

# A Hyper Meta-Heuristic Cascaded Support Vector Machines for Big Data Cyber-Security

G.A.Mylavathi, B.Srinivasan



**Abstract**— At an incredible speed, cyber security evolves in the ever-changing setting of attacks. Organisation processing of information inward and outward is huge in quantity and determining a threat amidst of information is challengeable. Late discovery of such instance is standstill challenge of the meticulous process. Thence, detection of intrusion and its prevention are rising challenge in Big data factors. the information inundation generally incorporate the Big data terms to dataset. The majorly focused issues are industrial oriented in big data challenge. Existing systems for big data cyber security problems are based on Online Support Vector Machines (OSVMs) framework. Bi-objective optimisation problem with primary objectives is designed as OSVMs configuration process for improving accuracy and less complexity of model. Here, a bi-objective optimization is implemented based on an Artificial Bee Colony (ABC). However, Online Support Vector Machines (OSVMs) has issue with computational complexity, and prematurity and local optimum is major problems in ABC algorithm. By overcoming this issue, developed research system designs an Ensemble Support Vector Machine (ESVM) framework for big data cyber security. Initially, the feature selection is done by using improved K-means clustering. Based on the selected features the intrusion detection and malware detection are performed using ESVM approach. In this proposed research work, a bi-objective optimization problem is designed as the ESVM configuration process for improving accuracy and less complexity of model and achieve its objectives. Cuckoo Search (CS) optimization algorithm is implemented for the bi-objective optimization. accuracy, precision, recall and f-measure are the parametric meters compared in proposed research attaining higher performance against existing approaches.

**Keywords:** Cyber security, Cuckoo Search (CS), Ensemble Support Vector Machine (ESVM) and improved K-means clustering.

## I. INTRODUCTION

Huge growth of threats in cyber security widely influences the network user. Even several systems in monitoring and screening are exist, cyber attacks also increases. Hence there is a need of security monitoring system strongly for huge

dataset in networking. A work against malware attacks are proposed in this research. IP address of attackers is retrieved incorporated dataset for producing data statistically. The characteristics of attacker's IP addresses can be extracted from our integrated datasets to generate statistical data. Each attribute weight is generated by the cyber security expert to produce a scoring system by detailing the log history.

Technologies of big data explores tremendously the Cyber security system [1].

Huge data processing, data storing in Big data are widely growing by technically. Conventional computing environment holds Big data requiring of a flexible storage space and performing data analytics by querying is faced a vast amount of challenges in dealing with huge volume of data against volume of Big Data positive feature. Several networks in social media like Google are much benefited with holding out vast information. Tremendous unsupervised data as uncategorised and a little supervised data are included in Big Data [2].

Several group, business and society used to retrieve, stock up and perform analytics. These sort of data are referred as "big data" on account of its velocity, veracity, volume, and -variety of data. It can be better defined as

- High volume—data quantity
- High velocity— rate of data creation
- High variety—data sort (both homogeneous and heterogeneous) .

Capturing of data, storing of data, and analysing of data are performing characteristics of innovated techniques. There are several sources of Big data [3]. It is mainly focused with its complexity and challenges of big data along its characteristics than mere its volume and its security. a bi-objective optimization is the primary aim of this research in formulating ESVM. The bi -objective functions are system accuracy and the system complexity. This work have introduced feature selection model for space search before clustering to improve the performance of detecting intrusion detection and malware detection.

Irrelevant features are eliminated by model of feature selection in function of the model more appropriate [4]. Minimum attributes are retrieved from this process in improvisation and understanding patterns thus increasing speed of learning stages.

data mining encloses crucial tasks as clustering. Several applications as business intelligence, image pattern recognition, biology, security, and Web search get succeed by Clustering. Data intrinsic structures are exploited and sorted in sub clusters.

Manuscript published on November 30, 2019.

\* Correspondence Author

**Mrs.G.A.Mylavathi\***, Assistant Professor of Computer Science, Gobi Arts & Science College, Gobichettipalayam, Tamil Nadu, India. (Email: mylavathiga@gmail.com)

**Dr.B.Srinivasan**, Associate Professor of Computer Science, Gobi Arts & Science College, Gobichettipalayam, Tamil Nadu, India. (Email: bsrini2561967@gmail.com)

© The Authors. Published by Blue Eyes Intelligence Engineering and Sciences Publication (BEIESP). This is an [open access](http://creativecommons.org/licenses/by-nc-nd/4.0/) article under the CC-BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>)

The major work is to split the group into subclusters with holding similar attributes of objects. Single cluster is a small subsets. In sequence they are clustered of objects based on interclass minimization and similarity of intraclass maximization. Like and unlike objects along with attributes of feature values determines distance measures. Objects are compared with other based on its

feature value similarity. Supreme distances, Manhattan distance, and Euclidean distance are the distance measures supports to classify the objects of numeric data. Globally, performing of analysis Cluster on dataset is huge task and therefore require several efficient algorithms to implement [5].

SVM is used for the classification which classify whether the unknown entry of user is a authenticated or not. Further this work brings a concept of Meta-heuristic scheme [6]. A heuristic optimization algorithm is developed to resolve a framework of metaheuristic, high-level problem-independent with designed strategies. Some instances are existing like metaheuristics include genetic/evolutionary algorithms, tabu search, simulated annealing, variable neighborhood search, (adaptive) large neighborhood search, and ant colony optimization, a metaheuristic framework or metaheuristic with low level heuristics and high level heuristic are referred as a heuristic optimization algorithm with problem-specific implementation. metaheuristic framework strategies are included in Metaheuristic algorithms, as method of optimization which is nature of heuristic.

