
딥러닝을 이용한 자연어 처리 입문

02. 텍스트 전처리(Text Preprocessing) 06~10

2021.06.20

집현전 초급반 2조: 이창훈(조장), 이문형, 유예림

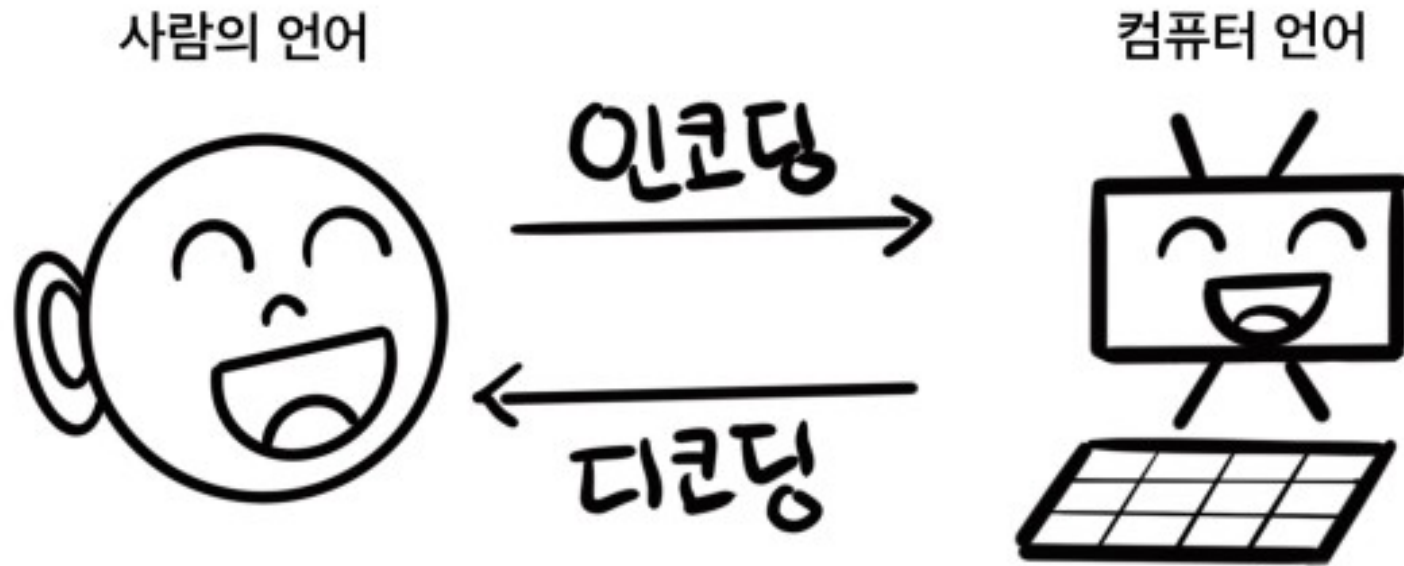
목차

1. 정수 인코딩(Integer Encoding)
2. 패딩(Padding)
3. 원-핫 인코딩(One-Hot Encoding)
4. 데이터의 분리(Splitting Data)
5. 한국어 전처리 패키지(Text Preprocessing Tools for Korean Text)

1. 정수 인코딩(Integer Encoding)

◆ 인코딩이란?

- 자연어를 기계가 이해할 수 있는 언어로 변환하는 과정



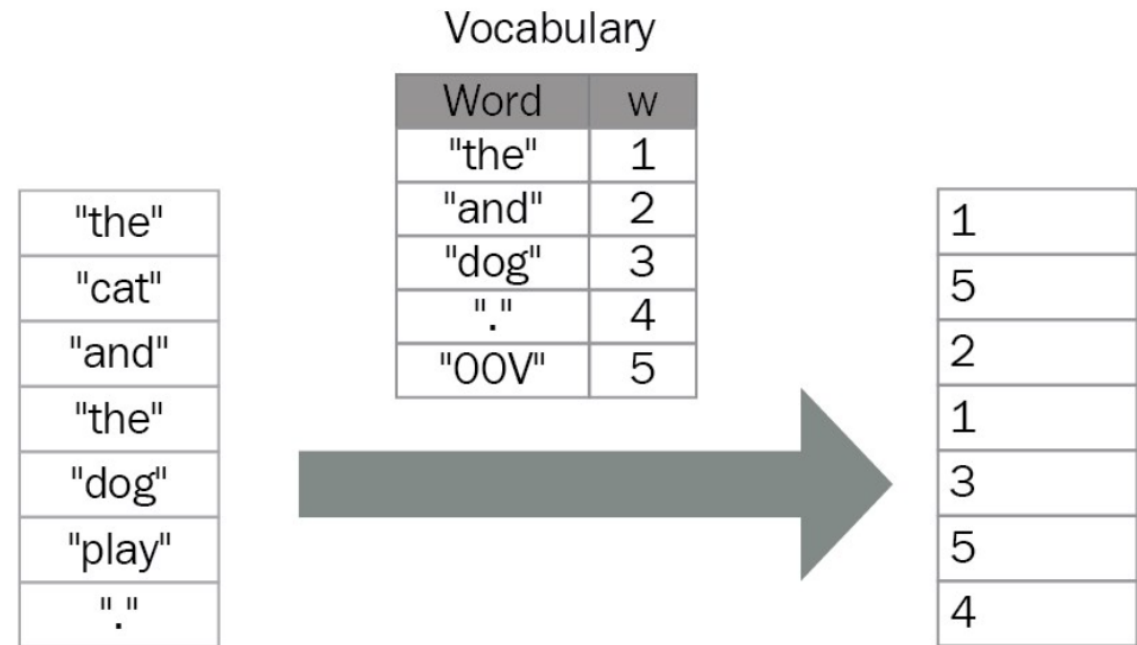
1. 정수 인코딩(Integer Encoding)

◆ 정수 인코딩이란?

- **word2index** 맵핑
- cleansing, tokenization, subword segmentation 등의 전처리 작업 이후에 모델에 넣기 전에 스해하

◆ 인덱스 부여 방법

- 1) 랜덤하게 부여함
- 2) 단어에 대한 **빈도수** 기준으로 정렬 후 부여함



1. 정수 인코딩(Integer Encoding)

◆ Word Representation

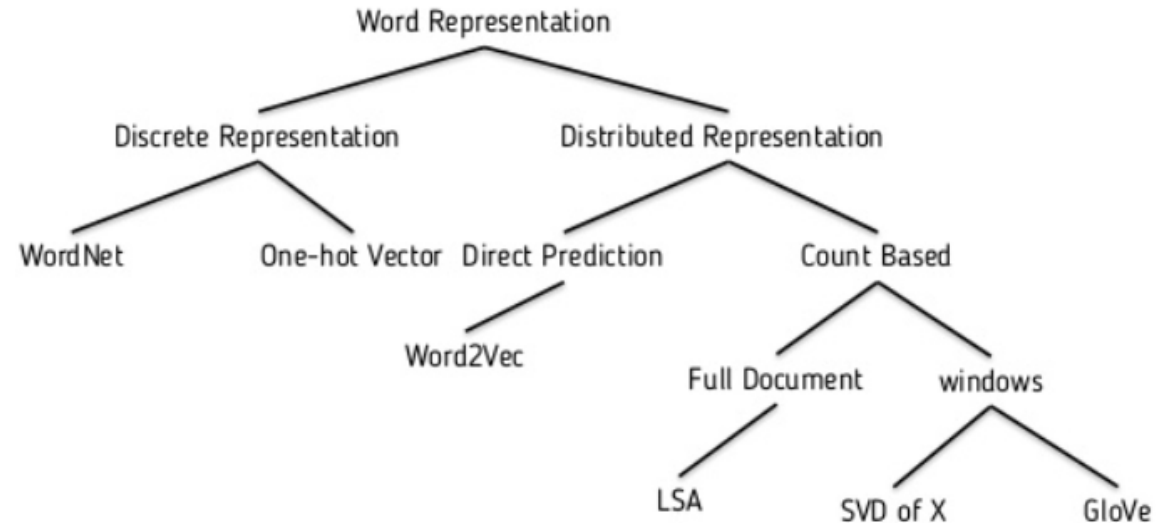
- 단어를 벡터로 **잘** 표현하는 방법으로 발전함



Language Representation 흐름도

- ❑ One Hot Encoding
- ❑ Word2Vec – CBOW, Skip-Gram
- ❑ Glove
- ❑ FastText
- ❑ Cove
- ❑ ELMO
- ❑ GPT
- ❑ BERT
- ❑ GPT-2
- ❑ ERNIE
- ❑ XLNET
- ❑ ERNIE 2.0
- ❑ ROBERTa

Word Representation 분류



1. 정수 인코딩(Integer Encoding)

◆ 정수 인코딩 방법

- 1) dictionary 사용하기
- 2) Counter 사용하기
- 3) NLTK의 FreqDist 사용하기
- 4) enumerate 이해하기
- 5) 케라스(Keras)의 텍스트 전처리

1. 정수 인코딩(Integer Encoding)

◆ 정수 인코딩 방법 1) dictionary 사용하기 - 라이브러리 및 데이터

```
[1] !pip install nltk
```

```
Requirement already satisfied: nltk in /usr/local/lib/python3.7/dist-packages (3.2.5)  
Requirement already satisfied: six in /usr/local/lib/python3.7/dist-packages (from nltk) (1.15.0)
```

```
[2] import nltk
```

```
nltk.download('punkt')  
nltk.download('stopwords')
```

```
from nltk.tokenize import sent_tokenize  
from nltk.tokenize import word_tokenize  
from nltk.corpus import stopwords
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...  
[nltk_data]   Unzipping tokenizers/punkt.zip.  
[nltk_data] Downloading package stopwords to /root/nltk_data...  
[nltk_data]   Unzipping corpora/stopwords.zip.
```

```
[3] text = "A barber is a person. a barber is good person. a barber is huge person. he Knew A Secret! The Secret He Kept is huge secret. Huge secret. His barber kept his word. a l
```

1. 정수 인코딩(Integer Encoding)

◆ 정수 인코딩 방법 1) dictionary 사용하기

- 문장 토큰화
- NLTK의 `sent_tokenize()` 함수는 `nltk.tokenize.punkt` 모듈의 `PunktSentenceTokenizer` 인스턴스를 사용해서 문장의 시작과 끝을 구분함
- 이미 다수의 유럽 언어에서 검증이 됨

문장 토큰화

```
text = sent_tokenize(text)
print(text)
```

['A barber is a person.', 'a barber is good person.', 'a barber is huge person.', 'he Knew A Secret!', 'The Secret He Kept is huge secret.', 'Huge secret.', 'His barber kept his word.', 'a barber &

1. 정수 인코딩(Integer Encoding)

- ◆ 정수 인코딩 방법 1) dictionary 사용하기
 - 정제와 단어 토큰화
 - NLTK의 `word_tokenize()` 함수는 TreebankWordTokenizer 클래스의 인스턴스에 `tokenize()` 함수를 호출하는 래퍼 함수임
 - 공백과 구두점을 이용해 단어를 분해하고, 구두점을 버리지는 않음

1. 정수 인코딩(Integer Encoding)

◆ 정수 인코딩 방법 1) dictionary 사용하기

- 정제와 단어 토큰화

[5] # 정제와 단어 토큰화

vocab = {} # 파이썬의 dictionary 자료형

```
sentences = []
```

```
stop_words = set(stopwords.words('english'))
```

```
for i in text:
```

```
sentence = word_tokenize(i) # 단어 토큰화를 수행합니다.
```

```
result = []
```

```
for word in sentence:
```

```
word = word.lower() # 모든 단어를 소문자화하여 단어의 개수를 줄입니다.
```

```
if word not in stop_words: # 단어 토큰화 된 결과에 대해서 불용어를 제거합니다.
```

```
if len(word) > 2: # 단어 길이가 2이하인 경우에 대하여 추가로 단어를 제거합니다.
```

```
result.append(word)
```

```
if word not in vocab:
```

```
vocab[word] = 0
```

```
vocab[word] += 1
```

```
sentences.append(result)
```

```
print(sentences)
```

[[['barber', 'person'], ['barber', 'good', 'person'], ['barber', 'huge', 'person'], ['knew', 'secret'], ['secret', 'kept', 'huge', 'secret'], ['huge', 'secret'], ['barber', 'kept', 'word'], ['barber

1. 정수 인코딩(Integer Encoding)

- ◆ 정수 인코딩 방법 1) dictionary 사용하기
 - vocab에는 중복을 제거한 단어와 각 단어에 대한 빈도수가 기록됨

```
[6] print(vocab)
```

```
{'barber': 8, 'person': 3, 'good': 1, 'huge': 5, 'knew': 1, 'secret': 6, 'kept': 4, 'word': 2, 'keeping': 2, 'driving': 1, 'crazy': 1, 'went': 1, 'mountain': 1}
```

```
[7] print(vocab["barber"]) # 'barber'라는 단어의 빈도수 출력
```

1. 정수 인코딩(Integer Encoding)

- ◆ 정수 인코딩 방법 1) dictionary 사용하기
 - 빈도수가 높은 순서대로 정렬함

```
[8] vocab_sorted = sorted(vocab.items(), key = lambda x:x[1], reverse = True)
    print(vocab_sorted)
```

```
[('barber', 8), ('secret', 6), ('huge', 5), ('kept', 4), ('person', 3), ('word', 2), ('keeping', 2), ('good', 1), ('knew', 1), ('driving', 1), ('crazy', 1), ('went', 1), ('mountain', 1)]
```

1. 정수 인코딩(Integer Encoding)

