

# 네트워크 침입 탐지를 위한 사례 기반 학습 기법\*

박 미 영\* 이 도 현\* 원 용 관\*\*

\*전남대학교 컴퓨터정보학부

\*\*전남대학교 정보통신공학부

{mypark, dhlee}@dbcore.chonnam.ac.kr

\*ykwon@chonnam.chonnam.ac.kr

## Instance-Based Learning for Intrusion Detection

Meeyoung Park\* Doheon Lee\* Yong Kwan Won\*\*

\*School of Computer & Information, Chonnam National University

\*\*School of Electronics & Computer Engineering, Chonnam National University

### 요 약

침입 탐지란 컴퓨터와 네트워크 자원에 대한 유해한 침입 행동을 식별하고 대응하는 과정이다. 점차적으로 시스템에 대한 침입 유형들이 복잡해지고 전문적으로 이루어지면서 빠르고 정확한 대응을 할 수 있는 시스템이 요구되고 있다. 이에 따라, 대용량의 데이터를 지능적으로 분석하여 의미있는 정보를 추출하는 데이터 마이닝 기법을 적용함으로써 지능적이고 자동화된 탐지를 수행할 수 있도록 한다. 본 논문에서는 학습 데이터를 각각 사례로 데이터베이스에 저장한 후, 실험 데이터가 입력되면 가장 가까운 거리에 있는 학습 데이터의 클래스로 분류하는 사례 기반 학습을 이용하여 빠르게 사용자의 이상 행위에 대해 판정한다. 그러나 많은 사례로 인해 기억 공간이 늘어날 경우 시스템의 성능이 저하되는 문제점을 고려하여, 빈발 에피소드 알고리즘을 수행하여 발견한 순차 패턴을 사례화하여 정상 행위 프로파일로 사용하는 순차패턴에 대한 사례 기반 학습을 제안한다. 이로써, 시스템 성능의 저하율을 낮추고 빠르게 정확하게 지능적인 침입 탐지를 수행할 수 있다.

### 1. 서론

침입 탐지란 컴퓨터와 네트워크 자원에 대한 유해한 침입 행동을 식별하고 대응하는 과정이다[1]. 인터넷을 통한 전자 상거래가 주요 서비스로 등장하면서 정보 유출, 전산망 침해, 정보 위변조 행위 등과 같은 침입 행위가 급격히 늘고 있어, 그에 따른 대응책들이 절실해지면서 한가지 해결책으로 침입 탐지 시스템의 필요가 증대되고 있다.

일반적으로 침입 탐지 기법은 오용 탐지(misuse detection) 기법과 이상 탐지(anomaly detection) 기법으로 분류된다. 오용 탐지는 이미 알려진 침입 행위를 이용하여 규칙을 생성한 후, 입력받은 데이터와 규칙이 일치하는 여부에 따라 침입으로 판정한다. 오용 탐지 기법은 전문가 시스템[2], 모델 기반 기법[3], 상태 전이 분석[4] 등이 있는데, 이 방법은 알려진 침입 유형에 대해서만 판정이 가능하다는 단점을 지니고 있다. 이상 탐지는 사용자의 정상 행위를 분석하여 정상 행위 패턴을 생성한 후, 새로운 입력과 사용자의 정상 행위 패턴을 비교하여 정상 행위에 어긋나는 경우 침입으로 판정한다. 이상 탐지 기법은 통계적인 방법, 인공지능, 데이터 마이닝 등을 이용한다.

통계적 분석의 기본적인 분석 방법은 파일 사용이나 CPU 사용량 등을 판정 요소로 하고 단기간의 사용자 행동을 장기간의 행동 패턴에 비교하는 것으로, 비교 값의 차가 크면 시스템 관리자에게 경보를 알려준다. 이 방법은 실시간 판정에서 사용하게 되는 프로파일 데이터를 최소화하고 주기적인 변경이나 유지보수가 불필요하다는 장점이 있다. 반면에 감사 데이터를

통계적인 수치 값으로 표현함으로써 데이터의 손실을 발생하거나 침입자가 침입 탐지 시스템을 학습시키는 침입 유형에는 대응하기가 어려운 단점이 있다[5][6][7].

