*A Novel Approach for Detecting Malicious Activities in Credit Card Transactions using K-Nearest Neighbour Algorithm to Improve Accuracy and Compared with Gaussian Naïve Bayes Algorithm*

*Section A-Research paper*

# A NOVEL APPROACH FOR DETECTING MALICIOUS ACTIVITIES IN CREDIT CARD TRANSACTIONS USING K-NEAREST NEIGHBOUR ALGORITHM TO IMPROVE ACCURACY AND COMPARED WITH GAUSSIAN NAÏVE BAYES ALGORITHM

**B. Madhumitha[1], V. Parthipan[2*]**

**Abstract**

**Aim:** The main aim of the research work is to create and build a novel fraud detection approach for Streaming Transaction Data, with the goal of analyzing historical customer transaction details and extracting behavioural patterns. Cardholders are divided into groups based on the volume of their transactions.
**Materials and Methods:** The categorizing is performed by adopting a sample size of n = 10 in K-Nearest Neighbour and sample size n = 10 in Gaussian Naive Bayes algorithms with G power in 80% and threshold 0.05%, CI 95% mean and standard deviation . For the implementation, the FraudTest dataset was used.
**Results :** The analysis of the results shows that the K-Nearest Neighbor has a high accuracy of (99.53) in comparison with Gaussian Naive Bayes algorithm (81.95). There is a statistically significant difference between the two groups with value p=0.005 (p< 0.05).
**Conclusion:** The results show that the K-Nearest Neighbor algorithm for detecting fraud in credit card transactions appears to generate better accuracy than Gaussian Naive Bayes algorithm.

**Keywords:** Novel Fraud Detection, Credit Card, Machine Learning, K-Nearest Neighbour, Gaussian Naïve Bayes, Streaming Transaction.

[1]Research Scholar, Department of Computer Science and Engineering, Saveetha School of Engineering, Saveetha Institute of Medical and Technical Sciences, Saveetha University, Tamil Nadu, India, PinCode: 602105
[2*]Project Guide, Department of Computer Science and Engineering, Saveetha School of Engineering, Saveetha Institute of Medical and Technical Sciences, Saveetha University, Tamil Nadu, India, PinCode: 602105

Eur. Chem. Bull. 2023, 12 (S1), 3559 – 3566

3559

*A Novel Approach for Detecting Malicious Activities in Credit Card Transactions using K-Nearest Neighbour Algorithm to Improve Accuracy and Compared with Gaussian Naïve Bayes Algorithm*

*Section A-Research paper*

## 1. Introduction

The purpose of the research work is to design and develop a novel fraud detection method for Streaming Transaction Data, to extract the behavioural patterns of customers by analyzing customers previous transaction records (Li et al. 2021). Credit card can be defined as a card that is given to customers by the banks which allows them to make payments within the credit limit and the customers can also withdraw cash in advance, for which they have to repay the bank within a given time. Credit card fraud can be defined as a fraud committed during the payment using this card. Credit card fraud can also lead to identity theft. This project is to propose a credit card fraud detection system using supervised learning algorithms (Dal Pozzolo et al. 2018). Credit card fraud can be done in a variety of ways. Fraudsters are highly skilled and quick-thinking individuals. This study uses the traditional approach to identify Application Fraud, which occurs when a person provides false information about himself in order to obtain a credit card. There's also unlawful use of Lost and Stolen Cards, which accounts for a sizable portion of credit card fraud (Dal Pozzolo et al. 2018; Seeja and Zareapoor 2014). The number of legitimate streaming transaction data greatly outnumbers the number of fraudulent transactions. To reduce their losses, banks and financial institutions have turned to novel fraud detection methods. Fraud detection systems can be used to detect fraudulent transactions in a variety of industries (Borzykowski 2013).

