

교통사고상황 텍스트 정보를 활용한 자연어처리 기반의 도로교통법 위반에 대한 판결예측 시스템

A Natural Language Processing-Based Judgment Prediction System for Road Traffic Act Violations Using Text Information on Traffic Accident Situations

민현식¹, 윤준영², 노병준^{*}

minun001@sch.ac.kr, ypjun100@sch.ac.kr, *powernoh@sch.ac.kr

목차

1. 서론
2. 자연어처리 기술을 활용한 판결 예측 시스템
 - 2.1 데이터 수집 및 전처리
 - 2.2 텍스트 임베딩
 - 2.3 판결 예측 모델
3. 실험 및 실험결과
 - 3.1 성능 평가 척도
 - 3.2 실험 결과
4. 결론
5. 참고 문헌

1. 서론

1. 서론

필요성

사법 취약계층은 전통적으로 법적 지식의 부족, 법률 접근성의 낮음, 경제적 어려움을 겪는 사람들로 정의되어 왔다. 이들은 법률의 복잡성과 지속적인 변화로 인하여 법의 보호를 충분히 받지 못하는 상황에 처해 있으며, 법적 지원 및 교육의 필요성이 높다[1].

사법 취약계층이 발생하는 주요 원인

1. 법률의 복잡성
2. 법률의 지속적인 변화
3. 법률 접근성의 부족.

개념 확장

확장된 사법 취약계층 정의

1. 전통적 취약계층 포함
2. 기술적 변화에 대응하는 데 어려움을 겪는 사람들 포함
3. 도로교통법의 변화에 따른 위법 행위자 포함

[1] 황승흠, "한국사회의 민사 법률구조의 이념과 현실", 법학논총, Vol. 21, No. 2, pp. 247-280, Feb. 2009.

[2] 김용훈, "제4차 산업혁명의 도래와 미래도로환경변화에 따른 도로교통법 진단 및 입법전략 방향", 법학연구소, Vol. 30, No. 2, pp. 163-190, Apr. 2022.

1. 서론

관련연구

기계 학습 모델에 기반한 형량 예측 시스템 연구

- 살인범죄와 같은 반사회적 범죄는 여론 및 사회적 요인에 영향을 받는 경우가 있는데 이런 사회적 요인들을 분석하여 사건에 적합한 형량이 무엇인지 예측하는 연구가 활발히 진행되고 있다.

사건 현장이나 상황을 고려하지 않는 한계

- 사회적 요인에 대한 영향을 고려한 형량 판결 예측 시스템은 존재하지만 **사건 현장이나 상황을 고려한 판결 예측 시스템은 미비한 상황이다.**

유럽 인권 협약 위반 여부를 예측하는 머신러닝 모델 연구

- 유럽 인권 재판소의 판결문을 분석하여 평균 75%의 정확도로 결과를 예측하였다.
- Support Vector Machine (SVM) Linear Classifier 사용

머신러닝의 한계

- 복잡한 법적 상황을 포착하기 위해 심층 신경망 (Deep Neural Networks), 자연어 처리 (Natural Language Processing) 기술이 요구되는 상황이다.**

아랍어 법률 문서를 분석하여 법원 판결을 예측 연구

- 문서에서 무작위로 문장을 선택하여 위치를 바꾸는 방법과 GloVe 와 FastText 모델의 어휘를 활용하여 원문의 임의 단어를 유사한 단어로 대체하는 방법으로 데이터를 증강
- CNN, BiLSTM, LSTM, GRU, BiGRU 등 다양한 딥러닝 모델을 활용하여 실험을 진행하였다.

임베딩 모델의 선택의 중요성

- 임베딩 선택이 딥 러닝 모델의 성능에 중요한 영향을 미칠 수 있음을 시사한다**

[3] 변재욱, 김한솔, 박미랑, 신종원, "기계 학습 모델에 기반한 형량 예측 시스템 연구", 한국범죄학, Vol. 12, No. 1, pp.3-17, May. 2018.

[4] Medvedeva, M., Vols, M. and Wieling, M., "Using Machine Learning to Predict Decisions of the European Court of Human Rights.", Artificial Intelligence and Law, Vol. 28, No. 2, pp. 237-66, June. 2020.

[5] Zahir, J., "Prediction of court decision from Arabic documents using deep learning", Expert Systems, Vol. 40, No. 6, pp. 1-16, July. 2023.

1. 서론

대용량 텍스트 자원을 활용한 한국어 형태소 임베딩의 모델별 성능 비교 분석

- 한국어 텍스트 데이터에 가장 적합한 모델을 찾기 위해 Word2Vec, GloVe, FastText를 비교 분석
- FastText는 임베딩 차원과 윈도우 크기에 따라 일관된 성능을 보여, 임베딩 차원이 200과 300일 때, 윈도우 크기가 10일 때 가장 좋은 결과를 나타냈다.

