

An Analysis of Supervised Tree Based Classifiers for Intrusion Detection System

Sumaiya Thaseen
School of Computing Science and Engineering, VIT
University,
Chennai, India
sumaiyathaseen@gmail.com

Ch. Aswani Kumar
School of Information Technology and Engineering, VIT
University,
Vellore, India
Cherukuri@acm.org

Abstract— Due to increase in intrusion incidents over internet, many network intrusion detection systems are developed to prevent network attacks. Data mining, pattern recognition and classification methods are used to classify network events as a normal or anomalous one. This paper is aimed at evaluating different tree based classification algorithms that classify network events in intrusion detection systems. Experiments are conducted on NSL-KDD 99 dataset. Dimensionality of the attribute of the dataset is reduced. The results show that RandomTree model holds the highest degree of accuracy and reduced false alarm rate. RandomTree model is evaluated with other leading intrusion detection models to determine its better predictive accuracy.

Keywords— Classification Models; Discretization; Feature Selection; Intrusion detection system; RandomTree;

I. INTRODUCTION

Intrusion detection research includes both misuse detection and anomaly detection. Misuse detection requires a learning algorithm trained by a dataset in which each instance is identified as a normal class or an anomaly. This algorithm cannot identify novel attacks not included in the training set but can learn the new attacks through a new training dataset[9]. Anomaly detection is the process of constructing models of normal network events and identifying the events that deviate from these models. Anomaly detection suffers from high false alarm rate as many unobserved normal events are also considered as anomalies [24] [25]. These limitations led to an increasing interest in intrusion detection techniques based upon data mining.

Data mining is used by many researchers to differentiate attack data from normal data by using various outlier detection methods[10]. Classification is used to identify the class of label of instances based on the features(attributes) in a dataset. Scholars have tested different classifier models to the intrusion detection problem, such as rule-based detection[11], Neural Networks [12,13,14], fuzzy logic [15], hidden markov model[16], random forest model [17] Bayesian analysis [18] and data mining[19] [20]. In this paper, we describe the use of machine learning techniques that aids the analysts to automatically generate rules used for computer network intrusion detection. We evaluate tree classifier algorithms to construct structured data. Given a structured data set, a list of attributes that describes every data element and a set of categories that partition the data, the tree classifier algorithms

determine which set of attributes most accurately categorizes the data.

Feature selection is an approach to remove redundant or irrelevant features from the data to increase classification accuracy and decrease computational costs[20]. It is the process of choosing a subset of original features that optimally reduces the feature space to evaluation criterion. The raw data is usually large so a subset of data set is used to create feature vectors that represent most of the information present in data. Feature selection methods typically fall into three broad categories: (1) Embedded (2) Wrapper and (3) filter methods [26]. Embedded methods are embedded in specific mining methods such as random forests in which the importance of each feature is estimated. Wrapper methods use feedback obtained from specific classifier to determine the quality of feature subset. Filter methods rely on the general characteristics of training data. Wrapper methods are more accurate than filter methods and they are also less expensive and do not rely on an explicit classifier. In this paper, we have considered two different feature selection approaches such as Consistency subset evaluation, correlation feature subset evaluation. Analysis is carried out using the following tree based classifiers ADTree, J48, LADTree, NBTree, RandomForest, RandomTree, REPTree. A supervised filter is applied at the preprocessing stage. NSL-KDD data set is used to train and test the tree based classifiers.

This paper is organized as follows: Section 2 reviews the related work. Section 3 discusses the materials and methods required for our analysis. Section 4 explains and discusses the experiments and results. Finally section 5 concludes the paper.

II. RELATED WORK

Classification task in intrusion detection system is widely discussed in the literature. Chebrolu et al [3] investigated the performance of two feature selection algorithms involving Bayesian networks (BN) and Classification Regression Trees (CRC), and developed the ensemble of both methods. Furthermore, Tsang et al. [4] used genetic-fuzzy rule mining approach to evaluate the importance of IDS features. Annur et al [5] implemented a hybrid statistical approach using data mining and decision tree classification to identify the false alarms. The proposed IDS detected threats and benign traffic in critical network applications. In [6], the authors proposed support vector machine approach to classify network anomalies. Two hybrid approaches for modeling IDS were proposed by

Abraham et al [7] namely hierarchical hybrid intelligent system model (DT-SVM) and an ensemble approach combining the base classifiers. The authors concluded that their research proved accurate for intrusion detection systems. An ensemble of ANNs, MARS and SVMs [8] proved superior to individual techniques for intrusion detection in terms of classification accuracy.

