The Data Set for this Research is collected from Kaggle which is an Open source Platform for getting Machine Learning Datasets. The Url for the datasets is mentioned below (ishansoni 2018). I got 10743 rows and 2 columns By combining the two datasets used in the Algorithms. In the Datasets, different dependent and independent Variables are Considered to Perform Machine Learning Techniques.

K-Nearest Neighbor Algorithm

The K-Nearest Neighbor method is one of the most fundamental Machine Learning algorithms and is based on the Supervised Learning methodology. The K-Nearest Neighbor approach assumes that new and existing data are comparable, and it assigns the new example to the category that is most similar to the existing categories. The K-Nearest Neighbor approach preserves all previously saved data and categorizes new data points based on their similarity. This means that when new data is available for filtering, the K-Nearest Neighbor approach can swiftly categorize it into an appropriate category. The K-Nearest Neighbor approach may be used for both regression and classification, while classification is more typically utilized. It is straightforward to build because all that is required is a probability calculation. We will forecast the Accuracy value of the K-Nearest Neighbor using the data. This is an iterative procedure that must be done for each data point in the dataset. Assume we already have a cleaned dataset that has been separated into training and testing data sets. To determine the distance among the testing dataset and each training dataset using (1) Euclidean distance metrics.

$$d(p, q) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + \dots + (p_i - q_i)^2 + \dots + (p_n - q_n)^2} \text{ ----- (1)}$$

Pseudocode for K-Nearest Neighbor

Input: Training dataset
Output: Classifier Accuracy
Read dataset as input
select 'k' samples from the total samples
 classifier. fit(x_train, y_train)
Among 'k' tokens calculate the node 'd' using the best split
Repeat 1 to 3 steps until many samples are reached
Build the K-classifier
Predict value using predict feature
prediction=model.predict(parameters, "")
Calculate vote for each predicted value
Get predicted accuracy
Get Test Results.

Multinomial Naive Bayes Algorithm

Multinomial Naive Bayes is a probabilistic learning method used in Natural Language Processing (NLP). Using the Bayes theorem, this approach guesses the tag of a text, such as an email or a news item. It computes the likelihood of each tag for a given sample and returns the tag with the highest likelihood. The Naive Bayes classifier is a group of algorithms that all follow the same basic principle: each feature being classified is unconnected to any other feature. One character's existence or absence has no bearing on the presence or absence of another. The formula for Naive Bayes efficiency and increase is used for text data analysis and multi-class scenarios (2). To understand how the Naive Bayes theorem works, you must first comprehend the Bayes theorem concept, as it is based on it. The Bayes theorem, developed by Thomas Bayes, states that previous knowledge of event-related circumstances does not affect the likelihood of an event occurring. It is calculated using the following equation:

$$P(A|B) = P(A) * P(B|A)/P(B) \text{ ----- (2)}$$

The probability of class A when predictor B is already provided.

$P(B)$ = prior probability of predictor B

$P(A)$ = prior probability of class A

$P(B|A)$ = occurrence of predictor B given class A probability

This formula helps in calculating the probability of the tags in the text.

Pseudocode for Multinomial Naive Bayes

Input: Training dataset

Output: Classifier accuracy

The first step is Data collection.

Pre-processing and text cleaning of the train data.

Fit the Training Data Set to the Multinomial Naive Bayes.

Now Predict the Results for test split data.

Define class

```
Def MultinomialNB()
```

```
if(condition satisfies)
```

```
    return accuracy
```

```
else
```

```
    return previous step
```

```
End
```

Create the Confusion Matrix find the Test Accuracy Results.

Get Test Results.

The platform used to evaluate the Machine learning Algorithm was Anaconda/Jupyter. The hardware used to perform the work is Intel(R) Core(TM) i7-10750H CPU @ 2.60GHz with a RAM size of 8 GB. The system type used was 64 bit, Windows OS, X64-based processor with an SSD of 256 GB. The Operating System used was Windows 10, and the tool used was JupyterLabs with the Python programming language. The testing procedure was to split the data into train and test data and then implement the Machine learning classifier to build and train a model on our data. After training, the predictions are made and the performance of the model is evaluated using the available metrics.

The dataset for Innovative spam prediction is collected from Kaggle. Data preprocessing was performed to gain some context about the data using Statistical Analysis techniques. Data cleaning methods such as removing unnecessary attributes, and contents and filling null values are done. The comparison of the K-Nearest Neighbor Algorithm and Multinomial Naive Bayes Algorithm with data exploration gives us some context and valuable insight into the dataset. The Spam Email Prediction with two widely spread classification algorithms in machine learning was selected K-Nearest Neighbor and Multinomial Naive Bayes. The algorithms will be trained with some data when the test data is given then it will predict the output whether the given email is spam or not. The testing data is used to give the predicted output and analyzes the data according to that.