◆ 정수 인코딩 방법 1) dictionary 사용하기

- 높은 빈도수를 가진 단어일수록 낮은 정수 인덱스를 부여함
- 빈도수 상위 n개의 단어만 사용함

```
[9] word_to_index = {}  
    i=0  
    for (word, frequency) in vocab_sorted :  
        if frequency > 1 : # 정제(Cleaning) 챕터에서 언급했듯이 빈도수가 적은 단어는 제외한다.  
            i=i+1  
            word_to_index[word] = i  
    print(word_to_index)
```

```
➡ {'barber': 1, 'secret': 2, 'huge': 3, 'kept': 4, 'person': 5, 'word': 6, 'keeping': 7}
```

```
[10] vocab_size = 5  
words_frequency = [w for w,c in word_to_index.items() if c >= vocab_size + 1] # 인덱스가 5 초과인 단어 제거  
for w in words_frequency:  
    del word_to_index[w] # 해당 단어에 대한 인덱스 정보를 삭제  
print(word_to_index)
```

```
{'barber': 1, 'secret': 2, 'huge': 3, 'kept': 4, 'person': 5}
```

1. 정수 인코딩(Integer Encoding)

◆ 정수 인코딩 방법 1) dictionary 사용하기

- **word_to_index**에 단어 집합에 존재하지 않는 단어들 (Out-Of-Vocabulary, OOV)를 새롭게 추가하고,
단어 집합에 없는 단어들은 'OOV'의 인덱스로 인코딩함

```
[11] word_to_index['OOV'] = len(word_to_index) + 1
```

```
[12] encoded = []  
for s in sentences:  
    temp = []  
    for w in s:  
        try:  
            temp.append(word_to_index[w])  
        except KeyError:  
            temp.append(word_to_index['OOV'])  
    encoded.append(temp)  
print(encoded)
```

```
[[1, 5], [1, 6, 5], [1, 3, 5], [6, 2], [2, 4, 3, 2], [3, 2], [1, 4, 6], [1, 4, 6], [1, 4, 2], [6, 6, 3, 2, 6, 1, 6], [1, 6, 3, 6]]
```

1. 정수 인코딩(Integer Encoding)

◆ 정수 인코딩 방법 2) Counter 사용하기

- 라이브러리 및 데이터
- 단어 집합을 만들기 위해서 문장의 경계인 [,]를 제거하고
단어들을 하나의 리스트로 만듦

```
[13] from collections import Counter
```

```
[14] print(sentences)
```

```
[['barber', 'person'], ['barber', 'good', 'person'], ['barber', 'huge', 'person'], ['knew', 'secret'], ['secret', 'kept', 'huge', 'secret'], ['huge', 'secret'], ['barber', 'kept', 'word'], ['barber',
```

```
[15] words = sum(sentences, [])
```

```
# 위 작업은 words = np.hstack(sentences)로도 수행 가능.
```

```
print(words)
```

```
['barber', 'person', 'barber', 'good', 'person', 'barber', 'huge', 'person', 'knew', 'secret', 'secret', 'kept', 'huge', 'secret', 'huge', 'secret', 'barber', 'kept', 'word', 'barber', 'kept', 'wor
```

1. 정수 인코딩(Integer Encoding)

◆ 정수 인코딩 방법 2) Counter 사용하기

- Counter(), most common()

```
[16] vocab = Counter(words) # 파이썬의 Counter 모듈을 이용하면 단어의 모든 빈도를 쉽게 계산할 수 있습니다.  
      print(vocab)
```

```
Counter({'barber': 8, 'secret': 6, 'huge': 5, 'kept': 4, 'person': 3, 'word': 2, 'keeping': 2, 'good': 1, 'knew': 1, 'driving': 1, 'crazy': 1, 'went': 1, 'mountain': 1})
```

```
[17] print(vocab["barber"]) # 'barber'라는 단어의 빈도수 출력
```

```
8
```

```
[18] vocab_size = 5  
      vocab = vocab.most_common(vocab_size) # 등장 빈도수가 높은 상위 5개의 단어만 저장  
      vocab
```

```
[('barber', 8), ('secret', 6), ('huge', 5), ('kept', 4), ('person', 3)]
```

```
[19] word_to_index = {}  
      i = 0  
      for (word, frequency) in vocab :  
          i = i+1  
          word_to_index[word] = i  
      print(word_to_index)
```

```
{'barber': 1, 'secret': 2, 'huge': 3, 'kept': 4, 'person': 5}
```


1. 정수 인코딩(Integer Encoding)

- ◆ 정수 인코딩 방법 3) NLTK의 FreqDist 사용하기
 - 라이브러리 및 FreqDist()

```
[20] from nltk import FreqDist  
import numpy as np
```

```
[21] # np.hstack으로 문장 구분을 제거하여 입력으로 사용 . ex) ['barber', 'person', 'barber', 'good' ... 중략 ...  
vocab = FreqDist(np.hstack(sentences))
```

```
[22] print(vocab["barber"]) # 'barber'라는 단어의 빈도수 출력
```

1. 정수 인코딩(Integer Encoding)

- ◆ 정수 인코딩 방법 3) NLTK의 FreqDist 사용하기
 - `most_common()`, `enumerate()`

```
[103] vocab_size = 5
      vocab = vocab.most_common(vocab_size) # 등장 빈도수가 높은 상위 5개의 단어만 저장
      vocab
```

```
[('barber', 8), ('secret', 6), ('huge', 5), ('kept', 4), ('person', 3)]
```

```
[104] word_to_index = {word[0] : index + 1 for index, word in enumerate(vocab)}
      print(word_to_index)
```

```
{'barber': 1, 'secret': 2, 'huge': 3, 'kept': 4, 'person': 5}
```

1. 정수 인코딩(Integer Encoding)

◆ 정수 인코딩 방법 4) enumerate 이해하기

- enumerate()는 순서가 있는 자료형(list, set, tuple, dictionary, string)을 입력으로 받아 인덱스를 순차적으로 함께 리턴함

```
[105] test=['a', 'b', 'c', 'd', 'e']  
      for index, value in enumerate(test): # 입력의 순서대로 0부터 인덱스를 부여함.  
          print("value : {}, index: {}".format(value, index))
```

```
value : a, index: 0  
value : b, index: 1  
value : c, index: 2  
value : d, index: 3  
value : e, index: 4
```

1. 정수 인코딩(Integer Encoding)

- ◆ 정수 인코딩 방법 5) 케라스(Keras)의 텍스트 전처리
 - 라이브러리 및 데이터
 - keras에서는 기본적인 전처리를 위한 도구들을 제공함
 - 자동으로 높은 빈도수를 가진 단어일수록 낮은 정수 인덱스를 부여함

```
[106] from tensorflow.keras.preprocessing.text import Tokenizer
```

```
[107] sentences=[['barber', 'person'], ['barber', 'good', 'person'], ['barber', 'huge', 'person'], ['knew', 'secret'], ['secret', 'kept', 'huge', 'secret'], ['huge', 'secret'], ['barber', 'kept', 'word
```

1. 정수 인코딩(Integer Encoding)

◆ 정수 인코딩 방법 5) 케라스(Keras)의 텍스트 전처리

- `fit_on_texts()`, `word_index`, `index_word`, `word_counts`, `text_to_sequences()`

```
[156] tokenizer = Tokenizer()  
tokenizer.fit_on_texts(sentences) # fit_on_texts()안에 코퍼스를 입력으로 하면 빈도수를 기준으로 단어 집합을 생성한다.
```

```
[157] print(tokenizer.word_index)  
  
{'barber': 1, 'secret': 2, 'huge': 3, 'kept': 4, 'person': 5, 'word': 6, 'keeping': 7, 'good': 8, 'knew': 9, 'driving': 10, 'crazy': 11, 'went': 12, 'mountain': 13}
```

```
[158] print(tokenizer.index_word)  
  
{1: 'barber', 2: 'secret', 3: 'huge', 4: 'kept', 5: 'person', 6: 'word', 7: 'keeping', 8: 'good', 9: 'knew', 10: 'driving', 11: 'crazy', 12: 'went', 13: 'mountain'}
```

```
[159] print(tokenizer.word_counts)  
  
OrderedDict([('barber', 8), ('person', 3), ('good', 1), ('huge', 5), ('knew', 1), ('secret', 6), ('kept', 4), ('word', 2), ('keeping', 2), ('driving', 1), ('crazy', 1), ('went', 1), ('mountain', 1)])
```

```
[160] print(tokenizer.texts_to_sequences(sentences))  
  
[[1, 5], [1, 8, 5], [1, 3, 5], [9, 2], [2, 4, 3, 2], [3, 2], [1, 4, 6], [1, 4, 6], [1, 4, 2], [7, 7, 3, 2, 10, 1, 11], [1, 12, 3, 13]]
```

1. 정수 인코딩(Integer Encoding)

- ◆ 정수 인코딩 방법 5) 케라스(Keras)의 텍스트 전처리
 - `tokenizer = Tokenizer(num_words = vocab_size + 1)`
 - `most_common()`과 같은 기능임
 - keras tokenizer가 숫자 0까지 단어 집합의 크기로

```
[112] vocab_size = 5  
      tokenizer = Tokenizer(num_words = vocab_size + 1) # 상위 5개 단어만 사용  
      tokenizer.fit_on_texts(sentences)
```

```
[113] print(tokenizer.word_index)
```

```
{'barber': 1, 'secret': 2, 'huge': 3, 'kept': 4, 'person': 5, 'word': 6, 'keeping': 7, 'good': 8, 'knew': 9, 'driving': 10, 'crazy': 11, 'went': 12, 'mountain': 13}
```

```
[114] print(tokenizer.word_counts)
```

```
OrderedDict([('barber', 8), ('person', 3), ('good', 1), ('huge', 5), ('knew', 1), ('secret', 6), ('kept', 4), ('word', 2), ('keeping', 2), ('driving', 1), ('crazy', 1), ('went', 1), ('mountain', 1)])
```

```
[115] print(tokenizer.texts_to_sequences(sentences))
```

```
[[1, 5], [1, 5], [1, 3, 5], [2], [2, 4, 3, 2], [3, 2], [1, 4], [1, 4], [1, 4, 2], [3, 2, 1], [1, 3]]
```

1. 정수 인코딩(Integer Encoding)

◆ 정수 인코딩 방법 5) 케라스(Keras)의 텍스트 전처리

- **word_counts, word_index**를 출력해도 같은 결과가 나오지만, **text_to_sequences()**를 하면 적용이 됨

```
[116] tokenizer = Tokenizer() # num_words를 여기서는 지정하지 않은 상태
      tokenizer.fit_on_texts(sentences)
```

```
[117] vocab_size = 5
      words_frequency = [w for w,c in tokenizer.word_index.items() if c >= vocab_size + 1] # 인덱스가 5 초과인 단어 제거
      for w in words_frequency:
          del tokenizer.word_index[w] # 해당 단어에 대한 인덱스 정보를 삭제
          del tokenizer.word_counts[w] # 해당 단어에 대한 카운트 정보를 삭제
      print(tokenizer.word_index)
      print(tokenizer.word_counts)
      print(tokenizer.texts_to_sequences(sentences))
```

```
{'barber': 1, 'secret': 2, 'huge': 3, 'kept': 4, 'person': 5}
OrderedDict([('barber', 8), ('person', 3), ('huge', 5), ('secret', 6), ('kept', 4)])
[[1, 5], [1, 5], [1, 3, 5], [2], [2, 4, 3, 2], [3, 2], [1, 4], [1, 4], [1, 4, 2], [3, 2, 1], [1, 3]]
```

1. 정수 인코딩(Integer Encoding)

◆ 정수 인코딩 방법 5) 케라스(Keras)의 텍스트 전처리

- **keras tokenizer**는 정수 인코딩 과정 중 단어 집합에 없는 단어(OOV)에 대해서는 아예 해당 단어를 제거해버림
- `tokenizer = Tokenizer(num_words = vocab_size + 2, oov_token = 'OOV')`

```
[118] vocab_size = 5
      tokenizer = Tokenizer(num_words = vocab_size + 2, oov_token = 'OOV')
      # 빈도수 상위 5개 단어만 사용. 숫자 0과 OOV를 고려해서 단어 집합의 크기는 +2
      tokenizer.fit_on_texts(sentences)
```

```
[119] print('단어 OOV의 인덱스 : {}'.format(tokenizer.word_index['OOV']))
```

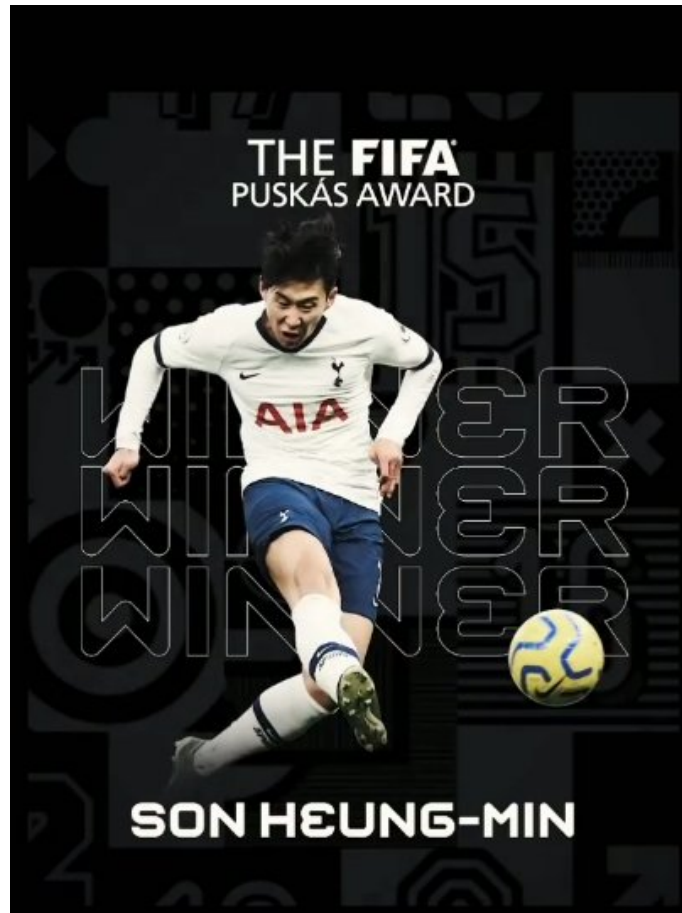
단어 OOV의 인덱스 : 1

```
[120] print(tokenizer.texts_to_sequences(sentences))
```

