



Improving the Accuracy of ODS

A. Venkatesh¹, P. Remanth Reddy², B. Vinay Kumar³, U. Pavan Achari⁴, M. Venkata Sunil⁵

^{1,2,3,4,5} Department of Computer Science and Engineering, Santhiram Engineering College, Nandyal

ABSTRACT

Intrusion detection could be very crucial for supplying security to extraordinary community domain names and is by and large used for finding and tracing the intruders. There are many issues with conventional intrusion detection models (IDS) including low detection functionality in opposition to unrecognised community attack, excessive fake alarm fee and insufficient evaluation functionality. Hence the foremost scope of the studies on this area is to broaden an intrusion detection version with progressed accuracy and decreased schooling time. This paper proposes a hybrid intrusion detection version via way of means of integrating the main component evaluation (PCA) and aid vector machine (SVM). The novelty of the paper is the optimization of kernel parameters of the SVM classifier the usage of automatic parameter choice approach. This approach optimizes the punishment factor (C) and kernel parameter gamma (γ), thereby enhancing the accuracy of the classifier and lowering the schooling and trying out time.

The experimental consequences acquired at the NSL-KDD and gurekddcup dataset display that the proposed technique plays higher with better accuracy, quicker convergence velocity and higher generalization. Minimum sources are ate up because the classifier enter calls for decreased characteristic set for maximum classification. A comparative evaluation of hybrid fashions with the proposed version is likewise performed.

Keywords: cross validation, dimensionality reduction, intrusion detection system, principal component analysis, radial basis function kernel, support vector machine

1. Introduction

Intrusion Detection Systems (IDS) are developed to become aware of unauthorized tries to access or control the laptop structures. IDS collects community facts to become aware of one of a kind types of malware and assaults in opposition to offerings and applications. IDS has been labeled into foremost categories, specifically signature primarily based totally detection and anomaly primarily based totally detection. In signature primarily based totally IDS, assault sample of intruders are modeled and the device will notify as soon as the suit is identified.

All recognized assaults are identified with decreased fake wonderful rate. Signature databases need to be up to date regularly for you to become aware of the brand new assault sample. However, anomaly detection structures create a profile of regular activity. Any sample that deviates from the regular profile is dealt with as an anomaly. Hence even unknown assault styles are identified with none guide intervention

Data mining strategies had been used these days withinside the improvement of intrusion detection models to reduce data overloading. These fashions extract the beneficial know-how with the aid of using searching for styles and relationships from the statistics collected, thereby enhancing choice making. Data mining technology inclusive of neural networks [1], naïve bayes networks [2], genetic algorithms [3], fuzzy logic [4] and guide vector machine [5] are used for category and sample popularity in lots of industries as they have got advanced the overall performance of the fashions that installation such algorithms. In category, the capabilities of newly gift items are examined and are assigned to one of the present set of classes. Classifier fashions advantage know-how from the schooling statistics and discover the elegance label for the brand new instances. Many supervised mastering fashions are used to clear up category problems.

Support Vector Machine (SVM) is one of the green strategies used because the generalization functionality is better even if the pattern training information is small. In the latest years, many hybrid smart structures had been proposed to enhance the accuracy in contrast to individual strategies.

Anomaly detection fashions have the issue of "curse of dimensionality" that's a completely critical issue. To triumph over this issue, an optimal characteristic subset must be acquired to imshow accuracy and get rid of noise. In the proposed model, a SVM classifier is mixed with PCA for figuring out the anomalies. PCA is one of the considerably used statistical strategies to lessen the dimensionality and SVM has the advantage of accomplishing best overall performance for the category of strange patterns.

The rest of the paper is prepared as follows: Section 2 discusses numerous device studying strategies, SVM strategies utilized in numerous intrusion detection fashions and latest hybrid strategies advanced integrating SVM and dimensionality reduction. The history of numerous strategies used withinside the version is disstubborn in Section 3. The proposed method is mentioned in Section 4. The experiments and consequences of the version are stated in Section 5. Section 6 incorporates the conclusion.

2. Background

In this segment we in brief evaluation the statistics mining strategies which are hired in our proposed model.

2.1 Scaling

Huge volumes of community visitors need to be processed for figuring out the anomalies and as a result type might not be accurate. Therefore, facts packets go through normalization method in which facts is sanitized. The cause of normalization is to file the facts to a numerous scale. Different strategies used for normalization are Z-score, Decimal and Min-Max scaling. The Min-Max normalization method is selected for the proposed version because it has much less range of misclassification errors [29] compared to different strategies. Min-Max normalization achieves a linear amendment at the unique facts. Normalization is done for a given range. To carry out mapping for a price v of a characteristic f inside range $[\min f, \max f]$ to a brand new range $[\text{new_min } f, \text{new_max } f]$, the calculation is given by

$$v' = \frac{(v - \min_f)(\text{new_max}_f - \text{new_min}_f) + \text{new_min}_f}{\max_f - \min_f} \quad (1)$$

where v' is the new value in the specified range. The benefit of this technique is that all values are concealed within certain ranges.

