

# 모바일 결제 환경에서의 데이터마이닝을 이용한 이상거래 탐지 시스템

한 희 찬,<sup>†</sup> 김 하나, 김 휘 강<sup>‡</sup>  
고려대학교 정보보호대학원

## Fraud Detection System in Mobile Payment Service Using Data Mining

Hee Chan Han,<sup>†</sup> Hana Kim, Huy Kang Kim<sup>‡</sup>  
Graduate School of Information Security, Korea University

### 요 약

전 세계적으로 스마트폰 보급률이 증가함에 따라 스마트폰을 이용한 다양한 결제 서비스들이 출시되었고, 모바일 결제로 금전을 탈취하는 사례도 증가하였다. 이미 금융권은 온/오프라인 환경에서 이상거래를 탐지하기 위한 다양한 보안 조치들을 마련하였지만, 모바일 결제 환경의 보안 솔루션이나 연구들은 미비한 실정이다. 모바일 결제는 소액 결제 위주의 결제 패턴을 보이고 결제 환경도 다르기 때문에 기존의 이상거래 탐지와는 다른 모바일에 특화된 이상거래 탐지가 필요하다. 이에 본 논문에서는 국내 PG사의 실제 모바일 결제 로그를 분석하고 데이터 마이닝 알고리즘을 이용한 모바일 결제에 특화된 이상거래 탐지 시스템을 제안하였다. 해당 시스템은 1단계 탐지 모듈에서 2가지 알고리즘을 사용해 빠른 속도로 정상거래를 판별하고, 2단계 탐지 모듈에서는 고도화된 3가지 알고리즘으로 이상거래를 정확히 탐지하도록 설계하였다. 그 결과 1초에 13건 이상의 거래를 93% 이상의 정확도로 판별할 수 있었다.

### ABSTRACT

As increasing of smartphone penetration over the world, various mobile payment services have been emerged and fraud transactions have drastically increased. Although many financial companies have deployed security solutions to detect fraud transactions in on/off-line environment, mobile payment services still lack fraud detection solutions and researches. The mobile payment is mainly comprised of micro-payments and payment environment is different from other payments, so mobile-specialized fraud detection is needed. In this paper, we propose a FDS (Fraud Detection System) based on data mining for mobile payment services. The method of this paper is applied to the real data provided by a PG (Payment Gateway) company in Korea. The proposed FDS consists of two phases; (1) the first phase is focused on classifying transactions at high speed (2) the second is designed to detect abnormal transactions with high accuracy. We could detect 13 transactions per second with 93% accuracy rate.

**Keywords:** Mobile Payment Service, Fintech, Fraud Detection System, Data Mining

## 1. 서 론

전 세계적인 스마트폰의 열풍으로 모바일 시장은

포화 상태가 되었으며, 스마트폰 보급률은 PC의 보급률을 넘어섰다. 독일의 시장조사업체 TNS Infratest에 따르면, 2015년 3월 한국의 스마트폰 보급률은 83%로 국민 10명 중 8명이 스마트폰을 사용하고 있는 것으로 나타났다[1]. 더불어 모바일 결제 시장 규모는 2013년 1분기 1조 1천억 원에서 2015년 1분기 5조 원으로 2년 만에 약 5배가량 증

Received(10. 26. 2016), Modified(12. 09. 2016),  
Accepted(12. 13. 2016)

<sup>†</sup> 주저자, huer2783@korea.ac.kr

<sup>‡</sup> 교신저자, cenda@korea.ac.kr(Corresponding author)

가하였다[2]. 이처럼 모바일 환경은 많은 사람들에게 익숙해졌고 모바일을 이용한 수많은 서비스들이 우리의 생활을 더욱 편리하게 만들어주고 있다.

대표적으로 핀테크를 들 수 있는데, 핀테크는 금융(finance)과 기술(technology)의 합성어로 금융과 IT의 융합을 통한 금융서비스 및 산업의 변화를 통칭한다. 기존의 금융기법에서 벗어나 모바일, SNS(Social Network Service), 빅데이터 등을 활용한 새로운 금융서비스를 제공하며 최근 사례로는 모바일뱅킹, 애플카드를 들 수 있다. 핀테크 환경에서는 비금융기업이 지급결제와 같은 금융서비스를 직접 제공하기도 하는데, 애플의 애플페이, 알리바바그룹의 자회사인 알리페이, 이베이의 자회사 페이팔 등이 있다. 국내에서는 카카오의 카카오페이, 삼성의 삼성페이 등이 서비스 중이며, 모바일과 인터넷에 특화된 인터넷전문은행인 K뱅크와 카카오�뱅크도 곧 출범될 예정이다. 이러한 기업들은 신속함과 편리함을 강점으로 내세워 스마트폰을 기반으로 한 모바일 결제를 제공하고 있다.

이처럼 모바일 환경이 보편화 될수록 이를 악용하는 사례 역시 증가하고 있는데 대표적인 예로 스미싱(smishing)을 들 수 있다. 스미싱은 문자메시지(SMS)와 피싱(phishing)의 합성어로, 지인이나 공공기관 등을 사칭한 문자메시지를 발송하고 메시지 내의 링크를 클릭할 경우 악성코드가 포함된 어플리케이션이 설치된다. 악성코드는 사용자의 스마트폰에 들어있는 계좌 정보, 공인인증서, 개인 정보 등을 획득한 뒤 금융 거래나 소액 결제를 통해 금전을 탈취한다. 스미싱에 의한 피해는 2014년 4,917건에서 2015년 1,120건으로 줄었지만, 피해 금액은 3억 4천만 원에서 17억 4천만 원으로 5배 이상 늘어났다[3].

