

# 스마트워크 환경에서 이상접속탐지를 위한 의사결정지원 시스템 연구\*

이 재 호,<sup>†</sup> 이 동 훈, 김 휘 강<sup>‡</sup>  
고려대학교 정보보호대학원

## Decision Support System to Detect Unauthorized Access in Smart Work Environment\*

Jae-Ho Lee,<sup>†</sup> Dong-Hoon Lee, Huy Kang Kim<sup>‡</sup>  
Graduate School of Information Security, Korea University

### 요 약

스마트워크 환경에서는 재택근무나 기업에서 구축한 스마트워크 센터, 휴대 가능한 모바일 단말기 등을 활용하여 원격 협업 환경을 구성하고 유연한 근무 환경을 조성하지만, 개인정보 및 업무상 중요 정보의 해킹, 노출 등의 위험성이 상존한다. 이러한 위험에 빠르게 대처하기 위해 기업 외부에서 일하는 직원이 내부망으로 접속할 때 사용하는 VPN(Virtual Private Network) 접속로그를 모니터링함으로써 직원들의 사용 패턴을 파악하고 비이상적인 행동을 탐지할 수 있다.

본 논문에서는 VPN 접속로그를 이용하여 기존의 로그 셋과 현재 접속의 유사도 측정 및 설문을 통한 적합한 시각화 방식을 제시하여 현재 접속의 정상 유무를 판단하는 시스템을 관리자에게 제공한다. 제안한 방법론을 통해 실제 기업환경에서 사용한 VPN 접속로그를 이용하여 실험을 한 결과 비정상 접속로그를 평균 88.7%로 추출할 수 있었으며, 관리자는 이 시스템을 이용하여 비정상적으로 접속하는 주체를 실시간으로 확인하여 대응할 수 있다.

### ABSTRACT

In smart work environment, a company provides employees a flexible work environment for tele-working using mobile phone or portable devices. On the other hand, such environment are exposed to the risks which the attacker can intrude into computer systems or leak personal information of smart-workers' and gain a company's sensitive information. To reduce these risks, the security administrator needs to analyze the usage patterns of employees and detect abnormal behaviors by monitoring VPN(Virtual Private Network) access log.

This paper proposes a decision support system that can notify the status by using visualization and similarity measure through clustering analysis. On average, 88.7% of abnormal event can be detected by this proposed method. With this proposed system, the security administrator can detect abnormal behaviors of the employees and prevent account theft.

**Keywords:** smart work, unauthorized access, decision support system, visualization, clustering.

접수일(2011년 12월 26일), 수정일(1차: 2012년 5월 18일, 2차: 2012년 6월 27일), 게재확정일(2012년 7월 30일)

\* 본 연구는 지식경제부 및 한국인터넷진흥원의 "고용계약형

지식정보보안 석사과정 지원사업"의 연구결과로 수행되었음

<sup>†</sup> 주저자, jh0814@korea.ac.kr

<sup>‡</sup> 교신저자, cenda@korea.ac.kr

## I. 서 론

최근 정보통신기술의 발달과 스마트폰, 태블릿 등 모바일 단말기의 확산으로 인해 언제 어디서나 일할 수 있는 새로운 근무 환경인 스마트워크가 도입되고 있다. 스마트워크란 정보통신기술을 이용하여 시간과 장소의 제약 없이 업무를 수행하는 유연한 근무 형태를 말한다[1]. 스마트워크 기술은 재택근무나 기업에서 구축한 스마트워크 센터, 휴대 가능한 모바일 단말기 등을 활용하여 원격 협업 환경을 구성하고 보다 유연한 근무 환경 기반을 만들 수 있다. 하지만, 개인정보 및 업무상 중요 정보의 해킹, 노출 등의 위험성이 상존하기 때문에 기업들의 스마트워크 도입 결정에 보안 문제가 걸림돌로 작용하고 있다. 주요 보안 위협으로는 단말기 분실 및 도난, 주요 정보 유출, 악성코드 감염, 네트워크 서버 해킹 등이 있다. 기업에서는 스마트워크를 도입할 때 발생할 수 있는 보안위협에 대한 하나의 해결방안으로 VPN을 도입하여 안전한 원격 네트워크 연결을 제공하고 있다. 그러나 직원의 계정을 탈취한 가장자(masquerader)의 공격이나 내부자가 의도적으로 정보를 노출하는 위협을 해결할 수는 없으며, 이를 해결하기 위한 방법으로 사용자의 행동을 분석하여 탐지하는 방법을 제시하고자 한다. 본 논문에서 제시하는 방법은 분실 모바일기기에 악의적으로 접속하는 경우 등 행동패턴이 이전과 상이할 경우 각 사용자는 고유한 행동 패턴을 가지고 있기 때문에 이러한 행동 패턴을 주로 접속하는 기기, 위치, 시간 등을 나타내는 접속 로그를 이용해서 발견 및 추출할 수 있다.

GS 칼텍스 개인정보 유출사고(2008), SK 하나카드 고객정보 유출 사고(2011) 등 최근 몇 년간 발생한 개인정보유출사고를 살펴보면 대부분 내부 직원에 의한 사고였으며, 2011년에 발생한 농협 전산망 마비 사태에서도 알 수 있듯이, 보안사고의 원인으로는 내부 관리 소홀이 가장 많은 비중을 차지하고 있다. 이러한 정보 유출 사고에 대응하기 위해 가장 중요한 것은 보안 관리자의 재빠른 탐지이다. 하지만 내부적인 요인에 의해 발생하는 경우, 탐지가 어려운 범주에 속한다. 또한, 현재 스마트워크가 증가함에 따라 모바일 기기를 이용하여 일을 하는 사용자는 증가하고 있지만 이에 대한 보안 및 관리에 대한 대책은 아직 미비하다. 본 논문에서는 스마트워크 환경에서 발생할 수 있는 위협에 빠르게 대처하기 위해 VPN 접속로그를 이용하여 이상 접속을 탐지하기 위한 방법론을 제시하고

자 한다. 이상 접속이란 각 사용자의 기존 접속 패턴을 벗어나는 접속을 의미하며, 분실 모바일기기에 악의적으로 접속하는 경우처럼 행동패턴이 이전과 상이할 경우를 예로 들 수 있다. 제안하는 시스템은 VPN 접속로그를 시각화하기 위한 설문을 수행하고 클러스터링 기법으로 유사도 측정을 하며, 이를 통해 관리자가 이상 접속을 탐지하여 접속하는 사용자에 대해 실시간으로 정상여부를 판단할 수 있다. 로그 시각화는 분석을 직관적으로 할 수 있도록 텍스트 형식의 데이터 셋을 그래프로 표현하는 것을 말하며, VPN 기존 로그 셋을 기기, 위치, 시간을 기준으로 표현하여 사용자의 행동 패턴을 직관적으로 표현한다. 클러스터링 기법을 이용한 유사도 측정 방법은 기존 로그를 정상으로 판단할 수 있을 만큼의 양을 가지고 있다는 전제하에 유사한 특징끼리 묶어 클러스터링한 후, 새로 입력된 로그와의 유사도를 숫자로 표현한다.

