



Scopus®



# Journal of Vibration Engineering



ISSN:1004-4523

Registered



SCOPUS



DIGITAL OBJECT  
IDENTIFIER (DOI)



GOOGLE SCHOLAR



IMPACT FACTOR 6.1



Our Website

[www.jove.science](http://www.jove.science)

# Enhancing a Dataset to Improve Anomaly Detection Using Machine Learning

V.Sankar, Research Scholar, Department of Computer Science, Bharathiar University, Coimbatore, Tamil Nadu, India

Dr.G.Zayaraz, Professor, Department of CSE, Pondicherry Engineering College, Pondicherry

**ABSTRACT** This paper discusses the data security of customer and business data. The business website receives requests from various sources such as customers, crawlers, bots, and hackers. The elimination of anomalies improves the effective usage of hardware, software, and network bandwidth. A machine learning technique is used to identify anomalies from each API request.

The quality of the dataset is important for better accuracy. This paper proposes to improve the quality of the training dataset using Quadratic Discriminant Analysis and Linear Discriminant Analysis models. The proposed method is expected to yield better accuracy, precision, and F1-score.

**Keywords** – *Quadratic Discriminant Analysis and Linear Discriminant Analysis, Web service security, Anomaly detection, Machine Learning.*

## I. INTRODUCTION

E-Commerce security ensures customer data are stored safely and all transactions are carried out without any compromise in data security. Cyberattacks on e-commerce may come in many forms and proper preventive measures are to be taken to ensure data protection. Some common types of attacks are phishing, monetary theft, credit card fraud, hacking, and misusing of intellectual property.

Any e-commerce platform ensures proper protection against these threats and prioritizes security measures for its customers. The machine learning mechanism is used to restrict anomaly access.

Analyzing datasets using machine learning can help in processing huge datasets and detect patterns that could be used to isolate attacks and malicious usages. For example, when a large number of requests is originating from a single user, the algorithm can detect this abnormal activity and can notify the admin regarding a potential attack.

This study discussed the train dataset noises and how dataset tuning improves the classifier model accuracy. Restricting the access of anomaly requests from various sources is also discussed.

The proposed algorithm applied to Quadratic Discriminant Analysis and Linear Discriminant Analysis models. It predicts anomalies with better precision, accuracy and F1 score.

## II. RELATED WORK

In computer science, anomaly detection refers to the techniques of finding specific data points that do not conform to the normal distribution of the data set. Companies from different sectors including manufacturing, automotive, healthcare, lodging, travelling, fashion, food, and logistics are investing a lot of resources in collecting big data and exploring the hidden anomalous patterns in them to facilitate their customers. In most of the cases, the collected data are streaming time series data and due to their intrinsic characteristics of periodicity, trend, seasonality, and irregularity, it is a challenging problem to detect point anomalies precisely in them (Mohsin Munir 2018).

Anomaly-based intrusion detection systems (IDSs) have been deployed to monitor network

activity and to protect systems and the Internet of Things (IoT) devices from attacks (or intrusions). The problem with these systems is that they generate a huge amount of inappropriate false alarms whenever abnormal activities are detected and they are not too flexible for a complex environment. The high-level rate of the generated false alarms reduces the performance of IDS against cyber-attacks and makes the task of the security analyst particularly difficult and the management of intrusions detection process computational expensive (Wajdi Alhakami 2019).

Naïve Bayes is a simple technique for classification. Naïve Bayes model could be used without accepting Bayesian probability or using any Bayesian methods. An analysis of the Bayesian classification problem showed that there are sound theoretical reasons behind the apparently implausible efficacy of types of classifiers. A advantage of Naïve Bayes is that it only requires a small number of training data to estimate the parameters necessary for classification. The fundamental property of Naïve Bayes is that it works on discrete value. If attribute values are continuous it is recommended to use Gaussian Naïve Bayes classifier (Shikha Agarwal 2019).

Machine learning methods have been widely used in the intrusion detection field. Classification, clustering, Markov chains have been studied extensively on classical systems. In modern systems, there are many CPUs and they are shared amongst software provided by the kernel scheduler. If there is more than usual demand for CPU resources, the tasks create an order and are in standby mode for processing. Standby regimes slow down the execution time of tasks, which results in the reduction of performance metrics. To improve performance metrics CPU usage needs to be analysed. In most cases, CPU usage is analysed in terms of process, flow or task. Another metric to be analysed to improve the performance metrics is the CPU utilisation. CPU utilisation is measured in time when the CPU is engaged in the processing interval of the

task and shown in the percentage. Memory access attempts also cause high CPU usage. When the interrupt attempts are made, the CPU interrupts to work and waits for the process to complete (Rasim M. Alguliyev 2019).

