

## 텍스트 데이터 다루기

- tensorflow 버전 : 2.4.1
- python 버전 : 3.8.5
- keras 버전 : 2.4.3

In [3]:

```
import keras
import tensorflow as tf
print(keras.__version__)
print(tf.__version__)
```

2.4.3  
2.4.1

## 텍스트 데이터

- 텍스트는 가장 흔한 시퀀스 형태의 데이터이다.
- 시퀀스 처리용 딥러닝 모델은 문서 분류, 감성 분석, 저자 식별, 질문 응답 등의 애플리케이션에 적합하다.
- 텍스트를 수치형 텐서로 변환하는 과정을 텍스트 벡터화(vectorizing text)라고 한다.
  - 방법1 : 텍스트를 단어로 나누고 각 단어를 하나의 벡터로 변환
  - 방법2 : 텍스트를 문자로 나누고 각 문자를 하나의 벡터로 변환
  - 방법3 : 텍스트에서 단어나 문자의 n-그램(n-gram)추출하여 각 n-그램을 하나의 벡터로 변환
    - n-gram은 연속된 단어나 문자의 그룹이다. 텍스트에서 단어나 문자를 하나씩 이동하면서 추출
- 텍스트를 나누는 이런 단위(단어, 문자, n-그램)를 토큰(token)이라 한다.
  - 토큰으로 나누는 작업을 토큰화(tokenization)이라 한다.

## 텍스트의 벡터화 과정

- 01 어떤 종류의 토큰화 적용
- 02 토큰에 수치형 벡터를 연결
- 03 벡터는 시퀀스 텐서로 묶여져서 심층 신경망의 입력으로 사용

## 토큰과 벡터를 연결하는 방법 2가지

- 원-핫 인코딩(one-hot encoding)
- 토큰 임베딩(token embedding) - 보통 단어에만 사용되면 단어(word embedding)이라 말한다.

## 원핫 인코딩 실습해 보기

- 모든 데이터에 대해서(문장)을 토큰으로 나눈다. 그리고 이를 사전에 하나의 단어와 인덱스로 만든다.
- 단어 사전이 갖는 인덱스 수만큼의 배열을 만든다.(값은 0으로 채운다.) np.zeros이용.
- 단어들을 인덱스 값(0,10)을 원핫 인코딩하여 해당 위치에 1을 채운다.

In [4]:



```
import numpy as np

# 하나의 원소가 샘플. 하나의 문장.
samples = ['The cat sat on the mat.', 'The dog ate my homework.']

# 문장에 있는 단어들을 가지고 인덱스(사전)을 구축함.
# 단어:인덱스
token_index = {}
for sample in samples:
    for word in sample.split():
        if word not in token_index:
            token_index[word] = len(token_index) + 1

max_length = 10
token_index
```

Out[4]:

```
{'The': 1,
 'cat': 2,
 'sat': 3,
 'on': 4,
 'the': 5,
 'mat.': 6,
 'dog': 7,
 'ate': 8,
 'my': 9,
 'homework.': 10}
```

In [5]:



```
# 단어 인덱스의 길이 확인
max(token_index.values()) + 1
```

Out[5]:

11

In [6]:



```
# 결과를 저장할 배열(0으로 이루어진 벡터)
# max_length는 사용할 단어수
results = np.zeros((len(samples), max_length, max(token_index.values()) + 1))
print(results)
```

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

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

In [7]:



```
token_index.get(word)
```

Out[7]:

10

In [8]:



```
for i, sample in enumerate(samples):
    for j, word in list(enumerate(sample.split()))[:max_length]:
        index = token_index.get(word)
        results[i, j, index] = 1. # 행, 열, 인덱스를 1로 치환
```

In [9]:



results

Out[9]:

```
array([[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., 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., 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.],
       [0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.]])
```

## 케라스를 활용한 원핫 인코딩 실습

- 케라스에서는 원본 텍스트 데이터를 단어 또는 문자 수준의 원-핫 인코딩으로 변환해 주는 유틸리티가 있다.

