## LSTM을 이용한 Text 속 감정 분류





## Park Ju Ho:An undergraduate in statistics Jung Sang Hoon: Mentor

Conclusion

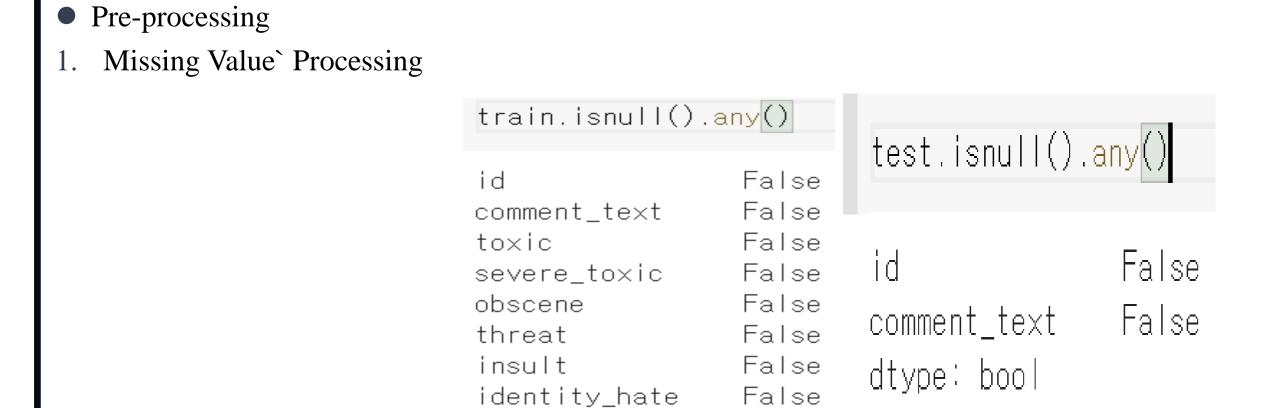
실히 획득함

에 대하여 고민 할 수 있음

## 1. INTRODUCTION

- 이 연구는 RNN의 한 종류인 LSTM을 이용하여 영어 자연어를 처리하는 기본적인 방법을 보여주고 있다.
- 과정은 크게 전처리와 모델 학습으로 나뉘고 최종적으로 loss와 accuracy를 확인한다.
- 전처리는 결측치처리,데이터 분리,토큰화,인덱싱,패딩의 순서로 이루어진다.
- 모델은 한 층의 RNN으로 가장 간단한 LSTM모델을 구현하였다.

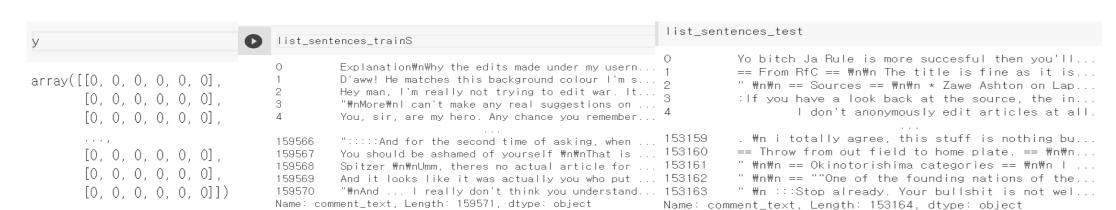
## 2. METHODS



dtype: bool

Train Data와 Test Data 모두 결측치가 없으므로 결측치를 처리 할 필요는 없음

2. Data Isolation



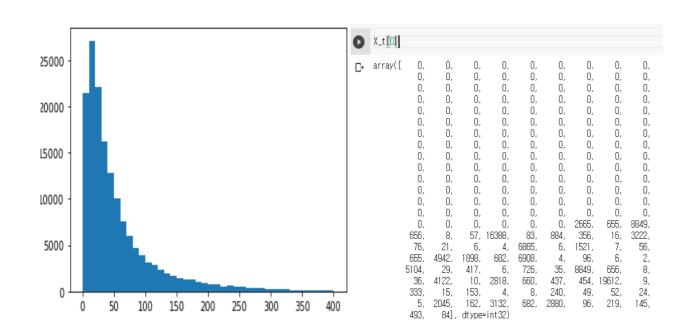
Y=감정 유무의 이진 데이터, train & test = text데이터

3. Tokenization & Indexing

print(list\_tokenized\_train[:1]) 75, 1, 126, 130, 177, 29, 672, 4511, 12052, 1116,