There are around 35 IEEE, Sciencedirect, Springer articles, and 30 Google scholar papers published in this domain over the past few years. For fraud detection, a variety of supervised and semi-supervised machine learning approaches are applied (Lamba 2020). However, we want to address three major issues with the card fraud dataset: severe class imbalance, the inclusion of labelled and unlabelled samples, and the ability to handle a large number of transactions. To detect fraudulent transactions in real-time datasets, several Supervised machine learning methods such as Decision Trees, Naive Bayes Classification, Least Squares Regression, Logistic Regression, and SVM are utilized. To train the behavioural aspects of normal and aberrant transactions, two approaches under random forests are utilized. Even while random forest produces decent results on small sets of data, it has certain issues when dealing with unbalanced data (Goyal and Sharma 2020). Potential fraud instances must be recognised in real time and labelled before the transaction is allowed or rejected, which is a typical uncertain domain. The incorporation of uncertainty aspects has an impact on an event processing engine's architecture and logic at all levels (Dal Pozzolo et al. 2018). The goal of this research is to develop a reliable and thorough deception financial detection model. The number of internet transactions has increased dramatically in recent years. Credit card transactions account for a significant share of these transactions. On the other hand, the growth of internet fraud is remarkable, which is mostly due to everyone's easy access to cutting-edge technology (Seeja and Zareapoor 2014). When compared to previous transactions made by the customer, card transactions are always unfamiliar. In the actual world, this unfamiliarity is known as concept drift problems, and it is a challenging problem to solve. Concept drift is a variable that evolves over time and in unexpected ways. This allows for the flexible design of event-driven systems with uncertainty characteristics from various domains. A first application was created in the field of detecting credit card fraud. Our preliminary findings are positive, indicating that adding uncertainty factors into the area of credit card fraud detection can result in significant gains (Baesens, Verbeke, and Van Vlasselaer 2015).

Our institution is passionate about high quality evidence based research and has excelled in various domains (Vickram et al. 2022; Bharathiraja et al. 2022; Kale et al. 2022; Sumathy et al. 2022; Thanigaivel et al. 2022; Ram et al. 2022; Jothi et al. 2022; Anupong et al. 2022; Yaashikaa, Keerthana Devi, and Senthil Kumar 2022; Palanisamy et al. 2022).Based on the literature survey, the methods utilized in prior articles produce less accurate outcomes. Previous studies used a smaller number of transactions to train the system, making it less efficient than the existing systems. The addition of more transactions improved the suggested systems efficiency. The project's purpose is to build and develop a novel fraud detection approach for Streaming Transaction data that analyses customers' historical transaction details and extracts behavioural patterns using the K-Nearest Neighbour algorithm.

## 2. Materials and Methods

The research work was performed in the OOAD Lab, Department of Computer Science and Engineering, Saveetha School of Engineering, Saveetha Institute of Medical and Technical Sciences, Chennai. Basically it is considered that two groups of classifiers are used, namely K-Nearest Neighbour and Gaussian

Eur. Chem. Bull. 2023, 12 (S1), 3559 – 3566

3560

*A Novel Approach for Detecting Malicious Activities in Credit Card Transactions using K-Nearest Neighbour Algorithm to Improve Accuracy and Compared with Gaussian Naïve Bayes Algorithm*

*Section A-Research paper*

Naive Bayes algorithms, which are used to classify the fraud in credit card transactions. Group 1 is the K-Nearest Neighbour algorithm with the sample size of N=10 and the Gaussian Naive Bayes algorithm is group 2 with sample size of N=10 and they are compared for more accuracy score and precision score values for choosing the best algorithm (Goyal and Sharma 2020).

The Pre- test analysis has been prepared using clinical.com by having a G power of 80% (Verma and Verma 2017) and threshold 0.05%, CI 95% mean and standard deviation. The fraudTest dataset was used in the study. This dataset was taken from the kaggle open source website. The K-Nearest Neighbour algorithm was chosen for implementation in this study, and it was compared to Gaussian Naive Bayes algorithm.

### K-Nearest Neighbour

In comparison to other supervised statistical pattern recognition fraud detection strategies, this is a supervised learning methodology that consistently produces good performance. Distance to find the least distant neighbours, some criterion to deduce a categorization from k-nearest neighbour, and the count of neighbours to label the new sample are three important aspects that influence its performance. This technique classifies any streaming transaction data by computing the transaction's least remote neighbour, and if that neighbour is tagged as fraudulent, the new transaction is likewise labelled as fraudulent. In this circumstance, Euclidean distance is a solid choice for calculating distances. This method is quick and produces fault alerts. Distance metric adjustment can help it perform better.