Quadratic discriminant analysis (QDA) is a widely used classification technique that generalizes the linear discriminant analysis (LDA) classifier to the case of distinct covariance matrices among classes. For the QDA classifier to yield high classification performance, an accurate estimation of the covariance matrices is required. Such a task becomes all the more challenging in high-dimensional settings, wherein the number of observations is comparable with the feature dimension (Houssem Sifaou 2020).

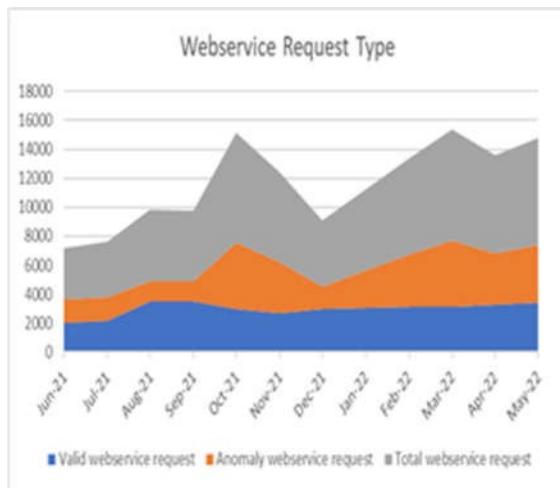
The performance of LDA-based classifiers depends heavily on accurate estimation of the class statistics, namely, the sample covariance matrix and class mean vectors. These statistics can be estimated with fairly high accuracy when the number of available samples is large compared to the data dimensionality. In practical high-dimensional data settings, the challenge is to cope with a limited number of available samples. In this case, the sample covariance estimates become highly perturbed and ill-conditioned resulting in severe performance degradation. In some practical situations, it occurs that the test data deviates from the training data model. For example, the training data and the test data might represent measurements obtained from non-identical devices. In such a case, the value of the regularisation parameter computed during the training phase may no longer be adequate, let alone be optimal (Alam Zaib 2021).

### III. PRELIMINARIES

JSON parameters are most commonly used to send and receive data in web services. JSON allows the transfer of complex data structures efficiently, which enables the transfer of large amounts of data with minimal effort and resources.

On analyzing the above survey, we find that around 48.58% of the web service requests come from valid users and the remaining 51.42% are found to be anomalous. Most of the requests come from invalid users, to protect against these potentially harmful requests. This proposed framework has been designed to detect and identify suspicious access patterns and automatically prevent unwanted web requests. It works by continuously monitoring the user activities and denies access for any anomalous request. It also provides a better and easier method to update rules and security policies to ensure the safety of the e-commerce web service. It also has the option to customize security settings and fine-tune control lists and restrict user activities, etc. This paper discusses a secure environment for operating an e-commerce business.

The e-commerce dataset used to train a machine learning model to find anomalies. The dataset was prepared from API requests from various sources, ranging from 50 to 200. Each API had different parameters. Our proposed method was applied before training the Quadratic Discriminant Analysis and Linear Discriminant Analysis models. The below API retrieves product details.



**Figure 1 Chart for Month wise Count of Webservices Request Type**

Dataset for get product details

```
{
  "Time": "2022-12-0600:00",
  "API": "resolve_product_detail",
  "data": [
    {
      "user_id": "retsSDWJwdw",
      "product": "Dolo 700",
      "slug": "dolo-700mg"
    }
  ],
  "ID": "frRFEwdswrd",
  "IP": "172.134.04.99",
  "Country": "India",
  "OperatingSystem": "Linux"
}
```

### 3.1 Quadratic Discriminant Analysis

Quadratic Discriminant Analysis is derived from the linear discriminant analysis. QDA classifies the result from two or more groups of dataset. It uses a quadratic function to classify the categories. QDA observes the difference of mean and covariance of each category. QDA provides better accuracy than LDA, with different covariance in the class.

$$\delta_k(x) = -\frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) + \log \pi_k$$

The Classification rule is

$$\hat{G}(x) = \arg \max_k \delta_k(x)$$

QDA works effectively to classify the boundaries of non-linear classes.