한국어 형태소 임베딩에
효과적인 모델

- 한국어 형태소 임베딩에 FastText 모델이 효과적이다.

본 연구의 목표

확장된 사법취약계층을 위한 FastText를 활용한 딥러닝 판결 예측 시스템을 제공한다.

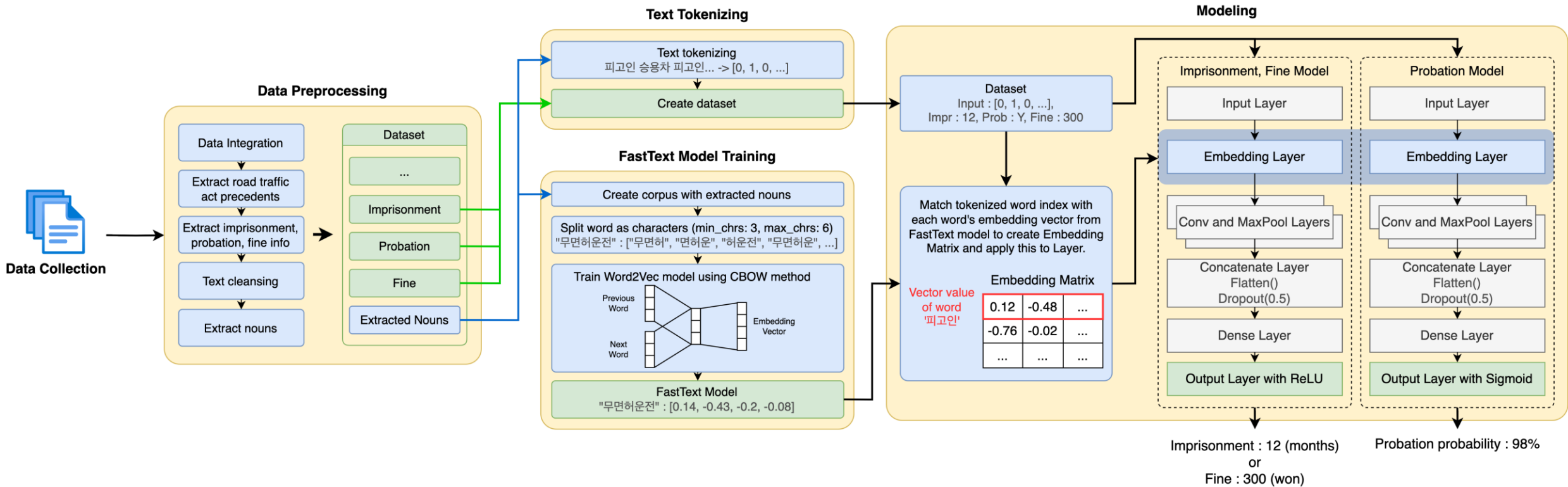
판결 예측 시스템은 사건 정보를 입력 시, 이를 딥러닝 모델을 통해 분석하고 최종적으로 예측된 형량, 집행유예 여부 및 벌금에 대한 정보를 제공한다.

판결 예측 시스템의 기대 부분

- 판결 예측 시스템은 확장된 사법 취약계층이 도로교통법을 이해하고 적용하는 데 도움을 줄 것으로 기대된다.
- 판결 예측 시스템을 통해 사건이 법적으로 어떻게 평가될지 미리 예측할 수 있게 된다.
- 이는 도로교통법 위반의 심각성을 일깨우는 계기가 되며, 도로교통법의 최신 동향을 파악하고, 법률적 지식을 습득하는 데 도움을 줄 수 있을 것이다.

