

# 효과적인 초동 수사를 위한 지도/준지도 학습 기반의 범죄 사건 정보 예측 및 용의자 후보군 특정 기술 연구

백명선, 박원주, 지중호\*, 신현정\*, 장광호\*\*, 이용태  
한국전자통신연구원, \*아주대학교, \*\*치안정책연구소

sabman@etri.re.kr, wjpark@etri.re.kr, baical77@ajou.ac.kr, shin@ajou.ac.kr,  
pathfinder@police.go.kr, ytleee@etri.re.kr

## A Study on the Prediction of Crime Case Information and Identification of Suspect Candidates based on Supervised/Semi-Supervised Learning Technique for Efficient Preliminary Investigation

Myung-Sun Baek, Wonjoo Park, Jong Ho Jhee\*, Hyunjung Shin\*, Kwangho Jang\*\*,  
Yong-Tae Lee  
ETRI, \* Ajou University, \*\*Police Science Institute

### 요 약

범죄사건 접수 초기에 신속하고 효과적으로 해당 범죄에 대응하기 위해서는 발생한 범죄 유형 및 용의자 후보군 대한 정보를 현장 파견 요원 및 수사관들에게 제공하는 것이 매우 효과적인 방법이다. 본 논문에서는 지도학습 기반의 인공지능 기술을 활용하여 접수된 범죄의 유형을 예측하는 기술을 설계한다. 설계된 기술은 범죄사건 접수 데이터인 텍스트 기반의 범죄사건 개요 정보를 활용하여 해당 범죄가 21 종의 중분류 기반 범죄유형 중 어느 유형에 해당하는지 예측하여 정보를 제공한다. 또한 준지도학습 기반의 범죄사건에 대한 유력한 용의자 후보군 특정 기술을 설계한다. 용의자 후보군 특정을 위해, 기존에 누적된 여러 종류의 치안데이터가 활용된다. 상기 이종의 데이터를 활용하여 사건과 인물 사이의 관계를 기반으로 하는 다계층 네트워크를 생성한 후, 신규 사건의 용의자/사건관련자 등을 준지도학습기반의 기계학습기법으로 예측한다. 이를 통하여 설계된 기술은 사건의 위험수준을 초기에 인지하여 신속한 초동대응방법을 도출하고, 유력한 용의자 후보군을 특정하는 것에 주요한 참고 자료로 활용될 수 있다.

### I. 서 론

최근 국민의 안전을 보장하기 위해 스마트 폴리싱 등 고도화된 치안기술개발을 위한 인공지능/빅데이터의 적용 기술이 활발하게 연구되고 있다 [1]-[3]. 본 논문에서는 누적된 범죄/수사데이터와 실시간으로 수집되는 다양한 치안데이터를 활용하여 지도학습/준지도학습을 수행하고, 이를 통해 신규 범죄사건 접수 시, 효과적인 초동 수사방안 수립이 가능한 인공지능 기반 치안 기술을 설계하고 성능을 검증한다. 첫번째 기술은 누적된 치안관련 통계데이터를 활용하여 인공지능 학습(지도학습)을 수행하고, 신규 접수되는 사건의 개요 정보(텍스트 정보)로부터 범죄의 유형 추론하는 기술이다. 해당 기술은 사건 내용을 포함하는 텍스트 데이터를 사용하여 중분류 기반 21 종의 범죄 유형을 고려하여 해당 유형중 한가지의 범죄 유형으로 범죄유형을 추론할 수 있다. 두번째 설계 기술은 누적된 여러 종류의 치안데이터를 활용하여 용의자 후보군을 특정하는 기술이다. 개별 사건에 대한 용의자 후보군을 추론하기 위해, 기존 사건 및 인물 별로 유사 사건/인물 계층 네트워크를 생성하고, 생성된 네트워크를 기반으로 용의자 후보군을 특정한다. 상기 개발된 기술을 활용하면 실제 치안현장에서 신규 사건 접수 시 사건유형/위험성 추론 결과를 참고하여 현장인력 파견/배치 등을 효과적으로 수행할 수 있으며 유력한 용의자 후보군 특정에 참고가 될만한 결과를 제공받을 수 있다.

