

# Leveraging Feature Ranking for System Fault Identification and Classification Using Machine Learning Algorithms



K. Pradheep Arumuham and B. Booba

**Abstract** Computing environments necessitate the prerequisite of ensuring performance pinnacle and entailing resilience during disruptions. Most often, system faults can interrupt the overall productivity and prolong operational time, thereby increasing computational complexities, along with leading to downtime, mitigated productivity, and detriments in terms of finances. The inevitability to swiftly effectuate precise system failure detection, stratification and further resolution, becomes crucial for effectively maintaining system veracity and evade unethical injections. This indagation pivots on analyzing the various attributes relevant to system fault processing, and to entail the data thresholding combined with feature extraction to efficiently identify and classify system failures using machine learning algorithms. The proposed study entails a multi-modal real-time feature evaluation from the database constructed using the primary attributes such as the upstream connection, response time, API latency, connection time, and transaction status. Feature ranking using variance, Region of Curve (ROC) and *T*-Test are incorporated to enhance accuracy of classification. The machine learning algorithms used in this research paper are the Efficient Linear Support Vector Machine (EL SVM), Naïve Bayes algorithm and Tri-layered Neural Network, and the performance accuracy rendered by each of the algorithms are scrutinized. The simulation results are carried out in MATLAB, and the results are procured successfully.

**Keywords** System fault identification and classification · Feature ranking · API latency · Efficient linear SVM · Naïve Bayes · Tri-layered neural network · MATLAB

---

K. P. Arumuham (✉)  
Research Scholar, VISTAS, stitch.sa, Chennai, India  
e-mail: [pradheepz@gmail.com](mailto:pradheepz@gmail.com)

Engineering, stitch.sa, Chennai, India

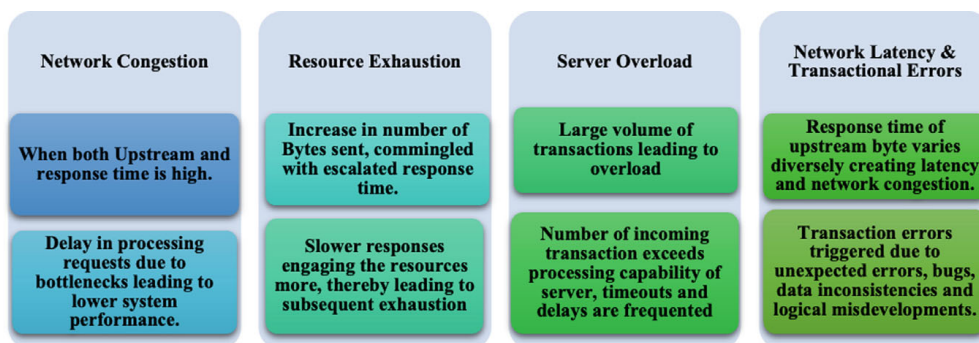
B. Booba  
Professor, VISTAS, Chennai, India

## 1 Introduction

The technological realm drives the need for paramount sustainability of systems through resilient working environments that further necessitate the explicit classification of system failures. System failures [1] can be triggered from diverse dysfunctions, malwares and network problems that can lead to compromised interruptions, financial constraints, and unsatisfied customer experiences. The prerequisite for timely and accurate classification of system failures is crucial for effective troubleshooting and proactive maintenance, thereby minimizing downtime and optimizing system performance. The use of machine learning algorithms [2] integrates the stratification pipeline to effectively leverage the right features procured from thresholding and feature ranking. Adapting the technique to monitor systems through the procured feature data, the distinct patterns indicative of system failures can be unsheathed. Most often system failures occur due to diverse reasons generally categorized into the following as shown in Fig. 1.