### Algorithm for K-Nearest Neighbour

Step 1:scaler = RobustScaler()
Step 2: X_train = scaler.fit_transform(X_train)
Step 3: X_test = scaler.transform(X_test)
Step 4: param_grid = {'n_neighbors': range(1,20)}
Step 5: clf = RandomizedSearchCV(KNeighborsClassifier(), param_grid)
Step 6: clf.fit(X_train,y_train)
Step 7: clf_pred = clf.predict(X_test)
Step 8: from sklearn.model_selection import cross_val_score
Step 9: scores = cross_val_score(clf,X_train,y_train,cv=10)
Step 10: print("accuracy scores",scores*100)
Step 11:Stop

### Gaussian Naive Bayes

The Bayes Theorem is used in this technique to compute the probability of a hypothesis and determine whether it is true or false. A classifier is used to calculate conditional probabilities for all possible classes and then place it in the class with the highest conditional probability for a given value of X. It is graphically represented as an acyclic directed graph, in which the samples are represented by the nodes of the graph and the dependencies between them are reflected by the directed edges. If there are no connecting edges between two variables, they are said to be independent. It also offers the joint probability distribution a specification and factorization.

### Algorithm for Gaussian Naive Bayes

Step 1: gnb = GaussianNB()
Step 2: gnb.fit(X_train,y_train)
Step 3: gnb_pred = gnb.predict(X_test)
Step 4: from sklearn.model_selection import cross_val_score
Step 5: scores = cross_val_score(gnb,X_train,y_train,cv=10)
Step 6: print("accuracy scores",scores*100)
Step 7: print(confusion_matrix(y_test,gnb_pred))
Step 8: print(classification_report(y_test,gnb_pred))
Step 9:Stop

The data collection of this research topic, detection and no detection data are observed, collected, and stored as a dataset. With the help of the device and the data, find the accuracy from the statistics tool or software.

The testing setup has all the components to do our test process. The testing setup has 2 types of configurations, Hardware configuration, and Software configuration. The Hardware configurations include Intel core i3 5[th] generation processor, 8 GB RAM (Random Access Memory), 64-bit Windows OS. The software configuration includes Windows OS. The language which is used to code the program is Python language.

### Statistical Analysis

IBM SPSS v23 is used for statistical analysis. The independent sample t-test calculation for analyzing equal variance, standard error, and levene's test are evaluated. The independent data sets are transaction id,cardholder id. The dependent values are amount, date, time. The independent T-test analysis is performed.

### 3. Results

Eur. Chem. Bull. 2023, 12 (S1), 3559 – 3566

3561

*A Novel Approach for Detecting Malicious Activities in Credit Card Transactions using K-Nearest Neighbour Algorithm to Improve Accuracy and Compared with Gaussian Naïve Bayes Algorithm*

*Section A-Research paper*

K-Nearest Neighbour is classified as a better algorithm because it has high output efficiency than Gaussian Naive Bayes. K-Nearest Neighbour proved with better accuracy than Gaussian Naive Bayes. Analysis of accuracy rate is done by varying the size of the datasets in Table 1. Table 2. shows the findings of the group statistics on all data variables with the count of N=10 and it calculates the mean, standard deviation, and standard error mean.Because it uses the clustering trick as a transformation approach, K-Nearest Neighbour (99.53) acquire superior accuracy and standard deviation when compared to Gaussian Naive Bayes (81.95). As a result of these transformations, it obtains the best boundary between the viable outcomes. Because of the relevance of equality of variance, the probability value states that the results in the research effort are significant and correlated with each other, the table demonstrates the difference in accuracy of both KNN and Gaussian Naive Bayes.The accuracy comparison of KNN and Gaussian Naive Bayes algorithms is shown in Fig. 1. The results of the independent sample t-test are shown in Table 3. Because of its efficient classification feature based on the clustering trick, the algorithm outperforms the Gaussian Naive Bayes algorithm.

## 4. Discussion

In this study, it was discovered that the K-Nearest Neighbour algorithm outperforms existing Gaussian Naive Bayes algorithm with an accuracy of 99.53% due to the consideration of more number of transactions, whereas existing Gaussian naive bayes consider less transactions (81.95%). The existing system considers less transactions, but the suggested method considers more number of transactions.

