# Analyzing Feature Selection Techniques for Machine Learning Based Anomaly Detection in IOT System.

## ABSTRACT:

As perhaps the most serious gamble in internet-based administrations, Distributed Denial of Service (DDoS) stays common. To slow the course of the client's entry, assailants can utilize a straightforward and effective DDoS assault technique. To perceive the DDoS assault, ML estimations are used. The regulated AI computations like k-closest neighbors(k-NN), calculated relapse, Naive Bayes, Decision tree, and arbitrary backwoods are used for acknowledgment and lightening of assault. By utilizing these five phases of information gathering, Pre-handling, includes designing, information parting, and element choice to find the best highlights in the given datasets (CICIDS-2017, NSL-KDD) with better precision. We accept that contrasting Machine Learning calculations regarding execution and pertinence will generally be important to get great exactness from the given datasets. It requires the best model to recognize noxious exercises as fast true to form and precisely. Different computations show unmistakable lead reliant upon the chosen highlights. The introduction of DDOS assault recognition is checked out and the best calculation is proposed. The paper portrays how to use highlight significance, connection network with heatmap, Principal Component Analysis (PCA), Linear Discriminant Analysis, and the Roulette wheel choice strategy is proposed to choose significant elements of the first assault information with the hereditary calculation.

## INTRODUCTION:

DDoS (Distributed Daniel of Service) is a sort of forswearing assault. For this situation, the DDoS identification model and the guarded framework programming depend on broad practice in the climate. Models of organization traffic can be utilized to prepare models, and afterward, track network assault exercises [1]. The outcomes show that the model essentially outflanks customary AI techniques.

Existing DDoS discovery strategies are essentially separated into edge recognition and scientific classification based on recognizable proof. Limit location frequently recognizes DDoS assaults and normal occasions by partitioning the edge of Internet highlights (e.g., source IP and objective IP). As of late, AI (ML) and information handling techniques play had a significant influence in recognizing and hence ordering meddlesome assaults. Normally, DDoS assaults are sent off with a botnet constrained by the aggressor. The botnet is typically planned with compromised motors.

The motivation behind this exploration is to more readily comprehend how every one of these AI calculations attempts to track down better precision between highlights. In this exploration paper, we will utilize Correlation-Based Feature Selection (CFS), a channel strategy for choosing the best component subset as per some assessment work, where highlights are viewed as restrictively free. With regards to expectation issues, the most generally utilized order-based AI calculations are Logistic Regression, Decision Tree, Random Forest, KNN, and Naive Bayes.

# RELATED WORKS:

ML techniques depend on regulated, solo, semi-managed, and support practice. These strategies are applied to recognize shrewd assaults and foster hearty guard instruments. The most well-known learning strategy is AI, which is observed in AI where the result is grouped given info utilizing a prepared informational collection. This strategy is likewise utilized in DDoS assault locations [2].

The Internet of Things (IoT) applications are moving around the world in general. Despite this, there are a variety of security risks. If there is insufficient information, the proposed technique considers recovering missing pieces as well as component reproduction. The NSL-KDD dataset was used for experiments that yielded a precision of 98.27% for 4-class identification, which was increased by increasing the number of hazy hubs [3].

The AI recognition process is being developed with extraordinary outcomes and perceptions in the security area of IoT. The proposed IDS is liable for DOS, information type examining, malignant control, vindictive activity, filtering, spying, and pernicious arrangement assault identification on different IoT destinations. Of all the observing calculations utilized, Random Forest gave the most elevated precision of 99.4% to recognize all assault types [4].

Numerous exemplary AI techniques, for example, the One-Class Support Vector Machine (SVM) and Principal Component Analysis (PCA), are presented as out-of-the-crate modules through the Azure Cloud Platform to speed up the advancement of universally useful inconsistency identification arrangements. The investigation of applying AI-based peculiarity location techniques to vertical plant divider frameworks for indoor environment control is critical to the Building and Environmental Research people group [5].

The paper, which utilizes a mix of brain association and an assist vector with machining, gives a location and portrayal technique for DDOS assaults on the media transmission network [6].

We contemplated Deep Neural Networks (DNN) to distinguish assaults in IoT. A sharp hinder recognition system should be constructed on the off chance that the constrained informational collection is to be gotten to [7].

The DDoS assault was done utilizing the ping of Death Strategy and distinguished the utilization of the AI approach by the WEKA instrument [8].

The outcome showed that the J48 arrangement had a high acknowledgment rate in the review. Gupta acquainted a strategy with forestall assaults in the cloud climate [9].

Commonplace AI estimations are taken on to accomplish the nave base, K-closest area (KNN), and support vector machine (SVM) to identify peculiar ways of behaving of data traffic. These three estimations are conflicting with their exhibitions and KNN has been viewed as more sensible than the other two [10].

Allomari subtleties and investigates the current botnet-based DDoS assaults and the harm brought about by DDoS assaults [11].

The information portrayed in this paper has assisted the accompanying analysts with actually safeguarding themselves from such assaults. Revathi and Malathi directed a nitty-gritty concentrate on the recognition of assaults utilizing different order strategies [12].

This proposes that it is essential to limit side effects before the arrangement. The NSL-KDD result is better when the element choice is utilized. In 2014, Aditya Harbola utilized arranging and voracious pursuit of exploration to channel NSL-KDD information. By utilizing this strategy, the choice of 20 significant elements and how much information can be split from the first. What's more, a high acknowledgment rate was gotten utilizing the KNN scientific categorization acknowledgment calculation [13].