### 3.2 Linear Discriminant Analysis

Linear Discriminant Analysis uses the technique of data reduction. LDA eliminates redundant data from the dataset and reduces the dimensionality of the dataset. The reduced dimension dataset enhances the between-class variance.

$$P_{lda} = \arg \max_P \frac{|P^T S_b P|}{|P^T S_w P|}$$

## IV. THE PROPOSED APPROACH

The proposed method tunes the businessdatasetforachievingbetteraccuracy,prediction, and F1 score. The proposed method improves the quality of the dataset. The outcome of this work gives better Quadratic Discriminant Analysis and Linear Discriminant Analysis results. This study proposes the following method to improve the dataset quality.

- Accept requests only from the e-commerce service area.
- Provided distinct user id to every guest user.
- Replace null user id by default user id.
- Replace null geolocation by default geolocation.

### Algorithm1

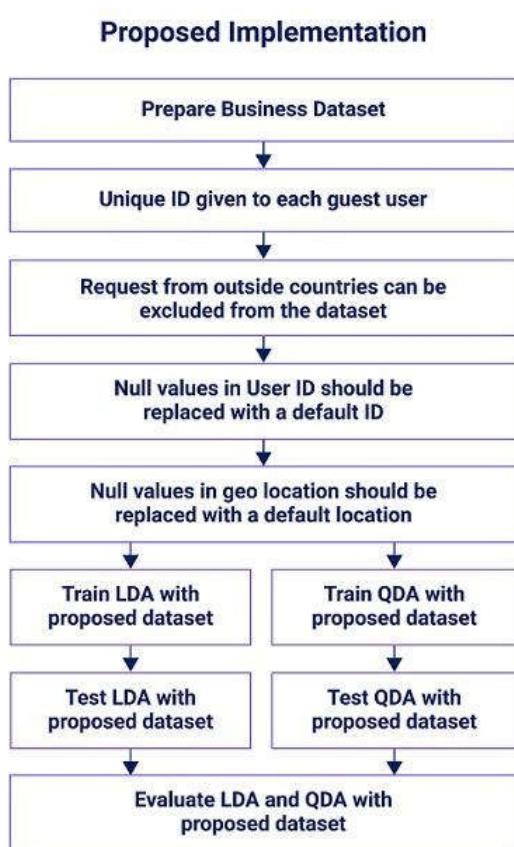
#### Training dataset algorithm

```
#Preparedataset
D={d1,d2,d3...dn}
#Acceptrequestsfromthee-
commerceservicearea
G=[xforxinDifx[geo_location]=country]
#Provideeveryguestuserwithadistinctuserid
I=[xforxinGifx[user_id]isguest]x[use
r_id]=uniquevalue
#ReplacenulluseridbydefaultuseridN =
[x for x in G if x[user_id] is
null]x[user_id]=defaultvalue
#Replacenullgeolocationbydefaultgeolocatio
n
L=[xforxinGifx[geo_location]isnull]x[ge
o_location]=defaultvalue
#Createadatasetwiththeabovecriteria
DS=[I ∩ N ∩ L]
#Splittingthedata
DS={x∈DS|80%astrainingdatasetand20%aste
ringdataset}
#Building the LDA
modelLDA=build_Ldamodel
(DS)
```

```
#Building the QDA
modelQDA=build_Qdamodel
(DS)
#Predicting the anomaly by
LDALDA_prediction =
predict_anomaly(LDA,DS)
#Predicting the anomaly by
QDAQDA_prediction=
predict_anomaly(QDA,DS)
#Evaluatetheresult
CompareperformanceofLDAandQDAwithProp
oseddataset
```

#### 4.1 Requests from countries outside of the E-Commerce service area is flagged

One easy but effective technique used to reduce unwanted traffic in an e-commerce platform is by restricting access to the requests that originate from outside its geographic serviceable area. This can ensure the safety of the platform from cyber attacks from other countries and also helps in reducing the traffic load of the website. This is typically done by identifying the origin of the IP address of the request and denying access to such a request that originates from outside the serviceable geo-location. In our case, all requests originating from outside India are restricted access, as our e-commerce service area is only in India.



