RESEARCH ARTICLE

# REDUCTION IN FALSE POSITIVE RATE BY COMBINING SVM AND KNN ALGO

## Sushil Kumar Mishra[1], Deepak Singh Rana[2]

[1]PG Student, Computer Science Engineering, Graphic Era Hill University, Uttarakhand, India

[2]Assistant Professor, Computer Science Engineering, Graphic Era Hill University, Uttarakhand, India

*ABSTRACT: With the emergence of internet, we have faced intrusion problem in computer system such as a file is self-replicating in nature and try to encrypt other file etc. This research work  is based on intrusion detection problem of network security. Its motive is to detect to detect a network behavior as normal or abnormal. In this research work, two different machine learning algorithms have been combined together to reduce its duplicity and redundancy to enhance accuracy of new developed algorithm. Its experimental result produces better result than other algorithm in terms of performance, accuracy and false positive. This combined algorithm has been applied on KDDCUP99 dataset to find better result by enhancing its performance, accuracy and reducing its false positive rate.*

*The objective of this proposed work is to predict the estimated false positive and accuracy with number of features and attributes by combining support vector machine (SVM) and K nearest neighbor (KNN) algorithm, to find the better and improved accuracy as compared to previous work done. As support vector machine (SVM) and K nearest neighbor (KNN) algorithms are combined to develop a new algorithm. In new developed algorithm, duplicity and redundancy are removed. By which, accuracy is enhanced and false positive rate is reduced.*

*Key Words:  KDDCUP99 dataset, Accuracy, False positive Rate.*

--------------------------------------------------------------------------***--------------------------------------------------------------------------

## 1.  INTRODUCTION

With the rapid growth of internet, there have increased in the number of attacks, intrusion detection system (IDS) has become a solution of information security. With the help of firewall we can facilitate some protection over a short period of time but they can not provide full protection. The main motive of intrusion detection system (IDS) is to help computer systems to deal with attacks.

An intrusion detection system (IDS) is able to monitor the behavior of all files those are coming in computer system if any file is suspicious or malicious so intrusion detection system (IDS) can detect that malicious file or attacks. Intrusion detection system (IDS) has created many models based on clustering to separate a normal and abnormal file. Intrusion detection system (IDS) provides the protection a network traffic from malicious file or virus, It basically maintains confidentiality and integrity of a computer system. Any unauthorized access of any data can not be done so secrecy of network traffic can be well maintained.

*On the basis of data intrusion detection system (IDS) are categorized to Host based intrusion detection system (HIDS) and Network based intrusion detection system (NIDS). In network based intrusion detection system (NIDS), only individual packet coming through the network are analyzed. The host based intrusion detection system (HIDS) analyzes above operation on a single host.*

*There are basically only two kinds of intrusion detection system used to detect intrusions in computer system*
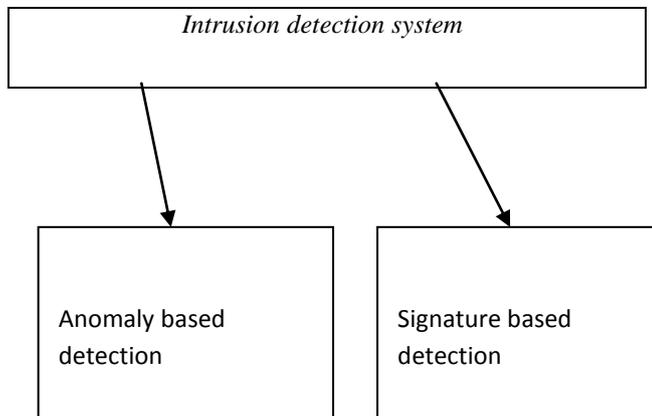


**Fig-1 Types of IDS**

## 1.1 Anomaly based detection

*Anomaly based detection system is based on heuristic rule which is used to monitor a normal and abnormal behaviors in a computer system. Anomaly based detection system contains many types of machine learning algorithm, these algorithm is used to identify a normal behavior and abnormal behavior in computer system.*
*The main drawback of anomaly based detection system is higher false positive rate and low accuracy.*

