

Word 2 Vec

- 주변에 같은 단어가 나타나는 단어일 수록 비슷한 벡터값을 가져야 한다.

문장 문맥에 따라 결정 X \rightarrow 코퍼스 $\xrightarrow{\text{학습}}$ 모델

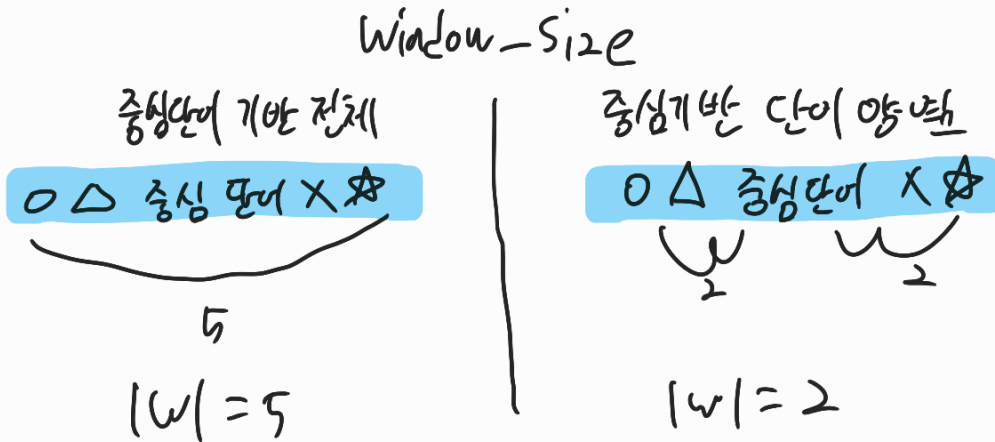
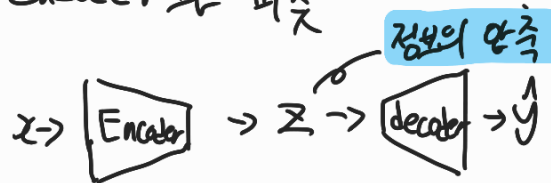
단어 \rightarrow fix된 임베딩 벡터

• CBOW, skip-gram

skip gram

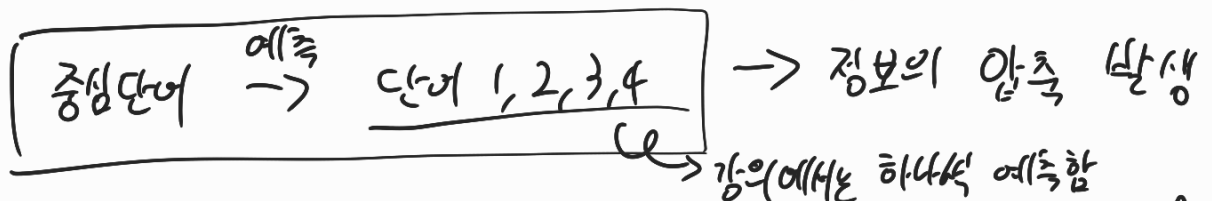
- 주변 단어 예측 \rightarrow 정보의 양측(임베딩) 가능
- Non-linear activation func 없음
- auto-encoder와 비슷

$|W| = \text{window-size}$



단어1, 단어2, 중심단어, 단어3, 단어4

classification



one-hot vector \rightarrow linear layer \rightarrow word embedding vector \rightarrow linear & softmax $\rightarrow \hat{y}$

장점 (before)

- 쉽다
- 빠르다
- 비교적 정확한 배제

단점 (현재)

- 느리다
- 출현빈도↓ 정확도↓

Glove

단어 \rightarrow [Encoder] \rightarrow Embedding vector \rightarrow Regression! [decoder] \rightarrow 각 단어 출현빈도 예측

- 출현빈도가 적은 단어에 대해서는 loss 기여도를 낮춤
↳ 출현빈도가 적은 단어에 대해 부정확해지는 단점을 보완

- 장점

- Skip-gram에 비해 빠르다

Glove: 전체 코퍼스에 대해 각 단어별 co-occurrence를 구한 후 regression 수행

- 출현빈도가 적은 단어에 대해서도 Good!