**Figure 2 Proposed AlgorithmImplementation**

#### 4.2 Provide every guest user with a distinctuserid

An e-commerce website receives different types of requests from different origins, and requests from users are among them. It may be from a registered user or a guest user, guest users are the types of users who are accessing the e-commerce platform but have not registered with it. They can access certain parts of the website, but some parts are restricted for them. It is important to track these users and their activities on the website. It is also important to include these users' data with the registered user's data so that the user behavior in the platform can be better understood. To track the guest user behaviors, they need to be identified individually. This can be achieved by providing the guest users with unique ids when they visit the first-time, which is known as the guestuserid. A unique guestuserid is provided

to each guest user so that all the guest user activities can be tracked.

#### 4.3 Replace nulluserid by defaultuserid

Null user IDs are the data in the dataset that does not have an ID for the user but all other information is present. This typically happens when a user id has not been assigned to a user or if the user id is lost when clearing the cache of the user's device. In any case, if we omit these data from the dataset, any analysis done on the dataset may become unreliable and inaccurate. So, to include these data in the analysis we are providing these null user ids with the default user id. This enables the null user id to be included in the dataset and return, it helps in including all the user records and behaviors are included in the dataset. Also, this ensures null user ids get valid id to be included in tracking, reporting, and other analysis. Adding the default user id to the null user id helps in increasing the accuracy and comprehensiveness of the dataset.

#### 4.4 Replace null geolocation with defaultgeolocation

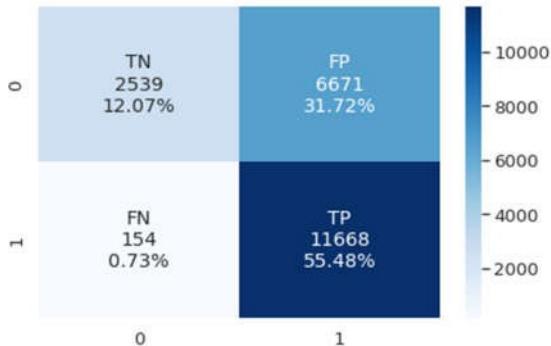
Null geo-location is when user data in a dataset has no geo-location in it. This typically happens when a location entered by a user is wrong or not identified or if the GPS coordinates of the user are not available or if the IP address cannot be geo-located. The null geolocation has to be changed to the default geolocation to include the particular user data to be included in the dataset. By replacing the null geolocation with the default geo-location, we are assigning an approximate location for the user for whom the true geo-location cannot be identified. This technique enables the system to complete the user data and helps in increasing the accuracy when analyzing the dataset. In our case, we use the default location as India.

## V. EXPERIMENTAL RESULT

The proposed dataset is tested with Quadratic Discriminant Analysis and Linear Discriminant Analysis classifiers. The experimental result is derived in the below section.

### 5.1 QuadraticDiscriminantAnalysiswithproposedalgorithm

The datashows the Confusion Matrix of QuadraticDiscriminantAnalysis.



**Figure3ConfusionMatrixofQuadraticDiscriminantAnalysiswithproposedalgorithm**

Theproposedalgorithmperformance metricsarecalculatedfromTruePositive, TrueNegative, FalsePositiveandFalseNegative.

$$\text{TruePositive(TP)}=11668$$

$$\text{TrueNegative(TN)}=2539$$

$$\text{False Positive (FP)} =$$

$$6671$$

$$\text{FalseNegative(FN)}=154$$

54

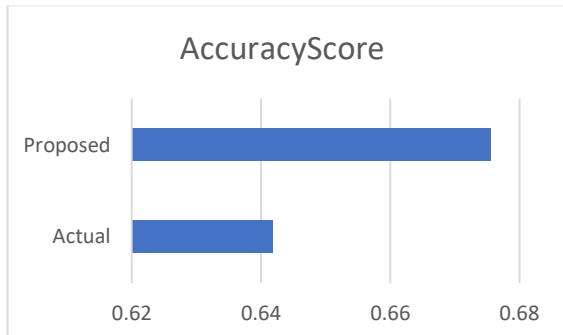
#### 5.1Accuracy

The accuracy is the percentage of correctpredictionsoutofthetotalpredictions.

$$\text{Accuracy} = \frac{(TP+TN)}{(TP+TN+FP+FN)}$$

$$\begin{aligned} \text{Accuracy} &= \frac{(11668+2539)}{(11668+2539+6671+154)} \\ &= \end{aligned}$$

$$\text{Accuracy}=0.6754944846$$



**Figure4AccuracyscoreofQDA**

TheAccuracyscorehasimprovedby5.24%.