## 1.2 Signature based detection

*Signature based detection system can detect only predefined attacks, these attack or virus patterns are stored in database if a similar pattern and similar situation occurs so it is declared as attack. In signature based detection system, previous pattern data is stored in database. The computer virus, those are discovered. Its signatures are created and stored in database. If any file comes in a computer network so its signatures are matched with database. If file matches with virus signature so it is declared a computer virus otherwise a normal file. The main drawback of signature based detection system is that it can not detect a new attack.*

## 2. EXPERIMENTAL PARAMETERS

*There are the following parameters to be evaluated during this experiment*

**Performance:** *Performance is basically used to deal with achieving a target in more correct manner.*
*Performance = (TP)/(TP)+(TN)*

**Accuracy:** *Accuracy is the degree of closeness of measurements of a quantity to that quantity's actual (true) value.*
*Accuracy = (TP+TN)/(TP+TN+FP+FN).*

**False Positive Rate:** *False positive rate is used to falsely detect a normal file as abnormal file.*
*False positive rate = (FP)/(FP+TN).*

**Negative Predictive Value:** *Negative predictive value is used to detect a abnormal file as normal file.*
*Negative predictive value = TN/(FN+TN).*

## 3. EVALUATION DATA SOURCES

*Several parameters is evaluated by the standard data set KDDCUP99 provided by MIT laboratory. In KDDCUP99 data set, there are several types of attacks those can be categorized as normal and abnormal data.*

*MIT Lincoln laboratory has established a computer network. About 9 days, which monitors network traffic, containing normal and abnormal data.*

*KDDCUP99 dataset contains normal, remote to user, buffer overflow, user to root, guess_pawd and probe type of attacks.*

*Denial of service: Denial of service (DOS) intrusion is an attack in which, correct information can not be made available to legitmate receiver.*

*Denial of service intrusion can slow down computer system.*

*User to Root (U2R): User to Root (U2R) intrusion ia an attack in which, an attacker can access client's password without any authentication and access secret information from computer by using stolen password.*

*Remote to User (R2U): In this Intrusion, attacker transmits a packet over network which is not legitmate for particular network which overcrowds traffic in the network. Remote to User (R2U) can restart a computer system repeatedly.*

*Probe: In this intrusion, attacker has the ability to monitor all information which are being transmitted in a network and can access that information without any authentication.*

## 4. COMBINING SVM AND KNN ALGORITHM

*Support vector machine (SVM) is a machine learning method used for classification in which, a hyperplane is created by support vector for separating normal and abnormal data from each other. Support vector machine (SVM) can be divided into two phases-*

1- ***Training phase***
2- ***Testing phase***

***1-Training phase****: Support vector machine (SVM) is basically used to learn a large set of pattern from dataset. In dataset, there are several types of homogeneous pattern and heterogeneous pattern which can give better classification scheme for separating normal and abnormal data.*

***2-Testing phase:*** *Testing phase basically depends upon training phase can be done by support vector machine. In testing phase, machine learning algorithm is applied on data sets to get desired result.*

***K nearest neighbor*** *algorithm is a kind of machine learning algorithm which can provide a mechanism to solve traveling salesman problem (TSP).*

