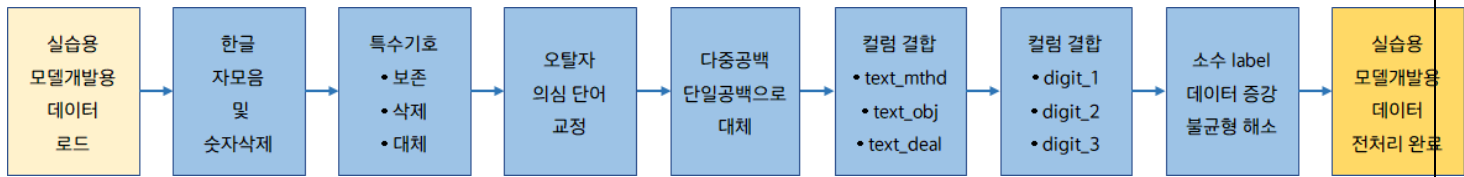


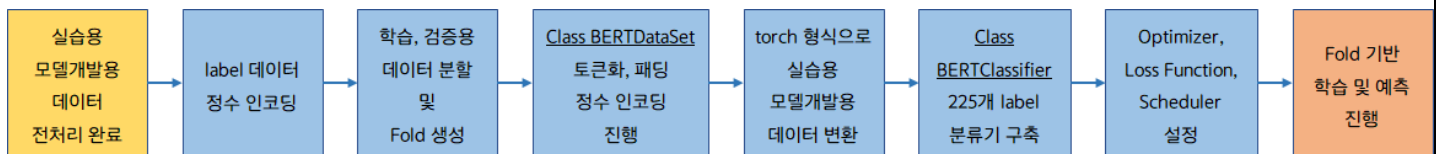
- 오탈자 의심 단어 직접 교정(케? -> 케익 / 숙박시? -> 숙박시설 등)
- text_obj / text_mthd / text_deal 3개의 column을 순서대로 결합
- digit_1 / digit_2 / digit_3 3개의 column 결합해 label column 생성해 예측할 target 데이터 생성
- 학습데이터 100만개 중에서 label 빈도수가 20개 미만인 label의 경우 데이터 개수 3배로 증강



[그림3 : 전처리 도식]

2) 사전학습 언어모델 – KoBERT (출처 : <https://github.com/SKTBrian/KoBERT>)

- KoBERT는 트랜스포머로 구현된 양방향 자연어 처리 모델인 BERT의 한국어 성능 한계를 극복하기 위해 위키피디아나 뉴스 등을 수집한 수백만 개의 한국어 문장으로 이뤄진 대규모 말뭉치(corpus)를 추가학습한 언어모델
- 본 팀은 SKTBrian팀이 제공하는 KoBERT 모델과 토큰라이저를 사용함.
- 사전학습된 KoBERT모델에 모델개발용 데이터를 전처리해 학습하는 fine-tuning하여 예측함



[그림4 : KoBERT 모델링 도식]

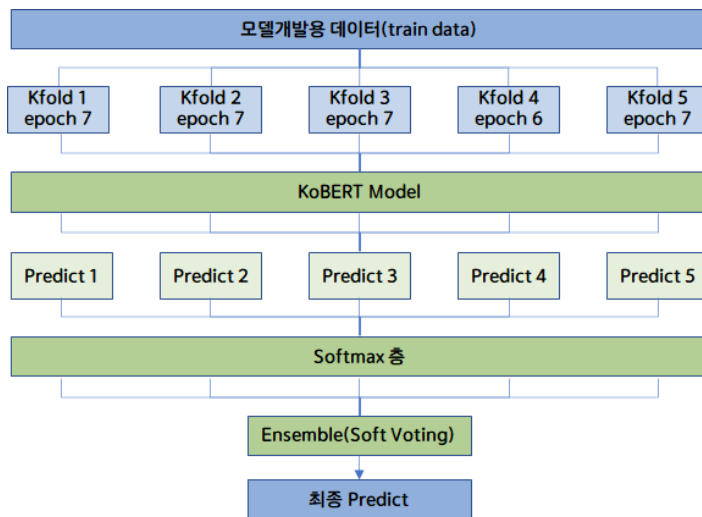
- 모델개발용 데이터를 StratifiedKfold를 이용해 학습, 검증용데이터로 분할
- Class BERTDataset을 정의해 전처리된 데이터가 KoBERT 모델의 입력가능 하도록 토큰화, 정수 인코딩, 패딩 진행
- KoBERT 모델의 8000여개 대규모 말뭉치(Corpus)를 기반으로 정수 인코딩
- 데이터의 각 문장을 토큰화하고, 토큰화된 문장의 최대 차원 수(max_seq_len)를 임의로 지정하고, max_seq_len보다 차원의 수가 작다면 1로 패딩해 모든 토큰화된 문장의 차원 통일함
- DataLoader를 이용해 torch 형식의 데이터로 변환
- Class BERTClassifier로 학습용 데이터의 225개의 label(소분류 기준으로 분류된 산업군 개수) 분류기를 정의함