Statistical Analysis

The IBM Spss is the Statistical Software Tool that is used for Spam Email data analysis. The IBM Statistical Tool can analyze the data and helps to create Graphs and Charts to display it quite easily. Before sending results into the Spss tool the Data sets are standardized and then the data is converted into arrays. The IBM tool can easily handle large data because it consists of a wide array of characteristics. The number of clusters required is pictured and analyzed and therefore the existing algorithms are obtained. It gives the Mean value for the Group statistics. The Group-1 and Group-2 Accuracy as shown in Table 1 the Different Test Sizes and their average accuracy values that are acquired after being tested with the K-Nearest Neighbor Classifier and Multinomial Naive Bayes Classifier with 10 Sample test sizes (Liu, Lu, and Nayak 2021). The Data Sets for the Spam Email Prediction are taken from the kaggle which consists of Both Dependent Variables and In-Dependent Variables in Table 2 and Table 3. The Statistical Comparison of The Spam Email Prediction using two Sample groups was done with the SPSS Version 25. The

Analysis was done using the Mean, Median, Independent T-Test, and Deviation. For each sample size of data, the Accuracy is deviating between 3% to 5 %. So that we finally sent all the Test sizes and also their Accuracy into the Spss tool and found the Average Accuracy values of the K-Nearest Neighbor Classifier and Multinomial Naive Bayes Classifier.

RESULT

In the proposed model, data is trained so that Machine learning can work properly. After applying the Multinomial Naïve Bayes algorithm, emails are taken as inputs which will give us the probabilistic index of that and will identify whether the Email is spam or not. This necessitates the development of a sensible method for detecting or identifying such spam emails, therefore saving a significant amount of time and memory space for the system. Spammers may easily create a false profile and email account by pretending to be a legitimate person in their spam emails. This paper will discuss machine learning algorithms and apply all of these algorithms to our data sets, and the best algorithm is selected for email spam detection with the highest precision and accuracy.

The Innovative Spam Prediction using K-Nearest Neighbor Algorithm gave us an accuracy of 91% and Multinomial Naive Bayes gave us an accuracy of 90% compared with their accuracy rate. Each algorithm was repeated 10 times, for each algorithm and the accuracy varies for different test sizes in decimals. The accuracy varies due to random changes in the test sizes of the algorithm as given in Table 1.

The observed values for the metrics of Group Statistics, the mean accuracy, and the standard deviation for the K-Nearest Neighbor are 90.187 and 0.79363. The Multinomial Naive Bayes Algorithm's mean accuracy is 87.997 and the standard deviation is 2.14172. The K-Nearest Neighbor also obtained a standard error mean rate of 0.25097 whereas the Multinomial Naive Bayes Algorithm obtained an error mean rate of 0.67727 as given in Table 2.

Then an independent sample test of 10 samples was performed, K-Nearest Neighbor obtained a mean difference of 2.1900 and a standard error difference of 0.72228. When compared to other algorithm's performance, the K-Nearest Neighbor performed better than the Multinomial Naive Bayes Algorithm and the significance value of 0.011 shows that our hypothesis is valid as given in Table 3.

It is called the Innovative Spam Prediction architecture. The architecture defines the steps which are performed to develop a spam email prediction. It consists of the steps as Data Pre-processing, Database, Data Extraction, Modeling Classifier, Implementation, and Predicted Accuracy.

The GGraph represents a bar chart of the simple bar mean accuracy, with the K-Nearest Neighbor Algorithm achieving an accuracy of approximately 91%, and the Multinomial Naive Bayes Algorithm achieving 90%. The 95% error bars represent the variation in the corresponding coordinates of the point. Independent t-tests were performed to compare the accuracy of the two algorithms and a statistically significant difference was noticed between the two algorithms $0.011 < 0.05$. When comparing the two algorithms the performance of the K-Nearest Neighbor Algorithm achieved a better performance than the Multinomial Naive Bayes Algorithm is given in Fig. 1.