2. 자연어처리 기술을 활용한 판결 예측 시스템

2. 자연어처리 기술을 활용한 판결 예측 시스템

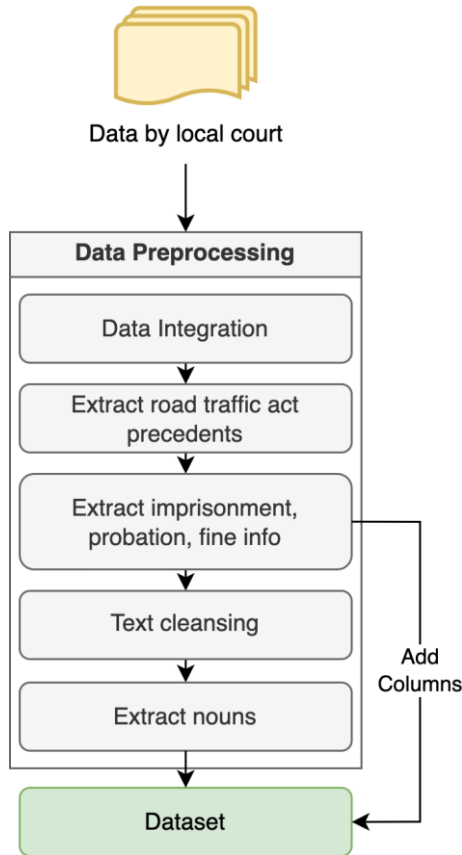


본 논문에서 제안하는 자연어처리 기술을 활용한 판결 예측 시스템 전체 구조도

1. 판례 데이터를 수집
2. 수집된 데이터를 가지고 데이터 전처리 진행
3. Keras에서 제공하는 토크나이저를 통해 추출된 명사들을 토큰화
4. 추출된 명사 단어를 글자단위로 분리시키고 이를 CBOW 방법을 적용한 Word2Vec 모델로 학습시켜 임베딩을 진행
5. 토큰 데이터와 FastText에서 추출된 각 단어들의 임베딩 벡터들을 일치시켜 임베딩 행렬을 생성하고 이를 예측 모델 내의 임베딩 레이어에 적용
6. 각 모델에는 각 단어에 대한 토큰값들이 들어가며 서로 다른 3개의 모델을 구성하여 최종 출력값으로 각각 징역, 집행유예 확률, 벌금 예측

2. 자연어처리 기술을 활용한 판결 예측 시스템

2.1 데이터 수집 및 전처리



- 데이터 수집 단계에서는 판결문 제공 서비스인 LBox[7]에서 도로교통법을 위반한 판결문 총 37,407개를 지방법원 12곳에서 수집하였다.
- 각 지방법원 별로 나뉘진 판결문들은 하나로 통합한 뒤 다음의 전처리 과정을 진행한다.
 1. 먼저, 통합된 판결문 중 도로교통법 위반이 아닌 판결문들을 삭제한다.
 2. 그리고 판결문 내의 최종 판결 부분에서 징역 집행유예 여부, 벌금에 대한 수치형 데이터를 추출한 뒤, 각각 징역, 집행유예 여부, 벌금 컬럼에 저장한다.
 3. 마지막으로 판결 요지 텍스트를 클렌징 처리를 한 후 명사만 추출하여 저장한다.
- 최종적으로 전처리 후 데이터 컬럼 구성은 [Table]과 같다.

[7] LBox [Online]. Available: <https://lbox.kr/case/%EB%8C%80%EB%B2%95%EC%9B%90/2022%EB%8F%847443> (downloaded 2023, Dec. 21)

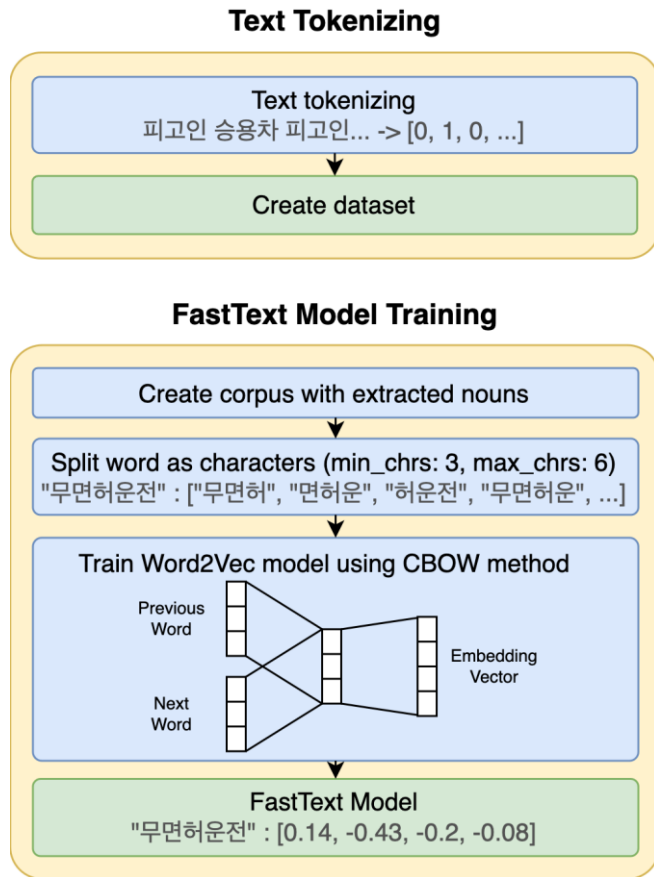
2. 자연어처리 기술을 활용한 판결 예측 시스템

2.1 데이터 수집 및 전처리

Variable name	Description	Example
Court Name	Local court	서울동부지방법원
Case Number	Case number	2023고단1234
Case Result	Simplified final judgement decision	징역 10월.집행유예 2년
Violation Groups	Violation law groups	['교통사고처리특례법위반', '도로교통법위반']
Violations	Name of violation law	['교통사고처리특례법위반(치상)', '도로교통법위반(무면허운전)', '도로교통법위반(음주운전)']
Judgement Date	Judgement date	2023. 7. 20.
Judgement Decision	Final judgement decision	피고인을 징역 10월에 처한다. 다만, 이 판결 확정일로부터 2년간 위 형의 집행을 유예한다.
Judgement Reason	Reason of judgement decision	피고인은 (차량번호 1 생략) G80 승용차의 운전업무에 종사하는 사람이다. 피고인은 2022. 12. 13. 01:15경...
Imprisonment	Imprisonment	10
Probation	Probation	2
Fine	Fine	0
Extracted Nouns	Nouns extracted from Judgement Reason	피고인 승용차 피고인 운전업무 운전면허 술...

