



INDIAN INSTITUTE OF TECHNOLOGY ROORKEE

Advanced Computer Network Project

REPORT ON

Detection of DDoS using RDF-SVM

Submitted To:-

Dr. P Sateesh Kumar
Asst. Prof.
Computer Science & Engg

Submitted By:-

Himanshu Gupta (18535012)
Rishabh Mangain (18535023)
Shubham Bhatia (18535024)
Tushar Rangari (18535028)
Vishal Arora (18535029)

Contents

1 Problem Statement:-.....	3
2 Brief of the protocol.....	3
3 Methods of implementation:-	4
3.1 Random forest:-	4
3.2 SVM Algorithm:-.....	4
3.3 Proposed RDF – SVM algorithm.....	5
3.4 Proposed Algorithm.....	6
4 Results Achieved:	7
5 Remarks on result:	8
6 References:	8

1 Problem Statement:-

DDoS attack detection using RDF-SVM. Previously DDoS attacks were easily misled by flash crowd traffic but new method can detect DDOS attack using RDF-SVM algorithm. This algorithm is made to exploit random forest which uses important feature and SVM to rescreen it which will prevent from removing features taken mistakenly. At last is obtained with optimal features which will high detection rate and recall rate. We have used KDD99 dataset to train and test which distinguish between DDoS attack traffic and normal traffic which is flash crowd. RDF-SVM algorithm has much higher detection rate as compared to CART, neural network, logistic regression and ADABOOST.

2 Brief of the protocol:-

Denial of Service (DoS) forces the victim to receive false traffic which makes him fail to receive legal requests, resulting in providing no services to legitimate users. DDoS attacks have generally two important factors, i.e., resource consumption and bandwidth consumption. Bandwidth and resource are consumed by the control of botnets to generate many forged packets. DoS attacks are implemented by hackers on one or more than one targets and exhaust victim's resources, making the victim incapable of providing normal network services.

As the technology is developing day by day, DDoS attacks have been increasing year by year, making difficult to perform commercial operations, and disturbing people's normal life. Statistics have shown that there were repeated attacks on DNS root servers by DDoS resulting in large number of paralyzed servers. GitHub also encountered large traffic from DDoS attacks.

Tracing back for fake source IP DDoS becomes difficult. Moreover, detection can be easily misled by flash crowd traffic which is quite similar to the DDoS attacks.

A new method is proposed to detect DDoS attacks which is based on RDF-SVM, i.e., random forest-support vector machines. Random forest is used to calculate variable importance, whereas SVM algorithm is used to rescreen features. This algorithm will obtain a subset with higher detection rate.

Following are some of the innovations:-

1. Rescreening and preventing from deletion of features which contribute to DDoS detection.
2. Attacks are suppressed before reaching target host.
3. Differentiate between flash crowd traffic and DDoS traffic.
4. Detection of both unknown and known attacks.

3 Methods of implementation:-

3.1 Random forest:-

This algorithm is an integrated machine learning which uses random sampling technique and node splitting technique to construct multiple decision trees randomly and gain final classification result by voting. Random forests are an ensemble learning technique for classification, regression as well as other tasks that use decision trees.

This algorithm uses several classification models $\{(X, \theta_k) | k = 1, 2, \dots\}$ of decision tree. $\{\theta_k\}$ is assumed to be independent and identically distributed random variable and k is number of decision trees. Following are some important highlights of this algorithm.

- N samples are selected randomly and then they are put back after the selection. This is done in order to train the decision tree.
- There are M attributes to each sample, when we split the decision tree node is split. Then m attributes are selected ($m \ll M$) and the information gain strategy is used to select an attribute as split attribute of the node.
- Each node in step 2 is split, until it can't be splitted further.
- Build decision trees without pruning
- Then the vote on the variable X_i is collected from decision trees, then the highest votes is set as label for classification.

Entropy tells about the disorder-ness. Information Entropy

$$\text{Entropy}(x) = \sum_{i=1}^c -P_i \log_2 P_i \quad (1)$$

X is a random variable of finite random discrete variable, P_i is the ratio of the random variable and the dataset C. The expectation of feature Y for x is:

$$\text{Entropy}_Y(x) = \sum_{j=1}^N \frac{|x_j|}{|x|} \text{Entropy}_Y(x) \quad (2)$$

The rise in information entropy leads to the rise in uncertainty of training samples. Information gain is as follows:

$$\text{Gain}(x, Y) = \text{Entropy}(x) - \text{Entropy}_Y(x) \quad (3)$$

Entropy is used as information gain. Feature is directly related to the information gain. Thus feature makes contributions to the detection rate.

3.2 SVM Algorithm:-

The SVM algorithm is a 2 step classification algorithm and it performs good in generalization and handling with imbalanced dataset.

For the training set

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}, y_i \in \{-1, 1\},$$

the basic step is to find a hyper-plane. We can represent the hyper plane by the linear equation

$$\omega^T x + b = 0 \quad (4)$$

ω is normal vector set, b is displacement. The distance r of any point x of sample space to hyper-plane (ω, b) is expressed as follows:

$$r = \frac{|\omega^T x + b|}{\|\omega\|} \quad (5)$$

We assume that the samples are trained correctly by the hyper-plane. For $(x_i, y_i) \in D$, then

$y_i = +1$, so $\omega^T x + b > 0$. If $y_i = -1$, then $\omega^T x + b < 0$.

