



## 1. 프로젝트 목표

- ❖ 수업시간에 배운 것을 의미 있는 실습을 통해 복습하자
- ❖ 자연어처리와 ML 과의 경계를 확인하자.

# 실습순서

1. 네이버 영화평 200,000개 에서 출현단어의 빈도수 확인
2. 학습을 시킨 후 긍정과 부정에 대한 문장입력으로 분류성능 확인
3. 네이버 영화평 200,000개에서 조사와 문장부호를 제거한 상태에서 출현단어의 빈도수 확인
4. 학습을 시킨 후 긍정과 부정에 대한 문장입력으로 분류성능 확인
5. 100% 분류를 통한 긍정 및 부정 단어사전 작성



## 2. 실습 – 학습의 방법 및 분류성능

❖ Keras를 이용한 학습 (Dense와 유닛 그리고 활성화 함수의 이해)

```
from tensorflow.keras import models
from tensorflow.keras import layers
from tensorflow.keras import optimizers
from tensorflow.keras import losses
from tensorflow.keras import metrics

model = models.Sequential()
model.add(layers.Dense(64, activation='relu', input_shape=(10000,)))
model.add(layers.Dense(64, activation='relu'))
model.add(layers.Dense(1, activation='sigmoid'))

model.compile(optimizer=optimizers.RMSprop(lr=0.001),
              loss=losses.binary_crossentropy,
              metrics=[metrics.binary_accuracy])

model.fit(x_train, y_train, epochs=10, batch_size=512)
results = model.evaluate(x_test, y_test)
```

```
predict_pos_neg("올해 최고의 영화! 세 번 넘게 봐도 질리지 않네요.")
predict_pos_neg("배경 음악이 영화의 분위기와 너무 안 맞았습니다. 몰입에 방해가 됩니다.")
predict_pos_neg("주연 배우가 신인인데 연기를 진짜 잘 하네요. 몰입감 ㅎㄷㄷ")
predict_pos_neg("믿고 보는 감독이지만 이번에는 아니네요")
predict_pos_neg("주연배우 때문에 봤어요")
```

[올해 최고의 영화! 세 번 넘게 봐도 질리지 않네요.]는 99.22% 확률로 긍정 리뷰이지 않을까 추측해봅니다.^^

[배경 음악이 영화의 분위기와 너무 안 맞았습니다. 몰입에 방해가 됩니다.]는 91.84% 확률로 부정 리뷰이지 않을까 추측해봅니다.^^;

[주연 배우가 신인인데 연기를 진짜 잘 하네요. 몰입감 ㅎㄷㄷ]는 96.14% 확률로 긍정 리뷰이지 않을까 추측해봅니다.^^

[믿고 보는 감독이지만 이번에는 아니네요]는 68.46% 확률로 부정 리뷰이지 않을까 추측해봅니다.^^;

[주연배우 때문에 봤어요]는 69.67% 확률로 부정 리뷰이지 않을까 추측해봅니다.^^;

Dense layer는 입력과 출력을 모두 연결하는 기능을 합니다. 예를 들어 입력 뉴런이 8개이고 출력이 8개 있다면 총 연결선은 64개 ( $8 \times 8 = 64$ ) 입니다.  
이렇게 입력뉴런과 출력뉴런을 모두 연결한다고 해서 전결합층이라고 불리고, 케라스에서는 Dense 라는 클래스로 구현되어 있습니다.

[https://tykimos.github.io/2017/01/27/MLP Layer Talk/](https://tykimos.github.io/2017/01/27/MLP%20Layer%20Talk/)

Sigmoid

입력값을 0과 1 사이의 값으로 변환하여 출력 (이진분류에 적합)하며, 레이어가 깊어질수록 적용하기 어려운 문제가 있다.

ReLU (rectified linear unit)

0이하의 값은 다음 레이어에 전달하지 않고, 0이상의 값은 그대로 출력하며, 주로 CNN을 학습 시킬 때 많이 사용됩니다.