논문의 구성은 먼저, 2장에서는 스마트워크의 정의와 보안위협, 시각화 분석 기술, 클러스터링의 기법에 대해 분석한다. 3장에서는 이상접속탐지를 위한 의사결정 시스템을 제안하며, 4장에서는 실험 방법 및 결과 분석, 5장에서는 본 논문에 대한 결론을 설명한다.

## II. 관련 연구

### 2.1. 스마트워크의 정의 및 보안위협

스마트워크는 이동환경에서도 언제 어디서나 편리하게, 효율적으로 업무에 종사할 수 있도록 하는 업무 환경으로, 각 지역 주거지 인근에 구축된 전용 시설인 스마트워크센터에서 근무하는 형태, 정보통신기술을 활용하여 자택에 업무공간을 마련하고, 업무에 필요한 시설을 구축한 환경에서 근무하는 형태, 스마트폰, PDA, 노트북 등을 이용하여 공간적 제약 없이 업무를 수행하는 이동 근무 형태로 구분할 수 있다[1]. 스마트워크는 현장에서의 신속한 업무처리를 통해 업무속도와 생산성이 향상되며, 원격협업을 통해 실시간 협업이 가능해져 신속한 의사결정과 빠른 문제해결이 가능해진다. 또한 근무형태의 유연화로 여성, 장애인, 고령자 등 근로취약계층의 취업기회 확대 등의 긍정적인 효과를 기대할 수 있다. 반면에 스마트워크 환경에 적용되는 새로운 정보통신기술로 인한 보안 취약성이 발생할 수 있다. 서로 다른 소속의 직원들이 공동으로 사용하는 스마트워크센터의 경우 문서의 유출이 발생할 수 있고, 모바일 기기를 이용한 이동 근무의 경우

(표 1) 스마트워크 환경에서의 보안 위협과 대책

구분	보안 위협	대책
단말기 보안	- 개인정보 유출 : 단말기 도난 및 분실로 인한 개인·업무정보 유출, 스마트폰 소유자가 악의적으로 업무정보를 외부로 유출	- MDM(Mobile Device Management) : 단말기 위치 추적, 데이터 백업 등을 지원 - DRM(Digital Rights Management) : 화면 캡처, 프린트 금지 등을 통한 업무자료 보호
응용 프로그램 및 플랫폼 부문	- 개인정보 유출 : 스마트폰의 블루투스 및 애플리케이션을 통해 SMS, 통화기록, 위치정보 유출 - 장치이용 제한 : 콘텐츠 삭제, 아이콘 변경을 통해 기기 사용 및 일부 기능을 마비시키는 공격 수행 - 악성코드 감염 : 플랫폼 취약점을 이용한 바이러스 감염으로 인한 원격 제어 및 비정상인 요금 등을 유도 - 모바일 DDoS(Distributed Denial of Service) : 감염된 좀비 단말기는 특정 사이트에 트래픽을 유발	- 안티 바이러스 백신 : 웹·바이러스 탐지 및 실시간 검사 수행 - 단말 가상화 : 업무와 인터넷을 별도 플랫폼으로 분리하여 실행 - PMS(Patch Management System) : 설치된 앱의 버전 관리 - 코드 서명 : 신뢰된 인증기관의 인증서를 이용한 플랫폼·앱 무결성을 검증
네트워크 및 서버 보안	- 무선 구간 네트워크 해킹 : 네트워크구간에서 패킷 가로채기, 상용인터넷망을 통한 해킹 가능 - 주요 정보 노출 : 업무 내용에 대한 도·감청을 통해 기밀정보 수집	- VPN : 단말기에서 외부 서버까지 트래픽 암호화를 통한 안전한 통신채널 제공 - Wireless IDS(Intrusion Detection System) 및 IPS(Intrusion Prevention System) : 비정상 트래픽 모니터링 및 차단 - 릴레이 서버 구축 : 단말기에서 내부 업무서버로의 직접연결을 차단하여 외부망과 내부망을 분리

단말기 분실 및 도난, 악성코드 감염으로 인해 공격자에 의한 내부 사설망의 접근으로 개인정보나 기밀 정보가 유출될 수 있다. [표 1]은 스마트워크 환경에서의 보안 위협과 대책을 보여준다[1][2].

[표 1]을 보면, 스마트워크 환경에서는 기존의 보안 위협에 근무환경의 변화에 따른 새로운 보안 위협이 추가된 사실을 알 수 있으며, 아직 이에 대한 대응책은 미비한 실정이다. 본 논문에서는 이러한 보안 위협에 빠르게 대처할 수 있는 방법 중 하나로 VPN 접속로그를 모니터링하여 관리자가 의사를 결정하는데 도움을 주는 시스템을 제안하고자 한다. 이 의사결정 시스템은 시각화 분석 기술과 클러스터링 기법으로 구성되어 있다.

## 2.2. 시각화 분석 기술

시각화란 사용자에게 더 효율적으로 정보를 전달하기 위하여 그래픽 요소를 활용하여 데이터가 정보로서의 의미가 생성되도록 직관적으로 형상화하는 것을 말한다[3]. 로그 분석은 로그 내에 사용자 각각의 고유한 패턴이 존재한다는 가정 하에 이루어지며, 방대한 양

의 로그에서 패턴을 찾을 수 있는 정보만을 필터링하는 가공 과정을 거친다. 그러나 로그 필터링 후에도 텍스트 형식의 로그를 통해 직원의 사용 패턴을 확인하는 것은 관리자가 패턴을 확인하는 능력이 떨어뜨릴 뿐만 아니라 일일이 모니터링해야 하는 불편함을 초래하는 문제점을 가지고 있다. 따라서 대용량의 보안 이벤트의 경우 관리자가 이상을 감지하기 어려우며, 오탐을 유발할 확률도 높다. 시각화 기반 분석 기술은 실시간으로 발생하는 보안 이벤트를 즉각 처리하여 관리자가 방대한 양의 데이터를 빠르게 이해하는 것을 가능하게 해주며, 알려지지 않은 이상 패턴을 표현해 줌으로써 관리자가 신속하게 상황에 대처할 수 있도록 도와준다.

본 논문에서 제안하는 시스템은 관리자에게 기업 외부에서 접속하는 주체가 정상적인 시도인지 아닌지 판별하기 위한 의사결정을 지원하는 시스템이며, 그 판별 기준으로 주체의 기존 접속 기록을 사용한다. 기존 접속 기록은 현재 접속 시도의 이상 유무를 관리자가 실시간으로 판단할 수 있도록 주체의 접속성향을 직관적으로 표현되어야 하며, 이러한 요구를 충족시키기 위해 시각화 분석 기술을 적용한다.