3) 모델 hyper-parameter 정보

실험을 통해 최적의 파라미터 및 손실함수 설정

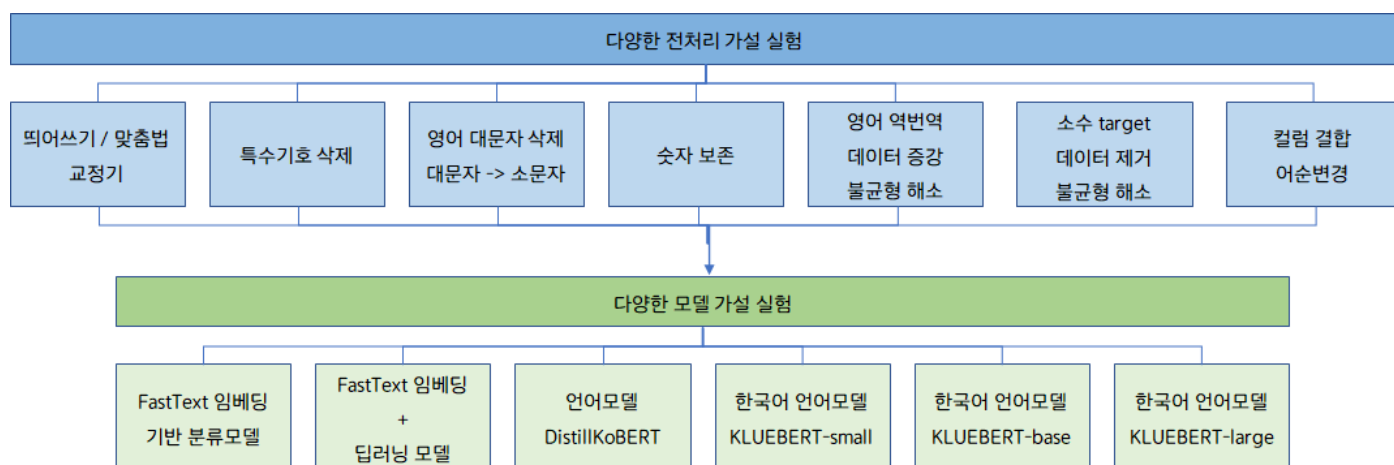
Loss Fuction	Drop rate	Batch size	Max seq len	Epoch	Hidden size	Max Grad Norm	Log Interval	Learning rate
Cross Entropy	0.2	64	64	6~7	768	1	200	5e-5

- StratifiedKfold : 특정 label에 편중된 불균형한 분포도를 가진 label(결정 클래스) 데이터 집합을 위한 KFold 방식
- Ensemble(Soft Voting) : 각 label 결정확률을 더한 평균 중 확률이 가장 높은 label을 최종 보팅 결과값으로 선정함



4) 기타 코드

다양한 전처리 및 모델 가설을 실험해보았고 성능이 낮아지는 경우 해당 항목을 배제함



신청자	소속/직위/팀명	건국대학교	성명	김유빈
	휴대전화	010-7164-6794	전자우편	kimyusintwo@gmail.com
제출일	2022년 4월 15일			