

IMDB - Binary Classification

NLP(Natural Language Processing)

Import Tensorflow & Keras

```
import warnings
warnings.filterwarnings('ignore')
```

- TensorFlow '1.x' Version 지정

```
%tensorflow_version 1.x
import tensorflow as tf

tf.__version__
```

```
TensorFlow 1.x selected.
'1.15.2'
```

- GPU 설정 확인

```
tf.test.gpu_device_name()

'/device:GPU:0'
```

- GPU 정보 확인

```
!nvidia-smi
```

```
Wed Mar 17 23:31:26 2021
```

| NVIDIA-SMI 460.56 Driver Version: 460.32.03 CUDA Version: 11.2 | | | | | | | | | |
|----------------------------------------------------------------|----------|---------------|------------------|-------------------|----------------------|------------|--|--|--|
| GPU | Name | Persistence-M | Bus-Id | Disp.A | Volatile Uncorr. ECC | | | | |
| Fan | Temp | Perf | Pwr:Usage/Cap | Memory-Usage | GPU-Util | Compute M. | | | |
| | | | | | | MIG M. | | | |
| 0 | Tesla T4 | Off | 00000000:00:04.0 | Off | 0 | | | | |
| N/A | 36C | P0 | 26W / 70W | 104MiB / 15109MiB | 0% | Default | | | |
| | | | | | | N/A | | | |

| Processes: | | | | | | | |
|------------|----|----|-----|------|--------------|------------|--|
| GPU | GI | CI | PID | Type | Process name | GPU Memory | |
| ID | ID | | | | | Usage | |
| | | | | | | | |

```
|=====|
|-----|
```

- Keras Version 확인

```
import keras
```

```
keras.__version__
```

```
Using TensorFlow backend.
'2.3.1'
```

▼ I. IMDB Data_Set Load & Review

▼ 1) Load IMDB Data_Set

- Word to Vector
- 전체 데이터 내에서 단어의 사용빈도에 따라 인덱스화
- 정수 인덱스 '11'은 11번째로 자주 사용된 단어를 나타냄
- num_words = 10000 : 인덱스 값 10000 이하의 단어만 추출
- 단어 인덱스 값이 10000을 넘지 않는 단어만 분석에 사용

```
from keras.datasets import imdb
```

```
(train_data, train_labels), (test_data, test_labels) = imdb.load_data(num_words = 10000)
```

```
Downloading data from https://s3.amazonaws.com/text-datasets/imdb.npz
17465344/17464789 [=====] - 1s 0us/step
```

```
max(max(W) for W in train_data)
```

```
9999
```

▼ 2) Visualization & Frequency(Optional)

- x - Histogram(리뷰 길이)

```
import matplotlib.pyplot as plt
```

```
print('리뷰 최대 길이 :', max(len(L) for L in train_data))
print('리뷰 평균 길이 :', sum(map(len, train_data))/len(train_data))
```

```
plt.figure(figsize = (9, 6))
plt.hist([len(L) for L in train_data], bins = 50)
plt.xlabel('Length of train_data')
```

```
plt.ylabel('Number of train_data')
plt.show()
```

- y - Frequency(0:부정, 1:긍정)

```
import numpy as np

unique_elements, counts_elements = np.unique(train_labels, return_counts = True)

print('Label 빈도수:')
print(np.asarray((unique_elements, counts_elements)))

Label 빈도수:
[[ 0  1]
 [12500 12500]]
```

▼ 3) Data Structure Review(Optional)

```
# 전체 train_data 개수
print(len(train_data))

# 첫번째 train_data 정보
print(len(train_data[0]))
print(train_data[0][0:10])
print(train_data[0].count(4))
print(train_labels[0])

25000
218
[1, 14, 22, 16, 43, 530, 973, 1622, 1385, 65]
15
1
```

▼ 4) Vector to Word(Optional)

- `get_word_index()`: 단어와 인덱스를 매핑한 사전
- 0, 1, 2: '패딩', '문서 시작', '사전에 없음'

```
word_index = imdb.get_word_index()

print(word_index)
```

- 인덱스와 단어 위치 변경

```
reverse_word_index = dict([(value, key) for (key, value) in word_index.items()])
```

```
print(reverse_word_index)
```

- 0번 영화 리뷰 디코딩(1:긍정)

```
decoded_review = ' '.join([reverse_word_index.get(i - 3, '?') for i in train_data[0]])
```

```
print(decoded_review)
print(train_labels[0])
```

```
? this film was just brilliant casting location scenery story direction everyone's really sui
1
```

- 1번 영화리뷰 디코딩(0:부정)

```
decoded_review = ' '.join([reverse_word_index.get(i - 3, '?') for i in train_data[1]])
```

```
print(decoded_review)
print(train_labels[1])
```

```
? big hair big boobs bad music and a giant safety pin these are the words to best describe th
0
```

▼ II. Tensor Transformation

▼ 1) X_train & X_test : (25000, 10000)