기술이 고도화됨에 따라 악성코드 역시 다양한 경로를 통해 유입되어 공격을 시도하고 있기 때문에, 앞으로는 소액 결제보다 큰 규모의 피해가 발생할 수 있다. 또한 신용카드 정보를 도용하거나 카드 위·변조, 도난·분실로 인한 카드 부정사용 피해 역시 무시할 수 없는 부분이다. 우리나라에서 2010년부터 2015년까지 발생한 카드 부정사용 피해 건수는 227,579건, 피해금액은 1,378억 원으로 대부분의 피해는 카드 회원의 부주의로 인해 발생하였다[4].

온라인 결제의 활성화로 인해 은행, 카드사 등 금융권에서는 이상거래 탐지 시스템 구축이 활발하게 이루어졌다. 간편결제를 서비스하는 대형IT 기업인

카카오의 카카오페이는 금융감독원의 보안 '가'군 인증을 받은 Mpay 보안방식을 적용하였고, 네이버의 네이버페이도 빅데이터를 이용한 이상거래 탐지 시스템을 구축하여 운영 중이다. 그러나 이상거래 탐지 시스템 구축에는 많은 비용이 소요되기 때문에 상대적으로 자본이 부족한 영세 업체들은 구축에 많은 어려움을 겪고 있다.

따라서 본 연구에서는 국내 PG(Payment Gateway)사의 실제 모바일 결제 로그를 이용해 이상거래를 효율적으로 탐지할 수 있는 방법을 연구하였다. 모바일 결제는 휴대폰으로 결제는 수행한다는 점이 기존의 금융거래와의 가장 큰 차이점이다. 이로 인해 휴대성이 높아져 언제 어디서나 결제를 수행할 수 있게 되었다. 그러나 아직까지 기존의 결제 환경에 익숙한 사람들은 보안에 대한 불안감이나 새로운 사용법을 익혀야 하는 점 때문에 모바일 결제로 쉽게 넘어오지 못하고 있다. 그래서 현재 모바일 결제를 이용하는 사용자들의 결제 패턴은 기존의 결제 패턴과는 다르다고 볼 수 있다. 따라서 모바일 결제 사용자 패턴 분석을 통한 모바일 환경에 적합한 이상거래 탐지가 필요하다.

개발한 프로그램은 모바일 결제 데이터를 학습해 모델을 만들고, 이를 이용해 실시간으로 결제 로그를 분석하여 이상거래를 탐지할 수 있다. 탐지 모듈을 2단계로 만들고 이중 검증을 통해 높은 탐지 정확도와 빠른 탐지 속도를 보장한다. 또한 탐지 시스템은 실제 결제 로그를 기반으로 학습하기 때문에 공격자의 공격 패턴이 달라져도 재학습을 통해 대응이 가능하며, 서버에서 동작하는 프로그램이기 때문에 고객의 휴대폰에 추가적인 보안 솔루션 등을 설치할 필요가 없다.

## II. 관련 연구

이상거래 탐지는 PG사, 은행, 카드사 등 여러 분야에서 다양한 방식으로 쓰이고 있다. 이에 따라 이상거래 탐지 관련 연구도 활발히 진행 중이다. Seong Hoon Jeong 등[5]은 이상거래 탐지 연구에 사용된 데이터 마이닝 알고리즘을 정리하고 이상거래 탐지 연구를 사용한 데이터셋, 알고리즘, 연구 관점으로 분류하여 정리하였다. 본 장에서는 선행 연구에서 사용한 결제 데이터에 따라 카테고리를 나누고 각 연구에 대해 정리하였다.

## 2.1 모바일 이상거래 탐지

모바일 환경에서 이상거래를 탐지하는 연구들은 전화 통화 기반의 데이터를 주로 이용하였고 모바일 결제 관련 연구는 아직 많이 연구되지 않았다. Constantinos S. Hilas 등[6]은 K-means clustering, Agglomerative clustering을 이용해 전화통화를 이용한 PIN(Personal Identification Number) 번호 입력 방식의 결제 데이터에서 이상거래를 탐지하고자 하였다. 고객의 개인 정보와 같은 민감한 정보를 배제하고 최소의 데이터를 이용하였다. Sharmila Subudhi 등[7]은 QS-SVM(Quarter-Sphere Support Vector Machine)을 이용해 정상 행위와 악성 행위를 분류하였다. 데이터셋은 헬싱키 대학에서 100대의 노키아 스마트폰에 소프트웨어를 설치해 9개월 동안 수집한 데이터를 사용하였으며, 통화 시간, 통화 유형, 위치, 시간, 주파수 정보를 실험에 이용하였다. QS-SVM은 90% 이상의 정확도와 10% 미만의 오답률을 보였다. Vincent S. Tseng 등[8]은 정상 전화와 사기 전화를 구분하기 위해 그래프 마이닝에 기반한 HITS(Hyperlink-Induced Topic Search) 알고리즘을 사용하였다. 실험 데이터는 Whoscall이라는 스팸·보이스피싱 차단 어플리케이션에서 제공받은 통화 관련 데이터셋을 이용하였다. 그러나 저자는 해당 방법론이 사용자의 주소록 데이터와 같은 추가적인 개인정보를 필요로 하기 때문에, 현업에 적용하기 위해서는 다른 피쳐를 사용해야 하는 문제가 있다고 하였다.

또한 모바일 어플리케이션에서 발생하는 이상거래에 대한 연구도 이루어졌다. Hee Yeon Min 등[9]은 모바일 뱅킹 앱을 이용하는 사용자의 입력 패턴 및 거래 패턴을 수집하고 SVM을 이용해 이상거래를 탐지하였다. 실험 데이터는 안드로이드 스마트폰에 가상 앱을 설치해 입력 패턴을 수집하였고, 데이터 마이닝 툴인 WEKA를 SVM 학습에 이용하였다. 실제 모바일 디바이스에서 정보를 수집하기 때문에 1ms보다 적은 빠른 탐지 시간과 98% 이상의 높은 정확도를 나타냈다.