From the overall, bi-objective optimisation problem as ESVM configuration process are done successfully here via hyper meta heuristic framework using Cuckoo Search algorithm. Configuration is generated by Cuckoo Search algorithm and process to the ESVM for cost function production. Thus accuracy is increased and complexity of system is decreased in this research proposed. Before classification, cluster based feature selection will takes place for classification improvisation in turn of reducing complexity of system.

## II. LITERATURE REVIEW

An approach developed by Gunantara et al [7] does work in pair path selection on more criteria as adhoc network. It analyze the meta-heuristic methods as of Ant Colony optimization (ACO), genetic algorithm (GA), and particle swarm optimization (PSO). Power consumption, signal-to-noise ratio (SNR), and load variance are many criteria used for analysis. Time consumption, performance as SNR, load variance, and power consumption on criteria leads to study analysis of experimental results iteratively with same values. At last, GA algorithm performs the path pair via ACO and PSO methods varying results.

An approach developed by Hassani and Jafarian [8] determines the cancer at early phase that helps in reducing risk factors via technique of hybrid classification. Women face breast cancer as commonly found widely. Hence this developed work of optimization in metaheuristic algorithms helps to determine the Fuzzy-ART parameters. But Fuzzy-ART fails to good prediction of cancer on availing data. But its result gets improved efficiently via optimization methods of evolution. Every techniques of hybrid classification for breast cancer are examined with data of

trained set provided through dataset of Wisconsin. accuracy of 97.80% is attained and specificity of 98.92% is attained by the developed approach of Hassani and Jafarian .

An approach developed by Teoh et al [9] consider a dataset in huge network inclusive of malware attackers. Hence it is possible to recognise the malware data on training with cyber security attack for exposing a practised system. The data extraction in statistical along with incorporated dataset extracts the of attacker's IP addresses characteristics. Each attribute weights are annotated by expert of cyber security and log history annotation builds a scoring system. In sequence, log system for cyber security is determined, classified and evaluated by specific semi supervise method. It unures that there is no possible attacks by Fuzzy K-Means (FKM) of segregation of data into 3 clusters and label manually small data (analyst intuition). Following that finally neural network classifier train the Multi-Layer Perceptron (MLP) to the base of data labelled manually. On performing such a kind, experimental results establishes more hopeful against cyber security logs that determining anomalies excluding labelling by analyst intuition. A Cyber security log by intuition establishes detection by fake.

An approach developed by Khorram and Baykan [10] attach on Particle Swarm Optimization (PSO), Ant Colony Optimization (ACO), Artificial Bee Colony (ABC), K-Nearest Neighbors (KNN) and Support Vector Machine (SVM). Attacks on networks are identified with appropriate feature set are determined by the above said algorithms while feature performance are helps to classify by KNN and SVM algorithms of feature selection. in this study, the training and testing are done by standard NSL-KDD dataset. Rate of accuracy as 98.9% is attained and false alarm 0.78% is attained by Feature selection to the base of ABC and KNN classifier algorithm accordingly, exhibits best as amidst of evaluated algorithms.

An approach developed by Saidala and Devarakonda [11] developed hybridizing WOA with Clonal Selection Algorithm for a new parallel meta-heuristic optimization. 23 standard mathematical benchmark functions are used for evaluation. Experimental results significantly perform the statistical scenario of proposed to determine the solutions of optimization for complexity models. Additionally, heart patient dataset are optimised using this technique. Accuracy, Precision, Recall, and F-measure are four evaluation metrics to determine the result against existing methods in prediction of heart disease. The simulated results exploits this developed method is better than many and competitive for predicting heart disease.

An approach developed by Tsai et al [12] establishes meta-heuristics approach on survey for healthcare system. Effective healthcare system is developed by researchers meta-heuristics work with roadmap. Healthcare system gets change by review of features updating at regular intervals. Big data challenges are overcome by a framework of learnable big data analytics in system of healthcare system for obtaining a solution of greater performance. At last, addressing of updation, modification, open problems, further trends of metaheuristics are done in systems of healthcare.

An approach developed by Bajpai and Dayanand [13] concentrates on security of big data while implementation. Automation of discovering hidden insights, decisions improvement and processing business are carried out. Visually able analyse and insight drawings are able to predict and haul threats of cyber security for security in Big data

analytics for huge data gathering. A cyber defence posture is obtained amidst of technologies in security. Network threats are represented by entering organizations for activity of pattern reorganization.

An approach developed by Sabar et al [14] designed bi-objective optimisation problem with SVM configuration process for obtaining better accuracy. Two conflicting objectives are focused with model complexity. For bi-objective optimisation, domain problem are independent to the proposal of a novel hyper heuristic framework. A low-level heuristics and high-level strategy are done in the proposed hyper-heuristic framework. Selection is accessed and controlled in performance of search by the high-level strategy and a new SVM configuration produces low-level heuristic. The SVM configurations of search space explore the low-level heuristics with varying several protocols significantly. Decomposition strength is incorporated to Pareto-based approaches for addressing optimisation of bi-objective and approximations of Pareto set of SVM configurations.

An approach developed by Hou et al [15] developed for detection of a new android malware system to the base of deep belief network depending on API call blocks. This research extracts small files on evaluating Application Programming Interface (API) calls and split the small code files into blocks in API calls. Deep learning framework (i.e., Deep Belief Network) is applied depending on block code generation and hence find any sort of unpredictable Android malware detection. Experimental results are produced by a real sample collection from Comodo Cloud Security Center.