The authors investigated the effectiveness of classification models in detecting credit card fraud and offered three classification models: decision tree, neural network, and logistic regression. The neural network and logistic regression outperform the decision tree among the three models (Seeja and Zareapoor 2014). Following a review of Bayesian theory, the nave bayes classifier and k-nearest neighbour classifier are implemented and applied to the credit card system dataset. Data mining applications, automated fraud detection, and adversarial detection are among the strategies used in this domain, according to a comprehensive survey . Another article discussed credit card fraud detection approaches such as Supervised and Unsupervised Learning. Despite their unexpected success in some areas, these methods and algorithms failed to provide

a long-term and consistent answer to fraud detection (Jurgovsky 2019). The research was cited for credit card fraud detection, and seven categorization algorithms were applied. To reduce the risk of the banks, they used decision trees and SVMs in this study (Dal Pozzolo et al. 2018). Artificial Neural Networks and Logistic Regression Classification Models are more useful in improving fraud detection ability, according to them. A similar study domain was reported in which they employed Outlier mining, Outlier detection mining, and Distance sum algorithms to accurately forecast fraudulent transactions in an emulation experiment of credit card transaction data from a single commercial bank (Goyal and Sharma 2020). There have also been attempts to move forward from an entirely other perspective. In the event of a fraudulent transaction, efforts have been made to improve the alert feedback interaction (Dorronsoro et al. 1997)). In the event of a fraudulent transaction, the authorized system will be notified, and a response will be delivered to refuse the current transaction. In the event of a fraudulent transaction, the authorized system will be notified, and a response will be delivered to refuse the current transaction (Borzykowski 2013).

Despite the fact that the suggested KNN algorithm outperformed the previous approach,this study has a few shortcomings. To improve the accuracy of the suggested technique, it might be evaluated on a real time dataset rather than an existing dataset. To achieve more efficient outcomes, this system can be implemented utilizing a variety of existing machine learning methods.

## 5. Conclusion

Credit card fraud detection system was successfully developed . The current study focused on machine learning algorithms, KNN over gaussian naive bayes for higher classification in detecting fraud. It can be slightly improved based on the random data sets analysis in future. The outcome of the study shows KNN 99.53% has higher accuracy than Gaussian Naive Bayes 81.95%.

**Declarations**
**Conflict of Interests:** No conflict of interest

**Authors Contribution:** Author BM was involved in data collection, data analysis, manuscript writing. Author VP was involved in the Action process, Data verification and validation, and Critical review of manuscript.

*A Novel Approach for Detecting Malicious Activities in Credit Card Transactions using K-Nearest Neighbour Algorithm to Improve Accuracy and Compared with Gaussian Naïve Bayes Algorithm*

*Section A-Research paper*

### 5. References

Anupong, Wongchai, Lin Yi-Chia, Mukta Jagdish, Ravi Kumar, P. D. Selvam, R. Saravanakumar, and Dharmesh Dhabliya. 2022. "Hybrid Distributed Energy Sources Providing Climate Security to the Agriculture Environment and Enhancing the Yield." *Sustainable Energy Technologies and Assessments*. https://doi.org/10.1016/j.seta.2022.102142.

Baesens, Bart, Wouter Verbeke, and Veronique Van Vlasselaer. 2015. *Fraud Analytics Using Descriptive, Predictive, and Social Network Techniques: A Guide to Data Science for Fraud Detection*. John Wiley & Sons.

Bharathiraja, B., J. Jayamuthunagai, R. Sreejith, J. Iyyappan, and R. Praveenkumar. 2022. "Techno Economic Analysis of Malic Acid Production Using Crude Glycerol Derived from Waste Cooking Oil." *Bioresource Technology* 351 (May): 126956.

Borzykowski, Bryan. 2013. *Debit Card Fraud Detection For Dummies (Custom)*. John Wiley & Sons.

Dal Pozzolo, Andrea, Giacomo Boracchi, Olivier Caelen, Cesare Alippi, and Gianluca Bontempi. 2018. "Credit Card Fraud Detection: A Realistic Modeling and a Novel Learning Strategy." *IEEE Transactions on Neural Networks and Learning Systems* 29 (8): 3784–97.

Dorronsoro, J. R., F. Ginel, C. Sgnchez, and C. S. Cruz. 1997. "Neural Fraud Detection in Credit Card Operations." *IEEE Transactions on Neural Networks / a Publication of the IEEE Neural Networks Council* 8 (4): 827–34.

Goyal, Yogita, and Anand Sharma. 2020. *Credit Card Fraud Detection and Analysis Through Machine Learning*.