## 2.2 금융 이상거래 탐지

금융권에서는 신용카드와 관련된 이상거래 탐지 연구가 가장 활발히 진행되었다. Seung Hyun

Kim 등[10]은 국내 카드사의 이상거래 탐지시스템에서 사용하는 스코어 방식과 룰 방식에 온라인 거래 정보를 추가 룰로 적용하였다. 그 결과 오프라인 기반 방식의 한계점을 극복하고 기존에는 탐지하지 못했던 온라인 부정거래를 탐지할 수 있었다. A. Prakash 등[11]은 MSHMM(Multiple Semi-Hidden Markov Model)에 Cuckoo search 알고리즘을 적용한 OMSHMM(Optimized Multiple Semi-Hidden Markov Model)을 제안하였다. OMSHMM은 신용카드 사기를 탐지하는 자동화 기법으로써 기존의 MSHMM보다 Precision과 Recall을 높이는 결과를 보였다. Yusuf Sahin 등[12]은 Cost-sensitive decision tree를 이용해 실제 신용카드 거래 데이터에서 부정사용을 탐지하였고 약 92%의 정확도를 보였다. K. RamaKalyani 등[13]은 Genetic algorithm을 이용해 신용카드의 부정 거래를 탐지하였다. 이 방법론은 카드 소지자의 거래 정보를 기반으로 유전과 진화를 수행하기 때문에 과거 거래에 대한 많은 데이터가 필요하다.

또한 은행 거래, 금융사기 데이터에서 이상치를 찾아내기 위한 연구도 수행되었다. Massoud Vadoodparast 등[14]은 KDA(K-means, DBSCAN, Agglomerative) clustering 모델을 만들어 은행에서 발생하는 이상거래를 탐지하고자 하였다. 논문에서는 약 360만 건의 실제 은행 거래 데이터와 32건의 이상거래 데이터를 사용하였다. 또한, 온라인 실시간 탐지와 오프라인 탐지 2가지 방법을 제시하였으며 이상거래 탐지율은 각각 68.75%, 81.25%를 보였다. Jae Hoon Park 등[15]은 의사결정나무를 사용한 정규화를 통해 기존의 탐지 룰을 개선하였다. 개선한 탐지 룰을 은행의 전자금융 사고 데이터에 적용하여 이상거래 여부를 판단하였다. Chengwei Liu 등[16]은 Random forest 알고리즘을 이용해 금융 사기를 분석하였다. CSMAR(China Stock Market & Accounting Research)의 데이터베이스 데이터를 정상 샘플과 사기 샘플로 나누어 실험에 이용하였다. 제안한 방법론은 Logistic, KNN(K-Nearest Neighbors), Decision tree, SVM과 결과를 비교해 성능을 평가하였고 Random forest 알고리즘이 88%의 정확도로 가장 좋은 성능을 보였다. 이상거래 탐지 연구에서 사용한 데이터와 알고리즘에 대하여 Table 1.에 정리하였다.

기존의 이상거래 탐지 연구들은 모바일 통화, 신

Table 1. Previous researches on fraud detection

Category	Method	Target
Mobile (6)-[9]	K-means clustering, Agglomerative clustering	Call detail record
	Quarter-Sphere SVM	Fraudulent call
	HITS	
	SVM	Mobile banking transaction
Finance (10)-[16]	Score & rule	Credit card fraud
	OMSHMM	
	Cost-sensitive decision tree	
	Genetic algorithm	
	KDA clustering	banking fraud
	Decision tree	bank e-finance accident
	Random forest	Financial fraud

용카드, 은행 등의 여러 환경에서 다양한 방법론을 기반으로 수행되었다. 그러나 모바일 분야에서 모바일 결제에 관련한 이상거래 탐지 연구는 많이 이루어지지 않았으며, 사용한 실험 데이터가 시뮬레이션 데이터인 경우가 많았다.

기존 연구는 시뮬레이션 데이터를 사용하였기 때문에 현실성이 떨어지는 단점이 있다. 본 연구에서는 이를 보완하기 위해 여러 알고리즘을 이용해 결과를 투표하는 방식으로 설계해 특정 알고리즘의 특징에 치우치지 않도록 하였고, 실제 국내 PG사의 모바일 결제 로그를 이용해 실험을 수행하였다.

### III. 이상거래 탐지 시스템

본 장에서는 이상거래 탐지에 필요한 전체적인 과정과 개발한 이상거래 탐지 시스템에 대하여 소개한다. 해당 시스템은 하나의 알고리즘을 사용하는 기존의 이상거래 탐지와는 다르게 다양한 알고리즘의 앙상블 기법을 통해 이상거래를 탐지한다. 또한 탐지 모듈을 2단계로 분리함으로써 탐지 속도를 올리고 정확도를 향상시키고자 하였다. 시스템 구현은 자바

를 이용해 CLI(Command Line Interface)로 구현하였고, 데이터 마이닝 알고리즘은 WEKA[17]에서 지원하는 자바 API를 사용하였다. 각 알고리즘은 트레이닝 데이터를 기반으로 학습해 결과를 모델로 만들어 파일로 저장한다. 탐지 모듈은 TCP 소켓을 이용해 데이터를 송수신하며 저장된 모델을 불러와 이상거래를 탐지하며 탐지 결과는 CSV(Comma Separated Value) 형식의 로그로 저장된다.