### 5.2 Precision

The precision is the metric that defines themodel's ability to predict the correct number ofyesoutcomesoutofthetotalyesoutcomespredicted.

$$\text{Precision} = \frac{TP}{(TP+FP)}$$

$$\text{Precision} = \frac{11668}{(11668+6671)}$$

$$\text{Precision}=0.6362397077$$



**Figure5PrecisionscoreofQDA**

ThePrecisionscorehasimprovedby12.53%.

### 5.3 Recall

Therecallcomparsthemodel'snumberof

correctyesoutcomespredictedwiththetotalnumbertofofactualyesoutcomes.

$$\text{Recall} = \frac{TP}{(TP+FN)}$$

$$\text{Recall} = \frac{11668}{(11668+154)} = 0.9869734394$$

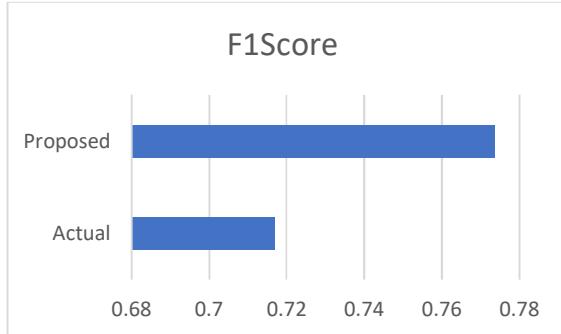
### 5.4 F1Score

F1 Score combines both precision and recallsscoresforfoolproofevaluation.

$$\text{F1Score} = \frac{2x(\text{Precision score} \times \text{Recall score})}{(\text{Precision score} + \text{Recall score})}$$

$$F1Score = \frac{2x(0.6362397077x0.9869734394)}{(0.6362397077+0.9869734394)}$$

$$F1Score=0.7737143994$$



**Figure6F1ScoreofQDA**

The F1 score has improved by 7.91%.

### 5.5 MissRate

The miss rate measures the incorrect predictions. The positive result that is predicted as negative is known as a false negative.

$$MissRate = \frac{FN}{FN + TP}$$

$$MissRate = \frac{154}{154+11668}$$

$$MissRate=0.01302656065$$

### 5.6 Fall-out

The fall-out rate measures the incorrect predictions. The negative results predicted as positive are known as a false positive.

$$Fall-out = \frac{FP}{FP+TN}$$

$$Fall-out = \frac{6671}{6671+2539}$$

$$Fall-out=0.7243213898$$

### 5.7 Specificity

Specificity is the metric used to find the true negatives; it is also referred to as the True Negative Rate.

$$Specificity = \frac{TN}{TN+FP}$$

$$Specificity = \frac{2539}{2539+6671}$$

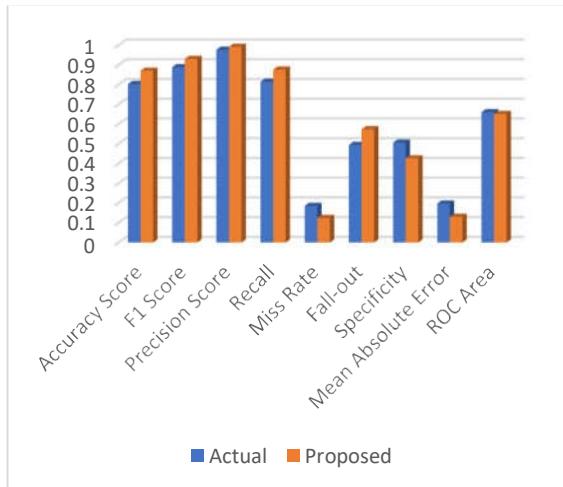
$$Specificity=0.2756786102$$

### 5.8 MeanAbsoluteError

The Mean absolute error value is 32.45% for the proposed algorithm. It is reduced by 3.4%. Quadratic Discriminant Analysis is a classifier model tested with the proposed algorithm. The results, shown in the below table, were compared using accuracy score, F1 score, precision score, recall, miss rate, fall-out, specificity and mean absolute error. The results demonstrate that the proposed algorithm is improving the anomaly detection performance.