문장을 토큰화 시킨 후 각 단어를 인덱싱 시켜 문장을 인덱싱 번호로 바꾼 작업

4. Padding



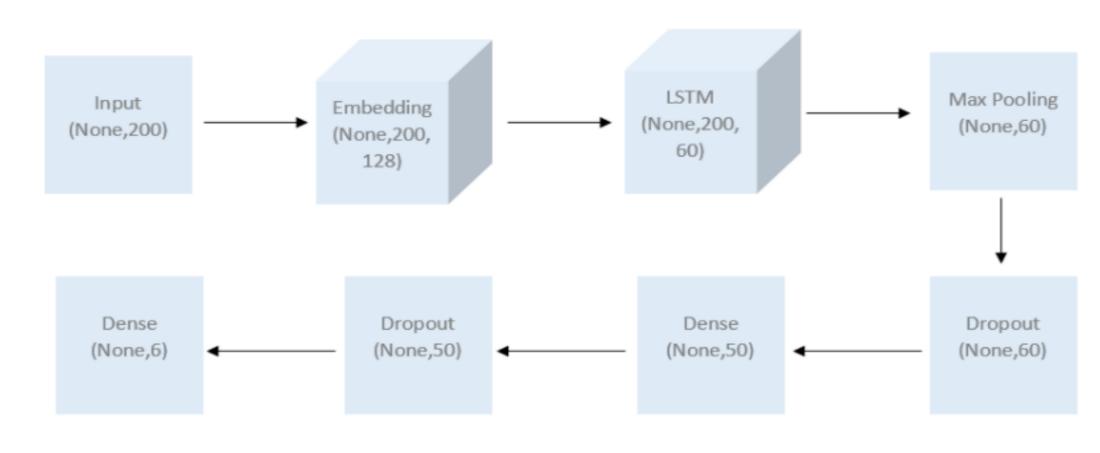
1. Maxlen을 40~50(그래프상 평균)으로 잡는다 => 정보가 손실됨(효율성 우선)

2. Maxlen을 더 크게 잡는다 => 학습이 오래 걸림(정확도 우선)

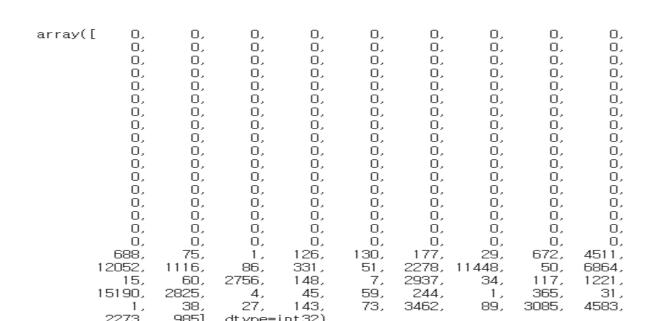
여기선 정확도를 좀 더 높이기 위하여 2번을 선택하여 maxlen을 200으로 잡고 padding을 진행시킴