데이터 마이닝이란 대용량의 실제 데이터로부터 잠재적으로 유용하지만 미리 알려지지 않은 목시적인 지식을 얻어내는 기법이다. 이 기술은 기존의 통계학적인 방법과 인공지능 분야를 기반으로 발전하였으나, 발생 가능한 여러 가지 패턴을 분석하고 예측할 수 있는 장점으로 다양한 분야에서 유용하게 이용되고 있다. 이 기법을 적용한 침입 탐지 시스템은 *syslog*와 같은 사용자 로그 정보, *tcpdump*, *BSM*(Basic Security Module) 감사 데이터 등의 대용량의 실제 데이터에 연관 규칙 탐사, 순차패턴, 분류, 군집화 등을 이용한다. 이것은 대용량의 데이터를 관리자가 직접 분석하고 탐지 규칙을 만들어 시스템을 보호하는 수작업관리 방식을 탈피하여 자동적이고 지능적이며, 보다 정확하게 탐지할 수 있도록 한다.

현재 데이터 마이닝 기법을 이용한 대표적인 방법에는 빈발 에피소드[8]와 사례 기반 학습을 이용한 방법[9] 등이 있다. 빈발 에피소드는 일정 기간 동안 발생한 사용자 명령 데이터에 빈발 에피소드 알고리즘을 적용하여 순차패턴 집합을 찾는다. 생성된 패턴 집합을 사용자의 정상 행위 프로파일로 정의한 후, 새로운 명령 시퀀스가 입력되면 사용자의 정상 행위 프로파일과 비교하여 유사도를 구한다. 유사도가 정상 행위 프로파일의 유사도 범위 내에 포함되는 경우 정상으로 판정하고, 그렇지 않으면 비정상으로 판정한다. 이 경우, 정상 행위임에도 불구하고 규칙과 정확히

\* 본 연구는 전남대학교 리눅스 시스템 보안 연구 센터를 통하여 정보통신부로부터 지원을 받았다.

일치하지 않으면 비정상으로 판정하는 거짓 탐지율(false positive)이 높아진다.

사례 기반 학습을 이용한 방법은 각 사용자별 히스토리 파일을 분석하여 명령어의 시퀀스를 수집하고, 각각을 사례(instance)로 만들어 데이터베이스에 저장한다. 관찰 대상이 되는 명령 시퀀스와 사용자의 정상 행위 프로파일 사례의 유사도를 적당한 거리 함수를 이용하여 측정한다. 유사도 값이 임계값 이상이면 정상으로 판정하고, 그렇지 않을 경우 비정상으로 판정한다. 그러나 사례 기반 학습의 기본적인 접근법은 몇 가지 단점을 가지고 있다. 첫째, 사례를 만들어가는 과정에서 데이터베이스가 커지게 되면 기억 공간과 유사도를 측정하는 과정의 계산 비용이 많이 요구된다. 둘째, 쓸모없는 데이터에 의해 분류 기준이 많이 발생하고 셋째, 유사도를 측정하기 위한 적절한 거리 함수를 선택하기가 어렵다.

본 논문에서는 위에 언급한 빈발 에피소드 방법과 사례 기반 학습을 이용한 방법이 갖는 장점을 취한 순차 패턴에 대한 사례 기반 학습을 제안한다. *icpdump* 데이터를 입력 값으로 사용하여, 빈발 에피소드 알고리즘을 통하여 찾아낸 에피소드를 사례화하고, 이를 각 사용자의 정상 행위 프로파일로 만든다. 새로 입력된 패턴과 정상 행위 프로파일의 각 사례 간의 유사도를 측정하여 임계값을 기준으로 정상과 비정상을 분류하여 판정한다. 이 방법을 적용함으로써, 사용자의 정상 행위에서 벗어난 행위를 빠르고 정확하게 탐지하고 기억 장소를 차지하는 양을 크게 줄여 침입 탐지 시스템의 성능을 높일 수 있다.

