

Vector Representation

언어를 컴퓨터가 알아들도록 숫자로 표현하기

Vector space representation

- One hot representation (Bag of Words model)
 - (row, column) 은 (문서, 단어) 해당하는 값은 단어의 중요도 혹은 빈도수를 의미

	기계	학습	은	텍스트	마이닝	는
Doc 1	3	2	5	0	0	0
Doc 2	0	0	0	3	5	5
...

Doc 1 = [(0, 3), (1, 2), (2, 5)]
 Doc 2 = [(3, 3), (4, 5), (5, 5)]



	0	1	2	3	4	5
Doc 1	3	2	5	0	0	0
Doc 2	0	0	0	3	5	5
...

Vector space representation

- One hot representation (Bag of Words model)

	0	1	2	3	4	5
Doc 1	3	2	5	0	0	0
Doc 2	0	0	0	3	5	5
...

- Column 개수 $|V|$ 는 문서 전체에서 등장한 단어 종류로, 매우 큼니다.
- 한 문서에 등장하는 단어의 개수는 적기 때문에, 대부분의 값이 0입니다 (Sparse vector).
- 문서에 등장한 단어를 쉽게 확인할 수 있어 해석이 쉽지만, 모든 단어는 다른 단어로 취급합니다. 단어간 유사성을 표현하기 어렵습니다.

Vector space representation

- Distributed representation
 - 단어/문서를 Word2Vec 등으로 **d차원 공간의 벡터**로 표현
 - 단어의 "**의미적 유사성**"을 벡터 공간에 표현
 - 벡터 공간의 거리가 가까운 단어 또는 문서는 의미가 비슷

'dog' = [0.31, -0.21, 2.01, 0.58, ...]

'cat' = [0.45, -0.17, 1.79, 0.61, ...]

'topic modeling' = [-2.01, 0.03, 0.22, 0.54, ...]

'dim. reduction' = [-1.88, 0.11, 0.19, 0.45, ...]

Vector space representation

- Distributed representation
 - 각 벡터는 "의미 공간"에서 좌표 역할



Document clustering

- 군집화 (Clustering)는 비슷한 데이터를 하나의 집합으로 그룹화합니다.
- 리뷰가 비슷한 영화들의 군집화 결과

cluster # 1	cluster # 2	cluster # 3
해무	인터스텔라	응답하라 1988
베를린	미스터 고	인턴
내가 살인범이다	다크 나이트	님아, 그 강을 건너지 마오
신세계	영웅: 샐러맨더의 비밀	카트
곡성(哭聲)	인셉션	인사이드 아웃
검은 사제들	트랜스포머 3	형
악마를 보았다	배틀쉽	비긴 어게인
용의자	스카이라인	두근두근 내 인생
감기	2012	라라랜드
감시자들	그래비티	반창꼬

Document clustering

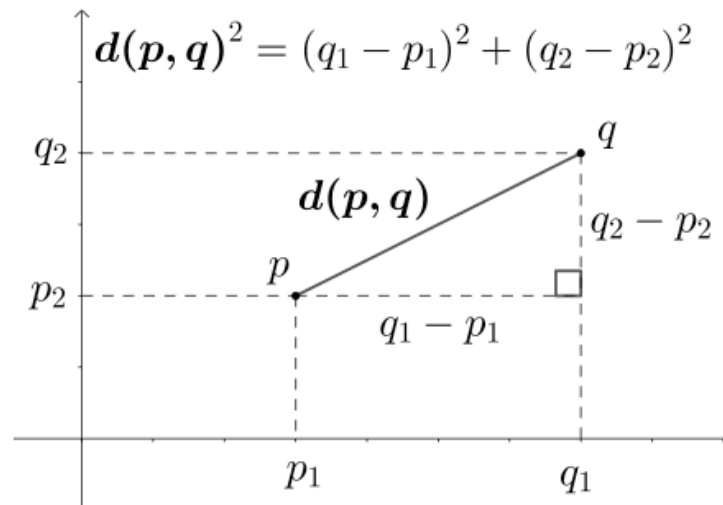
- 벡터 공간에서의 거리 척도로 Euclidean, Cosine 등이 이용됩니다.
- Euclidean distance 는 벡터의 크기(문서의 길이)에 영향을 받는다.
 - doc 1: 3 단어 / doc 2: 4 단어 / doc 3: 7 단어
 - 문서 1과 문서 3은 방향은 비슷하지만 단어 수가 달라서 거리가 멀다