DISCUSSION

The K-Nearest Neighbor has better accuracy than Multinomial Naive Bayes. The results are collected by performing multiple times for identifying different scales of accuracy rates. Independent samples t-tests are performed on the dataset. In this study of spam email prediction, the K-Nearest Neighbor Algorithm has an accuracy of

approximately 91%, which is higher than that of the Multinomial Naive Bayes Algorithm which is 90%. K-Nearest Neighbor has a better significance of 0.011 while using the independent samples T-test. The mean accuracy and standard deviation for the K-Nearest Neighbor Algorithm are 90.187 and 0.79363 using a missing value imputation and a machine learning model to get an accuracy of 91%. The Multinomial Naive Bayes Algorithm's mean accuracy is 87.997 and the standard deviation is 2.14172. In the paper, (Kontsewaya, Antonov, and Artamonov 2021) the K-Nearest Neighbor Algorithm obtained an accuracy of 90%, and (Sharaff and Rao 2020) the Multinomial Naive Bayes Algorithm achieved an accuracy of 89% accuracy. Based on the literature survey, it is evident that the K-Nearest Neighbor performs better than Multinomial Naive Bayes. By running independent sample tests in IBM's SPSS statistical program, it can be seen that the difference between the two algorithms is statistically significant at 0.011. The SPSS statistical program is also used to compute the mean and standard deviation.

Using IBM's SPSS statistical tool, independent sample analysis confirmed that the difference between the two methods is statistically significant at $0.011 < 0.05$. The mean and standard deviation are determined using the SPSS statistical tool. K-Nearest Neighbor outscored other algorithm classification accuracy by 91% percentage in this paper (Kontsewaya, Antonov, and Artamonov 2021).

The main limitation is that the attributes in the dataset contain fewer data to predict accuracy (%) for spam email classification. The more the independent and dependent variables the more accuracy will be improved. For future work, the dataset contains many attributes the classifier can work efficiently and can improve the prediction accuracy. Attributes like this can result in improved accuracy and exact precision values (Wood and Krasowski 2020). There exists a strong relationship between the content and the subject of the emails. With the help of this relationship, one can easily classify the documents. Positive value tells us how strongly that word belongs to the subject and negative tells how much it differs from a subject. With the help of a negative score also the accuracy of the classifier has been improved (Rafat et al. 2022) and this paper is to improve the relationship between the subject and the content of the email by identifying the most relevant words using evolutionary computation of Email.

CONCLUSION

These results were achieved through machine learning models such as Multinomial Naive Bayes, and K-Nearest Neighbors. In this paper, we have demonstrated that for the spam filtering method the most efficient algorithms are KNN and MNB given as they have the highest level of accuracy. These spammers target those who are unaware of these scams and have filtering issues. So, it is necessary to identify those spam emails that are fraudulent, this project will identify those spam by using machine learning techniques. The results can be used to create a more intelligent spam detection classifier by combining algorithms of filtering methods.

DECLARATIONS

Conflict of Interests

No conflict of interest in this manuscript.

Author Contribution

Author PCR was involved in data collection, data analysis, and manuscript writing. Author PSR was involved in the conceptualization, data validation, and critical review of the manuscript.

Acknowledgement

The authors would like to express their gratitude towards Saveetha School of Engineering, Saveetha Institute of Medical and Technical Sciences (Formerly known as Saveetha University) for providing the necessary infrastructure to carry out this work successfully.

Funding

We thank the following organizations for providing financial support that enabled us to complete the study.

1. Vee Eee Technologies Solutions Pvt. Ltd.
2. Saveetha University
3. Saveetha Institute of Medical and Technical Sciences
4. Saveetha School of Engineering