An approach developed by Nguyen et al [16] proposes a new approach on enhancing environmental security of mobile devices. It implements a novel system for determining the automatic malware intrusions. The designed system applies in the process of device operation by the approach of heuristic analysis for producing logs by mobile and holds the base in behaviours of user modelling. Achieving of awareness individually is majorly to be focused in cyber security and in socially though individual user behaviours play a role efficiently in social cyber security. Semantic formalization in a light-weight for classifying the collected raw log data are in formulation of taxonomy in physical and logical. in sequence, lemmatization, sliding windows, feature selection, etc., are used to data performance. In midway determination of malware attacks are done potentially by incremental machine learning mechanisms on account of task complexity.

### III. PROPOSED METHODOLOGY

This research paper stages the hyper meta heuristic framework applying Cuckoo search algorithm to produce the configuration to devise the ESVM in the function of bi-objective optimization. It is very much essential to enhance the performance of the system with regard to accuracy and system complexity and hence ESVM is applied

to recognise and categorize the malware detection and intrusion detection. For feature selection so as to lessen complication of classification by the way of reducing the features, an enhanced K-means clustering is employed. The overview of the proposed system is represented in figure 1.

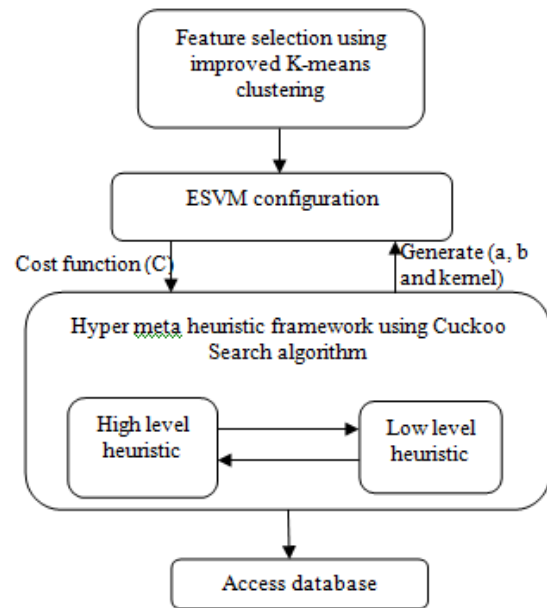


Figure 1. Overview of the proposed methodology

#### 3.1 Feature selection

Feature selection [17] is an indispensable methodology applied to decrease the dimensionality issue in data mining activity. Developing several data classification models established on the output from feature selection techniques helps to enrich the classification process's predictive performance. A cluster based feature selection utilizing improved K-means clustering algorithm is proposed in this paper.

The accountability of creating a collection of objects and arranging those objects is that the entities in an analogous group are much comparable with each other than to those objects existing in dissimilar groups (clusters). Clustering is exploratory information technique that organizes arranges the dataset information into a small number of groups. For the purpose of grouping the information, several grouping techniques are available. For different categories of information, different methods of computations need to be applied. For clustering analysis, K-means is the most commonly employed algorithm. Big data analytics incorporates various crucial data mining tasks including clustering that organizes the information into significant clusters taking in to consideration, the likeness or uniqueness between the objects. Towards evaluating the high dimensional datasets, primarily K-mean clustering algorithm is performed through the Hadoop and MapReduce for Cluster analysis technique. Once the unlabelled information is operated and exploited to group clusters of information, clustering is achieved in big data analytics. Data mining is an application programming approach that is used to evaluate the tremendously huge volume of information which is categorised in to structured, unstructured and semi-structured data.

The classification methods are strengthened by the supervised learning so as to execute the data and likewise the clustering algorithms are supported by unsupervised learning [18]. Following the assessment of traditional K-means algorithm few of the drawbacks like prominent error rate, variations in precision and some more setbacks are recognised. Taking in to account of these downsides, a novel approach is proposed with the aim of governing the disputes in conventional K-means algorithm. Primarily the information is pre-processed by the proposed method followed by the identifying the outliers from a group of related sets of information which is comprised of distinct elements. The information attained in the preceding stage is operated by a set of procedures and their results are analysed with the aid of suitable validation and verification techniques. Like so, an improved approach for enhancing the execution of K-means algorithm is proposed in this paper. The proposed research works ensures the effective resource utilization and improve the group leader selection procedure through which the issues of fluctuating precision and higher error rate is minimised basically. Furthermore the time consumption taken by the clustering procedure is much reduced by [19].

### 3.1.1 The improved k-means algorithm

Each datum point includes M estimations, which means that a data point can be considered as an ordered data set representing a record including M variables ( $V_{a1}, V_{a2}, \dots, V_{am}$ ). Among the M approximations of the data centres, p estimations are arbitrarily chosen. Restructure them in decreasing solicitation of necessity compatible with the essential merging of investigations, as ( $d_1, d_2, d_3, \dots, d_p$ ). The primarily estimation  $d_1$  is considered the fundamental estimation and rest are termed as optional estimations. The group size k is selected earlier.

Step 1: The variant is calculated as,  $H_i = \text{max}_i - \text{min}_i / k$ , for every dimension, where maximum value of  $i^{\text{th}}$  dimension is given by  $\text{max}_i$  and minimum value of  $i^{\text{th}}$  dimension by  $\text{min}_i$ .

Step 2: Applying the following set of circumstances, first the cluster is merged together. For any data point in case  $\text{min}_1 + j * H_1 \leq \text{value}_i < \text{min}_1 + (j+1) * H_1$ , then the data point of the group belongs to cluster j.

