

한국어 생의학 개체명 인식 성능 비교와 오류 분석

이재홍*

Performance Comparison and Error Analysis of Korean Bio-medical Named Entity Recognition

Jae-Hong Lee*

요 약

딥러닝 분야에서 트랜스포머 아키텍처의 출현은 자연어 처리 연구가 획기적인 발전을 가져왔다. 개체명 인식은 자연어 처리의 한 분야로 정보 검색과 같은 태스크에 중요한 연구 분야이다. 생의학 분야에서도 그 중요성이 강조되나 학습용 한국어 생의학 말뭉치의 부족으로 AI를 활용한 한국어 임상 연구 발전에 제약이 되고 있다.

본 연구에서는 한국어 생의학 개체명 인식을 위해 새로운 생의학 말뭉치를 구축하고 대용량 한국어 말뭉치로 사전 학습된 언어 모델을 선정하여 전이 학습시켰다. F1-score로 선정된 언어 모델의 개체명 인식 성능과 태그별 인식률을 비교하고 오류 분석을 하였다. 인식 성능에서는 KlueRoBERTa가 상대적 좋은 성능을 보였다. 태깅 과정의 오류 분석 결과 Disease의 인식 성능은 우수하나 상대적으로 Body와 Treatment는 낮았다. 이는 문맥에 기반하여 제대로 개체명을 분류하지 못하는 과분할과 미분할로 인한 것으로, 잘못된 태그들을 보완하기 위해서는 보다 정밀한 형태소 분석기와 풍부한 어휘사전 구축이 선행되어야 할 것이다.

ABSTRACT

The advent of transformer architectures in deep learning has been a major breakthrough in natural language processing research. Object name recognition is a branch of natural language processing and is an important research area for tasks such as information retrieval. It is also important in the biomedical field, but the lack of Korean biomedical corpora for training has limited the development of Korean clinical research using AI.

In this study, we built a new biomedical corpus for Korean biomedical entity name recognition and selected language models pre-trained on a large Korean corpus for transfer learning. We compared the name recognition performance of the selected language models by F1-score and the recognition rate by tag, and analyzed the errors. In terms of recognition performance, KlueRoBERTa showed relatively good performance. The error analysis of the tagging process shows that the recognition performance of Disease is excellent, but Body and Treatment are relatively low. This is due to over-segmentation and under-segmentation that fails to properly categorize entity names based on context, and it will be necessary to build a more precise morphological analyzer and a rich lexicon to compensate for the incorrect tagging.

키워드

Named Entity Recognition, BIO-Tagging, Transformers, Transfer Learning, Bio-Medical Corpus
개체명 인식, BIO-태깅, 트랜스포머, 전이 학습, 생의학 말뭉치

* 교신저자 : 전남도립대학교 보건의료과

• 접수일 : 2024. 06. 17

• 수정완료일 : 2024. 07. 15

• 게재확정일 : 2024. 08. 12

• Received : Jun. 17, 2024, Revised : Jul. 15 2024, Accepted : Aug. 12, 2024

• Corresponding Author : Jae-Hong Lee

Dept. of Health & Medical Science, Jeonnam State University

Email : hlee@dorip.ac.kr

I. 서론

지난 2017년 허깅페이스(Huggingface)의 트랜스포머¹⁾가 등장하며 딥러닝(deep learning)의 한 분야인 자연어 처리 분야에 획기적인 발전을 가져왔다. 트랜스포머는 기존 딥러닝 모델들과 달리 대용량 말뭉치(데이터 셋)로 사전 학습시킨 대규모 언어 모델(LLM, Large Language Model)들을 파생시키고 약간의 미세 조정만으로도 감정분석, 기계번역, 질의응답 등 특정 자연어 처리에 응용하기 쉽다는 장점이 있다[1].

