# TF-IDF calculation using Scikit-Learn

Ko, Youngjoong

Sungkyunkwan University nlp.skku.edu, nlplab.skku.edu



## **Contents**



- NLTK Part-Of-Speech(POS) Tag
- Count based Word Representation
  - Bag of Words(BoW)
  - TF-IDF
- ❖ Scikit-Learn
- Assignment

# **POS Tag**



#### POS Tag

- A POS tag is a label assigned to each word in a text to indicate the part of speech.
  - Ex) subject, verb, object, etc.
- In general, main components of sentences are subject, verb, object, and complement, and these are usually verbs or nouns.
- In this assignment, we use only verbs and nouns tags.
  - Verb: (VB, VBD, VBG, VBN, VBP, VBZ)
  - Noun: (NN, NNS, NNP, NNPS) (Details in page 4)

# POS Tag



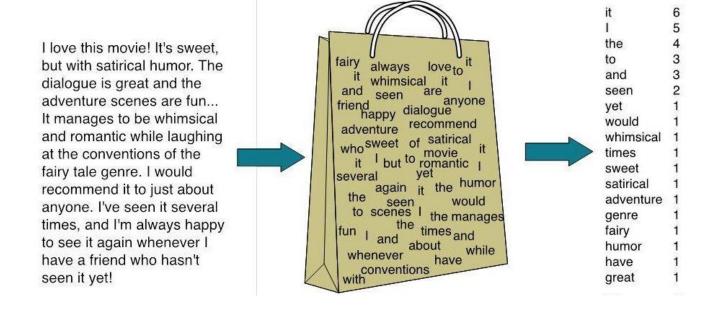
## ❖ NLTK POS Tag

Tag	Description	Example
CC	coordinating conjunction	and
CD	cardinal number	1, third
DT	determiner	the
EX	existential there	there is
FW	foreign word	d'hoevre
IN	preposition/subordinating conjunction	in, of, like
JJ	adjective	big
JJR	adjective, comparative	bigger
JJS	adjective, superlative	biggest
LS	list marker	1)
MD	Modal	could, will
NN	noun, singular or mass	Door
NNS	noun plural	Doors
NNP	proper noun, singular	John
NNPS	proper noun, plural	Vikings
PDT	Predeterminer	both the boys
POS	possessive ending	friend's

PRP	personal pronoun	I, he, it
PRP\$	possessive pronoun	my, his
RB	Adverb	however, usually
RBR	adverb, comparative	Better
RBS	adverb, superlative	Best
RP	Particle	give <i>up</i>
ТО	То	to go, to him
UH	Interjection	Uhhuhhuhh
VB	verb, base form	Take
VBD	verb, past tense	Took
VBG	verb, gerund/present participle	Taking
VBN	verb, past participle	Taken
VBP	verb, sing. present, non-3d	Take
VBZ	verb, 3rd person sing. Present	Takes
WDT	wh-determiner	Which
WP	wh-pronoun	who, what
WP\$	possessive wh-pronoun	Whose
WRB	wh-abverb	where, when



- Bag of Words(BoW)
  - A method of numerical expression of text data that focuses only on the frequency of words without considering the order of words.





- Bag of Words(BoW) based Document Features Extraction Method
  - Doc0: I/PRP am/VBP a/DT boy/NN
  - Doc1: I/PRP am/VBP a/DT girl/NN
    - 1) Remove duplicates from all words in Doc0 and Doc1, and extract verbs and nouns. List each word in column form. (Sorted by alphabet)

```
- [ 'am/VBP', 'boy/NN', 'girl/NN']
```

2) Give each of the following words a unique number (index).

```
- { 'am/VBP': 0, 'boy/NN': 1, 'girl/NN': 2 }
```

3) In each document, write the frequency in which the word appears in each number.

```
- Index: 0 1 2
```



- ❖ TF-IDF(Term Frequency-Inverse Document Frequency)
  - TF(Term Frequency)
    - : The number of times that term t occurs in document d ( $tf_{t,d}$ )
  - DF(Document Frequency)
    - : The number of documents that contain the term t ( $df_t$ )
  - IDF(Inverse Document Frequency)
    - : Inverse value of DF.