Step 3: Centroid data point is calculated for every one of the cluster. Centroid is calculated by taking the mean of all cluster points as expressed below:

$$\text{Centroid} = (\text{clp}_1 + \text{clp}_2 + \dots + \text{clp}_n) / n$$

Step 4: To compute the secondary dimension,  $2 \leq j \leq m$ , reiterate the following steps:

1. The outliers of every cluster is detected depending on the succeeding clauses.
2. For each data point if  $|\text{value}_{ij} - \text{value}_{cj}| > H_j$
3. Then the data point is considered as the outlier based on dimension j. Here,  $\text{value}_{ij}$  is the value of the  $i^{\text{th}}$  data point's  $j^{\text{th}}$  dimension, and  $\text{value}_{cj}$  is the value of the centroid's  $j^{\text{th}}$  dimension.
4. Distance from every centroid  $i^{\text{th}}$  to  $j^{\text{th}}$  dimension of outlier data point is computed.
5. As  $\text{distance} = \sum_{l=1}^{l=j} |\text{value}_{ij} - \text{value}_{cj}|$ .
6. Based on the least distance from all data points, select the cluster of the data point.

### 3.2 Ensemble SVM

Owing to the variety of the domains recognised from the follower of tweets, it is fascinate to analyse if a collection of machine learning algorithms is capable enough to utilize the various decision boundaries produced from the individual classifiers to deliberately merge the results of classification and thus a superior functioning is realized than is feasible with a single classifier. This research work concentrated on the potential of SVM assemblies [20] in categorizing the target audience from the followers list. In the succeeding sections, SVM is introduced in first which is then followed by two other sections illustrating the bootstrapping method and the algorithm applied for ensemble, respectively.

#### 3.2.1 SVM configuration

For two- or multi-class classification, the SVM is a supervised learning methodology and together with text categorisation it also has been effectively used in various applications. It pull out a known given set of  $\{+1, -1\}$  labelled training data through a hyperplane which is extremely distant from the positive and negative samples. This ideally parting hyperplane in the feature space relates to a non-linear decision boundary happening in the input space.

Consider a set of N distinct samples  $(x_i, y_i)$  with  $x_i \in \mathbb{R}^D$  and  $y_i \in \mathbb{R}^d$ . an SVM is modelled as

$$\sum_i a_i K(x, x_i) + b, i \in [1, N] \quad (1)$$

where  $K(x, x_i)$  is the kernel function, and a and b are the parameter and threshold of the SVM, respectively. The two goals to be improved ( $m=2$ ) can be expressed as follows:

$$\begin{aligned} \min_{s.t} F(x) &= |f_1(x), f_2(x)|, \quad f_1(x) = \text{error}, f_2(x) \\ &= NSV \end{aligned} \quad (2)$$

where  $f(x)$  is the cost function (C),  $\text{err}$  denotes the data sets number which are misclassified and  $NSV$  signifies the quantity of support vectors.

#### 3.2.2 Bootstrapping Using a Single SVM Model

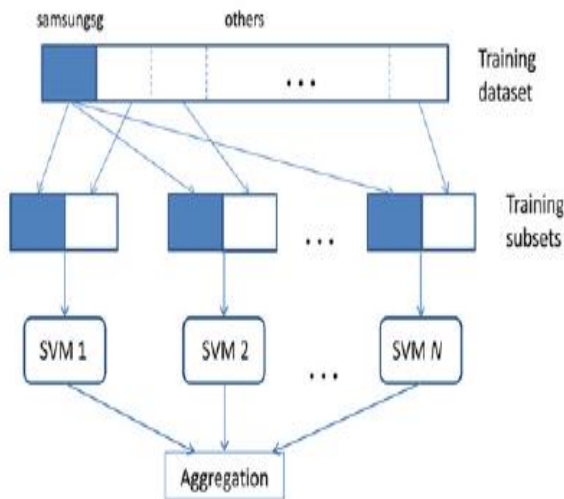
To tackle the issue of imbalance data via resampling of the minority class by means of replacement, bootstrapping [21] is a universally utilized method. Since our research be concerned about the temporal effect, the volume of tweets which we could attain is restricted to the number of tweets that are shared by the different owners within a period of 6-months. Therefore, it is impossible to gather enough samples to keep away from the difficulty of either running the risk of loss of information in the majority class or initiating unfairness in preference of the minority class.

As an alternative of traditional distributional assumptions, Bootstrap sampling bring into play a computation style, and it accepts a non-parametric approach to statistical inference with the purpose of achieving better estimation of the sample distribution instead of simply replicating the sample. A pseudocode description of bootstrapping is offered.

#### 3.2.3 Ensembles Using Multiple SVM Models

Since one of the attentions of this research study is to determine when the tweets that arise from the account owners could be utilized to find the target audience among the list of followers,

it is of importance to analyse if the ensemble of classifiers constructed using the training datasets of several domains can achieves superior than the general bootstrapping techniques defined prior. In spite of everything, the victory of an ensemble system hinge mostly on the variety of the classifiers that constitute the ensemble. The set of ensemble learning algorithms applied in this study comprises of bagging, stacking and majority vote. A common architectural depiction of the ensembles exploiting multiple SVM models is revealed in Figure 2. The approach of aggregation is dissimilar in each of the ensemble learning algorithms.



**Figure 2 A general architecture of the ensemble system using multiple SVM models**

Because of the various algorithms, several training datasets and configurations have been accepted for the function of manipulation of multiplicity from the different domains.