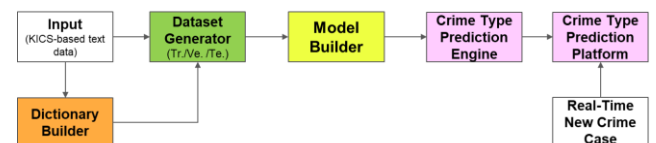


그림 1. 범죄유형 추론 시스템 구조도

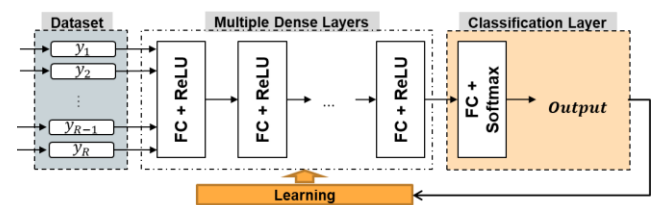


그림 2. 범죄유형 추정 엔진 구조도

### II. 범죄유형 추론 기술 설계

본 절에서는 텍스트 기반의 사건 접수데이터로부터 범죄 유형을 추론하는 기술을 설계한다. 기존에는 강력범죄 7 종에 대해 간략하게 유형을 추론한 바 있다 [2]. 그러나 강력범죄 7 종만을 추론하는 기존 기술은 다양한 실제 범죄 유형을 커버하지 못하므로, 현장적용에 한계가 있다. 따라서 본 논문에서는 범죄 중분류 기반의 21 종의 범죄 유형을 고려하여 범죄 유형 추론 기술을 설계하였다. 또한 범죄유형 추론 방식을 고도화 하여 기존 기술의 성능적 한계 또한 개선하였다. 본 절에서는 형사사법정보시스템 (KICS: Korea Information System of Criminal Justice Services)양식에 따른 사건관련

텍스트데이터를 고려한다. 해당 양식에서 텍스트데이터는 범죄 사건 개요를 약 150~500 자 내외로 기술하고 있다.

그림 1 은 제안된 범죄 유형 및 위험도 예측 시스템의 구조도이다. dictionary builder 블록은 텍스트 기반 데이터 소스에서 Feature Keyword 추출 알고리즘인 word rank 와 tf-idf 를 이용해 유의미한 단어를 추출하고 이를 통해 키워드 사전을 구축한다. dataset generator 블록은 상기 생성된 키워드 사전 및 입력된 KICS 데이터를 이용하여 데이터 셋을 생성한다. 생성된 데이터 셋은 training dataset, verification dataset, test dataset 이다. 상기 데이터셋을 활용하여 딥러닝 기반의 인공지능 추론기를 생성한다. 생성된 추론기는 training 데이터셋을 기반으로 훈련되고 verification 을 통해 검증된다. 그림 2 는 딥러닝 기반의 인공지능 추론기 구조를 보여준다. 그림에서와 같이 훈련기간 동안 training dataset 을 입력받아 추론기를 훈련하고 마지막으로 test dataset 을 통해 성능이 평가된다. 최종적으로 GUI 가 탑재된 플랫폼형태의 시스템이 개발되어 실시간으로 신규 사건내용을 입력받아 해당 범죄유형을 추론하는 것이 가능하다.

### III. 용의자 후보군 추론 기술

본 절에서 다루는 기술은 기존 누적된 다종의 사건데이터 (KICS, 성범죄자, 총포소지자, 112 신고데이터)를 기반으로 인물 및 사건 데이터를 추출하여 인물-사건 계층네트워크를 생성하고, 생성된 네트워크를 기반으로 사건 별 용의자 후보군을 추론하는 기술이다. 상기한 다양한 이종의 데이터들의 사건내용을 추출하여 관련된 사건을 연결하고 사건네트워크를 생성한다. 또한 인물정보를 사용하여 연관된 인물 기반의 인물네트워크를 생성한다. 이를 계층적으로 연결하여 인물-사건 계층네트워크를 그림 2 와 같이 생성한다. 상기 생성된 네트워크를 기반으로 그래프 기반 준지도학습 알고리즘 [4]을 통해 각 개별 사건에 따른 용의자 후보군을 추론하여 특정한다. 후보군을 추정하는 기술은 아래 수식과 같다.

$$f = (I + \mu_a L^{(intra)} + \mu_b L^{(inter)})^{-1} Y \quad (1)$$

상기 수식에서  $I$  는 identity matrix 이고  $L^{(inter)}$  와  $L^{(intra)}$  는 각각 네트워크 내 및 네트워크 사이의 weight matrix 기반 graph Laplacian [4]이다.  $\mu_a$  와  $\mu_b$  는 각각의 graph Laplacian 의 조절 계수이다. 또한  $Y$  는 기존 데이터 베이스에 저장된 알고있는 인물들의 정보이며  $f$  는 추정된 값이 된다. 추정된 값의 상위 값을 갖는 인물들을 해당 사건과 연관된 용의자 후보군으로 설정할 수 있다. 본 논문에서는 상위 1.5%에 속하는 인물들을 용의자 후보군으로 설정하였다.