본 연구에서는 자연어 처리 분야 중 정보 검색에 중요한 개체명 인식(NER, Named Entity Recognition)[2]을 위해 트랜스포머의 인코더 기반 모델 중에서 한국어 말뭉치로 사전 학습시킨 언어 모델들을 선정하고 한국어 생의학(bio-medical) 말뭉치를 구축하여 이를 전이 학습시킨 후에 선정된 언어 모델들의 인식성능을 비교하고 오류를 분석하여 향후 생의학 NER 인식 향상을 기하고자 한다.

본 논문은 II장에서 NER 개념과 생의학 NER 관련 연구를 살펴보고, III장에서는 한국어 생의학 NER 말뭉치를 구축한 후, 사전 학습된 언어 모델들을 선정한 후 전이 학습시켜 성능을 측정한 후, 각 모델들의 인식을 평가와 오류를 분석하고, IV장에서 결론을 맺는다.

II. NER 개념과 관련 연구

2.1 NER 개념

NER은 Named Entity(이름을 가진 개체)를 인식하는 것을 의미하며, ‘개체명 인식’이라고 한다³⁾. NER은 자연어 처리 태스크들 중 하나로, 텍스트 안의 키(key) 정보를 식별하여 이름, 위치, 조직 등과 같은 세트로 범주화(categorization)한다. NER를 수행하기 전에, 텍스트나 말뭉치를 토큰화, 덩어리로 묶기(chunking), 태깅(tagging)과 같은 전처리를 거치게 된다. 태깅에서는, 문장을 토큰 단위로 나누고 이 토큰들에 대해서 개체명 여부를 판별하여 적절한 태그

를 부여한다. 단일 토큰이 아닌 여러 개의 토큰 결합으로 하나의 개체명이 완성되는 경우, 여러 개의 토큰을 하나의 개체명으로 묶는 태깅 작업(chunking)이 필요하다. 태깅 작업에는 주로 BIO(Begin-Inside-Outside) 시스템⁴⁾이 사용되고 있다. 각 토큰에 부여된 태그 뒤에는 어떤 종류의 개체명인지를 식별하는 라벨이 부가된다.

표 1은 BIO 기반 태깅 작업의 예를 보여 준다. 이들은 향후 하나의 개체명으로 묶는 작업이 필요하다.

표 1. BIO 기반 태깅 작업의 예
Table 1. Example of a BIO-based tagging task

Token		Label
김길수	Gil-Soo Kim	PER-B
는	has been	o
급성	acute	Disease-B
상기도	upper respiratory	Disease-I
감염	infection	Disease-I
으로	-	o
지난	last	DAT-B
8월	August	DAT-I
부터	since	o
약물	medication	Treatment-B
치료	treatment	Treatment-I
를	-	o
받고	-	o
있다	-	o

2.2 관련 연구

생의학 분야 언어 모델은 트랜스포머를 기반으로 한 많은 파생 모델들이 출현했다. BioBERT는 BERT⁵⁾의 가중치를 그대로 사용하고 최소한의 구조 변경만을 하면서 대용량 생의학 말뭉치를 학습시켜 다양한 생의학 말뭉치를 대상으로 한 태스크에서 일반적인 말뭉치에서 학습된 BERT보다 좋은 성능을 보였다[3]. ClinicalBERT는 의무기록을 학습하여 입원 초기부터 퇴원까지 다양한 입원 시점에서 30일 내 재입원을 예측하기 위한 모델로 진단 예측, 사망 위험 추정, 재원기간 평가와 같은 태스크에도 적용할 수 있다고 보고되었다[4]. BioGPT는 GPT-2를 기반으로 15M

1) <https://arXiv.org/abs/1706.03762>

2) <https://www.letr.ai/blog/tech-20210723>

3) <https://www.letr.ai/blog/tech-230224>

4) <https://arxiv.org/abs/cmp-lg/9505040>

5) <https://arxiv.org/abs/1810.04805>

PubMed 요약들을 사전 학습시켰다⁶⁾. 그렇지만 BioBERT는 전문 의학용어를 충분히 반영하지 않고 있으며, ClinicalBERT는 재입원 예측에 초점을 맞추고 있고, BioGPT는 BERT기반 모델들과 달리 요약 생성, 지식 생성과 같은 생성 태스크에 적합하다. 이외에 BERT 기반 모델들의 최대 입력 시퀀스 길이 제한을 확장함으로써 좀 더 장기적인 의존성을 학습시켜 임상 개념의 더 풍부한 문맥을 다뤄 암 병기 예측에 활용된 Clinical-Longformer와 Clinical-BigBird도 있다⁵⁾.