보안 시각화 기술 분야는 네트워크 트래픽 정보 시각화 기술, 네트워크 공격 탐지 기술 등 다양한 연구가 진행되고 있다. Y.Livnat 등[4]은 다양한 alert의 what, when, where 속성을 연관 분석하여 사용자의 상황판단을 증진시키기 위한 시각화를 연구하였으며, 이 시스템은 사용자의 네트워크 이벤트의 상황판단력을 높이고, 네트워크 이상을 직접적으로 탐지할 수 있다. 최현상 등[5]은 웹, DDoS 공격 등의 패턴을 평행좌표계로 표현하여 공격을 탐지하는 PCAV (Parallel Coordinates Attack Visualizer)를 제안하였다.

하지만 보안 시각화 기술 분야는 시각화를 구현하는 환경이나 대상에 따라 구체적인 구현 및 디자인이 상이하다는 특성이 존재하기 때문에, 기 제안된 여러 시각화 분석 방법들은 본 논문에서 제안하는 환경에 직접적인 적용이 어렵다는 문제가 따른다. 따라서 본 논문에서는 스마트워크 환경에서의 VPN 접속로그 시각화에 특화된 시각화 방법 및 데이터 처리 알고리즘을 제시하고 실제 VPN 접속로그를 분석하여 최적화된 결과를 도출한다.

### 2.3. 클러스터링 기법

클러스터링이란 하나의 객체가 여러 속성을 갖고 있을 때, 유사한 속성들을 갖는 객체들을 묶어 전체의 객체들을 몇 개의 그룹으로 나누는 것을 말한다. 클러스터링 기법의 일반적인 절차는 데이터 샘플 중 특징(feature) 선택, 클러스터링 알고리즘 선택, 클러스터링 검증, 결과 도출 순으로 진행된다[6]. 특징을 선택할 경우에는 클러스터링을 수행한 결과, 서로 다른 클러스터에 속하는 값들이 구분 가능해야 하며, 알고리즘을 선택할 경우에는 클러스터링을 적용하는 분야와 데이터의 형식에 따라 적절하게 선택해야 한다. 따라서 각 클러스터의 밀집도가 높고, 클러스터 사이의 의존도가 낮은 결과를 도출하는 특징과 알고리즘을 선택해야 한다.

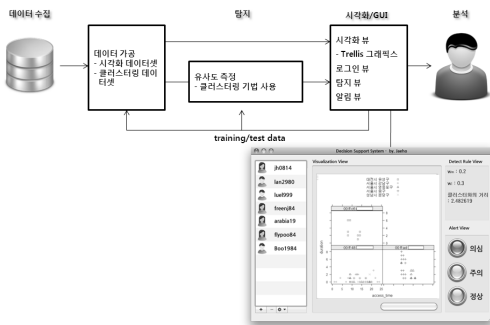
클러스터링 알고리즘은 크게 계층적 방식, 분할 방식, 밀도 기반 방식, 그리드 기반 방식, 모델 기반 방식으로 분류할 수 있다[7]. 분할 방식으로 가장 대표적인 알고리즘은 K-means로 군집의 개수 K를 결정하여 중심점과 객체와의 거리를 계산하여 클러스터링을 수행한다. PAM(Partitioning Around Medo-

ids)는 K-means와 유사한 방식이지만, 클러스터의 대표값을 평균으로 구하는 것이 아닌 대표 객체로 선택한다는 점에서 차이난다. 밀도 기반 방식은 거리에 기초하는 클러스터링이 아닌 밀도를 이용하여 임의의 형태가 있는 클러스터를 찾는 데 사용되며, 그리드 기반 방식은 객체공간을 격자구조로 이루어진 유한개의 공간으로 만들어 클러스터링을 하는 방식이다. 또한, 모델 기반 방식은 데이터가 특정한 확률밀도함수(Probability Density Function)들의 조합으로 생성된다고 가정하여, 주어진 객체로부터 특정한 확률 분포를 추정하고 그 결과에 따라 클러스터링을 수행하는 방법을 뜻한다. 마지막으로 신경망 접근 방식인 SOM(Self-Organizing Map)은 경쟁 학습을 통해 객체들을 상호 비교하며 스스로 클러스터를 조직하는 방법을 의미한다[8].

클러스터링 기법을 이용한 연구는 패턴인식, 웹 마이닝, 유전학, 지리학 등 다양한 분야에서 활발하게 진행되고 있다. 클러스터링 기법은 유사한 특징을 갖는 데이터가 하나의 클러스터를 이루며, 클러스터 간에는 뚜렷하게 특징이 다르다. 침입탐지분야에서는 침입탐지를 위한 새로운 클러스터링 기법이 제안되거나, 이상 행위 탐지 및 아웃라이어 탐지를 위해 클러스터링 기법이 사용되고 있다. 이상행위를 탐지하기 위해 W.Lee 등[9]는 데이터 마이닝 기술을 이용하여 연관 규칙과 빈도수 알고리즘으로 구성되어 있는 프레임워크를 제시하였으며, Y.Guan 등[10]은 K-means 알고리즘을 기반으로 하는 침입탐지를 위한 휴리스틱 클러스터링 기법을 제안하였다. 침입탐지 분야에서의 아웃라이어(outlier) 탐지는 네트워크의 비인가 접속을 탐지하는 것을 뜻하며, 이는 시스템 내에 악의적인 목적을 가지고 침입하는 공격일 수 있기 때문에 신속한 탐지가 필수적으로 요구된다[11].

본 논문에서는 사용자가 정상적으로 접속했다고 가정하는 데이터를 이용하기 때문에 이를 적절히 클러스터링하는 알고리즘을 선택하는 것이 중요하다. 또한, 각 클러스터의 개수 및 크기는 각 사용자의 행동 패턴을 나타내고 있기 때문에 클러스터의 크기가 작다고 해서 아웃라이어로 간주되는 것이 아니라, 이 또한 사용자의 고유한 특징이라고 인식한다. 따라서 제안하는 탐지 방법에서는 각 클러스터의 개수 및 크기를 고려하는 것이 아닌 사용자의 새로운 접속 기록과 인접한 클러스터와의 거리를 계산하여 이상 접속을 탐지한다.

III. 제안하는 의사결정시스템 방법론



(그림 1) 이상접속을 탐지하는 의사결정시스템의 순서도 및 구현된 응용프로그램 화면

제안하는 시스템은 다음과 같이 두 단계를 따라 효율적인 정보전달 및 의사결정시스템을 지원한다. 첫째, 선별된 정보를 사용자에게 가장 효과적으로 전달하는 방법에 따라 정보를 가공 및 처리하는 단계이다. 주어진 정보를 화면에 표현하는 방법은 다양하게 나타날 수 있지만, 본 논문에서는 시스템 관리자들에게 동일한 데이터를 여러 가지 방법으로 표현한 시각화를 전달하고, 이에 대한 선호도를 조사하여 최적의 표현 방식을 선정하였다. 둘째, VPN 접속로그로부터 유효한 정보를 추출하고, 이를 재가공한다. 이를 위해 클러스터링 알고리즘이 사용되며, 반복적인 실험을 통해 다수의 클러스터링 알고리즘 중에서 가장 효율적인 클러스터링 알고리즘을 선정하였다.