#### 3.1 데이터 소개 및 전처리

본 논문에서 사용한 데이터는 한국의 대표 PG사의 실제 모바일 결제 로그로, 2013년 1월 1일부터 2014년 12월 31일까지 총 2년 동안 수행되었던 정상거래 5,999,984건과 이상거래 15,148건으로 이루어져 있다. 정상거래는 사용자의 인증을 거쳐 정상적으로 결제가 완료된 거래를 의미하며, 이상거래는 결제까지 완료된 거래 중에서 Risk management system에서 이상거래로 판별한 거래를 의미한다. 2013년은 스미싱 범죄가 최고조에 달했던 해로 2012년에 비해 14배 이상 증가한 29,761건이 발생하였고 본 결제 로그는 해당 기간의 이상거래 중 많은 부분을 포함하고 있다[18].

결제 로그는 결제가 수행된 통신사 정보뿐만 아니라 휴대폰 번호, IP 주소 등 다양한 정보를 담고 있다. 거래 정보는 총 21개의 필드로 이루어져 있고 각 필드에 대한 정보는 Table 2.에 정리되어 있다. 괄호 안의 값은 각 필드의 데이터에 대한 예시이다. Trade ID는 각 거래에 부여되는 고유 번호를 의미하고, Cash type은 캐시구분 필드로 MA는 청구대행, MB는 복합을 의미한다. Pay mode는 거래형태로 4는 인증, 43은 자동결제이다. Service ID는 상품 타입(온라인게임, 가전, 오락 등), Merchant ID와 Enterprise ID는 상품을 판매한 상점과 법인의 고유 코드를 의미한다. 그러나 ID값과 휴대폰 번호, 이메일 주소와 같은 개인정보 관련 필드는 SHA-256 해쉬로 익명화 처리되어 있어 세부 내용은 파악할 수 없다. 통신사 정보는 데이터를 제공한 PG사의 요청으로 공란 처리하였다.

이상거래를 탐지하기에 앞서 데이터를 분석해 특징을 파악하는 과정이 필요하다. 결제 로그는 일반적인 데이터와 비교했을 때 정형 데이터에 가깝지만, 곧바로 데이터 마이닝에 적용하기에 부족하다. 데이터 분석을 통해 데이터의 특징이 잘 나타나도록 전처

Table 2. Field information

Trade ID (3675348247, ...)	Cash type (MA, MB, ...)
Pay mode (4, 43)	Status (C, None)
Phone number (SHA-256)	Merchant user ID (SHA-256)
Email flag (Y, N)	Email front address (SHA-256)
Email back address (SHA-256)	Encrypted IP address (SHA-256)
Approval client version (1001, B001, Bs1M, ...)	Encrypted user email (SHA-256)
Enterprise ID (SHA-256)	Telecommunication company ( - )
Merchant ID (SHA-256)	Service ID (SHA-256 → online game)
Product price (3000, 49500, ...)	Approval time (132541 → 13)
Approval day (20130101 → Tue)	Sales type (0, 1)
IP address(Country) (183.108.136.000 → KR)	Product price cluster (Price_1, Price_2, ...)

리를 하면 탐지율을 올릴 수 있다.

따라서 21개의 필드 중에서 중요도가 낮은 필드는 제외하고 필드의 정보를 재가공하는 과정을 수행하였다. 먼저 인증일자, 인증시간, IP주소 필드는 자세한 데이터를 그대로 사용할 경우 분류에 과적합이 생길 수 있기 때문에 추상화 과정을 거쳤다. 인증일자는 요일, 인증시간은 시간, IP 주소에서는 국가 정보만 추출하였다. 그리고 거래금액은 K-means 알고리즘을 통해 추가 분석을 수행하였다. 거래 금액에 따른 고객 그룹을 구별 하기 위해 K=5로 설정해 하위 그룹, 중간 그룹 2개, 상위 그룹 2개로 클러스터로 나누었다. 각 클러스터의 최대값, 최소값을 구한 결과 하위 그룹은 0원~9,371원, 중간 그룹1은 9,372원~29,689원, 중간 그룹2는 29,690원, 상위 그룹 1은 69,632원~149,465원, 상위 그룹2는 149,466원~691,250원으로 나타났다. 각 거래금액을 그룹과 매칭하고 매칭 결과를 Product price cluster 필드로 만들어 저장하였다. 서비스 아이디 필드는

SHA-256 해쉬가 적용되어 있으나 각 해쉬값과 매칭되는 카테고리 정보(온라인게임, 방송/신문, 영화 등)로 변환하였다.

### 3.2 피쳐 선택

가공이 끝난 22개의 필드 중 의미 있는 필드를 선별하기 위해 Weka에서 제공하는 Feature selection 알고리즘 중 12가지를 사용하였다. 사용한 알고리즘의 세부 내용은 Table 3.에 정리하였다.

Feature selection 알고리즘은 Evaluator와 Search method로 구성되어 있는데, Evaluator에서 각 속성을 평가하고 Search method에서 평가 결과를 기준으로 속성을 검색하거나 순위를 매긴다. Search method 중 Ranker를 사용한 알고리즘과 Ranker가 아닌 다른 method(greedy stepwise, ranksearch)를 사용한 알고리즘으로 분류를 나누었다. Ranker method를 사용한 Evaluator는 단일 속성 평가 방식, Non-ranker method를 사용한 Evaluator는 속성 종속 집합 평가 방식이다. 단일 속성 평가 방식은 전체 데이터에서 각 속성이 지니는 독립적인 성질을 파악하기 때문에 Ranker method를 이용해 속성의 순위를 나열하는 것이 적절하다. 종속 집합 평가 방식은 여러 속성들 사이의 연관성을 파악하기 때문에 속성의 집합 내부를 검색하는 과정으로 결과를 출력하는 방법이 적합해 Non-ranker method를 사용하였다.