**3.3.4 Random Sampling with Majority Vote**

To split the majority class’s dataset into numerous subsets through random sampling prior to joining with the minority class to develop a balanced training dataset for the purpose of classification is one of the plainest solutions to the issue of imbalance data. Then the Individual classifiers are aggregated or combined by gathering a straightforward majority vote of their decisions.

In the place of minority class, here random sampling is carried out on the majority class which is in conflict to the bootstrapping algorithm. To partition the majority dataset into similar size subsets of the minority classis the objective here. Up until the needed size is reached, very record existing in the majority class is selected randomly and positioned in a subset. This random selection procedure is recurred till all the subsets are generated. Within the same subset, the records are distinctive but then duplication of records are located among several subsets.

**3.3.5 The majority vote algorithm**

**Input:**

- D: training dataset with labels denoting C classes
- L: learning algorithm
- W: training dataset labels
- N: number of L employed

**Do n=1 to N**

1. Call L with Dn and obtainthe classifier Ln.
2. Compare Wn with Cn createdfrom Ln, update vote.
3. Aggregate vote to the ensemble.

**End**

**3.4 Framework of hyper meta heuristics algorithm using CS**

To resolve the multi objective function, the current work presented the frame work of hyper meta heuristics algorithm and optimization of the parameters is carried out by Cuckoo Search algorithm. The cuckoo search (cuckoo search CS) algorithm is a metaheuristic algorithm recommended in the latest years [22]; Cuckoo Search (CS) is a new swarm intelligent optimization algorithm. Introductory research establish that cuckoo search algorithm is simple, easy to implement, effective and possess fewer number of parameters [23].For the optimization of SVM parameter, Cuckoo search algorithm is capable of offering a new approach. It holds several advantages such as an effective search path, robust global search capacity, lesser parameters and also is strong when multi-objective issues are resolved. It is proposed in this study, a system of dynamic measurement error prediction intended for sensor grounded on a CS-optimized support vector machine.

Cuckoos used to lay their eggs in other birds' nests during the host birds leave the nest unprotected. In the course of this process, few of these eggs that are alike the host bird's eggs, hatch and raise into adult cuckoos. In case the host birds find that the eggs are not their own, they will force out the strange eggs or desert their own nest and discover a different location to make again a new nest. Every egg present in the nest denotes are result, and a cuckoo egg stand for a new solution. The goal of the CS algorithm is to make use of the new and possibly improved solutions (cuckoos) to substitute the not-so-good solutions existing in the nests. The CS algorithm encompasses the following three rules [24]:

- 1) Every cuckoo can lay only one egg (single solution) at oneinstance, and it positions the eggs in a nestthat is randomly chosen
- 2) The finest nest among these nests, having high class eggs, (solutions) will movefurther to the subsequent generation
- 3) The total amount of accessible host nests is fixed. With the probability of p<sub>a</sub>, a host bird can identify an alien egg. The host bird may either force out the egg or abandon the nest and beginto create a new nest in a different location in this case.

On the basis of the above mentioned three rules, the bird nest locations are updated by the CS algorithm. Its search path can be defined as follows:

$$X_i^{t+1} = X_i^t + \alpha \oplus L, \tag{3}$$

where X<sub>i</sub><sup>t</sup>denotes the position of the ith nest at iteration t. Entry-wise multiplicationis represented by the product ⊕, and the step size is denoted by α, which is bound bya normal distribution. Levy random search path is indicated by L, which can be expressed as follows:

$$L = 0.01 \times \frac{\mu}{|v|^{\beta}} \times (g_{best} - X_i^t), \tag{4}$$

where  $g$  best signifies the current best nest. When  $\mu, v$  is subject to a normal distribution,  $\mu \sim N(0, \delta_\mu^2), v \sim N(0, \delta_v^2)$ , and

$$\begin{cases} \delta_\mu = \left\{ \frac{\Gamma(1+\beta)\sin(\pi\beta/2)}{\Gamma[(1+\beta)/2]\beta 2^{(\beta-1)/2}} \right\}^{1/\beta} \\ \delta_v = 1 \end{cases} \quad (5)$$

Where  $\beta = 1.5$ .

In comparison with more meta-heuristic algorithms, the CS algorithm bears two advantages. The first one is that the CS algorithm can further effectively uphold the balance amongst the local search strategy and the effective investigation of the whole search space. The second is that the CS algorithm possess two parameters alone (population size,  $N$ , and the probability of egg detection,  $p_a$ ).  $p_a$  alone regulates the balance between random and local search once  $N$  is fixed. Since the CS algorithm has lesser number of parameters, its universality is superior.

So as to optimize the SVM parameters  $a, b$  and kernel values, the CS algorithm is operated as follows:

1. The cuckoo search algorithm is initialized and set the quantity of nests,  $N$ , the probability parameters,  $p_a$ , the maximum iterations,  $t_{max}$ , and the ranges of  $C$  for effective intrusion detection of classifier.

2. The nest positions are randomly generated applying  $q_i^0 = [x_1^0, x_2^0, \dots, x_n^0]^T$ . Every nest corresponds to a set of parameters  $(C)$ .

3. Assess the fitness value of every nest, identify the existing sound solution, and document the least fitness value and its corresponding location.

4. Hold the finest solutions from the preceding generation, and using Formula (3), update the other nests' position. Analyse the new position's fitness value at that moment.

5. In case the fitness value of the new generation is finer than that of the previous generation, then change the best solution of the previous generation and make record the position of the best nest.

6. As the probability of detection of egg, establish a random number. Then with  $p_a$ , compare it. Modify the position of the nest randomly so as to acquire a new set of positions, if  $random > p_a$ .