2. 순차 패턴에 대한 사례 기반 학습

2.1 빈발 에피소드

빈발 에피소드 알고리즘은 이벤트 시퀀스에서 서로 밀접하게 관련되어 빈번하게 발생하는 이벤트들의 순차적 패턴을 찾을 경우 사용한다.

에피소드는 일정 시간 동안 이벤트 시퀀스에서 서로 밀접하게 관련된 이벤트들의 집합이다[10].

[정의 1] 에피소드(episode)

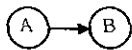
$e_i$ 는 이벤트일 때, 에피소드  $\varphi$ 는 이벤트들의 집합이며 다음과 같이 정의한다.

$$\varphi = (e_1, \dots, e_i),$$

예를 들어 다음과 같은 이벤트 시퀀스가 있다고 하자.

- (A, 123), (B, 125), (D, 140), (A, 150), (C, 151),
- (B, 155), (D, 210), (A, 220), (D, 222), (B, 225).

이 때, A, B, C, D는 다른 형태의 이벤트라고 하고, 숫자는 이벤트가 발생한 시간을 나타낼 때, (그림 1)과 같은 에피소드를 찾을 수 있다.



(그림 1) 에피소드 (a)

[정의 2] 최소 발생(minimal occurrence)

시간 간격  $[t_1, t_2]$  동안의 에피소드  $\varphi$ 의 발생은 다음 조건을 만족하면 최소 발생이다.

에피소드  $\varphi$ 가  $[t_1, t_2]$ 에 발생하고, 임의의 시간 간격  $[u, v]$   $C [t_1, t_2]$ 에서 발생하지 않는다.

위의 에피소드 (a)에 대해 가능한 최소 발생은  $[123, 125]$ ,  $[150, 155]$ ,  $[220, 225]$  등이다. 에피소드  $\varphi$ 의 최소 발생들의 집합은  $mo(\varphi)$ 로 표기한다.

[정의 3] 빈번도(frequency)

에피소드  $\varphi$ 의 빈번도,  $freq(\varphi) = |mo(\varphi)|$ 이다.

최소 빈번도,  $min\_fr$ 이 주어졌을 때 에피소드  $\varphi$ 의 빈번도가  $min\_fr$ 보다 크다면  $\varphi$ 는 빈번한 에피소드이다.

예를 들어,  $min\_fr$ 이 2인 경우, 에피소드 (a)는 빈번도가 3이고, 따라서 빈발 에피소드이다.

[정의 4] 빈발 에피소드 규칙(Frequent Episode Rule)

전체 집합 D의 이벤트들의 수에 대한  $mo(X)$ 의 비율을 지지도(support(X))라고 한다면, 빈발 에피소드 규칙은

$$X, Y \rightarrow Z, [c, s, w]$$

으로 표현된다. 이 때, 지지도  $s = support(XUYUZ)$  이고, 신뢰도  $c = support(XUYUZ) / support(XUY)$  이다. 단, 발생 간격의 너비는  $w$ 보다 작아야 한다.

예를 들어, <표 1>의 각 사용자의 일정 기간동안의 네트워크 패킷 정보를 수집하여 빈번하게 발생하는 패턴을 찾아 빈발 에피소드 규칙으로 나타내면 <표 2>와 같다.

<표 1> 네트워크 연결 레코드

시간	간격	서비스	출발지 호스트	목적지 호스트	Flag
1.1	0	http	Spoofed1	Victim	S0
1.1	0	http	Spoofed2	Victim	S0
1.1	0	http	Spoofed3	Victim	S0
1.1	0	http	Spoofed3	Victim	S0
...	...	...	...	...	...
10.1	2	ftp	A	B	SF
12.3	1	Sntp	B	D	SF