REFERENCES

- Akinyelu, A. A. 2021. "Advances in Spam Detection for Email Spam, Web Spam, Social Network Spam, and Review Spam: ML-Based and Nature-Inspired-Based Techniques." <https://doi.org/10.3233/JCS-210022>.
- Benin, S. R., S. Kannan, Renjin J. Bright, and A. Jacob Moses. 2020. "A Review on Mechanical Characterization of Polymer Matrix Composites & Its Effects Reinforced with Various Natural Fibres." *Materials Today: Proceedings* 33 (January): 798–805.
- Dada, E., J. Bassi, H. Chiroma, S. Abdulhamid, A. O. Adetunmbi, and Opeyemi Emmanuel Ajibuwa. 2019. "Machine Learning for Email Spam Filtering: Review, Approaches and Open Research Problems." <https://www.ncbi.nlm.nih.gov/pubmed/31211254>.
- Faris, Hossam, Ala' M. Al-Zoubi, Ali Asghar Heidari, Ibrahim Aljarah, Majdi M. Mafarja, Mohammad A. Hassonah, and H. Fujita. 2019. "An Intelligent System for Spam Detection and Identification of the Most Relevant Features Based on Evolutionary Random Weight Networks." <https://doi.org/10.1016/J.INFFUS.2018.08.002>.
- Gaurav, Devottam, Sanju Mishra Tiwari, Ayush Goyal, Niketa Gandhi, and Ajith Abraham. 2019. "Machine Intelligence-Based Algorithms for Spam Filtering on Document Labeling." *Soft Computing* 24 (13): 9625–38.
- Gudipani, Ravi Kumar, Mohammad Khursheed Alam, Santosh R. Patil, and Mohamed Isaqali Karobari. 2020. "Measurement of the Maximum Occlusal Bite Force and Its Relation to the Caries Spectrum of First Permanent Molars in Early Permanent Dentition." *The Journal of Clinical Pediatric Dentistry* 44 (6): 423–28.
- Hossain, Fahima, M. N. Uddin, and Rajib Kumar Halder. 2021. "Analysis of Optimized Machine Learning and Deep Learning Techniques for Spam Detection." <http://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=9422508>.
- ishansoni. 2018. "SMS Spam Collection Dataset." Kaggle. October 6, 2018. <https://kaggle.com/ishansoni/sms-spam-collection-dataset>.
- Ismail, Safaa S. I., Romany F. Mansour, Rasha M. Abd El-Aziz, and Ahmed I. Taloba. 2022. "Efficient E-Mail Spam Detection Strategy Using Genetic Decision Tree Processing with NLP Features." *Computational Intelligence and Neuroscience* 2022 (March): 7710005.
- Kontsewaya, Yuliya, Evgeniy Antonov, and Alexey Artamonov. 2021. "Evaluating the Effectiveness of Machine Learning Methods for Spam Detection." 6074. EasyChair. https://easychair.org/publications/preprint_open/Th28.
- Kumar, N., Sanket Sonowal, and Nishant. 2020. "Email Spam Detection Using Machine Learning Algorithms." <http://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=9183098>.
- Liu, Xiaoxu, Haoye Lu, and Amiya Nayak. 2021. "A Spam Transformer Model for SMS Spam Detection." *IEEE Access*. <https://doi.org/10.1109/access.2021.3081479>.
- Maguluri, Lakshmana Phaneendra, R. Ragupathy, Sita Rama Krishna Buddi, Vamshi Ponugoti, and Tharun Sai Kalimil. 2019. "Adaptive Prediction of Spam Emails : Using Bayesian Inference." 2019 3rd International Conference on Computing Methodologies and Communication (ICCMC). <https://doi.org/10.1109/iccmc.2019.8819744>.
- Nalini, Devarajan, Jayaraman Selvaraj, and Ganesan Senthil Kumar. 2020. "Herbal Nutraceuticals: Safe and Potent Therapeutics to Battle Tumor Hypoxia." *Journal of*

- Cancer Research and Clinical Oncology 146 (1): 1–18.
- Rafat, Khan Farhan, Qin Xin, Abdul Rehman Javed, Zunera Jalil, and Rana Zeeshan Ahmad. 2022. "Evading Obscure Communication from Spam Emails." *Mathematical Biosciences and Engineering: MBE* 19 (2): 1926–43.
- Reddy, Poornima, Jogikalmat Krithikadatta, Valarmathi Srinivasan, Sandhya Raghu, and Natanasabapathy Velumurugan. 2020. "Dental Caries Profile and Associated Risk Factors Among Adolescent School Children in an Urban South-Indian City." *Oral Health & Preventive Dentistry* 18 (1): 379–86.
- Sathish, T., and S. Karthick. 2020. "Gravity Die Casting Based Analysis of Aluminum Alloy with AC4B Nano-Composite." *Materials Today: Proceedings* 33 (January): 2555–58.
- Sathish, T., D. Bala Subramanian, R. Saravanan, and V. Dhinakaran. 2020. "Experimental Investigation of Temperature Variation on Flat Plate Collector by Using Silicon Carbide as a Nanofluid." In *PROCEEDINGS OF INTERNATIONAL CONFERENCE ON RECENT TRENDS IN MECHANICAL AND MATERIALS ENGINEERING: ICRTMME 2019*. AIP Publishing. <https://doi.org/10.1063/5.0024965>.
- Sharaff, Aakanksha, and Ulligaddala Srinivasa Rao. 2020. "Towards Classification of Email through Selection of Informative Features." <http://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=9071488>.
- Sivasamy, Ramesh, Potu Venugopal, and Rodrigo Espinoza-González. 2020. "Structure, Electronic Structure, Optical and Magnetic Studies of Double Perovskite Gd₂MnFeO₆ Nanoparticles: First Principle and Experimental Studies." *Materials Today Communications* 25 (December): 101603.
- "Spam Filtering Using a Logistic Regression Model Trained by an Artificial Bee Colony Algorithm." 2020. *Applied Soft Computing* 91 (June): 106229.
- Venu, Harish, and Prabhu Appavu. 2021. "Experimental Studies on the Influence of Zirconium Nanoparticle on Biodiesel–diesel Fuel Blend in CI Engine." *International Journal of Ambient Energy* 42 (14): 1588–94.
- Wood, Kelly E., and Matthew D. Krasowski. 2020. "Academic E-Mail Overload and the Burden of 'Academic Spam.'" *Academic Pathology* 7 (January): 2374289519898858.