Term frequency		Document frequency		Normalization Normalization	
n (natural)	$tf_{t,d}$	n (no)	1	n (none)	1
I (logarithm)	$1 + \log(tf_{t,d})$	t (idf)	$\log \frac{N}{\mathrm{df_t}}$	c (cosine)	$\frac{1}{\sqrt{w_1^2 + w_2^2 + + w_M^2}}$
a (augmented)	maxt(CIt,d)	p (prob idf)	$max\{0, log \tfrac{N-\mathrm{df}_t}{\mathrm{df}_t}\}$	u (pivoted unique)	1/u
b (boolean)	$\begin{cases} 1 & \text{if } \operatorname{tf}_{t,d} > 0 \\ 0 & \text{otherwise} \end{cases}$			b (byte size)	$1/\mathit{CharLength}^{lpha}, \ lpha < 1$
L (log ave)	$\frac{1 + \log(tf_{t,d})}{1 + \log(ave_{t \in d}(tf_{t,d}))}$				



#### ❖ TF (Term Frequency)

Doc0: I/PRP am/VBP a/DT boy/NN

Doc1: I/PRP am/VBP a/DT girl/NN

Doc2: who/WP is/VBZ a/DT boy/NN

-  $tf_{t,d}$ : The number of times that term t occurs in document d ( $tf_{t,d}$ )

	am/VBP	is/VBZ	boy/NN	girl/NN
Doc0	1	0	1	0
Doc1	1	0	0	1
Doc2	0	1	1	0



- ❖ IDF (Inverse Document Frequency)
  - Doc0: I/PRP am/VBP a/DT boy/NN
  - Doc1: I/PRP am/VBP a/DT girl/NN
  - Doc2: who/WP is/VBZ a/DT boy/NN
    - $log \frac{N}{df_t}$ : Inverse value of DF.
    - N: Total number of document
    - $df_t$ : The number of documents that contain the term t  $(df_t)$

Ex) am/VBP = 
$$log \frac{3}{2} = 0.176$$
, girl/NN=  $log \frac{3}{1} = 0.477$ 

	am/VBP	is/VBZ	boy/NN	girl/NN
IDF	0.176	0.477	0.176	0.477



TF-IDF(Term Frequency Inverse Document Frequency)

$$-tf_{t,d} \times log \frac{N}{df_t}$$

Ex) am/VBP : 
$$(1 \times 0.176) = 0.176$$

	am/VBP	is/VBZ	boy/NN	girl/NN
Doc0	0.176	0	0.176	0.477
Doc1	0.176	0	0.176	0
Doc2	0	0.477	0	0



#### Install scikit-learn

- · Google, 'Colab'
  - Packages installed by default

#### Anaconda

- If your virtual environment is activated,
  - \$ conda install scikit-learn
- If your virtual environment is deactivated
  - \$ conda install --name [virtual\_environment\_name] scikit-learn
- Python environment
  - \$ pip install -U scikit-learn



#### ❖ Scikit-Learn Package

- One of the most widely-used Python package for data science and machine learning.
  - Benchmark dataset
  - Preprocessing
  - Supervised learning
  - Unsupervised learning
  - Evaluation and selection



## Preprocessing with Scikit-Learn

- Subpackage sklearn.feature\_extraction and sklearn.feature\_extraction.text provide the following subsystems for document preprocessing:
- The following classes create vectors after automatically sorting each token.
  - DictVectorizer
    - Creates a BOW vector using a dictionary consisting of counts of each word.
  - CountVectorizer
    - Generates word tokens and then creates a BOW vector by counting the number of each token.
  - TFidfVectorizer
    - It is similar to CountVectorizer but creates a BOW vector using weights of words calculated by TF-IDF.



#### DictVectorizer

- DictVectorizer is provided by the feature\_extraction subpackage.
- Creates a BOW vector using a dictionary consisting of counts of each word.

```
1 from sklearn.feature extraction import DictVectorizer
 3 Bag Of Word = DictVectorizer(sparse=False)
 4 Frequency = [{'am':1, 'boy':1}, {'am':1, 'girl':1}]
 5 output = Bag_Of_Word.fit_transform(Frequency)
 7 print (output)
[[1, 1, 0,]
[1. 0. 1.]]
 1 print(Bag_Of_Word.feature_names_)
['am', 'boy', 'girl']
 1 print(Bag_Of_Word.transform({'girl':2, 'female':3}))
[[0. 0. 2.]]
```



#### CountVectorizer

- CountVectorizer performs following tasks:
  - 1. Converting the document to a list of tokens.
  - 2. Counting the frequency of each token in each document.
  - 3. Converting each document to a BOW vector.