<표 2> 빈발 에피소드 규칙

빈발 에피소드 규칙	의미
(service = http, flag=S0, dst_host = victim), (service = http, flag = S0, dst_host = victim) → (service = http, flag = S0, dst_host = victim)	flag = S0 인 victim 호스트로의 http연결이 3%의 지지도를 가지고 발생한 후, 2초 내에 3번째 http연결이 93%의 신뢰도를 가지고 발생함 [0.93, 0.03, 2]

2.2 빈발 에피소드를 이용한 사례 기반 학습

사례 기반 학습(Instance-Based Learning)은 학습 데이터를 각각 사례로 데이터베이스에 저장한 후, 실험 데이터가 입력 되면 사례와 실험 데이터 간의 거리를 측정하여 유사도를 구하고, 가장 가까운 거리에 있는 학습 데이터의 클래스로 분류하는 방법이다.

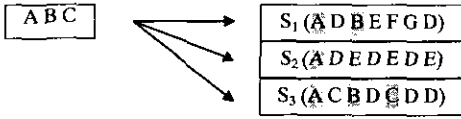
본 논문에서는 2.1에서 구한 빈발 에피소드를 사례화한 후, 새로운 이벤트 시퀀스와 각 사례간의 유사도를 측정한다.

[정의 5] 빈발 에피소드와 시퀀스간의 유사도

$E = (e_1, e_2, \dots, e_k)$ 가 빈발 에피소드이고, 시퀀스  $S = \{s_1, s_2, \dots, s_j\}$ 는 새로운 이벤트 시퀀스일 때, 유사도는 다음과 같이 정의된다.

$$Sim(E, S) = \frac{|E \cap S|}{|E|}$$

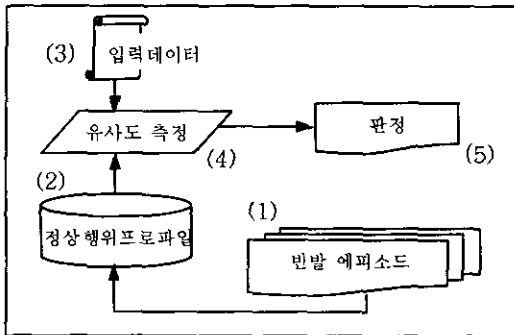
예를 들어, 빈발 에피소드 규칙 R이  $A, B \rightarrow C, [0.93, 0.25, 30]$  일 경우, 입력 시퀀스가 다음과 같다면, 유사도는 각각  $Sim(S_1) = 2/3$ ,  $Sim(S_2) = 1/3$ ,  $Sim(S_3) = 3/3$  이다.



(그림 2) 유사도 측정

2.3 수행 단계

순차 패턴에 대한 사례 기반 학습의 전체 수행 단계는 크게 두 단계로 구성된다. 먼저, [정의 4]를 이용하여 도출한 빈발 에피소드를 각각 사례화하여 정상 행위 프로파일로 저장한다. 다음 단계는 새로운 이벤트 시퀀스가 입력되면, [정의 5]를 이용하여, 정상 행위 프로파일의 각 사례와 시퀀스 간의 유사도를 측정한다. (그림 2)는 순차패턴에 대한 사례 기반 학습의 전체적인 흐름도이다.



(그림 2) 순차패턴에 대한 사례 기반 학습의 흐름도

[1단계] 프로파일 생성 과정

- (1) 일정 기간의 네트워크 패킷 데이터를 마이닝하기 위하여 tcpdump 파일에서 특성 항목을 선택한다.
- (2) (1) 에서 추출한 데이터를 입력 데이터로 사용하여 빈발 에피소드 알고리즘을 수행한다.
- (3) (2)에서 빈번하게 발생하는 순차 패턴을 각각 사례화하여 정상 프로파일로 저장한다.

[2단계] 이상 행위 판정 과정

- (1) 새로운 입력 시퀀스에서 특성항목을 추출한다.
- (2) 정상 프로파일의 패턴과 시퀀스를 차례대로 각각 비교한다.
- (3) 유사도가 임계값 이상인 사례가 나타날 때까지 (2)를 반복한다.
- (4) 비교하는 시퀀스에 임계값 이상인 사례가 존재하지 않을 경우, 해당 시퀀스를 비정상 행위로 판정한다.