## Background:

Unknown location is the identification of interesting occasions, subtleties, or assents that contrast essentially from standard activities or examples. Inconsistencies in information are otherwise called standard redirections, anomalies, clamor, sundries, and special cases.

There are three fundamental classes of peculiarity identification courses that are not regulated, or semi-endlessly administered. Fundamentally, the right issue identification situation depends on the markers accessible in the dataset. Ways of recognizing checked peculiarities require a total arrangement of informational indexes of "typical" and "uncommon" markers to work with the section calculation. This sort of style additionally incorporates preparing for grouping.

Observing machine proficiency makes a future model utilizing a preparation set named with basic and sporadic examples. The most well-known managed styles incorporate Bayesian organizations, k-close to neighbors, choice trees, regulated brain organizations, and SVMs.

Unaided styles don't request natively constructed naming for preparing information. All things considered, they felt that main a little, genuinely unique chance of organization business was horrendous and unprecedented.

The principal thought is to work on the capacity to distinguish highlights as a powerful method for diminishing the size of elements without lessening the simplicity of articulation to expand the effectiveness of mental examination, which comprises two primary advances: include choice and component extraction. What's more, even though highlight determination strategies are broadly utilized and accomplish great outcomes, overt repetitiveness because of the relationship between's elements isn't considered, so include extraction techniques ought to be presented. A more viable technique called chief can lessen estimations by wiping out part relationship (PCA) connection. Direct Discriminant Analysis (LDA) is a typical innovation for layered decrease issues as a pre-handling venture for AI and example division tasks.

The writing recommends that calculated relapse, nave base, KNN (K prompt neighbors), choice tree, and irregular woodland designs have shown shrewd execution for answers to scientific classification issues. During this review, the exhibition rate location of DDoS assaults from their order bunch competitors, strategic relapse, nave base, KNN (K closest neighbors), choice tree, and irregular woods designs was examined.

## Methodology:

A hereditary calculation is a system that drives organic development to take care of restricted and uncontrolled streamlining issues. The hereditary calculation chooses people from the ongoing populace as guardians at each stage and uses them to make them up-and-coming age of youngsters. GA keeps up with the number of inhabitants in chromosomes (arrangements) related to wellness values [14]. Guardians are chosen to mate given their wellness by delivering posterity through a regenerative arrangement. Thus, the most appropriate arrangements are offered more chances for propagation with the goal that the posterity acquires characteristics from each parent.

After the underlying populace is arbitrarily created, the calculation is created by the administrator:

select wellness is the same decision for endurance;

Hybrid showing intercourse between people;

The transformation presents irregular changes [16].

# ROULETTE WHEEL SELECTION:

The roulette wheel, created by Holland, was the best option technique. The likelihood of Pi for every individual is characterized by the way that Fi is equivalent to the wellness of individual I. The utilization of roulette wheel choice restricts the hereditary calculation to amplification since assessment work arrangements should be planned to completely requested values above + ℜ [14]. Expansions, for example, windows and scaling are proposed to permit minimization and antagonism.

All individuals of the accompanying age are chosen to utilize the roulette wheel technique. In a hereditary calculation, this is a basic choice procedure. The relative wellness (individual wellness to add up to wellness proportion) of every individual is utilized to plan the roulette wheel. The individual hybrid likelihood is determined by isolating a singular's wellness by the absolute populace wellness.

Pi is the likelihood of every chromosome, which is the chromosome recurrence isolated by the all-out wellness of the roulette wheel. Accept that the choice system for roulette wheels is like the diagram above. Every individual has wellness esteem and the circle is the amount of every one of them. Thus, your possibilities observing a potential not entirely settled by your total wellness. The fitter individual has more top on the haggle, hence, bound to land before a decent point/pointer when the wheel is turned. Subsequently, the probability of picking an individual relies straightforwardly upon their wellness [15]. It has less time intricacy when run equally. No wellness evaluating or arranging is required.

A straightforward determination strategy appoints the likelihood of choice Pj to every individual j in light of its wellness esteem. Ci-1 <U (0, 1) ≤ Ci however N is a progression of irregular numbers created and contrasted and the aggregated and replicated into the new populace. There are various strategies for appointing probabilities to people: roulette wheel, straight positioning, and mathematical positioning.

Given the likelihood that an individual will be chosen as a parent for the hybrid

$$P(i) = \frac{f(i)}{\sum_{i} f(i)}$$

$f(i)$: fitness of individual $i$
$P(i)$: probability that individual $i$ is chosen

## Algorithm 1 Pseudocode for a Genetic Algorithm

1: t ← 0;

2: InitPopulation[P(t)]: {Initializes the population}

3: EvalPopulation[P(t)]; {Evaluates the population}

4: **While** not termination **do**

5:     P'(t) ← Variation|P(t)); {Creation of new solutions}

6:     EvalPopulation[P'(t)]; {Evaluates the new solutions}

7:     P(t+1) ← ApplyGeneticOperators[P'(t)∪Q]; {Next generation pop.}

8:     t ← t+1;

9: **End While**

## Algorithm 2 Pseudocode of Roulette Wheel Selection

**Function** index=Roulette Wheel Selection (**FV**)

1: **Begin**

2:    **FV** = [fv,fv2, ..., fvN]^T is the unsorted fitness vector;

       // bold font indicates vector form

3:    [**sFV, id** ] = sort(**FV**, 'ascend');    // **sFV** is the sorted fitness vector

                                // **id** is the index of sorted individuals

4:    Let **p** = [pi]^T,i=1, 2, ...,N    // probability vector

5:    ∀i,pi = fvi/Σ^Ni=1 fvi;