[[2, 6], [2, 1, 6], [2, 4, 6], [1, 3], [3, 5, 4, 3], [4, 3], [2, 5, 1], [2, 5, 1], [2, 5, 3], [1, 1, 4, 3, 1, 2, 1], [2, 1, 4, 1]]

1. 정수 인코딩(Integer Encoding)

- ◆ 실습 : 유튜브 동영상(손흥민 푸스카스상) 댓글 정수 인코딩



1. 정수 인코딩(Integer Encoding)

◆ 실습 : 유튜브 동영상(손흥민 푸스카스상) 댓글 정수 인코딩

- 동영상 댓글 크롤링
(3929개)

```
In [14]: 1 # 손흥민 골 댓글  
2 dfs[9].head(20)
```

Out [14]:

	0	text	vote	time
0	{"cid": "UgxydcFfKsrXwwH1X-N4AaABAg", "text": ...	"금호타이어 광고 제대로 했네"	"0"	"3주 전"
1	{"cid": "UgxQZfMv3XCAFAHAUbN4AaABAg", "text": ...	"What a amazing goal"	"0"	"3주 전"
2	{"cid": "UgyRHAhvQqg6zil2hLR4AaABAg", "text": ...	"Rancak Bana..!!!"	"0"	"3주 전"
3	{"cid": "UgyILPpkWGthnm8nAOt4AaABAg", "text": ...	"Show...!!!"	"0"	"3주 전"
4	{"cid": "UgzsZCepVjnMouQ8WjF4AaABAg", "text": ...	"Well deserved"	"0"	"2주 전"
5	{"cid": "Ugyz1nSywaPSpWHo9T94AaABAg", "text": ...	"Superrr"	"0"	"2주 전"
6	{"cid": "UgxI_0aVVMGUj9zT0f4AaABAg", "text": ...	"Barcelona vs PSG Messi rocket 🚀"	"0"	"2주 전"
7	{"cid": "UgzeS6lVo4SpWTNenFN4AaABAg", "text": ...	"🍌🍌🍌🍌🍌🍌🍌🍌"	"0"	"2주 전"
8	{"cid": "UgwOEBcb9IMLg76czcJ4AaABAg", "text": ...	"Suarez snd Arrascaetas were both much better ...	"0"	"2주 전(수정됨)"
9	{"cid": "Ugxknx0_gv--W7jza6F4AaABAg", "text": ...	"Better than Maradona"	"0"	"2주 전"
10	{"cid": "UgxCAiciw6jidKDXZfn4AaABAg", "text": ...	"아시아 최고의 축구 선수"	"0"	"2주 전"
11	{"cid": "Ugzy-KnVawoR9G1g_5l4AaABAg", "text": ...	"VIDA"	"0"	"2주 전"
12	{"cid": "UgwQv_kaVR6w-B0hFB94AaABAg", "text": ...	"In any other country he would have been unfai..."	"2"	"2주 전"
13	{"cid": "UgxTMw94B4L_re8Lnvd4AaABAg", "text": ...	"Nice one"	"0"	"2주 전"
14	{"cid": "Ugybribs6TnThQ9SQRAaABAg", "text": ...	"🤔👉👈🍌🍌🍌🍌🍌🍌"	"0"	"1주 전"
15	{"cid": "Ugy9N7fLIKbSJDFDWxt4AaABAg", "text": ...	"Son needs to leave Tottenham or he will never..."	"0"	"1주 전(수정됨)"
16	{"cid": "Ugwnlkx6LRu7H5KsC094AaABAg", "text": ...	"O cara canta sbt.joga na Inglaterra e ainda g..."	"3"	"1주 전"
17	{"cid": "Ugw8h6qzunzVe8mELaB4AaABAg", "text": ...	"Era de suarez para mi pero no puedo negar que..."	"0"	"6일 전"
18	{"cid": "UgwQgGK5y0LNutoVOG14AaABAg", "text": ...	"OGM he is world class!"	"0"	"1일 전"
19	{"cid": "UgwnVbiHxi5CDNqif_p4AaABAg", "text": ...	"amazing goal but didn't deserve a puskas award"	"0"	"14시간 전"

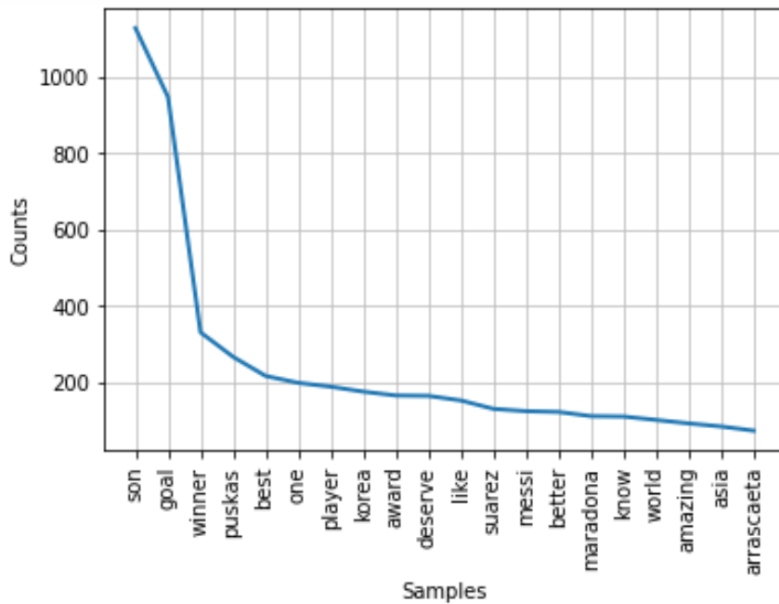
1. 정수 인코딩(Integer Encoding)

- ◆ 실습 : 유튜브 동영상(손흥민 푸스카스상) 댓글 정수 인코딩
 - 간단한 텍스트 전처리
 - NLTK의 word_tokenize()를 활용하여 토큰화를 수행함
 - 단어의 토큰화는 단어 전체, 명사, 형용사, bigram, trigram에 대하여 진행함
 - 불용어 제거 (영어, 특수문자, 사용자 정의)
 - 대소문자 통일 및 단어 길이 2이하 제거
 - 빈도수 TOP 100의 결과를 비교하여 단어 치환을 수행함
(유사단어, 고유명사, 복수형 등)
(예시)
son heung min, sonny, sonaldo, ... → son
deserved, deserve, ... → deserve
jumo, south korea, korean, ... → korea
girl, woman, female, ... → women 등

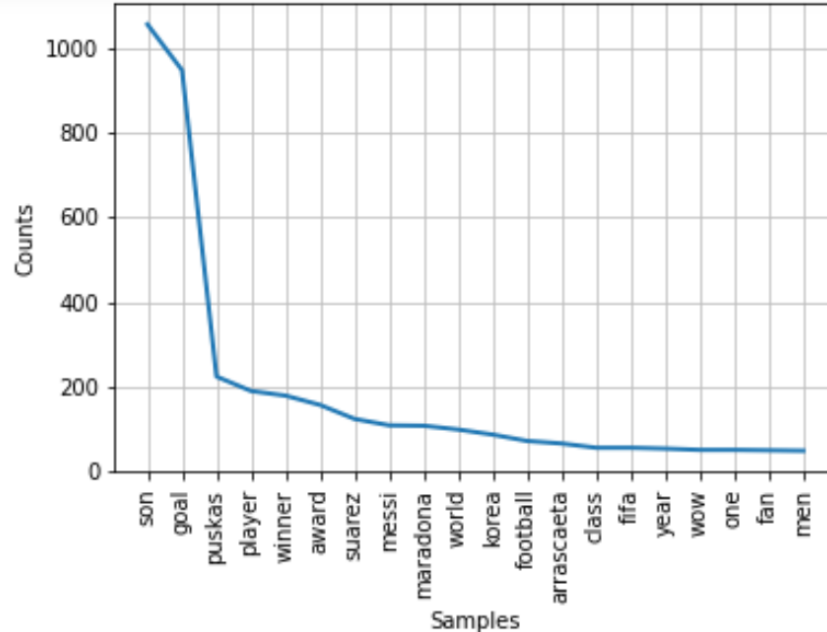
1. 정수 인코딩(Integer Encoding)

◆ 실습 : 유튜브 동영상(손흥민 푸스카스상) 댓글 정수 인코딩 - 댓글 빈도수 분석

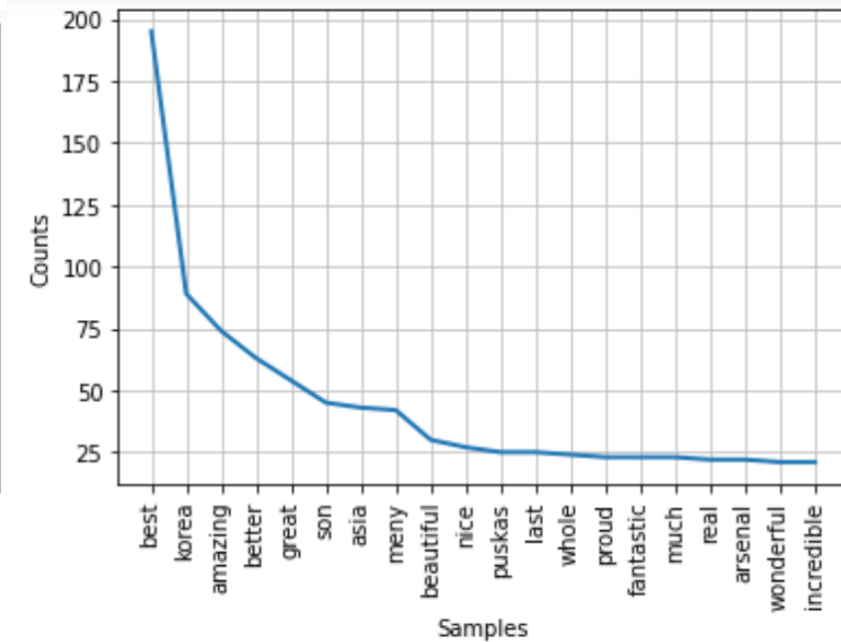
모든 단어(unigram) TOP 20



명사 TOP 20



형용사 TOP 20



1. 정수 인코딩(Integer Encoding)

- ◆ 실습 : 유튜브 동영상(손흥민 푸스카스상) 댓글 정수 인코딩
 - 댓글 빈도수 분석
 - 손흥민을 지칭하는 단어들이 가장 많이 나옴
 - winner, best, award, deserve가 높게 나타나는 것을 보아, 손흥민이 1등을 했다는 것을 짐작해볼 수 있음
 - korea, asia가 높은 것을 보아 손흥민의 국적을 알 수 있음
 - messi, maradona를 보아 손흥민이 메시, 마라도나를 연상시킨다는 말이 나왔음을 짐작해봄

1. 정수 인코딩(Integer Encoding)

- ◆ 실습 : 유튜브 동영상(손흥민 푸스카스상) 댓글 정수 인코딩
- 워드 클라우드



1. 정수 인코딩(Integer Encoding)

◆ 실습 : 유튜브 동영상(손흥민 푸스카스상) 댓글 정수 인코딩

- 정수 인코딩(NLTK의 FreqDist)