이와 같이 2가지 타입의 Search method와 12가지의 Evaluator를 사용하여 속성의 독립적인 특

Table 3. Feature selection algorithm

Category	Evaluator
Ranker	ChiSquaredAttributeEval
	FilteredAttributeEval
	GainRatioAttributeEval
	InfoGainAttributeEval
	OneRAttributeEval
	SymmetricalUncertAttributeEval
	ReliefFAttributeEval
Non-ranker	CfsSubsetEval
	ClassifierSubsetEval
	WrapperSubsetEval
	ConsistencySubsetEval
	FilteredSubsetEval

징과 중속적인 특징을 분석하고 피처를 선별하였다. Ranker는 중요한 필드 순으로 순위를 출력하고 Non-ranker의 경우 중요 필드만 결과로 출력한다. 일부 Non-ranker 알고리즘은 2~3개의 적은 수의 필드만 출력하였기 때문에 결과를 그대로 사용하기 어려워 결과를 종합적으로 정리하는 과정이 필요하였다. 따라서 Ranker의 결과 중 상위 10개 필드와 Non-ranker에서 출력한 필드의 빈도를 측정하여 빈도가 높은 피처를 뽑아 최종 선정 피처로 사용하고 하였다.

그러나 Feature selection 알고리즘은 기계적인 계산을 통해 보편적인 선택을 도와주는 방법이기 때문에 의도와 다른 피처가 선택될 가능성이 있다. 따라서 전문가의 경험에 의한 주관적인 의견을 반영하는 휴리스틱 기법도 함께 적용하는 것이 중요하다. 따라서 정상거래와 이상거래의 각 필드를 구성하는 데이터의 비율을 분석하였다.

정상거래와 이상거래의 차이가 드러나는 필드는 대표적으로 인증 일자(Approval day)와 거래 금액(Product price)이 있다. 인증 일자는 사용자가 거래를 위해 인증을 수행한 요일을 의미하는 필드로, Fig.1.에 정상거래와 이상거래가 수행되었을 때의 각 요일 비율을 정리하였다.

정상거래의 인증 일자는 주말에 비해 평일의 비율이 높았고 평일 중에서도 월요일, 화요일, 수요일의 비율이 높았다. 이는 일반적인 상품 배송이 2~3일 정도 소요되기 때문에 주말보다 평일, 평일 중에서도 앞 요일에 거래를 많이 요청하는 것으로 유추할 수 있다. 그러나 이상거래의 경우는 정상거래의 비율과 완전히 다른 패턴을 보이고 있다. 이상거래는 평일보다 주말에 많이 이루어진 것으로 보아 주말에 이상거래를 수행하기 수월했던 것으로 유추할 수 있다.

Fig.2.의 거래 금액에서 정상거래는 대부분이 1

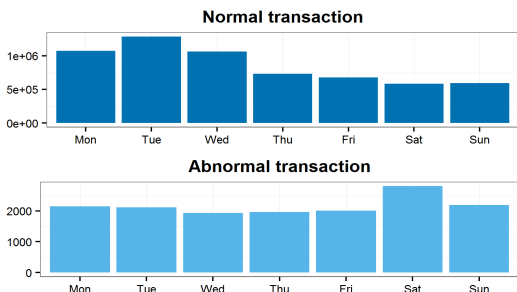


Fig. 1. Distribution of approval day

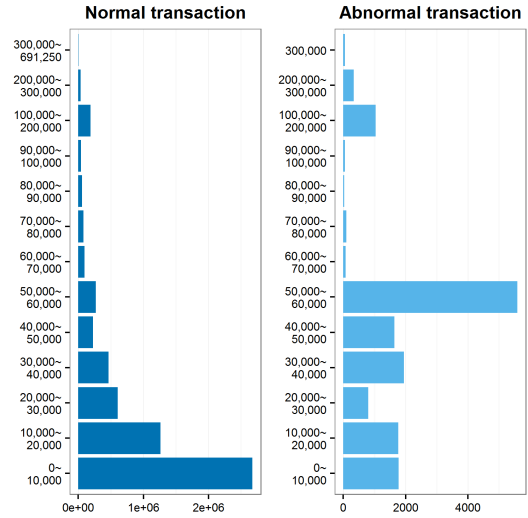


Fig. 2. Distribution of product price

만원 미만의 소액 결제였고 최대 결제 금액은 691,250원이었다. 그러나 이상거래는 정상거래에 비해 소액 결제의 비율이 낮은 편이었고 오히려 5만원 대의 거래가 높았다. 또한 최대 결제 금액은 300,000원이었다. 이는 금전 탈취를 목적으로 하는 이상거래의 특징을 보여준다. 그러나 최대 금액이 정상거래에 비해 낮은 이유는 모바일 결제가 30만 원 이상 결제 시에 공인인증서를 요구하기 때문으로 유추된다.

이와 같이 휴리스틱 기반의 데이터 분석을 통해 기계적인 피처 선택 알고리즘으로 분석할 수 없는 부분들을 보완하였고, 두 가지 기법에서 선정된 피처 중 중복으로 선정된 필드를 최종 피처로 사용하였다. 최종 피처는 총 10개로 Table 4.에 정리하였다.

Table 4. Final feature

Enterprise ID	Merchant ID
Product price	Service ID
Approval day	Approval time
IP address(Country)	Sales type
Product price cluster	Telecommunication company

### 3.3 시스템 구조

이상거래 탐지 시스템은 1단계 탐지 모듈과 2단계 탐지 모듈로 구성되어 있다. 먼저 1단계 탐지 모듈

에서 모바일 결제 데이터를 받아 전처리를 수행한다. 전처리 과정에서는 결제 데이터의 사용하지 않는 필드를 제외한 후에 요일, 시간 등의 값을 가공하고 거래금액 필드에서 거래금액 구간을 계산한다. 전처리가 끝난 데이터는 2가지의 데이터마이닝 알고리즘을 이용해 이상거래 여부를 판별한다. 2가지 알고리즘 모두 해당 거래를 이상거래로 판별할 경우 2단계 탐지 모듈로 보내 추가 분석을 수행하고, 정상거래로 판별된 거래는 정상적으로 승인되어 거래가 완료된다. 2단계 탐지 모듈은 1단계 탐지 모듈로부터 전처리된 데이터를 받았기 때문에 바로 3가지 알고리즘을 이용해 분석을 수행한다. 3가지 알고리즘 중 2가지 이상의 알고리즘이 이상거래로 판별할 경우 최종적으로 이상거래로 탐지한다. 위에서 설명한 탐지 과정을 Fig.3.에 정리하였다.