Fast text

facebook!

skip-gram의 upgrade version

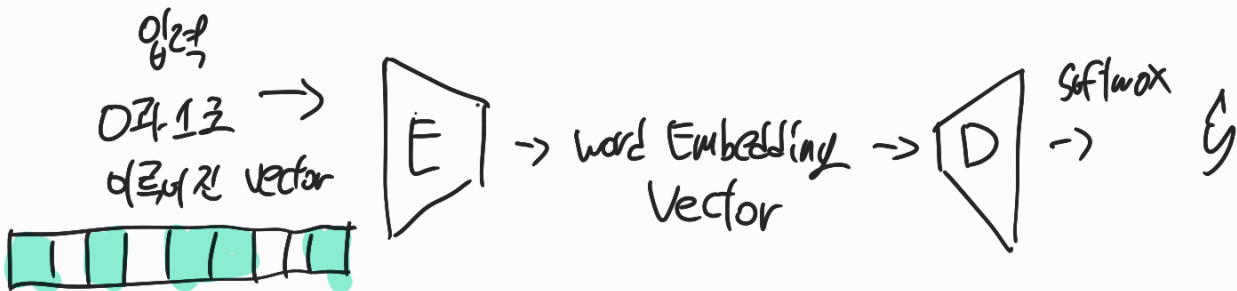
기존 skip-gram \rightarrow 저빈도 단어에 대한 학습 & oov에 대한 대처 어려움

fast text

- 1) 단어를 subword로 분절
- 2) **skip-gram** 활용 (subword Embedding vector \times ^{단어} 주변 context vector)
- 3) 주변 단어에 대한 확률 출력 학습

비율: Subword에 대한 Embedding vector의 한이 word Embedding vector

입력 vector \rightarrow ex) where = { <wh, whe, her, ere, re> }



Embedding vector summary

- 여러 알고리즘이 훨씬 뛰어나게 하는 것보다
· 구현 쉽고 빠른 오픈소스 사용이 best
- 2개의 다른 알고리즘(ex. glove fasttext) concat 하기도