A network intrusion detection system based on hidden naïve bayes classification was proposed by Levent et al [1] wherein the data mining model can be applied to intrusion detection data that has high correlated features and high volume of network data stream. A decision tree using wrapper approach was proposed by Siva et al [2]. Neuro tree approach implemented by the authors resulted in better detection accuracy. Optimal features were evolved maximizing specificity and sensitivity of Intrusion Detection Systems

III. MATERIALS AND METHODS

A. Data Set

NSL-KDD data set is used for our analysis as it overcomes the problems suggested in [22]. This dataset is believed to be an effective data set benchmarked to compare different intrusion detection methods. The number of training and testing data set in NSL-KDD are reasonable. Hence the experiments are carried out on the complete data set rather than randomly selecting a smaller portion. NSL-KDD has the following advantages in comparison with the original KDD data set

1. Redundant records are not included in the training set as a result classifiers are not biased towards frequent records.
2. Better detection rate on the records as there are no
3. data set for training and testing purposes. duplicate records in the test sets.
4. Evaluation results carried out by different researchers will be comparable as it is affordable to run the experiments given the reasonable size of training and test data set. Hence our experiments will be carried out using NSL-KDD 20% labeled data set for training and testing purposes.

Figure 1 shows a conceptual framework of Intrusion Detection System. The framework receives regular and anomalous traffic pattern as input and performs supervised discretization. The discretized output is sent to each of the tree based classifying models to classify the anomalies based on a feature selection technique. A decision analyzer is used to classify the normal and anomalous instances of our dataset. Finally a performance evaluation is done based on performance and error metrics.

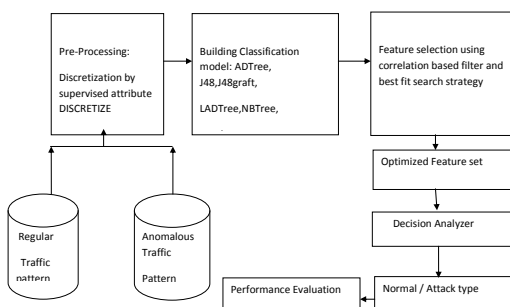


Figure 1. Conceptual Framework of Intrusion Detection System

B. Data Pre-processing

In order to make the data set more efficient, supervised filtering is used to reduce the data set. After streamlining the raw data, a small sample data set needs to be chosen representing the entire data set. Supervised attribute filtering technique called Discretize has been applied to improve the performance of the IDS. Discretization is used to partition or convert continuous attributes features or variables to nominal attributes, features or variables. Better models can be produced by discretization of continuous attributes.

C. Classification

Tree based classifiers such as ADTree, C4.5, J48graft, LADTree, NBTree, RandomTree, RandomForest, REPTree are used for analysis in our paper. Previous studies show that tree based classifiers have a superior performance in comparison with rule based and function based classifiers[1] [17]

I) ADTree

Alternating Decision Tree is one of the machine learning methods for classification. It generalizes decision trees and has connectivity with boosting. ADTree consist of decision nodes and prediction nodes. Decision nodes specify a predicate condition. Prediction nodes contain a single number. ADTrees always have prediction nodes as both root and leaves.

II) C4.5

C4.5 [24] is a statistical classifier that constructs decision trees using information entropy. The attributes of the data set is carefully chosen such that it effectively splits set of samples into subsets enhancing any one of the class. Splitting criterion is the difference in entropy. Highest information gain attribute is identified to make decision. A pruned or unpruned C4.5 decision tree is generated using J48graft decision tree

III) LADTree

LADTree learning algorithm applies induces an alternating decision tree by applying the logistic boosting algorithm. Each training instance has a working response with respect to mean value of instances [23]. The working response has to fit to the mean value of instances in a particular subset by reducing the least square value between them.

IV) NBTree

NBTrees are shown to scale large databases and outperform decision trees and naïve Bayesian classifiers. This approach is viable for inducing classifiers as many attributes are considerable for classification; attributes need not be independent; importance of interpretability of classifier;

V) RandomTree

A random tree is a tree formed by stochastic process. Types of random trees include Uniform spanning tree, Random minimal spanning tree, Random binary tree, Random recursive tree, Treap, Rapidly exploring random tree, Brownian tree, Random forest and branching process.

