

SK네트웍스 Family AI과정 6기

데이터 전처리 학습된 인공지능 모델

□ 개요

- 산출물 단계 : 데이터 전처리
- 평가 산출물 : 이진 분류 모델 (학습된 인공지능 모델 1)
- 제출 일자 : 2025.03.17.
- 깃허브 경로 : [git@github.com:SKNETWORKS-FAMILY-AICAMP/SKN06-FINAL-4Team.git](https://github.com/SKNETWORKS-FAMILY-AICAMP/SKN06-FINAL-4Team.git)
- 작성 팀원 : 김동명

사용 데이터 및 기술	<ul style="list-style-type: none">• KcELECTRA pre-trained model, NSMC로 Fine tuning *NSMC(Naver Sentiment Movie Corpus, train/test = 150000/50000)• meta data는 네이버리뷰, 왓챠피디아 리뷰(최소 30개, 최대 150개)• Ray tune, Grid search 으로 하이퍼 파라미터 탐색• grad_accum_steps : 4• 적용된 하이퍼 파라미터 Batch_size, Weight_decay, Warm_up, early stop, Learning_rate• AdamW optimizer 사용• sklearn, matplotlib 으로 eval data 생성
상황	<ul style="list-style-type: none">• 1. GPU, 모델 선정 LSTM 보다 성능이 좋은 LLM 모델 활용, VRAM 용량 문제로 Distil-bert 및 KcELECTRA 검토, 최종 모델 선정(KcELECTRA), GPU 는 L4(테스트), L40S(최종학습) 사용 *GPU : L40s (18000+ CUDA cores, 48GB VRAM, 16 vCPU)• 2. 라이브러리 활용 모델 학습을 확인하기 위해 MLflow, Ray tune 을 검토, 실시간으로 확인이 용이하고 VS code 상에서 바로 확인이 편리한 Ray 사용• Validation set을 적용해도 학습이 지나치게 빠르게 종료되고 과대적합 문제가 발생 -> Warm up 적용, 일반화 성능 향상을 위해 Weight decay 적용• 메모리 부담을 줄이기 위해 Grad accumulation 4 적용(최종 Batch size = 8)
수행 결과	<ul style="list-style-type: none">• 하이퍼 파라미터 탐색 범위 Learning rate : $1e-3 \sim 5e-6$ batch_size : 8, 16, 32 epoch max : 10 Weight decay : 0.1, 0.01, 0.001 Warm up : 0.1, 0.2, 0.3, 0.4
보완점	<ul style="list-style-type: none">• NSMC 의 학습데이터에 오류가 확인됨(더 나은 train data 필요)• 특수문자, 이모티콘에 대한 해석이 불분명• 장르에 따라 꺼림칙함, 무서움, 공포와 같은 내용을 부정으로 판단

수 행 결 과		DistilBERT	KoBERT	KcELECTRA
Hyperparameter	Learning_rate	6.73e-05	5.52e-05	6.45e-05
	Batch_size	16	16	8
	epochs	5	2	4
Strategy	Warm-up weight decay	0.2 0.1	0.2 0.01	0.3 0.01
Performance	Train loss		0.271	0.177
	Validation loss	0.328	0.297	0.2359
	Validation accuracy	0.861	0.876	0.9107

- 평가 산출물 : 다중 분류 모델 (학습된 인공지능 모델 2)
- 제출 일자 : 2025.03.17.

사용 데이터 및 기술	<ul style="list-style-type: none"> ● KcELECTRA pre-trained model, KOTE로 Fine tuning한 가중치 사용 *KOTE(Korean Online Contents Emotions, train/test = 45000/5000) ● meta data는 네이버리뷰, 왓차피디아 리뷰(최소 30개, 최대 150개)
상황	<ul style="list-style-type: none"> ● 논문에 따르면 하나의 텍스트(리뷰)에서 2개 이하의 감정이 포함될 확률 99%, 3개 이하의 감정이 포함될 확률 86%로 분석 결과 감정의 강도를 수치로 표현 -> 감정 강도 평균을 리뷰의 감정 포인트로 계산
수행 결과	<ul style="list-style-type: none"> ● 대부분의 경우 만족할만한 성능
보완점	<ul style="list-style-type: none"> ● 댓글에 감정표현이 잘 나타나지 않는 경우 ● 특수문자, 이모티콘에 대한 해석이 불분명 ● 감정이 나타나지 않아 “없음”으로 표기된 경우