

# Intrusion Detection Using KDD Dataset

Divya Chaganti  
Saili Sahasrabuddhe  
Ashwin Viswanathan  
Vignesh Miriyala  
Sanil Sinai Borkar

December 16, 2015

## 1 Introduction

Our project is to use a data mining approach to develop an intrusion detection system (IDS) using the KDD dataset [1].

## 2 Feature Extraction

The original KDD dataset contains 42 features including the label. Since working this size of data may not be feasible always from the perspective of a real-time IDS, we had to reduce the number of features used. This dimensionality reduction was achieved using Principal Component Analysis (PCA).

34 of the 42 features of the dataset were continuous on which PCA was applied. The result of applying PCA is placed in the file *PCA\_Result.txt*.

Out of all the 34 continuous features, we picked up the principal components that contained positive weights for most of the *error rate* indicating features. Out of the 34 candidate principal components, we narrowed down it to 2 principal components, and selected the one that had more positive weights. This led us to finalize our principal component which is specified below:

Table 1: Final Feature Set  
Description

Feature	Description
num_access_files	number of operations on access control files
error_rate	% of connections that have “SYN” errors
srv_error_rate	% of connections that have “SYN” errors
srv_error_rate	% of connections that have “REJ” errors
same_srv_rate	% of connections to the same service
diff_srv_rate	% of connections to different services
dst_host_diff_srv_rate	
dst_host_srv_diff_host_rate	
dst_host_error_rate	
dst_host_srv_error_rate	
dst_host_error_rate	

$$\begin{aligned}
Comp23 = & 0.124 * num\_access\_files - 0.131 * error\_rate + 0.141 * srv\_error\_rate - \\
& 0.490 * srv\_error\_rate + 0.107 * same\_srv\_rate + 0.126 * diff\_srv\_rate - \\
& 0.154 * dst\_host\_diff\_srv\_rate + \\
& 0.263 * dst\_host\_srv\_diff\_host\_rate - 0.122 * dst\_host\_error\_rate + 0.162 * \\
& dst\_host\_srv\_error\_rate + 0.717 * dst\_host\_error\_rate
\end{aligned}$$

Based on this principal component, 11 features were selected out of the available 41 as shown in Table 1 as per the description given in [2].

Since the IDS needs to tag each network packet as either malicious (attack packet) or benign (normal packet), the type of attack is not significant but only the presence of an attack packet is. Therefore, as a pre-processing step, the normal packets were labeled as ‘normal’ and everything else was tagged as an ‘attack’ packet.

### 3 Classification

Decision Tree classifier was used to classify the instances. Out of the 494022 packets, 80% of it was used as training data (395218 records), and the remaining 20% was used for testing (98803 records).

The time taken to train a decision tree classifier model was around **14 seconds** on an average, and accuracy was **95.5318%**. The summary of the classifier model is given below:

=== Summary ===

Correctly Classified Instances	377559	95.5318 %
Incorrectly Classified Instances	17659	4.4682 %
Kappa statistic	0.8471	
Mean absolute error	0.083	
Root mean squared error	0.2037	
Relative absolute error	26.2432 %	
Root relative squared error	51.2282 %	
Coverage of cases (0.95 level)	99.8133 %	
Mean rel. region size (0.95 level)	81.0894 %	
Total Number of Instances	395218	

=== Confusion Matrix ===

a	b	<-- classified as
316384	1011	a = attack.
16648	61175	b = normal.

The prediction for the 98803 records was done in **0.2 seconds** with an accuracy of **95.46%**. The confusion matrix for the prediction is given below:

Confusion Matrix and Statistics

	Reference	
Prediction attack. normal.		
attack.	79084	4223
normal.	264	15232

Accuracy : 0.9546  
95% CI : (0.9533, 0.9559)  
No Information Rate : 0.8031  
P-Value [Acc > NIR] : < 2.2e-16

Kappa : 0.8445  
McNemar's Test P-Value : < 2.2e-16

Sensitivity : 0.9967  
Specificity : 0.7829  
Pos Pred Value : 0.9493  
Neg Pred Value : 0.9830  
Prevalence : 0.8031  
Detection Rate : 0.8004  
Detection Prevalence : 0.8432

Balanced Accuracy : 0.8898

'Positive' Class : attack.

## 4 Association

## 5 Clustering

## 6 Conclusion

## References

- [1] “KDD Cup 1999 Data,” <http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html>, accessed: 2015-11-28.
- [2] “KDD Cup 1999 Feature Set,” <http://kdd.ics.uci.edu/databases/kddcup99/task.html>, accessed: 2015-11-28.