By monitoring and analyzing corresponding features, proactive steps in failure detection [3] can be constructed in order to substantially mitigate risks [4]. Some steps entail optimizing network configurations, scaling resources to handle increased loads, implementing caching mechanisms to reduce latency and corroborating development of logic with appropriate testing methods [5]. This research pivots on different machine learning classifiers that aid in training and testing the complex relationships between system metrics and the respective cause of failure. The key contributions of this research include the thresholding of the data to comprehend the success and failures of the system, while corroborating the thresholding result with the entailed tri-classifier model. The real-time data is procured to assess the performance of the systems, while cognizing the need for diversity of features to enhance the precision of the stratification.

This paper is structured with Sect. 2 elaborating the literature review in relevance to identification of system failures. Section 3 explicates the methodology used in this indagation, followed by the results and conclusion in Sects. 4 and 5, respectively.



**Fig. 1** Common triggers for system failures

## 2 Empirical Review

Waqar et al. [1], in the paper titled “Fault identification, classification, and localization in microgrids using superimposed components and Wigner distribution function” delineated about the resilience of microgrids, and the integration of Distributed Energy Resources (DERs) into disseminated grids that has become less compatible with the evolution of the former. The protection of microgrids due to its popular utilization is an issue that is addressed in this study. The study explicates the use of superimposed components and the Wigner distribution function (WDF). The protection methodology utilizes the alienation coefficient and WDF to effectively detect faults in the microgrid. The introduction of superimposed positive sequence reactive power (SPSQ) aids in stratifying and identifying detriments in both grid-based and islanded operable platforms thereby enhancing reliability and augmented security of microgrids.

Stravani et al. [6] explicated the encountering of faults in the orifice flowmeter for detection of flow in a resilient manner. Transient data procured from computational fluid dynamics serves as the base for this indagation. The proposed study entails a model-based approach to construct a second-order transfer function that contrives residue for detection of faults in the flowmeter. The Luenberger Observer for Residue Generation otherwise called as the LPV observer is juxtaposed with the working of a neural network to identify the performance of the systems. While the time of fault detection is equaled by both models, the real-time fault data is better identified in LPV observers as negated by the neural network models. The inclusion of entropy is also better handled by the LPV observers, with the neural network model showcasing dissemination in accuracy when noise in data is entailed.

Zhang et al. [7], proposed a study that involved the diagnosis of fault [8] in electronic devices due to external faults through the utilization of power electronic converter. The study pivots on signal injections that identifies the impact by line inference signals. This thereby explicitly locates line faults through harmonic impedance that further computes the distance metrics for unambiguous fault identification.

## 3 Proposed Methodology

This research proposes the classification of successful and unsuccessful systems in terms of their resilience to failure. The dataset incorporated for this study consists of real-time data procured from working environments from stitch.sa, with a total of 65,535 data consisting of various system attributes. System testing using the proposed algorithmic models are tested in stitch.sa, thereby further enhancing the accuracy of system fault classification and substantially bolstering system uptime in a pragmatic level of utilization. The training and testing set are stratified as 70% and 30%, respectively. Pre-processing of data is implemented by cleaning of outliers, and entropy analysis. Z-Score median absolute deviation is used for normalizing

the data to enable better accuracy in the further processes. Feature extraction is effectuated after effective pre-processing, and ranking of the attributes renders better comprehension for initiating classification. The classification learner in MATLAB is used for simulating classifier accuracy efficacy for the incorporated study. The overall architecture for this indagation is depicted in Fig. 2 as illustrated:

The process of data pre-processing entails various phases as shown in figure (Fig. 3).

The missing value imputation for this study is effectuated through K-Nearest Neighbor imputation [9]. The distance is chosen to be the Euclidean measure. The “knnimpute” function in MATLAB is used for filling the missing values through the process of identifying the mean of the specific attribute. Post the process of missing value imputation, the next step of outlier removal is implemented by identifying the deviations from the mean. This process entails a threshold value, beyond which the values are chosen as outliers. In this study, the threshold value is set to 3. The computation of outliers also entails the calculation of standard deviation in order to comprehend the absolute deviation. Therefore, the below formula is used to identify the outliers using the standard deviation and mean as shown in equation.

$$\text{Oulier} = \text{abs}(\text{attribute value} - \text{attribute mean}) > (\text{threshold} * \text{std}_{\text{data}}) \quad (1)$$