<https://yeomko.tistory.com/39>

## 2. 실습 – 학습의 방법 및 분류성능

❖ Keras를 이용한 학습 (Dense와 유닛 그리고 활성화 함수의 이해)

```
model = models.Sequential()
model.add(layers.Dense(64, activation='relu', input_shape=(10000,)))
model.add(layers.Dense(64, activation='relu'))
model.add(layers.Dense(1, activation='sigmoid'))

model.compile(optimizer=optimizers.RMSprop(lr=0.001),
              loss=losses.binary_crossentropy,
              metrics=[metrics.binary_accuracy])

model.fit(x_train, y_train, epochs=10, batch_size=512)
results = model.evaluate(x_test, y_test)

Epoch 1/10
293/293 [=====] - 7s 23ms/step - loss: 0.3856 - binary_accuracy: 0.8352
Epoch 2/10
293/293 [=====] - 6s 21ms/step - loss: 0.3158 - binary_accuracy: 0.8644
Epoch 3/10
293/293 [=====] - 6s 21ms/step - loss: 0.2924 - binary_accuracy: 0.8778
Epoch 4/10
293/293 [=====] - 6s 20ms/step - loss: 0.2724 - binary_accuracy: 0.8885
Epoch 5/10
293/293 [=====] - 6s 20ms/step - loss: 0.2535 - binary_accuracy: 0.8975
Epoch 6/10
293/293 [=====] - 6s 20ms/step - loss: 0.2338 - binary_accuracy: 0.9073
Epoch 7/10
293/293 [=====] - 6s 20ms/step - loss: 0.2144 - binary_accuracy: 0.9170
Epoch 8/10
293/293 [=====] - 6s 20ms/step - loss: 0.1959 - binary_accuracy: 0.9252
Epoch 9/10
293/293 [=====] - 6s 20ms/step - loss: 0.1793 - binary_accuracy: 0.9321
Epoch 10/10
293/293 [=====] - 6s 20ms/step - loss: 0.1647 - binary_accuracy: 0.9380
1563/1563 [=====] - 4s 2ms/step - loss: 0.4180 - binary_accuracy: 0.8520

> print MI
results
[0.4179568298718449, 0.852080252532959]

> print MI

def predict_pos_neg(review):
    token = tokenize(review)
    tf = tf.nn.embedding_lookup(embeddings, token)
    data = np.expand_dims(np.asarray(tf).astype('float32'), axis=0)
    score = float(model.predict(data))
    if(score > 0.5):
        print("{}는 {:.2f}% 확률로 긍정 리뷰이지 않을까 추측해봅니다.^^\n".format(review, score * 100))
    else:
        print("{}는 {:.2f}% 확률로 부정 리뷰이지 않을까 추측해봅니다.^^\n".format(review, (1 - score) * 100))

> print MI

predict_pos_neg("올해 최고의 영화! 세 번 넘게 봐도 질리지 않네요.")
predict_pos_neg("배경 음악이 영화의 분위기를 너무 안 맞았습니다. 물집에 방해가 됩니다.")
predict_pos_neg("주연 배우가 신인인데 연기를 진짜 잘 하네요. 물입감 ㅎㅎ")
predict_pos_neg("맞고 보는 감쪽이지만 이번에는 아니네요")
predict_pos_neg("주연배우 때문에 봤어요")

[올해 최고의 영화! 세 번 넘게 봐도 질리지 않네요.]는 98.86% 확률로 긍정 리뷰이지 않을까 추측해봅니다.^^
[배경 음악이 영화의 분위기를 너무 안 맞았습니다. 물집에 방해가 됩니다.]는 73.25% 확률로 부정 리뷰이지 않을까 추측해봅니다.^^
[주연 배우가 신인인데 연기를 진짜 잘 하네요. 물입감 ㅎㅎ]는 99.71% 확률로 긍정 리뷰이지 않을까 추측해봅니다.^^
[맞고 보는 감쪽이지만 이번에는 아니네요]는 76.21% 확률로 긍정 리뷰이지 않을까 추측해봅니다.^^
[주연배우 때문에 봤어요]는 94.73% 확률로 부정 리뷰이지 않을까 추측해봅니다.^^
```

시간이 오래 걸리지만 input\_shape가 추측확률을 결정짓는 요소가 되기 때문에 10,000을 유지

확률 0.85 로써 조심스러운 수치가긴 하지만 단어사전 구분 용 85%에서 다시 80%를 선택하는 작업

비교용 데이터

```
# 명사, 동사, 형용사, 부사 등 의미있는 단어를 남긴다.
okt = Okt()
def str_filter(review):
    ....str_result = ""
    ....for t in okt.pos(review, norm=True, stem=True):
    ....    if t[1] == "Josa" or t[1] == "Punctuation":
    ....        str_result = str_result
    ....    else:
    ....        str_result = str_result + "." + t[0]
    ....return str_result
```