한국어 생의학 개체명 인식은 생의학 말뭉치 구축의 미비로 많은 연구가 보고되지 않았다. [6]에서는 네이버의 한국어 온라인 QA 서비스에서 임상관련 질의응답들로부터 4개 진료과목에 대해 말뭉치를 만든 후, 상병(DZ), 증상(SX), 신체 부위(BP)의 3개 태그를 부여한 후, BiLSTM-CRF(Bidirectional Long Short Term Memory-Conditional Random Fields) 모델과 다국어 버전 BERT 모델로 개체명 인식을 한 결과 BERT가 좋은 성능을 보여 주었다. [7]에서는 네이버 NLP Challenge 2018의 개체명 말뭉치⁷⁾에서 용어(TERM) 관련 TRM 태그를 포함한 데이터들을 제거하고 대신에 ChatGPT로 생성한 생의학 문장을 자동 및 수동으로 태깅하여 기존 13개에 3개(질병: DISEASE, 신체:BODY, 처치:TREATMENT) 태그를 추가한 16개 태그로 이루어진 한국어 바이오-의학 말뭉치(KBMC)를 만들었다. KM-BERT⁸⁾는 임상 연구 자료, 건강 정보 뉴스 기사, 의학 서적의 3종류의 약 6백만 문장과 116백만 단어들로 구성된 한국어 의학 말뭉치를 만들고, KR-BERT⁸⁾ 문자 토큰라이저와 가중치를 활용하여 학습시켰다.

III. 인식 성능 평가와 오류 분석

3.1 실험용 말뭉치

본 연구에서는 생의학 문장을 포함한 말뭉치를 구성하기 위해 네이버 개체명 말뭉치와 한국어 생의학 말뭉치(KBMC)^[7]를 활용하였다. 네이버 개체명 말뭉

치 총 82,393개 문장으로 14종류의 개체명을 사용하고 태깅 시스템은 BIO를 적용하고 있다. KBMC 말뭉치는 총 6,147문장으로 되어 있다, 본 연구에서는 KBMC 말뭉치를 7:3의 비율로 훈련용과 테스트용으로 분할하고, 여기에 네이버 NER 말뭉치(훈련용 81,000문장, 테스트용 9,000문장)를 각각 추가한 다음, TRM 태그를 가진 문장을 훈련용에서 11,099개, 테스트용에서 1,310개를 제거하여 훈련용 74,201개, 테스트용 9,540개의 한국어 임상 말뭉치(KCC : Korean Clinical Corpus)를 새로 구축하였다.

표 2는 네이버 개체명 말뭉치와 KBMC에 사용된 태그들을 보여 준다.

표 2. 네이버 개체명 말뭉치와 KBMC의 태그들
Table 2. Named Entity Tags for NAVER NER corpus and KBMC

No	Category	Corpus	
		NAVER	KBMC
1	person	PER	PER
2	location	LOC	LOC
3	organization	ORG	ORG
4	date	DAT	DAT
5	time	TIM	TIM
6	quantity	NUM	NUM
7	artifact	AFW	AFW
8	event	EVT	EVT
9	animal	ANM	ANM
10	plant	PLT	PLT
11	material	MAT	MAT
12	study field	FLD	FLD
13	civilization	CVL	CVL
14	term	TRM	-
15	disease	-	Disease
16	body	-	Body
17	treatment	-	Treatment

KCC 말뭉치로 전이 학습된 언어 모델들의 태깅 예측을 평가하기 위해 예측용 데이터를 작성하였다. 예측용 데이터는 ChatGPT-4를 활용하여 프롬프트에 질병, 증상, 처치 관련 전문 용어를 포함한 문장들을 생성하도록 지시하여 60개 문장을 생성시켰다.

