

# 결합손실함수를 이용한 112 신고 데이터의 분류 성능 향상

이은결<sup>1,2</sup>, 박현호<sup>1</sup>, 박영수<sup>1</sup>, 변성원<sup>1</sup>  
<sup>1</sup>한국전자통신연구원, <sup>2</sup>광주과학기술원

doryeon514@gm.gist.ac.kr, hyunhopark@etri.re.kr, yspark@etri.re.kr, wbyon@etri.re.kr

## Classification Performance Improvement of 112 Report Data By Using Combined Loss Function

Eungyeol Lee<sup>1,2</sup>, Hyunho Park<sup>1</sup>, Young-Su Park<sup>1</sup>, Sungwon Byon<sup>1</sup>

<sup>1</sup>Electronics and Telecommunications Research Institute,  
<sup>2</sup>Gwangju Institute of Science and Technology

### 요약

경찰이 범죄현장에 적절히 대응하기 위해, 112 신고데이터의 분석을 통한 위험도와 긴급도의 정확한 분류는 중요하다. 본 논문은 Cross Entropy 와 Mean Square Error 함수가 결합된 결합손실함수를 사용하는 112 신고데이터의 KcELECTRA 기반 위험도와 긴급도의 분류 기법을 소개하고 성능향상을 검증한다. 결합손실함수를 이용한 112 신고데이터의 위험도와 긴급도 분류를 통해, 경찰의 범죄현장 대응력을 향상시키고 그에 따른 범죄현장 피해를 최소화할 수 있을 것이다.

### I. 서론

범죄 현장에 적절하게 대응하기 위해서는 범죄현장의 상태에 따른 경찰의 적절한 대응능력과 경찰인력의 적절한 배치가 중요하다[1]. 경찰에게 범죄현장에 대한 적절한 대응방안을 제공하고 범죄현장에 적절한 경찰인력의 배치를 위해, 112 신고 분석을 통한 범죄현장의 위험도와 긴급도를 정확히 파악하는 것이 필요하다. 112 신고 데이터를 KoBERT 를 이용하여 학습하여, 범죄현장의 위험도를 분류하는 시스템이 연구되었다[2].

본 논문에서는 결합손실함수를 사용하여 기존의 112 신고 데이터 분석 시스템에 비해 위험도 분류 성능을 향상시키고, 긴급도를 제공할 수 있는 방안을 제공한다. 본 논문에서 분류 성능을 향상시키기 위해 KoBERT 보다 우수한 KcELECTRA[3]에 Dense 층을 결합하여 사용하였다. 그리고, 본 논문에서 위험도와 긴급도와 같이 단계적인 값 형태의 분류를 위해 Cross Entropy 손실함수와 회귀분석에 특화된 Mean Square Error 손실함수를 결합하여 사용하는 방안을 제안한다. 본 논문에서 KcELECTRA 와 결합손실함수를 이용한 위험도와 긴급도 분류성능을 측정하여, 결합손실함수의 분류성능 향상을 검증한다. 향상된 위험도와 긴급도 분류성능을 통해, 경찰의 적절한 범죄현장 대응을 지원할 것이다.

### II. 결합손실함수를 이용한 위험도 및 긴급도 분류 모델

그림 1 은 본 논문에서 제안하는 KcELECTRA 와 결합손실함수를 이용한 112 신고 데이터 분류모델을 보여준다. 112 신고 데이터의 위험도 및 긴급도 분류를 위한 사전학습 모델로 KcELECTRA 를 사용하였다. 또한, 본 논문의 손실함수로서 Cross Entropy 손실함수와 Mean Square Error 손실함수를 결합한 손실함수를 사용하였다. 여기에 성능 향상을 위한 Dense

층 1 개와 ReLU 활성화 함수, 그리고 분류를 위한 Dense 층 1 개와 분류를 위한 Softmax 함수를 쌓아 더 효과적으로 Label 을 분류할 수 있도록 하였다.

