

프로젝트 Wrap Up

1. 프로젝트 개요

본 프로젝트에서는 2022 베이징 동계 올림픽과 관련된 위키 원시 말뭉치를 활용해 자연어처리 관계 추출 태스크에 쓰이는 주석 코퍼스를 제작했습니다. 프로젝트의 의의는 한국어 및 다른 언어에서의 자연어처리 데이터셋의 유형 및 포맷이 어떠한지, 그리고 데이터셋을 구축하는 일반적인 프로세스가 무엇인지 학습하는 것입니다.

2. 프로젝트 팀 구성 및 역할

이름	역할
공통	가이드라인 작성, Entity•Relation 정의, 파일럿 및 메인 어노테이션
강나경	카테고리별 문장 split, fleiss-kappa 계산
김산	가이드라인 FAQ 작성
김현지	데이터셋 전처리, 가이드라인 이미지 제작
정민지	모델 Fine-tuning, 데이터셋 분석
최지연	여러 개의 파일을 카테고리별로 분류하여 통합

3. 프로젝트 수행 절차 및 방법

3-1. 데이터셋 소개

- 베이징 동계 올림픽 관련 위키 데이터
- 총 43개의 문서, 문장 1,693개로 구성
- 예시



2022년 동계 올림픽은 2022년 2월 4일부터 2월 20일까지 중화인민공화국 베이징에서 열린 동계 올림픽이다.

3-2. 데이터 제작 과정



3-3. 마주한 문제점 및 해결 방안

1. 여러 주제로 소분류 된 데이터셋
 - 대분류 베이징 동계 올림픽에서 소분류 2022, Building, Organization, Sports로 분할
 - 소분류된 카테고리 내에서 등장하는 관계 선정
2. 여러 명의 작업자간 작업물 일치