그림 1은 ChatGPT-4로 생성한 예측용 데이터의 일부이다.

6) <https://arxiv.org/abs/2210.10341>

7) <https://github.com/naver/nlp-challenge>

8) <https://arxiv.org/abs/2008.03979>

당뇨병은 대사 장애의 일종으로 신체가 인슐린을 제대로 사용하지 못하거나 충분한 인슐린을 생산하지 못하는 질환 고혈당성 질환은 심혈관계에 여러 합병증을 일으킬 수 있으며, 그 중 뇌혈관 질환은 뇌졸중으로 이어질 수 있는 심각한 결과는 전염성이 높은 질병으로 주로 폐를 감염시키지만, 때로는 중의 다른 부분도 영향을 미칠 수 있습니다. 폐렴과 인플루엔자는 호흡기 감염의 일종으로, 특히 면역력이 약한 사람들에게 심각한 건강 문제를 일으킬 수 있습니다. 골다공증은 주로 노인에게 발생하는 질환으로 뼈의 밀도와 질이 저하되어 골절 위험이 증가합니다. 관상동맥 및 영골 장애는 아르트로스스로도 알려져 있으며, 이는 관절의 연골이 손상되어 발생하는 진행성 질환입니다. 피부암과 흉선암 피부암 발생하는 영종 반종으로 다양한 원인에 의해 발생할 수 있으며 종종 가려움증과 발진을 동반한 신생물 신장 기능이 정지되는 질환으로 치료는 수술적 제거는 투사 필요한 단계에 이를 수 있습니다. 상부 호흡기 감염은 흔히 감기로 알려져 있으며, 이는 바이러스에 의해 주로 발생하는데, 적절한 관리가 이루어지지 않다면 중증의 약화를 가져와 골절 위험을 증가시키므로, 환습과 비타민 D가 풍부한 식사와 함께 규칙적인 운동이 만약 당신이 당뇨병을 진단받았다면, 정기적인 혈당 검사와 적절한 식사 계획을 통해 질병을 관리할 수 있습니다. 피부암은 중앙형에서 가장 흔한 악성 신생물 중 하나로, 조기 발견과 적절한 치료는 생존율을 크게 향상시킬 수 있으며 고혈당성 질환은 심장 질환 및 뇌졸중과 같은 다른 심각한 건강 문제로 이어질 수 있으므로, 정기적인 의료 검진이 필요는 매우 중요이 심할 수 있으며, 때로는 수술적 치료가 필요할 수도 있습니다. 골다공증은 뼈의 약화를 가져와 골절 위험을 증가시키므로, 칼슘과 비타민 D가 풍부한 식사와 함께 규칙적인 운동이 류마티스염은 주로 어린이와 청소년에게 영향을 미치는 질병으로, 심각한 경우 심장 손상을 초래할 수 있습니다. 폐렴은 특히 고령이나 면역체계가 약한 사람들에게 심각한 합병증을 초래할 수 있는 호흡기 질환입니다. 화상과 부상은 즉각적인 의료 조치를 필요로 하는 손상으로, 제대로 된 치료가 이루어지지 않으면 영구적인 흉터나 패암과 같은 카르시노마는 조기 발견이 중요하기 때문에 정기적인 건강 검진을 받아야 합니다. 육종은 골격근이나 연조직에서 발생하는 악의 일종으로, 다른 종양과는 다르게 진행됩니다. 고혈당성 질환은 심장, 뇌, 신장에 다양한 합병증을 일으킬 수 있으니 관리가 필수적입니다. 뇌혈관 질환은 뇌졸중을 유발할 수 있어 즉각적인 의료 조치가 필요합니다. 심근경색 심장의 일종의 혈관 질환으로, 심장 근육에 산소가 부족하게 됩니다. 호흡기 감염은 특히 공공장소에서 손 씻기 등의 예방 조치가 중요합니다. 당뇨병은 복부 내장이 복박을 돕고 나오는 질환이므로 수술로 교정해야 합니다. 흉선이 있는 피부는 자외선 노출을 피해야 하며 적절한 치료가 필요합니다.