## 2. 실습 – 동일한 학습과 분류성능

❖ Keras를 이용한 학습 (Dense와 유닛 그리고 활성화 함수의 이해)

```
from tensorflow.keras import models
from tensorflow.keras import layers
from tensorflow.keras import optimizers
from tensorflow.keras import losses
from tensorflow.keras import metrics

model = models.Sequential()
model.add(layers.Dense(64, activation='relu', input_shape=(10000,)))
model.add(layers.Dense(64, activation='relu'))
model.add(layers.Dense(1, activation='sigmoid'))

model.compile(optimizer=optimizers.RMSprop(lr=0.001),
              loss=losses.binary_crossentropy,
              metrics=[metrics.binary_accuracy])

model.fit(x_train, y_train, epochs=10, batch_size=512)
results = model.evaluate(x_test, y_test)
```

조사를 제거하면 조금 더 긍정에 가까운 확률이 나옴

```
predict_pos_neg("올해 최고의 영화! 세 번 넘게 봐도 질리지 않네요.")
predict_pos_neg("배경 음악이 영화의 분위기랑 너무 안 맞았습니다. 물입에 방해가 됩니다.")
predict_pos_neg("주연 배우가 신인인데 연기를 진짜 잘 하네요. 물입감 ㅎㅎ")
predict_pos_neg("밀고 보는 감독이지만 이번에는 아니네요")
predict_pos_neg("주연배우 때문에 봤어요")
```

[올해 최고의 영화! 세 번 넘게 봐도 질리지 않네요.]는 99.53% 확률로 긍정 리뷰이지 않을까 추측해봅니다.^^

[배경 음악이 영화의 분위기랑 너무 안 맞았습니다. 물입에 방해가 됩니다.]는 92.94% 확률로 부정 리뷰이지 않을까 추측해봅니다.^^;

[주연 배우가 신인인데 연기를 진짜 잘 하네요. 물입감 ㅎㅎ]는 99.48% 확률로 긍정 리뷰이지 않을까 추측해봅니다.^^

[밀고 보는 감독이지만 이번에는 아니네요]는 58.15% 확률로 부정 리뷰이지 않을까 추측해봅니다.^^;

[주연배우 때문에 봤어요]는 74.49% 확률로 부정 리뷰이지 않을까 추측해봅니다.^^;

## 2. 실습 – 동일한 학습과 분류성능

### ❖ 조사가 없을 경우의 분류성능

```
model = models.Sequential()  
model.add(layers.Dense(64, activation='relu', input_shape=(10000,)))  
model.add(layers.Dense(64, activation='relu'))  
model.add(layers.Dense(1, activation='sigmoid'))  
  
model.compile(optimizer=optimizers.RMSprop(lr=0.001),  
              loss=losses.binary_crossentropy,  
              metrics=[metrics.binary_accuracy])  
  
model.fit(x_train, y_train, epochs=10, batch_size=512)  
results = model.evaluate(x_test, y_test)  
  
Epoch 1/10  
293/293 [=====] - 7s 22ms/step - loss: 0.3972 - binary_accuracy: 0.8266  
Epoch 2/10  
293/293 [=====] - 7s 23ms/step - loss: 0.3272 - binary_accuracy: 0.8603  
Epoch 3/10  
293/293 [=====] - 7s 22ms/step - loss: 0.3055 - binary_accuracy: 0.8717  
Epoch 4/10  
293/293 [=====] - 7s 22ms/step - loss: 0.2893 - binary_accuracy: 0.8811  
Epoch 5/10  
293/293 [=====] - 7s 22ms/step - loss: 0.2726 - binary_accuracy: 0.8890  
Epoch 6/10  
293/293 [=====] - 7s 22ms/step - loss: 0.2557 - binary_accuracy: 0.8971  
Epoch 7/10  
293/293 [=====] - 7s 23ms/step - loss: 0.2391 - binary_accuracy: 0.9046  
Epoch 8/10  
293/293 [=====] - 7s 22ms/step - loss: 0.2225 - binary_accuracy: 0.9117  
Epoch 9/10  
293/293 [=====] - 7s 22ms/step - loss: 0.2068 - binary_accuracy: 0.9193  
Epoch 10/10  
293/293 [=====] - 7s 22ms/step - loss: 0.1922 - binary_accuracy: 0.9256  
1563/1563 [=====] - 5s 3ms/step - loss: 0.4393 - binary_accuracy: 0.8380
```