6:    **p = (1− p)/Σp**;    /* to assign high probability to individuals

      with low numerical value as fitness for minimization problem */

7:    r = rand(0,1);    // random number between 0 to 1

8:    i = 1;

9:    **While** r > 0

10:       r = r - Pi; i++;

11:    **End While**
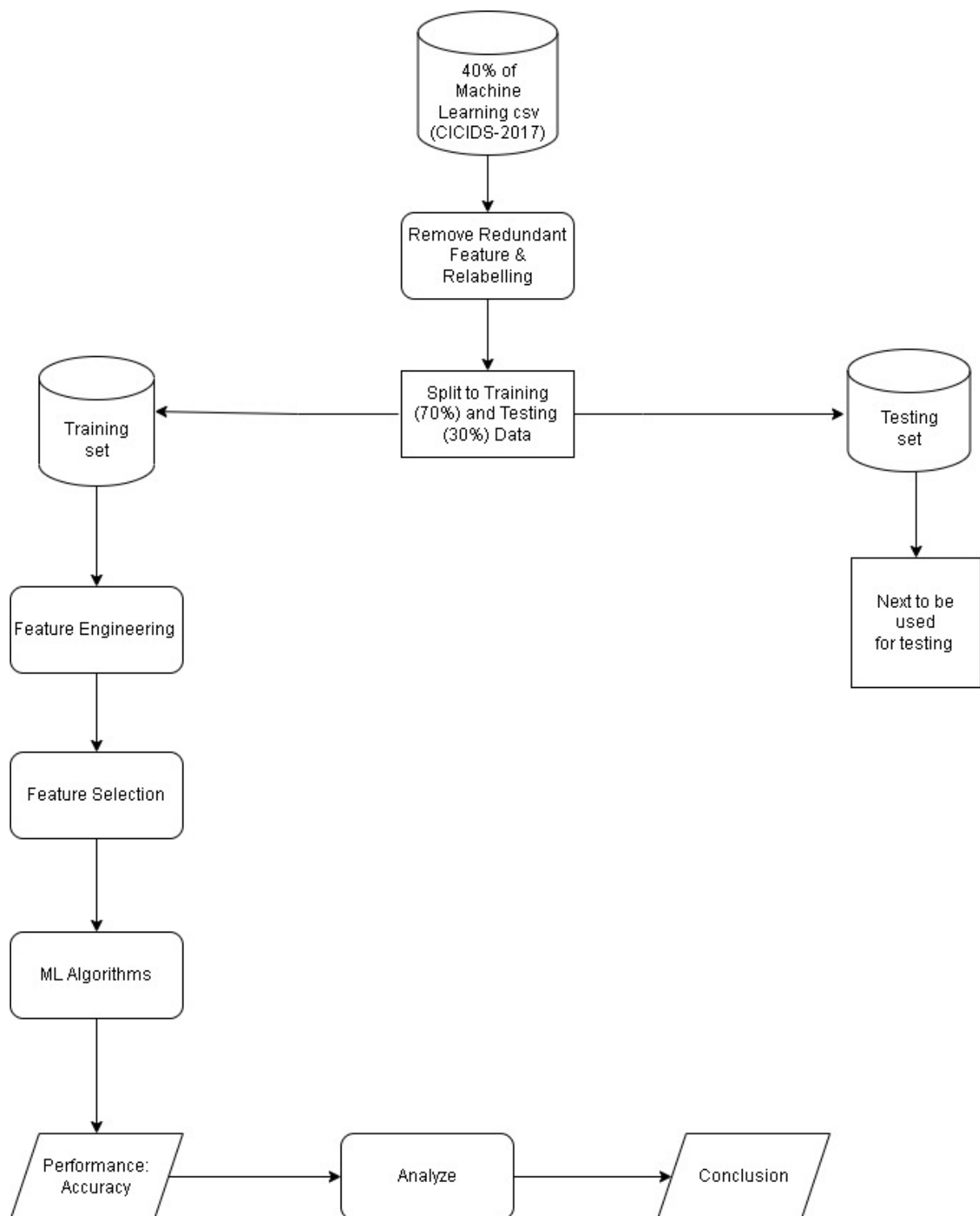
12:    index ← id(i − 1);