*By applying K nearest neighbor algorithm (KNN), false positive rate can be evaluated but it gives higher false positive rate which is inversely proportional to accuracy.*
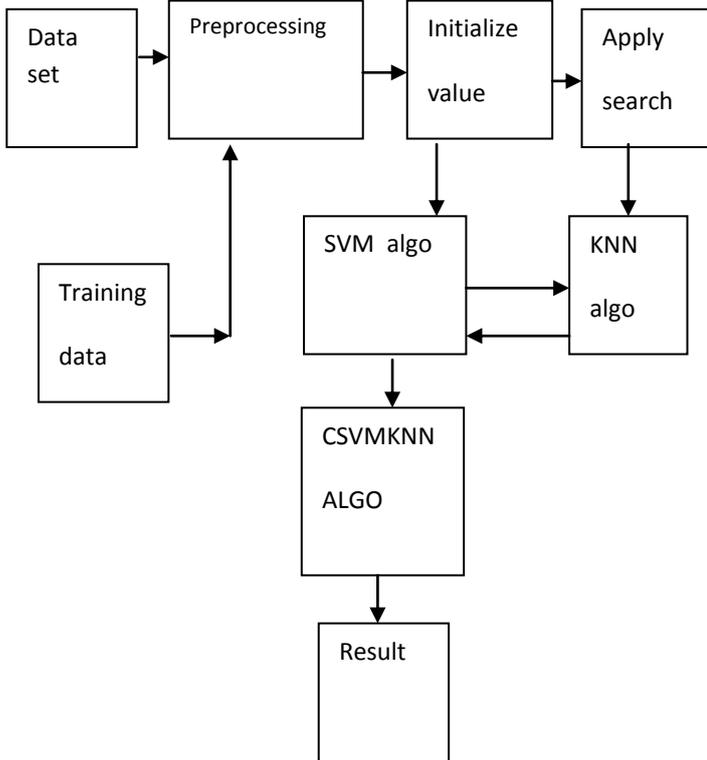
```
┌──────┐     ┌──────────────┐     ┌──────────┐     ┌──────────┐
│ Data │────▶│ Preprocessing │────▶│ Initialize│────▶│  Apply   │
│ set  │     │              │     │  value   │     │  search  │
└──────┘     └──────────────┘     └──────────┘     └──────────┘
                    ▲                   │                 │
                    │                   ▼                 ▼
┌──────────┐        │            ┌──────────┐     ┌──────────┐
│ Training │        │            │ SVM  algo│────▶│   KNN    │
│          │────────┘            │          │◀────│  algo    │
│  data    │                     └──────────┘     └──────────┘
└──────────┘                           │
                                       ▼
                                ┌──────────┐
                                │ CSVMKNN  │
                                │  ALGO    │
                                └──────────┘
                                       │
                                       ▼
                                ┌──────────┐
                                │  Result  │
                                └──────────┘
```

*Fig- 2  IDS using CSVMKNN*

*Algorithm1: SVM with KNN clustering*

*Input: Use training data set containing normal and abnormal data (Class type).*
*Output: classifier.*
*1 start*
*2 Randomly select data  from different class;*
*3 Generate of support vector machine classifier;*
*4 **While** number of iteration to add data to data set*
*5  Make available support vector to create hyperplane;*
*6 Hyperplane separate normal and abnormal data;*
*7 **Apply KNN clustering***
*8  Add points to the training dataset*
*9 Retrain support vector machine classifier using update     training dataset*
*10   **end***

*The active learning support vector machine process is introduced at this stage; support vector machine (SVM) is used to perform support vector machine training rely on various training data sets.The main goal of active learning support vector machine is how to choose the training data set for each training step. The separate hyperplane is used to tackle such issue; hyperplanes are created by support vectors. The hyperplane is tailored gradually by adding data points between maximum margins after each support vector machine training phase. Hence it produces more efficient selection strategy of data points.*

*Algorithm2: Training in CSVMKNN*

*Input: Training data set.*
*Input: T-Number of iteration.*
*Input: D-Maximum detection rate.*
*Output: Support vector machine (SVM) and K nearest neighbor (KNN) Classifier.*
*1 Start*
*2 Normalize data;*

*3 Let DR is detection rate , initially 0;*
*4 **While** DR < D do**
*5 **for i =1,…….. , T do***
*6 **SVM training phase;***
*7 **KNN Clustering Phase;***
*8 **end***
*9 Construct classifiers;*
*10 Do testing to updated DR;*
*11 **end***
*12 **end***

## 5.   EXPERIMENTAL  SETUP  AND RESULT

*Support Vector Machine (SVM) Classifier:*

*Support vector machine (SVM) is a group of supervised method used to apply on dataset (KDDCUP99) for classification and regression through which, false positive rate can be evaluated.*

