



도배 하자 질의 응답 처리

DACON 한솔데코 시즌2 AI 경진대회

EURON 5기 중급 하자하자!팀 | 김유민 박은혜 이아영

목차

01 문제 정의 및
배경

02 데이터 수집 및
전처리

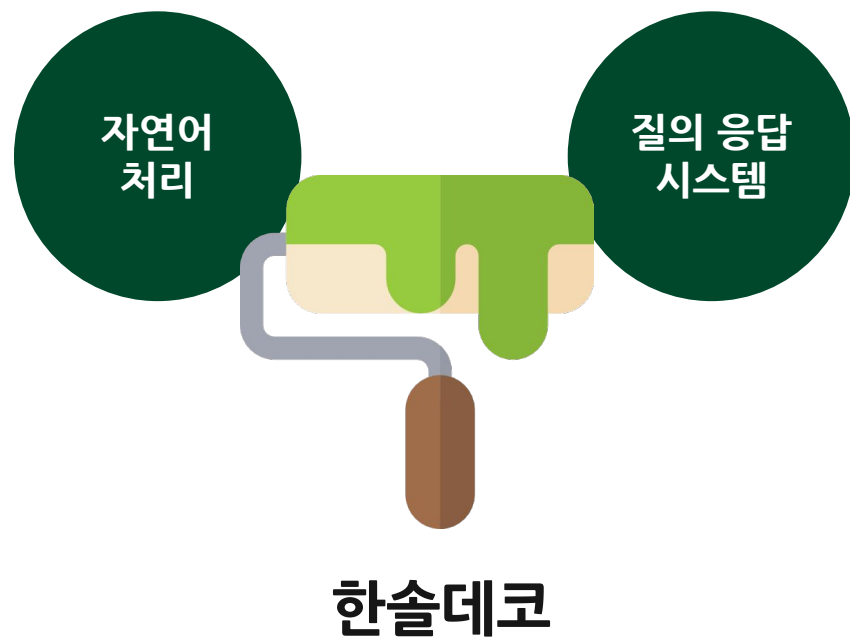
03 모델 개요 및 학습
방법

04 결과 분석 및
평가 지표

05 시스템 구현 계획 및
기술적 이슈

06 결론 및 제언

문제 정의 및 배경



데이터 수집 및 전처리

- 데이터 수집

공모전 주최측에서 제공한 train, test 데이터 사용

train data

	id	질문_1	질문_2	category	답변_1	답변_2	답변_3	답변_4	답변_5
0	TRAIN_000	면진장치가 뭐야?	면진장치에 사용되는 주요 기술은 무엇인가요?	건축구조	면진장치란 지반에서 오는 진동 에너지를 흡수하여 건물에 주어진 진동을 줄여주는 진동...	면진장치란 건물의 지반에서 발생하는 진동 에너지를 흡수하여 건물을 보호하고, 진동을...	면진장치란 지반으로부터 발생하는 진동 에너지를 흡수하여 건물에 전달되는 진동을 줄여...	면진장치는 건물의 지반으로부터 오는 진동 에너지를 흡수하여 건물에 전달되는 진동을 ...	면진장치는 건물에 오는 지반 진동의 영향을 최대한으로 흡수하여 건물에 전달되는 진동...

test data

	id	질문
0	TEST_000	방청 페인트의 종류에는 어떤 것들이 있는지 알고 계신가요? 또한, 원목사이딩을 사용...
1	TEST_001	도배지에 녹은 자국이 발생하는 주된 원인과 그 해결 방법은 무엇인가요?
2	TEST_002	큐블럭의 단점을 알려주세요. 또한, 압출법 단열판을 사용하는 것의 장점은 무엇인가요?
3	TEST_003	철골구조를 사용하는 고층 건물에서, 단열 효과를 높이기 위한 시공 방법은 무엇이 있...
4	TEST_004	도배지의 완전한 건조를 위해 몇 주 동안 기다려야 하나요?