13:**End**



**Fig. 1** The Process of Evaluation-Based Algorithms
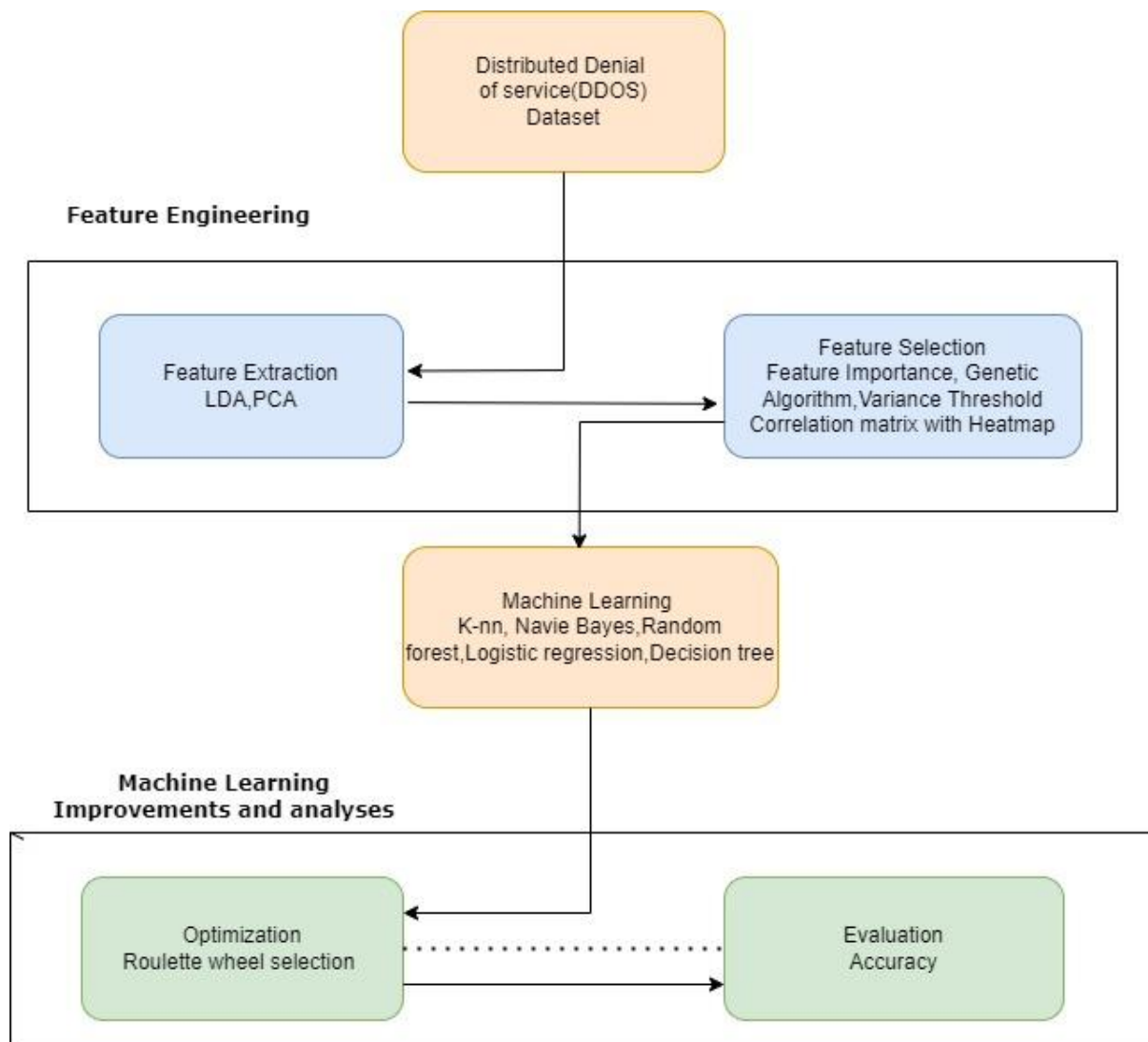
**Fig. 2** Strategic-Level framework for DDoS attack