- 약 4000개의 단어 중에서 빈도수가 3 이상인 단어의 수를

num_words(=1060)으로 설정함

- 전미하 저처리 데이터 수직 등으로 단어 히스토그램을 줄여야 함

```
In [35]: 1 # 정수 인코딩 (NLTK)
2
3 vocab = top3_fd[0].most_common(1060)
4 word_to_index = {word[0] : index + 1 for index, word in enumerate(vocab)}
5 word_to_index['OOV'] = 0
6
7 sorted_list = sorted(word_to_index.items(), key=lambda item: item[1])
8 word_to_index = {k: v for k, v in sorted_list}
9
10 print(word_to_index)
```

```
{'OOV': 0, 'son': 1, 'goal': 2, 'winner': 3, 'puskas': 4, 'best': 5, 'one': 6, 'player': 7, 'korea': 8, 'award': 9, 'deserve': 10, 'I
like': 11, 'suarez': 12, 'messi': 13, 'better': 14, 'maradona': 15, 'know': 16, 'world': 17, 'amazing': 18, 'asia': 19, 'arrascaeta':
20, 'football': 21, 'scored': 22, 'filfa': 23, 'great': 24, 'wow': 25, 'even': 26, 'think': 27, 'get': 28, 'run': 29, 'would': 30, 'cl
ass': 31, 'lol': 32, 'year': 33, 'nice': 34, 'arsenal': 35, 'ball': 36, 'fan': 37, 'men': 38, 'see': 39, 'solo': 40, '2020': 41, 'eve
n': 42, 'burnley': 43, 'really': 44, 'vote': 45, 'doubt': 46, 'well': 47, 'meny': 48, 'proud': 49, 'super': 50, 'sead': 51, 'never':
52, 'team': 53, 'every': 54, 'kolasinac': 55, 'much': 56, 'yes': 57, 'could': 58, 'salah': 59, 'dribbling': 60, 'time': 61, 'fast': 6
2, 'way': 63, 'running': 64, 'got': 65, 'ronaldo': 66, 'beautiful': 67, 'seen': 68, 'newflow': 69, 'absolutely': 70, 'another': 71,
'love': 72, 'must': 73, 'gol': 74, 'skill': 75, 'say': 76, 'times': 77, 'speed': 78, 'league': 79, 'que': 80, 'legend': 81, 'congratu
lations': 82, 'still': 83, 'past': 84, 'play': 85, 'look': 86, 'definitely': 87, 'people': 88, 'two': 89, 'real': 90, 'whole': 91, 'd
ribble': 92, 'also': 93, 'last': 94, 'defenders': 95, 'score': 96, 'fantastic': 97, 'dont': 98, 'watch': 99, 'right': 100, 'everyon
e': 101, 'level': 102, 'overrated': 103, 'video': 104, 'let': 105, 'make': 106, 'give': 107, 'gon': 108, 'assist': 109, 'south': 110,
'crazy': 111, 'first': 112, 'winnning': 113, 'made': 114, 'prize': 115, 'wonderful': 116, 'game': 117, 'tottenham': 118, 'easy': 11
9, 'winners': 120, 'fans': 121, 'back': 122, 'incredible': 123, 'kick': 124, 'no.1': 125, 'far': 126, 'without': 127, 'years': 128,
'spurs': 129, 'yeah': 130, 'sure': 131, 'del': 132, 'nothing': 133, 've': 134, '👍': 135, 'defender': 136, 'come': 137, 'socce
r': 138, 'dribbled': 139, 'nosonated': 140, 'century': 141, 'ran': 142, 'puskas': 143, 'crowd': 144, 'history': 145, 'congrats': 146,
'stop': 147, 'long': 148, 'else': 149, 'others': 150, 'agree': 151, 'mean': 152, 'style': 153, 'coys': 154, 'special': 155, 'box': 15
6, 'watched': 157, 'defending': 158, 'pace': 159, 'field': 160, 'king': 161, '.....': 162, 'anyone': 163, 'life': 164, 'diego': 165,
'next': 166, 'literally': 167, 'respect': 168, 'top': 169, 'maybe': 170, 'makes': 171, 'robbed': 172, 'something': 173, 'actually': 1
74, 'done': 175, 'nosonee': 176, 'already': 177, 'golazo': 178, 'second': 179, 'course': 180, 'brilliant': 181, 'premier': 182, 'simi
lar': 183, 'voted': 184, 'weak': 185, 'simply': 186, 'wonder': 187, 'haha': 188, 'sprint': 189, 'taker': 190, 'played': 191, 'believ
```

```
In [36]: 1 encoded = []
2 for s in top3_tokens[0]:
3     temp = []
4     for w in s:
5         try:
6             temp.append(word_to_index[w]) # 단어를 인덱스로 치환
7         except:
8             temp.append(word_to_index['OOV'])
9     encoded.append(temp)
10
11 print(encoded)
```

```
[[0, 912], [18, 2], [0, 0], [294], [47, 10], [0], [244, 0, 13, 913], [0], [12, 0, 0, 56, 14, 2], [14, 15], [0, 0], [0], [266, 30, 0,
453, 679, 168, 454, 0, 0], [34, 6], [0], [1, 338, 0, 118, 52, 3, 267, 268, 53, 914, 7, 0], [0, 0, 0, 0, 680, 0, 0, 0, 0], [385, 12, 2
45, 269, 0, 0, 80, 222, 178], [0, 17, 31], [18, 2, 10, 4, 9], [0, 1], [1, 295, 681, 2, 339, 0, 109, 0, 246, 2, 915, 0], [0, 0, 0, 0],
[0, 0], [916, 132, 682, 917, 0, 0], [918], [0], [1, 386], [0, 0, 0, 0, 0, 0, 0, 0], [296, 5, 109, 387], [20], [0, 0, 161, 19], [19, 1
25, 0], [12, 10], [1], [0, 919], [920, 0], [340, 6, 223, 0, 341, 8], [4, 9, 141], [297, 84, 89, 7, 99, 4, 102], [0], [683, 1], [109,
298, 0], [43, 224, 223, 342, 14, 53, 1, 103, 0, 225, 455], [0, 0, 0], [0], [2, 12, 2, 1], [299, 538, 94, 185, 539, 0, 921, 270, 922,
923, 684, 226], [25], [0, 0, 0, 0], [0], [0, 0], [1], [0, 685, 456], [154], [1, 540, 186, 0], [0, 0, 0, 0, 0, 0, 74, 0], [1], [5, 126,
142, 300, 186, 0, 0, 686, 11, 95, 133, 155, 687, 924, 7], [1], [0, 0, 0, 118, 1, 89, 2, 6, 109, 925, 457, 225, 541], [388], [0, 0, 54,
53, 30, 343, 85], [119, 54, 185, 7, 0, 458, 96, 11, 686, 11, 95, 2, 10, 4], [247, 344], [542, 688, 542, 1, 542, 542, 688, 688], [90,
227, 926, 38, 25], [0], [0, 0, 0, 0, 0], [927, 928, 0, 11, 0, 927, 928, 0, 0, 0], [35, 81, 10, 9, 119, 0], [0, 389], [15, 2, 0], [1
0], [1, 120, 4, 60, 84, 95, 187, 48, 4, 13, 10], [82], [543, 14, 2, 17, 21, 91, 0], [0, 43, 0, 0, 0, 0, 0, 689, 0, 0, 74, 0, 0], [34
3, 0, 1], [49, 50, 1, 0, 690], [1, 1, 929], [49, 19], [110, 8, 49, 1], [188, 1], [13, 11, 459, 0], [247, 344], [57], [691, 0, 544],
[0, 0, 543, 545, 460, 461, 58, 143, 9], [204, 301, 156], [345, 13, 62, 1], [0, 2, 1], [346, 0, 0, 0], [1, 930], [32], [17, 692, 70, 1
7, 692], [1, 1, 0, 205, 5, 7], [24, 7, 18, 2], [1, 693], [0, 694, 302, 2, 35, 71, 546, 2], [0, 931, 13, 2], [72, 0], [1, 1, 19, 206,
8, 0, 118, 347, 19, 462, 271, 169, 17, 170, 6, 223, 462, 6, 17], [456, 695, 10], [14, 15, 2, 248], [12, 2, 14], [111, 0, 162], [932,
0, 0, 920, 0, 0, 0, 932], [0, 933, 547, 463], [249, 104], [0], [0, 0], [0, 0], [303, 301, 156, 171, 390], [169, 17], [169, 696, 1
7], [32], [0], [1, 540], [81], [934, 934, 0], [0, 0, 189], [348, 935], [304, 105, 90, 349, 179, 11, 464, 272, 305, 112, 53, 21, 387,
91, 17, 391, 189, 2, 0, 3, 4], [934, 0], [113, 1], [100, 24, 2, 4, 9], [5, 2, 306], [0, 697, 936, 0, 12, 0, 0], [350, 1, 1, 392, 1, 1,
49, 1, 1, 1, 937, 190], [0], [465, 393, 172], [938, 30, 0], [34, 6, 1, 34, 6, 0, 6, 1, 939, 71, 0, 0], [191, 0, 351, 0], [0], [8, 9
0], [940, 548, 0, 2, 11, 302, 0, 2, 172, 941, 29, 540], [942, 0], [943, 144, 0, 0, 192], [0, 11, 0], [698, 394, 699, 65, 0, 250, 2, 2
```

1. 정수 인코딩(Integer Encoding)

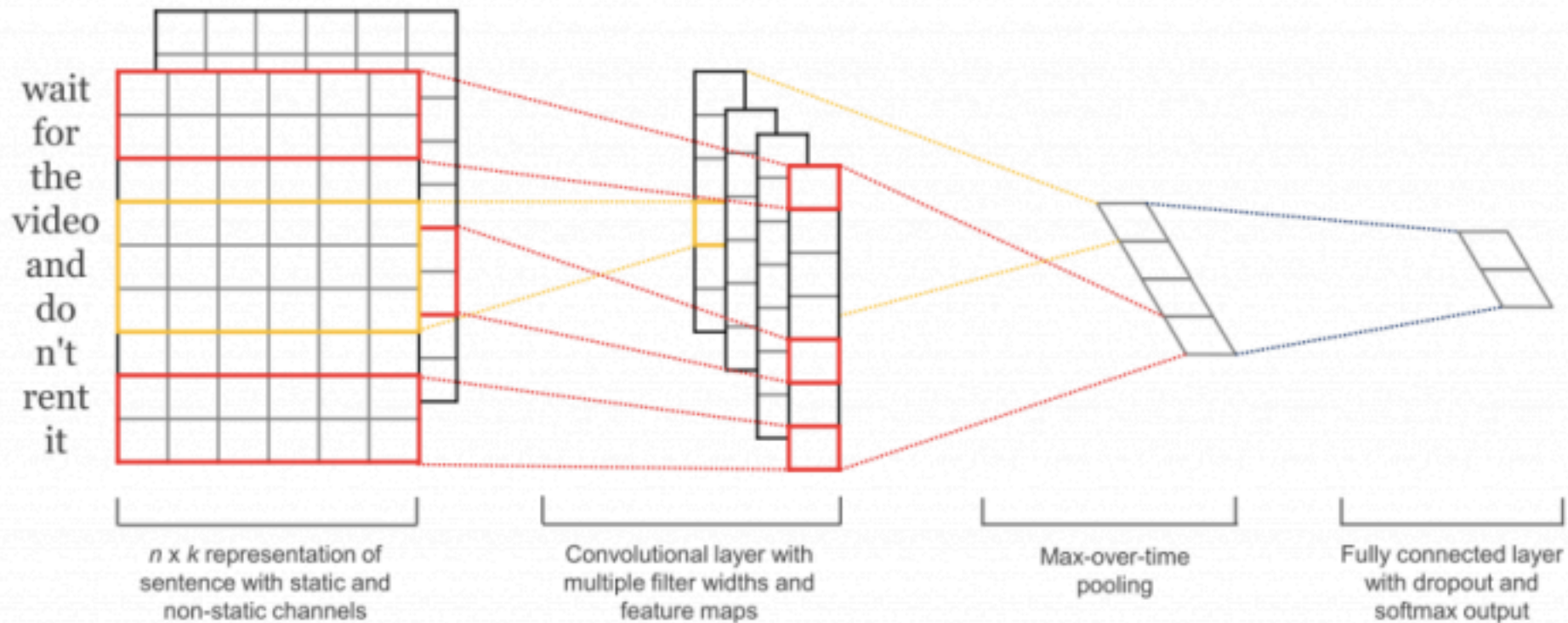
◆ 실습 : 유튜브 동영상(손흥민 푸스카스상) 댓글 정수 인코딩 - 정수 인코딩(Keras tokenizer)

```
In [37]: 1 # 정수 인코딩 (keras tokenizer)
2
3 from tensorflow.keras.preprocessing.text import Tokenizer
4
5 vocab_size = 1060
6 tokenizer = Tokenizer(num_words = vocab_size + 2, oov_token = '00Y')
7 tokenizer.fit_on_texts(top3_tokens[0])
8 print(tokenizer.word_index)
```