그림 1. ChatGPT-4를 이용하여 생성된 문장들
Fig. 1 Sentences generated using ChatGPT-4

3.2 실험 환경

한국어 생의학 개체명 인식을 위한 실험 환경으로 구글의 코랩(Colab Pro)을 사용하였으며, 런타임 환경으로 GPU(Graphics Processing Unit) 유형 L4를 활용하여 학습 속도를 높였다. 개체명 인식 성능 비교를 위해 허깅페이스에 공개된 한국어 언어 모델⁹⁾ 중 한국어 감정 분류와 개체명 인식에서 좋은 성능을 보였던[9,10] KlueBERT¹⁰⁾, KlueRoBERTa¹¹⁾와 LMkorBERT¹²⁾를 선정하였고, 추가로 음절 단위 양방향 워드피스 토큰라이저를 사용하는 KR-BERT¹³⁾를 포함시켰다. 이 모델들은 각각 Tensorflow나 Pytorch로 작성되어 있지만 본 연구에의 전이 학습에는 Tensorflow 버전으로 변환하여 코딩하여 Pytorch 버전으로 개발된 모델의 일부 하이퍼파라미터들이 반영되지 않을 수 있다.

그림 2는 생의학 NER 성능 비교를 위한 실험 절차를 보여 준다. 1) 대용량 한국어 말뭉치로 사전 학습된 언어 모델들을 선정하고, 2) KCC 말뭉치를 훈련 데이터와 테스트 데이터로 분할한 후, 3) 훈련 데이터로 전이 학습시킨 다음, 4) 테스트 데이터를 입력하여 예측된 태깅 정보를 얻고 이를 정답과 비교하여

9) <https://arxiv.org/abs/2112.03014>
10) <https://huggingface.co/klue/bert-base>
11) <https://huggingface.co/klue/roberta-large>
12) <https://github.com/kiyoungkim1/LMkor>
13) <https://github.com/snunlp/KR-BERT-MEDIUM/blob/main/README.md>

F1-score를 얻는다. 5) ChatGPT-4를 이용하여 생성한 예측용 데이터를 학습된 언어 모델에 입력하여 태깅된 출력으로 모델들의 예측 성능을 평가한다.

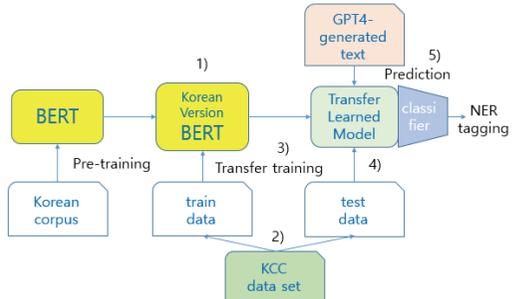


그림 2. BIO-NER 성능 비교를 위한 실험 절차
Fig. 2 Experimental Procedures for Comparing BIO-NER Performance

실험에 사용된 언어 모델들은 tensorflow/keras로 전이 학습시켰으며 선정된 언어 모델들과 하이퍼파라미터들은 표 3과 표 4와 같다.

