

## 워드 임베딩

단어를 벡터로 표현

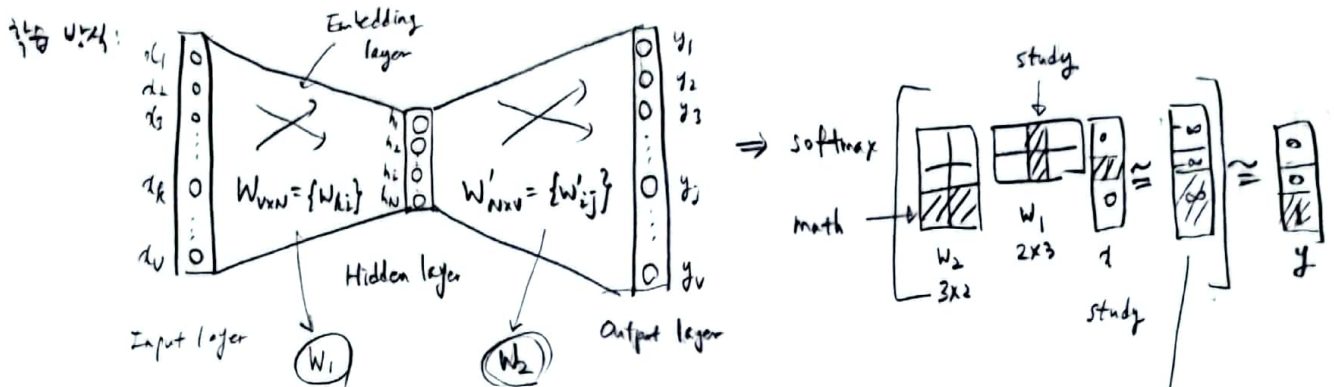
ex) "cat"과 "kitty"를 유사한 단어이므로 비슷한 벡터로 표현. 같은 거리  
 "hamburger"는 "cat", "kitty"와 유사하지 않으므로 다른 벡터로 표현 먼 거리  
 이런 식으로 공간에 좌표값을 부여

## Word2Vec

같은 문맥에서 나타난 인접한 단어들 간의 차이가 비슷한 것. 라는 가정

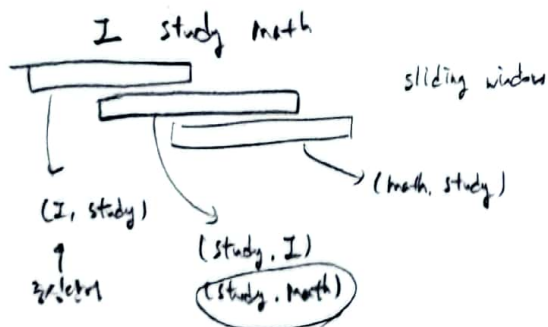
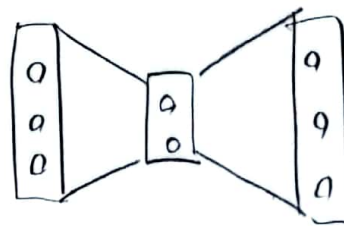
ex) (The cat purrs)  
 (This cat hunts mice)

Idea: 한 단어가 주변에 등장하는 단어를 통해 의미를 알 수 있다  
 주어진 학습 데이터로 비언어적 특징 단어 주변에 나타나는 단어들의 확률 분포 예측



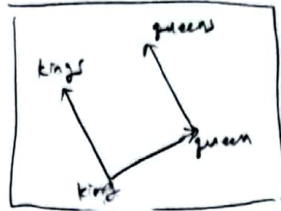
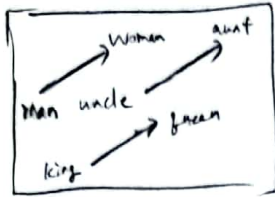
$$y = \text{softmax}(W_2 W_1 x)$$

ex) Sentence: "I study math"  
 Vocabulary: {"I", "study", "math"}  
 Input: "study" [0, 1, 0] one-hot  
 Output: "math" [0, 0, 1]



Input의  $W_1$  단어의 벡터  
 Output의  $W_2$  단어의 벡터의 내적  
 최대한 커리도록. 내적은 작-비대응  
 (비대응이 좋다)

항상 Word2Vec은 단어들 간의 시미플릭 관계를 통해 임베딩의 정렬에 관한 상동 같은 관계는 같은 벡터로 표현



$$\text{ex) } \text{vec}[\text{queen}] - \text{vec}[\text{king}] = \text{vec}[\text{woman}] - \text{vec}[\text{man}]$$

$$\text{ex) } \text{vec}[\text{사촌}] - \text{vec}[\text{남자}] + \text{vec}[\text{여자}] = \text{vec}[\text{처사촌DC}]$$

Word intrusion detection.

→ 단어들이 주어졌을 때 나머지 단어가 의의가 가장 상이한 단어는 무엇인가?

ex) staple hammer saw drill

유음으로 머리 계산 수 평행

평행 머리가 가장 큰 것을 선택

Word2Vec은 NLP의 다수 영역에서 쓰는 방법

- 단어 유사도
- 기계 번역
- part-of-speech (pos) 태깅
- Named entity recognition (NER, 고유명사 인식)
- 강제 분석
- 전공화
- semantic lexicon building ...

## Glove

word2vec라 달리 사전에 각 단어들의 동시 등장 빈도수를 계산하여  
단어 임베딩 벡터의 내적 값과 사전에 계산된 값의 차이를 제곱하여 합을

$$J(\theta) = \frac{1}{2} \sum_{i,j=1}^N f(p_{ij}) (u_i^T v_j - \log p_{ij})^2$$

합을 속도가 빠르고 작은 corpus에서도 잘 동작  
(구독 예상은 중!)