The Z-score normalization using Median Absolute Deviation (MAD) [10] is resilient to outliers as contrasted to the conventional methods. The process of normalization is implemented by subtracting the median of each point and further dividing

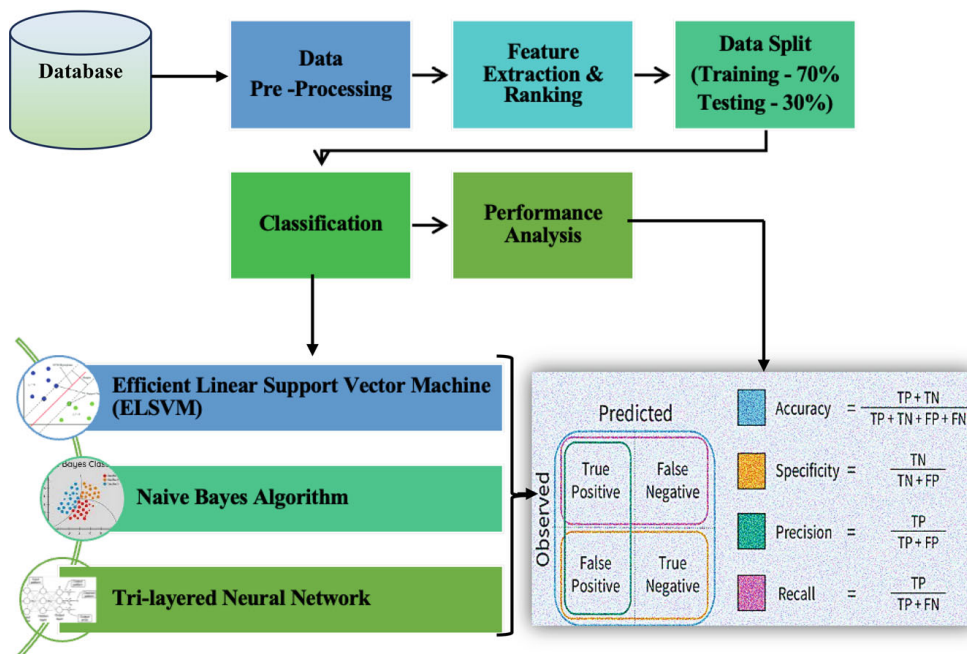


Fig. 2 Proposed architectural flow



**Fig. 3** Steps implemented for pre-processing data

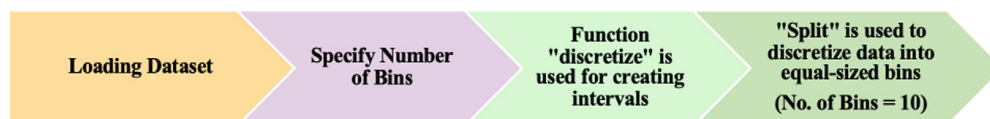
the points with their corresponding median absolute deviation. The results of the Z-score Median Absolute Deviation implemented on the Fault Tolerance dataset is shown in Sect. 4.

Subsequent to the process of normalization, binning using the split algorithm is performed in order to discretize the data, augment classifier accuracy through the analysis of bins formed for each attribute. The sequence of steps implemented in binning is as shown in Fig. 4.

The process of feature extraction and ranking is used to ensure that the classification of data renders high precision while effectively scrutinizing the attributes that render pinnacle of information as compared to those that do not. The feature extraction and ranking [11] are effectuated through the Region of Curve (ROC), variance of data loaded and values sorted post the  $P$ -value for explicitly procuring the  $T$ -Test results.

Classification of data is initiated by splitting the dataset into 70% for training and 30% for testing. Three different classifier methods are implemented to analyze the performance of the models, and further to comprehend the efficacy derived from the initial pre-processing phases.