	Term 1	Term 2	Term 3	Term 4	Term 5
Doc 1	1	1	1		
Doc 2			2	1	1
Doc 3	2	2	2		1

Euclidean Distance

- 두 점 사이에 줄을 긋고, 그 줄의 길이를 계산하는 것!

$$\text{distance}(A, B) = \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2 + \dots + (a_n - b_n)^2}$$



2D에서 euclidean distance

Cosine Similarity

- Cosine similarity

$$\text{similarity} = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$

- $\cos((3, 0, 2), (1, 2, 0)) = \frac{3 \times 1 + 0 \times 2 + 2 \times 0}{\sqrt{3^2 + 2^2} \times \sqrt{1^2 + 2^2}} = \frac{3}{\sqrt{13 \times 5}}$

- Vector A 와 Vector B
 - 같은 방향(0) = 1
 - 완전히 반대 방향(180) = -1
 - 서로 독립적(90) = 0

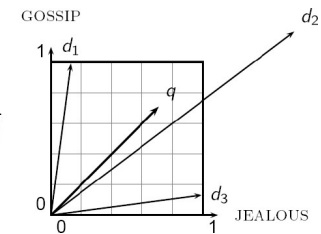
Document clustering

- Bag of words model 에는 Euclidean 보다 Cosine 이 적절합니다.
 - 두 문서에 공통으로 등장한 단어에 대해서만 유사성을 판단합니다.
 - Cosine은 문서 길이에(벡터의 크기,norm) 영향을 받지 않습니다
 - Sparse representation 에서는 벡터의 방향이 가장 중요합니다.

	Term 1	Term 2	Term 3	Term 4	Term 5
Doc 1	1	1	1		
Doc 2			2	1	1
Doc 3	2	2	2		1

Euclidean(d1, d2) = Euclidean(d1, d3)이지만, d1과 d3이 공통된 단어가 많기 때문에 더 비슷

$$Sim(D_1, D_2) = \frac{\sum_i x_{1i} x_{2i}}{\sqrt{\sum_j x_j^2} \sqrt{\sum_k x_k^2}}$$

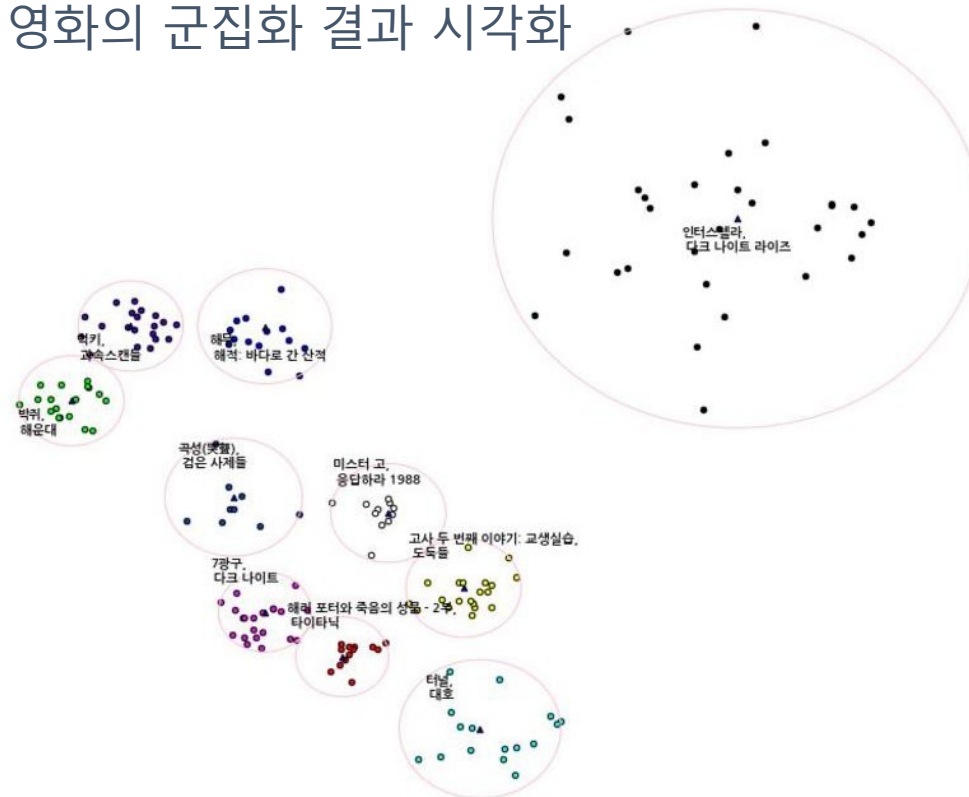


Document clustering

- Document distance/similarity 를 계산할 때에는 Cosine 이 적합합니다.
 - 문서 표현에 distributed representation 을 이용한다 하더라도 벡터의 방향이 가장 중요합니다.
 - Logistic regression, Neural network 등의 머신 러닝 알고리즘도 벡터 방향이 큰 영향을 미칩니다.