In [10]:



```
from keras.preprocessing.text import Tokenizer

samples = ['The cat sat on the mat.', 'The dog ate my homework.']
```

- tensorflow의 버전이 다를 경우 에러 발생할 수 있음. (tf 2.5.x에서는 에러 발생 가능성 있음.
  - AttributeError: module 'keras.utils.generic\_utils' has no attribute 'populate\_dict\_with\_module\_objects'
  - pip uninstall tensorflow
  - pip install tensorflow==2.4.1

## 2-1 단어들에 대한 인덱스 값 얻기

- fit\_on\_texts(text) # 코퍼스를 입력으로 하면 빈도수를 기준으로 단어 집합 생성.

In [11]:

```
# 가장 빈도가 높은 1,000개의 단어만 선택하도록 Tokenizer 객체를 만듭니다.
tokenizer = Tokenizer(num_words=1000)

# 단어 인덱스를 구축합니다.
tokenizer.fit_on_texts(samples)

# 문자열을 정수 인덱스의 리스트로 변환합니다.
sequences = tokenizer.texts_to_sequences(samples)
sequences
```

Out[11]:

```
[[1, 2, 3, 4, 1, 5], [1, 6, 7, 8, 9]]
```

## 2-2 직접 원핫 이진 벡터 표현 얻기

- 인덱스된 값을 1000여개의 단어의 해당 단어의 위치에 매핑

In [12]:

```
one_hot = tokenizer.texts_to_matrix(samples, mode='binary')
one_hot.shape, one_hot
```

Out[12]:

```
((2, 1000),
 array([[0., 1., 1., ..., 0., 0., 0.],
        [0., 1., 0., ..., 0., 0., 0.]])
```

In [13]:

```
# 계산된 단어 인덱스를 구합니다.
word_index = tokenizer.word_index
print(word_index)
print('Found %s unique tokens.' % len(word_index))
```

```
{'the': 1, 'cat': 2, 'sat': 3, 'on': 4, 'mat': 5, 'dog': 6, 'ate': 7, 'my': 8, 'home
work': 9}
Found 9 unique tokens.
```

## 03 해싱 기법을 이용한 원핫 인코딩-변형된 형태

- 원핫 인코딩의 변종중의 하나는 원-핫 해싱(one-hot hashing)기법입니다.
- 고유한 토큰의 수가 너무 커서 다루기 어려울 때 사용.
  - 각 단어에 명시적으로 인덱스를 할당하고 이 인덱스를 딕셔너리에 저장하지 않고, 단어를 해싱한다.

In [14]:



```
samples = ['The cat sat on the mat.', 'The dog ate my homework.']

# 단어를 크기가 1,000인 벡터로 저장합니다.
# 1,000개(또는 그이상)의 단어가 있다면 해싱 충돌이 늘어나고 인코딩의 정확도가 감소될 것입니다
dimensionality = 1000
max_length = 10

results = np.zeros((len(samples), max_length, dimensionality))
for i, sample in enumerate(samples):
    for j, word in list(enumerate(sample.split()))[:max_length]:
        # 단어를 해싱하여 0과 1,000 사이의 랜덤한 정수 인덱스로 변환합니다.
        index = abs(hash(word)) % dimensionality
        results[i, j, index] = 1.
```

In [15]:



```
results.shape, results
```

Out[15]:

```
((2, 10, 1000),
array([[0., 0., 0., ..., 0., 0., 0.],
       [0., 0., 0., ..., 0., 0., 0.],
       [0., 0., 0., ..., 0., 0., 0.],
       ...,
       [0., 0., 0., ..., 0., 0., 0.],
       [0., 0., 0., ..., 0., 0., 0.],
       [0., 0., 0., ..., 0., 0., 0.]],

      [[0., 0., 0., ..., 0., 0., 0.],
       [0., 0., 0., ..., 0., 0., 0.],
       [0., 0., 0., ..., 0., 0., 0.],
       ...,
       [0., 0., 0., ..., 0., 0., 0.],
       [0., 0., 0., ..., 0., 0., 0.],
       [0., 0., 0., ..., 0., 0., 0.]])])
```