데이터 수집 및 전처리

데이터 전처리

for문을 활용해 질문과 답변이 1:1 대응이 되도록 데이터 포매팅

question	answer
0 면진장치가 뭐야?	면진장치란 지반에서 오는 진동 에너지를 흡수하여 건물에 주는 진동을 줄여주는 진동 ...
1 면진장치가 뭐야?	면진장치란 건물의 지반에서 발생하는 진동 에너지를 흡수하여 건물을 보호하고, 진동을...
2 면진장치가 뭐야?	면진장치란 지반으로부터 발생하는 진동 에너지를 흡수하여 건물에 전달되는 진동을 줄여...

토큰화/인코딩

데이터셋의 질문과 답변을 모델의 입력으로 사용하기 위해 텍스트를 토큰화하고 숫자로 변환

```
[tensor([[ 9411,  8265, 37765, 46651,  7991,   406,     1,  9411,  8265, 20725,
          7374,  9027,  7599,  9023, 14472, 15898, 14820, 33220, 36928, 10764,
          9166, 11818, 28037, 10090, 15898, 34062, 20725, 21154]])],
```

데이터 수집 및 전처리

- 데이터 증강

train data 양을 늘려 모델의 정확도를 개선하기 위해 2가지 방법으로 진행하여 데이터 증강

- **MLM data augmentation**

Masked Language Modeling 방식으로 학습 모델을 학습한 후에 새로운 문장의 일부에 마스킹을 적용하고 인퍼런스를 적용해 마스킹된 부분에 알맞는 새로운 토큰을 후보로 생성

기존 문장	증강된 문장
면진장치란 지반에서 오는 진동 에너지를 흡수하여 건물에 주는 진동을 줄여주는 진동 격리장치입니다.	면진장치란 지반에서 나오는 진동 에너지를 전달하여 건물에서의 진동을 줄여주는 진동의장치다.

데이터 수집 및 전처리

- BERT augmentation

Bert based 모델을 활용하여, 의미상 자연스러운 토큰을 삽입하거나 대체(masking, insertion)하는 형식으로 문장 augmentation 수행

Masking	기존	면진장치란 건물의 지반에서 발생하는 진동 에너지를 흡수하여 건물을 보호하고, 진동을 줄여주는 장치입니다. 주로 지진이나 기타 지반의 진동으로 인한 피해를 방지하기 위해 사용됩니다.
	증강	면진장치란 건물의 지반에서 발생하는 진동 에너지를 막아 건물을 보호하고, 최대한 줄여주는 데 주로 지진이나 기타 지반의 진동으로 인한 피해를 방지하기 위해 사용됩니다.
Insertion	기존	면진장치란 지반에서 오는 진동 에너지를 흡수하여 건물에 주는 진동을 줄여주는 진동 격리장치입니다.
	증강	면진장치란 지반에서 전해 오는 진동 에너지를 흡수하여 건물에 주는 진동을 최대한 줄여주는 진동 격리장치입니다.

모델 개요 및 학습 방법

KoGPT 모델

gpt-3 기반 한국어 언어 생성모델

- 한국어를 사전적, 문맥적으로 이해하고 이용자가 원하는 결과값을 보여줌
- 60억개의 매개변수와 2000억개 토큰(token)의 한국어 데이터를 바탕으로 구축
- 주어진 텍스트의 다음 단어 예측 가능

```
model = GPT2LMHeadModel.from_pretrained('skt/kogpt2-base-v2')
```


모델 개요 및 학습 방법

KoGPT 모델 하이퍼 파라미터

[하이퍼파라미터 선택]

- learning_rate(학습률) : $2e-5$
- epoch(에포크) : 10
- batch_size = 4

[최적화 알고리즘]

- AdamW로 설정 : 가중치 감쇠(weight decay)를 적용하여 더욱 효과적인 학습을 할 수 있도록 함

결과 분석 및 평가 지표

평가 산식 : Cosine Similarity

- 내부 곱 공간의 두 벡터간의 유사성 측정
- 생성된 답변을 **Sentence Transformer** 모델을 이용하여 512 차원의 Embedding Vector로 변환한 후, 변환된 벡터와의 코사인 유사도 계산
- (코사인 유사도 값이 0보다 작은 경우 0으로 간주)