2.2 Principal Component Analysis (PCA)

PCA is one of the widely used statistical techniques with inside the discipline of records mining to lessen dimensionality and to perceive records factors with the very best feasible variance [30] [31]. Lakhina et al. [32] hired PCA to distinguish community site visitors records into regular and anomalous sub regions. In this method, the focal point is on detection of quantity primarily based totally anomalies in origin-destination go with the drift aggregated in spine networks and it's miles a crucial issue inside numerous IDS systems today. The PCA technique identifies anomalous site visitors quantity on a selected hyperlink via way of means of evaluating it with beyond values. Thus, PCA separates hyperlink site visitors measurements into sub areas representing regular and odd site visitors. The final results of the PCA is to venture a function area onto a smaller subspace that represents records via way of means of decreasing the size of function area. This reduces computational expenses and the mistake of parameter estimation.

The standard PCA approach can be summarized in six simple steps:

- (i) Determine the covariance matrix of the normalized d -dimensional dataset.
- (ii) Determine the eigenvectors and eigenvalues of the covariance matrix.
- (iii) Sort the eigenvalues in descending order.
- (iv) Select the k eigenvectors that correspond to the k largest eigenvalues where k is the number of dimensions of the new feature subspace.
- (v) Construct the projection matrix from the k selected eigenvectors.
- (vi) Transform the original dataset to build a new k -dimensional feature space.

2.3 Support Vector machine Classification Model

Support Vector Machines (SVMs) are a hard and fast of supervised gaining knowledge of strategies especially used for classification, outlier detection and regression. The essential blessings of SVM are:

- Efficient results in high dimensional spaces
- Helpful wherein the quantity of dimensions is higher than the quantity of data samples
- Efficient usage of memory as SVM uses only a subset of training points in the decision making function
- Different kernel functions can be used for the decision function. Common kernels are available, but we can also develop custom kernels

2.3.1. Linear Support Vector Machine

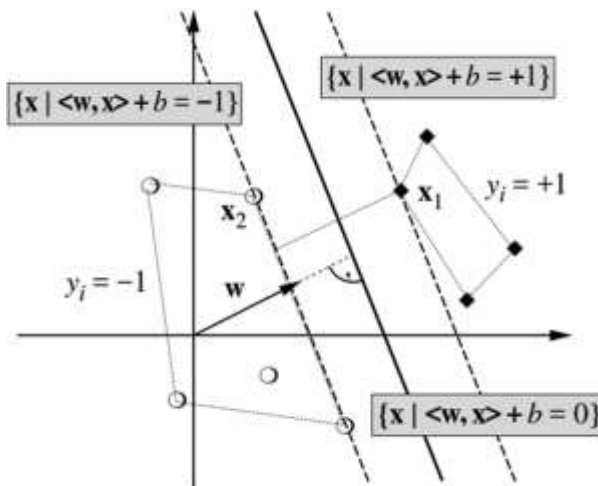
Consider the categorization of lessons that may be separated in a linear style as proven in Figure 1 [34]. Figure 1 indicates that the hyperplane for the linear classifier is of the form $(w \cdot x + b) = \text{zero}$ having the most margin (each expressions $\hat{w} \cdot \hat{x}$ and $w \cdot x$ denote the scalar made from vectors, i.e. constitute the equal operation). The classifier is defined via way of means of the set of pairs (w, b) , wherein w is a weight vector and b is the bias, that may specify the inequality for any pattern x_i within the education set, whilst y_i represents the magnificence label:

$$\begin{aligned} w \cdot x_i + b &\geq +1 \text{ if } y_i = +1 \\ w \cdot x_i + b &\leq -1 \text{ if } y_i = -1 \end{aligned} \quad (2)$$

Here, in which w is normalized with admire to a fixed of factors x such that: $\min_i |w \cdot x_i| = 1$. Minimizing $\|w\|_2$ situation to equation (2) and representation of constraints in a compact shape are as follows

$$\begin{aligned} y_i (w \cdot x_i + b) &\geq +1 \\ y_i (w \cdot x_i + b) - 1 &\geq 0 \end{aligned} \quad (3)$$

Figure 1. Linear support vector machine [34].