#### CountVectorizer

[0 1 0 1]]

```
1 from sklearn.feature_extraction.text import CountYectorizer
 2 import nltk
 3 from nltk.tokenize import word_tokenize
 5 corpus = ['l am a boy',
             'Lamagirl'
             'Who is a boy'l
 9 POS_corpus = list()
10
11 for sentence in corpus:
       pos_token = nltk.pos_tag(word_tokenize(sentence))
       POS_corpus.append(' '.join([t[0] for t in pos_token if t[1] in ['NN', 'NNS', 'NNP', 'NNPS', 'YB', 'YBD', 'YBB', 'YBN', 'YBP', 'YBZ']]))
13
14
15 print (POS_corpus)
['am boy', 'am girl', 'is boy']
 1 vect = CountYectorizer()
 2 vect.fit_transform(POS_corpus)
 3 print(vect.vocabulary_)
{'am': 0, 'boy': 1, 'girl': 2, 'is': 3}
1 print(vect.transform(POS corpus).toarray().tolist())
[[1 1 0 0]
 [1 0 1 0]
```



#### TfidfVectorizer

- When TfidfVectorizer creates a vector from TF-IDF, the vector is normalized using L2 Norm.
  - L2 Norm:  $||x||_2 := \sqrt{x_1^2 + x_2^2 + \dots + x_n^2}$  where  $x = (x_1, x_2, \dots, x_n)$
- TfidfVectorizer uses  $idf = log(\frac{N+1}{df_t+1}) + 1$  to calculate TF-IDF.

$$idf = log\left(\frac{N+1}{df_t+1}\right) + 1$$

- The constant "1" is added to the numerator and denominator of  $log\left(\frac{N}{df_t}\right)$  to prevent dividing by zero.
- The constant "1" is added to  $log\left(\frac{N+1}{df_t+1}\right)$  to prevent terms which occur in all documents in a training set from being entirely ignored.



#### TfidfVectorizer

```
1 from sklearn.feature extraction.text import TfidfVectorizer
 2 import nltk
 3 from nltk.tokenize import word_tokenize
 5 corpus = ['I am a boy',
            'l am a girl',
             'Who is a boy']
 9 POS_corpus = list()
10
11 for sentence in corpus:
      pos token = nltk.pos tag(word tokenize(sentence))
      POS_corpus.append(' '.join([t[0] for t in pos_token if t[1] in ['NN', 'NNS', 'NNP', 'NNPS', 'YB', 'YBD', 'YBB', 'YBB', 'YBP', 'YBZ']]))
13
14
15 print (POS_corpus)
['am boy', 'am girl', 'is boy']
 1 tfidfvect = TfidfVectorizer()
 2 tfidfvect.fit transform(POS corpus)
 3 print(tfidfvect.vocabulary_)
{'am': 0, 'boy': 1, 'girl': 2, 'is': 3}
1 print(tfidfvect.transform(POS_corpus).toarray().tolist())
[[0.70710678 0.70710678 0.
 [0.60534851 0. 0.79596054 0.
 [0.
            0.60534851 0. 0.79596054]
```

# **Assignment**



## Assignment

- 1) The attached JSON file in assignment 2 consists of 300 articles with multiple sentences.
  - See page 21 for JSON input file format.
- 2) Calculate TF-IDF vector of 300 articles using only nouns and verbs from given data.
  - Verbs (VB, VBD, VBG, VBN, VBP, VBZ) and nouns (NN, NNS, NNP, NNPS)
  - See page 4 for POS tags
- 3) You MUST create an output text file.
  - See page 22 for an example of output file.
- 4) You can use only Scikit-Learn for calculating TF-IDF and PyKomoran for morpheme analysis respectively.

# **Assignment**



#### Submit File

- 1) Python code file (.py) (python version 3.x)
  - Format: "Student Number\_Name\_TFIDF.py".
    - Ex) "2020000000\_홍길동.py"
- 2) TEXT file (.txt)
  - Format: "Student Number\_Name\_TFIDF.txt".

# **JSON Input File**



```
"business":
                "Ad sales boost Time Warner profit Quarterly profits at US media ...
                "Dollar gains on Greenspan speech The dollar has hit ...
                                                                                        Article
Article Topic
                (중략)
               "GM, Ford cut output as sales fall ...
               "Ebbers denies WorldCom fraud ...
                (중략)
               "politics": [
                "Labour plans maternity pay rise...
                "Watchdog probes e-mail deletions...
               (중략)
                "Blair stresses prosperity goals...
                "Guantanamo man 'suing government'...
```

JSON input file includes 300 articles on business, politics, and tech topics.

# An Example of Output File



```
Article
             Student Number_Name_TFIDF.txt
                                           Tab
              (business1)
                       0.4311
                                 0.0
                                           0.0
                                                    ...(하략)..
             0.0 0.0
                                      0.0
                                                0.0
              (business2)
              0.0 0.0 0.0 0.0 0.1334 0.0 0.0 ...(하략)...
Article Topic
              (중략)
              (중략)
              (politics1)
              0.0 0.0
                      0.0 0.0 0.0 0.0 0.3311 ...(하략)...
              (중략)
              (politics 100)
                            0.3311
                                                     0.0 ...(하략)...
              0.1176
                       0.0
                                      0.0
                                           0.0
                                                0.0
```

- This is an example.
- This can be different with actual answers.