VI) RandomForest

RandomForest [17] is an ensemble classifier consisting of many decision trees and outputs the mode of classes as output by individual trees. This method combines bagging idea and random selection of features independently to construct a collection of decision trees with controlled variation. It is one of the highly accurate classifier for many datasets. Thousands of input variables can be handled without any variable deletion.

VII) REPTree

REPTree builds a decision/regression tree using information gain/variance and prunes it using reduced error backfitting. Numeric attributes are only sorted. Missing values are also dealt by splitting the corresponding instances into parts.

D. Feature Evaluation

Feature evaluation of attributes is performed using the following methods (1) correlation based feature selection (2) consistency based filter. Correlation based feature selection selects the feature sets containing features highly correlated with class and uncorrelated with each other. Redundant features are removed because of high correlation with remaining feature set. Consistency based filter method (Dash & liu,2003; Huan & Setiono,1997) generates in each round a random subset. A comparison of random subset is done with the current best subset .If new subset is more consistent than current best subset, latter is replaced by new set(Hall,1999).

IV. EXPERIMENTS AND RESULTS

NSL-KDD data set is preprocessed by extracting 41 features from the tcpdump data. It includes the full training set, the 20% training set and test set. The full training set has 125973 connections. The training set has 25192 instances containing the normal and anomaly classes. The test set contains 11850 connections.

The experiments were carried out using default values of the parameters for all tree based classifiers in the Waikato Environment for Knowledge Analysis (WEKA)[5]. We calculate the detection rate and false positive rate by changing the label to 1 or 0 where a '1' specifies anomaly class and a '0' specifies normal class. We improve the performance by employing a supervised attribute filtering technique namely Discretize during the preprocessing phase. Classification model is built by the following 8 classifiers: ADTree,C4.5, J48graft,LADTree,NBTree , RandomTree, RandomForest, REPTree.Feature selection were carried out using consistencysubset(CONS) evaluator and correlationfeature subset (CFS) evaluator. Table 2 shows the combination of features selected for CONS and CFS.CFS selects 8 features from the 41 feature data set whereas CONS selects 10 features from the entire data set. Table 1 shows the different combination of features selected by CFS and CONS attribute evaluator. As the performance metric results were very close for both the evaluators we tend to show the results for CFS as it has a higher performance with reduced feature set. We applied a 10-fold cross validation to accurately represent the training data and build classifier models. Training data is divided into 10 subsets of equal size randomly. In every iteration one of the

subsets is selected for testing and remaining 9 subsets are used to train the classifier.

Feature Evaluator	Features Selected
CFS	protocol_type dst_bytes num_root
CONS	duration service src_bytes dst_bytes hot count diff_srv_rate dst_host_srv_count dst_host_same_srv_rate dst_host_diff_srv_rate dst_host_same_src_port_rate dst_host_rerror_rate

a) A. Classifier Performance Metrics

b) We evaluate our classifiers by measuring their performance.All measures of performance are based on the following values resulted from the training and test set.These values are true positives (tp) ,false positives (fp) ,true negatives (tn) and false negatives (fn).Every application has different statistics computed from the entries.

c)

d) I) Accuracy $a : (tp + tn) / n$;

Proportion of the total number of predictions that were correct.

e) II) Precision $p : tp / (tp + fp)$;

Proportion of the predicted positive cases that were correct

f) III) Recall $r : tp / (tp + fn)$;

Proportion of positive cases that were correctly identified.

g) IV) F- measures: $p + r / (2 pr)$

Average of the information retrieval precision and recall metrics.

h) Table 2 gives the detailed information of the execution time required to build the model in training and test data set.From table 2 it is inferred that RandomTree model and REPTree classifies the training and test instances very quickly in comparison to other models. Tables 5 and 6 depict the various performance metrics evaluated for the training and test data set.It is inferred from table 5 that RandomForest model has the highest precision and lowest false alarm rate.It is inferred from table 6 that NBTree and RandomForest has the highest precision whereas NBTree has the lowest false alarm rate.