Jothi, K. Jeeva, K. Jeeva Jothi, S. Balachandran, K. Mohanraj, N. Prakash, A. Subhasri, P. Santhana Gopala Krishnan, and K. Palanivelu. 2022. "Fabrications of Hybrid Polyurethane-Pd Doped ZrO2 Smart Carriers for Self-Healing High Corrosion Protective Coatings." *Environmental Research*. https://doi.org/10.1016/j.envres.2022.113095.

Jurgovsky, Johannes. 2019. *Context-Aware Credit Card Fraud Detection*.

Kale, Vaibhav Namdev, J. Rajesh, T. Maiyalagan, Chang Woo Lee, and R. M. Gnanamuthu. 2022. "Fabrication of Ni–Mg–Ag Alloy Electrodeposited Material on the Aluminium Surface Using Anodizing Technique and Their Enhanced Corrosion Resistance for Engineering Application." *Materials Chemistry and Physics*. https://doi.org/10.1016/j.matchemphys.2022.125900.

Lamba, Harshit. 2020. *Credit Card Fraud Detection In Real-Time*.

Li, Yuening, Zhengzhang Chen, Daochen Zha, Kaixiong Zhou, Haifeng Jin, Haifeng Chen, and Xia Hu. 2021. "Automated Anomaly Detection via Curiosity-Guided Search and Self-Imitation Learning." *IEEE Transactions on Neural Networks and Learning Systems* PP (September). https://doi.org/10.1109/TNNLS.2021.3105636.

Palanisamy, Rajkumar, Diwakar Karuppiah, Subadevi Rengapillai, Mozaffar Abdollahifar, Gnanamuthu Ramasamy, Fu-Ming Wang, Wei-Ren Liu, Kumar Ponnuchamy, Joongpyo Shim, and Sivakumar Marimuthu. 2022. "A Reign of Bio-Mass Derived Carbon with the Synergy of Energy Storage and Biomedical Applications." *Journal of Energy Storage*. https://doi.org/10.1016/j.est.2022.104422.

Ram, G. Dinesh, G. Dinesh Ram, S. Praveen Kumar, T. Yuvaraj, Thanikanti Sudhakar Babu, and Karthik Balasubramanian. 2022. "Simulation and Investigation of MEMS Bilayer Solar Energy Harvester for Smart Wireless Sensor Applications." *Sustainable Energy Technologies and Assessments*. https://doi.org/10.1016/j.seta.2022.102102.

Seeja, K. R., and Masoumeh Zareapoor. 2014. "FraudMiner: A Novel Credit Card Fraud Detection Model Based on Frequent Itemset Mining." *TheScientificWorldJournal* 2014 (September): 252797.

Eur. Chem. Bull. 2023, 12 (S1), 3559 – 3566

3563

*A Novel Approach for Detecting Malicious Activities in Credit Card Transactions using K-Nearest Neighbour Algorithm to Improve Accuracy and Compared with Gaussian Naïve Bayes Algorithm*

*Section A-Research paper*

Sumathy, B., Anand Kumar, D. Sungeetha, Arshad Hashmi, Ankur Saxena, Piyush Kumar Shukla, and Stephen Jeswinde Nuagah. 2022. "Machine Learning Technique to Detect and Classify Mental Illness on Social Media Using Lexicon-Based Recommender System." *Computational Intelligence and Neuroscience* 2022 (February): 5906797.

Thanigaivel, Sundaram, Sundaram Vickram, Nibedita Dey, Govindarajan Gulothungan, Ramasamy Subbaiya, Muthusamy Govarthanan, Natchimuthu Karmegam, and Woong Kim. 2022. "The Urge of Algal Biomass-Based Fuels for Environmental Sustainability against a Steady Tide of Biofuel Conflict Analysis: Is Third-Generation Algal Biorefinery a Boon?" *Fuel*. https://doi.org/10.1016/j.fuel.2022.123494.

Verma, Priyam, and J. Verma. 2017. *Determination of Sample Size and Power Analysis with G\*Power Software: Step-Wise Illustrated Manual for Research Scholars*.