9	<input checked="" type="checkbox"/>	수비수 안셀 살렐라가 첫 홈 개막 경기에서 선제골 을 기록하여 이는 이 경기장의 첫 골 로 기록되었다.	안셀 살렐라 - 이 경기장의 첫 골 (per - rst) 가능할까요?	현지) 선제골...만 경기 결과...라고 해야하지 않을까 싶지만 "선제골"이 경기 결과 (승패, 메달 등) 인가...? 싶네요 (부덕:혼란)	민지) 음 선제골을 보통 경기의 결과로 보지는 않는 것 같아요..! 경기 결과 선제골을 했습니다 라는 걸 못들어본 거 같아요
10	<input checked="" type="checkbox"/>	베이징 수도 체육관: 2004년 10월 17일, 전미 농구 협회(NBA)의 2004-05시즌 시범 경기 를 수도 실내 체육관에서 진행했다.	시범경기(이벤트)의 태깅 범위에 "전미 농구 협회"가 포함될까요?	ks) ~의 ~~ 가 애매한 것 같은데 저는 웬만하면 ~의까지 포함하고 있던 했습니다!	지연) '시범경기'라고 하면 그냥 일반명사인데, 앞에 수식을 붙여줘서 특정 경기를 의미하는게 되는거니, 포함하는게 더 적절할 것 같다고 생각합니다!
11	<input checked="" type="checkbox"/>	준결승 1조에서 경기를 치른 대한민국 선수 황태현과 2조에서 경기를 치른 대한민국 선수 이준서가 각각 조 1위, 조 2위로 결승선을 통과 했다.	이 문장에서 보면 결승선을 통과한 것이 경기 결과인 것 같습니다. 그런데 다음 문장에서는 실격당했다 라는 내용이 나옵니다! 이 문장만 가지고 판단할 때는 조 1위, 조 2위가 경기결과라고 생각할 수 있을까요?	현지) 황태현-조 1위, 이준서-조 2위 로 태깅하면 될 것 같습니다. 본 문장에서 실격이라는 단어가 나오지 않으니깐요! (외부 지식 활용으로 생각할 수 있을 것 같음!)	
12	<input checked="" type="checkbox"/>	베이징 수도 체육관: 2022년, 2022년 동계 올림픽 쇼트트랙, 피겨스케이팅 종목 경기를 개최하고 있다.	이 경우 저희가 사전에 정의한 심포 등으로 인해 분할되는 케이스에 해당되어 "피겨스케이팅"은 이벤트 태깅 대상이 아닐까요?	민지) 피겨스케이팅만을 이벤트로 보기는 애매한 것 같습니다! 2022년 동계 올림픽 쇼트트랙은 이벤트인 것 같지만요, 그리고 말씀해주신대로 심포 분할 케이스인것 같습니다!	
13	<input checked="" type="checkbox"/>	한국의 김연아	김연아(per) - 한국(org) 가능하겠죠!..??	현지) 넵! 가능하다고 생각합니다!	
14	<input type="checkbox"/>	캐나다 중부 지역에서 컬링 하던 사람들은 19세기 초까지 스톤 보다는 찰 을 종종 사용했다.	1. (컬링-아이템-찰) 가능할까요? 19세기 초까지 사용한거라 지금은 사용하지 않는거라서 시점 생각하면 아닌것 같기도 해서 고민입니다! 2. 여기서의 "스톤"이 컬링의 (현재 사용하는)아이템이라는걸 유추할 수 있을까요?	현지) 찰...은 따로 아이템은 아닌 것 같습니다! 구성 요소? 가 더 적절할 것 같아요 ㅋㅋㅋㅋㅋ 여기서 스톤은... 돌을 의미하는 걸까요? 뭔가 문장내 뒤임스가 돌 vs 찰 의 느낌이네요	
15	<input type="checkbox"/>	컬링: 스톤 핸들 에 불빛이 녹색이면 적법하게 투구했음을 말해주며 (보류: 아이템 유지?)	"스톤 핸들"은 아이템 vs 전문 용어, 당신의 선택은? (두둥)	민지) 전문용어..? 요 근데 애매한 것 같습니다!ㅠ 이거 다시 논의해보면 좋을 것 같습니다!	지연) 저는 아이템이라고... 생각했습니다. 손잡이만 따로 떼기도 하더라고요
16	<input checked="" type="checkbox"/>	베이징 지구 는 주로 동계 올림픽의 빙상 종목 의 경기를 치를 예정이고, 개막식 과 폐막식 이 열릴 예정이다.	베이징 지구가 건물은 아니지만, 장소의 용도가 나오기도 하는 것 같습니다! (그냥 발견, 건물:용도 관계의 확장이 필요한가?에 대한 궁금증)		
17	<input checked="" type="checkbox"/>	컬링: 이 대회에서 이듬해 2007년 월드컵 챔피언 이 된 켈리 스콧 이 1998년 월드컵 챔피언인 동메달리스트였던 캐시 킹 을 상대로 한 경기의 한 엔드에서 8점을 기록했다.	"2007년 월드컵 챔피언"은 켈리 스콧의 경기 결과에 해당되지 않겠죠? (FAQ 업적 참고) (캐시 킹은 저렇게 태깅했습니다!)	민지) - 우선, 캐시킹에 대해서 <1998년 월드컵 챔피언인 동메달> 이렇게 경기결과 태깅하기로 했던 것으로 기억합니다! - 월드컵 챔피언도 경기결과 아닐까요..?	지연) 챔피언이 '우승자'라는 뜻이어서... '우승'만을 떼어서 태깅할 수 있다면 딱 좋겠는데, 이 경우는 좀 애매한듯하네요ㅠㅠ 어렵군요! 그리고 동메달 부분은 민지님 의견과 같습니다!
18	<input checked="" type="checkbox"/>	2022년 동계 올림픽 빅 에어 종목 개최를 위해 빅 에어 쇼우강 경기장 을 2018년 착공했다. (건물:용도 없애지는거 예약)	해당 문장은 "이벤트: 개최 장소" 일까요? "건물:용도" 일까요? 아래 중에서 적합한 태깅은 무엇일까요? 2022년 동계 올림픽 빅 에어 종목 개최를 위해 빅 에어 쇼우강 경기장 을 2018년 착공했다. > evt:holding_place 2022년 동계 올림픽 빅 에어 종목 개최를 위해 빅 에어 쇼우강 경기장 을 2018년 착공했다. > bid:spo(1) 2022년 동계 올림픽 빅 에어 종목 개최를 위해 빅 에어 쇼우강 경기장 을 2018년 착공했다. > bid:spo(2)	민지) 1번과 2번 중에는 2번이 더 적절하다고 생각합니다! 근데 이벤트:개최 장소, 건물:용도가 명확하게 구분되지 않는 것 같기도 합니다!	지연) 이벤트인 경우 '이벤트:개최 장소', 스포츠인 경우 '건물:용도'로 하기로 했던걸로 기억합니다. 이 문장의 경우는 이벤트인듯 하면서 스포츠인듯도 해서 어렵군요TTTT 문맥상으로는 개최를 '위해' 착공한것이니 스포츠라고 보고 '건물:용도'로 태깅하는 것이 좋을 듯 합니다(2번)
19	<input checked="" type="checkbox"/>	빅 에어 쇼우강: 2022년 동계 올림픽의 스노보드 빅 에어 종목과 프리스타일 스키 빅 에어 종목을 이 경기장에서 진행한다. (건물:용도 없애짐)	12번이랑 동일한 질문 프리스타일 스키 빅 에어 종목은 태깅 대상이 아니겠죠?	지연) 넵, 앞의 '2022년 동계 올림픽' 부분이 '스노보드 빅 에어 종목'과 '프리스타일 스키 빅 에어 종목'을 모두 수식하는 걸로 보여서 현지님 말씀대로 뒤에 오는 '프리스타일 스키 빅 에어 종목'은 태깅 대상이 아닌 것 같아요 근데 저희 '종목'은 제외하기로 했던 것 같은데.. 그럼 '2022년 동계 올림픽의 스노보드 빅 에어' 부분만 태깅해야 하지 않을까요!	

- 동일한 작업 결과물을 얻기 위해 헛갈릴 수 있는 entity의 영역과 relation 관계 같은 경우는 slack과 구글 스프레드시트를 활용하여 팀원들과 의견을 정리한 후, 방향을 명확하게 잡았습니다.

3-4. Relation 정의 기준 및 방식 소개

- Relation set의 구성 및 정의, 가이드라인 작성, 파일럿 및 메인 어노테이션, 그리고 간단한 모델 Fine-tuning의 과정을 통해 실제 데이터 제작의 workflow를 경험해 볼 수 있습니다.
- 이 과정에서 정밀한 가이드라인 제작의 중요성과 inter-annotator agreement(IAA)의 개념을 체득할 수 있습니다.

4. 프로젝트 수행 결과

4-1. 가이드라인 소개

가이드라인

4-2. 데이터셋 분석 결과

스포츠에서 사용하는 용어에 대한 관계

각 관계의 데이터 개수 분포



