

Intro & Word Vectors

"Language is a social system,
constructed and interpreted by people."

trained on text data, neural machine translation is quite good

How do we represent the meaning of a word?

Commonest linguistic way of thinking of meaning:

Signifier (Symbol) \Leftrightarrow Signified (idea or thing)

= denotational semantics word \leftrightarrow thing = pairing

Common NLP solution: 'WordNet' - a thesaurus containing lists of
synonym sets and hypernyms (is a rel)
문제점 ←

- Nuance 반영 어려움
- 신조어 합침
- 주관적
- 인력 소모 大
- word similarity 정확히 계산 불가

단어들이 discrete symbols로 여겨지면서 one-hot vectors 사용

ex. 'motel' & 'hotel' are orthogonal (직교) in one-hot vec.
BUT one-hot vectors는 '유사성' 개념이 없음

Solution: Learn to encode similarity in the vectors themselves

Distributional semantics : 문맥의 미론

A word's meaning is given by
the words that frequently appear close-by

When a word w appears in a text, its **context** is the set of words that appear nearby

Word vectors : word embeddings, word representations

- build a dense vector for each word,
chosen so that it is similar to vectors of words
that appear in similar contexts

단어 차원을 **embed**

Word2vec

Idea: large corpus of text

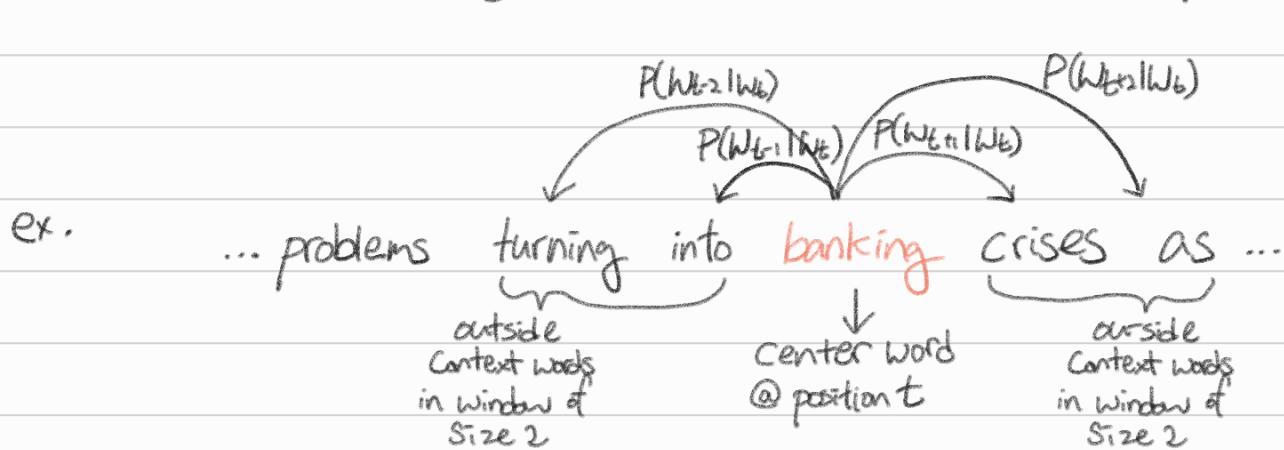
vector for each word

center word c , context words o

use the similarity of the word vectors for $c \& o$

to calculate the probability of o given c

keep adjusting the word vectors to maximize this probability $P(w_{t+1}|w_t)$



For each position $t=1, \dots, T$, predict context words within a window of fixed size m , given center word j

$$\text{Likelihood} = L(\theta) = \prod_{t=1}^T \prod_{\substack{-m \leq j \leq m \\ j \neq 0}} P(w_{t+j} | w_t; \theta)$$

The objective function $J(\theta)$ is the (average) negative log likelihood
loss, cost

$$J(\theta) = -\frac{1}{T} \log L(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log P(w_{t+j} | w_t; \theta)$$

Minimizing $J(\theta)$ = Maximizing predictive accuracy

How to calculate " $P(w_{t+j} | w_t; \theta)$ "?

$$\text{Word vector} \begin{cases} v_w & (\text{when } w \text{ is } c) \\ u_w & (\text{when } w \text{ is } o) \end{cases}$$

$$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$

- ② Exponential makes anything positive
- ① Dot product compares similarity of o and c (e.g., v_c)
- ③ Normalize over entire vocab to give probability distribution

이거 바로 Softmax function $\mathbb{R}^n \rightarrow (0, 1)^n$ 여기 1

$$\text{Softmax}(x_i) = \frac{\exp(x_i)}{\sum_{j=1}^n \exp(x_j)} = p_i$$

Soft: still assigns some probability to smaller x_i

Max: amplifies probability of largest x_i

Optimization done by gradient descent

$$\theta \in \mathbb{R}^{2dV}$$

v_c 에 대해 미분:

$$\frac{\partial}{\partial v_c} \log \frac{\exp(u_0^T v_c)}{\sum_{w=1}^V \exp(u_w^T v_c)}$$

$$\frac{\partial}{\partial v_c} \log \cancel{\exp(u_0^T v_c)} - \frac{\partial}{\partial v_c} \log \sum_{w=1}^V \exp(u_w^T v_c)$$

↓

↓ chain rule

$$\frac{\partial}{\partial v_c} u_0^T v_c = u_0$$

$$\frac{\sum_{x=1}^V \exp(u_x^T v_c) u_x}{\sum_{w=1}^V \exp(u_w^T v_c)}$$