### IV. 개발 기술의 성능 검증

본 절에서는 설계된 두가지 치안 기술의 성능을 F1 score 를 활용하여 검증한다. 표 1 은 범죄유형 추론 기술의 성능을 검증한 결과이다. 본실험에서는 II 절에서 설명한 바와 같이 총 21 종의 범죄 유형을 포함하는 5000 개의 치안 데이터를 활용하였다. 데이터셋 구성은 다음과 같다.

- training: 3000, verification: 1000 test: 1000

표 1 에서와 같이 각 범죄 유형별 소폭의 성능 차이는 있지만, 평균 범죄유형 추론 성능은 0.89 점으로 우수한 추론 성능을 제공하는 것을 확인할 수 있다.

그림 4 는 용의자 후보군 추론 기술의 성능검증 결과이다. 사건 및 인물관련 실험환경은 다음과 같다.

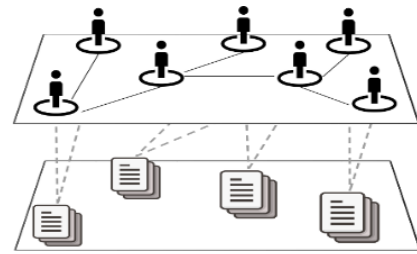


그림 3. 누적 치안데이터 연계 기반 인물-사건네트워크

표 1. 범죄유형 추론 기술 성능 검증 결과

	범죄유형	f1-score	samples
1	절도	0.98	192
2	손괴	0.91	80
3	공갈	0.98	62
4	약취유인	0.95	28
5	상해	0.83	82
6	폭행	0.85	107
:	:	:	:
21	성폭속범죄	0.91	5
	전체 평균 성능	0.89	

```
num= 995 , idx= 1903 , f1_score= [0.82348582]
num= 996 , idx= 4232 , f1_score= [0.82348582]
num= 997 , idx= 2059 , f1_score= [0.82348582]
num= 998 , idx= 2673 , f1_score= [0.82348582]
num= 999 , idx= 1968 , f1_score= [0.82348582]
Averaged F1 Score= 0.8234796471348657
completed
```

그림 4. 용의자 후보군 추론 기술 성능 검증 결과

- 전체사건: 5,539 건
- 전체인물: 12,334 명
  - 피의자: 6,084 명, 피해자: 5,729 명, 참고인: 521 명

실험방법은 용의자 후보 추론 결과 내의 상위 분위(percentile) 1.5%내에 용의자가 포함되었는지 여부로 성능을 검증하였다. 임의의 1000 건의 사건을 추출하여 용의자후보군을 추론 하였으며 그림 4 와 같이 전체 사건의 평균 F1 score 는 0.82 점으로 범죄수사에 참고할만한 점수를 제공하였다.

### ACKNOWLEDGMENT

이 논문은 2021 년도 정부(과학기술정보통신부)의 재원으로 정보통신기획평가원의 지원을 받아 수행된 연구임 (No. 2018-0-00440, 위험 상황 초기 인지를 위한 ICT 기반의 범죄 위험도 예측 및 대응 기술 개발).

### 참 고 문 헌

- [1] 장광호, “스마트치안”, 2020. 06.
- [2] 백명선 외 4, “효과적 대응을 위한 기계학습 기반의 범죄 유형 및 범죄 위험스코어 예측 기술 연구”, 2020 한국통신학회 하계 종합학술대회, 2020. 08.
- [3] 백명선 외 2, “인공지능 기반 지능형 범죄 위험 예측 및 대응 기술을 활용한 스마트 치안 기술 연구”, 2020 한국인공지능학회, 2020. 12.
- [4] M. Kim, D.-g. Lee, H. Shin, “Semi-supervised learning for hierarchically structured networks,” Pattern Recognition, vol 95, pp. 191-200, Nov. 2019.