{'00Y': 1, 'son': 2, 'goal': 3, 'winner': 4, 'puskas': 5, 'best': 6, 'one': 7, 'player': 8, 'korea': 9, 'award': 10, 'deserve': 11, 'like': 12, 'suarez': 13, 'messi': 14, 'better': 15, 'maradona': 16, 'know': 17, 'world': 18, 'amazing': 19, 'asia': 20, 'arrascaeta': 21, 'football': 22, 'scored': 23, 'fifa': 24, 'great': 25, 'wow': 26, 'even': 27, 'think': 28, 'get': 29, 'run': 30, 'would': 31, 'class': 32, 'lol': 33, 'year': 34, 'nice': 35, 'arsenal': 36, 'ball': 37, 'fan': 38, 'men': 39, 'see': 40, 'solo': 41, '2020': 42, 'ever': 43, 'burnley': 44, 'really': 45, 'vote': 46, 'doubt': 47, 'well': 48, 'meny': 49, 'proud': 50, 'super': 51, 'sead': 52, 'never': 53, 'team': 54, 'every': 55, 'kolašinac': 56, 'much': 57, 'yes': 58, 'could': 59, 'salah': 60, 'dribbling': 61, 'time': 62, 'fast': 63, 'way': 64, 'running': 65, 'got': 66, 'ronaldo': 67, 'beautiful': 68, 'seen': 69, 'newflow': 70, 'absolutely': 71, 'another': 72, 'love': 73, 'must': 74, 'gol': 75, 'skill': 76, 'say': 77, 'times': 78, 'speed': 79, 'league': 80, 'que': 81, 'legend': 82, 'congratulations': 83, 'still': 84, 'past': 85, 'play': 86, 'look': 87, 'definitely': 88, 'people': 89, 'two': 90, 'real': 91, 'whole': 92, 'dribble': 93, 'also': 94, 'last': 95, 'defenders': 96, 'score': 97, 'fantastic': 98, 'dont': 99, 'watch': 100, 'right': 101, 'everyone': 102, 'level': 103, 'overrated': 104, 'video': 105, 'let': 106, 'make': 107, 'give': 108, 'gon': 109, 'assist': 110, 'south': 111, 'crazy': 112, 'first': 113, 'winning': 114, 'made': 115, 'prize': 116, 'wonderful': 117, 'game': 118, 'tottenham': 119, 'easy': 120, 'back': 121, 'incredible': 122, 'kick': 123, 'no': 124, 'far': 125, 'without': 126, 'years': 127, 'spurs': 128, 'yeah': 129, 'sure': 130, 'del': 131, 'nothing': 132, 've': 133, 'defender': 134, 'come': 135, 'soccer': 136, 'dribbled': 137, 'nosonated': 138, 'century': 139, 'puskas': 140, 'crowd': 141, 'history': 142, 'congrats': 143, 'stop': 144, 'long': 145, 'else': 146, 'others': 147, 'agree': 148, 'mean': 149, 'style': 150, 'coys': 151, 'special': 152, 'box': 153, 'watched': 154, 'defending': 155, 'pace': 156, 'field': 157, 'king': 158, 'anybody': 159, 'lifes': 160, 'diane': 161, 'next': 162, 'literally': 163, 'respect': 164, 'top': 165, 'maybe': 166, 'makes': 167, 'robbed': 168, 'tottenham': 169, 'easy': 170, 'winners': 171, 'fan': 172, 'back': 173, 'incredible': 174, 'kick': 175, 'no': 176, 'far': 177, 'without': 178, 'years': 179, 'spurs': 180, 'yeah': 181, 'sure': 182, 'del': 183, 'nothing': 184, 've': 185, 'defender': 186, 'come': 187, 'soccer': 188, 'dribbled': 189, 'nosonated': 190, 'century': 191, 'puskas': 192, 'crowd': 193, 'history': 194, 'congrats': 195, 'stop': 196, 'long': 197, 'else': 198, 'others': 199, 'agree': 200, 'mean': 201, 'style': 202, 'coys': 203, 'special': 204, 'box': 205, 'watched': 206, 'defending': 207, 'pace': 208, 'field': 209, 'king': 210, 'anybody': 211, 'lifes': 212, 'diane': 213, 'next': 214, 'literally': 215, 'respect': 216, 'top': 217, 'maybe': 218, 'makes': 219, 'robbed': 220, 'tottenham': 221, 'easy': 222, 'winners': 223, 'fan': 224, 'back': 225, 'incredible': 226, 'kick': 227, 'no': 228, 'far': 229, 'without': 230, 'years': 231, 'spurs': 232, 'yeah': 233, 'sure': 234, 'del': 235, 'nothing': 236, 've': 237, 'defender': 238, 'come': 239, 'soccer': 240, 'dribbled': 241, 'nosonated': 242, 'century': 243, 'puskas': 244, 'crowd': 245, 'history': 246, 'congrats': 247, 'stop': 248, 'long': 249, 'else': 250, 'others': 251, 'agree': 252, 'mean': 253, 'style': 254, 'coys': 255, 'special': 256, 'box': 257, 'watched': 258, 'defending': 259, 'pace': 260, 'field': 261, 'king': 262, 'anybody': 263, 'lifes': 264, 'diane': 265, 'next': 266, 'literally': 267, 'respect': 268, 'top': 269, 'maybe': 270, 'makes': 271, 'robbed': 272, 'tottenham': 273, 'easy': 274, 'winners': 275, 'fan': 276, 'back': 277, 'incredible': 278, 'kick': 279, 'no': 280, 'far': 281, 'without': 282, 'years': 283, 'spurs': 284, 'yeah': 285, 'sure': 286, 'del': 287, 'nothing': 288, 've': 289, 'defender': 290, 'come': 291, 'soccer': 292, 'dribbled': 293, 'nosonated': 294, 'century': 295, 'puskas': 296, 'crowd': 297, 'history': 298, 'congrats': 299, 'stop': 300, 'long': 301, 'else': 302, 'others': 303, 'agree': 304, 'mean': 305, 'style': 306, 'coys': 307, 'special': 308, 'box': 309, 'watched': 310, 'defending': 311, 'pace': 312, 'field': 313, 'king': 314, 'anybody': 315, 'lifes': 316, 'diane': 317, 'next': 318, 'literally': 319, 'respect': 320, 'top': 321, 'maybe': 322, 'makes': 323, 'robbed': 324, 'tottenham': 325, 'easy': 326, 'winners': 327, 'fan': 328, 'back': 329, 'incredible': 330, 'kick': 331, 'no': 332, 'far': 333, 'without': 334, 'years': 335, 'spurs': 336, 'yeah': 337, 'sure': 338, 'del': 339, 'nothing': 340, 've': 341, 'defender': 342, 'come': 343, 'soccer': 344, 'dribbled': 345, 'nosonated': 346, 'century': 347, 'puskas': 348, 'crowd': 349, 'history': 350, 'congrats': 351, 'stop': 352, 'long': 353, 'else': 354, 'others': 355, 'agree': 356, 'mean': 357, 'style': 358, 'coys': 359, 'special': 360, 'box': 361, 'watched': 362, 'defending': 363, 'pace': 364, 'field': 365, 'king': 366, 'anybody': 367, 'lifes': 368, 'diane': 369, 'next': 370, 'literally': 371, 'respect': 372, 'top': 373, 'maybe': 374, 'makes': 375, 'robbed': 376, 'tottenham': 377, 'easy': 378, 'winners': 379, 'fan': 380, 'back': 381, 'incredible': 382, 'kick': 383, 'no': 384, 'far': 385, 'without': 386, 'years': 387, 'spurs': 388, 'yeah': 389, 'sure': 390, 'del': 391, 'nothing': 392, 've': 393, 'defender': 394, 'come': 395, 'soccer': 396, 'dribbled': 397, 'nosonated': 398, 'century': 399, 'puskas': 400, 'crowd': 401, 'history': 402, 'congrats': 403, 'stop': 404, 'long': 405, 'else': 406, 'others': 407, 'agree': 408, 'mean': 409, 'style': 410, 'coys': 411, 'special': 412, 'box': 413, 'watched': 414, 'defending': 415, 'pace': 416, 'field': 417, 'king': 418, 'anybody': 419, 'lifes': 420, 'diane': 421, 'next': 422, 'literally': 423, 'respect': 424, 'top': 425, 'maybe': 426, 'makes': 427, 'robbed': 428, 'tottenham': 429, 'easy': 430, 'winners': 431, 'fan': 432, 'back': 433, 'incredible': 434, 'kick': 435, 'no': 436, 'far': 437, 'without': 438, 'years': 439, 'spurs': 440, 'yeah': 441, 'sure': 442, 'del': 443, 'nothing': 444, 've': 445, 'defender': 446, 'come': 447, 'soccer': 448, 'dribbled': 449, 'nosonated': 450, 'century': 451, 'puskas': 452, 'crowd': 453, 'history': 454, 'congrats': 455, 'stop': 456, 'long': 457, 'else': 458, 'others': 459, 'agree': 460, 'mean': 461, 'style': 462, 'coys': 463, 'special': 464, 'box': 465, 'watched': 466, 'defending': 467, 'pace': 468, 'field': 469, 'king': 470, 'anybody': 471, 'lifes': 472, 'diane': 473, 'next': 474, 'literally': 475, 'respect': 476, 'top': 477, 'maybe': 478, 'makes': 479, 'robbed': 480, 'tottenham': 481, 'easy': 482, 'winners': 483, 'fan': 484, 'back': 485, 'incredible': 486, 'kick': 487, 'no': 488, 'far': 489, 'without': 490, 'years': 491, 'spurs': 492, 'yeah': 493, 'sure': 494, 'del': 495, 'nothing': 496, 've': 497, 'defender': 498, 'come': 499, 'soccer': 500, 'dribbled': 501, 'nosonated': 502, 'century': 503, 'puskas': 504, 'crowd': 505, 'history': 506, 'congrats': 507, 'stop': 508, 'long': 509, 'else': 510, 'others': 511, 'agree': 512, 'mean': 513, 'style': 514, 'coys': 515, 'special': 516, 'box': 517, 'watched': 518, 'defending': 519, 'pace': 520, 'field': 521, 'king': 522, 'anybody': 523, 'lifes': 524, 'diane': 525, 'next': 526, 'literally': 527, 'respect': 528, 'top': 529, 'maybe': 530, 'makes': 531, 'robbed': 532, 'tottenham': 533, 'easy': 534, 'winners': 535, 'fan': 536, 'back': 537, 'incredible': 538, 'kick': 539, 'no': 540, 'far': 541, 'without': 542, 'years': 543, 'spurs': 544, 'yeah': 545, 'sure': 546, 'del': 547, 'nothing': 548, 've': 549, 'defender': 550, 'come': 551, 'soccer': 552, 'dribbled': 553, 'nosonated': 554, 'century': 555, 'puskas': 556, 'crowd': 557, 'history': 558, 'congrats': 559, 'stop': 560, 'long': 561, 'else': 562, 'others': 563, 'agree': 564, 'mean': 565, 'style': 566, 'coys': 567, 'special': 568, 'box': 569, 'watched': 570, 'defending': 571, 'pace': 572, 'field': 573, 'king': 574, 'anybody': 575, 'lifes': 576, 'diane': 577, 'next': 578, 'literally': 579, 'respect': 580, 'top': 581, 'maybe': 582, 'makes': 583, 'robbed': 584, 'tottenham': 585, 'easy': 586, 'winners': 587, 'fan': 588, 'back': 589, 'incredible': 590, 'kick': 591, 'no': 592, 'far': 593, 'without': 594, 'years': 595, 'spurs': 596, 'yeah': 597, 'sure': 598, 'del': 599, 'nothing': 600, 've': 601, 'defender': 602, 'come': 603, 'soccer': 604, 'dribbled': 605, 'nosonated': 606, 'century': 607, 'puskas': 608, 'crowd': 609, 'history': 610, 'congrats': 611, 'stop': 612, 'long': 613, 'else': 614, 'others': 615, 'agree': 616, 'mean': 617, 'style': 618, 'coys': 619, 'special': 620, 'box': 621, 'watched': 622, 'defending': 623, 'pace': 624, 'field': 625, 'king': 626, 'anybody': 627, 'lifes': 628, 'diane': 629, 'next': 630, 'literally': 631, 'respect': 632, 'top': 633, 'maybe': 634, 'makes': 635, 'robbed': 636, 'tottenham': 637, 'easy': 638, 'winners': 639, 'fan': 640, 'back': 641, 'incredible': 642, 'kick': 643, 'no': 644, 'far': 645, 'without': 646, 'years': 647, 'spurs': 648, 'yeah': 649, 'sure': 650, 'del': 651, 'nothing': 652, 've': 653, 'defender': 654, 'come': 655, 'soccer': 656, 'dribbled': 657, 'nosonated': 658, 'century': 659, 'puskas': 660, 'crowd': 661, 'history': 662, 'congrats': 663, 'stop': 664, 'long': 665, 'else': 666, 'others': 667, 'agree': 668, 'mean': 669, 'style': 670, 'coys': 671, 'special': 672, 'box': 673, 'watched': 674, 'defending': 675, 'pace': 676, 'field': 677, 'king': 678, 'anybody': 679, 'lifes': 680, 'diane': 681, 'next': 682, 'literally': 683, 'respect': 684, 'top': 685, 'maybe': 686, 'makes': 687, 'robbed': 688, 'tottenham': 689, 'easy': 690, 'winners': 691, 'fan': 692, 'back': 693, 'incredible': 694, 'kick': 695, 'no': 696, 'far': 697, 'without': 698, 'years': 699, 'spurs': 700, 'yeah': 701, 'sure': 702, 'del': 703, 'nothing': 704, 've': 705, 'defender': 706, 'come': 707, 'soccer': 708, 'dribbled': 709, 'nosonated': 710, 'century': 711, 'puskas': 712, 'crowd': 713, 'history': 714, 'congrats': 715, 'stop': 716, 'long': 717, 'else': 718, 'others': 719, 'agree': 720, 'mean': 721, 'style': 722, 'coys': 723, 'special': 724, 'box': 725, 'watched': 726, 'defending': 727, 'pace': 728, 'field': 729, 'king': 730, 'anybody': 731, 'lifes': 732, 'diane': 733, 'next': 734, 'literally': 735, 'respect': 736, 'top': 737, 'maybe': 738, 'makes': 739, 'robbed': 740, 'tottenham': 741, 'easy': 742, 'winners': 743, 'fan': 744, 'back': 745, 'incredible': 746, 'kick': 747, 'no': 748, 'far': 749, 'without': 750, 'years': 751, 'spurs': 752, 'yeah': 753, 'sure': 754, 'del': 755, 'nothing': 756, 've': 757, 'defender': 758, 'come': 759, 'soccer': 760, 'dribbled': 761, 'nosonated': 762, 'century': 763, 'puskas': 764, 'crowd': 765, 'history': 766, 'congrats': 767, 'stop': 768, 'long': 769, 'else': 770, 'others': 771, 'agree': 772, 'mean': 773, 'style': 774, 'coys': 775, 'special': 776, 'box': 777, 'watched': 778, 'defending': 779, 'pace': 780, 'field': 781, 'king': 782, 'anybody': 783, 'lifes': 784, 'diane': 785, 'next': 786, 'literally': 787, 'respect': 788, 'top': 789, 'maybe': 790, 'makes': 791, 'robbed': 792, 'tottenham': 793, 'easy': 794, 'winners': 795, 'fan': 796, 'back': 797, 'incredible': 798, 'kick': 799, 'no': 800, 'far': 801, 'without': 802, 'years': 803, 'spurs': 804, 'yeah': 805, 'sure': 806, 'del': 807, 'nothing': 808, 've': 809, 'defender': 810, 'come': 811, 'soccer': 812, 'dribbled': 813, 'nosonated': 814, 'century': 815, 'puskas': 816, 'crowd': 817, 'history': 818, 'congrats': 819, 'stop': 820, 'long': 821, 'else': 822, 'others': 823, 'agree': 824, 'mean': 825, 'style': 826, 'coys': 827, 'special': 828, 'box': 829, 'watched': 830, 'defending': 831, 'pace': 832, 'field': 833, 'king': 834, 'anybody': 835, 'lifes': 836, 'diane': 837, 'next': 838, 'literally': 839, 'respect': 840, 'top': 841, 'maybe': 842, 'makes': 843, 'robbed': 844, 'tottenham': 845, 'easy': 846, 'winners': 847, 'fan': 848, 'back': 849, 'incredible': 850, 'kick': 851, 'no': 852, 'far': 853, 'without': 854, 'years': 855, 'spurs': 856, 'yeah': 857, 'sure': 858, 'del': 859, 'nothing': 860, 've': 861, 'defender': 862, 'come': 863, 'soccer': 864, 'dribbled': 865, 'nosonated': 866, 'century': 867, 'puskas': 868, 'crowd': 869, 'history': 870, 'congrats': 871, 'stop': 872, 'long': 873, 'else': 874, 'others': 875, 'agree': 876, 'mean': 877, 'style': 878, 'coys': 879, 'special': 880, 'box': 881, 'watched': 882, 'defending': 883, 'pace': 884, 'field': 885, 'king': 886, 'anybody': 887, 'lifes': 888, 'diane': 889, 'next': 890, 'literally': 891, 'respect': 892, 'top': 893, 'maybe': 894, 'makes': 895, 'robbed': 896, 'tottenham': 897, 'easy': 898, 'winners': 899, 'fan': 900, 'back': 901, 'incredible': 902, 'kick': 903, 'no': 904, 'far': 905, 'without': 906, 'years': 907, 'spurs': 908, 'yeah': 909, 'sure': 910, 'del': 911, 'nothing': 912, 've': 913, 'defender': 914, 'come': 915, 'soccer': 916, 'dribbled': 917, 'nosonated': 918, 'century': 919, 'puskas': 920, 'crowd': 921, 'history': 922, 'congrats': 923, 'stop': 924, 'long': 925, 'else': 926, 'others': 927, 'agree': 928, 'mean': 929, 'style': 930, 'coys': 931, 'special': 932, 'box': 933, 'watched': 934, 'defending': 935, 'pace': 936, 'field': 937, 'king': 938, 'anybody': 939, 'lifes': 940, 'diane': 941, 'next': 942, 'literally': 943, 'respect': 944, 'top': 945, 'maybe': 946, 'makes': 947, 'robbed': 948, 'tottenham': 949, 'easy': 950, 'winners': 951, 'fan': 952, 'back': 953, 'incredible': 954, 'kick': 955, 'no': 956, 'far': 957, 'without': 958, 'years': 959, 'spurs': 960, 'yeah': 961, 'sure': 962, 'del': 963, 'nothing': 964, 've': 965, 'defender': 966, 'come': 967, 'soccer': 968, 'dribbled': 969, 'nosonated': 970, 'century': 971, 'puskas': 972, 'crowd': 973, 'history': 974, 'congrats': 975, 'stop': 976, 'long': 977, 'else': 978, 'others': 979, 'agree': 980, 'mean': 981, 'style': 982, 'coys': 983, 'special': 984, 'box': 985, 'watched': 986, 'defending': 987, 'pace': 988, 'field': 989, 'king': 990, 'anybody': 991, 'lifes': 992, 'diane': 993, 'next': 994, 'literally': 995, 'respect': 996, 'top': 997, 'maybe': 998, 'makes': 999, 'robbed': 1000, 'tottenham': 1001, 'easy': 1002, 'winners': 1003, 'fan': 1004, 'back': 1005, 'incredible': 1006, 'kick': 1007, 'no': 1008, 'far': 1009, 'without': 1010, 'years': 1011, 'spurs': 1012, 'yeah': 1013, 'sure': 1014, 'del': 1015, 'nothing': 1016, 've': 1017, 'defender': 1018, 'come': 1019, 'soccer': 1020, 'dribbled': 1021, 'nosonated': 1022, 'century': 1023, 'puskas': 1024, 'crowd': 1025, 'history': 1026, 'congrats': 1027, 'stop': 1028, 'long': 1029, 'else': 1030, 'others': 1031, 'agree': 1032, 'mean': 1033, 'style': 1034, 'coys': 1035, 'special': 1036, 'box': 1037, 'watched': 1038, 'defending': 1039, 'pace': 1040, 'field': 1041, 'king': 1042, 'anybody': 1043, 'lifes': 1044, 'diane': 1045, 'next': 1046, 'literally': 1047, 'respect': 1048, 'top': 1049, 'maybe': 1050, 'makes': 1051, 'robbed': 1052, 'tottenham': 1053, 'easy': 1054, 'winners': 1055, 'fan': 1056, 'back': 1057, 'incredible': 1058, 'kick': 1059, 'no': 1060, 'far': 1061, 'without': 1062, 'years': 1063, 'spurs': 1064, 'yeah': 1065, 'sure': 1066, 'del': 1067, 'nothing': 1068, 've': 1069, 'defender': 1070, 'come': 1071, 'soccer': 1072, 'dribbled': 1073, 'nosonated': 1074, 'century': 1075, 'puskas': 1076, 'crowd': 1077, 'history': 1078, 'congrats': 1079, 'stop': 1080, 'long': 1081, 'else': 1082, 'others': 1083, 'agree': 1084, 'mean': 1085, 'style': 1086, 'coys': 1087, 'special': 1088, 'box': 1089, 'watched': 1090, 'defending': 1091, 'pace': 1092, 'field': 1093, 'king': 1094, 'anybody': 1095, 'lifes': 1096, 'diane': 1097, 'next': 1098, 'literally': 1099, 'respect': 1100, 'top': 1101, 'maybe': 1102, 'makes': 1103, 'robbed': 1104, 'tottenham': 1105, 'easy': 1106, 'winners': 1107, 'fan': 1108, 'back': 1109, 'incredible': 1110, 'kick': 1111, 'no': 1112, 'far': 1113, 'without': 1114, 'years': 1115, 'spurs': 1116, 'yeah': 1117, 'sure': 1118, 'del': 1119, 'nothing': 1120, 've': 1121, 'defender': 1122, 'come': 1123, 'soccer': 1124, 'dribbled': 1125, 'nosonated': 1126, 'century': 1127, 'puskas': 1128, 'crowd': 1129, 'history': 1130, 'congrats': 1131, 'stop': 1132, 'long': 1133, 'else': 1134, 'others': 1135, 'agree': 1136, 'mean': 1137, 'style': 1138, 'coys': 1139, 'special': 1140, 'box': 1141, 'watched': 1142, 'defending': 1143, 'pace': 1144, 'field': 1145, 'king': 1146, 'anybody': 1147, 'lifes': 1148, 'diane': 1149, 'next': 1150, 'literally': 1151, 'respect': 1152, 'top': 1153, 'maybe': 1154, 'makes': 1155, 'robbed': 1156, 'tottenham': 1157, 'easy': 1158, 'winners': 1159, 'fan': 1160, 'back': 1161, 'incredible': 1162, 'kick': 1163, 'no': 1164, 'far': 1165, 'without': 1166, 'years': 1167, 'spurs': 1168, 'yeah': 1169, 'sure': 1170, 'del': 1171, 'nothing': 1172, 've': 1173, 'defender': 1174, 'come': 1175, 'soccer': 1176, 'dribbled': 1177, 'nosonated': 1178, 'century': 1179, 'puskas': 1180, 'crowd': 1181, 'history': 1182, 'congrats': 1183, 'stop': 1184, 'long': 1185, 'else': 1186, 'others': 1187, 'agree': 1188, 'mean': 1189, 'style': 1190, 'coys': 1191, 'special': 1192, 'box': 1193, 'watched': 1194, 'defending': 1195, 'pace': 1196, 'field': 1197, 'king': 1198, 'anybody': 1199, 'lifes': 1200, 'diane': 1201, 'next': 1202, 'literally': 1203, 'respect': 1204, 'top': 1205, 'maybe': 1206, 'makes': 1207, 'robbed': 1208, 'tottenham': 1209, 'easy': 1210, 'winners': 1211, 'fan': 1212, 'back': 1213, 'incredible': 1214, 'kick': 1215, 'no': 1216, 'far': 1217, 'without': 1218, 'years': 1219, 'spurs': 1220, 'yeah': 1221, 'sure': 1222, 'del': 1223, 'nothing': 1224, 've': 1225, 'defender': 1226, 'come': 1227, 'soccer': 1228, 'dribbled': 1229, 'nosonated': 1230, 'century': 1231, 'puskas': 1232, 'crowd': 1233, 'history': 1234, 'congrats': 1235, 'stop': 1236, 'long': 1237, 'else': 1238, 'others': 1239, 'agree': 1240, 'mean': 1241, 'style': 1242, 'coys': 1243, 'special': 1244, 'box': 1245, 'watched': 1246, 'defending': 1247, 'pace': 1248, 'field': 1249, 'king': 1250, 'anybody': 1251, 'lifes': 1252, 'diane': 1253, 'next': 1254, 'literally': 1255, 'respect': 1256, 'top': 1257, 'maybe': 1258, 'makes': 1259, 'robbed': 1260, 'tottenham': 1261, 'easy': 1262, 'winners': 1263, 'fan': 1264, 'back': 1265, 'incredible': 1266, 'kick': 1267, 'no': 1268, 'far': 1269, 'without': 1270, 'years': 1271, 'spurs': 1272, 'yeah': 1273, 'sure': 1274, 'del': 1275, 'nothing': 1276, 've': 1277, 'defender': 1278, 'come': 1279, 'soccer': 1280, 'dribbled': 1281, 'nosonated': 1282, 'century': 1283, 'puskas': 1284, 'crowd': 1285, 'history': 1286, 'congrats': 1287, 'stop': 1288, 'long': 1289, 'else': 1290, 'others': 1291, 'agree': 1292, 'mean': 1293, 'style': 1294, 'coys': 1295, 'special': 1296, 'box': 1297, 'watched': 1298, 'defending': 1299, 'pace': 1300, 'field': 1301, 'king': 1302, 'anybody': 1303, 'lifes': 1304, 'diane': 1305, 'next': 1306, 'literally': 1307, 'respect': 1308, 'top': 1309, 'maybe': 1310, 'makes': 1311, 'robbed': 1312, 'tottenham': 1313, 'easy': 1314, 'winners': 1315, 'fan': 1316, 'back': 1317, 'incredible': 1318, 'kick': 1319, 'no': 1320, 'far': 1321, 'without': 1322, 'years': 1323, 'spurs': 1324, 'yeah': 1325, 'sure': 1326, 'del': 1327, 'nothing': 1328, 've': 1329, 'defender': 1330, 'come': 1331, 'soccer': 1332, 'dribbled': 1333, 'nosonated': 1334, 'century': 1335, 'puskas': 1336, 'crowd': 1337, 'history': 1338, 'congrats': 1339, 'stop': 1340, 'long': 1341, 'else': 1342, 'others': 1343, 'agree': 1344, 'mean': 1345, 'style': 1346, 'coys': 1347, 'special': 1348, 'box': 1349, 'watched': 1350, 'defending': 1351, 'pace': 1352, 'field': 1353, 'king': 1354, 'anybody': 1355, 'lifes': 1356, 'diane': 1357, 'next': 1358, 'literally': 1359, 'respect': 1360, 'top': 1361, 'maybe': 1362, 'makes': 1363, 'robbed': 1364, 'tottenham': 1365, 'easy': 1366, 'winners': 1367, 'fan': 1368, 'back': 1369, 'incredible': 1370, 'kick': 1371, 'no': 1372, 'far': 1373, 'without': 1374, 'years': 1375, 'spurs': 1376, 'yeah': 1377, 'sure': 1378, 'del': 1379, 'nothing': 1380, 've': 1381, 'defender': 1382, 'come': 1383, 'soccer': 1384, 'dribbled': 1385, 'nosonated': 1386, 'century': 1387, 'puskas': 1388, 'crowd': 1389, 'history': 1390, 'congrats': 1391, 'stop': 1392, 'long': 1393, 'else': 1394, 'others': 1395, 'agree': 1396, 'mean': 1397, 'style': 1398, 'coys': 1399, 'special': 1400, 'box': 1401, 'watched': 1402, 'defending': 1403, 'pace': 1404, 'field': 1405, 'king':