7. In Step 6 identify the best nest position. When the maximum limit of iteration is attained, halt the process of searching, and provide as output the best position to realise the optimal parameter value ( $a, b$  and kernel); or else, go back to Step 3.

Thus the parameters of SVM are optimized by the proposed framework, accuracy is enhanced and the system complexity is diminished. The result analysis based on the proposed and existing method are discussed in detail in the following section.

#### IV. RESULTS AND DISCUSSION

In this cyber security protection technique, the comparison of the algorithms is done against the benchmarked ones to ensure the accuracy and model complexity. framework of Google Funf influences the Sher Lock data [25-26] collection agent. Media Lab of MIT develops a framework for processing data in the Framework of Funf Open Sensing specifically for mobile devices. Fun f was not developed for intensified frequent feature monitoring, like computing statistics on motion sensors. Hence, framework of pipeline processing are needed to get modify on stability and

robustness steadiness. Also, framework incorporated with probes gathers information on every running application . Physical or virtual data sources are obtained by Sensors (e.g., exterior temperature or memory intake). In general, two sort of sensors as PUSH and PULL are utilized.

Event-based PUSH sensors performs the sensation of SMS arrival or screen gets on. Gathered PULL sensors accelerate the device or sample of CPU. Sher Lock's gathering of data get stores in a temporary basis with format of JSON on the volunteer's device in the text format file. The file is zipped to ~50MB if size exceeds 500MB. Next, when the end user bond to Wi-Fi, and in sequence, temporarily stored the zip files on device are get into the server. Sher Lock dataset for cyber security research: like App Profiling & Malware Detection. malware detection and profiling of App are performed by activity of applications implicitly. Also, the dataset has several contextual features like the device location, movement, and utilization of battery, are used for improving malicious threat detection. Results of the newly introduced ESVMs classifier are then compared with Hyper-heuristic Support Vector Machines (HH-SVM), [27], Random Forest (RF) [28] and Gaussian Naive Bayes Tree (GNBT) [29] and Online Support Vector Machines (OSVMs) [30]. The typical performance metrics like accuracy, recall, precision, f-measure are considered for experimental analysis in proposed approach. table 1 exploits the metrics used for two class classification task based to confusion matrix.

| Classes |     | Predicted |    |
|---------|-----|-----------|----|
|         |     | Yes       | No |
| Actual  | Yes | TP        | FP |
|         | No  | FN        | TN |

Table 1: Confusion Matrix

Recall is defined as the ratio of actual positives that are predicted positive. Recall is defined as a function of the rightly classified cases (TPs) and the misclassified cases (FNs).

$$Recall = TP / (TP + FN) \quad (6)$$

Precision is defined as a measure of the accuracy given that a particular class has been predicted and specifies the proportion of the positives detected, which are actually right. Precision is the ratio of predicted positives that are real positive

$$Precision = TP / (TP + FP) \quad (7)$$

F1-score measure is a uniformly balanced precision and recall. The proportion of test precision and test recall are focused in score evaluation. It can be considered to be a average weight of precision and recall, in which best value of 1 is attained by F1-score and 0 as worst score.

$$F1\text{-score} = 2 * precision * recall / (precision + recall) \quad (8)$$

Classification accuracy is one among the most generic evaluation techniques used for measuring the system performance. It is utilized as the chief parameter of criteria to assess the performance of classification systems. The more the classification accuracy, the superior would be the system performance. The merit of this measure is in its simplicity and is computed with the expression  $Accuracy = (TN + TP) / (TN + FP + FN + TP)$  (9)



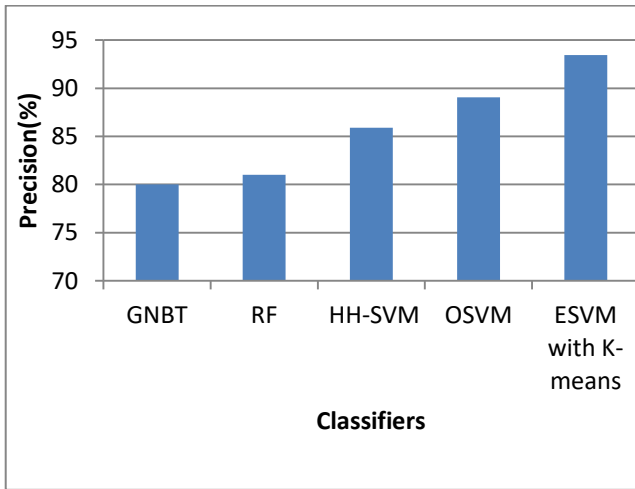


Figure 3 Precision Results of Classifiers Vs Comparison

the comparison results of precision in developed or proposed ESVM classifier along K-means, and the available classifiers like HH-SVM and OSVM classifiers illustrated in Figure 3. Figure 3 illustrates that the newly introduced ESVM classifier with K-means yields greater precision of 93.45%, while other classifiers like GNB, RF, HH-SVM and OSVM yield 80%, 81%, 85.89% and 89.05% correspondingly.