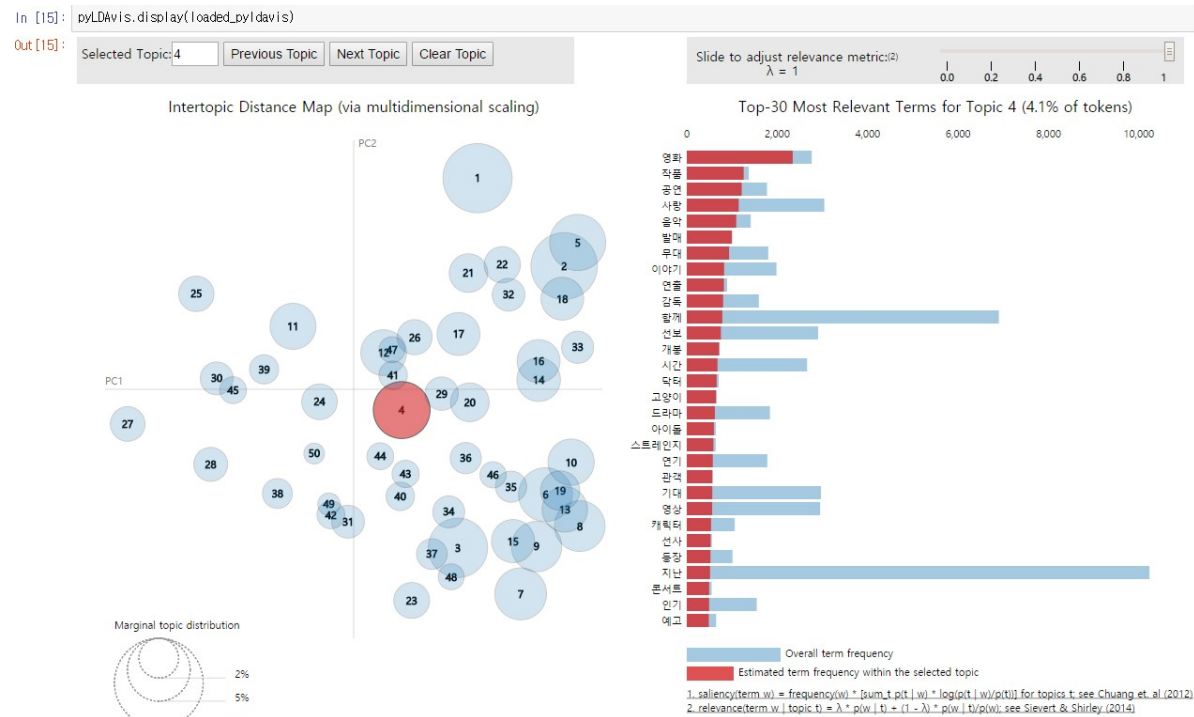
Word/Document Visualization

- 고차원의 벡터를 2차원으로 압축함으로써, 벡터 공간을 설명합니다.
- 리뷰 기반 영화의 군집화 결과 시각화



Word/Document Visualization

- 고차원의 벡터를 2차원으로 압축함으로써, 벡터 공간을 설명합니다.
- pyLDAVis 를 이용한 토픽 모델링의 시각화

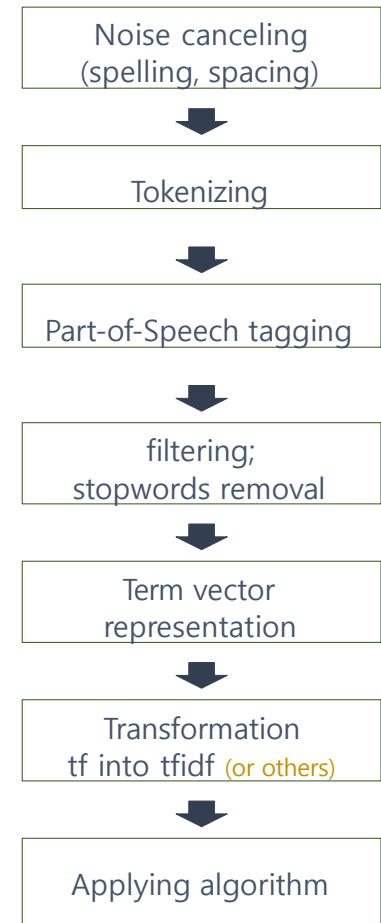


Text Processing

Feat. TF-IDF and some other techniques

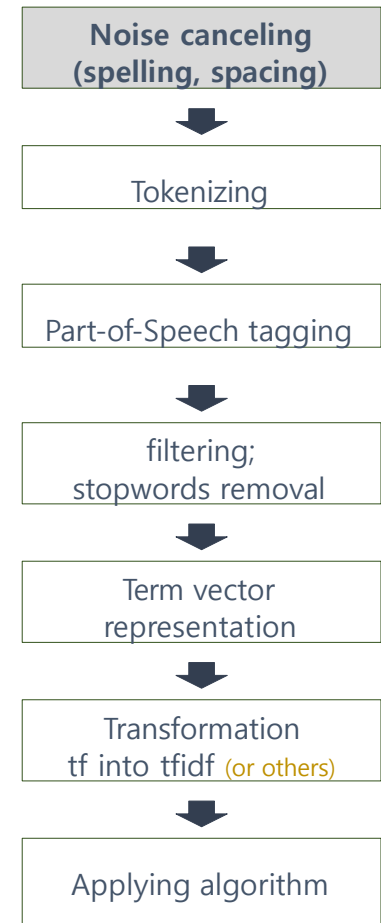
Framework

15



Spelling

- 한 단어가 다르게 적힌다면,
 - 같은 의미로 여러개 다른 단어로 표현
 - Bag of words model가 더 큰 sparse representation
