

The Research on Intrusion Feature Selection Algorithm Based on Particle Swarm Optimization

¹Wang Yuanzhi and ²Ge Wengeng

¹School of Computer and Information, Anqing Normal College, Anqing, 246011, China

²School of Information Engineering, Huanghuai University, Zhumadian, 463000, China

Abstract: High-dimensional intrusion detection data concentration information redundancy results in low processing velocity of intrusion detection algorithm. Accordingly, the current study proposes an intrusion feature selection algorithm based on Particle Swarm Optimization (PSO). Analyzing the features of the relevance between network intrusion data allows the PSO algorithm to optimally search in a featured space and autonomously select effective feature subset to reduce data dimensionality. Experimental results reveal that algorithm can effectively eliminate redundancy and reduce intrusion feature selection time to effectively increase the detection velocity of the system while ensuring detection accuracy rate.

Keywords: Feature relevance, intrusion detection, intrusion feature selection, optimization searching, particle swarm optimization

INTRODUCTION

Intrusion detection technology is an important research direction for network security. The presence of any attack behavior against security policy in the network or system is detected by analyzing network flow or system audit record, which generates corresponding strategies, compensates system bugs and fills system functions. However, during intrusion detection, the data size collected by the detector is large and with different features. Some of these features are not related with detection features, which reduce classification or cluster accuracy and significantly increase complicity in study, training time and space, thereby affecting the efficiency of the algorithm operation. Therefore, methods on how to increase the detection velocity of the system without compromising the detection accuracy rate have received much attention. Previous studies reported that the intrusion feature selection could maintain the integrity of the original network data and eliminate redundancy features to increase system detection velocity. Based on current intrusion feature selection algorithms, heuristic search strategies exist, including forward search, reverse search and sequential search, among others. Several scholars employed a method that combines expansion matrix theory and genetic algorithm, among others (Kennedy and Eberhart, 2001; Prasad, 2000; Prasad *et al.*, 2003a, b).

For intrusion feature selections in a large-scale data set, these search strategies have large computing resource degradation, slow convergence rate and high time complexity. Although the feature subsets obtained from non-search strategies in Literature (Gopal, 2003) have relatively low time complexities, these subsets

have relatively more redundancy features, which affect classification accuracy. The current study aims to improve the defects in existing algorithms by proposing an Intrusion Feature Selection Algorithm (IFSA) based on Particle Swarm Optimization (PSO), which introduces feature relevance analysis to direct PSO with a faster convergence rate in feature space and to realize the self-adaption and self-optimization of the intrusion feature selection.

FEATURE RELEVANCE ANALYSIS

Relevance selection (David, 2000) is widely used in machine learning and statistics relevance analysis to evaluate relevance between features; in addition, relevance selection significantly affects the efficiency of intrusion feature selection. The relevance between two random variables is generally measured through entropy (Eberhart and Shi, 2001) defined in information theory.

Definition 1 (entropy): For a dispersed feature X with concentrated data, the possible dereference might be $\{x_1, x_2, \dots, x_n\}$, the corresponding probability distribution is $\{p(x_1), p(x_2), \dots, p(x_n)\}$ and the entropy definition of X is:

$$H(X) = -\sum_{i=1}^n p(x_i) \log p(x_i) \quad (1)$$

Entropy expresses information size of feature X . The smaller the entropy, the more asymmetrical the distribution of data dereferencing in X will be. The more data of some or several values X values for, the

less data of other values it values. If all data of X obtain the same value, then the entropy of X is 0 and the information contained in the feature is 0; that is, no available information of such feature is present for the data in the data set. By contrast, the bigger the entropy of X, the more symmetric distribution of data value and the more information it contains.

Definition 2 (joint entropy): For dispersed features X and Y in a data set, if the joint probability they value for x_i and y_i , respectively, is $p(x_i, y_i)$, then the Y joint entropy of X is defined as:

$$H(X, Y) = - \sum_{i=1}^n \sum_{j=1}^m p(x_i, y_j) \log p(x_i, y_j) \quad (2)$$

The concept of joint entropy is based from entropy, which describes the information size shared by the two random variables. The bigger the value, the greater the relevance between the two variables will be. If joint entropy between the two variables is 0, then they are independent.