3.1. 시스템 구성

VPN 접속 로그 기록을 이용한 의사결정시스템은 [그림 1]에서 보듯이 탐지에 필요한 데이터 가공 처리, 관리자에게 시각화와 유사도 측정값을 제공하는 이상 접속 탐지, 그래픽 기반의 인터페이스로 구분할 수 있다. 이상 접속을 탐지하기 위한 데이터 처리 과정에서는 두 가지의 전제가 필요하다. 첫째, 기존 로그가 정상으로 판단할 수 있을 만큼의 양이 있어야 한다. 둘째, 측정된 로그는 시각화와 유사도 측정에 필요한 속성만을 선별하는 로그전처리(preprocessing) 작업을 거쳐야 한다. 따라서 일정기간동안 정상접속로그를 축적하여 가공하는 전체 과정이 끝나면 새로 들어오는 접속을 실시간으로 탐지할 수 있다. 탐지 방법으로는 시각화와 클러스터링 기법을 이용한 유사도 측정 방식을 사용한다. 시각화는 접속주체의 기

존 로그 기록을 통해 표현된 접속 성향과 새로 들어온 접속로그를 그래프로 표현한다. 유사도 측정은 기존의 로그들을 클러스터링한 결과의 대푯값과 새로 들어온 접속 로그간의 거리 계산을 통해 유사도를 측정하며, 이는 각 사용자별 패턴을 시각화를 통해 쉽게 유추할 수 없을 경우를 보완하기 위해 추가적인 정보를 제공하기 위함이다. 거리가 임계값보다 작을 경우 기존 로그와 비슷하다고 판단하여 정상으로 분류하며, 임계값보다 클 경우 주의 혹은 의심으로 분류하여 관리자에게 알린다. 그래픽 기반의 인터페이스는 시각화 뷰(visualization view), 탐지 뷰(detection view), 로그인 뷰(login view), 알림 뷰(alert view)로 구성되어 있다. 시각화 뷰는 선택한 접속주체의 접속 성향에 대한 그래프를 보여주는 뷰로 관리자는 기존 접속 패턴과 지금 접속한 로그와 얼마나 연관이 있는지 신속히 확인할 수 있다. 탐지 뷰는 클러스터링과 점의 사이의 거리를 이용해서 유사도를 측정한 값을 보여주는 뷰로써 로그인 접속정보의 위험도를 관리자에게 제공한다. 알림 뷰는 탐지 뷰의 위험도를 정상/주의/의심으로 분류하여 관리자가 접속로그의 위험도를 직관적으로 알 수 있도록 도움을 준다. 로그인 뷰는 최근에 로그인한 주체를 보여주는 뷰이며, ID를 선택하면 주체에 대한 시각화 뷰, 탐지 뷰, 알림 뷰가 갱신되어 관리자에게 이상접속정보를 제공한다.

3.2. 시각화 방법

이상접속을 탐지하려면 우선 직원 개개인의 고유한 행동 패턴을 파악하는 것이 선행되어야 한다. 예를 들어, 스마트워크센터나 재택에서 근무하는 직원의 접속 로그를 살펴보면 IP 주소는 주로 같은 대역이고, 접속 시각 및 체류시간 또한 거의 일정한 값을 갖는 특징이 보일 것이다. 반면에 휴대용 기기로 주로 접속하는 직원의 경우에는 IP 주소와 접속 시각은 다양한 값을 지니며, 체류시간은 대부분 작은 값을 갖는 특징을 볼 수 있을 것이다. 이외에도 특정 요일 및 시간대에만 접속을 하거나 출장을 정기적으로 가는 직원의 경우 특정 지역에서 주기적으로 접속하는 등 각 직원은 고유한 접속 패턴을 가진다. 이러한 행동 패턴은 여러



(그림 2) VPN 접속로그의 예시(예 : OpenVPN(12)의 로그)

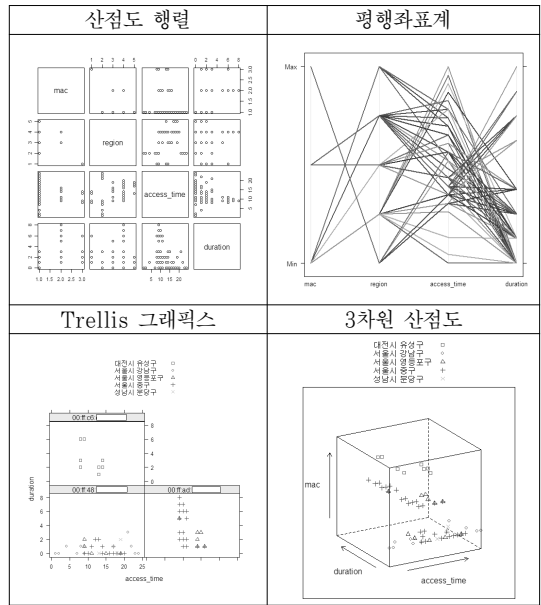
형태로 나타날 수 있으며, 관리자는 그래프를 통해서 텍스트 형식의 로그보다 직관적으로 패턴을 확인할 수 있고, 실시간으로 접속하는 직원의 접속기록이 기존 접속 패턴과 얼마나 비슷한지 시각적으로 확인할 수 있다.

VPN 접속로그는 [그림 2]과 같은 형태를 띠고 있으며, 이를 가공하여 시각화에 필요한 속성들을 추출해야 한다. 속성은 기본적으로 접속주체 ID, 접속기기 MAC Address값, 접속 IP 주소 및 포트번호, 접속 날짜 및 시각, 종료 날짜 및 시각으로 구성되어 있으며, 이 정보를 통해 관리자는 누가, 언제, 어디서, 어떠한 기기로 접속했는지 확인이 가능하다. 본 논문에서는 [표 2]와 같이 접속기기 MAC Address값, 접속 IP 주소, 접속 시각, 접속 체류시간을 이용하여 각 직원의 접속 패턴을 시각화하고자 한다.

시각화를 수행하기 위한 고려사항으로는 속성의 개수가 4개이므로 다차원을 표현할 수 있고, 각 사용자의 이용 성향이 직관적으로 드러나는 그래프이어야 한다. 따라서 본 논문에서 제안하는 시각화를 위한 그래프는 다차원 그래프이어야 하고, 관리자가 그래프를 통해 쉽게 접속주체의 이용성향을 파악할 수 있도록 가시성이 뛰어나야 한다. 다차원의 데이터를 시각화하는데 사용하는 그래프로는 대표적으로 산점도 행렬, 평행좌표계, Trellis 그래픽스, 3차원 산점도 등이 있다. [그림 3]은 openVPN 로그를 기반으로 수집한 동일한 데이터 셋을 산점도 행렬, 평행 좌표계, Trellis 그래픽스, 3차원 산점도 그래프로 각각 시각화하여 표현한 결과이다.