TABLES AND FIGURES

Table 1. Accuracy Values for the Algorithms. The Data Accuracy for the K-Nearest Neighbor (Group-1) and Multinomial Naive Bayes (Group-2) with different Test sizes has been taken. In these different Test Sizes, the Accuracy value for K-Nearest Neighbor is 91.03 and the Multinomial Naive Bayes is 90.35.

S No	Test Size	Group-1 Accuracy	Group-2 Accuracy
1	0.2	90.55	90.32
2	0.25	90.72	90.35
3	0.3	90.53	89.97
4	0.35	89.84	89.31
5	0.4	90.85	88.64
6	0.45	91.03	87.94
7	0.5	90.65	87.36
8	0.55	89.6	86.47
9	0.6	89.65	85.58
10	0.7	88.45	84.03

Table 2. Group Statistics the mean accuracy and standard deviation for K-Nearest Neighbor (KNN) are 90.187 and 0.79363 and For Multinomial Naive Bayes(MNB) Algorithm is 87.9970 and 2.14172.

Group Statistics					
	KNN, MNB	N	Mean	Std. Deviation	Std. Error Mean
Accuracy	KNN	10	90.1870	.79363	.25097
	MNB	10	87.9970	2.14172	.67727

Table 3. Independent Samples Test. Independent t-tests were performed to compare the accuracy of the two algorithms and a statistically significant difference was noticed between the two algorithms $0.011 < 0.05$ and Std. Error Difference is noticed as .72228.

		Levene's Test for Equality of Variances		t-test for Equality of Means		
		F	Sig.	Std. Error Difference	95% Confidence Interval of the Difference	
					Lower	Upper
Accuracy	Equal variances assumed	7.934	.011	.72228	.67255	3.70745
	Equal variances not assumed			.72228	.60748	3.77252

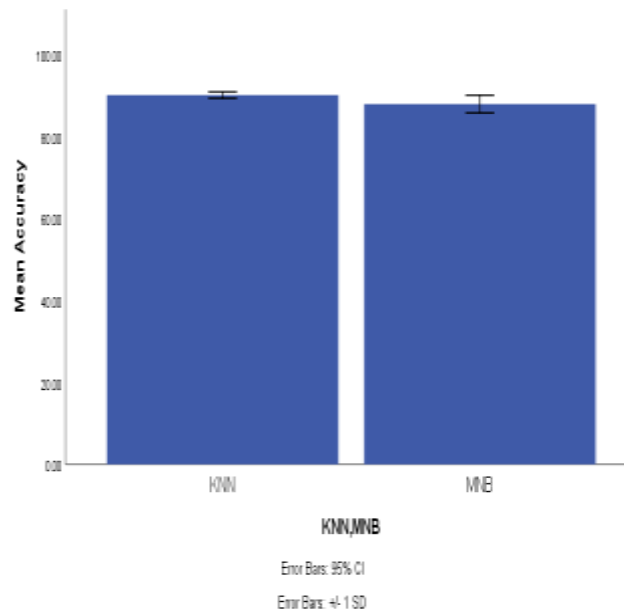


Fig. 1. Simple Bar Mean of Accuracy by K-Nearest Neighbor(KNN) and Multinomial Naive Bayes(MNB), the bar chart representing the comparison of mean accuracy of K-Nearest Neighbor Algorithm is 90.1870 and Multinomial Naive Bayes Algorithm is 87.9970. X-Axis: K-Nearest Neighbor Algorithm vs Multinomial Naive Bayes Algorithm. Y-Axis: Mean accuracy. The error bars are 95% for both algorithms. The Standard Deviation Error Bars are +/- 1 SD.