IFSA: IFSA employs PSO with a faster convergence rate to search in feature space. IFSA introduces relevance analysis to guide algorithm search and realizes self-adaption and self-optimization of intrusion feature selection.

Intrusion feature selection definition and principle:

Definition 3 (feature subset): Feature subset is the new attribute set obtained after the irrelevant and redundant attributes entered into the attribute set are canceled.

Definition 4 (intrusion feature selection): Intrusion feature selection (also called attribute selection or feature extraction) recognizes and selects an effective attribute subset to describe an effective mode in a relatively large data set that generally contains redundant and irrelevant attributes. During feature subset selection, the algorithm generally selects the effective attribute set in the smallest scale using the following principle:

- Classification accuracy should not significantly decrease
- Classification distribution should maintain concord as much as possible before and after the intrusion feature selection

PSO algorithm: Based on the research result of bird flock foraging behavior (Shi and Eberhart, 1998), Kenedy and Eberhart (2001) proposed the PSO algorithm in 1995. The PSO algorithm, with a fast execution velocity and good resistance to dimension changes, immediately attracted much attention. The algorithm is described as follows:

In a D-dimensional target search space, assuming that m particles constitute a community, in which the i^{th} particle is expressed as a D-dimensional vector, $x_i = (x_{i1}, x_{i2}, \dots, x_{iD})^T$ where $i = 1, 2, \dots, m$, then the location of the i^{th} particle in a D-dimensional search space is x_i and the location of each particle is a potential solution. When x_i is substituted in a target function, then its adaptive value is calculated. Whether x_i is good or bad is identified and measured according to fitness. The “flight” velocity of the i^{th} particle is also a D-dimensional vector, given by $v_i = (v_{i1}, v_{i2}, \dots, v_{iD})^T$. To date, the optimal location of the i^{th} particle is written as $p_g = (p_{g1}, p_{g2}, \dots, p_{gD})^T$. PSO operates the particles demonstrated according to the following two formulas (Kennedy, 2000):

$$v_{id}^{n+1} = \omega v_{id}^n + c_1 r_1 (P_{id}^n - x_{id}^n) + c_2 r_2 (P_{gd}^n - x_{id}^n) \quad (3)$$

$$x_{id}^{n+1} = x_{id}^n + \eta v_{id}^{n+1} \quad (4)$$

where,

- d = 1, 2, ..., D
- i = 1, 2, ..., m
- m = The colony scale
- ω = The inertia weight, which is the coefficient of keeping the original velocity
- c_1 = The optimal value of weight coefficient in history that particles tracking themselves, which expresses the cognition of particles to themselves. Thus, it is named “cognition” and generally set up as 2
- c_2 = The weight coefficient of the optimal value that particles track the community, which expresses the particles cognition to the entire community and is thus named “social knowledge” or usually as “society” and is generally set up as 2
- r_1, r_2 = The random numbers between [0, 1]
- η = The coefficient added before velocity during location update, which is named constraint factor and is generally set up as 1
- n = 1, 2, ... are iterations

Binary Particle Swarm Optimization (BPSO)

algorithm: Kennedy and Eberhart (2001) proposed the BPSO version of the PSO algorithm (Angeline, 1998), which fueled such algorithm into a combinatorial optimization field. BPSO applies the binary coding form and restricts each dimension x_i and p_i as 1 or 0 in BPSO model and velocity v_i is not under such restriction. The Velocity Sigmoid function is used to express the possibility of location state change:

$$s(v) = 1/(1 + e^{-v}) \quad (5)$$

The velocity update in Eq. (3) of BPSO is not changed. By contrast, the location Eq. (4) is changed as follows:

If $(\text{rand} () < s (v_{id}^n))$ then $x_{id}^{n+1} =$; Els:

$$x_{id}^{n+1} = 0; \tag{6}$$

where, $\text{rand} ()$ is the random number between (0, 1) and the maximum velocity V_{\max} is used to restrain the possibility that x_i might be 0 or 1. The key in using PSO for optimization problems is the fitness function selection, which embodies the relationship between practical problems and optimal algorithms.