- `vectorize_sequence()` 정의
- 크기는 10000이고 모든 원소가 0인 행렬 생성
 - `np.zeros(len(sequences), dimension))`
- 값이 존재하는 인덱스의 위치를 1로 지정
 - `enumerate()`
 - `results[i, sequence] = 1.0`

```
import numpy as np

def vectorize_sequences(sequences, dimension = 10000):
    results = np.zeros((len(sequences), dimension))
    for i, sequence in enumerate(sequences):
        results[i, sequence] = 1.0
    return results
```

- `enumerate()` - Example

```
r = np.zeros((5, 10))
v = [1, 3, 5, 7, 9]

for i, v in enumerate(v):
    r[i, v] = 1.0

print(r)

[[0.  1.  0.  0.  0.  0.  0.  0.  0.  0.]
 [0.  0.  0.  1.  0.  0.  0.  0.  0.  0.]
 [0.  0.  0.  0.  0.  1.  0.  0.  0.  0.]
 [0.  0.  0.  0.  0.  0.  0.  1.  0.  0.]
 [0.  0.  0.  0.  0.  0.  0.  0.  1.  0.]]
```

- `vectorize_sequence()` 적용

```
X_train = vectorize_sequences(train_data)
X_test = vectorize_sequences(test_data)

X_train.shape, X_test.shape

((25000, 10000), (25000, 10000))
```

- Transformation Check

```
print(X_train[0][:21])
print(X_train[0][9979:])

print(X_test[0][:21])
print(X_test[0][9979:])

[0.  1.  1.  0.  1.  1.  1.  1.  1.  1.  0.  0.  1.  1.  1.  1.  1.  1.  1.  0.]
[0.  0.  0.  0.  0.  0.  0.  0.  0.  0.  0.  0.  0.  0.  0.  0.  0.  0.  0.  0.]
[0.  1.  1.  0.  1.  1.  1.  1.  1.  1.  1.  0.  0.  1.  1.  0.  1.  0.  0.  0.]
[0.  0.  0.  0.  0.  0.  0.  0.  0.  0.  0.  0.  0.  0.  0.  0.  0.  0.  0.  0.]
```

▼ 2) `y_train` & `y_test`

```
y_train = np.asarray(train_labels).astype(float)
y_test = np.asarray(test_labels).astype(float)

print(y_train[:21])
print(y_test[:21])

[1.  0.  0.  1.  0.  0.  1.  0.  1.  0.  1.  0.  0.  0.  0.  1.  1.  0.  1.  0.]
[0.  1.  1.  0.  1.  1.  1.  0.  0.  1.  1.  0.  0.  0.  1.  0.  1.  0.  0.  1.]
```

▼ 3) Train vs. Validation Split

```
X_valid = X_train[:10000]
partial_X_train = X_train[10000:]

y_valid = y_train[:10000]
partial_y_train = y_train[10000:]

partial_X_train.shape, partial_y_train.shape, X_valid.shape, y_valid.shape

((15000, 10000), (15000,), (10000, 10000), (10000,))
```

▼ III. IMDB Keras Modeling

▼ 1) Model Define

- 모델 신경망 구조 정의

```
from keras import models
from keras import layers

imdb = models.Sequential()
imdb.add(layers.Dense(16, activation = 'relu', input_shape = (10000,)))
imdb.add(layers.Dense(16, activation = 'relu'))
imdb.add(layers.Dense(1, activation = 'sigmoid'))
```

WARNING:tensorflow:From /tensorflow-1.15.2/python3.7/tensorflow_core/python/ops/resource_variable_ops.py:1432: tf.nn.nn_ops.nn_init is deprecated and will be removed in a future version.

Instructions for updating:

If using Keras pass *_constraint arguments to layers.

- 모델 구조 확인

```
imdb.summary()
```

Model: "sequential_1"

| Layer (type) | Output Shape | Param # |
|-----------------|--------------|---------|
| dense_1 (Dense) | (None, 16) | 160016 |
| dense_2 (Dense) | (None, 16) | 272 |
| dense_3 (Dense) | (None, 1) | 17 |

Total params: 160,305

Trainable params: 160,305
Non-trainable params: 0

2) Model Compile

- 모델 학습방법 설정

```
imdb.compile(loss = 'binary_crossentropy',
             optimizer = 'rmsprop',
             metrics = ['accuracy'])
```

WARNING:tensorflow:From /tensorflow-1.15.2/python3.7/tensorflow_core/python/ops/nn_impl.py:18:
Instructions for updating:
Use tf.where in 2.0, which has the same broadcast rule as np.where