Every hyperplane (\mathbf{w}, b) is a classifier that separates all patterns from the training set. To deal with non separable case, the problem is rewritten as: Minimize

$$\|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i \quad (4)$$

(where C is a regularization parameter: small C allows constraints to be easily ignored, large C makes constraints hard to ignore) with respect to

$$y_i (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i, \text{ where } \xi_i \geq 0 \quad (5)$$

Hence the decision function is of the form

$$f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b \quad (6)$$

The fundamental equation is obtained by

$$\min P(\mathbf{w}, b) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum H_1(y_i f(\mathbf{x}_i)) \quad (7)$$

where $P(\mathbf{w}, b)$ represents the primal formulation to minimize the training error, and H_1 would identify the number of errors.

2.4. Cross Validation

Cross validation is a version assessment method that makes use of a partial facts set for education with a view to be utilized by the learner. Certain facts are eliminated earlier than the education begins. The facts already eliminated is used as a check set to degree overall performance of the version on the brand new facts.

In K-fold move validation, the statistics set is split into okay subsets. In each iteration, one of the okay subsets is taken into consideration because the take a look at set and the alternative okay - 1 subsets are taken into consideration because the schooling set. The benefit of this approach is that the statistics factor may be withinside the take a look at set as a minimum as soon as and withinside the schooling set okay - 1 times.

3. Proposed Work

3.1 Preprocessing of Dataset

This paper analyzes the National Scientific Laboratory–Knowledge Discovery and Data Mining (NSL-KDD) dataset [35] and gurekd- dcup [36] for the experiments. Preprocessing is finished with the aid of using normalization of the discrete at- tributes into non-stop ones with the aid of using Min-Max tech- nique on each datasets. Every community statistics has forty one attributes in which 34 attributes are continu- ous and seven attributes are discrete in nature. Pre- processing is finished at the dataset after which the statistics is split into schooling and check sets

3.2 Principal Component Analysis

Figure four suggests the little by little evaluation of feature vector era the use of PCA. A normalized characteristic matrix is acquired after preprocessing and is fed as enter to achieve the imply and covariance of the man or woman features. Eigen- vectors are generated for each characteristic and the best eigenvalues and respective eigenvectors are retained (withinside the vector set) to achieve the most suitable characteristic subset.

Algorithm for Obtaining the Optimal Feature Subset Using PCA

Input (Training Set, Test Set)

Output (Optimal Training Set, Optimal Test Set) Step 1: Determine the size of training and test data. Step 2: Scale the training and test data.

Step 3: Subtract from each feature x its respective mean m

$$m = \frac{\sum_{k=1}^n x_k}{n}$$

wherein x_k specifies the individual element of x and n denotes the number of elements.

Step 4: Determine the covariance matrix C

$$C = \frac{X(m)XT(m)}{n}$$

where $X(m)$ represents the feature matrix after sub- tracting the respective means, XT is the transpose matrix, and n is the total number of elements.

Step 5: Determine the eigenvectors v_j and reign values λ_j of the covariance matrix C

$$Cv_j = \lambda_j v_j \quad j = 1, \dots, p, p \leq n$$

Step 6: Obtain a feature vector using a set of reigen- values $(\lambda_1, \lambda_2 \dots \lambda_p)$ and respective eigenvectors, where λ_1 is the highest eigenvalue. Select k such eigenvectors that match the largest k eigenvalues in the set. The computational complexity of the PCA is $O(p^2n + p^3)$ p is the variety of functions and n is the variety of statistics points. Covariance matrix computation is $O(p^2n)$ and the corresponding eigenvalue decomposition is $O(p^3)$.

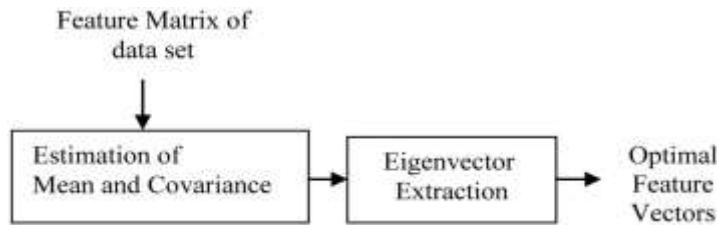


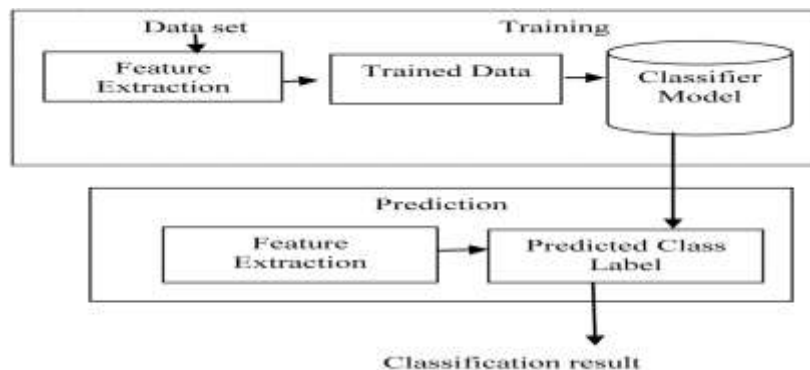
Figure 4. Optimal feature subset generation using PCA.