Vickram, Sundaram, Karunakaran Rohini, Krishnan Anbarasu, Nibedita Dey, Palanivelu Jeyanthi, Sundaram Thanigaivel, Praveen Kumar Issac, and Jesu Arockiaraj. 2022. "Semenogelin, a Coagulum Macromolecule Monitoring Factor Involved in the First Step of Fertilization: A Prospective Review." *International Journal of Biological Macromolecules* 209 (Pt A): 951–62.

Yaashikaa, P. R., M. Keerthana Devi, and P. Senthil Kumar. 2022. "Algal Biofuels: Technological Perspective on Cultivation, Fuel Extraction and Engineering Genetic Pathway for Enhancing Productivity." *Fuel*. https://doi.org/10.1016/j.fuel.2022.123814.

Alhazmi, Abrar Hassan, and Nojood Aljehane. 2020. "A Survey Of Credit Card Fraud Detection Use Machine Learning." 2020 International Conference on Computing and Information Technology (ICCIT-1441). https://doi.org/10.1109/iccit-144147971.2020.9213809.

Azhan, Mohammed, and Shazli Meraj. 2020. "Credit Card Fraud Detection Using Machine Learning and Deep Learning Techniques." 2020 3rd International Conference on Intelligent Sustainable Systems (ICISS). https://doi.org/10.1109/iciss49785.2020.9316002.

Bodepudi, Hariteja. 2021. "Credit Card Fraud Detection Using Unsupervised Machine Learning Algorithms." International Journal of Computer Trends and Technology. https://doi.org/10.14445/22312803/ijctt-v69i8p101.

Dornadula, Vaishnavi Nath, and S. Geetha. 2019. "Credit Card Fraud Detection Using Machine Learning Algorithms." Procedia Computer Science. https://doi.org/10.1016/j.procs.2020.01.057.

Huang, Jiayi. 2020. "Credit Card Transaction Fraud Using Machine Learning Algorithms." Proceedings of the 2019 International Conference on Education Science and Economic Development (ICESED 2019). https://doi.org/10.2991/icesed-19.2020.14.

Lacruz, Francisco, and Jafar Saniie. 2021. "Applications of Machine Learning in Fintech Credit Card Fraud Detection." 2021 IEEE International Conference on Electro Information Technology (EIT). https://doi.org/10.1109/eit51626.2021.9491903.

Oghenekaro, L. U., and C. Ugwu. 2016. "A Novel Machine Learning Approach to Credit Card Fraud Detection." International Journal of Computer Applications. https://doi.org/10.5120/ijca2016909316.

Popat, Rimpal R., and Jayesh Chaudhary. 2018. "A Survey on Credit Card Fraud Detection Using Machine Learning." 2018 2nd International Conference on Trends in Electronics and Informatics (ICOEI).

Eur. Chem. Bull. 2023, 12 (S1), 3559 – 3566

3564

*A Novel Approach for Detecting Malicious Activities in Credit Card Transactions using K-Nearest Neighbour Algorithm to Improve Accuracy and Compared with Gaussian Naïve Bayes Algorithm*

*Section A-Research paper*

### Tables and Figures

Table 1. Comparing accuracy values with the different sample sizes. It Represents the malicious activities in credit card transactions, the accuracy of KNN (99.53), and the GNB algorithm (81.95).

| S.No. | K-Nearest Neighbour Accuracy(%) | Gaussian Naive Bayes Accuracy(%) |
|---|---|---|
| 1 | 99.27 | 81.71 |
| 2 | 99.35 | 81.98 |
| 3 | 99.39 | 83.05 |
| 4 | 99.31 | 82.47 |
| 5 | 99.09 | 82.84 |
| 6 | 99.41 | 82.40 |
| 7 | 99.30 | 82.39 |
| 8 | 99.52 | 82.32 |
| 9 | 99.50 | 81.95 |
| 10 | 99.40 | 81.68 |

Table 2. Group Statistics of K-nearest neighbor with Gaussian naive bayes by grouping the iterations with Sample size 10, Mean = 99.3540 , Standard Derivation = .12340 , Standard Error Mean = 0.03902. Descriptive Independent Sample Test of Accuracy and Precision is applied for the dataset in SPSS. Here it specifies Equal variances with and without assuming a T-Test Score of two groups with each sample size of 10.

| | Algorithm | N | Mean | Std.Deviation | Std.Error Mean |
|---|---|---|---|---|---|
| Accuracy | KNN | 10 | 99.5340 | .12340 | .03902 |
| | GNB | 10 | 81.9530 | .41625 | .13163 |