Coding mode: The essence of the intrusion feature selection is the selection of N features to form the subset from M features. Therefore, each feature can be defined as a one-dimensional dispersed binary variable of the particle and the M features constitute M-dimensional binary spaces of the particles. For each particle, if the i^{th} digit is 1, then the i^{th} feature is selected; if vice versa, then the feature is not selected. Therefore, each particle represents a different feature subset; that is, a candidate set. For example, if particle $i = 100110$, then features 1, 3 and 5 are selected and the feature subset is $\{1, 3, 5\}$:

Fitness function: During intrusion feature selection, the selection of fitness evaluation function is of prime importance. Although scholars have provided several different suggestions, such as distance evaluation and relevance evaluation, among others, a standard measurement that is generally accepted is still not available to date. The principal idea of the relevance evaluation method applied in the present study lies in selecting an attribute subset, wherein much relevance exists between each attribute and property, but no internal relevance to eliminate irrelevant and repeated attributes. The relationship between the two attributes A and B can be measured using the symmetric uncertainty, given by:

$$F(A, B) = 2 * \frac{H(A) + H(B) - H(A, B)}{H(A) + H(B)} \tag{7}$$

Relevance-based attribute selection decides goodness of an attribute set, which is measured as follows:

$$f(A, C) = \left(\sum_j F(A_j, C) \right) / \left(\sqrt{\sum_i \sum_j F(A_i, A_j)} \right) \tag{8}$$

where, C is the property and I and j that include all sets attributes comprise the attribute set. In the fitness function of particle swarm denoted in Eq. (8), bigger values produce higher particle fitness.

Algorithm description: According to the above analysis, the key steps of algorithm in the present study are as follows:

- Step 1:** Load training data set and set initialization parameter.
- Step 2:** Generate original colony at random to generate random initialized velocity for each particle and set the individual extremum pbest of the particle and global extremum gbest of the colony.
- Step 3:** Evaluate fitness value of each particle based on Eq. (8).
- Step 4:** Compare fitness value with the best location pbest of each particle. If it is more optimal than the pbest, then it is the best current location pbest.
- Step 5:** Compare the fitness value of each particle with the best location gbest that the colony has ever been. If it is more optimal than the gbest, then it is the optimal location of the colony and the reference number of the gbest is reset.
- Step 6:** Update particle velocity and location according to Eq. (3), (5) and (6).
- Step 7:** If iterations reach the maximum value, proceed to Step 8, or repeat Step 3.
- Step 8:** Convert the optimal location of the colony to the corresponding feature subset and returns.

EXPERIMENTAL PROCESS AND RESULT ANALYSIS

Experimental data set: To verify the effectiveness of the algorithm, the experiment data applies KDD99 data set. The test data set collects operational data of the analog network for five weeks. The data of the first and third weeks are training data with no attack cases to train the abnormal detection system. The data of the second week are test data, including 43 attack cases (64 attack methods in the entire test data) to train the misused detection system. The data of the fourth and fifth weeks are test data, including 5 million records in the data set with 41 features (attribute) in each record. The data volume in the KDD99 data set is very large; therefore, the data set needs to be sampled to decrease data volume for convenient algorithm testing. The data from the training and test data sets were randomly sampled. The randomly sampled data were combined, which formed training data subset (a total of 21836) and test data subset (a total of 36715) for experimental use. Algorithm parameters were set as follows: $D = 41$, $m = 30$, $\omega = 0.9$, $c_1 = c_2 = 2.0$, $V_{\max} = 4.0$ and 50 iterations.

Experimental program: The following three experiments were designed to verify the performance of IFSA more clearly and correctly.