TABLE II. DETAILED INFORMATION OF THE DATASET

Algorithm	Time taken to build model in training data (in secs)	Time taken to build model in test data (in secs)
ADTree_DIS_CFS	18.91	5.14
C4.5_DIS_CFS	4.01	1.34
J48graft_DIS_CFS	5.93	2.16
LADTree_DIS_CFS	86.4	20.76
NBTREE_DIS_CFS	71.55	37.8
RandomTree	0.36	0.0321
RandomForest	3.3	1.4
REPTree	1.19	0.37

TABLE IV. ERROR MEASURES FOR EACH OF THE SUPERVISED ALGORITHMS IN TEST DATASET

Algorithm	MAE	RMSE	KS
ADTree_DIS_CFS	0.1613	0.2397	0.7631
C4.5_DIS_CFS	0.0389	0.1552	0.9022
J48graft_DIS_CFS	0.0412	0.163	0.8924
LADTree_DIS_CFS	0.1103	0.215	0.8057
NBTREE_DIS_CFS	0.0243	0.1322	0.9252
RandomTree	0.0321	0.0321	0.8926
RandomForest	0.0348	0.1405	0.9157
REPTree	0.0369	0.1471	0.9075

C. B. Error Metrics

The following are the different error metrics used to evaluate each algorithm.

I) Mean Absolute Error(MAE):

a) Average of absolute errors and specifies how the predicted values are different from the real values. The closer the predicted value to the real value, the smaller the MAE.

II) Root Mean Squared Error (RMSE):

Average difference between prediction value and real value. It is more sensitive to outliers in comparison to MAE.

3) Kappa Statistics (KS) :

Statistical measure which specifies the consistency between the predicted and real value in a dataset. The higher the value the more consistency between predicted and real values resulting in a better performance of the algorithm. Table 3 and 4 show the various error metrics analyzed in the training and test data set. It is inferred from table 3 that RandomTree has the least MAE and highest Kappa Statistic value whereas from Table 4 it is inferred that NBTree has the least MAE and highest Kappa Statistic but as the time required to build a model using NBTree is high. RandomTree is an appropriate model for classifying network anomalies in a minimal span of time with higher accuracy.

TABLE III. ERROR MEASURES FOR EACH OF THE SUPERVISED ALGORITHMS IN TRAINING DATASET

Algorithm	MAE	RMSE	KS
ADTree_DIS_CFS	0.0627	0.1357	0.9625
C4.5_DIS_CFS	0.0064	0.0651	0.9911
J48graft_DIS_CFS	0.0062	0.0644	0.9914
LADTree_DIS_CFS	0.0444	0.1319	0.9539
NBTREE_DIS_CFS	0.0033	0.0533	0.9932
RandomTree	0.0049	0.0697	0.9901
RandomForest	0.006	0.0478	0.9949
REPTree	0.0068	0.0653	0.9908

TABLE V. RESULTS OF EFFICIENCY MEASURES FOR THE SEVEN CLASSIFIERS IN THE TRAINING DATA SET

Algorithm	Recall	Precision	F-measure	False Alarm rate
ADTree_DIS_CFS	0.981	0.982	0.981	0.02
C4.5_DIS_CFS	0.996	0.996	0.996	0.004
J48graft_DIS_CFS	0.996	0.996	0.996	0.004
LADTree_DIS_CFS	0.977	0.977	0.977	0.024
NBTREE_DIS_CFS	0.997	0.997	0.997	0.004
RandomTree	0.995	0.995	0.995	0.005
RandomForest	0.997	0.997	0.997	0.003
REPTree	0.995	0.995	0.995	0.005

TABLE VI. RESULTS OF EFFICIENCY MEASURES FOR THE SEVEN CLASSIFIERS IN THE TEST DATA SET

Algorithm	Recall	Precision	F-measure	False Alarm rate
ADTree_DIS_CFS	0.934	0.933	0.932	0.224
C4.5_DIS_CFS	0.971	0.971	0.971	0.078
J48graft_DIS_CFS	0.969	0.968	0.968	0.093
LADTree_DIS_CFS	0.943	0.942	0.943	0.146
NBTREE_DIS_CFS	0.978	0.978	0.978	0.048
RandomTree	0.968	0.968	0.968	0.074
RandomForest	0.975	0.975	0.975	0.059
REPTree	0.972	0.973	0.972	0.061