결과 분석 및 평가 지표 - 증강 전 데이터

- 증강한 데이터는 훈련 시간이 너무 오래걸려 비효율적이라고 판단하여 우선 원래 train 데이터로 훈련 진행하며 비교

1차 시도

```
Epoch 1 - Avg Loss: 2.8414: 100%|██████████| 6440/6440 [08:35<00:00, 12.49it/s]
Epoch 1/10, Average Loss: 2.8414353972358732
Epoch 2 - Avg Loss: 1.7499: 100%|██████████| 6440/6440 [08:17<00:00, 12.93it/s]
Epoch 2/10, Average Loss: 1.7499371493778983
Epoch 3 - Avg Loss: 1.1237: 100%|██████████| 6440/6440 [08:20<00:00, 12.88it/s]
Epoch 3/10, Average Loss: 1.1237117445385605
Epoch 4 - Avg Loss: 0.7592: 100%|██████████| 6440/6440 [08:32<00:00, 12.56it/s]
Epoch 4/10, Average Loss: 0.7592249020633305
Epoch 5 - Avg Loss: 0.5347: 100%|██████████| 6440/6440 [08:19<00:00, 12.90it/s]
Epoch 5/10, Average Loss: 0.5346679408242059
Epoch 6 - Avg Loss: 0.3969: 100%|██████████| 6440/6440 [08:21<00:00, 12.85it/s]
Epoch 6/10, Average Loss: 0.3969358843988421
Epoch 7 - Avg Loss: 0.3167: 100%|██████████| 6440/6440 [08:25<00:00, 12.74it/s]
Epoch 7/10, Average Loss: 0.3167410741918807
Epoch 8 - Avg Loss: 0.2641: 100%|██████████| 6440/6440 [08:26<00:00, 12.71it/s]
Epoch 8/10, Average Loss: 0.2641438113820525
Epoch 9 - Avg Loss: 0.2307: 100%|██████████| 6440/6440 [08:22<00:00, 12.81it/s]
Epoch 9/10, Average Loss: 0.23065063478367298
Epoch 10 - Avg Loss: 0.2051: 100%|██████████| 6440/6440 [08:20<00:00, 12.87it/s]
Epoch 10/10, Average Loss: 0.2050533873121607
```

하이퍼파라미터 튜닝 + 배치정규화

```
Epoch 1 - Avg Loss: 1.7021: 100%|██████████| 1610/1610 [04:46<00:00, 5.62it/s]
Epoch 1/10, Average Loss: 0.4255206729592009
Epoch 2 - Avg Loss: 1.0646: 100%|██████████| 1610/1610 [04:45<00:00, 5.64it/s]
Epoch 2/10, Average Loss: 0.2661546274325492
Epoch 3 - Avg Loss: 0.7697: 100%|██████████| 1610/1610 [04:45<00:00, 5.65it/s]
Epoch 3/10, Average Loss: 0.19242231228521892
Epoch 4 - Avg Loss: 0.5883: 100%|██████████| 1610/1610 [04:44<00:00, 5.66it/s]
Epoch 4/10, Average Loss: 0.14707822467877257
Epoch 5 - Avg Loss: 0.4556: 100%|██████████| 1610/1610 [04:46<00:00, 5.63it/s]
Epoch 5/10, Average Loss: 0.1138881599183501
Epoch 6 - Avg Loss: 0.3580: 100%|██████████| 1610/1610 [04:46<00:00, 5.61it/s]
Epoch 6/10, Average Loss: 0.08950460070963973
Epoch 7 - Avg Loss: 0.2943: 100%|██████████| 1610/1610 [04:46<00:00, 5.63it/s]
Epoch 7/10, Average Loss: 0.0735697320286439
Epoch 8 - Avg Loss: 0.2472: 100%|██████████| 1610/1610 [04:45<00:00, 5.63it/s]
Epoch 8/10, Average Loss: 0.06179665763463293
Epoch 9 - Avg Loss: 0.2132: 100%|██████████| 1610/1610 [04:45<00:00, 5.63it/s]
Epoch 9/10, Average Loss: 0.053311259651221105
Epoch 10 - Avg Loss: 0.1877: 100%|██████████| 1610/1610 [04:45<00:00, 5.63it/s]
Epoch 10/10, Average Loss: 0.04693056394617935
```