The Efficient Linear Support Vector Machine (ELSVM) [12] is used as a classification algorithm [13] with the selection of cross-validation to effectively utilize the generated feature matrices and label vectors. This model is implemented through the



**Fig. 4** Binning using median absolute deviation

“fitsvm” function, with the kernel selection set to indicate linearity. The regularization parameter [14] is adjusted to ridge with automated lambda ridge strengthening. The convergence criteria are constrained by utilizing the beta tolerance levels, which is set to 0.0011 for the incorporated dataset. One-to-one multi-class coding is used to comprehend the binary classifier [15] model for complex datasets entailing large quantity of data. The training and testing of the model are effectuated to explicitly comprehend the classification accuracy, and the results procured from the algorithm is delineated in the subsequent section.

The Naïve Bayes algorithm [16] works on the basis of probabilistic classification with the combination of Bayes theorem and naïve assumption. This algorithm is categorized into two different distributions to fit the ensemble data of the “FaultTolerance” dataset that comprises of both numerical and categorical data. The Gaussian function is utilized as the distribution function for numeric predictors, and the probabilistic function of the same is as follows:

$$\text{GDf}(x_i|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} * e^{-\frac{(x_i-\mu)^2}{2\sigma^2}} \quad (2)$$

where  $\mu$  is the mean,  $\sigma^2$  is the variance computed for each feature in the dataset.

The multivariate normal distribution is used for categorical predictors in the dataset with equation:

$$\text{MVf} = (x|\mu, \Sigma) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} * e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)} \quad (3)$$

where  $x$  is the feature vector,  $\mu$  is the mean,  $\Sigma$  is the covariance matrix computed for each feature in the dataset, thereby identifying the likelihood of the data to further categorize them into classes. The results procured from this algorithm renders the highest accuracy performance for the incorporated dataset due to the ensemble distribution method used for the categories of data in the “FaultTolerance” dataset.

The tri-layered neural network [17] model is built with three fully-connected layers, with effective standardization performed on the data. The bias for layers is randomized, with the weights for each layer set to 30 for the first layer, 25 for the second, and the third layer comprising of 15. The Sigmoid activation function is used to diminish non-linearity in the network. The iteration limit is set to 1000, with the regularization strength equalized to 0.1. The classifier utilizes both forward and backpropagation to compute loss and update weights and bias for the model, respectively. The results of the model obtained after training and testing is depicted in the consecutive section.

## 4 Results

The results procured from the simulations are depicted as below. Figures 5, 6, 7 and 8 explicate the method of pre-processing implemented after explicit data visualization. Figure 5 delineates the cleaning process through outlier removal using mean detection method. Further to cleaning, the normalization through Z-score method is effectuated and is represented in Fig. 6. Figures 7 and 8 show the process of binning implemented through the splitting algorithm for various attributes to subsequently enhance classifier accuracy.

Figure 9 shows the process of feature extraction and ranking depicted through three methods such as ROC, variance, and *T*-Test. The feature “http\_status” renders the highest information and forms to be crucial feature for stratification accuracy enhancement.

Figures 10, 11 and 12 depict the various classification algorithms implemented in this research paper. The training and validation accuracies, along with the parallel