TABLE VII. VALIDATION RESULTS FOR THE OVERALL CLASSIFIER PERFORMANCE RANKED BY ACCURACY

Algorithm	Accuracy	Error rate
RandomTree_DIS_CFS	0.9974	0.254
NBTree_DIS_CFS	0.9962	0.3374
J48graft_DIS_CFS	0.9957	0.4287
C4.5_DIS_CFS	0.9955	0.4406
REPTree_DIS_CFS	0.9954	0.4565
RandomTree_DIS_CFS	0.995	0.4922
ADTree_DIS_CFS	0.9813	0.1865
LADTree_DIS_CFS	0.977	0.229

TABLE VIII. TEST RESULTS FOR THE OVERALL CLASSIFIER PERFORMANCE RANKED BY ACCURACY

Algorithm	Accuracy	Error rate
NBTree DIS CFS	0.9776	0.0223
RandomTree DIS CFS	0.9749	0.025
REPTree DIS CFS	0.9724	0.0275
C4.5 DIS CFS	0.9722	0.0287
J48Graft DIS CFS	0.9686	0.031
RandomForest DIS CFS	0.968	0.031
LADTree DIS CFS	0.9428	0.0571
ADTree DIS CFS	0.9344	0.0655

TABLE IX. COMPARISON OF THE OVERALL CLASSIFIER PERFORMANCE SORTED BY ACCURACY

Model	Accuracy
RandomTree DIS CFS	0.9749
HNB PKI INT [1]	0.9372
JRip[27]	0.9230
NBTree[27]	0.9228
LBk[27]	0.9222
SVM[28]	0.9218
J48[27]	0.9206
MLP[27]	0.9203
Decision Table[27]	0.9166
SMO[27]	0.9165
BayesNet[27]	0.9062
OneR[27]	0.8931
Naïve Bayes[27]	0.7832

C. Discussion

We summarize the obtained results from the evaluation conducted in the previous sections. The results indicate that the execution time of RandomTree algorithm is higher for training and test data sets in comparison with NBTree, LADTree and ADTree. The error of the predicted values for ADTree, LADTree are higher in comparison with RandomTree and J48graft. RandomTree performs more accurately in determining the values with respect to data samples.

From the accuracy perspective, RandomTree, RandomForest and REPTree detects anomalies powerfully with fewer false rate alarms compared to other algorithms.

To sum up, from the execution and accuracy point of view, RandomTree can be identified as the best choice for analysis and detection model among all the other classifier algorithms. Table 7 represents the validation results based on the training data and Table 8 represents the test data result sorted by accuracy. It is inferred from Table 7 and Table 8 that RandomTree and NBTree have the highest accuracy when experiments were carried out on training and test data set. Finally Table 9 compares the different classification models with our approach based on accuracy.

Based on the results of our studies, the random tree model with supervised discretize method and correlation feature selection method has good predictive accuracy. Random Tree

model provides an advantage that with a reduced feature set a better predictive performance in the intrusion detection set.

V. CONCLUSION

The main goal of this paper is to evaluate eight tree based classifier algorithms to classify network events based on supervised discretization and feature selection. We summarized the experiments conducted using NSL-KDD data set and explored the models based on performance and error metrics resulting in higher accuracy and decreased resource utilization.

We compared the performance of Random Tree model with leading NBTree and JRip approaches of intrusion detection system. The results of our experimental study shows that classification model integrated with discretization and feature selection method results in better accuracy, error rate and reduced false alarm rate. Considering the advantage and simplicity of RandomTree model, it can also be applied for dependent attributes such as NSL-KDD intrusion detection dataset.