- 1차 시도 때보다 average loss가 0.0174 낮아짐
- 리더보드 제출 결과 : 0.5966 -> 0.6105로 높아짐

결과 분석 및 평가 지표

하이퍼파라미터 튜닝 및 데이터 증강 후

Epoch 1 - Avg Loss: 2.2926: 100%|██████████| 25152/25152 [18:02<00:00, 23.24it/s]
Epoch 1/10, Average Loss: 2.2925947730585876
Epoch 2 - Avg Loss: 1.2044: 100%|██████████| 25152/25152 [18:08<00:00, 23.10it/s]
Epoch 2/10, Average Loss: 1.204397246482102
Epoch 3 - Avg Loss: 0.8654: 100%|██████████| 25152/25152 [18:22<00:00, 22.82it/s]
Epoch 3/10, Average Loss: 0.865386580531906
Epoch 4 - Avg Loss: 0.6960: 100%|██████████| 25152/25152 [18:15<00:00, 22.97it/s]
Epoch 4/10, Average Loss: 0.6959855962807132
Epoch 5 - Avg Loss: 0.5907: 100%|██████████| 25152/25152 [18:19<00:00, 22.87it/s]
Epoch 5/10, Average Loss: 0.5907449588482031
Epoch 6 - Avg Loss: 0.5154: 100%|██████████| 25152/25152 [18:22<00:00, 22.80it/s]
Epoch 6/10, Average Loss: 0.5154345666593265
Epoch 7 - Avg Loss: 0.4568: 100%|██████████| 25152/25152 [18:32<00:00, 22.60it/s]
Epoch 7/10, Average Loss: 0.4567984238119722
Epoch 8 - Avg Loss: 0.4107: 100%|██████████| 25152/25152 [18:21<00:00, 22.83it/s]
Epoch 8/10, Average Loss: 0.41068613466069265
Epoch 9 - Avg Loss: 0.3721: 100%|██████████| 25152/25152 [18:29<00:00, 22.66it/s]
Epoch 9/10, Average Loss: 0.37213897139545526
Epoch 10 - Avg Loss: 0.3402: 100%|██████████| 25152/25152 [18:23<00:00, 22.80it/s]
Epoch 10/10, Average Loss: 0.34020161796118303

- 증강한 데이터로 동일하게 튜닝했을 때 리더보드에서 더 좋은 결과를 보임
- 제출결과 : 0.6356

시스템 구현 계획 및 기술적 이슈

계획

- 다양한 하이퍼파라미터 및 파인튜닝 방법을 시도하여 모델을 개선시켜 loss를 줄이고 모델의 성능을 더 높일 예정

기술적 이슈

- 증강된 데이터로 학습하는데 시간이 매우 오래 걸려 epoch를 키워서 테스트하는데 어려움을 겪고 있음 이로 인해 하이퍼파라미터 조정이 어려움
- 증강한 데이터의 정제 작업

결론 및 제언

- 하이퍼파라미터 튜닝 및 데이터를 증강한 결과 1차 시도보다 높은 성능을 보임
- 증강된 데이터로 학습할 때에 학습시간이 너무 오래 걸려 우선 제공된 train 데이터로 여러 번 하이퍼파라미터를 수정한 후 증강된 데이터에 시도하는 중
- 앞으로 최종 코드 제출일까지 (3/11) 계속해서 파인튜닝을 진행할 예정