### 3.4 알고리즘

탐지 모듈에서 사용한 알고리즘은 총 5가지 (C4.5, Naïve bayes, CART, SMO, Random forest)로, 데이터 마이닝 분야에서 널리 쓰이는 알고리즘을 선별하였다.

#### 3.4.1 1단계

1단계 탐지 모듈에서는 빠른 탐지를 위해 과정이 직관적이고 속도가 빠른 알고리즘 2가지를 선택하여 사용하였다.

##### 3.4.1.1 Decision tree(C4.5)

Decision tree는 의사결정트리로 불리며 분류

알고리즘으로 널리 쓰이고 있다. 데이터의 속성을 트리 구조로 분할하고 정보이론(information theory)을 이용하여 데이터 집합을 분류한다. Decision tree는 분류 규칙과 결과를 시각화하여 과정을 쉽게 이해할 수 있어 탐지 결과를 사용자에게 명확하게 제시하고 설명할 수 있다. 또한 계산량이 적어 많은 컴퓨팅 작업을 필요로 하지 않기 때문에 빠른 속도로 분석을 수행할 수 있다. 본 시스템에서는 Decision tree 기반의 여러 알고리즘 중 C4.5 알고리즘을 사용하였다. C4.5는 기본적인 Decision tree 알고리즘인 ID3에서 수치형 속성 처리를 해결하고 과적합을 가지치기(pruning)로 해결하는 등의 몇 가지 단점을 보완한 알고리즘이다.

##### 3.4.1.2 Naïve bayes

1단계 탐지 모듈에서는 규칙에 의한 분류인 Decision tree와는 다른 확률 기반의 분류 알고리즘인 Naïve bayes 알고리즘도 사용하였다. Naïve bayes 분류는 특성들 사이의 독립을 가정하는 베이즈 정리를 사용한 분류 알고리즘으로 문서 분류에서 널리 쓰이고 있다. 해당 알고리즘은 가정과 공식 자체가 간단하기 때문에 속도가 빠르고 분류에 필요한 트레이닝 데이터가 적어도 좋은 성능을 보인다.

#### 3.4.2 2단계

2단계 탐지 모듈은 1단계 탐지에서 이상거래로 탐지된 거래를 추가로 판별하기 위해 보다 고도화된 알고리즘 3가지를 사용하였다.

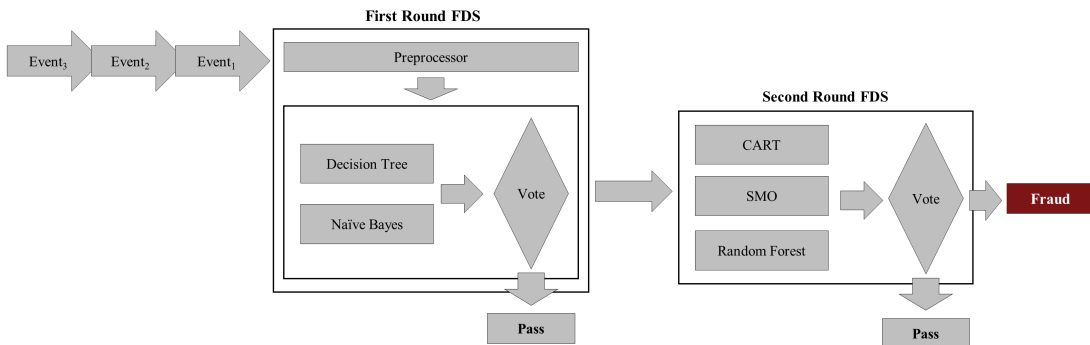


Fig. 3. System architecture

### 3.4.2.1 CART(Classification And Regression Trees)

CART는 통계학 분야에서 개발된 Decision tree 알고리즘 중의 하나로 Gini index와 이진 트리를 기반으로 동작한다. 엔트로피 매트릭스를 사용해 후보 트리를 여러 개 생성하고 그 중에서 최적의 트리를 찾는다. 또한 데이터를 내부적으로 트레이닝과 테스트로 구분해 사용함으로써 트리의 과적합을 줄인다.

### 3.4.2.2 SMO(Sequential Minimal Optimization)

SMO는 SVM의 최적화 문제를 효율적으로 풀기 위해 고안된 알고리즘이다. SVM은 패턴 인식, 자료 분석을 위한 지도 학습 알고리즘으로, 데이터를 공간으로 사상시키고 초평면(hyperplane)을 만들어 분류하는 방식이다. 비선형 분류의 경우는 커널 함수를 사용해 차원을 변환하여 분류를 수행한다. SMO는 SVM의 QP(Quadratic Programming)의 최적화 문제를 작은 QP 문제로 줄여 SVM의 장점인 낮은 과적합과 높은 정확도를 그대로 유지하고 메모리 사용을 줄이고 분석 속도를 높였다.

### 3.4.2.3 Random forest

Random forest는 Decision tree의 앙상블 기법으로, 여러 개의 트리를 만들어 학습하고 이를 결합해 최종적인 결과를 도출한다. 이를 통해 노이즈 데이터와 결측치 데이터에 대한 패턴 분류의 일반화율을 높이고 Decision tree가 특정 학습 데이터에 민감하거나 불안정한 성질을 보완해 과적합을 줄였다.