2. 자연어처리 기술을 활용한 판결 예측 시스템

2.2 텍스트 임베딩

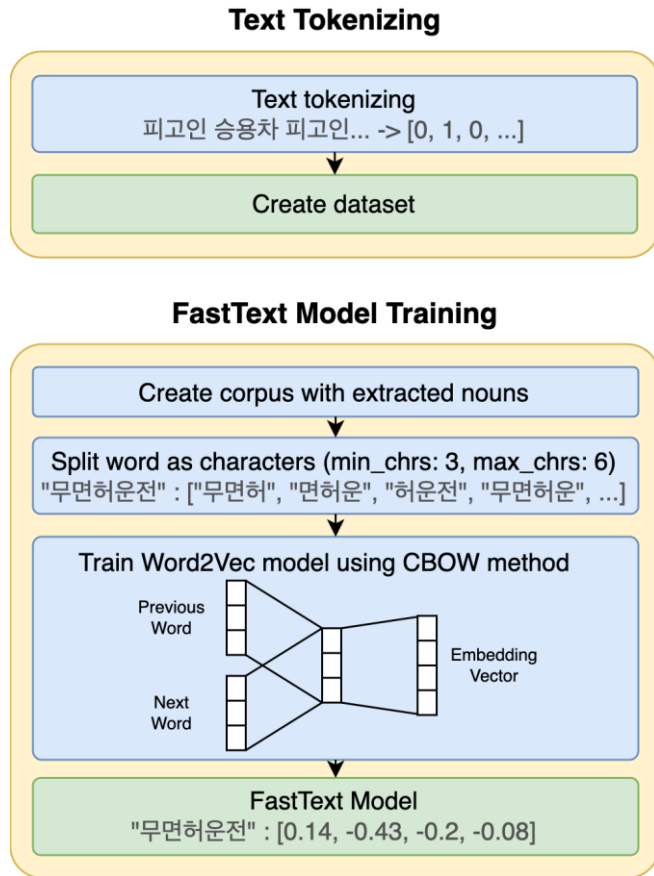


- FastText는 Word2Vec를 기반하지만, 문자 단위의 N-gram을 사용하여 단어를 표현함으로써 사전에 정의되지 않은 단어를 처리하는데 용이하다.
- 특히 FastText를 사용하면 단어 뿐만 아니라 단어 내의 문자 레벨 정보까지 표현되므로 단어 수가 부족한 환경에서 큰 효과를 낼 수 있다[8].
- 본 연구에서는 최소 3개 이상의 글자들로 단어를 분할한 뒤 이를 CBOW 방법을 적용한 Word2Vec 모델을 통해 학습시킨다.

[8] Umer, M., et al. "Impact of convolutional neural network and FastText embedding on text classification", Multimedia Tools and Applications: An International Journal, Vol. 82, No. 4, pp. 5569–5585, 2023.

2. 자연어처리 기술을 활용한 판결 예측 시스템

2.2 텍스트 임베딩



텍스트 임베딩 과정에서는 모델에 실질적으로 입력될 데이터를 가공하는 텍스트 토크나이징 과정과 추출된 명사를 벡터 형태로 변화시키기 위한 FastText 모델 학습 과정이 포함된다.

Text Tokenizing 과정

- Text Tokenizing 과정에서는 추출된 명사를 Keras에서 제공하는 토크나이저를 활용해 토큰화 시킨 후 이 데이터와 정답 데이터인 징역, 집행유예 여부, 벌금을 합쳐 데이터셋을 구성한다.