*Support vector machine classifier is a method used to create a hyper plane by using support vector between data points. The main function of hyper planes is to separate normal and abnormal data from each other. Support vector machine (SVM) is also used for pattern identification. These methods have some specific application in speech analysis, image analysis etc.*
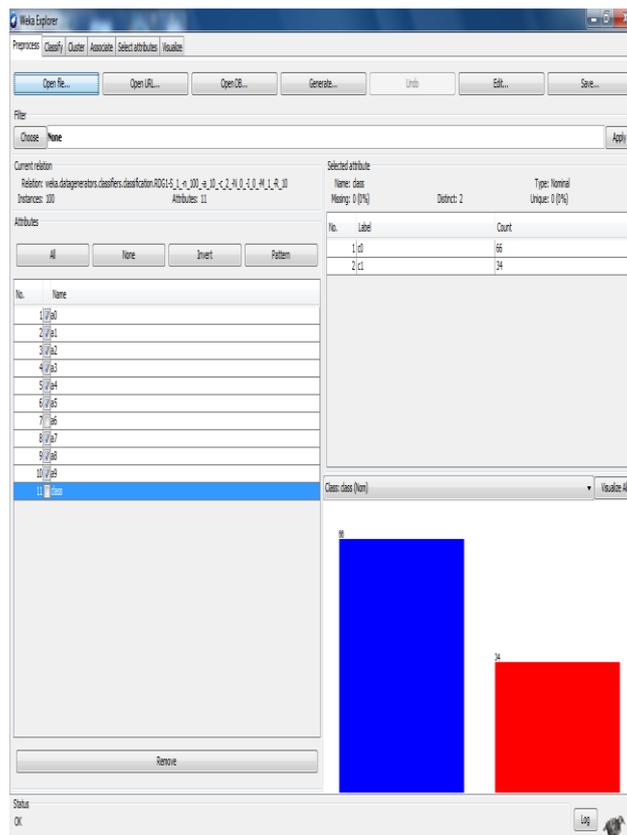


***Fig-3 SVM classifier***

*The optimal separation hyper plane is a linear classifier used with a maximum margin for finite pattern.*
*Suppose classification of different classes of same pattern is linear separable.*
*Linear separable hyper plane is w . x + b=0.*

$$\begin{cases} \mathbf{w} \cdot \mathbf{x}_i + b > +1 & \text{if} \quad y_i = +1 \\ \mathbf{w} \cdot \mathbf{x}_i + b < -1 & \text{if} \quad y_i = -1 \end{cases}$$

*These equations may be expressed in short form as*

$$y_i (\mathbf{w} \cdot \mathbf{x}_i + b) \geq +1$$

*Alternative*

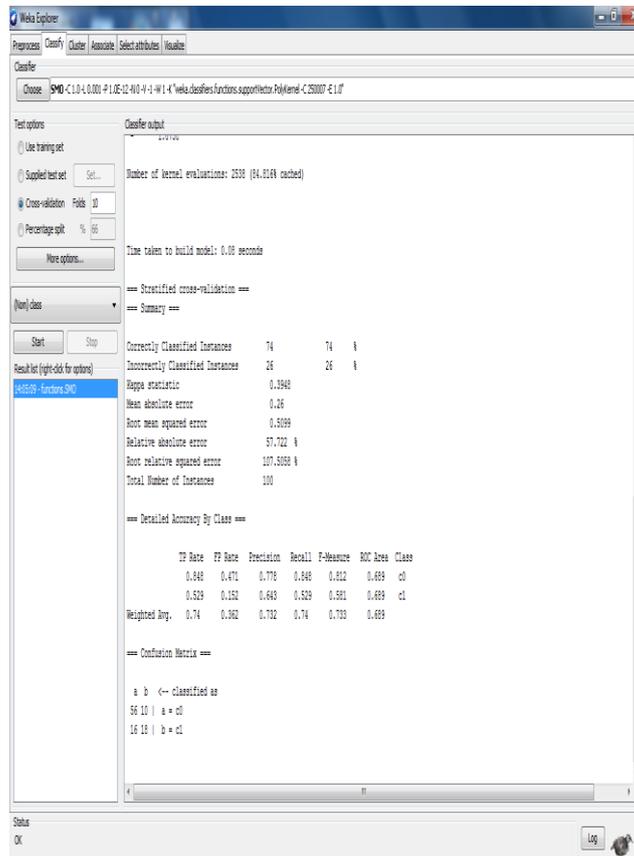$$y_i (\mathbf{w} \cdot \mathbf{x}_i + b) - 1 \geq 0$$