Modeling

0. 모델링 설계도

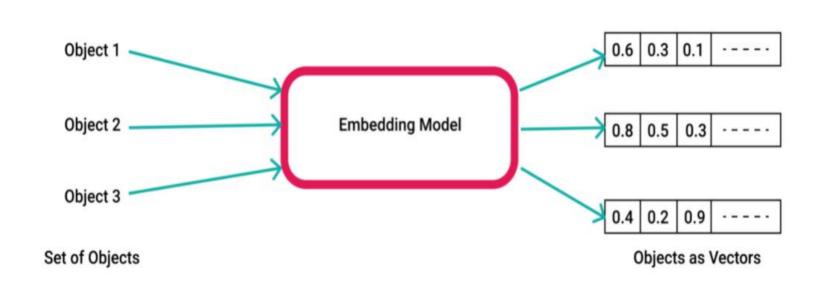


1. Input



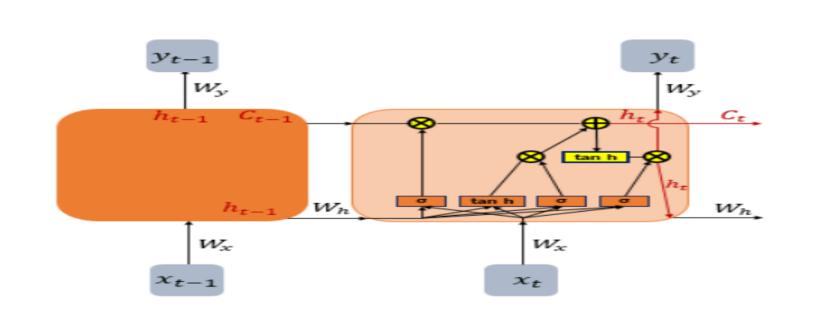
Maxlen의 길이가 200이므로 input의 차원은 200으로 설정함

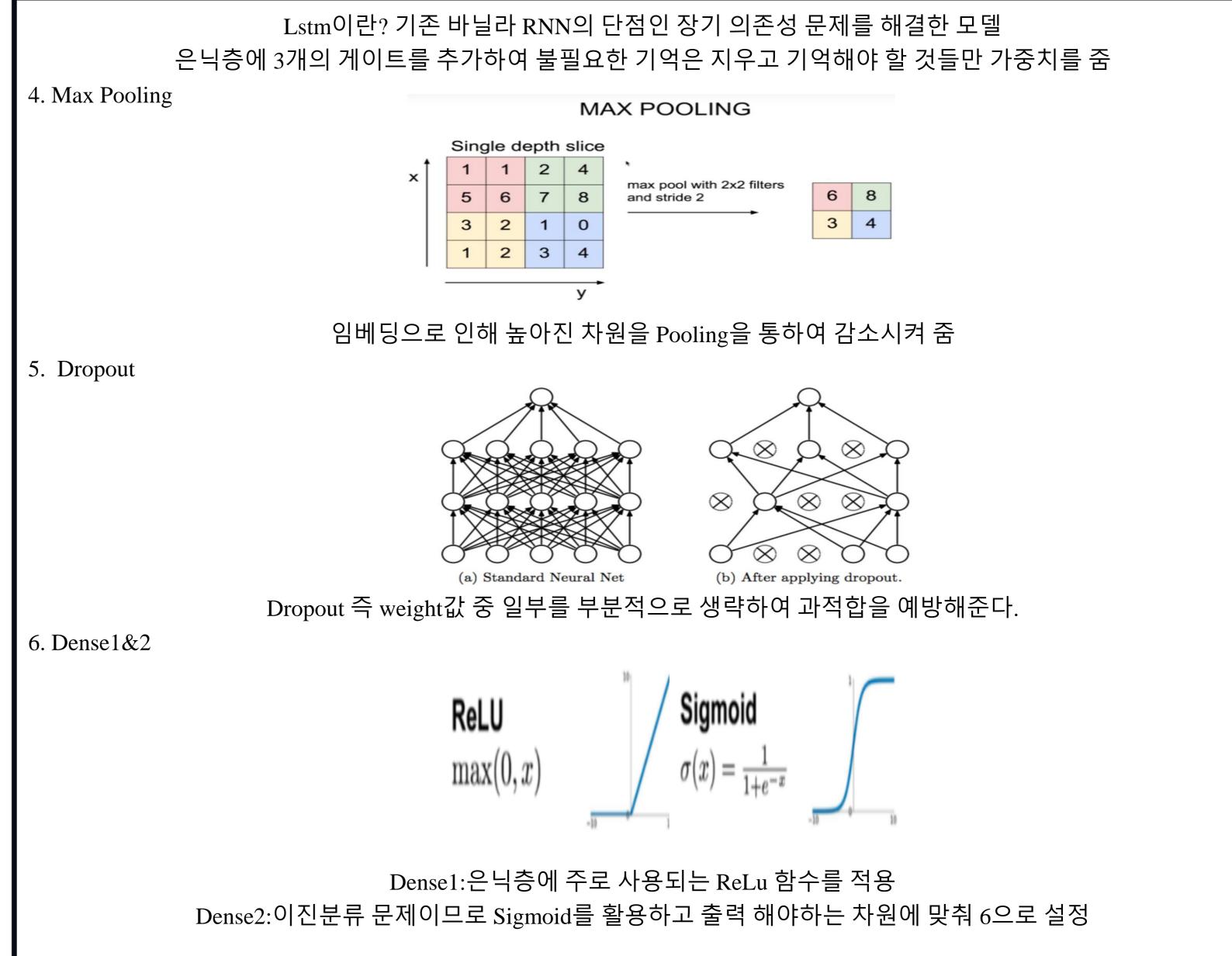
2. Embedding

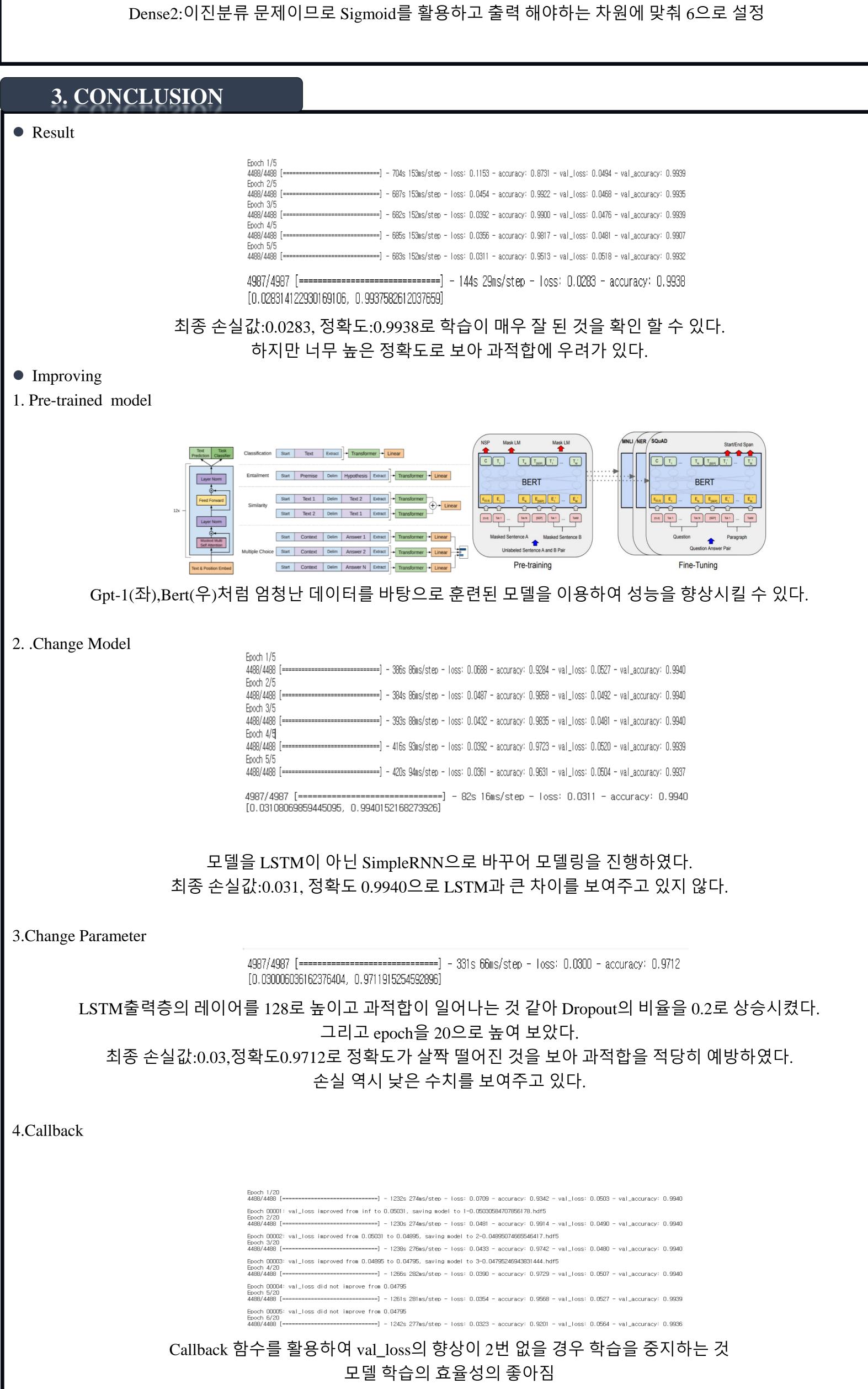


One-hot 인코딩 된 벡터를 임베딩함으로써 One-hot의 단점인 높은 차원을 해소할 수 있음 뿐만 아니라 임베딩 된 벡터를 이용 벡터의 성질(거리 등)을 그 대로 적용 시킬 수 있는 장점이 있음

3. LSTM







이 연구를 통하여 가장 간단한 형태의 자연어 분류 모델을 만들어 봄으로써 모델의 자연어 처리 과정과 및 향상 방안

고민한 향상 방안은 크게 4가지이며 가장 좋은 방법은 Pre-trained model로 예상되나 장비의 문제로 실행하지 못 함

그것을 제외하곤 변수 변환 및 모델 변환을 적용해 보았으나 큰 성과를 거두진 못하였고 callback은 통한 효율성은 확