2. 패딩(Padding)

- ◆ 왜 패딩을 하지?
자연어 처리시 동일 길이를 묶어서 연산하기 좋음



2. 패딩(Padding)

◆ Numpy로 패딩

```
import numpy as np
from tensorflow.keras.preprocessing.text import Tokenizer
```

```
sentences = [['barber', 'person'], ['barber', 'good', 'person'], ['barber', 'huge', 'person'], ['knew', 'secret'], ['secret', 'kept', 'huge', 'secret'], ['huge', 'secret'], ['barber', 'kept', 'word'], ['barber', 'kept', 'word'], ['barber', 'kept', 'secret'], ['keeping', 'keeping', 'huge', 'secret', 'driving', 'barber', 'crazy'], ['barber', 'went', 'huge', 'mountain']]
```

```
tokenizer = Tokenizer()
tokenizer.fit_on_texts(sentences) # fit_on_texts()안에 코퍼스를 입력으로 하면 빈도수를 기준으로 단어 집합을 생성한다.
```

```
encoded = tokenizer.texts_to_sequences(sentences)
print(encoded)
```

```
[[1, 5], [1, 8, 5], [1, 3, 5], [9, 2], [2, 4, 3, 2], [3, 2], [1, 4, 6], [1, 4, 6], [1, 4, 2], [7, 7, 3, 2, 10, 1, 1], [1, 12, 3, 13]]
```

2. 패딩(Padding)

◆ Numpy로 패딩

```
max_len = max(len(item) for item in encoded)
print(max_len)
```

```
for item in encoded: # 각 문장에 대해서
    while len(item) < max_len: # max_len보다 작으면
        item.append(0)
```

```
padded_np = np.array(encoded)
padded_np
```

```
array([[ 1,  5,  0,  0,  0,  0,  0],
       [ 1,  8,  5,  0,  0,  0,  0],
       [ 1,  3,  5,  0,  0,  0,  0],
       [ 9,  2,  0,  0,  0,  0,  0],
       [ 2,  4,  3,  2,  0,  0,  0],
       [ 3,  2,  0,  0,  0,  0,  0],
       [ 1,  4,  6,  0,  0,  0,  0],
       [ 1,  4,  6,  0,  0,  0,  0],
       [ 1,  4,  2,  0,  0,  0,  0],
       [ 7,  7,  3,  2, 10,  1, 11],
       [ 1, 12,  3, 13,  0,  0,  0]])
```

2. 패딩(Padding)