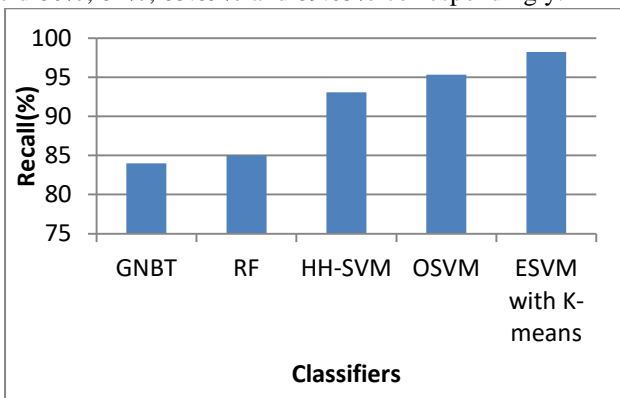


Figure 4 Recall Results Comparison vs. Classifiers

the comparison results of recall of the novel ESVM along K-means classifier, and the available classifiers like HH-SVM and OSVM classifiers correspondingly are illustrated in Figure 4. Figure 4 illustrates that the novel ESVM with K-means classifier yields greater recall of 98.25%, while other classifiers like GNB, RF, HH-SVM and OSVM yield 84%, 85%, 93.08% and 95.34% correspondingly.

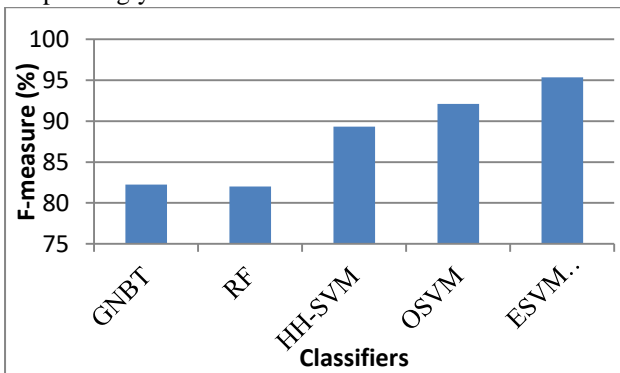


Figure 5 F-Measure Results Comparison vs. Classifiers

the comparison results of the f-measure of proposed ESVM along K-means classifier, and the available classifiers like HH-SVM and OSVM classifiers correspondingly are illustrated in Figure 5. Figure 5 reveals that the novel ESVM with K-means classifier yields a higher f-measure of 95.36%, while other classifiers like GNB, RF, HH-SVM and OSVM renders 82.25%, 82%, 89.34% and 92.09% correspondingly.

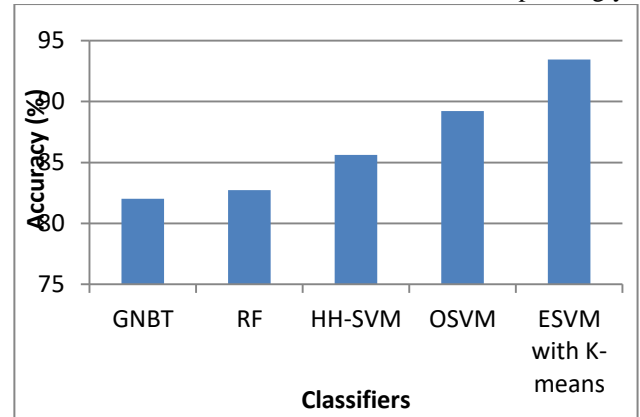


Figure 6 Accuracy Results Comparison vs. Classifiers

Figure 6 illustrates the results of the accuracy comparison of the newly introduced ESVM with K-means classifier, and the available classifiers like HH-SVM and OSVM classifiers correspondingly. Figure 6 reveals that the novel ESVM with K-means classifier yields a higher accuracy of 93.45%, while other classifier including GNB, RF, HH-SVM and OSVM renders 82.02%, 82.74%, 85.63% and 89.21 % correspondingly.

### V CONCLUSION AND FUTURE WORK

Intrusion detection and malware detection are eliminated in this research proposal by s Ensemble of Support Vector Machine (ESVM) for classification ensuring cyber security. Accuracy is improved in ESVM's bi-objective configuration. High and low level are considered in this research. Cuckoo Search algorithm is implemented to produce performance of ESVM parameters. Data are segregated by the SVM configuration. Accuracy and complexity of model are handled easily by means of the heuristics parameters efficiently. The goals are obtained and compared against several algorithms that performs well of operation and prove that the online support vector machine attains better result amid of all techniques. Data are considered individually so several combinations of data are possessed by the heuristic technique. Kernel methods are proposed in this developed framework and bi-objective optimisation is used for selection process effectively on account of dealing a huge data. Thus it achieves efficient result safeguard the protection from attackers by the independent optimization of model. standard checking are enhanced in future in the same model of real life scenarios and the several goals of safeguarding from threat of attackers. The proposed research could be enriched of potential model to contest any sort of attacks and dangers.6.

