

Word Embedding

1. 자연어를 좌표평면 위에 표기 하기 위한 방법.

✓ one-hot encoding

단어를 하나의 숫자 형태로 표현

표현된 단어는 해당 위치의 Index에만 1을 둔 나머지는 0으로 표현

ex) I am a student $\Rightarrow [0 \ 1 \ 2 \ 3]$
 $\Rightarrow \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 1 \end{bmatrix}$

단어의 표현 (representation) 이 너무 sparse 해짐

\Rightarrow 불필요한 저장 공간 확보 \rightarrow Data는 크지만 정작 제공하는 정보 \downarrow

차원의 저주 (Curse of Dimension)

\Rightarrow 단어의 의미와 관련된 정보를 제공 하는 것이 불가능

\rightarrow 단어간 유사도 확인 역시 불가능

✓ Word2Vec

· 단어의 의미를 벡터공간에 임베딩

· 단어의 주변 단어를 활용해 의미를 파악

ex) 나는 소를 좋아해. 소와 돼지는 비슷한 의미라는 것을
나는 돼지고기를 좋아해. 유추 가능

· Word embedding의 성능 검증.

① WordSim 353 (사람이 미리 만든 Dataset)

단어 사이의 코사인 유사도를 사용하여 단어간 유사도 점수를

비교한다. (0.1 이상의 경우 잘 학습된 임베딩 모델)
② Semantic & Syntactic analogy
의미적 · 문법적 분석을 통한 확인.

· 학습되지 않은 단어에 대해 OOV 문제가 생김.

✓ FastText

- FastText의 학습 방법은 기존의 Word2Vec과 같다.
- OOV 문제 해결을 위해 단어를 n -gram 단위로 분할하여 학습

ex) <orange>

"오렌지"

→ (<o, or, ra, an... , <ora, ... , orange >

↑ 마지막은 원래의 단어.

orange의 단어 벡터는 n -gram 으로 분할된 토큰들의 합이 된다.

따라서 oranges 라는 새로운 단어가 들어와도 비슷한 의미로 해석이 가능!!

- 한국어의 경우. 자소단위를 사용하는 경우가 많다.