3.3 Support Vector Machines

Figure five indicates the tiers of SVM classifier for predicting the elegance label of the community traffic. The tiers are divided into phases: Training and prediction.

In education segment, the function matrix is fed in to the classifier version to iden- tify the elegance label. The checking out segment obtains the getting to know regulations from education segment to iden- tify the sample of the unknown traffic

Figure 5. Predicting the class label using support vector machine.



4. Conclusion

This paper proposes an intrusion detection model integrating Principal Component Analysis (PCA) and Support Vector Machines (SVM) using RBF kernel. Dimensionality reduction using PCA removes noisy attributes and retains the optimal attribute subset. SVMs construct classification models based on training data obtained from PCA. Optimization of SVM parameters C and γ for RBF kernel by proposed automatic parameter selection technique reduces the training and testing time and produces better accuracy. Two different datasets NSL-KDD and gurekddcup were applied to the model to analyze the performance. The experimental results indicate that the classification accuracy of the proposed model outperforms other classification techniques using SVM as the classifier with PCA as the dimensionality reduction technique. Minimum resources are consumed as the classifier input requires reduced feature set and thereby minimizes training and testing overhead time.

References

- [1]. G. Wang et al., "A new approach to intrusion detection using artificial neural networks and fuzzy clustering", *Expert Syst. Appl.*, vol. 37, pp. 6225–6232, 2010. <http://dx.doi.org/10.1016/j.eswa.2010.02.102>
- [2]. W. Wang and R. Battiti, "Identifying intrusions in computer networks with principal Component analysis", in *Proceedings of the First International Conference on Availability Reliability and Security (ARES'06)*, 2006, pp. 270–279. <http://dx.doi.org/10.1109/ARES.2006.73>
- [3]. K. Shafi and H. A. Abbass, "An adaptive genetic based signature learning system intrusion detection", *Expert Syst. Appl.*, vol. 36, no. 10, pp. 12036–12043, 2009. <http://dx.doi.org/10.1016/j.eswa.2009.03.036>
- [4]. S. Srinoy et al., "Anomaly based intrusion detection using fuzzy rough clustering", in *Paper Presented at the International Conference on Hybrid Information Technology (ICHIT'06)*, 2006, pp. 329–334.
- [5]. L. Khan et al., "A new intrusion detection system using support vector machines and hierarchical clustering", *Journal on very large databases*, vol. 16, no. 4, pp. 507–521, 2007.
- [6]. M. V. Mahoney and P. K. Chan, "Learning non stationary models of normal network traffic for detecting network attacks", in *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining KDD'02*, New York, NY, USA, pp. 61–72.
- [7]. M. Mahoney and P. Chan, "Learning models of network traffic for detecting novel attacks", Florida Institute of Technology, Technical report CS-2001–2, 2002.
- [8]. L. Weijun and L. Zhenyu, "A method of SVM with normalization in intrusion detection", *Procedia Environmental Sciences*, pp. 256–262, 2011. <http://dx.doi.org/10.1016/j.proenv.2011.12.040>
- [9]. A. C. Catania et al., "An Autonomous labeling approach to support vector machines algorithms for network traffic anomaly detection", *Expert Systems with Applications*, pp. 1822–1829, 2012.
- [10]. Snort [Online]. Available: <https://www.snort.org>
- [11]. C. A. Kumar, "Analysis of Unsupervised Dimensionality Reduction Techniques", *Computer Science and Information Systems*, vol. 6, no. 2, pp. 217–227, 2009. <http://dx.doi.org/10.2298/CSIS0902217K>
- [12]. S. Thaseen and C. A. Kumar, "An Analysis of supervised tree based classifiers for intrusion detection systems", in *2013 International Conference on Pattern Recognition, Informatics and Mobile Engineering (PRIME)*, Salem, 2013, pp. 294–299.
- [13]. F. Kuang et al., "A novel hybrid KPCA and SVM with GA model for intrusion detection", *Applied Soft Computing*, pp. 178–184, 2014. <http://dx.doi.org/10.1016/j.asoc.2014.01.028>
- [14]. L. S. Thaseen and C. A. Kumar, "Intrusion Detection Model using fusion of PCA and optimized SVM", in *Proceedings of 2014 International Conference on Computing and Informatics (IC3I)*, Mysore, India, 2014, pp. 879–884.
- [15]. S. Thaseen and C. A. Kumar, "Intrusion Detection Model using fusion of chi-square feature selection and multi class SVM", *Journal of King Saud University – Computer and Information Sciences*, in press, 2016.