detection
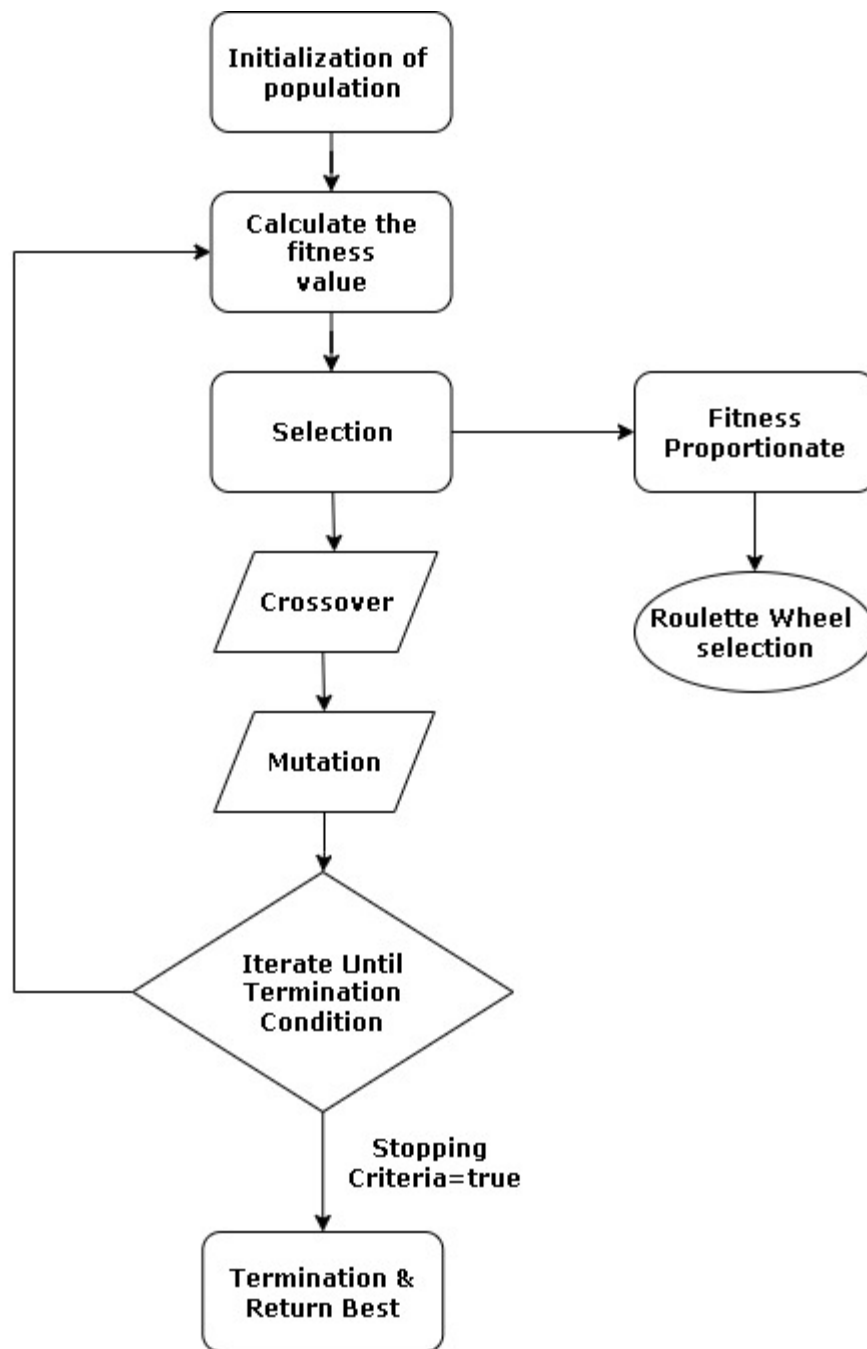
**Fig. 3** Steps involved in genetic algorithm
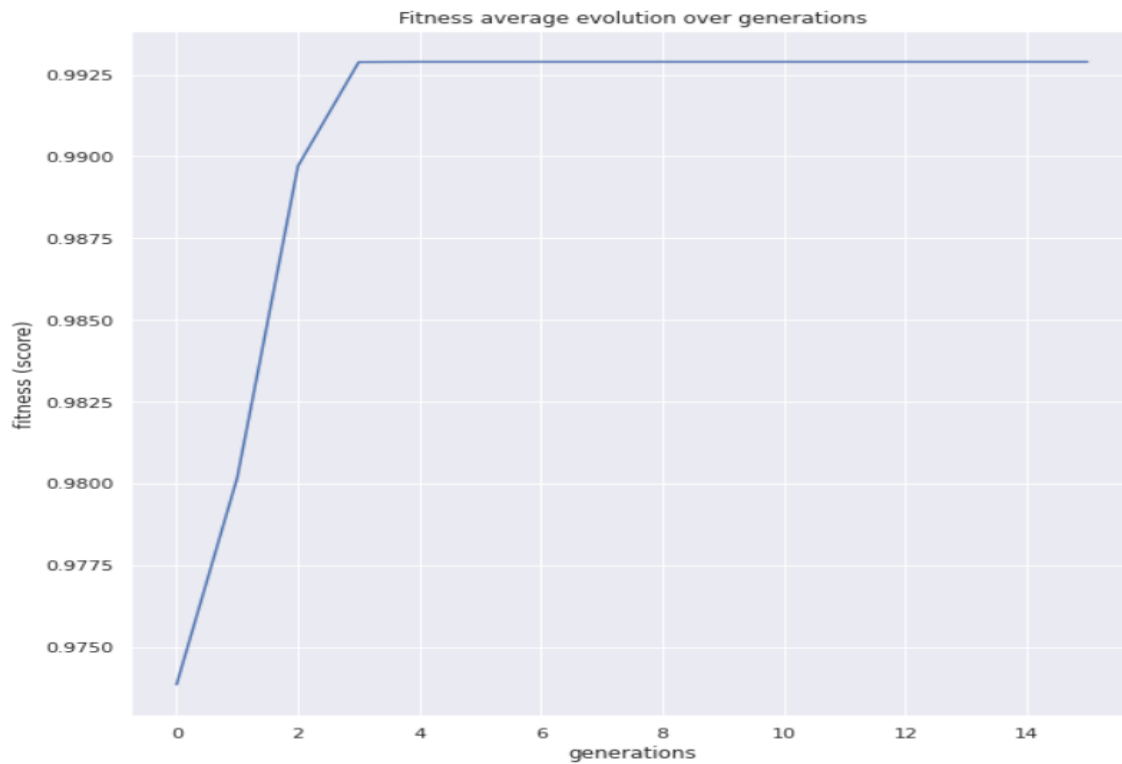
# Evaluation of result:



**Fig. 4** Fitness Evolution Plotting of CICIDS2017 Dataset
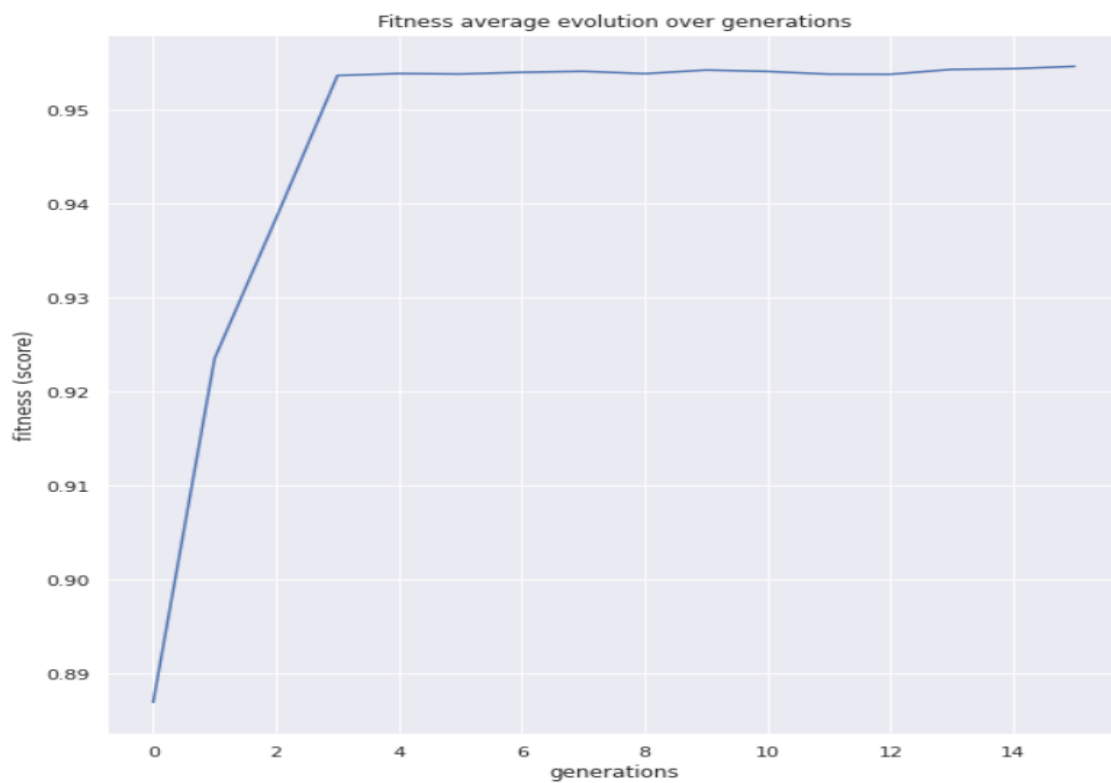


**Fig. 5** Fitness Evolution Plotting of NSL-KDD Dataset
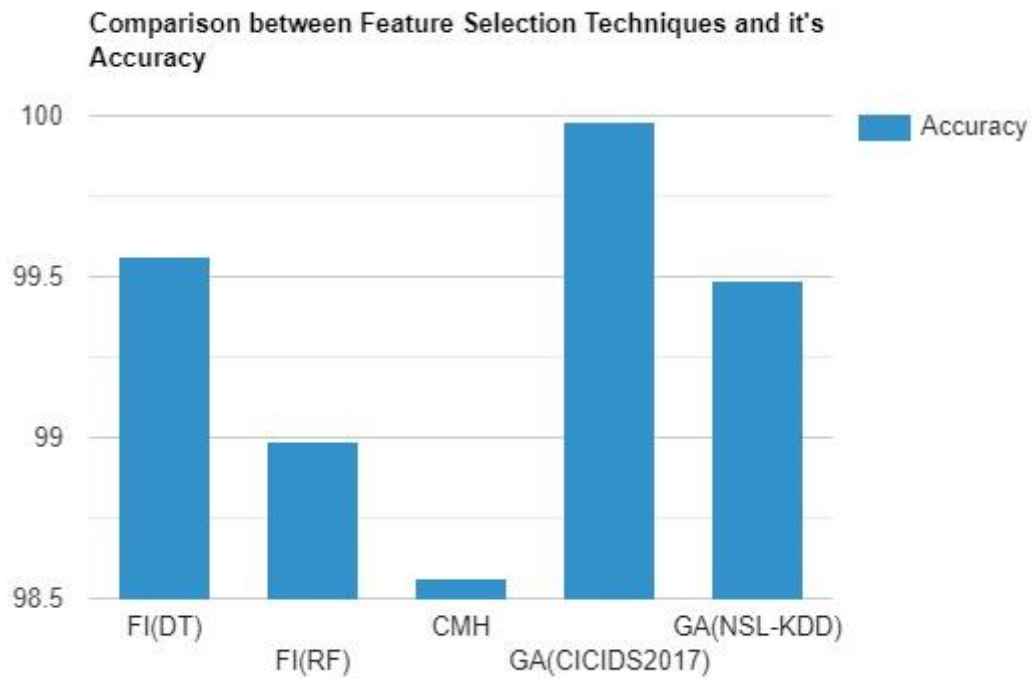
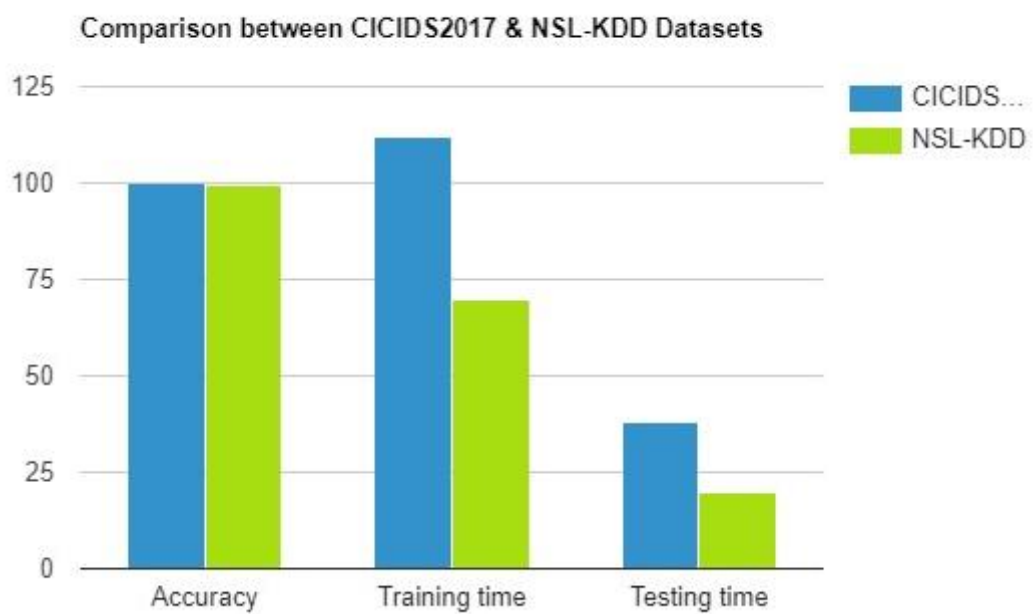**Fig. 6** Comparison between feature selection techniques and accuracy