[표 2] 시각화 모델링 파라미터

파라미터	설명
접속기기 MAC Address 값	기기의 고유한 정보이다. ex) 00-24-1D-00-00-00
접속 IP	접속 IP를 통해 접속 위치를 확인할 수 있다. 사용자가 주로 같은 IP 대역에서 접속할 경우 고정 근무, 여러 IP 대역에서 접속할 경우 이동 근무를 하거나 휴대용 기기로 작업하는 직원이라는 것을 유추할 수 있다. ex) 192.68.3.2
접속 시각	사용자가 주로 어느 시간대에 접속하는지 확인할 수 있다. ex) 8:54:03
접속 체류시간	접속 체류시간이 짧은 경우 이동 근무형, 긴 경우에는 주로 재택근무나 스마트워크 센터 근무형이라는 것을 유추할 수 있다. ex) 2:03:46

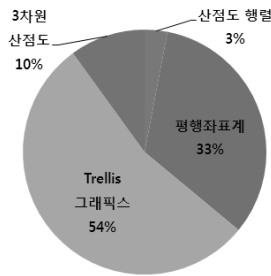


[그림 3] 다차원 그래프의 종류

산점도 행렬은 데이터가 여러 개의 변수를 갖고 있을 때, 이들 간에 어떤 상관성이 있는지, 가능한 모든 순서쌍에 대해 그림을 그려주는 그래프이다. 평행좌표계는 스크린 상에 세로축을 데이터필드의 크기만큼 배열하고 필드의 데이터 값을 선으로 연결하여 각 필드간의 관계를 육안으로 판별할 수 있다. Trellis 그래픽스는 여러 개의 패널이 행과 열의 구조로 정렬되어 있는 것이 일반적인 모습으로, 비교적 간단하게 여러 개의 그래프를 잘 정렬해서 배치할 수 있다[13]. 3차원 산점도는 산점도를 3차원으로 나타낸 것이다.

이 다차원 그래프들 중에서 사람들이 어떤 그래프를 가장 보기 쉬워하고 패턴이 잘 드러난다고 생각하는지 확인하기 위해 설문조사를 실시하였다. 일반적으로 통계분석에 많이 사용되는 다차원 그래프를 후보로 선정하여 정보보호 관련 연구원 50명을 대상으로 각 그래프를 읽는 방법을 알려준 후, “다음 중 어느 그래프가 사용자의 VPN 접속 패턴을 잘 나타내고 있는가? 그렇게 생각한 이유는 무엇인가?”, “당신이 관리자라면 어느 그래프를 선택하겠는가? 그렇게 생각한 이유는 무엇인가?”의 문항을 통해 어떤 그래프가 가장 보기 편하고 패턴을 파악하기 적합한지 참가자들의 의견을 조사하였으며, 조사 결과는 [그림 4]와 같다.

설문 조사 결과 실험 참가자의 54%가 Trellis 그래픽스를 선택하였으며, 평행좌표계가 그 뒤를 따랐다. Trellis 그래픽스는 다른 그래프에 비해 상당



(그림 4) 다차원 그래프 선호도

라벨 별로 보기 편했다는 의견이 대다수였으며, 평형좌표계는 양이 적을 때 이해가 가장 빠르겠지만, 고유한 패턴을 찾기 어려울 것이라고 답하였다. 산점도 행렬의 경우에는 처음 그래프를 숙지하는 데 시간이 오래 걸리지만, 그 후에 두 변수간의 관계를 정확하게 알아볼 수 있다는 의견이 있었으며, 3차원 산점도의 경우 입체적이라 가시성이 떨어지지만, 밀집도 확인은 쉬울 것이라는 의견이 있었다. 실험 결과를 통해 최종 선택된 Trellis 그래픽스는 다중 패널이며, 여러 변수 사이의 관계를 규명하는 데 매우 유용하다. Trellis 그래픽스는 관리자가 새로 접속하는 주체의 ID를 선택했을 경우 주체의 접속성향과 새로 입력된 로그와의 관계를 나타낸다. 상단 띠 라벨은 MAC Address를 나타내며, 이를 기준으로 로그들을 분류해서 시각화한다. 가로축은 접속 시각, 세로축은 접속 체류시간, 색은 접속 지역을 나타내어 주체가 자주 접속하는 시간대와 머물러 있는 정도 및 자주 사용하는 기기를 직관적으로 확인할 수 있다.

### 3.3. 클러스터링 기법을 이용한 유사도 측정

본 논문에서 제안하는 시스템은 VPN 접속로그를 시각화하여 이상접속을 판별할 뿐만 아니라 클러스터링 기법을 이용해서 관리자에게 현재 접속한 로그와 기존 로그 데이터 셋과의 거리를 측정하여 이상접속유무를 판단할 수 있는 객관적인 지표를 제공한다. 클러스터링은 주어진 데이터 셋의 속성을 분석하여 비슷한 패턴을 가진 데이터를 묶는 방법으로, 이번 장에서는 접속주체의 기존 접속정보를 유사한 것끼리 분류하고 새로 접속한 정보가 정상인지 아닌지를 판단하는 방법을 제시한다.

클러스터링 기법을 이용해서 유사도를 측정하기 위해서는 먼저 기존 데이터 셋을 유사한 것끼리 클러스

터링 하는 알고리즘 및 특징벡터를 선정해야 한다. 그리고 새로운 접속 정보가 입력되었을 경우 가중치 유클리드 거리(Weighted Euclidean distance)를 이용하여 군집의 대푯값과 거리를 측정하여 이상탐지 유무를 확인한다.

#### 3.3.1 클러스터링 기법 및 속성 선정

클러스터링 기법을 선정할 때, VPN 접속로그 데이터 셋을 유사한 것끼리 적절히 배치하여 클러스터링 하는 알고리즘을 선정해야 하며, 관리자가 새 접속이 올 때마다 최대한 신속히 분석해서 대응해야 하기 때문에 알고리즘의 시간 복잡도를 고려해야 한다. 객체가 적절한 클러스터에 배치되었는지를 확인하고, 클러스터의 개수를 정하기 위해서 실루엣 너비(Silhouette width)를 이용한다[14]. 실루엣 너비는 각 객체의 밀집도와 의존도를 모두 고려하는 검증 방법으로 클러스터링 결과를 비교 및 평가하거나 클러스터링의 개수를 결정하기 위해 사용되며, 아래와 같이 정의된다.

$$s(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))} \quad (1)$$

$a(i)$ 는 객체  $i$ 가 속한 클러스터의 모든 객체들과의 평균 거리를 나타내며,  $b(i)$ 는 객체  $i$ 와 객체  $i$ 가 속하지 않는 클러스터와의 최소 평균 거리를 말한다. 클러스터 내의 모든 객체에 대한  $s(i)$ 의 평균값을 클러스터의 실루엣 너비라 하고, 각 클러스터의 실루엣 너비의 평균값을 전체 클러스터의 실루엣 너비라 한다.  $s(i)$ 는 -1 부터 1사이의 값을 가지며, 1에 근사하는 값이 계산될수록 주어진 데이터 셋을 명확하게 클러스터링할 수 있는 알고리즘으로 간주된다.