 - 미등록단어 (Out of vocabulary) 문제 발생
- 사전에 존재하는 올바른 단어로 수정합니다
 - Edit distance(String distance) 비교를 통해 수정



Spelling

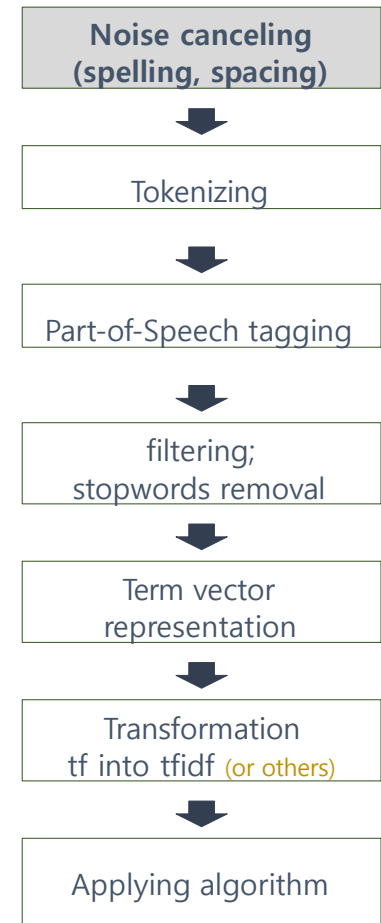
- 오타자는 수기로 입력된 데이터에서 주로 발생
- 오타자이외 은어, 비속어 등 올바른 표현으로 고쳐야하는 경우 종종 발생

데이터
제조외
제조, 도매, 부동산
건설업
조립금속제조, 기타화학제조
서비스 도소매
편의점, 담배
소매업.서비스업. 부동산업
식음료
제조 및 도소, 부동산업
제조업

사전		
가구내 고용활동	보험	운수
가스	부동산	원료재생
개인	사업시설관리	음식점
건설	사업지원	임대
과학	사회복지	임업
광업	서비스	자가소비생산활동
교육	소매	전기
국제	수도사업	전문
금융	수리	정보
기술	숙박	제조
기타	스포츠	증기
농업	어업	출판
단체	여가	폐기물처리
도매	영상	하수처리
방송통신	예술	협회
보건	외국기관	환경복원