▶ \*분 ML  
results

[0.4392968714237213, 0.838039942398071]

▶ \*분 ML

```
def predict_pos_neg(review):  
    token = tokenize(review)  
    tf = term_frequency(token)  
    data = np.expand_dims(np.asarray(tf).astype('float32'), axis=0)  
    score = float(model.predict(data))  
    if(score > 0.5):  
        print("{}")는 {:.2f}% 확률로 긍정 리뷰이지 않을까 추측해봅니다.^^\n".format(review, score * 100))  
    else:  
        print("{}")는 {:.2f}% 확률로 부정 리뷰이지 않을까 추측해봅니다.^^;\n".format(review, (1 - score) * 100))
```

▶ \*분 ML

```
predict_pos_neg("올해 최고의 영화! 세 번 넘게 봐도 질리지га 않네요.")  
predict_pos_neg("배경 음악이 영화의 분위기를 너무 안 맞았습니다. 물임에 방해가 됩니다.")  
predict_pos_neg("주연 배우가 신인인데 연기를 진짜 잘 하네요. 물임값 * c c")  
predict_pos_neg("믿고 보는 감독이지만 이번에는 아니네요")  
predict_pos_neg("주연배우 때문에 봤어요")
```

[올해 최고의 영화! 세 번 넘게 봐도 질리지га 않네요.]는 99.53% 확률로 긍정 리뷰이지 않을까 추측해봅니다.^^

[배경 음악이 영화의 분위기를 너무 안 맞았습니다. 물임에 방해가 됩니다.]는 92.94% 확률로 부정 리뷰이지 않을까 추측해봅니다.^^;

[주연 배우가 신인인데 연기를 진짜 잘 하네요. 물임값 \* c c]는 99.48% 확률로 긍정 리뷰이지 않을까 추측해봅니다.^^

[믿고 보는 감독이지만 이번에는 아니네요]는 58.15% 확률로 부정 리뷰이지 않을까 추측해봅니다.^^;

[주연배우 때문에 봤어요]는 74.49% 확률로 부정 리뷰이지 않을까 추측해봅니다.^^;

→ 확률 0.84 로써 반복되는 조사가 없어지니 분류의 정확률이 떨어짐

→ 긍정적인 확률이 근소한 차이로 올라감



### 3. 100% 긍정, 부정 단어 추출 (단어사전 작성)

- ❖ nltk.Textvocabulary().most\_common를 이용한 출현 빈도별 토큰분석
- ❖ 명사, 동사, 형용사, 부사를 중심으로 결과비교
- ❖ 80% 이상 긍정단어와 80%이상 부정단어 추출하여 사전 파일 작성

80% 이상 긍정적인 단어사전  
(조사가 포함된 파일로 제작)

명사  
동사  
형용사  
부사  
조사  
문장부호 등



80% 이상 부정적인 단어 사전  
(조사가 제거된 파일로 제작)

### 3. 100% 긍정, 부정 단어 추출 (단어사전 작성)

- ❖ nltk.Textvocabulary().most\_common를 이용한 출현 빈도별 토큰분석
- ❖ 명사, 동사, 형용사, 부사를 중심으로 결과비교
- ❖ 80% 이상 긍정단어와 80%이상 부정단어 추출하여 사전 파일 작성

```
def make_dictionary(word):
    token = tokenize(word)
    tf = term_frequency(token)
    data = np.expand_dims(np.asarray(tf).astype('float32'), axis=0)
    score = float(model.predict(data))
    if(score > 0.8):
        score = score * 100
        spoint = str(score)
        strReturn = token[0] + "/" + spoint
    else:
        strReturn = "fault"

    return strReturn

# print (len(set(text.tokens)))
checkwords = set(text.tokens)
arrcheckwords = list(checkwords)
# print ((arrcheckwords[0].split("/"))[0])
positive_dictionary = open("./data/positive_dictionary.txt", "w", encoding="utf-8")
for t in arrcheckwords:
    result = make_dictionary((t.split("/"))[0])
    if result != "fault":
        positive_dictionary.write(result+"\n")
positive_dictionary.close()
```

