

Intrusion Detection with Genetic Algorithms and Fuzzy Logic

Emma Ireland

Division of Science and Mathematics
University of Minnesota, Morris
Morris, Minnesota, USA

December 2013
UMM CSci Senior Seminar Conference

The Big Picture



Outline

- 1 Background
- 2 Using Fuzzy Genetic Algorithms
- 3 Using Genetic Algorithms
- 4 Conclusions

Outline

- 1 Background
 - Types of Networking Attacks
 - Detection Methodologies
 - Data Sets - KDD99 and RLD09
 - Rules
 - Genetic Algorithms
 - Determining the Accuracy of an Algorithm

2 Using Fuzzy Genetic Algorithms

3 Using Genetic Algorithms

4 Conclusions

Types of Networking Attacks

- Denial of Service (DoS): attacker makes a machine inaccessible to a user by making it too busy to serve legitimate requests.
- Remote to User (R2L): attacker tries to gain access to things a local user would have on the machine.
- User to Root (U2R): attacker starts out with access on the machine and then tries to gain root access to the system.
- Probe: attacker examines a machine in order to collect information about weaknesses or vulnerabilities that in the future could be used to compromise the system.

Detection Methodologies

- Signature-based detection: compares well-known patterns of attacks that are already in the intrusion detection system against captured events in order to identify a possible attack.
- Anomaly-based detection: looks for patterns of activity that are rare and uncommon.

KDD99

- Generated by simulating a military network environment in 1999.
- Has long been a standard data set for intrusion detection.
- Data was processed into 5 million records.
 - A record is a sequence of TCP packets, between which data flows to and from a source IP address to a target IP address.
- Data in the set is classified as normal or attack activity.
- Uses 41 features, which are properties of a record that are used to describe the activity and help to distinguish normal connections from attacks.

Some Features of KDD99

- duration: length of the normal or attack activity in seconds.
- num_failed_logins: number of failed login attempts.
- root_shell: returns 1 if root shell is obtained, else returns 0.

RLD09

- RLD09 was created because KDD99 is 14 years old.
- Data was captured from a university in Bangkok, Thailand.
- 17 different types of attacks (divided into denial of service and probe attacks), and 12 features.

Rules

- A commonly used approach for detecting intrusions and to differentiate between normal connections and attacks is to use rules.
- If-Then format: If (*condition*) then (*consequence*).
 - The condition is composed of one or more features, and the consequence says if it is an intrusion or not.
 - If *duration* = 4 then *intrusion*.

Genetic Algorithms

- Search technique used to find solutions to problems.
- Mutation: random bits in an individual, or possible solution, are changed.
- Selection: individuals that have a better fitness are chosen over the other individuals.
- Fitness function: determines the quality of a particular individual.
- Crossover: Two individuals swap one of their characteristics with the other to form two new individuals.

Determining the Accuracy of an Algorithm

		Predicted	
		Not Attack	Attack
Actual	Not Attack	True Negative (TN)	False Positive (FP)
	Attack	False Negative (FN)	True Positive (TP)

- Detection rate (DR): number of intrusions detected by the system divided by the total number of intrusions that happen.

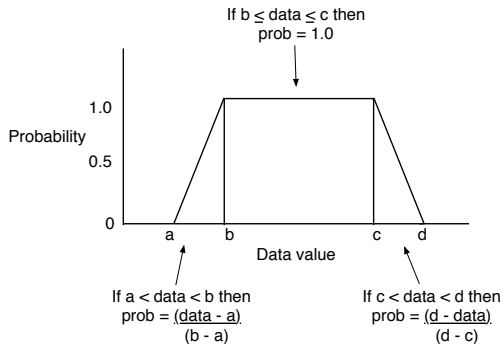
Outline

- 1 Background
- 2 Using Fuzzy Genetic Algorithms
 - Fuzzy Algorithm
 - Algorithm Overview
 - Experimental Design and Results
- 3 Using Genetic Algorithms
- 4 Conclusions

Fuzzy Logic

- Attacks on systems do not always have a fixed pattern, so fuzzy logic is used to detect patterns that have a behavior that is between normal and unusual.
- Fuzzy logic rules are similar to the rules described before, except that consequence is a certainty factor.
 - If (*duration* = 6) then (*probability of it being an attack is 50%*).

Measuring the Probability of a Record Being an Attack



Example:

- Feature: duration
(length of the activity in seconds).
- $a=1$, $b=3$, $c=5$, $d=7$
- The length of the activity is 6 seconds
(between c and d).
- $\text{prob} = \frac{d - \text{data}}{d - c} = \frac{7 - 6}{7 - 5} = 0.5$

Encoding of Features and Rules

- The four parameters are encoded into blocks.
- Each block is a feature with values between 0.0 and 7.0.

010	011	100	101
a=2	b=3	c=4	d=5

- A rule has 12 blocks of features, at the end is the type of attack.

010	011	100	101	010	011	101	111	DoS
a=2	b=3	c=4	d=5	a=2	b=3	c=5	d=7	
Block 1					Block 12				Type

Algorithm Overview

- The algorithm generates rules, improves rules, then rules are used to classify the data.
- One record is passed into a rule.
- Each feature in a record is matched to one block of the rule.
- The parameters of each block measure the probability of an attack using the trapezoidal fuzzy rule shape.
- The probabilities of each block are then compared with a threshold to determine if the record represents an attack or normal behavior.

Algorithm

```

for each record do
  for each rule do
    for each feature do
      prob = fuzzy();
      totalprob = totalprob + prob;
    end for
    if totalprob > threshold then
      class is attack;
    else
      class is normal;
    end if
  end for
  find  $A$ ,  $B$ ,  $\alpha$ , and  $\beta$ 
end for
calculate fitness
crossover(), mutation()

```

Fitness function:

$$\frac{\alpha}{A} - \frac{\beta}{B}$$

A : # of attack records.