Open-source

- Gensim
- Glove
- fasttext

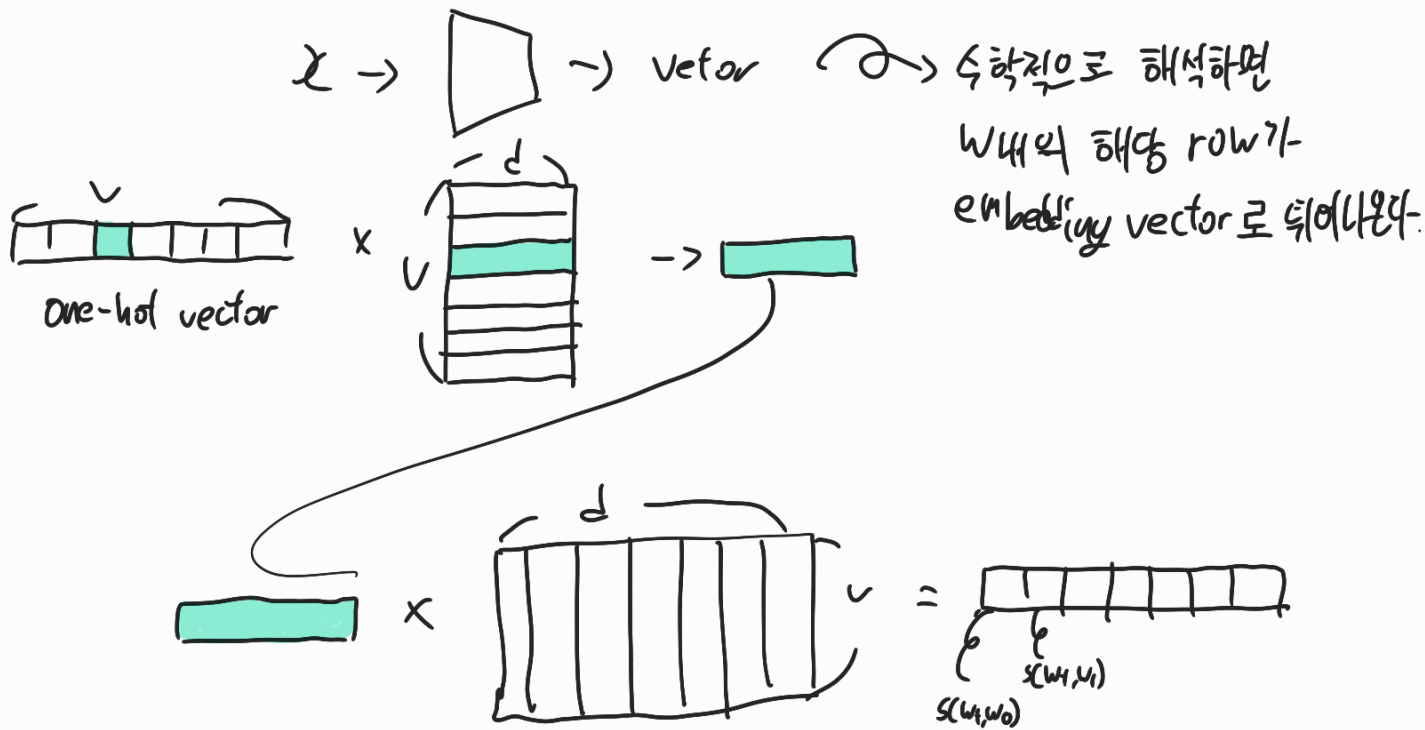
Skip-gram

$$\sum_{t=1}^T \sum_{c \in C_t} \log P(w_c | w_t)$$

Maximize (즉 최대화)
 즉 위 단어
 현재 단어
 t에서 윈도우 사이즈 만큼의 즉 위 단어
 w_t가 주어졌을 때 w_c의 확률

$$P(w_c | w_t) = \frac{e^{s(w_t, w_c)}}{\sum_{i=1}^{|V|} e^{s(w_t, w_i)}}$$

softmax
 where $s(w, w') = u_w^T v_{w'}$



Negative Sampling (Softmax 디지레(용))

$$\log(1 + e^{-s(w_t, w_c)}) + \sum_{n \in N} \log(1 + e^{s(w_t, w_n)})$$

Minimize
 (정답) (노답)
 정답

Glove

* regression

$$\Theta = \underset{\Theta \in \Theta}{\operatorname{argmin}} \sum_{x \in X} \underbrace{f(x)}_{\text{penalty}} \times \underbrace{\| \underbrace{w'}_{\text{decoder}} \underbrace{w}_x - \log \underbrace{c_x}_{\text{co-occurrences with } x \text{ vector in other word } y} \|_2^2}_{\text{MSE}}$$

$$f(x) = \begin{cases} (\text{count}(x)/\text{thres})^\alpha & \text{if } \text{count}(x) < \text{thres} \\ 1 & \text{otherwise} \end{cases}$$

각각 나타나는 단어 loss를 줄이며 적은 단어에 대한 과적합을 줄임

Fast Text

input: sum of sub-word one-hot vector

$$\sum_{t=1}^T \sum_{c \in C_t} \log p(w_c | w_t), \quad p(w_c | w_t) = \frac{e^{s(w_t, w_c)}}{\sum_{j=1}^{|V|} e^{s(w_t, w_j)}}$$

$$s(w, w') = \sum_{g \in g_w} z_g^T v_{w'}$$

where g_w is a set of subword n-grams

$z_g^T =$