◆ 케라스 전처리 도구로 패딩하기

```
from tensorflow.keras.preprocessing.sequence import pad_sequences
```

```
encoded = tokenizer.texts_to_sequences(sentences)
print(encoded)
```

```
[[1, 5], [1, 8, 5], [1, 3, 5], [9, 2], [2, 4, 3, 2], [3, 2], [1, 4, 6], [1, 4, 6], [1, 4, 2], [7, 7, 3, 2, 10, 1, 1], [1, 12, 3, 13]]
```

```
padded = pad_sequences(encoded)
padded
```

```
array([[ 0,  0,  0,  0,  0,  1,  5],
       [ 0,  0,  0,  0,  1,  8,  5],
       [ 0,  0,  0,  0,  1,  3,  5],
       [ 0,  0,  0,  0,  0,  9,  2],
       [ 0,  0,  0,  2,  4,  3,  2],
       [ 0,  0,  0,  0,  0,  3,  2],
       [ 0,  0,  0,  0,  1,  4,  6],
       [ 0,  0,  0,  0,  1,  4,  6],
       [ 0,  0,  0,  0,  1,  4,  2],
       [ 7,  7,  3,  2, 10,  1, 11],
       [ 0,  0,  0,  1, 12,  3, 13]], dtype=int32)
```

2. 패딩(Padding)

◆ 케라스 전처리 도구로 패딩하기 (post padding)

```
padded = pad_sequences(encoded, padding = 'post')
```

```
padded
```

```
array([[ 1,  5,  0,  0,  0,  0,  0],
       [ 1,  8,  5,  0,  0,  0,  0],
       [ 1,  3,  5,  0,  0,  0,  0],
       [ 9,  2,  0,  0,  0,  0,  0],
       [ 2,  4,  3,  2,  0,  0,  0],
       [ 3,  2,  0,  0,  0,  0,  0],
       [ 1,  4,  6,  0,  0,  0,  0],
       [ 1,  4,  6,  0,  0,  0,  0],
       [ 1,  4,  2,  0,  0,  0,  0],
       [ 7,  7,  3,  2, 10,  1, 11],
       [ 1, 12,  3, 13,  0,  0,  0]], dtype=int32)
```

```
(padded == padded_np).all()
```

2. 패딩(Padding)

◆ 케라스 전처리 도구로 패딩하기 (max len)

```
padded = pad_sequences(encoded, padding = 'post', maxlen = 5)  
padded
```

```
array([[ 1,  5,  0,  0,  0],  
       [ 1,  8,  5,  0,  0],  
       [ 1,  3,  5,  0,  0],  
       [ 9,  2,  0,  0,  0],  
       [ 2,  4,  3,  2,  0],  
       [ 3,  2,  0,  0,  0],  
       [ 1,  4,  6,  0,  0],  
       [ 1,  4,  6,  0,  0],  
       [ 1,  4,  2,  0,  0],  
       [ 3,  2, 10,  1, 11],  
       [ 1, 12,  3, 13,  0]], dtype=int32)
```

2. 패딩(Padding)

◆ 케라스 전처리 도구로 패딩하기 (last index)

```
last_value = len(tokenizer.word_index) + 1 # 단어 집합의 크기보다 1 큰 숫자를 사용
print(last_value)
```

```
padded = pad_sequences(encoded, padding = 'post', value = last_value)
padded
```

```
array([[ 1,  5, 14, 14, 14, 14, 14],
       [ 1,  8,  5, 14, 14, 14, 14],
       [ 1,  3,  5, 14, 14, 14, 14],
       [ 9,  2, 14, 14, 14, 14, 14],
       [ 2,  4,  3,  2, 14, 14, 14],
       [ 3,  2, 14, 14, 14, 14, 14],
       [ 1,  4,  6, 14, 14, 14, 14],
       [ 1,  4,  6, 14, 14, 14, 14],
       [ 1,  4,  2, 14, 14, 14, 14],
       [ 7,  7,  3,  2, 10,  1, 11],
       [ 1, 12,  3, 13, 14, 14, 14]], dtype=int32)
```

3. 원-핫 인코딩(One-Hot Encoding)

- ◆ 컴퓨터 연산 효율을 위해 문자 보다는 숫자로 처리
- ◆ Vocabulary: 단어 집합
- ◆ 원-핫 인코딩: 표현하고 싶은 단어 인덱스는 1, 그 외 단어는 0으로 표현, 원-핫 벡터(One-Hot vector)

3. 원-핫 인코딩(One-Hot Encoding)

```
from konlpy.tag import Okt
okt=Okt()
token=okt.morphs("나는 자연어 처리를 배운다")
print(token)
```

```
['나', '는', '자연어', '처리', '를', '배운다']
```

```
def one_hot_encoding(word, word2index):
    one_hot_vector = [0]*(len(word2index))
    index=word2index[word]
    one_hot_vector[index]=1
    return one_hot_vector
```

```
one_hot_encoding("자연어",word2index)
```

```
[0, 0, 1, 0, 0, 0]
```

3. 원-핫 인코딩(One-Hot Encoding)

◆ 케라스를 이용한 원-핫 인코딩

```
text="나랑 점심 먹으러 갈래 점심 메뉴는 햄버거 갈래 갈래 햄버거 최고야"
```

```
from tensorflow.keras.preprocessing.text import Tokenizer  
from tensorflow.keras.utils import to_categorical
```

```
text="나랑 점심 먹으러 갈래 점심 메뉴는 햄버거 갈래 갈래 햄버거 최고야"
```

```
t = Tokenizer()  
t.fit_on_texts([text])  
print(t.word_index) # 각 단어에 대한 인코딩 결과 출력.
```

```
{'갈래': 1, '점심': 2, '햄버거': 3, '나랑': 4, '먹으러': 5, '메뉴는': 6, '최고야': 7}
```

```
sub_text="점심 먹으러 갈래 메뉴는 햄버거 최고야"  
encoded=t.texts_to_sequences([sub_text])[0]  
print(encoded)
```

```
[2, 5, 1, 6, 3, 7]
```

3. 원-핫 인코딩(One-Hot Encoding)

◆ 케라스를 이용한 원-핫 인코딩

```
one_hot = to_categorical(encoded)
print(one_hot)
```

```
[[0. 0. 1. 0. 0. 0. 0. 0.] #인덱스 2의 원-핫 벡터
 [0. 0. 0. 0. 0. 1. 0. 0.] #인덱스 5의 원-핫 벡터
 [0. 1. 0. 0. 0. 0. 0. 0.] #인덱스 1의 원-핫 벡터
 [0. 0. 0. 0. 0. 0. 1. 0.] #인덱스 6의 원-핫 벡터
 [0. 0. 0. 1. 0. 0. 0. 0.] #인덱스 3의 원-핫 벡터
 [0. 0. 0. 0. 0. 0. 0. 1.] #인덱스 7의 원-핫 벡터]
```

3. 원-핫 인코딩(One-Hot Encoding): 한계

- ◆ 단어의 갯수 만큼 저장 공간 필요
- ◆ 단어 유사도를 표현하지 못함
- ◆ 카운트 기반 벡터화 방법: LSA, HAL
- ◆ 예측 기반 벡터화: NNLM, RNNLM, Word2Vec, FastText
- ◆ 카운트 기반 + 예측 기반 벡터화: GloVe

4. 데이터의 분리(Splitting Data)

◆ 지도 학습

정답지[스팸 여부]

문제지[메일 내용]

- 훈련데이터(training set)을 통해 컴퓨터를 학습시키고, 이를 바탕으로 새로운 데이터(test data)를 분류/예측하는 방식
- 훈련 데이터에 '레이블(label) - 원하는 답'이 포함됨

받은편지함 메일보기 스팸편지함

전체 선택 답장 전체 답장 전달 읽음 표시 휴지통으로 편지함 정리 편지 이동 수신 거부

보낸 사람

날짜: 2019.11.14

날짜	메일 종류	모드	전송 결과	제목	발신자	크기	필터링 정보
13 16:44:57	스팸	차단		[광고] Creative Clo...	adobelapacmid...	33.0 K	제목, (광고), 포...
13 13:51:22	스팸	차단		[광고][학지사_인파...	noreply@hakjisa...	11.3 K	제목, (광고), 포...
13 09:27:31	스팸	차단		[광고][웨비나] 201...	no-reply@direct...	12.3 K	제목, (광고), 포...
12 23:56:40	스팸	차단		Read Attached Omi...	info@lwypolnocy...	197.9 K	발신자 도메인, i...
12 10:20:24	스팸	차단		[광고][한국학술경...	returnmail@cheo...	7.0 K	제목, (광고), 포...
12 09:17:30	스팸	차단		[광고][인싸이트] ...	noreply@hakjisa...	15.5 K	제목, (광고), 포...

Total : 6개

삭제 파일로 저장 메일 복구 허용 차단

받은편지함 메일보기 스팸편지함

전체 선택 답장 전체 답장 전달 읽음 표시 휴지통으로 편지함 정리 편지 이동 수신 거부

보낸 사람

날짜: 2019.11.14

날짜	메일 종류	모드	전송 결과	제목	발신자	크기	필터링 정보
13 16:44:57	스팸	차단		[광고][학지사_인파...] 신규 상품 3종 출시 기념 이벤트 - 선착순	adobelapacmid...	33.0 K	제목, (광고), 포...
13 13:51:22	스팸	차단		[광고][학지사_인파...	noreply@hakjisa...	11.3 K	제목, (광고), 포...
13 09:27:31	스팸	차단		[광고][웨비나] 201...	no-reply@direct...	12.3 K	제목, (광고), 포...
12 23:56:40	스팸	차단		Read Attached Omi...	info@lwypolnocy...	197.9 K	발신자 도메인, i...
12 10:20:24	스팸	차단		[광고][한국학술경...	returnmail@cheo...	7.0 K	제목, (광고), 포...
12 09:17:30	스팸	차단		[광고][인싸이트] ...	noreply@hakjisa...	15.5 K	제목, (광고), 포...

“어떻게 새 메일이 ‘스팸’으로 분류되는가?”

4. 데이터의 분리(Splitting Data)

(예시) 20,000개의 행을 가진 데이터의 분리

텍스트(메일의 내용) x	레이블(스팸 여부) y
당신에게 드리는 마지막 혜택! ...	스팸 메일
내일 볼 수 있을지 확인 부탁...	정상 메일
...	...
(광고) 멋있어질 수 있는...	스팸 메일



<훈련 데이터> [학습용: 18,000개]

X_train : 문제지 데이터

y_train : 문제지에 대한 정답 데이터

<테스트 데이터> [시험용: 2,000개]

X_test : 시험지 데이터.

y_test : 시험지에 대한 정답 데이터.

- 문제지를 다 공부한 후 실력을 평가하기 위해 테스트용으로 일부 문제와 정답지를 빼놓는 것! (여기에서는 2,000개를 분리한다고 가정)
 - 분리시 x(문제지)와 y(정답지)의 맵핑 관계를 유지해야 함
- 기계는 훈련 데이터를 보면서 어떤 내용일 때 정상/스팸 메일인지 규칙 도출 → 학습 종료 후 X_test에 대해 정답 예측
- 기계가 예측한 답 vs. 실제 정답인 y_test 비교 → 기계의 **정확도(Accuracy)**

4. 데이터의 분리(Splitting Data)

◆ X와 Y 분리하기

1) zip 함수를 이용하는 방법

```
x,y = zip(['a',1], ['b', 2], ['c', 3])
```

```
print(x)
```

```
('a', 'b', 'c')
```

```
print(y)
```

```
(1, 2, 3)
```

- zip(): 동일한 개수를 갖는 시퀀스 자료형에서 각 순서에 등장하는 원소들끼리 묶어주는 역할
→ 리스트의 리스트 구성에서 x와 y 분리시 유용

4. 데이터의 분리(Splitting Data)

2) 데이터 프레임 이용하기

```
import pandas as pd
```

```
values = [['당신에게 드리는 마지막 혜택!', 1],  
          ['내일 볼 수 있을지 확인 부탁드립니다...', 0],  
          ['도연씨. 잘 지내시죠? 오랜만입...', 0],  
          ['(광고) AI로 주가를 예측할 수 있다!', 1]]  
columns = ['메일 본문', '스팸 메일 유무']
```

```
df = pd.DataFrame(values, columns=columns)
```

```
df
```

메일 본문 스팸 메일 유무

0	당신에게 드리는 마지막 혜택!	1
1	내일 볼 수 있을지 확인 부탁드립니다...	0
2	도연씨. 잘 지내시죠? 오랜만입...	0
3	(광고) AI로 주가를 예측할 수 있다!	1

분리

```
x=df['메일 본문']
```

```
y=df['스팸 메일 유무']
```

```
print(x)
```

```
0      당신에게 드리는 마지막 혜택!  
1      내일 볼 수 있을지 확인 부탁드립니다...  
2      도연씨. 잘 지내시죠? 오랜만입...  
3      (광고) AI로 주가를 예측할 수 있다!  
Name: 메일 본문, dtype: object
```

```
print(y)
```

```
0    1  
1    0  
2    0  
3    1  
Name: 스팸 메일 유무, dtype: int64
```

- 데이터 프레임(df)은 열의 이름을 통해 각 열에 접근이 가능함을 이용

4. 데이터의 분리(Splitting Data)

3) Numpy를 이용하는 방법

```
import numpy as np
ar = np.arange(0,16).reshape((4,4))
print(ar)
```

```
[[ 0  1  2  3]
 [ 4  5  6  7]
 [ 8  9 10 11]
 [12 13 14 15]]
```

```
x=ar[:, :3]
print(x)
```

```
[[ 0  1  2]
 [ 4  5  6]
 [ 8  9 10]
 [12 13 14]]
```

```
y=ar[:,3]
print(y)
```

```
[ 3  7 11 15]
```

4. 데이터의 분리(Splitting Data)

◆ 테스트 데이터 분리하기

1) 사이킷런을 이용하는 방법

```
import numpy as np
from sklearn.model_selection import train_test_split
x, y = np.arange(10).reshape((5, 2)), range(5) # 실습을 위해 임의로 x와 y가 이미 분리된 데이터 생성
print(x)
print(list(y)) # 레이블 데이터
```

```
[[0 1]
 [2 3]
 [4 5]
 [6 7]
 [8 9]]
[0, 1, 2, 3, 4]
```