3. 결론

점차적으로 시스템에 대한 침입 유형들이 복잡해지고 전문적으로 이루어지면서 빠르고 정확한 대응을 할 수 있는 시스템이 요구되고 있다. 기존에는 침입 탐지를 위해 통계학적인 기법이나 인공지능을 이용한 방법들이 다양하게 연구되었다. 그러나 이러한 연구들은 평균값 등을 사용하므로 데이터 손실을 유발하거나 부정확한 판정을 한다는 단점이 있다. 특히, 실시간으로 침입을 탐지할 경우 시스템의 성능이 저하되는 문제점 또한 간과할 수 없다. 이러한 문제점들을 해결하기 위해 본 논문에서는

빠르게 분류하기 위해 학습 데이터를 일반화하지 않고 각각 사례로 데이터베이스에 저장한 후, 실험 데이터가 입력되면 사례와 실험 데이터 간의 거리를 측정하여 유사도를 구하고, 가장 가까운 거리에 있는 학습 데이터의 클래스로 분류하는 사례 기반 학습을 이용하였다. 그러나 많은 사례로 인해 기억 공간이 늘어날 경우 시스템의 성능이 저하되는 문제점을 고려하여, 빈발 에피소드 알고리즘을 수행하여 발견한 순차 패턴을 사례화하여 정상 행위 프로파일로 사용하는 순차 패턴에 대한 사례 기반 학습을 제안하였다.

시스템 관리자가 상황에 따른 적절한 유사도의 임계값을 설정한다면 시스템 성능의 저하율을 낮추고, 빠르게 정확한 침입 탐지를 수행할 수 있다.

4. 참고 문헌

- [1] E. Amoroso, *Intrusion Detection*, Intrusion.Net, Sparta, New Jersey, 1999
- [2] Sandeep Kumar, *Classification and Detection of Computer Intrusions* PhD. Thesis, Purdue University, August 1995.
- [3] T. D. Garvey and T. F. Lunt. "Model based Intrusion Detection," *In Proceedings of the 14th National Computer Security Conference*, pages 372-385, October 1991.
- [4] K. Igun, R. Kemmerer, and P. Porras. "State Transition Analysis: A RuleBased Intrusion Detection System," *IEEE Transactions on Software Engineering*, 21(3), Mar. 1995.
- [5] H.S. Javitz and A. Valdes. "The NIDES statistical component description and justification," *Technical report, Computer Science Laboratory, SRI International, Menlo Park, California*, March 1994.
- [6] P.A. Porras and P.G. Neumann. "EMERALD: Event Monitoring Enabling Responses to Anomalous Live Disturbances," *In Proceedings of the Nineteenth National Computer Security Conference*, pages 353--365, Baltimore, Maryland, 22-25 October 1997. NIST/NCSC
- [7] T. Lunt, H. Javitz, A. Valdes et al. "A Real-time Intrusion-Detection Expert System (IDES)," *SRI International Technical Report, SRI Project 6784*, February 28, 1992.
- [8] W.Lee et al., "A Data Mining Framework for Building Intrusion Detection Models," *Proc. of IEEE Symposium on Security and Privacy*, pp.120-132, 1999.
- [9] T.Lane et al., "Temporal Sequence Learning and Data Reduction for Anomaly Detection," *ACM Trans. On Information and System Security*, Vol.2, No.3, pp.295-331, 1999.
- [10] Heikki Mannila, Hannu Toivonen, and A. Inkeri Verkamo. "Discovering frequent episodes in sequences," *In Proc. of the Int'l Conference on Knowledge Discovery in Databases and Data Mining*, Montreal, Canada, August 1995.
- [11] H. Mannila and H. Toivonen. "Discovering generalized episodes using minimal occurrences," *In Proceedings of the Second Int'l Conference on Knowledge Discovery and Data Mining*, pp. 146 - 151, Portland, Oregon, August 1996.