표 3. 실험에 사용된 언어 모델들
Table 3. Language Models used in the Experiment

Base Model	Derived model (Pre-trained language model)	Checkpoint
BERT	KlueBERT	klue/bert-base
	KlueRoBERTa	klue/roberta-base
	KR-BERT	snunlp/KR-BERT-MEDIUM
	LMkorBERT	kykim/bert-kor-base

표 4. 언어 모델들을 전이 학습하기 위한 하이퍼파라미터들

Table 4. Hyper-parameters for Transfer-Learning Language Models

Batch size	32(training)	64 or 1024(evaluating)
Learning rate	5e-5	
Optimizer	Adam	
Epoch	2~3	
Loss function	BinaryCrossentropy	

예측 데이터는 전이 학습된 언어 모델에 입력되기 전에 한국어 형태소 분석기로 전처리하였다. 형태소 분석기로는 KoNLPy(Korean NLP in Python)¹⁴⁾에서 지원하는 Hannanum, Kkma, Komoran, Okt, Mecab 중 상대적으로 실행시간과 형태소 분석 성능에서 좋은 성능을 보이며¹⁵⁾, 윈도우 환경을 지원하는 Okt¹⁶⁾를 선정하였다.

3.3 성능 평가와 오류 분석

본 연구에서는 생의학 개체명 인식률을 평가하기 위해 정밀도(precision)와 재현율(recall)을 조화 평균한 F1-score를 이용하였다. 3.1절에서 구축한 KCC 말뭉치를 대상으로 한 성능 측정 결과는 표 5와 같다. [10]의 한국어 개체명 인식 성능 비교에서와 같이 여러 한국어 언어 모델 중 KlueBERT와 KlueRoBERTa의 F1-score가 상대적으로 좋은 성능을 보여 주고 있다. 실험에 사용한 각 언어 모델들에 대해 Body, Disease, Treatment 태그별 F1-score는 표 6과 같다.

표 5. 한국어 BIO 개체명 인식 성능 비교
Table 5. Performance Comparison of Korean Biomedical Named Entity Recognition

Derived model (Pre-trained language model)	F1-score
KlueBERT	82.16
KlueRoBERTa	83.82
LMkorBERT	82.00
KR-BERT	79.30

표 6. 언어모델별 생의학 태그들의 F1-score
Table 6. F1-score for Biomedical tags by Language Model

	Body	Disease	Treatment
KlueBERT	0.81	0.86	0.77
KlueRoBERTa	0.81	0.86	0.79
LMkorBERT	0.79	0.86	0.76
KR-BERT	0.73	0.81	0.69

표 6의 태그별 인식 결과에서 Disease가 Body와 Treatment 보다 상대적으로 좋은 성능을 보였다. [6]

의 태그별 인식에서도 DZ(상병)가 SX(증상)와 BP(신체 부위) 보다 좋은 F1-score를 보였는데, 이는 많은 상병이 하나 이상의 단어로 이루어진 데 반해, SX와 BP는 보통 한 단어로 이루어져 발생한 태그 사이의 불균형으로 분석하고 있다.

표 7은 예측용 데이터를 학습된 언어 모델에 입력하여 태깅된 출력 중 잘못된 예들을 보여준다. ‘관절염 및 연골장애’ 문장에서 ‘잇’은 질병 관련 단어가 아닌 접속사지만 ‘Disease-I’로 태깅된 것을 보여 주고 있다. 마찬가지로 ‘관절의 연골이’에서도 단어를 연결하는 조사 ‘의’가 ‘Disease-I’로 오분류되었다. ‘급성상기도염’은 ‘급성’과 ‘상기도염’으로 분할되어야 하지만 ‘급성’, ‘상’기도’, ‘염’으로 하나의 질병명이 여러 개로 과분할되었다. 이에 반해, ‘고혈압성 질환’은 (‘고혈압’, ‘Disease-B’), (‘성’, ‘Disease-I’), (‘질환’, ‘Disease-I’)으로 과분할되었으나 모두 ‘Disease’로 태깅되었다. ‘몸’, ‘손’과 같이 한 단어로 이루어진 신체 부위가 ‘Body’ 태그가 아닌 ‘오’로 올바르게 태깅되지 않은 경우가 많이 발견되었다. 특히, ‘손씻기’는 ‘손’과 ‘씻기’로 분할되지 못해 ‘오’로 정확하게 태깅되지 않았다. 이런 오분류나 미분류는 문맥에 기반하여 제대로 개체명을 분류하지 못해서 발생한다.