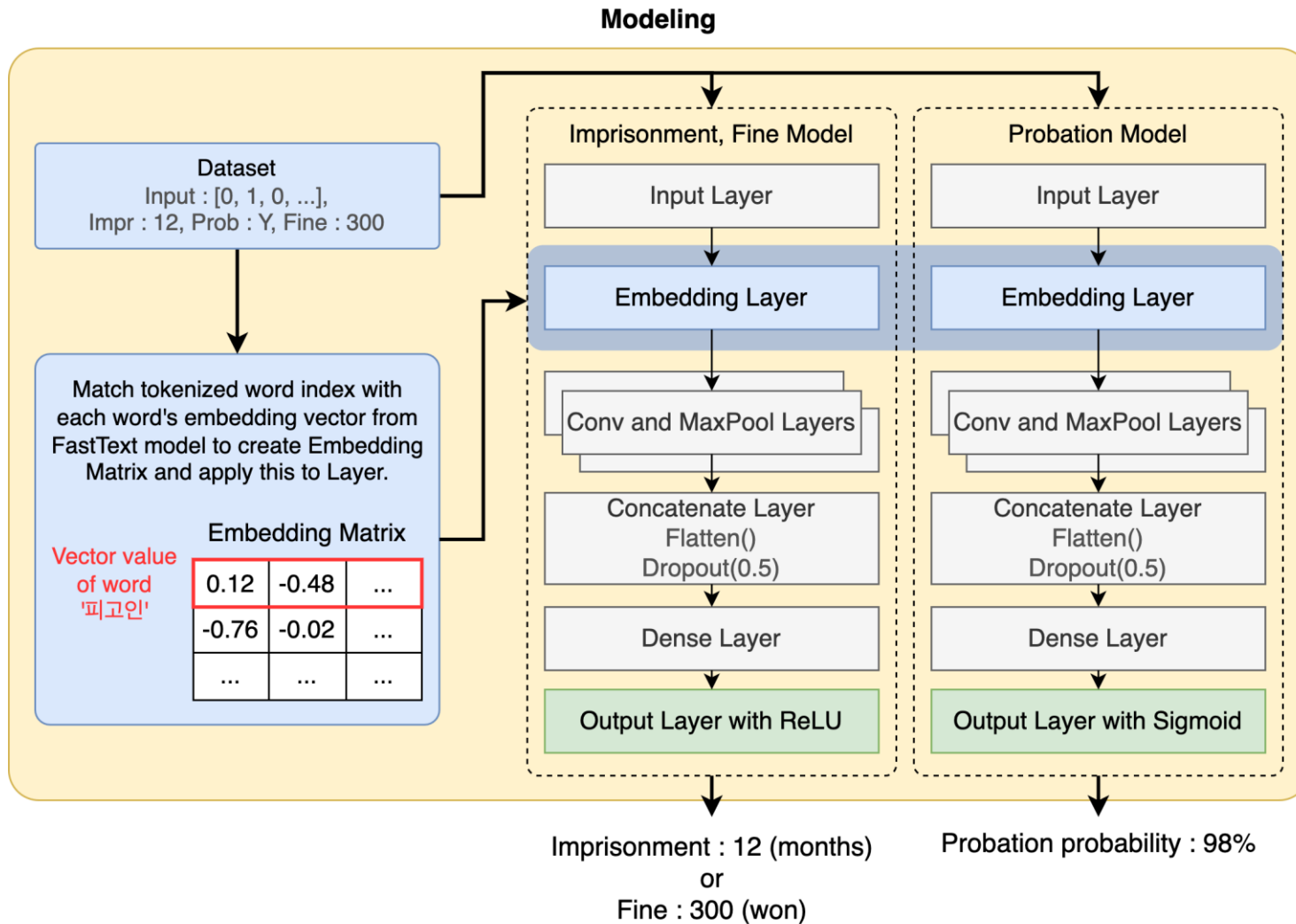
FastText 모델 학습 과정

- FastText 모델 학습 과정에서는 추출된 명사를 그대로 이용해 코퍼스를 구성한 뒤, 각 단어들을 최소 3개 이상의 문자를 가진 문자열로 분할한다.
- 분할된 문자들을 이용하여 CBOW 방식을 적용한 Word2Vec 모델을 학습시키는데, 여기서 CBOW는 현재 단어 주위의 단어들을 토대로 현재 단어에 대한 벡터값을 생성하는 방식이다.
- 최종적으로 모델 학습 이후, 하나의 단어에 대한 임베딩 벡터가 생성되며 이 벡터가 실질적으로 모델에 입력된다.

[8] Umer, M., et al. "Impact of convolutional neural network and FastText embedding on text classification", Multimedia Tools and Applications: An International Journal, Vol. 82, No. 4, pp. 5569–5585, 2023.