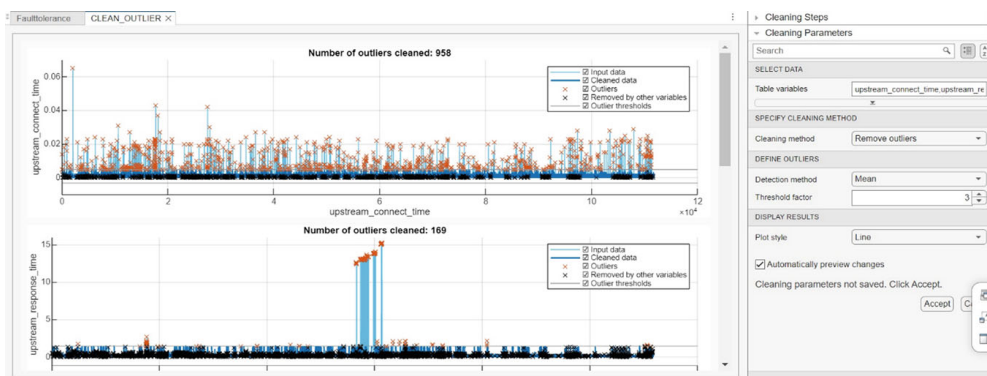


Fig. 5 Outlier removal for each attribute

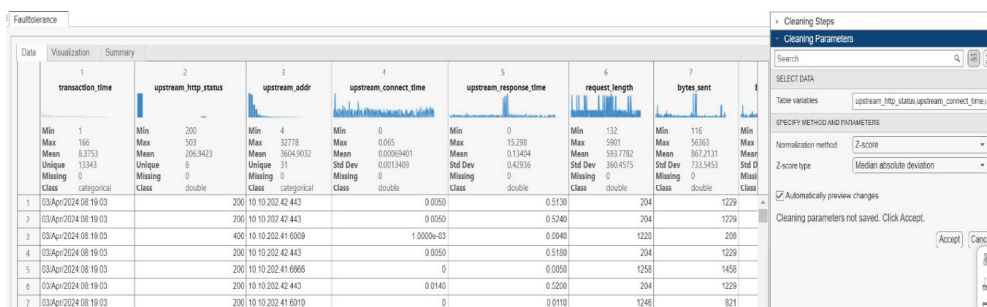


Fig. 6 Z-Score normalization

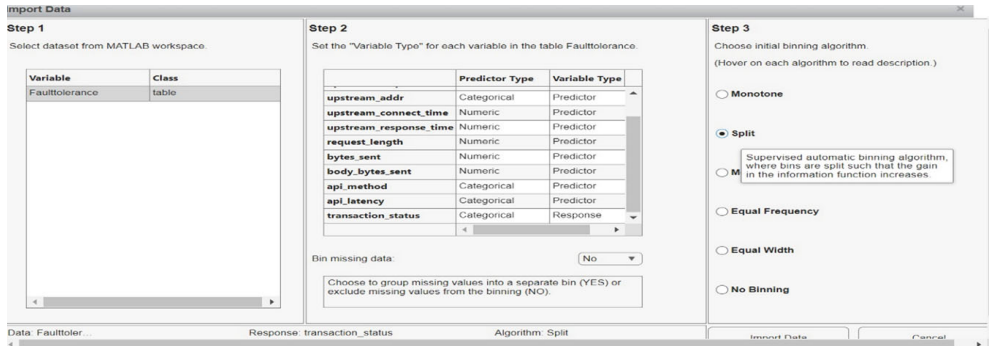


Fig. 7 Importing data for split binning method

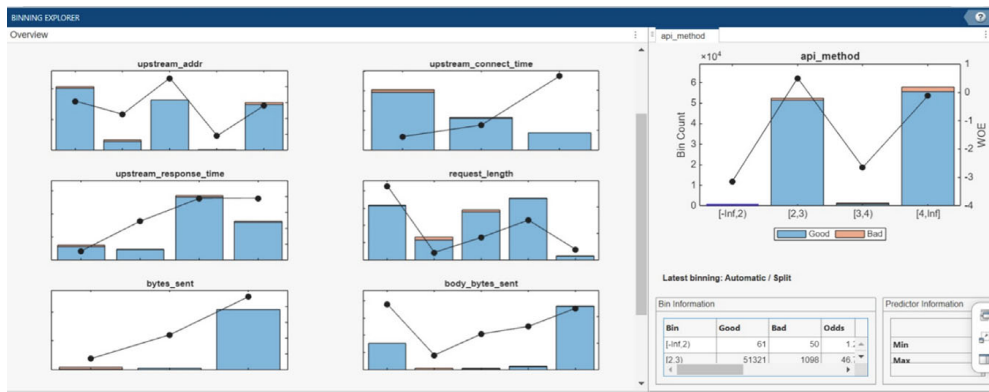


Fig. 8 Split binning method implemented for each attribute of incorporated dataset