Then we obtain the most suitable variables ω and b , thus the optimize SVM problem is as follows

$$\max_{\omega, b} \frac{1}{2\|\omega\|^2} \quad (6)$$

$$\text{s.t. } y_i(\omega^T x_i + b) \geq 1, i = 1, 2, \dots, m \quad (7)$$

For dual hyper-plane problem the optimum solution $Q(a)$ is :

$$Q(a) = \sum_{i=1}^n a_i - \frac{1}{2} \sum_{i,j=1}^n a_i a_j y_i y_j (x_i \cdot x_j) \quad (8)$$

Finally optimal classification function $f(x)$ is

$$f(x) = \text{sgn}\{\sum_{svm} a_i y_i (x_i \cdot x) - b\} \quad (9)$$

3.3 Proposed RDF – SVM algorithm

The $F(x_i)$ function define as

$$F(x_i) = \begin{cases} 1, & x_i \text{ is positive feature} \\ -1, & x_i \text{ is negative feature} \end{cases} \quad (10)$$

We obtain 2 values. $P_j^{(1)}$ denotes the training and prediction with deletion and $p_j^{(2)}$ denotes prediction without deletion. The value sum_P denotes the difference

$$sum_P = P_j^{(1)} - P_j^{(2)} \quad (11)$$

Then main difference is found for n samples by the equation

$$Mean_sumpre_diff = (\sum_j (P_j^{(1)} - P_j^{(2)})) / n, j = 1, 2, \dots \quad (12)$$

If $Mean_sumpre_diff > t$, then $X_i \in C$, the feature will be retained. Otherwise, $X_i \notin C$, then it will be removed.

3.4 Proposed Algorithm

Algorithm 1 RDF-SVM algorithm

Input: $\{x_i | i = 1, 2, 3 \dots, N\}, S_{train}, S_{test}$

Output: Optimal feature subset C

Steps:

1: initial $\omega \leftarrow \phi$, $sum_P \leftarrow 0$, $Mean_sumpre_diff \leftarrow 0$

2: for each x_i do

3: $\omega_{x_i} \leftarrow \text{Randomforest}(x_i)$;

4: end for

5: $W \leftarrow \text{Sort}(\omega_{x_i})$;

6: for each $\omega_{x_i} \in W$ do

7: if $\omega_{x_i} > \alpha$ then

8: $A_index \leftarrow x_i$;

9: else

10: $B_index \leftarrow x_i$;

11: end for

12: end for

13: for each $X_i \in B_index$ do

14: for $j \in (1, n)$ do

15: $P_j^{(1)} = \text{SVM}(y_i, x_i, y_{train}, X_{test})$;

16: $P_j^{(2)} = \text{SVM}(y_i, A_{x_i}, y_{train}, A_{X_{test}})$;

17: $sum_P = sum_P + P_j^{(1)} - P_j^{(2)}$;

18: $j++$;

19: end for

20: $Mean_sumpre_diff = sum_P / n$;

21: if $Mean_sumpre_diff > t$ then

22: $A_index \leftarrow x_i$;

23: else

24: $\text{remove}(x_i)$;

25: end for

26: return A_index ;

4 Results Achieved:

Analysis on KDD99 dataset

We have the KDD99 dataset and it consists of 42 features and these features belong to normal and attacked dataset. We further have 2 sets

1. KDD Test+ set which contains 39 kinds of attacks which is 17 kinds more than KDD Train+.
2. KDD Train+ set which contains 22 kinds of attacks.

42 features are:

“duration”, “protocol_type”, “service”, “flag”, “src_bytes”, “dst_bytes”, “land”, “wrong_fragment”, “urgent”, “hot”, “num_failed_logins”, “logged_in”, “num_compromised”, “root_shell”, “su_attempted”, “num_root”, “num_file_creations”, “num_shells”, “num_access_files”, “num_outbound_cmds”, “is_host_login”, “is_guest_login”, “count”, “srv_count”, “error_rate”, “srv_error_rate”, “error_rate”, “srv_error_rate”, “same_srv_rate”, “diff_srv_rate”, “srv_diff_host_rate”, “dst_host_count”, “dst_host_srv_count”, “dst_host_same_srv_rate”, “dst_host_diff_srv_rate”, “dst_host_same_src_port_rate”, “dst_host_srv_diff_host_rate”, “dst_host_error_rate”, “dst_host_srv_error_rate”, “dst_host_error_rate”, “dst_host_srv_error_rate”, “label”]

. Then we have applied RDF_SVM algorithm on the test data to find the best features (src_bytes, dst_bytes, logged_in, duration, protocol_type, and num_file_creations) out of these which can detect DDOS attack.

The RDF-SVM algorithm prevents from removing the feature subset from which classification is promoted. We can detect unknown and known attacks, differentiate real ip attacks from random ip attacks and flash crowd more effectively as compared to other methods.

After applying Random Forest Algorithm, the detected features out of all 42 are with their importance values

Detected Feature	Importance value
Duration	0.013450231133692056
src_bytes	0.5563474527516045
dst_bytes	0.32320092844470516
logged_in	0.03535644214871085
Hot	0.010177058552346224
num_access_files	0.0009165270390737588
srv_count	0.05868177970868924
is_guest_login	0.001869580221178301

We took all those features whose importance value is greater than a (0.20). A_index contains

Detected Feature	Importance Value
src_bytes	0.5563474527516045
dst_bytes	0.32320092844470516

On applying SVM classifier on the remaining features we captured one more feature which got missed out in Random forest classification.

Detected Feature	Importance Value
src_bytes	0.5563474527516045
dst_bytes	0.32320092844470516
logged_in	0.03535644214871085

5 Remarks on result:

We have Random Forest on data set KDD99 with the importance set threshold value 0.22 and SVM on remaining features with the threshold value of 0.0000000000005. Using these features we have tried to differentiate between DDoS attack and flash crowd traffic.

6 References:

<https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=8089926>