표 7 잘못된 태깅의 예
Table 7. Examples of incorrect tagging

Incorrect	Correct
(‘관절염’, ‘Disease-B’), (‘잇’, ‘Disease-I’), (‘연골’, ‘Disease-I’), (‘장애’, ‘Disease-I’)	(‘관절염’, ‘Disease-B’), (‘잇’, ‘O’), (‘연골’, ‘Disease-I’), (‘장애’, ‘Disease-I’)
(‘관절’, ‘Body-B’), (‘의’, ‘Body-I’), (‘연골’, ‘Body-I’), (‘이’, ‘O’)	(‘관절’, ‘Body-B’), (‘의’, ‘O’), (‘연골’, ‘Body-I’), (‘이’, ‘O’)
(‘급성’, ‘O’), (‘상’, ‘Disease-I’), (‘기도’, ‘Disease-I’), (‘염’, ‘Disease-I’)	(‘급성’, ‘Disease-B’), (‘상’, ‘Disease-I’), (‘기도’, ‘Disease-I’), (‘염’, ‘Disease-I’)
(‘홍반’, ‘Disease-B’), (‘이’, ‘Disease-I’), (‘있는’, ‘O’)	(‘홍반’, ‘Disease-B’), (‘이’, ‘O’), (‘있는’, ‘O’)
(‘때로는’, ‘O’), (‘몸’, ‘O’), (‘의’, ‘O’)	(‘때로는’, ‘O’), (‘몸’, ‘Body-B’), (‘의’, ‘O’)
(‘손씻기’, ‘O’)	(‘손’, ‘Body-B’), (‘씻기’, ‘O’)
(‘당뇨병’, ‘Disease-B’), (‘은’, ‘O’), (‘대사’, ‘Disease-B’), (‘장애’, ‘O’), (‘의’, ‘O’)	(‘당뇨병’, ‘Disease-B’), (‘은’, ‘O’), (‘대사’, ‘Disease-B’), (‘장애’, ‘Disease-I’), (‘의’, ‘O’)
(‘약성’, ‘Disease-B’), (‘신’, ‘Disease-I’), (‘생물’, ‘O’)	(‘약성’, ‘Disease-B’), (‘신’, ‘Disease-B’), (‘생물’, ‘Disease-I’)
(‘호흡기’, ‘O’), (‘질환’, ‘O’)	(‘호흡기’, ‘Disease-B’), (‘질환’, ‘Disease-I’)

14) <http://konlpy.org/en/latest>
15) <http://konlpy.org/ko/v0.6.0/morph>
16) <https://github.com/open-korean-text/open-korean-text>

이런 오분류를 줄이기 위해서는 과분류나 미분류된 서브단어(subword)들을 하나의 단어로 인식되도록 보다 정밀한 형태소 분석기와 풍부한 어휘사전 구축이 선행되어야 할 것이다.

V. 결론

개체명 인식은 자연어 처리의 한 분야로 정보 검색, 질의응답, 챗봇 등에 매우 중요한 연구 과제로 트랜스포머 기반 대규모 언어 모델들의 등장과 함께 활발한 연구가 진행되고 있다.

본 연구에서는 한국어 생의학 개체명 인식을 위해 네이버 개체명 말뭉치와 한국어 생의학 개체명 말뭉치인 KBMC를 기반으로 실험용 말뭉치를 구축하고, 선정된 한국어 BERT 파생 모델들에 전이 학습시켜 태깅 인식 성능을 비교하였다. 인식 성능에 있어서는 KlueRoBERTa가 상대적으로 좋은 성능을 보였으나 처리 시간이 타 모델들보다 많이 소요되었다. 전이 학습시킨 언어 모델들을 대상으로 ChatGPT-4로 생성한 문장들로 태깅을 예측하여 오류 분석을 하였다. 태그별 인식에서는 Disease가 Body와 Treatment보다 상대적으로 좋은 결과를 보였다. 이는 많은 질병명이 하나 이상의 단어로 이루어진 데 반해, 신체와 처치는 보통 한 단어로 이루어져 발생한 태그 사이의 불균형으로 보인다. 오류 분석 결과, 대부분의 오분류는 과분할과 미분할로 인한 것이었다. 이를 보완하기 위해서는 문맥을 고려한 형태소 분석과 더욱 충실한 생의학 어휘사전 구축이 필요하다.