2. 자연어처리 기술을 활용한 판결 예측 시스템

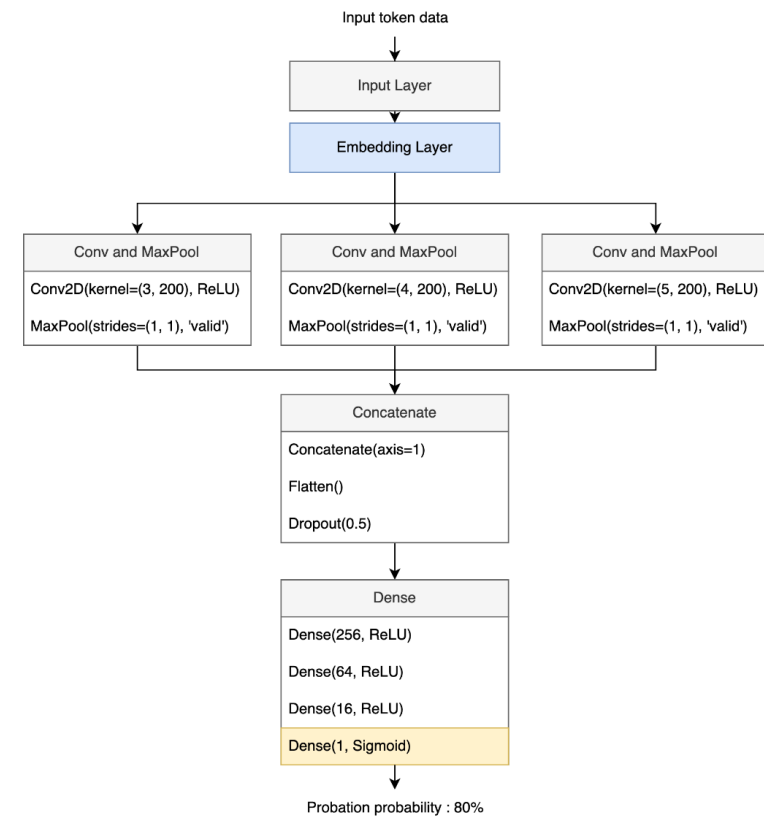
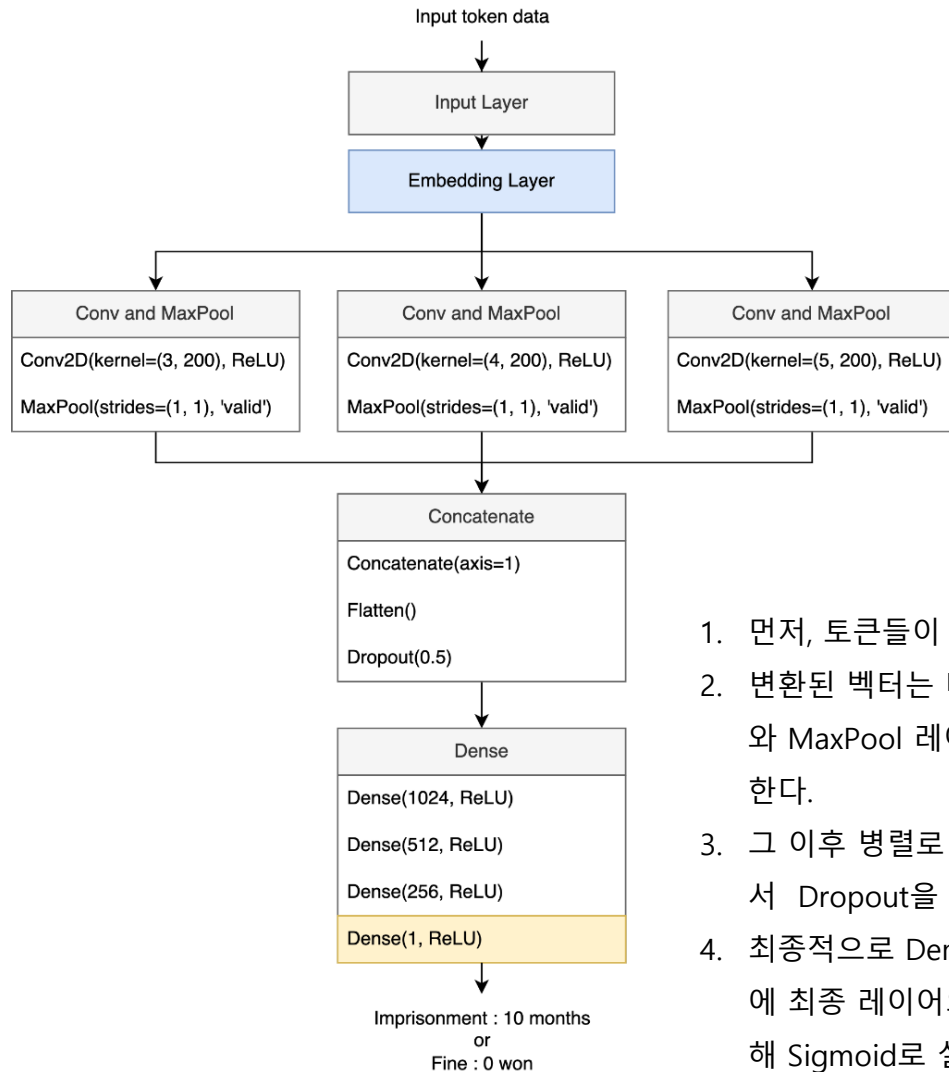
2.3 판결 예측 모델