***Fig-4 SVM FPR***

**Table 1 with Experiment 1 (SVM)**

| False Positive Rate (FPR) |
| --- |
| 0.362 |

*K Nearest Neighbor (KNN) Classifier is used to summarize training dataset which can be used to find out new observation by reducing number of comparison.*
*K nearest neighbor (KNN) algorithm is applied on dataset to find out its false positive rate to evaluate higher accuracy.*
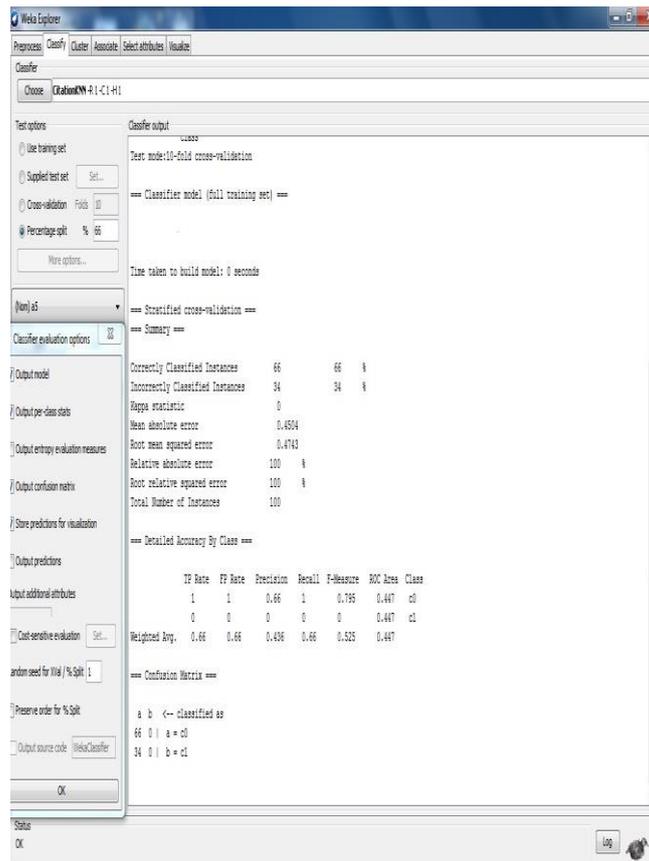


**Fig-5 KNN classifier**

**Table 2 with Experiment 2 (KNN)**

| False Positive Rate |
| --- |
| 0.66 |

*CSVMKNN classifier contains all attributes of support vector machine (SVM) and K nearest neighbor (KNN) algorithms. CSVMKNN algorithm is applied on dataset to produce accuracy, false positive rate.*

*False positive rate (FPR) can be reduced in CSVMKNN algorithm, by removing duplicate and redundant data from CSVMKNN algorithm.*
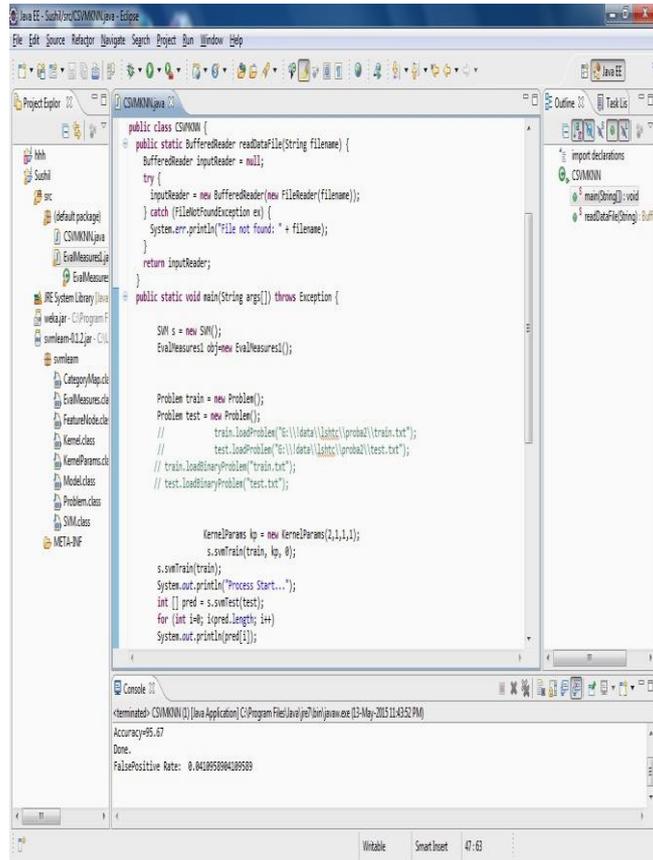
***Fig-6 CSVMKNN Classifier***