B : # of normal records.

α : # of attack records
correctly identified as
attack.

β : # of normal records
incorrectly classified as
attack.

Experiments

- A variety of experiments were run. Two experiments used just RLD09. Three experiments used both the RLD09 and KDD99 in order to compare how the fuzzy GA would perform on both.

Experiments Using Only RLD09

Experiment 1

- Fuzzy GA was used to create DoS and probe detection rules, then the rules were verified with known attack types.
- Two steps in the training process: find a DoS rule, find a probe rule. Both of these rules were then used together in the testing process to identify attacks from the testing data set.
- 10,000 records were used for the training set, 26,500 records were used for the test set.

Experiments Using Only RLD09

Experiment 1 Results

	Attack	Normal	Total	FP(%)	FN(%)	DR(%)
DoS Training	1499	8501	10000	1.46	47.50	91.64
Probe Training	2496	7504	10000	1.83	15.38	94.79
Testing	10500	16000	26500	1.13	4.10	97.92

Experiments Using Only RLD09

Experiment 2

- Attacks were pulled out of the training set and kept for unknown data testing. This was to test that the fuzzy GA could detect unknown attacks.
- Used fuzzy GA and a decision tree algorithm.
- For each test case there were 13 attack types plus normal activity that were in the training data set. Three attack types were used for the unknown testing data set.

Experiments Using Only RLD09

Experiment 2 Results

Test Case	Unknown Attacks	Decision Tree DR (%)	Fuzzy Genetic DR (%)
1	Adv Port Scan (Probe)	Avg =	Avg =
	Ack Scan (Probe)	98.33	100
	Xmas Tree (Probe)		
2	UDP Flood (DoS)	Avg =	Avg =
	Host Scan (Probe)	46.65	99.80
	UDP Scan (Probe)		
3	Jping (DoS)	Avg =	Avg =
	Syn Scan (Probe)	99.70	98.75
	Fin Scan (Probe)		
4	UDP Flood (DoS)	Avg =	Avg =
	RCP Scan (Probe)	70.35	98.15
	Fin Scan (Probe)		
5	Http Flood (DoS)	Avg =	Avg =
	RCP Scan (Probe)	99.94	97.50
	Fin Scan (Probe)		

Experiments Using Both RLD09 and KDD99

Experiment 1

- Used fuzzy GA to classify normal activity and attacks from KDD99 and RLD09.

Data set	Attack	Normal	FP (%)	FN (%)	DR (%)
KDD99	160,117	39,337	0.13	1.55	98.72
RLD09	10,500	16,000	1.14	3.39	97.97

Experiments Using Both RLD09 and KDD99

Experiment 2

- Used the fuzzy GA to classify types of attacks in KDD99.
- 158,597 records of DoS attacks. 1,500 records of probe attacks.

Test	Attack	Type	FP (%)	FN (%)	DR (%)
1	Back	DoS	85.33	0.00	16.56
2	Smurf	DoS	0.76	0.10	99.73
3	Portssweep	Probe	6.40	0.00	93.66
4	Satan	Probe	0.74	3.75	99.22

Experiments Using Both RLD09 and KDD99

Experiment 3

- Used the fuzzy GA to classify types of attacks in RLD09.
- 6,400 records of DoS attacks. 10,400 records of probe attacks.

Test	Attack	Type	FP (%)	FN (%)	DR (%)
1	Smurf	DoS	0.02	0	99.98
2	UDP Flood	DoS	11.06	0	89.59
3	Ackscan	Probe	0.03	0	99.97
4	Synscan	Probe	0.65	4.2	99.24

Outline

- 1 Background
- 2 Using Fuzzy Genetic Algorithms
- 3 Using Genetic Algorithms
 - Algorithm Overview
 - Experimental Design and Results
- 4 Conclusions

Algorithm Overview

- The system is divided into 2 phases: a precalculation phase and a detection phase.
- In the precalculation phase, a set of chromosomes are created using training data. It takes network data as an input and outputs a set of chromosomes. Then this set of chromosomes is used in the detection phase.
- In the detection phase, selection, crossover, and mutation occur, and then the type of the data (whether it is an attack or normal behavior) is predicted.

Experimental Design

- KDD99 data set.
- Used only the numerical features of KDD99 (34 out of 41 total features).
- Training set: 494,021 records (396,741 of them are attacks).
Test set: 311,029 records (250,436 of them are attacks).

Table: Number of records

	Training	Testing
Normal	97,280	60,593
Probe	4,107	4,166
DoS	391,458	229,853
U2R	52	228
R2L	1,124	16,189
Total	494,021	311,029

Results

		Predicted					% Correct
		Normal	Probe	DoS	U2R	R2L	
Actual	Normal	42138	1421	15835	486	713	69.5
	Probe	398	2963	654	2	149	71.1
	Dos	921	432	228489	1	10	99.4
	U2R	146	21	8	43	10	18.9
	R2L	11191	578	3398	141	881	5.4
% Correct		76.9	54.7	92.0	6.4	50.0	

Outline

- 1 Background
- 2 Using Fuzzy Genetic Algorithms
- 3 Using Genetic Algorithms
- 4 Conclusions**

Conclusions

- The use of genetic algorithms and fuzzy logic in intrusion detection are effective ways of detecting attacks.
- The fuzzy genetic algorithm had a higher detection rate than a decision tree algorithm in most cases.
- Fuzzy genetic algorithms are good at detecting unknown attacks.

Thanks!

Thank you for your time and attention!

Questions?

References