- 판결 예측 모델을 구성하기 앞서 임베딩 과정에서 나온 각 명사 단어에 대한 토큰들과 FastText 모델 내의 각 단어에 대한 임베딩 벡터를 매칭시켜 임베딩 행렬을 생성한다. 이 행렬은 모델 내의 임베딩 레이어에 적용되며 모델의 입력으로 토큰들이 들어왔을 때 해당 레이어에서 각 토큰을 해당 단어의 임베딩 벡터로 변환해 주는 역할을 수행한다.
- 판결 예측 모델은 CNN 모델을 기반해 판결문 내 영향력있는 단어들에 대한 특징을 추출하여 징역, 집행유예, 벌금을 예측하는 3개의 모델로 구성한다. 각 모델의 구조는 [Figure]와 같다.

2. 자연어처리 기술을 활용한 판결 예측 시스템

2.3 판결 예측 모델



1. 먼저, 토큰들이 입력되면 임베딩 레이어에서 각 토큰을 임베딩 벡터로 변환시킨다.
2. 변환된 벡터는 다양한 특징들을 추출하기 위해 3개의 서로 다른 하이퍼파라미터를 가진 Convolutional 레이어와 MaxPool 레이어에 병렬로 입력된다. 이러한 병렬 과정은 하나의 데이터에 대해 다양한 특징을 추출할 수 있도록 한다.
3. 그 이후 병렬로 연산된 값들은 Concatenate 레이어를 통해 하나로 통합된 후 1차원으로 변환되며, 이 과정에서 Dropout을 적용시켜 모델의 복잡성을 해소시킨다.
4. 최종적으로 Dense 레이어들을 차례대로 전파되면서 결과값이 도출된다. 이때 징역, 벌금 모델은 회귀 모델이기 때문에 최종 레이어의 활성화 함수를 ReLU로 설정하였고, 집행유예 모델은 집행유예 여부에 대한 확률을 출력하기 위해 Sigmoid로 설정하였다.

3. 실험 및 실험 결과

3. 실험 및 실험 결과

3.1 성능 평가지표

$$\left\{ \begin{array}{l} \text{징역} : 144 \text{ (12년)} \times 5\% = 7.2 \text{ (개월)} \quad (1) \\ \text{벌금} : 3000 \text{ (만원)} \times 5\% = 150 \text{ (만원)} \quad (2) \end{array} \right.$$

Evaluation Index	Tolerance
Imprisonment Error Rate	7.2 months
Probation Error Probability	5%
Fine Error Rate	1,500,000 won

징역 오차율 지표 (Imprisonment Error Rate)

- 징역 오차율 지표 (Imprisonment Error Rate Index)(1)는 교통사고 치사 후 유기 도주 사례에 적용되는 최고 가중형 징역 기간인 12년(144개월)[9]의 5%에 해당하는 7.2개월을 오차 범위로 설정한다. 이 지표는 징역예측 모델이 징역 기간을 얼마나 정확하게 예측하는지를 평가하는 역할을 한다. 징역 예측 모델의 예측이 오차 범위 안에 들어오면, 모델의 예측이 옳은 것으로 간주한다.

집행유예 오차확률 지표 (Probation Error Probability)

- 집행유예 오차확률 지표 (Probation Error Probability Index)는 집행유예 모델의 출력 확률값과 실제 집행유예의 부여 여부를 0과 1로 표현한 값에 대한 허용 오차확률을 나타낸다. 예측 확률과 실제 확률에 대한 오차가 5% 이내이면 정답이라 간주한다.

벌금 오차율 지표 (Fine Error Rate)

- 벌금 오차율 지표 (Fine Error Rate Index)(2)는 위험운전 교통사고, 어린이 교통사고, 교통사고 후 도주 사례에 부과될 수 있는 최대 벌금 3000만원[9]의 5%, 즉 150만원을 오차 범위로 정한다. 이 지표는 벌금 예측 모델이 벌금 금액을 얼마나 정확하게 예측하는지를 평가하는 데 사용된다. 예측된 벌금 금액이 5% 오차 범위 이내이면 정답이라 간주한다.

[9] 양형위원회 교통범죄 양형기준 [Online]. Available:

https://sc.scourt.go.kr/sc/krsc/criterion/criterion_35/traffic_change_01.jsp (downloaded 2023, Dec. 21)

3. 실험 및 실험 결과

3.2 실험 결과

Model	Accuracy
Imprisonment	96.13%
Probation	94.85%
Fine	95.49%

- 판결 예측에 사용된 모델은 징역, 집행유예, 벌금 예측으로 총 3개이며, 4.1에서 제안한 성능 평가 지표를 기반으로 모델에 대한 정확도를 계산한다. 각 모델에 대한 정확도는 징역 예측 모델은 96%, 집행유예 예측 모델은 95%, 벌금 예측 모델은 95%를 보였다.
- 이 정확도를 통해 징역 예측 모델의 경우 96%의 확률로 예측값과 실제값 간의 오차가 7.2개월 이내임을 알 수 있고, 집행유예 예측 모델의 경우 94%의 확률로 오차확률이 5% 이하, 벌금 예측 모델의 경우 94%의 확률로 오차가 150만원 이내였다는 것을 알 수 있다.
- 이러한 높은 정확도는 3개의 모델이 판결을 예측하는 데 있어 높은 신뢰성을 가지며, 실제 법률적 상황에서 유용하게 적용될 수 있음을 나타낸다.