기존 연구 [4]에 따르면, BERT 기반 사전학습모델을 fine tuning 하면 효과적인 텍스트 분류 모델을 만들 수 있다. 이에 따라 본 논문에서는 BERT 기반 사전학습 모델로 KcELECTRA 를 사용하였다. KcELECTRA 는 ELECTRA 초기 모델에서 한국어 온라인 뉴스 댓글을 수집하여 추가로 학습한 모델으로, 2022 년 기준 NER, NSMC, Question Pair 등 많은 Task 에서 가장 좋은 성능을 보인 모델이다[3].

본 논문의 모델이 수행하는 태스크(task)는 112 신고 텍스트 데이터를 통하여 단계적인 값을 가진 위험도와

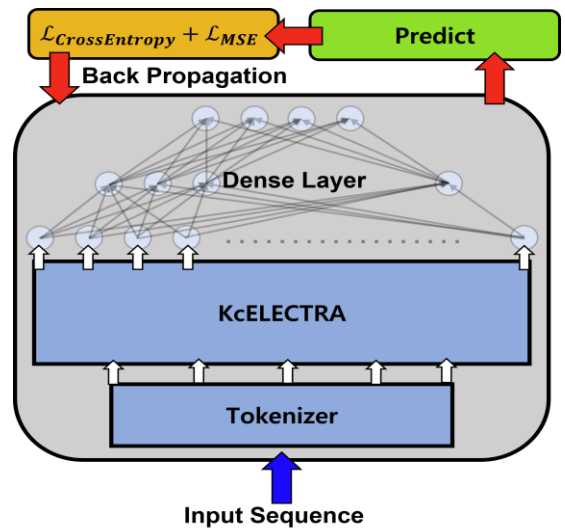


그림 1. KcELECTRA 와 결합손실함수를 이용한 112 신고 데이터 분류모델

긴급도를 분류한다. 보통 다중 분류 모델에서는 Cross Entropy 손실 함수를 사용한다. 레이블이 단계적으로 분포되어 있는 특성상, 단순히 Cross Entropy 손실 함수를 사용하는 것보다 Cross Entropy 와 Mean Square Error 를 일대일 선형 결합한 손실 함수를 사용하는 것이 분류성능 향상 측면에서 더 효과적일 것이라 판단하여, 본 논문에서 Cross Entropy 손실함수와 Mean Square Error 손실함수를 결합한 손실함수를 사용하였다.

### III. 데이터 및 모델 평가

112 신고 데이터의 개수는 총 16372 개이며 데이터마다 위험도와 긴급도가 레이블링 되어있다. 위험도는 피의자 행위 정도에 따라 D1 에서 D5 까지 다섯 단계로 레이블링 되었고, 긴급도는 얼마나 빨리 경찰 지원이 필요한 지에 따라 C0 에서 C3 까지 네 단계로 나누었다. C4 의 긴급도에 대한 데이터는 확보하지 못해 C4 긴급도의 데이터는 학습에 사용하지 못하였다. 데이터 개수 중 90%인 14734 개는 훈련 데이터로 사용되었으며 나머지 10%인 1638 개는 시험 데이터로 사용되었다. 별도의 검증 데이터는 사용되지 않았다.

소개한 데이터를 통해 본 논문에서 제안한 모델 및 다양한 비교군 모델을 학습했다. 학습은 총 5 번씩 진행했다. 최대 18 epochs 까지 학습했으며 그 중 시험 데이터의 예측 정확도가 가장 높은 epoch 의 결과값의 평균을 계산했다.

아래 표 1 은 다양한 모델들의 112 신고 텍스트 데이터 분류 결과이다. KcELECTRA\_v1 는 KcELECTRA 와 레이블로 분류하는 하나의 Dense 층만 있는 가장 기본적인 형태의 베이스 모델이다. KcELECTRA\_v2 모델은 KcELECTRA\_v1 모델에 결합 손실 함수를 적용한 모델이다. KcELECTRA\_v3 모델은 KcELECTRA\_v1 모델에 Dense 층을 하나 추가한 응용 모델이다. KcELECTRA\_v4 모델은 본 논문에서 최종적으로 제시한 KcELECTRA\_v3 모델에 결합 손실 함수를 적용한 모델이다.