## IV. 실험

개발한 이상거래 탐지 시스템이 이상거래와 정상거래를 효과적으로 탐지하는지 확인하기 위해 결제 로그 중 일부를 샘플링하여 실험을 진행하였다. 정상/이상거래 결제 로그에서 알고리즘 학습을 위한 트레이닝 데이터 5,000건, 탐지율 측정을 위한 테스트 데이터 5,000건을 전체 기간에서 랜덤하게 추출하였다.

먼저, 트레이닝 데이터의 정상/이상거래 비율에 따른 탐지 성능의 차이를 알아보기 위해 트레이닝 비율을 변경하면서 테스트를 진행하였다. 1단계 탐지 모듈은 일반적인 트레이닝 비율인 8:2(정상/이상)으

로 설정하였다. 그러나 2단계 탐지 모듈은 1단계에서 이상거래로 판별된 거래에 대해서만 판별을 수행하기 때문에 트레이닝의 비율을 8:2, 5:5, 2:8로 변경하면서 차이를 살펴보았다. Table 5.의 결과를 보면 이상거래의 비율이 높아질수록 이상거래에 과적합되어 미탐률(false negative rate)은 낮아지지만 오탐률(false positive rate)이 높아지는 것을 확인할 수 있다. 그리고 3가지 실험 모두 93% 이상의 탐지율을 보였다. 따라서 가장 낮은 오탐률을 보이는 8:2 비율을 선정해 2차 실험을 진행하였다.

2차 실험은 1차 실험에서 선정한 비율을 사용해 추가적으로 정확도를 측정하였다. 실험 환경을 실제 거래 환경과 가장 유사하도록 설정하기 위해 결제 로그의 전체 기간에서 랜덤하게 추출한 정상거래 4,000건, 이상거래 1,000건 총 5,000건의 테스트 데이터를 사용하였다. 또한 특정 데이터에 따라 탐지율의 변화가 있는지 확인하기 위해 10번의 샘플링을 통해 10번을 실험을 진행하였다. 실험 결과는 Table 8.에 정리하였다. 본 논문에서 제안한 이상거래 탐지 시스템은 1단계 탐지 모듈에서 92% 이상의 탐지율과 4% 이하의 오탐률로 이상거래를 탐지하였다. 그리고 2단계 탐지 모듈의 추가 검증을 통해 약 1%의 오탐률 감소와 약 0.6%의 탐지율 상승 효과를 보여 최종적으로 93% 이상의 탐지율과 3% 이하의 오탐률로 이상거래를 탐지해낼 수 있었다.

기존의 이상거래 탐지 연구는 약 80%에서 90% 정도의 정확도를 보였으나 해당 시스템은 기존의 연구들과 비슷하거나 더 높은 성능을 보였다. 더불어 기존 연구는 정확도에만 초점을 맞춰 설계하였으나, 본 시스템은 오탐률을 낮추고 탐지 시간도 최소화하고자 하였다. 그 결과, 5,000건의 거래를 분석하는데 6분 17초가 소요되어 1초에 13건 이상의 거래를 분석하는 것이 가능하였다.

그러나 해당 시스템을 실제 서비스에 적용하기에는 부족한 부분들이 존재한다. 모바일 결제 로그의 일부 필드가 개인정보보호 이슈로 인해 익명화 처리가 되어 있어 분석이 어렵고 활용도가 떨어지는 문제

Table 5. First experiment result

Normal/ Abnormal	Accuracy	False negative	False positive
8:2	93.28%	3.90%	2.82%
5:5	93.30%	3.58%	3.12%
2:8	93.22%	3.54%	3.24%



Table 6. Second experiment result

		Accuracy	False negative	False positive
Avg	1st round	92.58%	3.74%	3.69%
	2nd round	93.19%	4.03%	2.78%
Max	1st round	92.94%	4.00%	3.90%
	2nd round	93.34%	4.14%	3.00%
Min	1st round	92.30%	3.42%	3.44%
	2nd round	93.00%	3.92%	2.58%

가 있었다. 기업 내에서 익명화 처리가 되기 전의 정보를 분석하고 그 결과를 이용해 이상거래 탐지를 수행한다면 더 높은 정확도로 탐지를 수행할 수 있을 것으로 사료된다. 분석 속도 측면에서는 해당 시스템의 전처리, 저장 과정을 줄여 분석을 최적화하고, 시스템을 병렬로 동시에 구동함으로써 단점을 보완할 수 있다.

## V. 결 론

스마트폰을 이용한 간편결제 환경이 보편화되면서 모바일 결제 규모는 더욱 증가하고 있는 추세이다. 그러나 모바일 결제 환경에 적합한 이상거래 탐지 연구는 1차원적인 수준에 머무르고 있다. 단순한 탐지 탐지 규칙이 알려질 경우 쉽게 우회될 수 있고, 탐지 규칙을 업데이트하기 위해서는 분석가의 추가적인 분석을 통해 주관적인 판단 기준을 반영해야 한다는 단점이 있다.

이에 본 논문에서는 데이터 마이닝 알고리즘을 이용해 모바일 결제 환경에서 효율적으로 이상거래를 탐지하는 시스템을 제안하였다. 탐지 모듈을 2단계로 나눔으로써 1단계 탐지에서 모바일 결제의 장점인 빠른 속도를 보장하고, 2단계 탐지를 통해 오탐률을 낮춰 안정성을 높이는 효과를 보였다. 또한 해당 시스템은 트레이닝 데이터를 이용해 각 알고리즘을 학습하기 때문에 새로운 이상거래 패턴을 트레이닝 데이터에 적용해 새로운 모델을 만들어 업데이트할 수 있는 확장성이 있다.

또한 모바일 결제 로그를 분석하여 모바일 결제

환경에 적합한 이상거래 탐지를 구현하였으며, 실제 데이터에서도 높은 정확도와 낮은 오탐률을 보여 현업에서 운영 중인 이상거래 탐지 시스템 연구에 활용될 수 있을 것으로 기대한다.