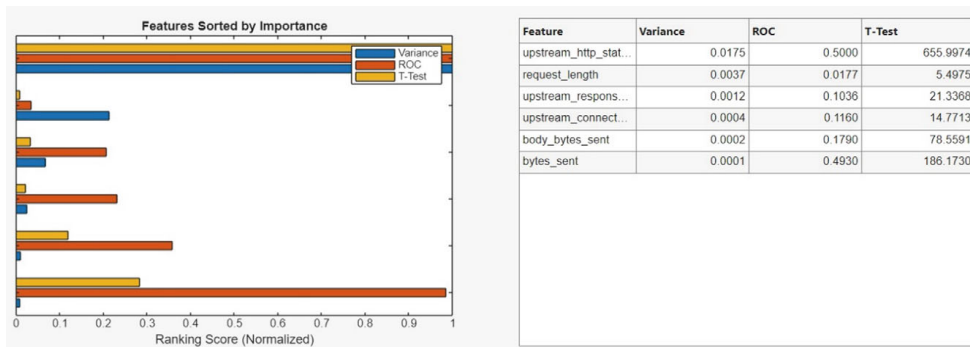


Fig. 9 Feature extraction and ranking

coordinate chart for each observation in the study is defined through the set of figures shown.

The performance chart of each algorithm used for stratification is illustrated through Fig. 13 shown. The results clearly depict that the Naïve Bayes algorithm supersedes the other two classifiers, and may necessitate optimal processing time.