**Table 3. Independent Samples T-test for accuracy of KNN shows significance value achieved is p=0.005 (p<0.05), which shows that two groups are statistically significant. Mean Difference = 17.07500 and confidence interval = (.14876- 16.76248).**

| | Levene's Test for Equality of Variances | t -test for Equality of Means |
|---|---|---|
| | | |

Eur. Chem. Bull. 2023, 12 (S1), 3559 – 3566

3565

*A Novel Approach for Detecting Malicious Activities in Credit Card Transactions using K-Nearest Neighbour Algorithm to Improve Accuracy and Compared with Gaussian Naïve Bayes Algorithm*

*Section A-Research paper*

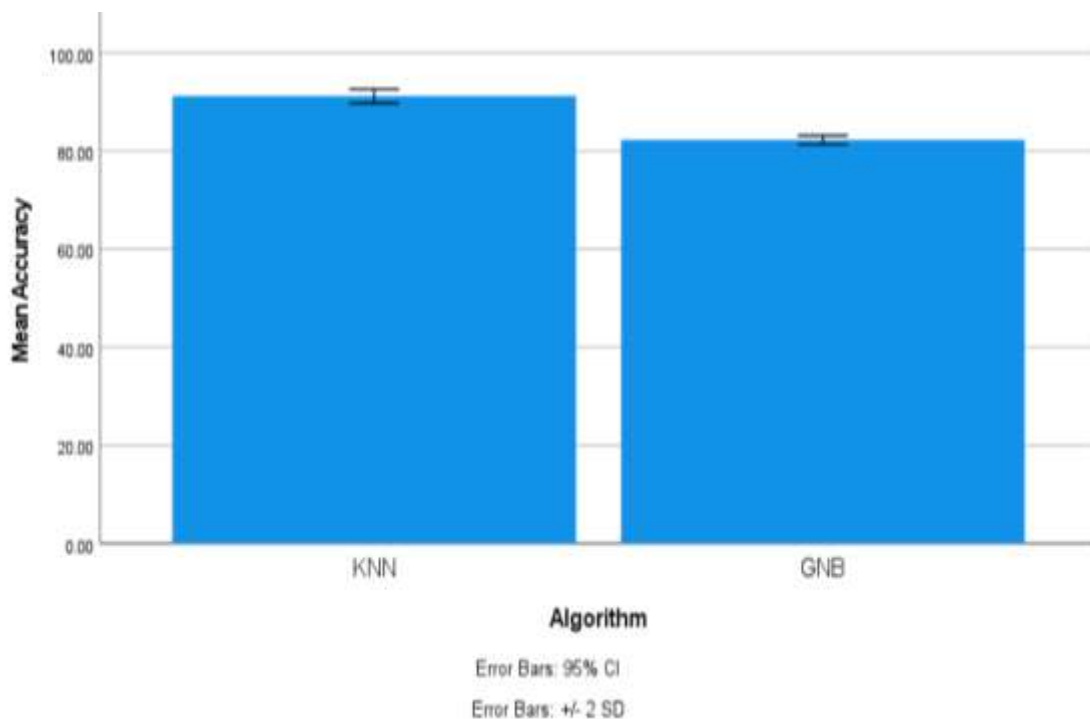| | | F | sig. | t | df | sig.(2-tailed) | Mean Difference | Std.Error Difference | 95% Confidence interval of the Difference | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | | | | | | | Lower | Upper |
| Accuracy | Equal variances assumed | 10.515 | .005 | 114.786 | 18 | .000 | 17.07500 | .14876 | 16.76248 | 17.38752 |
| | Equal variances not assumed | | | 114.786 | 10.323 | .000 | 17.07500 | .14876 | 16.74495 | 17.40505 |



Fig. 1. Comparison of KNN over GNB in terms of mean accuracy. It explores that the mean accuracy (99.53) is better than GNB (81.95) and the standard deviation is moderately improved KNN slightly lower than the GNB. Graphical representation of the bar graph is plotted using groupid as X-axis KNN vs GNB, Y-Axis displaying the error bars with a mean accuracy of detection +/- 2 SD.

Eur. Chem. Bull. 2023, 12 (S1), 3559 – 3566

3566