

KoNLPy (한글어 처리 package) ("코엔엘파이")

- i) kKma: 코리안 (서울대)
- ii) Hanmanum: 한남대 (kaist)
- iii) Twitter
- iv) Mecab: 일본어용 형태소 분석기를 한글에 적용
- v) Komoran: 코로안

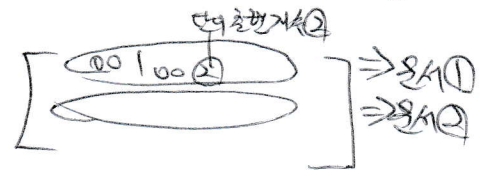
Sakit-Learn의 문서 처리 기능

- Bow (Bag of Words) : 문서를 숫자 Vector로 변환하는 기본적인 방법.

x_{ij} = 문서 d_i 의 단어 w_j 의 출현 빈도

↓

$x_{ij} = \begin{cases} 0, & \text{단어 } w_j \text{가 문서 } d_i \text{에 없으면} \\ 1, & \text{단어 } w_j \text{가 문서 } d_i \text{에 있으면} \end{cases}$



- Skat-Learn의 문서 처리 기능

↗ "단어 사용 빈도"

i) DictVectorizer : 각 단어의 수를 세어 놓은 사전에서 Bow 벡터를 만든다.

ii) CountVectorizer : ^{토큰화 (tokenizer) 후 CountVectorizer} 문서 집합에서 단어 토큰을 생성하고, 각 단어의 수를 세어 Bow 인코딩한 Vector를 만든다.

① 토큰화 (tokenizer) 후 CountVectorizer
② 각 문서를 Bow 인코딩한 Vector로 변환
③ 각 문서에서 토큰의 출현 빈도를 세다

iii) TfidfVectorizer : TF-IDF 방식으로 단어의 가중치를 조절된 Bow (Term Frequency-Inverse document frequency) 벡터를 만든다. → 단어의 수로 계산된 TF, 모든 문서에 공통된 단어가 있으면 가중치 감소

iv) HashVectorizer : (해시 함수) hash function을 사용해 단어 메모리와 빠른 속도로 Bow 벡터를 만든다.

★ ii) CounterVectorizer 변수

① Stop-words → 무시할 단어 목록

② Analyzer, tokenizer, token-pattern → 사용한 토큰 생성기를 선택

③ N-gram → 토큰의 크기를 결정. (1-그램: 토큰 하나만 단어 사용, 2-그램: 연속된 토큰 두 개를 하나의 단어로 사용)

④ max_df, min_df → 빈도수 (하나 자주 나오거나, 너무 안나오면 무시)