Experiment 1: The performance of the intrusion detection models based on all 41 features was compared with that of the feature subset model based on intrusion feature selection in detection time and precision. The intrusion feature selection methods proposed in the current study were first applied in the randomly

sampled data set to obtain the corresponding feature subsets. Subsequently, intrusion detection models were constructed in the feature subsets based on all 41 features and after the intrusion feature selection in the training set.

Experiment 2: The performance of intrusion detection models using the proposed algorithm was compared with intrusion detection models using Genetic Algorithm (GA) and Relief algorithm in detection time and detection precision. The IFSA proposed in the current study was firstly applied in the data set to obtain corresponding feature subsets and construct intrusion detection models based on the feature subset. The GA and Relief algorithms were then used for the intrusion feature selection to the same training data set.

Experiment 3: Support Vector Mechanism (SVM) was used as the classifier. The classification error score generated by sample classification of the three feature subsets were compared after the intrusion feature selection in Program 2.

Experiment result analysis: The IFSA algorithm in the current study was applied in the experimental training data set after random sampling. Table 1 demonstrates the generated feature subset.

Table 1 shows that different feature subsets for different attacks were obtained after the intrusion feature selection. For each attack type, data sets for intrusion detection model were constructed. These data sets have not gone through intrusion feature selection and feature subsets after intrusion feature selection demonstrated in Table 1. The performance of each intrusion detection model in detection time and detection precision was compared. The results from the comparison are shown in Table 2.

Table 2 shows that the performance of intrusion detection models using the algorithm proposed in the current study is obviously better than the intrusion detection models that have not gone through intrusion feature selection in detection precision and detection time. IFSA, GA and Relief algorithms were used in the training data subset for experimental use after random sampling. The experimental results are presented in Table 3 and 4.

Finally, SVM was used as the classifier. Feature subsets obtained after the IFSA, GA and Relief algorithms performed intrusion feature selection to the experimental training data set as the classification sample. The experimental results are demonstrated in Table 5 and 6.

Table 4 shows that no significant difference exists among the feature subsets as classification sample generated through IFSA algorithm and the basic GA algorithm proposed in the current study in terms of classification precision, which is higher than the Relief algorithm in precision rate. In terms of intrusion feature

Table 1: Feature subset obtained via the IFSA algorithm

| Attack type | Feature subset |
|-------------|---|
| Dos | Dos, protocol_type, src_bytes, count, dst_host_same_src_rate |
| Probe | Probe, duration, service, src_bytes, dst_bytes, count, dst_host_diff_src_rate |
| R2L | R2L, duration, service, src_bytes |
| U2R | U2R, duration, service, src_bytes, root_shell, dst_host_count |
| Normal | Protocol_type, service, src_bytes, count, dst_host_count |

Table 2: Comparison in detection rate and detection time before and after the intrusion feature selection

| Attack type | Classification precision | | Detection time (sec) | |
|-------------|--------------------------|----------------|----------------------|----------------|
| | All features | Feature subset | All features | Feature subset |
| Dos | 83.5% | 98.4% | 1.09 | 0.25 |
| Probe | 87.4% | 97.5% | 1.23 | 0.36 |
| R2L | 85.7% | 98.7% | 1.16 | 0.28 |
| U2R | 84.2% | 97.1% | 0.94 | 0.17 |
| Normal | 86.3% | 98.8% | 1.13 | 0.24 |

Table 3: Comparison of classification precision among IFSA, GA and relief algorithms

| Attack type | Classification precision | | |
|-------------|--------------------------|-------|--------|
| | IFSA | GA | Relief |
| Dos | 98.4% | 97.3% | 95.6% |
| Probe | 97.5% | 96.7% | 96.4% |
| R2L | 98.7% | 98.9% | 94.7% |
| U2R | 97.1% | 95.3% | 95.3% |
| Normal | 98.8% | 97.2% | 95.8% |