명사  
동사  
형용사  
부사  
조사  
문장부호 등

```
def make_dictionary(word):
    token = tokenize(word)
    tf = term_frequency(token)
    data = np.expand_dims(np.asarray(tf).astype('float32'), axis=0)
    score = float(model.predict(data))
    if(score < 0.2):
        score = (1-score) * 100
        spoint = str(score)
        strReturn = token[0] + "/" + spoint
    else:
        strReturn = "fault"

    return strReturn

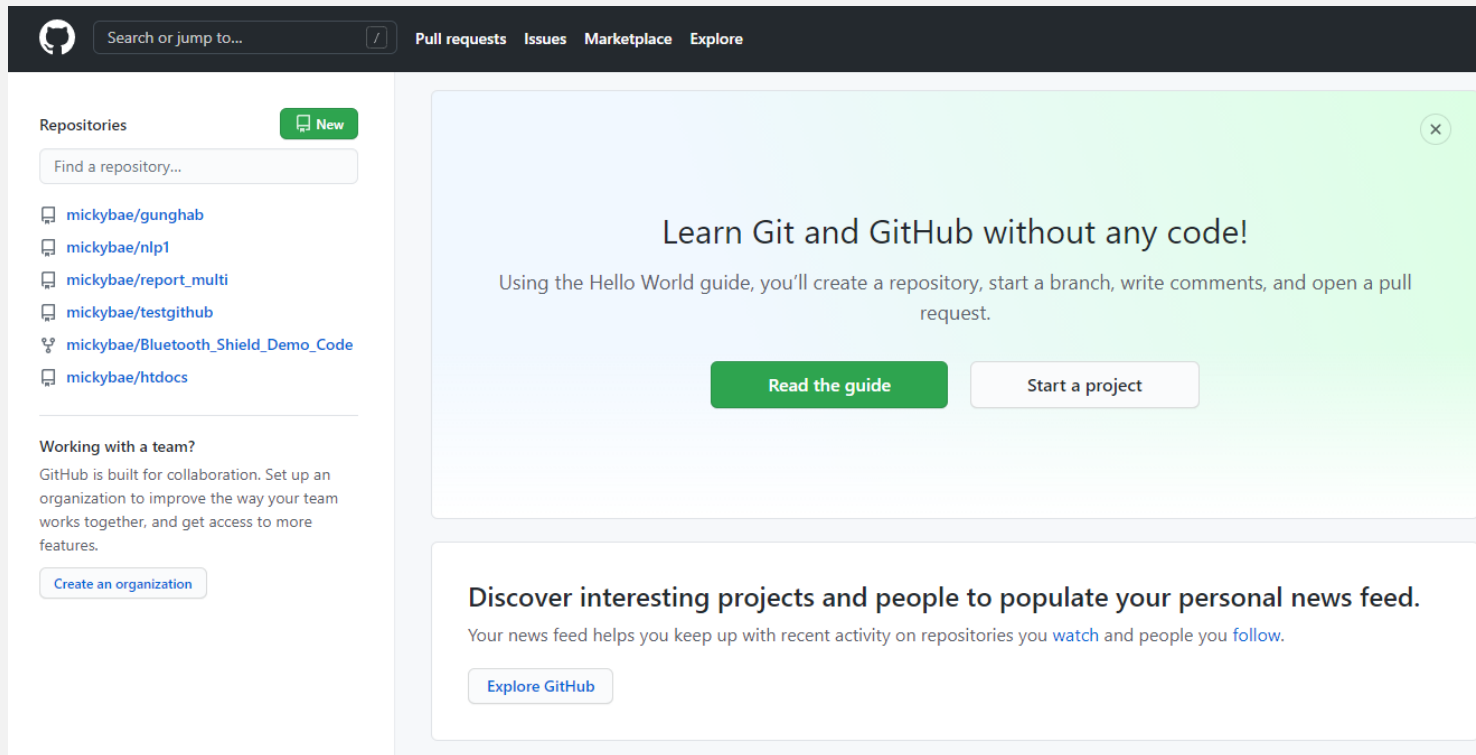
# print (len(set(text.tokens)))
checkwords = set(text.tokens)
arrcheckwords = list(checkwords)
# print ((arrcheckwords[0].split("/"))[0])
negative_dictionary = open("./data/negative_dictionary.txt", "w", encoding="utf-8")
for t in arrcheckwords:
    result = make_dictionary((t.split("/"))[0])
    if result != "fault":
        negative_dictionary.write(result+"\n")
negative_dictionary.close()
```

positive\_dictionary.txt

negative\_dictionary.txt

## 4. 사용소스 및 데이터 공개

❖ Git-hub 사용하기



## 온라인 공유 게시

- 프로젝트 개요
- 사용소스
- 원본데이터
- 긍정사전
- 부정사전

감사합니다.