**Fig. 7** Comparison between CICIDS2017 & NSL-KDD Datasets

**TABLE 1 Accuracy of ML Algorithms of CICIDS2017 Dataset**

| Algorithms | Accuracy |
|---|---|
| Decision Tree | 100% |
| Logistic Regression | 98.79% |

**TABLE 2 Accuracy of ML Algorithms of NSL-KDD Dataset**

| Algorithms | Accuracy |
|---|---|
| Decision Tree | 99.97% |
| Logistic Regression | 86.41% |

**TABLE 3 Feature Selection Techniques with Accuracy of CICIDS2017 Dataset**

| Feature Selection | Algorithm | Accuracy |
|---|---|---|
| Principal Component Analysis(PCA) | Naïve Bayes | 96.8% |
| | Logistic Regression | 96.8% |
| Linear Discriminant Analysis(LDA) | Knn | 99.74% |
| | Logistic Regression | 97.35% |

**TABLE 4 Parameters assumed in evolutionary algorithm**

| Parameter | Value |
|---|---|
| Number of generations | 15 |
| Chromosome population Size | 60 |
| Crossover probability | 0.6 |
| Crossover type | Single |
| Mutation probability | 0.05 |
| Mutation type | Uniform |
| Evolutionary algorithm | eaMuPlusLambda |

**TABLE 5 Fitness Evolution Table of CICIDS2017 Dataset**

| Generation | Fitness(accuracy) |
|:---:|:---:|
| 0 | 0.973849 |
| 1 | 0.98018 |
| 2 | 0.989705 |
| 3 | 0.992883 |
| 4 | 0.992891 |
| 5 | 0.992891 |
| 6 | 0.992891 |
| 7 | 0.992891 |
| 8 | 0.992891 |
| 9 | 0.992891 |
| 10 | 0.992891 |

**TABLE 6 Fitness Evolution Table of NSL-KDD Dataset**

| Geneartion | Fitness(accuracy) |
|:---:|:---:|
| 0 | 0.886952 |
| 1 | 0.923541 |
| 2 | 0.938454 |
| 3 | 0.953646 |
| 4 | 0.953857 |
| 5 | 0.953804 |
| 6 | 0.95399 |
| 7 | 0.954096 |
| 8 | 0.953844 |
| 9 | 0.954228 |
| 10 | 0.954082 |

**TABLE 7 Features Selected By Roulette Wheel with accuracy of CICIDS2017 Dataset**

| Feature Selection | Algorithm | Feature Selected | Accuarcy |
|---|---|---|---|
| Feature importance | Decision Tree | Destination port<br>Subflow Fwd Packets<br>Fwd IAT Min<br>Packet Length Variance<br>Init_Win_bytes_backward<br>Bwd Packets/s<br>Fwd Packet Length Max<br>Flow IAT Min<br>Bwd IAT Total<br>Bwd IAT std | 99.98% |
| | Random Forest | Destination port<br>Avg Fwd Segment Size<br>Fwd Packet Length Std<br>Packet Length Variance<br>Fwd Packet Length Max<br>Average Packet Size<br>Packet Length Std<br>Packet Length Mean<br>Max Packet Length<br>Flow IAT Min | 98.34% |
| Roulette Wheel | Decision Tree | Destination Port<br>Flow IAT Min<br>Fwd IAT Total<br>Fwd IAT Min<br>Bwd URG Flags<br>CWE Flag Count<br>Fwd Header Length.1<br>Subflow Bwd Bytes<br>Init_Win_bytes_forward<br>Init_Win_Bytes_backward | 99.88% |

**TABLE 8 Features Selected By Roulette Wheel with accuracy of NSL-KDD Dataset**

| Feature Selection | Algorithm | Feature Selected | Accuarcy |
|---|---|---|---|
| Roulette Wheel | Decision Tree | service<br>flag<br>src_bytes<br>is_host_login<br>same_srv_rate<br>diff_srv_rate<br>dst_host_count<br>dst_host_srv_diff_host_rate<br>dst_host_error_rate<br>level | 99.49% |

**TABLE 9 Comparision of different selection methods**

| Feature Selection Methods | Number Of Features | Features Selected | Accuarcy |
|---|---|---|---|
| Feature-Importance(Decision Tree) | 10 | 28,26,20,7,38,68,43, 25,63,1 | 99.56% |
| Feature-Importance(Random Forest) | 10 | 20,40,41,42,53,7,43, 10,54,1 | 98.99% |
| Correlation matrix with heatmap | 35 | 74,53,54,55,27,29,30,13,14, 51,18,19,56,22,24,25,10,77, 78,40,41,42,43,45,65,66,64, 4,6,69,26,21,37,75,63 | 98.56% |
| GA-RouleteeWheel(CICIDS2 | 10 | 1,20,21,25,34,50,56,66,67,68 | 99.98% |
| GA-RouleteeWheel(NSL-KDD) | 10 | 3,4,5,21,29,30,32,37,38,43 | 99.49% |

## CONCLUSION:

In this work, a methodology for including a decision is called Feature Selection for CICIDS2017. Interruption Dataset is acquainted and it was capable with characterize various kinds of DoS/DDoS assaults. A web-based answer for identifying DoS/DDoS assaults was presented in this examination. As per the testing results, when we use highlight significance with the choice tree it is giving 99.98% precision. Whenever we use include significance with the irregular backwoods it is given 98.34%. A choice tree is giving better precision with highlighted significance. Whenever we use head part examination with credulous Bayes it is giving 96.8% exactness. At the point when we use head part examination with strategic relapse, it is giving 96.8%. Both the strategic relapse and the nav base give similar exactness with the examination of the significant parts. Whenever we utilize direct separation examination with calculated relapse it gives 97.35%. Whenever we utilize straight segregation examination with KNN it gives 99.74%.KNN offers better precision with direct separation investigation. After investigating the CICIDS2017 and NSL-KDD datasets, the elements were diminished from 79 to 10 and the information on the CICIDS2017 dataset was decreased to around 13%. Highlights have been decreased from 43 to 10 and information has been diminished to practically 24% of the first IN NSL-KDD dataset.