1/3만 테스트 데이터로 지정(test_size, train_size 중 1개만 기재 가능)

```
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.33, random_state=1234)
print(x_train)
print(x_test)    x : 배열 또는 데이터프레임
```

랜덤하게 훈련 데이터와 테스트 데이터 분리

```
print(y_train)
print(y_test)    y : 레이블

[[2 3]
 [4 5]
 [6 7]]
[[8 9]
 [0 1]]

[1, 2, 3]
[4, 0]
```

- `Train_test_split()`: 학습용 데이터와 테스트용 데이터 분리 함수

4. 데이터의 분리(Splitting Data)

2) 수동으로 분리하는 방법

```
import numpy as np
x, y = np.arange(0,24).reshape((12,2)), range(12)
print(x)
print(list(y))
```

```
[[ 0  1]
 [ 2  3]
 [ 4  5]
 [ 6  7]
 [ 8  9]
 [10 11]
 [12 13]
 [14 15]
 [16 17]
 [18 19]
 [20 21]
 [22 23]]
[0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11]
```

```
n_of_train = int(len(x) * 0.8) # 데이터의 전체 길이의 80%에 해당하는 길이값 산출
n_of_test = int(len(x) - n_of_train) # 전체 길이에서 80%에 해당하는 길이를 뺀 값 산출
# (소수점 처리로 인한 데이터 누락 방지 위해 0.2를 곱하지 않음에 유의)

print(n_of_train)
print(n_of_test)
```

```
9
3
```

분리

```
print(x_test)
print(list(y_test))
```

```
[[18 19]
 [20 21]
 [22 23]]
[9, 10, 11]
```

```
print(x_train)
print(list(y_train))
```

```
[[ 0  1]
 [ 2  3]
 [ 4  5]
 [ 6  7]
 [ 8  9]
 [10 11]
 [12 13]
 [14 15]
 [16 17]]
[0, 1, 2, 3, 4, 5, 6, 7, 8]
```

5. 한국어 전처리 패키지(Text Preprocessing Tools for Korean Text)

◆ PyKoSpacing

- 딥러닝 모델을 활용한 한국어 띄어쓰기 패키지
- Spacing() 활용

```
!pip install git+https://github.com/haven-jeon/PyKoSpacing.git
```

```
sent = '김철수는 극중 두 인격의 사나이 이광수 역을 맡았다. 철수는 한국 유일의 태권도 전승자를 가리는 결전의 날을 앞두고 10년간 함께 훈련한 사형인 유연재(김광수 분)를 찾으러 속세로 내려온 인물이다.'
```

```
new_sent = sent.replace(" ", '') # 띄어쓰기가 없는 문장 임의로 만들기  
print(new_sent)
```

김철수는극중두인격의사나यी광수역을맡았다.철수는한국유일의태권도전승자를가리는결전의날을앞두고10년간함께훈련한사형인유연재(김광수분)를찾으러속세로내려온인물이다.

```
from pykospadding import Spacing  
spacing = Spacing()  
kospacing_sent = spacing(new_sent)  
  
print(sent)  
print(kospacing_sent)
```

김철수는 극중 두 인격의 사나이 이광수 역을 맡았다. 철수는 한국 유일의 태권도 전승자를 가리는 결전의 날을 앞두고 10년간 함께 훈련한 사형인 유연재(김광수 분)를 찾으러 속세로 내려온 인물이다.
김철수는 극중 두 인격의 사나이 이광수 역을 맡았다. 철수는 한국 유일의 태권도 전승자를 가리는 결전의 날을 앞두고 10년간 함께 훈련한 사형인 유연재(김광수 분)를 찾으러 속세로 내려온 인물이다.

5. 한국어 전처리 패키지(Text Preprocessing Tools for Korean Text)

◆ Py-Hanspell

- 네이버 한글 맞춤법 검사기를 활용한 한국어 맞춤법 + 띄어쓰기 교정 패키지
- spell_checker.check() 활용

```
!pip install git+https://github.com/ssut/py-hanspell.git
```

```
from hanspell import spell_checker

sent = "맞춤법 틀리면 외 안돼? 쓰고싶은대로쓰면돼지 "
spelled_sent = spell_checker.check(sent)

hanspell_sent = spelled_sent.checked
print(hanspell_sent)
```

맞춤법 틀리면 왜 안돼? 쓰고 싶은 대로 쓰면 되지

```
spelled_sent = spell_checker.check(new_sent)

hanspell_sent = spelled_sent.checked
print(hanspell_sent)
print(kospacing_sent) # 앞서 사용한 kospacing 패키지에서 얻은 결과
```

김철수는 극 중 두 인격의 사나이 이광수 역을 맡았다. 철수는 한국 유일의 태권도 전승자를 가리는 결전의 날을 앞두고 10년간 함께 훈련한 사형인 유연재(김광수 분)를 찾으러 속세로 내려온 인물이다.
김철수는 극중 두 인격의 사나이 이광수 역을 맡았다. 철수는 한국 유일의 태권도 전승자를 가리는 결전의 날을 앞두고 10년간 함께 훈련한 사형인 유연재(김광수 분)를 찾으러 속세로 내려온 인물이다.

5. 한국어 전처리 패키지(Text Preprocessing Tools for Korean Text)

◆ SOYNLP

- 비지도 학습으로 데이터에 자주 등장하는 단어를 단어로 추출하는 단어 토큰나이저
 - 기존 형태소 분석기가 미등록 단어(예: 신조어)를 제대로 분리하지 못하는 점 보완
 - 단어 점수 표(응집 확률, 브랜칭 엔트로피 활용)를 바탕으로 동작

```
!pip install soynlp
```

<훈련 데이터 활용하여 학습한 후, 응집 확률 및 브랜칭 엔트로피 계산>

```
import urllib.request
from soynlp import DoublespaceLineCorpus
from soynlp.word import WordExtractor
```

```
urllib.request.urlretrieve("https://raw.githubusercontent.com/lovit/soynlp/master/tutorials/2016-10-20.txt", filename="2016-10-20.txt")
# 학습에 필요한 한국어 문서 다운로드
```

```
('2016-10-20.txt', <http.client.HTTPMessage at 0x7f6414117890>)
```

```
corpus = DoublespaceLineCorpus("2016-10-20.txt") # 훈련 데이터를 다수의 문서로 분리
len(corpus)
```

30091

5. 한국어 전처리 패키지(Text Preprocessing Tools for Korean Text)

- wordExtractor.extract() 활용: 전체 코퍼스에 대한 응집 확률, 브랜칭 엔트로피 단어 점수표 산출

```
word_extractor = WordExtractor()  
word_extractor.train(corpus)  
word_score_table = word_extractor.extract()
```

```
training was done. used memory 1.467 Gb  
all cohesion probabilities was computed. # words = 223348  
all branching entropies was computed # words = 361598  
all accessor variety was computed # words = 361598
```

- 응집 확률(cohesion probability): 문자열을 문자 단위로 분리하여 하위 문자열을 만드는 과정에서, 왼쪽부터 순서대로 문자를 추가하면서 각 문자열이 주어졌을 때 그 다음 문자가 나올 확률을 계산하여 누적곱을 한 값

$$\text{cohesion}(n) = \left(\prod_{i=1}^{n-1} P(c_{1:i+1} | c_{1:i}) \right)^{\frac{1}{n-1}}$$

→ 이 값이 클수록 전체 코퍼스에서 단어로 등장할 가능성이 높음

5. 한국어 전처리 패키지(Text Preprocessing Tools for Korean Text)

$$\text{cohesion}(n) = \left(\prod_{i=1}^{n-1} P(c_{1:i+1} | c_{1:i}) \right)^{\frac{1}{n-1}}$$

- 예시) ‘반포한강공원에’(길이: 7)에 대해 각 하위 문자열의 응집 확률 구하기

- $\text{cohesion}(2) = P(\text{반포}|\text{반})$
- $\text{cohesion}(3) = \sqrt[2]{(P(\text{반포}|\text{반}) \cdot P(\text{반포한}|\text{반포}))}$
- $\text{cohesion}(4) = \sqrt[3]{(P(\text{반포}|\text{반}) \cdot P(\text{반포한}|\text{반포}) \cdot P(\text{반포한강}|\text{반포한}))}$
- $\text{cohesion}(5) = \sqrt[4]{(P(\text{반포}|\text{반}) \cdot P(\text{반포한}|\text{반포}) \cdot P(\text{반포한강}|\text{반포한}) \cdot P(\text{반포한강공}|\text{반포한강}))}$
- $\text{cohesion}(6) = \sqrt[5]{(P(\text{반포}|\text{반}) \cdot P(\text{반포한}|\text{반포}) \cdot P(\text{반포한강}|\text{반포한}) \cdot P(\text{반포한강공}|\text{반포한강}) \cdot P(\text{반포한강공원}|\text{반포한강공}))}$
- $\text{cohesion}(7) = \sqrt[6]{(P(\text{반포}|\text{반}) \cdot P(\text{반포한}|\text{반포}) \dots \text{중략} \dots \cdot P(\text{반포한강공원}|\text{반포한강공}) \cdot P(\text{반포한강공원에}|\text{반포한강공원}))}$

```
word_score_table["반포한"].cohesion_forward
```

```
0.08838002913645132
```

```
word_score_table["반포한강"].cohesion_forward
```

```
0.19841268168224552
```

```
word_score_table["반포한강공"].cohesion_forward
```

```
0.2972877884078849
```

```
word_score_table["반포한강공원"].cohesion_forward
```

```
0.37891487632839754
```

```
word_score_table["반포한강공원에"].cohesion_forward
```

```
0.33492963377557666
```


5. 한국어 전처리 패키지(Text Preprocessing Tools for Korean Text)

- 브랜칭 엔트로피(branching entropy): 주어진 문자열에서 다음 문자가 얼마나 등장할 수 있는지를 판단하는 척도
("다음 문자 예측을 위해 헛갈리는 정도")

```
word_score_table["디스"].right_branching_entropy
```

```
1.6371694761537934
```

```
word_score_table["디스플"].right_branching_entropy
```

```
-0.0
```

```
word_score_table["디스플레이"].right_branching_entropy
```

```
-0.0
```

```
word_score_table["디스플레이"].right_branching_entropy
```

```
3.1400392861792916
```

‘디스플레이’ 뒤에 조사나 다른 단어가 다양하게 올 수 있음!

→ 전체 코퍼스에서 하나의 단어가 끝나면 그 경계 부분부터 이 값이 증가

5. 한국어 전처리 패키지(Text Preprocessing Tools for Korean Text)

1) L 토크나이저: 점수가 가장 높은 L 토크를 찾아 분리

```
from soynlp.tokenizer import LTokenizer

scores = {word:score.cohesion_forward for word, score in word_score_table.items()}
l_tokenizer = LTokenizer(scores=scores)
l_tokenizer.tokenize("국제사회와 우리의 노력들로 범죄를 척결하자", flatten=False)

[('국제사회', '와'), ('우리', '의'), ('노력', '들로'), ('범죄', '를'), ('척결', '하자')]
```

2) 최대 점수 토크나이저: 점수가 가장 높은 글자 시퀀스를 순차적으로 탐색

- 띄어쓰기가 되어 있지 않는 문장 처리시 유용

```
from soynlp.tokenizer import MaxScoreTokenizer

maxscore_tokenizer = MaxScoreTokenizer(scores=scores)
maxscore_tokenizer.tokenize("국제사회와우리의노력들로범죄를척결하자")

['국제사회', '와', '우리', '의', '노력', '들로', '범죄', '를', '척결', '하자']
```

5. 한국어 전처리 패키지(Text Preprocessing Tools for Korean Text)

3) soynlp.normalizer: 반복되는 문자 정제

```
from soynlp.normalizer import *
```

```
print(emoticon_normalize('악ㅋㅋㅋㅋ이영화존잼쓰ㅠㅠㅠㅠ', num_repeats=2))  
print(emoticon_normalize('악ㅋㅋㅋㅋㅋㅋㅋㅋ이영화존잼쓰ㅠㅠㅠㅠ', num_repeats=2))  
print(emoticon_normalize('악ㅋㅋㅋㅋㅋㅋㅋㅋㅋㅋ이영화존잼쓰ㅠㅠㅠㅠ', num_repeats=2))  
print(emoticon_normalize('악ㅋㅋㅋㅋㅋㅋㅋㅋㅋㅋㅋㅋ이영화존잼쓰ㅠㅠㅠㅠ', num_repeats=2))
```

```
아ㅋㅋ영화존잼쓰ㅠㅠ  
아ㅋㅋ영화존잼쓰ㅠㅠ  
아ㅋㅋ영화존잼쓰ㅠㅠ  
아ㅋㅋ영화존잼쓰ㅠㅠ
```

```
print(repeat_normalize('와하하하하하하하하하하', num_repeats=2))  
print(repeat_normalize('와하하하하하하하', num_repeats=2))  
print(repeat_normalize('와하하하하하', num_repeats=2))
```

```
와하하하  
와하하하  
와하하하
```

5. 한국어 전처리 패키지(Text Preprocessing Tools for Korean Text)

◆ Customized KoNLPy

- 사용자 사전 추가가 매우 쉬운 패키지
- add_dictionary() 활용

```
pip install customized_konlpy
```

```
from konlpy.tag import Twitter
twitter = Twitter()
twitter.morphs('은경이는 사무실로 갔습니다.')
```

```
/usr/local/lib/python3.7/dist-packages/konlpy/tag/_okt.py:16: UserWarning: "Twitter" has changed to "Okt" since KoNLPy v0.4.5.
  warn('"Twitter" has changed to "Okt" since KoNLPy v0.4.5.')
['은', '경이', '는', '사무실', '로', '갔습니다', '.']
```

```
twitter.add_dictionary('은경이', 'Noun')
```

```
twitter.morphs('은경이는 사무실로 갔습니다.')
```

```
['은경이', '는', '사무실', '로', '갔습니다', '.']
```