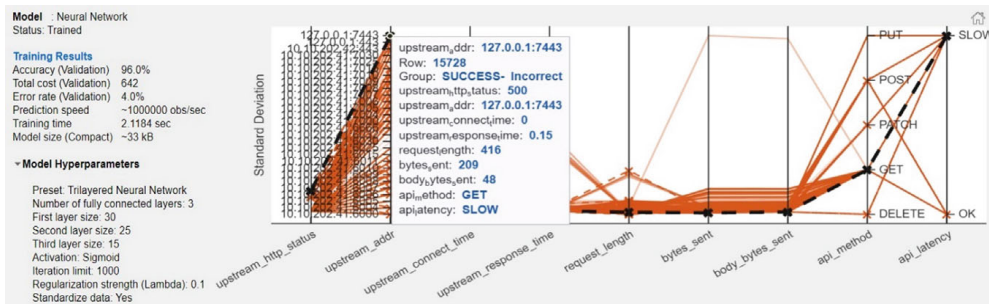


Fig. 10 Tri-layered neural network

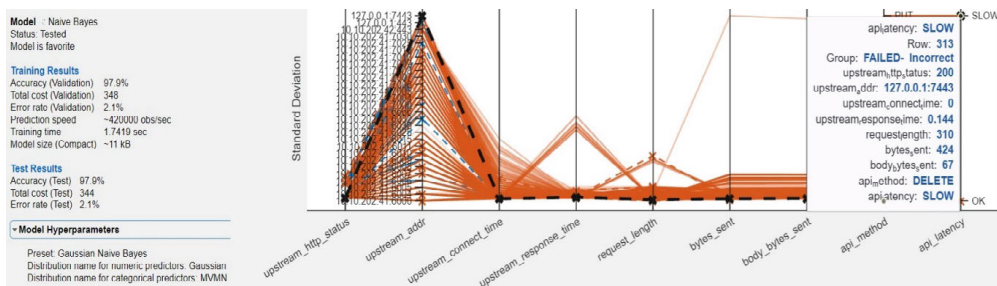


Fig. 11 Naïve Bayes algorithm

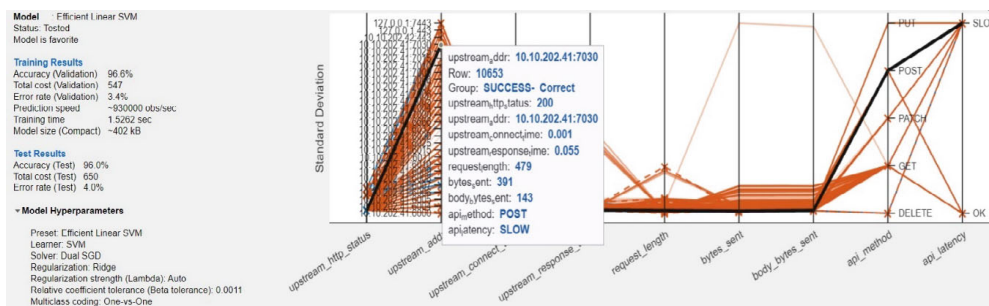
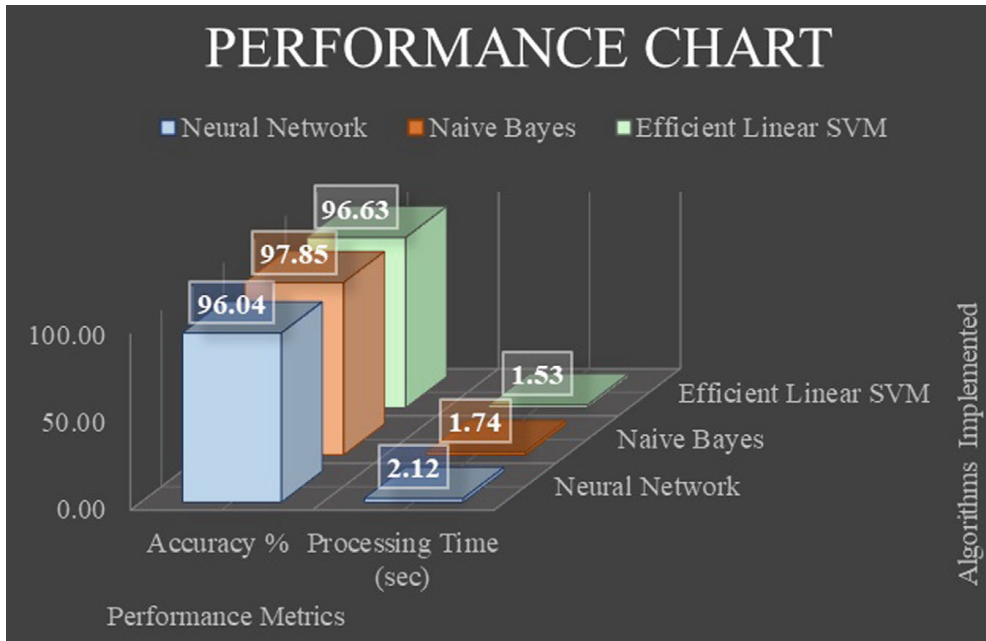


Fig. 12 Efficient linear support vector machine



**Fig. 13** Performance of classification algorithms

## 5 Conclusion

The pursuit of operatively functional system fault identification is a continuing endeavor that coerces collaboration across disciplines, innovation in methodologies, and incessant comprehensive cognizance of techniques. System behavior and development of robust fault identification strategies, mandates the understanding of vulnerabilities, while efficaciously scrutinizing the causal attributes. This study utilizes the real-time database comprising of various attributes that can trigger successful parsing, while bottlenecks in certain feature values can lead to failure of the system. The pivotal aspect of this indagation combines the various phases from pre-processing of the data to effective feature extraction and ranking to explicitly entail classification algorithms for procuring performance accuracy. This research effectuates a comprehensive pre-processing methodology to effectively clean the data and implement binning that further augments classifier accuracy. The feature ranking also considers the ROC, variance and *T*-Test to unambiguously identify the “http\_status” to render highest information gain. Classification through the ELSVM, Naïve Bayes and Tri-Layer Neural Network renders proximal classifier accuracies, with the naïve bayes accounting with the highest classification accuracy of 97.85%, followed by the ELSVM and Neural Network with 96.63% and 96.04%, respectively. This research renders a clear purview of stratification accuracy of each classifier algorithm, but future work can delve into clustering and stratification of specific system failures with advanced deep learning algorithmic models.