표 1. 모델에 따른 위험도, 긴급도 예측 정확도 비교

모델	위험도(acc)	긴급도(acc)	평균(acc)
KoBERT[2]	0.828	0.722	0.775
KcELECTRA_v1	0.842	0.785	0.813
KcELECTRA_v2	0.889	<b>0.859</b>	0.874
KcELECTRA_v3	0.852	0.816	0.834
<b>KcELECTRA_v4</b>	<b>0.893</b>	<b>0.859</b>	<b>0.876</b>

공통된 모델 학습 환경은 다음과 같다. Optimizer 는 AdamW 함수를 사용했다. Dropout 은 0.3, batch 의 크기는 32, learning rate 는  $5e-6$  이다. KcELECTRA\_v2 모델과 KcELECTRA\_v4 모델에서 중간의 Dense 층의 선형 계수는 256 이다.

모든 KcELECTRA 기반 모델이 기존보다 좋은 결과를 얻을 수 있었다. 특히 KcELECTRA\_v4 모델은 112 신고 텍스트 데이터의 위험도, 긴급도를 분류하는

Task 에서 가장 좋은 성능을 보였다. 또한 결합 손실 함수를 사용한 모델은 그렇지 않은 모델보다 상당한 성능 향상을 보였다.

### IV. 결론

본 논문에서는 112 신고 텍스트 데이터 분류 모델의 성능 향상 기법을 세 가지 제시했다. 첫째, Cross Entropy 손실 함수를 사용하는 대신 Cross Entropy 함수와 Mean Square Error 함수를 선형결합한 결합 손실 함수를 사용한다. 둘째, 기존에 제안한 방법에서 사용된 사전 학습 모델인 KoBERT 대신 KcELECTRA 를 사용한다. 셋째, 기존 분류를 위한 1 개의 Dense 층 사이에 추가적인 Dense 층을 삽입한다.

위 세 기법을 통하여 112 신고 텍스트 데이터 분류 모델의 성능을 향상시킬 수 있었다. 또한 결합 손실 함수는 단계적 특성의 다중 라벨 분류 태스크에서 기존의 연구 [2]에서 KoBERT 모델을 사용했을 때 평균 정확도 0.775 보다 약 0.1 높은 정확도의 평균 정확도 0.876 의 향상된 성능을 보였다.

단계적 레이블 분류 문제에서 결합손실함수의 사용이 성능을 향상시킬 수 있는지 여부에 대해 더 많은 사례를 통해 검증하는 과정이 필요하다. 향후 추가적인 연구를 통해, 더욱 다양한 종류의 데이터에 대한 학습에서의 결합손실함수의 사용에 따른 성능향상 여부를 검증할 것이다. 또한 수식적인 접근 방법을 통해 결합손실함수의 단계적 레이블 분류 문제의 분류 성능 향상을 증명할 예정이다.

### ACKNOWLEDGMENT

이 논문은 2024 년도 정부(경찰청)의 재원으로 지원받아 수행된 연구결과임 [내역사업명: 112 긴급출동 의사결정 지원 시스템/ 연구개발과제번호: PR08-03-000-21]

### 참 고 문 헌

- [1] 정웅, "지역경찰 수사사무의 범위 및 적정 인력 연구", 치안정책연구소 책임연구보고서 2020-10, 2021.
- [2] 박현호, 권은정, 변성원, 이민정, 박영수, 정의석, "112 신고 데이터 기반 위험도 분석 시스템 구현", 한국통신학회 2023 인공지능 학술대회 논문집, pp. 51-53.
- [3] KcELECTRA: Korean comments ELECTRA, <https://github.com/Beomi/KcELECTRA>
- [4] Sun, C., Qiu, X., Xu, Y., and Huang, X., "How to fine-tune bert for text classification?" Chinese Computational Linguistics: 18th China National Conference, CCL 2019, Kunming, China, October 18- 20, 2019, Proceedings 18. Springer International Publishing, 2019.