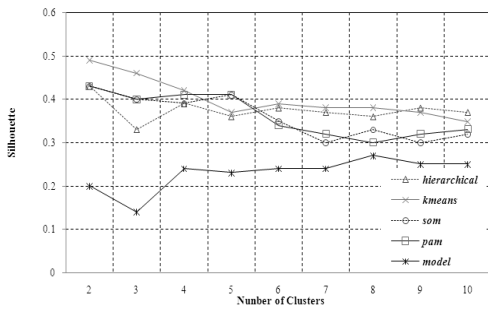
[그림 5]부터 [그림 8]은 각 알고리즘과 속성별로 주어진 데이터 셋의 특징벡터를 클러스터링 하고, 그 결과를 실루엣 너비로 나타낸 그래프이다. 모든 그래프에서 특히 Hierarchical 알고리즘과 K-means 알고리즘으로 분류된 결과가 상대적으로 높은 실루엣 너비를 보여주고 있다. 하지만 Hierarchical 알고리즘의 경우 특정한 구간에서 불안정한 결과를 도출해 내는 것으로 나타나며, 그에 반해, K-means 알고리즘이 가장 VPN 접속로그 데이터를 안정적으로 분류하는 것을 확인할 수 있다. 또한, K-means의 시간 복잡도는  $O(NK\bar{a})$  로 다른 알고리즘들보다 낮은 복잡

도를 가지고 있어, [6] VPN 접속로그를 이용하여 실시간으로 이상접속을 분석하는데 적당한 알고리즘을 최종적으로 K-Means로 선정하였다. 또한, <MAC, region, day, access\_time>의 클러스터의 개수가 3일 때, 실루엣 너비가 다른 속성들을 가지는 실루엣 너비보다 큰 값을 가지므로 이를 특징벡터로 선정하여 클러스터링을 하는데 사용하기로 한다.

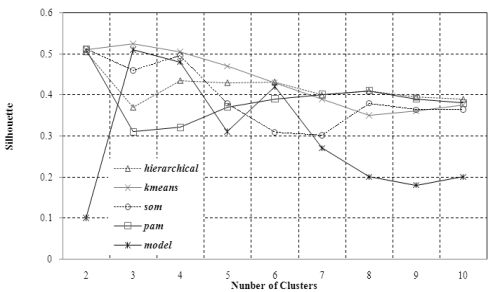
3.3.2. 이상 접속 탐지 방법

클러스터링 후에 새로 접속하는 이벤트의 이상 유무를 판단하기 위해 가중치 유클리드 거리를 사용하여 유사도를 측정하며, 거리에 따라 위험도를 정상/주의/의심으로 분류하여 새로운 접속 로그와 클러스터사이의 거리가 가까울수록 정상, 멀수록 의심으로 판단한다. 기존 데이터 셋의 각 클러스터를 대표하는 값이  $Y = (y_1, y_2, y_3, y_4)$  이고, 새로운 접속로그에 대한 정보를  $X = (x_1, x_2, x_3, x_4)$  라고 주어졌을 때, 두 정보의 유클리드 거리는 아래와 같이 정의될 수 있다.

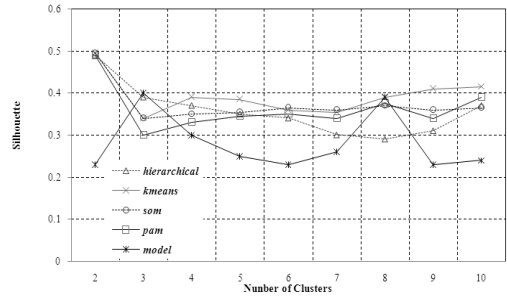
$$d(X, Y) = \sqrt{w_m(x_1 - y_1)^2 + w_i(x_2 - y_2)^2 + (x_3 - y_3)^2 + (x_4 - y_4)^2} \quad (2)$$



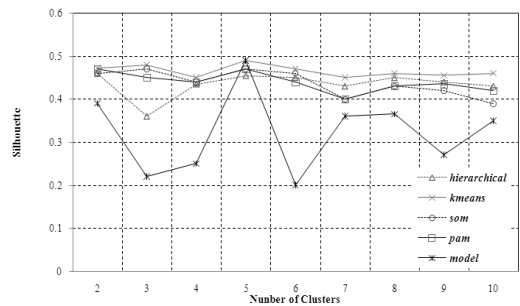
(그림 5) <MAC Address, region, day, access\_time, duration>



(그림 6) <MAC Address, region, day, access\_time>



(그림 7) <MAC Address, region, day, duration>



(그림 8) <MAC Address, region, access\_time, duration>

이때  $x_1$ 과  $y_1$ 는 MAC Address정보를 나타내는 특징벡터의 정보이고,  $x_2$ 과  $y_2$ 은 IP에 대응하는 정보이다. 이 정보들은 접속로그의 이상을 탐지하는데 있어 1차적으로 매우 중요한 척도가 될 수 있기 때문에, 그에 대한 가중치를 부여하여 이상 탐지의 효율을 향상시킬 수 있다. 예를 들어,  $w_m$ 은 새로 입력되는 MAC Address가 가지는 가중치이며 아래와 같이 계산된다.

$$w_m = \frac{N_u}{N_{total}} \quad (3)$$

이때  $N_u$ 는 기존의 MAC Address의 값과 새로 입력된 MAC Address정보가 다른 개수이고,  $N_{total}$ 은 전체 MAC Address의 수이다. 가령, 기존에 존재하는 MAC Address의 값이 다음과 같이  $[MAC_a, MAC_b, MAC_b, MAC_c]$  이고,  $MAC_b$ 가 새로 기록된 로그의 MAC Address정보일 때, 가중치  $w_m = 2/4 = 0.5$ 처럼 계산될 수 있다. 다시 말해, 새로 입력되는 MAC Address가 기존의 정보에 관측되던 내용과 같을수록 가중치는 낮아지게 되어 결과적으로



거리  $d(X, Y)$ 를 짧아지게 만든다.  $w_i$ 는 새로 입력되는 IP의 가중치를 나타내며 아래와 같이 계산된다.