References

- [1] J. Lee, O. Kwon, "Performance Assessment of Machine Learning and Deep Learning in Regional Name Identification and Classification in Scientific Documents," *J. of The Korea Institute of Electronic Communication Sciences*, vol. 19, no. 2, 2024, pp. 389-396.
<https://doi.org/10.13067/JKIECS.2024.19.2.389>
- [2] J. L. A. Sun, J. Han, and C. Li, "A Survey on Deep Learning for Named Entity Recognition," *IEEE Trans. on Knowledge and Data Engineering*, vol. 34, no. 1, 2020, pp. 50-70.
<https://doi.org/10.48550/arXiv.1812.09449>
- [3] J. Lee, W. Yoon, S. Kim, D. Kim, S. Kim, C. So, and J. Kang, "BioBERT: a pre-trained biomedical language representation model for biomedical text mining," *Bioinformatics*, vol. 36, issue 4, 2019, pp. 1234-1240.
<https://doi.org/10.48550/arXiv.1901.08746>
- [4] K. Huang, J. Altsaar, and R. Ranganath, "ClinicalBERT: Modeling Clinical Notes and Predicting Hospital Readmission", In *Proc. of CHIL(Conference on Health, Inference, and Learning)'20 Workshop*, Toronto, ON, Canada, April 2-4, 2020.
<https://doi.org/10.48550/arXiv.1904.05342>
- [5] Y. Li, R. Wehbe, F. Ahmed, H. Wang, and Y. Lao, "A Comparative study of pretrained language models for long clinical text", *J. of the American Medical Informatics Association*, vol. 30, no. 2, 2023, pp. 340-347.
<https://doi.org/10.48550/arXiv.2301.11847>
- [6] Y. Kim and T. Lee, "Korean clinical entity recognition from diagnosis text using BERT," *BMC Medical Informatics and Decision Making*, vol. 20, supplement 7, 2020.
<https://doi.org/10.1186/s12911-020-01241-8>
- [7] S. Byun, J. Hong, S. Park, D. Jang, J. Seo, M. Kim, C. Oh, and H. Shin, "Korean Bio-Medical Corpus (KBMC) for Medical Named Entity Recognition," In *Pro. of the 2024 Joint Int. Conf. on Computational Linguistics, Language Resources and Evaluation*, Torino, Italy, 2024, pp. 9941-9947.
<https://doi.org/10.48550/arXiv.2403.16158>
- [8] Y. Kim, J. Kim, J. Lee, M. Jang, Y. Yum, S. Kim, U. Shin, Y. Kim, H. Joo, and S. Song, "A pre-trained BERT for Korean medical natural language processing," *Scientific Reports*, vol. 12, no. 1, 2022.
<https://doi.org/10.1038/s41598-022-17806-8>

- [9] J. Lee, "Comparison of Sentiment Classification Performance of for RNN and Transformer- Based Models on Korean Reviews, " *J. of The Korea Institute of Electronic Communication Sciences*, vol. 18, no. 4, 2023, pp. 693-700.
<https://doi.org/10.13067/JKIECS.2023.18.4.693>
- [10] J. Lee, "Performance Comparison and Error Analysis of Transformer-based Korean Named Entity Recognition," *J. of Jeonnam State University*, vol. 25, 2023, pp. 293-302.

저자 소개

이재홍(Jae-Hong Lee)



1999년 충남대학교 대학원 컴퓨터공학과(공학박사)
1988년-1993년 국방과학연구소 연구원
2000년-현재 전남도립대학교 보건의료과 교수
※ 관심분야 : 자연어 처리, AI