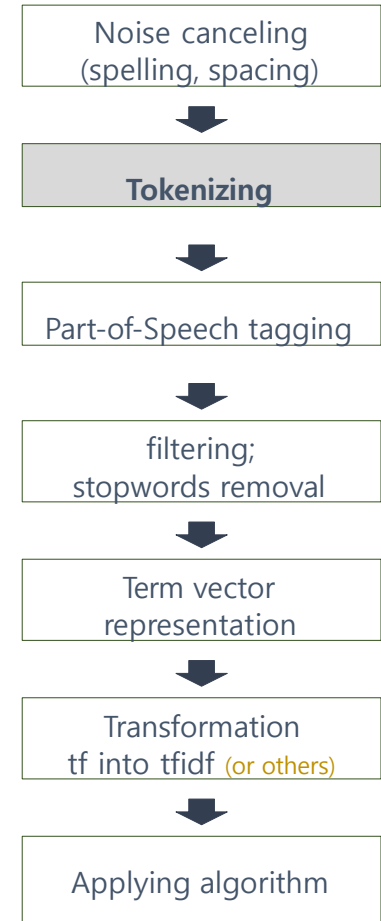
spacing

- 한국어의 어절은 띄어쓰기로 구분됩니다
 - 띄어쓰기 오류는 자연어처리의 정확도와 계산 시간에 영향을 줍니다
- 띄어쓰기 오류에 대응할 수 있는 토크나이저 혹은 띄어쓰기 오류 교정이 필요합니다



Tokenizing

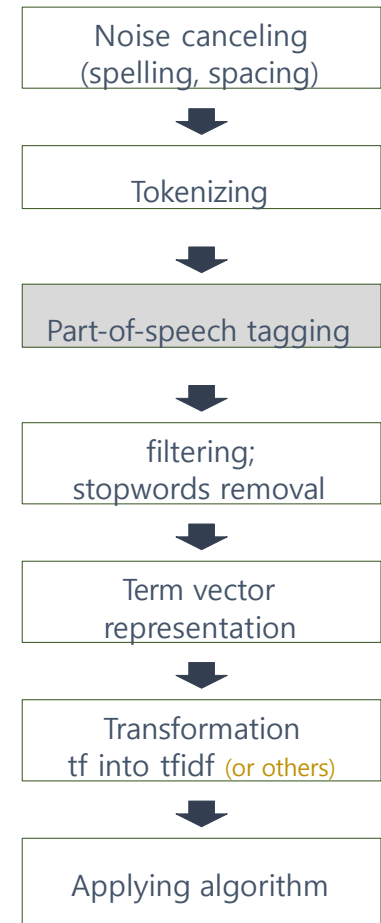
- 토크나이징은 어절에서 단어를 나누는 것입니다
 - [토크나이징, 은, 어절, 에서, 단어,를, 나누는, 것, 입니다]
- 정확히는 "문장"을 "토큰"으로 나누는 것입니다
 - 토큰은 n-gram, 어절, 단어, phrase 등으로 목적에 따라 다르게 정의
 - 토큰을 나누는 다양한 방법이 학습에 영향을 줌



Part-of-Speech tagging

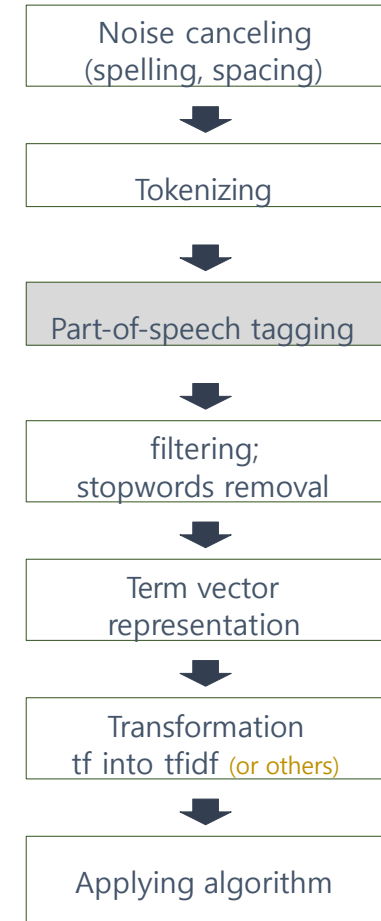
- 품사 판별은 주어진 단어의 품사를 구분합니다

- [토크나이징, 은, 어절, 에서, 단어, 를, 나누는, 것, 입니다] →
[(토크나이징, 명사),
(은, 조사),
(어절, 명사),
(에서, 조사),
(단어, 명사),
(를, 조사),
(나누는, 동사),
(것, 명사),
(입니다, 형용사)]



Morphological analysis

- 형태소 분석은 단어의 형태소를 인식
 - 형태소는 단어를 구성하는 최소단위
 - 품사 판별: "입니다" → 형용사
 - 형태소 분석: "입니다"
→ 이/형용사어근 + ㅂ니다/어미
- 형태소 분석을 바탕으로 단어의 품사를 추정