## References

1. Waqar H, Ali Bukhari SB, Wadood A, Albalawi H, Mehmood KK (2021) Fault identification, classification, and localization in microgrids using superimposed components and Wigner distribution function. *Front Energy Res.* <https://doi.org/10.3389/fenrg.2024.1379475>
2. Fahim SR, Sarker SK, Muyeen SM, Sheikh MRI, Das SK (2020) Microgrid fault detection and classification: machine learning based approach, comparison, and reviews. *Energies* 13:3460
3. Pirmani SK, Mahmud MA (2023) Advances on fault detection techniques for resonant grounded power distribution networks in bushfire prone areas: identification of faulty feeders, faulty phases, faulty sections, and fault locations. *Elect Power Syst Res* 220:109265
4. Li Y, Lin J, Niu G, Wu M, Wei X (2021) A Hilbert-Huang transform-based adaptive fault detection and classification method for microgrids. *Energies* 14:5040. <https://doi.org/10.3390/en14165040>
5. Chauhan P, Gupta CP, Tripathy M (2022) A novel adaptive protection technique based on rate-of-rise of fault current in DC microgrid. *Electr Power Syst Res* 207:107832. <https://doi.org/10.1016/j.epsr.2022.107832>
6. Sravani V, Venkata SK (2023) Detection of sensor faults with or without disturbance using analytical redundancy methods: an application to orifice flowmeter. *J Sens* 23(14):6633. <https://doi.org/10.3390/s23146633>
7. Zhang C, Wang H, Wang Z, Li Y (2023) Active detection fault diagnosis and fault location technology for LVDC distribution networks. *Int J Electr Power Energy Syst* 148:108921
8. Garramiola F, Poza J, Madina P, Del Olmo J, Ugalde G (2020) A hybrid sensor fault diagnosis for maintenance in railway traction drives. *Sensors* 20:962. <https://doi.org/10.3390/s20040962>
9. Fadlil A, Dikky PM (2023) Single imputation using statistics-based and K nearest neighbor methods for numerical datasets. *Ing des Syst d'Inform* 28(2):451. <https://doi.org/10.18280/isi.280221>
10. Kappal S (2019) Data normalization using median median absolute deviation MMAD based Z-score for robust predictions versus min–max normalization. *Lond J Res Sci Nat Formal.* <https://doi.org/10.13140/RG.2.2.32799.82088>
11. Gaudioso M, Gorgone E, Labbé M, Rodríguez-Chía AM (2017) Lagrangian relaxation for SVM feature selection. *Comput Oper Res.* <https://doi.org/10.1016/j.cor.2017.06.001>
12. Hemapriya CK, Suganyadevi MV, Krishnakumar C (2020) Detection and classification of multi-complex power quality events in a smart grid using Hilbert-Huang transform and support vector machine. *Electr Eng* 102:1681–2170
13. Atik C, Kut RA, Yilmaz R, Birant D (2023) Support vector machine chains with a novel tournament voting. *J Electron* 12(11):2485. <https://doi.org/10.3390/electronics12112485>
14. Trávez CJ, Quilumba FL, Transform AW (2018) Wavelet transform and support vector machine-based current-only directional overcurrent relay for transmission line protection. *IEEE PES Trans Distrib Conf Exhib Lat Am* 2:1–5
15. Chongya S, Kang Y, Alexander P, Jin L (2020) Rank-based chain-mode ensemble for binary classification. *Int J Comput Syst Eng* 14:153–158
16. Aker E, Othman ML, Aris I, Emmanuel O, Abdul Wahab NI, Hizam H (2020) Transmission line fault identification and classification with integrated FACTS device using multiresolution analysis and naïve Bayes classifier. *Int J Power Electron Drive Syst (IJPEDS)* 11(2):907. <https://doi.org/10.11591/ijpeds.v11.i2.pp907-913>
17. Srinivasa Rao Y, Ravi Kumar G, Kesava Rao G (2017) A new approach for classification of fault in transmission line with combination of wavelet multi resolution analysis and neural networks. *Int J Power Electron Drive Syst (IJPEDS)*. 8(1):505–512