$$w_i = s \left( \frac{N_u}{N_{total}} \right) \quad (4)$$

이때  $s$ 는 IP주소 체계의 각 서브 넷에 따른 분류를 구체화한 것으로, IP주소의 A 클래스가 일치할 경우 0.75, B 클래스가 일치할 경우 0.5, C 클래스가 일치할 경우 0.25, 모두 일치하지 않을 경우 1을 부여한다.  $N_u$ 는 새로 입력된 IP주소의 클래스와 기존 로그 셋에 있는 IP주소의 클래스가 다른 개수이고,  $N_{total}$ 은 전체 IP주소의 개수이다. 만약 기존 로그 셋의 IP주소가 [168.28.1.8, 173.126.43.2, 192.84.2.18, 192.28.3.139, 168.28.31.231]이고, 192.28.3.8이 새로 입력된 로그의 IP정보일 때,  $w_i = 0.25 * 4/5 = 0.2$ 의 값을 가진다. 결국 새로 입력되는 IP주소가 기존에 존재하는 IP와 일치하는 것이 많고, 그 중에서 IP의 클래스가 A보다는 C일수록 가중치가 낮아지게 되어 결과적으로 거리  $d(X, Y)$ 는 짧아지게 된다.

위와 같은 방식으로 클러스터의 개수만큼 가중치 유클리드 거리를 측정된 후 가장 작은 값을 최종 유사도로 결정하며, 이는 새로운 접속기록과 가장 짧은 거리의 클러스터 내에 있는 객체들의 행동이 가장 유사하여 이 클러스터에 포함될 확률이 높다는 것을 나타낸다. 또한, 가중치 유클리드 거리를 이용하여 구한 값이 정상인지 아닌지 판별하기 위한 기준으로 임계값을 설정한다. 임계값을 결정할 때에는 군집의 대푯값과 다량의 정상 데이터들의 거리를 반복적으로 계산한 후에 특정 값까지의 확률의 합을 의미하는 누적확률분포 (Cumulative Probability Distribution Function)를 이용하여 구한다. 이때 X축은 군집의 대푯값과 새로 입력된 값의 거리, y축은 해당 거리까지의 확률의 합을 나타내며, 기울기가 0에 가까워지는 점의 X좌표(거리)를 임계값으로 선정한다. 다시 말해, 기울기가 0이 되는 부분의 X좌표는 정상이라고 간주할 수 있는 거리의 최댓값을 의미하므로 이를 기준으로 하여 이상접속의 유무를 확인하는 것이 가능하다.

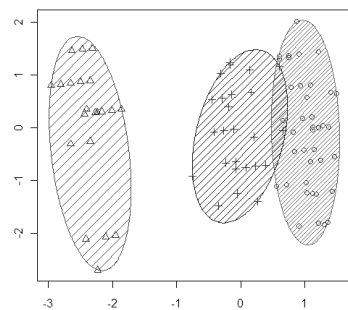
#### IV. 실험 및 결과 분석

제안하는 기법을 실제 기업환경에 적용하고 그 유효성을 판단하기 위해, 6개월 동안 특정 기업의 구성 직원들을 선정하고 정상적인 접속 로그를 수집하여 실

험을 진행하였다. 실험 데이터로는 VPN을 자주 사용하는 주체 5명의 접속로그 각각 500개를 사용하며, 각 주체 당 400개의 로그는 Training data로 사용하여 이를 클러스터링하고, 여분의 100개의 로그는 Test data로 사용하여 접속 주체 및 다른 사용자의 로그를 제대로 구분해서 탐지하는 지 확인한다.[15] 실험에 사용되는 VPN 접속로그의 속성은 <MAC, IP, Day, Access time>으로 구성되어 있으며, MAC Address값은 명칭 값으로 거리개념을 갖고 있지 않기 때문에 {0, 1, 2, ...}와 같은 수치 형식으로 변환한다. IP의 경우 IP는 특정 의미가 없는 숫자로 이루어져있기 때문에 각 시/도/구를 구별하는 지역코드로 변환하여 클러스터링 작업을 수행하며, 후에 서브넷 클래스의 일치여부를 확인하여 일정한 가중치를 부여한다. 이러한 방식으로 클러스터링한 예로 주체 1명의 Training data는 [그림 9]와 같은 클러스터를 이루었으며 각 클러스터는 다른 클러스터와 확연한 구별이 가능하다.

실험은 Microsoft Windows 7, Intel core i3, 2GHz, 2GB RAM 환경에서 공개 소프트웨어 통계 패키지인 R-project를 사용하였으며, 실험 절차는 다음과 같다.

- 1) 공개 통계 프로그램인 R-project를 이용하여 Training data를 K-means알고리즘으로 클러스터링한다.
- 2) 각 클러스터의 대푯값을 추출한다.
- 3) 2)에서 추출된 각 클러스터의 대푯값과 Test data와의 가중치 유클리드 거리를 측정된 후, 가장 작은 값을 유사도로 선정한다.
- 4) 유사도가 임계값보다 작으면 정상, 크면 비정상으로 탐지한다. 비정상 중에서도 위험도를 주의와 의심으로 분류하여 거리가 임계값의 두 배일 경우에는 의심으로 통보한다.



(그림 9) user1의 클러스터링결과 그래프

5) 1)~4)까지의 과정을 Test data set의 크기만큼 반복한다.

실험을 진행한 결과, 접속주체의 데이터 셋에 다른 사용자의 접속로그를 입력하여 이상접속으로 탐지한 결과는 [표 3]과 같이 평균 88.7%로 나타났다. Training data가 user3이고, user5의 데이터로 실험을 했을 때, 모든 접속을 이상으로 판단하였으며, 가장 높은 위험도를 나타내는 '의심'이라는 결과를 얻었다. 이는 다른 사용자의 로그가 클러스터링된 주체의 로그와 확연히 차이가 나기 때문에 좋은 결과를 얻었다고 판단할 수 있다. 한편, Training data가 user3이고, user4의 데이터로 실험을 했을 경우에는 상대적으로 낮은 탐지율을 보이고 있는데, 이는 특징 벡터의 모든 요소의 값이 각 클러스터의 대푯값 중 하나와 거의 같을 경우 결과적으로 짧은 거리를 산출하여 정상으로 탐지하게 되기 때문이다. 다시 말해, 서로 성향이 다른 사람의 로그인 경우에는 탐지율이 높지만, 그렇지 않은 경우에는 탐지율이 낮아지는 것을 확인할 수 있었다. 실제로 공격자가 실제 사용자의 사용 패턴에 대해 사전에 파악하지 못하는 경우가 대부분이기 때문에 좋은 결과를 얻을 수 있을 것이라 예상된다. 추가적으로 내부의 악의적인 이용자가 특정 사용자의 패턴에 맞게 접속을 시도할 경우에는 탐지하지 못할 가능성이 존재하기 때문에 이러한 경우에는 OTP와 같은 부차적인 인증절차를 추가하여 탐지하도록 해야 할 것이다.