## REFERENCES

1. Teoh, T.T., Zhang, Y., Nguwi, Y.Y., Elovici, Y. and Ng, W.L., 2017, July. Analyst intuition inspired high velocity big data analysis using PCA ranked fuzzy k-means clustering with multi-layer perceptron (MLP) to obviate cyber security risk. In 2017 13th International Conference on Natural Computation, Fuzzy Systems and Knowledge Discovery (ICNC-FSKD) (pp. 1790-1793).
2. Senthamil Selvi, R. and Valarmathi, M.L., Enabling data security in data using vertical split with parallel feature selection using meta heuristic algorithms. *Concurrency and Computation: Practice and Experience*, p.e5248.
3. Anita, M. and Kumar, M.S., 2017. Security issues related to query phasing using metaheuristic algorithm. *International Journal of Engineering Science Inventio*. pp.22-27.
4. Peralta D, del Río S, Ramírez-Gallego S, Triguero I, Benítez JM, Herrera F. Evolutionary feature selection for big data classification: a mapreduce approach. *Math Probl Eng*. 2015;501. Article ID 246139.
5. Priya, T.M., and Saradha, A., 2018. An Improved K-means Cluster algorithm using Map Reduce Techniques to mining of inter and intra cluster data in Big Data analytics. *International Journal of Pure and Applied Mathematics*. Vol19 (7), pp. 679-690.
6. Dehestani, D., Eftekhari, F., Guo, Y., Ling, S.S., Su, S. and Nguyen, H.T. Online support vector machine application for model based fault detection and isolation of HVAC system. *International Journal of Machine Learning and Computing* 1(1) (2011) 66-72.
7. Gunantara, N., Putra, N. and Nyoman, I.D., 2019. The Characteristics of Metaheuristic Method in Selection of Path Pairs on Multicriteria Ad Hoc Networks. *Journal of Computer Networks and Communications*.
8. Hassani, K. and Jafarian, K., 2016. An intelligent method for breast cancer diagnosis based on fuzzy ART and metaheuristic optimization. In *XIV Mediterranean Conference on Medical and Biological Engineering and Computing 2016* (pp. 200-204). Springer, Cham.
9. Teoh, T. T., Nguwi, Y. Y., Elovici, Y., Ng, W. L., & Thiang, S. Y. (2018). Analyst intuition inspired neural network based cyber security anomaly detection. *International journal of innovative computing information and control*, 14(1), pp.379-386.
10. Khorram, T. and Baykan, N.A., 2018. Feature selection in network intrusion detection using metaheuristic algorithms. *International Journal of Advance Research, Ideas and Innovations in Technology*, 4(4), pp.704-710.
11. Saidala, R.K. and Devarakonda, N.R., 2017, April. A new parallel metaheuristic optimization algorithm and its application in CDM. In *2017 2nd International Conference for Convergence in Technology (I2CT)* (pp. 667-674). IEEE.
12. Tsai, C.W., Chiang, M.C., Ksentini, A. and Chen, M., 2016. Metaheuristic algorithms for healthcare: Open issues and challenges. *Computers & Electrical Engineering*, 53, pp.421-434
13. Bajpai, A. and Dayanand, A.A., 2018. Big Data Analytics in Cyber Security.
14. Sabar, N.R., Yi, X. and Song, A., 2018. A bi-objective hyper-heuristic support vector machines for big data cyber-security. *IEEE Access*, 6, pp.10421-10431.
15. Hou, S., Saas, A., Ye, Y., & Chen, L. (2016, June). Droiddelver: An android malware detection system using deep belief network based on api call blocks. In *International Conference on Web-Age Information Management* (pp. 54-66). Springer, Cham.
16. Nguyen, G., Nguyen, B. M., Tran, D., & Hluchy, L. (2018). A heuristics approach to mine behavioural data logs in mobile malware detection system. *Data & Knowledge Engineering*, 115, pp.129-151.
17. Lin, H.Y., 2013. Feature selection based on cluster and variability analyses for ordinal multi-class classification problems. *Knowledge-Based Systems*, 37, pp.94-104.
18. Krzanowski, Wojtek J. *Statistical principles and techniques in scientific and social research*. 2007, Oxford University Press on Demand.
19. Bikkur, Thulasi. *A Novel Multi-Class Ensemble Model for Classifying Imbalanced Biomedical Datasets*. 2017, IOP Conference Series: Materials Science and Engineering, vol. 225, no. 1. IOP Publishing.
20. Lo, S.L., Chiong, R. and Cornforth, D., 2015. Using support vector machine ensembles for target audience classification on Twitter. *PLoS one*, 10(4), p.e0122855.
21. Xin-She, Y and Deb, S, 2009. "Cuckoo Search via Levy flights," in *Proceedings of World Congress on Nature & Biologically Inspired Computing*. India: IEEE Publications, pp. 210-214
22. Yang, X.-S and Deb, S, 2014. "Cuckoo search: recent advances and applications," *Neural Computing and Applications*, vol. 24, no. 1, pp. 169-174.
23. Long, W., Liang, X., Huang, Y and Chen, Y, 2014. "An effective hybrid cuckoo search algorithm for constrained global optimization," *Neural Comput. Appl.*, vol. 25, no. 3, pp. 911-926.
24. Yang, X.-S. and Deb, S, 2013. "Multiobjective cuckoo search for design optimization," *Comput. Oper. Res.*, vol. 40, no. 6, pp. 1616-1624.
25. <http://bigdata.ise.bgu.ac.il/sherlock/#/>
26. <http://bigdata.ise.bgu.ac.il/sherlock/#/download>
27. Sabar, N.R., Yi, X. and Song, A, 2018. A Bi-objective Hyper-Heuristic Support Vector Machines for Big Data Cyber-Security. *IEEE ACCESS* 6 (2018)1-11.
28. Luba Gloukhov, Cody Wild, and David Reilly. *Malware classification: Distributed data mining with spark*. In *Association for the Advancement of Artificial Intelligence*, pages 1-6. [www.aaai.org](http://www.aaai.org), 2015.
29. Mahbod Tavallae, Ebrahim Bagheri, Wei Lu, and Ali A Ghorbani. A detailed analysis of the kdd cup 99 data set. In *Computational Intelligence for Security and Defense Applications, 2009. CISDA 2009. IEEE Symposium on*, pages 1-6. IEEE, 2009.
30. Mylavathi and Srinivasan, "A Meta-Heuristic Online Support Vector Machines for Big Data Cyber-Security", *Jour of Adv Research in Dynamical & Control Systems*, Vol. 11, 01-Special Issue, 2019