4. 결론

4. 결론

- 본 연구에서 개발된 시스템은 자연어처리(Natural Language Processing) 기술을 활용한 CNN 모델을 통해 징역, 집행유예 확률, 벌금 예측을 수행하고 있으며, 모델은 각각 96.13%, 94.85%, 95.49%의 높은 정확도를 달성했다. 이러한 성능은 법률적 상황에서의 유용한 적용 가능성을 보여주며, 인공지능 기술이 법률 분야에서 중요한 역할을 할 수 있음을 시사한다.
- 본 연구에서는 징역 오차율 지표(Imprisonment Error Rate), 집행유예 오차확률 지표(Probation Error Probability), 벌금 오차율 지표(Fine Error Rate) 등 새로운 성능 평가 지표들을 도입했다. 이 지표들은 모델의 예측 정확도를 객관적이고 실제 법률적 상황에 맞게 평가하는 데 중요한 역할을 한다.

Discussions

- 본 연구에서는 판결 예측을 위한 모델로 CNN을 사용하였지만 최근 각광받고 있는 BERT나 GPT와 같은 언어 모델을 사용하지 않은 점과 더 방대한 데이터를 수집하여 모델을 훈련시키지 않은 점이 본 논문의 한계점이다. 이후에는 다양한 모델을 사용하여 비교함과 동시에 데이터셋 증진을 통한 판결 예측 모델의 정확도 향상을 통해 법률 분야에서 인공지능 기술의 응용 범위를 넓히고, 사법적 불평등을 해소하는 방안에 대한 연구를 계속해서 진행할 것이다.

5. 참고 문헌

5. 참고 문헌

- [1] 황승흠, "한국사회의 민사 법률구조의 이념과 현실", 법학논총, Vol. 21, No. 2, pp. 247-280, Feb. 2009.
- [2] 김용훈, "제4차 산업혁명의 도래와 미래도로환경변화에 따른 도로교통법 진단 및 입법전략 방향", 법학연구소, Vol. 30, No. 2, pp. 163-190, Apr. 2022.
- [3] 변재욱, 김한솔, 박미랑, 신종원, "기계 학습 모델에 기반한 형량 예측 시스템 연구", 한국범죄학, Vol. 12, No. 1, pp.3-17, May. 2018.
- [4] Medvedeva, M., Vols, M. and Wieling, M., "Using Machine Learning to Predict Decisions of the European Court of Human Rights.", Artificial Intelligence and Law, Vol. 28, No. 2, pp. 237-66, June. 2020.
- [5] Zahir, J., "Prediction of court decision from Arabic documents using deep learning", Expert Systems, Vol. 40, No. 6, pp. 1-16, July. 2023.
- [6] 이다빈, 최성필, "대용량 텍스트 자원을 활용한 한국어 형태소 임베딩의 모델별 성능 비교 분석", 정보과학회논문지, Vol. 46, No. 5, pp. 413-418, May. 2019.
- [7] LBox [Online]. Available: <https://lbox.kr/case/%EB%8C%80%EB%B2%95%EC%9B%90/2022%EB%8F%847443> (downloaded 2023, Dec. 21)
- [8] Umer, M., et al. "Impact of convolutional neural network and FastText embedding on text classification", Multimedia Tools and Applications: An International Journal, Vol. 82, No. 4, pp. 5569-5585, 2023.
- [9] 양형위원회 교통범죄 양형기준 [Online]. Available: https://sc.scourt.go.kr/sc/krsc/criterion/criterion_35/traffic_change_01.jsp (downloaded 2023, Dec. 21)
- [10] Jabbar, A., et al. "An Analytical Analysis of Text Stemming Methodologies in Information Retrieval and Natural Language Processing Systems", IEEE Access, Vol. 11, pp. 113681-133702, Jan. 2023.
- [11] Lee, S., et al. "Extracting Fallen Objects on the Road From Accident Reports Using a Natural Language Processing Model-Based Approach", IEEE Access, Vol. 11, pp. 139521-139533, Jan. 2023.
- [12] Graham, S.G., Soltani, H. and Isiaq, O., "Natural language processing for legal document review: categorising deontic modalities in contracts", Artificial Intelligence and Law, pp. 1-22, No. v. 2023.