**Table1ResultComparisonofQuadraticDiscriminantAnalysiswithproposedalgorithm**

Metrics	QuadraticDiscriminantAnalysis	Proposed algorithm
AccuracyScore	64.18%	67.55%
F1Score	71.69%	77.37%
PrecisionScore	56.54%	63.62%
Recall	97.96%	98.70%
MissRate	2.04%	1.30%
Fall-out	64.95%	72.43%
Specificity	35.05%	27.57%
Mean Absolute Error	35.82%	32.45%



**Figure 7 Performance Metrics of Quadratic Discriminant Analysis with proposed algorithm**

## 5.2 Linear Discriminant Analysis with proposed algorithm

The proposed algorithm performance metrics are calculated from True Positive, True Negative, False Positive and False Negative.

$$\text{True Positive (TP)} =$$

$$18187 \text{ True Negative (TN)}$$

$$= 113 \text{ False Positive (FP)}$$

$$= 152 \text{ False Negative (FN)}$$

$$= 80$$

The data shows Confusion Matrix of Linear Discriminant Analysis.



**Figure 8 Confusion Matrix of Linear Discriminant Analysis with proposed algorithm**

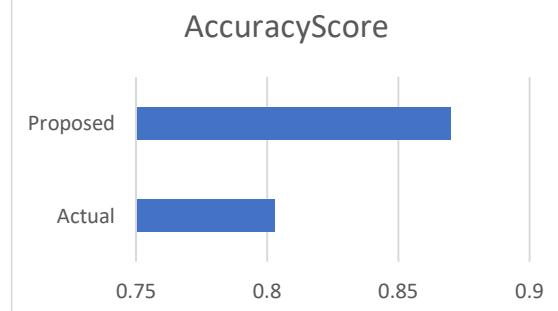
### 5.1 Accuracy

The accuracy is the percentage of correct predictions out of the total predictions.

$$\text{Accuracy} = \frac{(TP+TN)}{(TP+TN+FP+FN)}$$

$$\text{Accuracy} = \frac{(18187+113)}{(18187+113+152+2580)}$$

$$\text{Accuracy} = 0.870103$$



**Figure 9 Accuracy score of LDA**

The Accuracy score has improved by 8.37%.

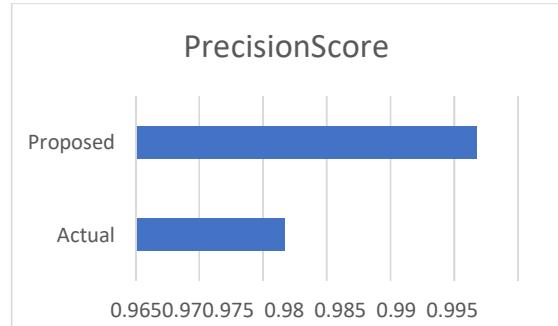
### 5.2 Precision

The precision is the metric that defines the model's ability to predict the correct number

of yes outcomes out of the total yes outcomes predicted.

$$\text{Precision} = \frac{TP}{(TP+FP)}$$

$$\text{Precision} = \frac{18187}{(18187+152)} = 0.991712$$



**Figure 10 Precision score of LDA**

The Precision score has improved by 1.54%.

### 5.3 Recall

The recall compares the model's number of correct yes outcomes predicted with the total number of actual yes outcomes.

$$TP$$

$$Recall = \frac{TP}{TP+FN}$$

$$Recall = \frac{18187}{(18187+2580)} = 0.875764434$$

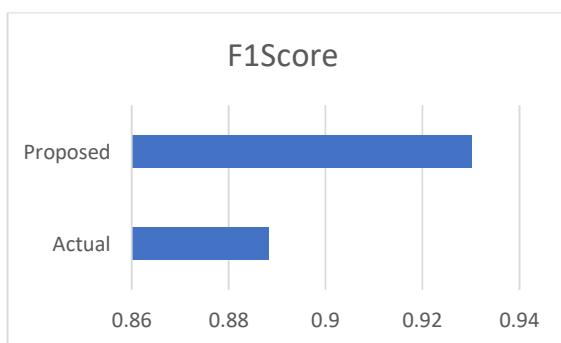
### 5.4 F1 Score

F1 Score combines both precision and recall scores for foolproof evaluation.

$$F1Score = \frac{2x(precision score \times recall score)}{(precision score + recall score)}$$

$$F1Score = \frac{2x(0.991712 \times 0.875764434)}{(0.991712 + 0.875764434)}$$

$$F1Score = 0.9301385977$$



**Figure 11 F1 Score of LDA**

The F1 score has improved by 4.71%.