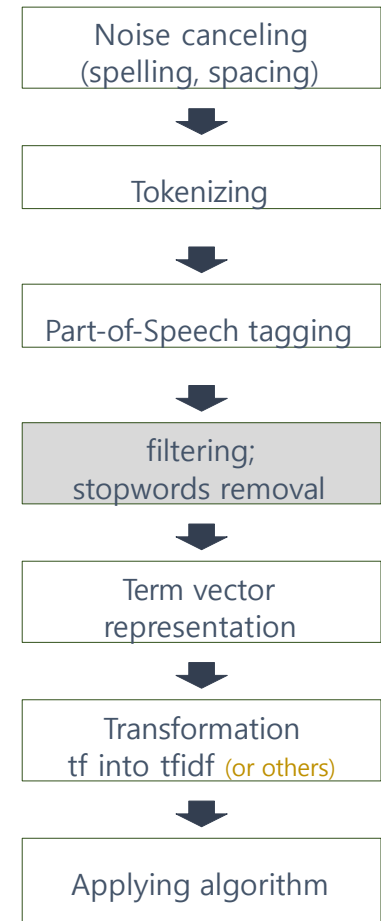


Stopwords removal

- Bag of words model 의 불필요한 단어를 제거합니다

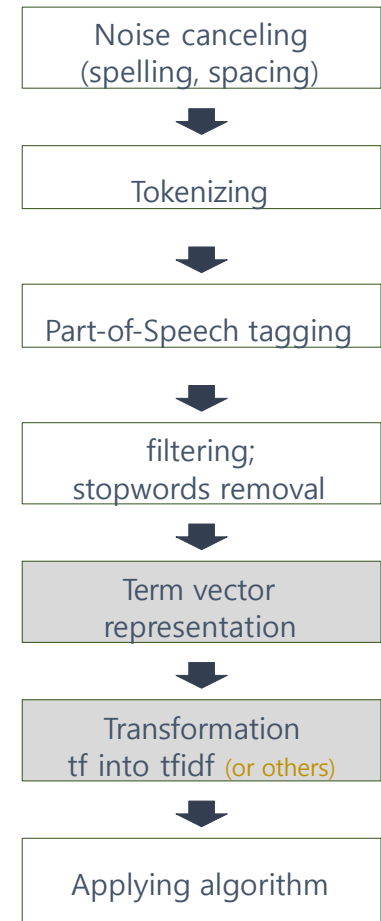
- 거의 등장하지 않는 단어
- -은, -는과 같은 조사(영어에서는 a, the, am, are, ...)
- 키워드 추출을 위한 명사 선택
- 특수문자

" \ , / , ' , \diamond , \ , " , / , ' , . , \ " , : , \Delta , \bullet , / , \blacksquare , (,) , \ " , >> , ` , / , - , ~ , = , \cdot , < , > , . , ? , ! , [,] , ... , \blacklozenge , \% "



Term vector representation

- Term weighting
 - $(i, j) = \text{weight}$
 - Term frequency vector 는 문서 i 에서의 단어 j 의 중요도를 단어의 빈도수로 표현
 - (i, j) 의 중요도는 정의하기 나뉘며, 반드시 TF 혹은 TF-IDF 를 이용해야 하는 것도 아닙니다



Term vector representation

- TF-IDF는 Information Retrieval 을 위하여 제안된 term weighting 입니다
 - 흔하게 등장하는 단어의 영향력을 줄입니다.

$$\text{TF-IDF}(w, d) = \text{TF}(w) \times \log\left(\frac{N}{\text{DF}(w)}\right)$$

TF: 단어 w 가 문서 d 에서 등장한 빈도 수

DF: 단어 w 가 등장한 문서의 개수

N : 문서집합에서 문서의 총 개수.

- $\text{DF}(w) = N$ 이면, 그 단어는 정보력이 없기 때문에 $\text{TF-IDF}(w, d) = 0$

Term vector representation

- Document frequency (df) 가 큰 단어는 정보력이 적습니다
 - 어디에나 등장하는 토큰은 정보력이 없음을 의미
 - 무의미하거나 문법적인 역할(조사)
 - 흔하게 등장하기 때문에 문서 간 거리를 계산할 때에도 무시

Word	Document frequency
1위	50
야구팀은	500
엘지 트윈스	1000
-은, -는	10,000

Term vector representation