***Table 3 with Experiment 3 (CSVMKNN)***

| *False Positive Rate* |
|---|
| *0.041* |

***Table 4 with Comparison of False Positive Rate***

| Evaluation Measure | SVM | KNN | CSVMKNN |
|---|---|---|---|
| False positive Rate | 0.362 | 0.66 | 0.041 |

*CSVMKNN algorithm generates lesser false positive rate than Support vector machine (SVM) and K nearest neighbor (KNN) algorithms..*

## 6. CONCLUSIONS

*In this research work, a new machine learning algorithm based on classification is introduced by combining support vector machine (SVM) and K nearest neighbor (KNN) algorithms for solving intrusion detection problem.*
*In order to achieve higher accuracy, duplicity and redundancy are removed by reducing false positive rate.*

## REFERENCES

*[1]. pgale, Robert, Sheodoor schote, rengin and Christopher kruegel."A Literature analysis on automated malware analysis technique ."*
*[2]. Pargas, Rob Jonathan jarcy, Eleazar Aguirre Anaya , Samon Galeana Huerta and Alba Felix Moreno Hernandez,"Security controls for Android" In Computational Aspects of Social Networks (CASoN), 2012 Fourth International Conference on, pp.212-216,IEEE,2012*

*[3]. Blasing, Thomas, Leonid Batyuk, A-D.Schmidt, Seyit Ahmet Camtepe, and Sahin Albayrak." An android application sandbox system for suspicious software detection" In Malicious and Unwanted Software (MALWARE), 2010 5<sup>th</sup> International Conference on ,pp. 55-62 IEEE, 2010.*

*[4]. Johnson Ryan,  Zhaohui Wang , Corey Gagnon and Angelos  Stavrou." Analysis of Android Applications' Permissions. " In Software Security and Reliability Companion(SERE-C),2012 IEEE Sixth International Conference on, pp. 45 - 46.IEEE,2012.*

*[5]. Susan M. B. and Rayford B.V. (2000). Intrusion  detection via fuzzy data mining, Proceedings of the 12th Annual Canadian Information Technology,Ottawa, Canada, June 19-23, 2000, PP.109-122.*

*[6]. A Detailed Analysis of the KDD CUP 99 Data Set, Mahbod Tavallaee, Ebrahim Bagheri, Wei Lu, and Ali A.*

*[7]. Susan M. B. and Rayford B.V. (2000). Intrusion detection via fuzzy data mining, Proceedings of the 12th Annual Canadian Information Technology,Ottawa, Canada, June 19-23, 2000, PP.109-122.*

*[8]. A Detailed Analysis of the KDD CUP 99 Data Set, Mahbod Tavallaee, Ebrahim Bagheri, Wei Lu, and Ali A.*

*[9]. Blasing, Thomas, Leonid Batyuk, A-D.Schmidt, Seyit Ahmet Camtepe, and Sahin Albayrak." An android application sandbox system for suspicious software detection" In Suspicious and Unwanted Software (MALWARE), 2013 5<sup>th</sup> International Conference on ,pp. 55-62 IEEE, 2013.*

**1045**