3) Model Fit

- 약 1분

```
%%time
```

```
Hist_imdb = imdb.fit(partial_X_train, partial_y_train,
                    epochs = 50,
                    batch_size = 512,
                    validation_data = (X_valid, y_valid))
```

```
Epoch 23/50
15000/15000 [=====] - 1s 68us/step - loss: 0.0059 - accuracy: 0.9
Epoch 24/50
15000/15000 [=====] - 1s 67us/step - loss: 0.0018 - accuracy: 0.9
Epoch 25/50
15000/15000 [=====] - 1s 66us/step - loss: 0.0042 - accuracy: 0.9
Epoch 26/50
15000/15000 [=====] - 1s 66us/step - loss: 0.0011 - accuracy: 0.9
Epoch 27/50
15000/15000 [=====] - 1s 66us/step - loss: 0.0024 - accuracy: 0.9
Epoch 28/50
15000/15000 [=====] - 1s 66us/step - loss: 7.4511e-04 - accuracy:
Epoch 29/50
15000/15000 [=====] - 1s 66us/step - loss: 6.1528e-04 - accuracy:
Epoch 30/50
15000/15000 [=====] - 1s 66us/step - loss: 0.0031 - accuracy: 0.9
Epoch 31/50
15000/15000 [=====] - 1s 67us/step - loss: 3.1138e-04 - accuracy:
Epoch 32/50
15000/15000 [=====] - 1s 66us/step - loss: 2.5170e-04 - accuracy:
Epoch 33/50
15000/15000 [=====] - 1s 66us/step - loss: 0.0022 - accuracy: 0.9
Epoch 34/50
15000/15000 [=====] - 1s 68us/step - loss: 1.5121e-04 - accuracy:
```

```

15000/15000 [-----] - 1s 66us/step - loss: 1.5121e-04 - accuracy:
Epoch 35/50
15000/15000 [=====] - 1s 66us/step - loss: 1.2341e-04 - accuracy:
Epoch 36/50
15000/15000 [=====] - 1s 67us/step - loss: 6.0734e-04 - accuracy:
Epoch 37/50
15000/15000 [=====] - 1s 67us/step - loss: 8.5682e-05 - accuracy:
Epoch 38/50
15000/15000 [=====] - 1s 67us/step - loss: 6.8835e-05 - accuracy:
Epoch 39/50
15000/15000 [=====] - 1s 67us/step - loss: 5.5887e-05 - accuracy:
Epoch 40/50
15000/15000 [=====] - 1s 66us/step - loss: 0.0018 - accuracy: 0.9
Epoch 41/50
15000/15000 [=====] - 1s 66us/step - loss: 3.7081e-05 - accuracy:
Epoch 42/50
15000/15000 [=====] - 1s 67us/step - loss: 2.9623e-05 - accuracy:
Epoch 43/50
15000/15000 [=====] - 1s 66us/step - loss: 2.4563e-05 - accuracy:
Epoch 44/50
15000/15000 [=====] - 1s 66us/step - loss: 1.7579e-05 - accuracy:
Epoch 45/50
15000/15000 [=====] - 1s 66us/step - loss: 8.5200e-04 - accuracy:
Epoch 46/50
15000/15000 [=====] - 1s 67us/step - loss: 1.9188e-05 - accuracy:
Epoch 47/50
15000/15000 [=====] - 1s 67us/step - loss: 1.0265e-05 - accuracy:
Epoch 48/50
15000/15000 [=====] - 1s 69us/step - loss: 8.0907e-06 - accuracy:
Epoch 49/50
15000/15000 [=====] - 1s 67us/step - loss: 6.4136e-06 - accuracy:
Epoch 50/50
15000/15000 [=====] - 1s 67us/step - loss: 8.8788e-04 - accuracy:
CPU times: user 52.3 s, sys: 1.32 s, total: 53.7 s
Wall time: 59.5 s

```

▼ 4) 학습 결과 시각화

- Loss Visualization

```

import matplotlib.pyplot as plt

epochs = range(1, len(Hist_imdb.history['loss']) + 1)

plt.figure(figsize = (9, 6))
plt.plot(epochs, Hist_imdb.history['loss'])
plt.plot(epochs, Hist_imdb.history['val_loss'])
plt.title('Training & Validation Loss')
plt.xlabel('Epochs')
plt.ylabel('Loss')
plt.legend(['Training Loss', 'Validation Loss'])
plt.grid()
plt.show()

```


- Accuracy Visualization

```
import matplotlib.pyplot as plt

epochs = range(1, len(Hist_imdb.history['accuracy']) + 1)

plt.figure(figsize = (9, 6))
plt.plot(epochs, Hist_imdb.history['accuracy'])
plt.plot(epochs, Hist_imdb.history['val_accuracy'])
plt.title('Training & Validation Accuracy')
plt.xlabel('Epochs')
plt.ylabel('Accuracy')
plt.legend(['Training Accuracy', 'Validation Accuracy'])
plt.grid()
plt.show()
```

▼ 5) Model Evaluate

- Loss & Accuracy

```
loss, accuracy = imdb.evaluate(X_test, y_test)

print('Loss = {:.5f}'.format(loss))
print('Accuracy = {:.5f}'.format(accuracy))
```

25000/25000 [=====] - 2s 71us/step
 Loss = 1.75344
 Accuracy = 0.83964

▼ 6) Model Predict

```
imdb.predict_classes(X_test)

array([[0],
       [1],
       [1],
       ...,
       [0],
       [0],
       [0]], dtype=int32)

#

#

#
```

The End

#

#

#