Table 4: Comparison among IFSA, GA and relief algorithms in classification detection time

| Attack type | Detection (sec) | | |
|-------------|-----------------|------|--------|
| | IFSA | GA | Relief |
| Dos | 0.25 | 0.35 | 0.41 |
| Probe | 0.36 | 0.43 | 0.48 |
| R2L | 0.28 | 0.38 | 0.47 |
| U2R | 0.17 | 0.21 | 0.33 |
| Normal | 0.24 | 0.36 | 0.42 |

Table 5: IFSA, GA, and relief classification precision comparison using the SVM classifier

| Data set | Classification precision | | |
|---------------------|--------------------------|-------|--------|
| | IFSA | GA | Relief |
| Training subset | 98.4% | 97.3% | 95.6% |
| Experimental subset | 97.5% | 96.7% | 96.4% |

Table 6: IFSA, GA, relief detection time comparison using the SVM classifier

| Data set | Detection time (sec) | | |
|-----------------|----------------------|------|--------|
| | IFSA | GA | Relief |
| Training subset | 0.25 | 0.35 | 0.41 |
| Test subset | 0.46 | 0.59 | 0.53 |

selection time, the method proposed in the present study exhibited the best performance.

Combining the results of Experiments 1, 2 and 3 reveals that the IFSA algorithm effectively reduces feature dimension of data information. In addition, the performance of intrusion detection models based on IFSA algorithm is better than intrusion detection

models that have not gone through intrusion feature selection in detection time and detection precision. Compared with GA and Relief algorithms, IFSA algorithm can significantly reduce time complexity of intrusion feature selection under high classification precision and effectively shorten intrusion feature selection time.

CONCLUSION

The current study proposes the algorithm IFSA based on PSO. Based on the analysis of the relevance among all the features in network intrusion data, IFSA uses the optimization searching of PSO algorithm in all feature spaces, conducts guidance search according to relevance and can select effective feature subsets through self-adaptation and self-optimization to reduce data dimensionality. The result of intrusion detection models for the experimental verification to KDD99 data set reveals that the intrusion feature selection method proposed in the current study can ensure effective increasing system detection performance while also ensuring detection after being used in the intrusion detection system. Therefore, IFSA is better than current intrusion feature selection methods in terms of detection time and classification precision.

ACKNOWLEDGMENT

The study is supported by the Natural Science Key project of Anhui province Education Department of China under Grant No. KJ2008A18ZC, KJ2010A232.

REFERENCES

- Angeline, P.J., 1998. Using selection to improve particle swarm optimization. *Proceeding of IEEE International Conference Evolutionary Computation*, pp: 84-89.
- David, E.G., 2000. *Genetic Algorithms in Search, Optimization and Machine Learning*. Pearson Education Asia Ltd., New Delhi.
- Eberhart, R.C. and Y. Shi, 2001. Particle swarm optimization: Developments, applications and resources. *Proceeding of the Congress on Evolutionary Computation*. IEEE Service Centre, Piscataway, NJ, Seoul, Korea.
- Gopal, M., 2003. *Control System Principles and Design*. Tata McGraw Hill Publishing Co. Ltd., New Delhi.
- Kennedy, J., 2000. Stereotyping: Improving particle swarm performance with cluster analysis. *Proceeding of International Conference on Evolutionary Computation*, pp: 1507-1512.
- Kennedy, J. and R.C. Eberhart, 2001. *Swarm Intelligence*. Morgan Kaufman Publishers, California.
- Prasad, R., 2000. Pade type model order reduction for multivariable systems using routh approximation. *Comp. Elec. Eng.*, 26(6): 445-459.
- Prasad, R., S.P. Sharma and A.K. Mittal, 2003a. Improved pade approximants for multivariable systems using stability equation method. *J. Inst. Eng. (India)*, 84: 161-165.
- Prasad, R., S.P. Sharma and A.K. Mittal, 2003b. Linear model reduction using the advantages of mikhailov criterion and factor division. *J. Inst. Eng.*, 84: 7-10.
- Shi, Y. and R.C. Eberhart, 1998. A modified particle swarm optimizer. *Proceeding of IEEE International Conference on Evolutionary Computation*. IEEE Press, Piscataway, NJ, pp: 69-73.