REFERENCES

- [1]. Levent Koc, Thomas A. Mazzuchi, and Shahrar. m. Sarkani, "A network intrusion detection system based on a hidden naïve bayes multiclass classifier", Expert Systems with Applications, 39(18), pp.13492-13500, 2012.
- [2]. Siva S. Sivatha Sindhu, S. Geetha and A. Kannan, "Decision tree based light weight intrusion detection using a wrapper approach", Expert Systems with Applications, 39(1), pp.129-141, 2012.
- [3]. Chebroly S., Abraham A., Thomas J. P., "Feature deduction and ensemble", Design of Intrusion Detection System, Computer and Security, 24(4), pp.295-307, 2005
- [4]. Tsang, C. H., Kwong, S., & Wang, H., "Genetic-fuzzy rule reordering in mining approach and evaluation of feature selection techniques for anomaly intrusion detection", Pattern Recognition, 40(9), pp. 2373-2391, 2007.
- [5]. N.B. Annur, H. Sallehudin, A. Gani and O. Zakari, "Identifying false alarm for network intrusion detection system using hybrid data mining and decision tree", Malaysian Journal of Computer Science, 21(2), pp.101-115, 2008.
- [6]. John Mill and Atsushi Inoue, "Support vector classifiers and network intrusion detection", In Proceedings of the IEEE International Conference on fuzzy systems, vol.1, pp.407-410, 2004.
- [7]. S. Peddabachigari, A. Abraham, C. Grosan and J. Thomas, "Modeling IDS using hybrid intelligent systems", Journal of network and computer applications, 30(1), pp.114-132, 2007.
- [8]. S. Mukkamala, A. H. Sung and A. Abraham, "Intrusion detection using an ensemble of intelligent paradigms", Journal of network and computer applications, 28(2), pp.167-182, 2005.
- [9]. Kumar and Spafford, "A pattern matching model for misuse intrusion detection", In Proceedings of the 18th National Systems Security Conference, pp.150-158, 1995.
- [10]. Aleksandar Lazarevic, Levent Ertöz, Aysel Ozgur, Vipin Kumar, Jaideep Srivastava, "A comparative study of anomaly detection schemes in network intrusion detection", Proceedings of the third SIAM International Conferences on data mining, San Francisco, CA, 2003.
- [11]. Lunt T. F., "Real-time intrusion detection expert system", Final Technical Report, February 28, 1992.
- [12]. Xin Yu. (2004). "Improving TCP performance over mobile ad hoc networks by exploiting cross layer information awareness", ACM Proceedings of the 10th Annual International Conference on Mobile Computing and networking, pp:231-244.
- [13]. Zhang Z., HIDE: a hierarchical network intrusion detection system using statistical preprocessing and neural network classification", In

- Proceedings of the IEEE Workshop on information assurance and security, pp,85-90,2001.
- [14]. Bridges,S.M., & Vaughn R.B, "Fuzzy data mining and genetic algorithms applied to intrusion detection", In Proceedings of the 23rd National Information Systems Security Conference (NISCC),pp.13-31,2000.
- [15]. Bo, Hui-Ye &Yu-Hang, "HMMs (Hidden Markov Models) based on Anomaly intrusion detection method", In paper presented at the International conference on machine learning and cybernetics,2005.
- [16]. Zhang J., Zulkernine,M., & Haque,A."Random-forests based network Intrusion detection systems, IEEE Transactions on Systems, Man and cybernetics, Part C:Applications and Reviews, 38(5), 649- 659, 2008.
- [17]. Barbara D.,Wu,N., & Jajodia,S, "Detecting novel network intrusions using Bayes estimators", In Paper presented at the first SIAM conference on data mining,Chicago,2001.
- [18]. Lee, W.Stolfo,S.J., & Mok,K.W, "A data mining framework for building Intrusion detection models", In Proceedings of the IEEE symposium on security and privacy, pp.120-132,1999.
- [19]. Wu., S & Yen.E, "Data mining based intrusion detectors", Expert SystemswithApplications,3(1),pp.5605-5612,2009.
- [20]. Blum, A., & Langley, P, "Selection of relevant features and examples in machine learning", Artificial Intelligence,97(1),45-271,1997.
- [21]. M. Tavallae, E. Bagheri, W. Lu, and A. Ghorbani, "A Detailed Analysis of the KDD CUP 99 Data Set," Second IEEE Symposium on Computational Intelligence for Security and Defense Applications(CISDA) , 2009.
- [22]. Geoffrey Holmes, Bernhard Pfahringer, Richard Kirkby, Eibe Frank, Mark Hall, " Multiclass alternating decision trees," Proceedings of the 13th European Conference on Machine Learning", 161-172,2002.
- [23]. J.R.Quinlan, "Improved use of continuous attributes in C4.5",Journal of Artificial Intelligence Research,vol.4,pp.77-90,1996.
- [24]. D.E.Denning."An Intrusion Detection Model", IEEE Transactions on Software Engineering,vol.13,222-232, 1987.
- [25]. H.S. Javitz, and A. Valdes, The NIDES Statistical Component: Description and Justification, Technical Report, Computer Science Laboratory, SRI International, 1993.
- [26]. I. Guyon, A. Saffari, G. Dror, and G. Cawley, "Model selection: Description and Justification, Technical Report, Computer Science Learning, vol. 11, pp. 61-87, 2010.