## V. 결 론

현재 국내 기업의 스마트워크 도입은 대기업을 중심으로 진행되고 있지만 직원의 만족도, 업무 생산성 등의 긍정적인 측면을 고려할 때 스마트워크가 확산될 잠재력은 매우 클 것으로 전망된다. 하지만, 자료·정보 유출과 같은 정보보호에 대한 부정적인 인식은 스마트워크가 확산되는 데 걸림돌이 되고 있다.

본 논문에서는 스마트워크 환경에서 발생할 수 있는 내부 보안 위협에 빠르게 대처하기 위한 하나의 해결책으로 VPN 접속로그를 이용하여 이상 접속을 탐지할 수 있는 방법론을 제시하였다. 이 시스템은 크게 시각화 방식과 클러스터링 기법으로 유사도를 측정하는 방식으로 구성되어 있으며, 비정상접속을 실시간으로 탐지할 수 있게 함으로써 정보 유출, 계정 탈취로 인한 공격자 침입과 같은 보안 위협에 빠르게 대처할 수 있다.

이 시스템은 스마트워크의 전반적인 보안을 증대시키는 부가적인 솔루션으로 사용될 수 있으며, 접속 로그만을 이용해서 이상을 탐지하기 때문에 목표 시스템의 결함에 대한 정보 없이 탐지 가능하여 특정 취약점이나 결함을 관찰할 필요가 없고, 접속자들의 PC에 Agent를 설치할 필요가 없는 강점이 있다. 더불어 시각화를 통해 직관적인 판단을 도와, 신속한 대응이 이루어 질 수 있도록 구현하였다는 점에서 활용도가 높다 할 수 있다. 향후 업무 시 접근하는 파일의 패턴, 네트워크 트래픽 등 많은 요소에 대한 이상접속 탐지 기법과 보다 많은 사용자의 수를 포함한 로그를 활용한 연구를 진행할 계획이다.

[표 3] 실험 결과

Training data	Test data	이상접속 탐지율 (1-False Negative <sup>†</sup> )	위험도가 '의심'일 경우	False Positive <sup>‡</sup>
user1	user2	0.93	0.28	0.08
	user3	0.89	0.06	
user2	user3	0.94	0.17	0.13
	user4	1.00	0.94	
user3	user4	0.72	0.07	0.06
	user5	1.00	1.00	
user4	user5	0.74	0.23	0.16
	user1	0.95	0.26	
user5	user1	0.82	0.14	0.02
	user2	0.88	0.13	
평균		0.887	0.328	0.09

<sup>†</sup> False Negative는 주체와 다른 로그가 들어왔을 경우 탐지하지 못한 비율

<sup>‡</sup> False Positive는 Training data의 주체와 같은 데이터를 넣었을 때 이상이라고 탐지한 비율

참고문헌

- [1] “기업을 위한 스마트워크 도입, 운영 가이드북”, 방송통신위원회 한국정보화진흥원, pp. 07-13, Oct. 2010
- [2] 이형찬, 이정현, 손기욱, “스마트워크 보안 위협과 대책”, *정보보호학회지*, 21(3), pp. 12-21, 2011년
- [3] 오병근, 강성중, *정보 디자인 교과서*, pp. 99-123 2008년
- [4] Y.Livnat, J.Agutter, S.Moon, F. Erbacher, S.Foresti, “A Visualization Paradigm for Network Intrusion Detection”, *Proceedings of the 2005 IEEE Workshop on Information Assurance and Security*, pp. 92-99, June. 2005
- [5] H.Choi and H.Lee, PCAV: Internet Attack Visualization on Parallel Coordinates, vol. 3783, *ICICS 2005*, pp. 454 - 466, 2005.
- [6] R.Xu, “Survey of Clustering Algorithms”, *IEEE Transactions On Neural Networks*, vol. 16, no. 3, pp. 645-678, May 2005
- [7] J.Han and M.Kamber, *Data Mining: Concepts and Techniques*, Morgan Kaufmann, 2000
- [8] 오일석, *패턴인식*, 교보문고, 2008
- [9] W.Lee and S.J. Stolfo, “Data Mining Approaches for Intrusion Detection”, *7th USENIX Security Symposium*, pp. 79-94, Apr. 1998
- [10] Y.Guan, A. Ghorbani, “Y-means : A Clustering Method For Intrusion Detection”, *Canadian Conference on Electrical and Computer Engineering*, pp. 1-4, May 2003
- [11] V.J.Hodge, J.Austin, “A Survey of Outlier Detection Methodologies”, *Artificial Intelligence Review*, vol. 22, no. 2, pp. 85-126, Nov. 2004
- [12] OpenVPN, <http://openvpn.net>
- [13] 박동련, “R에 의한 통계그래픽스 : 강의 내용 및 방법의 논의”, *응용통계연구* 20(3), pp. 619-634, 2007년
- [14] P.J. ROUSSEEUW, Silhouettes: a graphical aid to the interpretation and validation of cluster analysis, *Journal of Computational and Applied Mathematics*, vol. 20, pp.53-65, Nov. 1987
- [15] M.Schonlau, W.Dumouchel, W.Ju, “Computer Intrusion: Detecting Masquerades”, *Statistical Science*, vol. 16, no. 1, pp. 58-74, Feb. 2001

---

 〈著者紹介〉
 

---



이 재 호 (Jae-Ho Lee) 학생회원  
 2010년 2월: 한국기술교육대학교 인터넷미디어학부 학사  
 2012년 2월: 고려대학교 정보보호대학원 금융보안학과 석사  
 2012년 3월~현재: 하나아이엔에스 재직  
 <관심분야> 금융보안, 네트워크 보안



이 동 훈 (Dong-Hoon Lee) 종신회원  
 1983년 8월: 고려대학교 경제학과(학사)  
 1987년 12월: Oklahoma University 전산학 대학원(공학석사)  
 1992년 5월: Oklahoma University 전산학 대학원(공학박사)  
 1992년 8월: 단국대학교 전자계산학과 전임강사  
 1993년 3월~1997년 2월: 고려대학교 전산학과 조교수  
 1997년 3월~2001년 2월: 고려대학교 전산학과 부교수  
 2001년 2월~현재: 고려대학교 정보보호대학원 교수  
 관심분야 : 암호프로토콜, 암호이론, USN 이론, 키 교환, 익명성 연구, PET 기술



김 휘 강 (Huy Kang Kim) 종신회원  
 1998년 2월: KAIST 산업경영학과 학사  
 2000년 2월: KAIST 산업공학과 석사  
 2009년 2월: KAIST 산업및시스템공학과 박사  
 2004년 5월~2010년 2월: 엔씨소프트 정보보안실장, Technical Director  
 2010년 3월~현재: 고려대학교 정보보호대학원 조교수  
 <관심분야> 온라인게임 보안, 네트워크 보안, 네트워크 포렌식