Both proficiency and exactness were improved with the proposed include choice strategy. The acknowledgment rate isn't just expanded by the number of highlights. Contingent upon whether the chosen highlights are pertinent to the occasion. The proposed ID strategy centers around

further developing the ID model by choosing significant elements. The registering asset of the location part saves time as the determination cost increments.

## FUTURE WORK:

This work gives the establishment for a few future frameworks and more trials may likewise be directed to incorporate the variety of the AI calculations, for example, regulated, unaided, and semi-managed models across various DDoS Attack related datasets. Include determination is a to a great extent open investigation region and we accept that crossover techniques for highlight choice utilizing different methodologies of factual instruments and multi goals, for example, arbitrary pursuit or some extra calculation can be the best method of element determination for DDoS assault discoveries. Our Future work is to examine ways for further developing the disclosure pace of different assaults. It'll incorporate examining DDoS assaults given administration weaknesses like adding numerous class grouping and setting up prudent steps. Changes in advancement methods like subterranean insect province enhancement calculation, and PSO (molecule swarm improvement) make more investigations in hereditary calculations, contrast, and different datasets like Caida, SDN, and so forth. Change the choice interaction and examination with various information like changing the number. Over the ages, we desire to work with the analytical cycle by participating in the endeavors of scientists.

## REFERENCES:

[1] Pawar, Mohan V., and J. Anuradha. "Network security and types of attacks in the network." Procedia Computer Science 48 (2015): 503-506.

[2] Aamir, Muhammad, and Syed Mustafa Ali Zaidi. "Clustering-based semi-supervised machine learning for DDoS attack classification." Journal of King Saud University-Computer and Information Sciences 33.4 (2021): 436-446.

[3] Aversano, Lerina, et al. "Effective Anomaly Detection Using Deep Learning in IoT Systems." Wireless Communications and Mobile Computing 2021 (2021).

[4] Tyagi, Himani, and Rajendra Kumar. "Attack and Anomaly Detection in IoT Networks Using Supervised Machine Learning Approaches." Rev. d'Intelligence Artif. 35.1 (2021): 11-21.

[5] Liu, Yu, et al. "Anomaly detection based on machine learning in IoT-based vertical plant wall for indoor climate control." Building and Environment 183 (2020): 107212.

[6] Says, S. "DDOS attack detection in a telecommunication network using machine learning." Journal of Ubiquitous Computing and Communication Technologies (ACCT) 1.01 (2019): 33-44.

[7] Choudhary, Sarika, and Nishtha Kesswani. "Analysis of KDD-Cup'99, NSL-KDD and UNSW-NB15 datasets using deep learning in IoT." Procedia Computer Science 167 (2020): 1561-1573.

[8] Pande, Sagar, et al. "DDOS detection using machine learning technique." Recent Studies on Computational Intelligence. Springer, Singapore, 2021. 59-68.

[9] Gupta, B.B., Badve, O.P.: Taxonomy of dos and DDoS attacks and desirable defense mechanism in a cloud computing environment. Neural Comput. Appl. 28(12), 1–28 (2017)

[10] Prakash, Aditya, and Rojalina Priyadarshini. "An intelligent software-defined network controller for preventing distributed denial of service attack." 2018 Second International Conference on Inventive Communication and Computational Technologies (ICICCT). IEEE, 2018.

[11] Alomari, E., Manickam, S., Gupta, B.B., Karuppayah, S., Alfaris, R.: Botnet-based distributed denial of service (DDoS) attacks on web servers: classification and art. Int. J. Comput. Appl. 49(7),24–32 (2012)

[12] Revathi, S., Malathi, A.: A detailed analysis on the all-kdd dataset using various machine learning techniques for intrusion detection.In: International Journal of Engineering Research and Technology(2013)

[13] Harbolt, A., Harbola, J., Vaisla, K.S.: Improved intrusion detection in DDoS applying feature selection using rank & score of attributes in the kdd-99 data set. In: International Conference on Computational Intelligence and Communication Networks, pp. 840–845 (2014)

[14] Pencheva, Tania, Krassimir Atanassov, and Anthony Shannon. "Generalized net model of selection function choice in genetic algorithms." Recent Advances in Fuzzy Sets, Intuitionistic Fuzzy Sets, Generalized Nets, and Related Topics 2 (2010): 193-201.

[15] Shyalika, C. (2019, March 2). Genetic algorithms -selection. Medium. Retrieved April 22, 2022, from https://medium.datadriveninvestor.com/genetic-algorithms-selection-5634cfc45d78

[16] Barekatain, Behrang, Shahrzad Dehghani, and Mohammad Pourzaferani. "An energy-aware routing protocol for wireless sensor networks based on a new combination of genetic algorithm & k-means." Procedia Computer Science 72 (2015): 552-560.