그러나 해당 시스템은 낮은 오탐률에 비해 아직 높은 미탐률을 보이고 있는데, 미탐은 기업의 입장에서 고객의 피해보상과 직결되기 때문에 낮춰야 할 필요가 있다. 이에 향후 미탐률을 낮추기 위한 추가적인 연구가 필요하다.

## References

- [1] Digieco, [http://www.digieco.co.kr/KTFront/report/report\\_issue\\_trend\\_view.action?board\\_seq=10712&board\\_id=issue\\_trend](http://www.digieco.co.kr/KTFront/report/report_issue_trend_view.action?board_seq=10712&board_id=issue_trend), Jan. 2016.
- [2] BNK Finance Research Institute, <http://www.bnkfg.com/download?seq=3632>, Sep. 2015.
- [3] Kyung Wook Min, <http://blog.naver.com/minkyungwook/220800480732>, Aug. 2016.
- [4] Byung Do Min, <http://badmin.net/220449860352?Redirect=Log&from=postView>, Aug. 2015.
- [5] Seong Hoon Jeong, Hana Kim, Youngsang Shin, Taejin Lee, and Huy Kang Kim, "A Survey of Fraud Detection Research based on Transaction Analysis and Data Mining Technique," *Journal of the Korea Institute of Information Security and Cryptology*, vol. 25, no. 6, pp. 1525-1540, Dec. 2015.
- [6] Hilas, Constantinos S., Paris A. Mastorocostas, and Ioannis T. Rekanos. "Clustering of Telecommunications User Profiles for Fraud Detection and Security Enhancement in Large Corporate Networks: A case Study," *Applied Mathematics & Information Sciences*, vol. 9, no. 4, pp. 1709-1718, Jul. 2015.
- [7] Subudhi, Sharmila, and Suvasini Panigrahi. "Quarter-Sphere Support Vector Machine for Fraud Detection in

- Mobile Telecommunication Networks," *Procedia Computer Science*, vol. 48, pp. 353-359, 2015.
- [8] Tseng, V.S., Ying, J.C., Huang, C.W., Kao, Y., and Chen, K.T., "FraudDetector: A Graph-Mining-based Framework for Fraudulent Phone Call Detection," *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 2157-2166, Aug. 2015.
- [9] Hee Yeon Min, Jin Hyung Park, Dong Hoon Lee, and In Seok Kim, "Outlier Detection Method for Mobile Banking with User Input Pattern and E-finance Transaction Pattern," *Journal of Internet Computing and Services*, vol. 15, no. 1, pp. 157-170, Feb. 2014.
- [10] Seung-Hyun Kim, Huy Kang Kim, and Eun jin Kim. "A Study on the Improvement of FDS Effectiveness over the Case Analysis on E-commerce Credit-card Fraud-to-sales," *Journal of Knowledge Information Technology and Systems (JKITS)*, vol. 10, no. 6, pp. 723-734, Dec. 2015.
- [11] Prakash, A., and C. Chandrasekar, "An optimized multiple semi-hidden markov model for credit card fraud detection," *Indian Journal of Science and Technology*, vol. 8, no. 2, pp. 165-171, Jan. 2015.
- [12] Sahin Yusuf, Serol Bulkan, and Ekrem Duman, "A cost-sensitive decision tree approach for fraud detection," *Expert Systems with Applications*, vol. 40, no. 15, pp. 5916-5923, 2013.
- [13] RamaKalyani, K., and D. UmaDevi, "Fraud detection of credit card payment system by genetic algorithm," *International Journal of Scientific & Engineering Research*, vol. 3, no. 7, Jul. 2012.
- [14] Vadoodparast, Massoud, and Abdul Razak Hamdan, "Fraudulent Electronic Transaction Detection Using Dynamic KDA Model," *International Journal of Computer Science and Information Security*, vol. 13, no 3, pp. 90-99, Mar. 2015.
- [15] Jae Hoon Park, Huy Kang Kim, and Eun jin Kim, "Effective Normalization Method for Fraud Detection Using a Decision Tree," *Journal of The Korea Institute of Information Security & Cryptology*, vol. 25, no. 1, pp. 133-146, Feb. 2015.
- [16] Chengwei Liu, Yixiang Chan, Syed Hasnain Alam Kazmi, and Hao Fu, "Financial Fraud Detection Model: Based on Random Forest," *International journal of economics and finance*, vol. 7, no. 7, pp. 178-188, Jun. 2015.
- [17] Hall. M., Frank. E., Holmes. G., Pfahringer. B., Reutemann. P., and Witten. I.H., "The WEKA data mining software: an update," *ACM SIGKDD explorations newsletter*, vol. 11, no. 1, pp. 10-18, Jun. 2009.
- [18] The Fact, <http://news.tf.co.kr/read/life/1399255.htm>, Aug. 2014.

### 〈 저자 소개 〉



한 희 찬 (Hee Chan Han) 학생회원  
 2015년 2월: 서울시립대학교 수학과 졸업  
 2015년 3월~현재: 고려대학교 정보보호학과 석사과정  
 <관심분야> 데이터마이닝, 데이터 분석, IoT 보안



김 하 나 (Hana Kim) 학생회원  
 2013년 2월: 서울여자대학교 정보보호학과 졸업  
 2015년 8월: 고려대학교 정보보호학과 석사  
 2015년 9월~현재: 고려대학교 정보보호학과 박사과정  
 <관심분야> 온라인게임 보안, 데이터 마이닝



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~2014년 12월: 고려대학교 정보보호대학원 조교수  
 2015년 1월~현재: 고려대학교 정보보호대학원 부교수  
 <관심분야> 온라인게임 보안, 네트워크 보안, 네트워크 포렌식