### 5.5 Miss Rate

The miss rate measures the incorrect predictions. The positive result that is predicted as negative is known as a false negative.

$$MissRate = \frac{FN}{FN+TP}$$

$$MissRate = \frac{2580}{2580+18187}$$

$$MissRate = 0.124235566$$

### 5.6 Fall-out

The fallout rate measures the incorrect predictions. The negative results predicted as positive are known as a false positive.

$$Fall-out = \frac{FP}{FP+TN}$$

$$Fall-out = \frac{152}{152+113}$$

$$Fall-out = 0.5735849057$$

### 5.7 Specificity

Specificity is the metric used to find the true negatives; it is also referred to as the True Negative Rate.

$$Specificity = \frac{TN}{TN+FP}$$

$$Specificity = \frac{113}{113+152}$$

$$Specificity = 0.4264150943$$

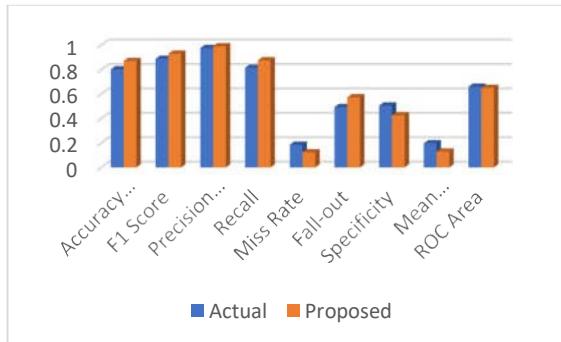
### 5.8 Mean Absolute Error

The Mean absolute error value is 12.99% for the proposed algorithm. It is reduced by 6.73%.

**Table 2 Result Comparison of Linear Discriminant Analysis with proposed algorithm**

Metrics	Linear Discriminant Analysis	Proposed algorithm
AccuracyScore	80.28%	87.01%
F1Score	88.83%	93.01%
PrecisionScore	97.67%	99.17%
Recall	81.45%	87.58%
MissRate	18.55%	12.42%
Fall-out	49.39%	57.36%
Specificity	50.61%	42.64%
MeanAbsoluteError	19.72%	12.99%

Linear Discriminant Analysis is a classifier model tested with the proposed algorithm. The results, shown in the above table, were comparing using accuracy score, F1 score, precision score, recall, miss rate, fall-out, specificity and mean absolute error. The result demonstrates that the proposed algorithm is improving the anomaly detection performance.



**Figure 12 Performance Metrics of Linear Discriminant Analysis with proposed algorithm**

## VI CONCLUSION

The proposed algorithm tested with Quadratic Discriminant Analysis and Linear Discriminant Analysis models. It improves Accuracy score, Precision score and F1 score.

Quadratic Discriminant Analysis with proposed algorithm. The Accuracy score has improved by 5.24%. The Precision score has improved by 12.53%. The F1 score has improved by 7.91%.

Linear Discriminant Analysis with proposed algorithm. The Accuracy score has improved by 8.37%. The Precision score has improved by 1.54%. The F1 score has improved by 4.71%.

## VII FUTURE WORK

The proposed algorithm can be used with different classifiers. They are Gaussian Naive Bayes, Logistic Regression, MLP Classifier, Adaboost Classifier, Random Forest Classifier, Decision Tree Classifier, and K-Nearest Neighbors.

Train the models with different domain datasets. The wide variety of training datasets predicts with better accuracy. The dataset can be retrieved from different seasons and different geographical location people data.

## REFERENCES

- [1] BayuAdhiTama,LewisNkenyereye,S. M.RiazulIslam, and Kyung-SupKwak(2020). An Enhanced Anomaly Detection in WebTraffic Using a Stack of Classifier Ensemble. IEEE Access, Digital Object Identifier 10.1109/ACCESS.2020.2969428, pages 24120–24134.
- [2] WajdiAlhakami, Abdullah Alharbi, SamiBourouis, RoobaeaAlroobaea, and Nizar Bouguila(2019). Network Anomaly Intrusion Detection Using a Nonparametric Bayesian Approach and Feature Selection. IEEE Access, Digital Object Identifier 10.1109/ACCESS.2019.2912115, pages 52181–52190.
- [3] HaowenXu, Wenxiao Chen, NengwenZhao, Zeyan Li, Jiahao Bu, Zhihan Li, YingLiu, Youjian Zhao, Dan Pei, Yang Feng, JieChen, ZhaogangWang, HonglinQiao(2018). Unsupervised Anomaly Detection via Variational Auto-Encoder for Seasonal KPIs in Web Applications, Cornell University.
- [4] Ya Su, YoujianZhao, ChenhaoNiu, RongLiu, WeiSun, DanPei(2019). Robust Anomaly Detection for Multivariate Time Series through Stochastic Recurrent Neural Network. Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining
- [5] VenkatN.Gudivada, AmyApon, and Junhu aDing(2017). Data Quality Considerations for Big Data and Machine Learning: Going Beyond Data Cleaning and Transformations. International Journal on Advances in Software 10.1, pages 1-20.
- [6] AlamZaib, TarigBallal, ShahidKhattak, and TareqY.Al-Naffouri(2021), A Doubly

RegularizedLinearDiscriminantAnalysisClassifierWithAutomaticParameterSelection.IEEEAcess,DigitalObjectIdentifier10.1109/ACCESS.2021.3068611, pages51343–51354.

[7] HoussemSifaou, AblaKammoun, and Mohamed-SlimAlouini(2020), High-Dimensional Quadratic Discriminant AnalysisUnderSpikedCovarianceModel.IEEE Access,DigitalObjectIdentifier10.1109/ACCESS.2020.3004812, pages117313–117323.

[8] Chen Zhuang, Hongbo Zhao, Chao Sun, and Wenquan Feng(2019), Detection andClassificationofGNSSSignalDistortionsBased on Quadratic Discriminant Analysis.IEEEAccess,DigitalObjectIdentifier 10.1109/ACCESS.2020.2965617, pages25221–25236.

[9] V. Sankar, Dr. G. Zayaraz(2020). JSONandXMLdatasecurityperformanceofEncryptedFileProcessinWebServices. Published in International Journal of CurrentScience(IJCPUB), Volume10,Issue4 .

[10] Ravi Chandra Jammalamadaka, SharadMehrotra (2006). Querying Encrypted XMLDocuments. Published International DatabaseEngineeringandApplicationsSymposium(IDEAS'06)

[11] RA.K.Saravanaguru,GeorgeAbraham,Krishnakumar Venkatasubramanian,KiransinhBorasia(2013). SecuringWebServicesUsingXMLSignatureandXMLEncryption. School of Computer Science andEngineering,VITUniversity,Vellore,India.

[12] Hoi Ting Poon and Ali Miri(2015). ComputationandSearchoverEncryptedXML Documents. Department of ComputerScienceRyersonUniversityToronto, Ontario, Canada,978-1-4673-7278-7/15

[13] GuYue-sheng, Ye Meng-tao, Gan Yong(2010). WebServicesSecurityBasedon

XML Signature and XML Encryption.Journalofnetworks,Vol.5, No.9.

[14] Nithin N and Anupkumar M Bongale(2012). XBMRSA:A New XML EncryptionAlgorithm. Proceedings of Information andCommunicationTechnologies.WorldCongress,pp567-571,2012.

[15] NithinNandHarshitha.K.S(2014). Analysis of Symmetric algorithm for XMLdocument security. International Journal ofInnovations in Engineering and Technology(IJIET), Vol.3Issue4.

[16] AamerNadeem (2005). A PerformanceComparison of Data Encryption Algorithms.IEEE.

[17] RaviVarma1,Dr.G.VenkatRamiReddy (2014). Schema Based Parallel XMLParser:AFastXMLParserDesignedforLarge XML Files. International Journal ofComputerScience andMobile Computing. Vol.3Issue.8.

[18] Shikha Agarwal, BalmukundJha, TisuKumar,ManishKumar,PrabhatRanjan(2019). Hybridof NaiveBayesandGaussian Naive Bayes for Classification: A Map Reduce Approach. International JournalofInnovativeTechnologyandExploringEngineering(IJITEE)ISSN:2278-3075,Volume-8,Issue-6S3.

[19] RasimM.Alguliyev,RamizM.Aliguliyev and FarganaJabbarAbdullayeva(2019). Hybridisat ionofclassifiersforanomaly detection in big data. InternationalJournalofBigDataIntelligence.10 .1504/IJBDI.2019.10018528