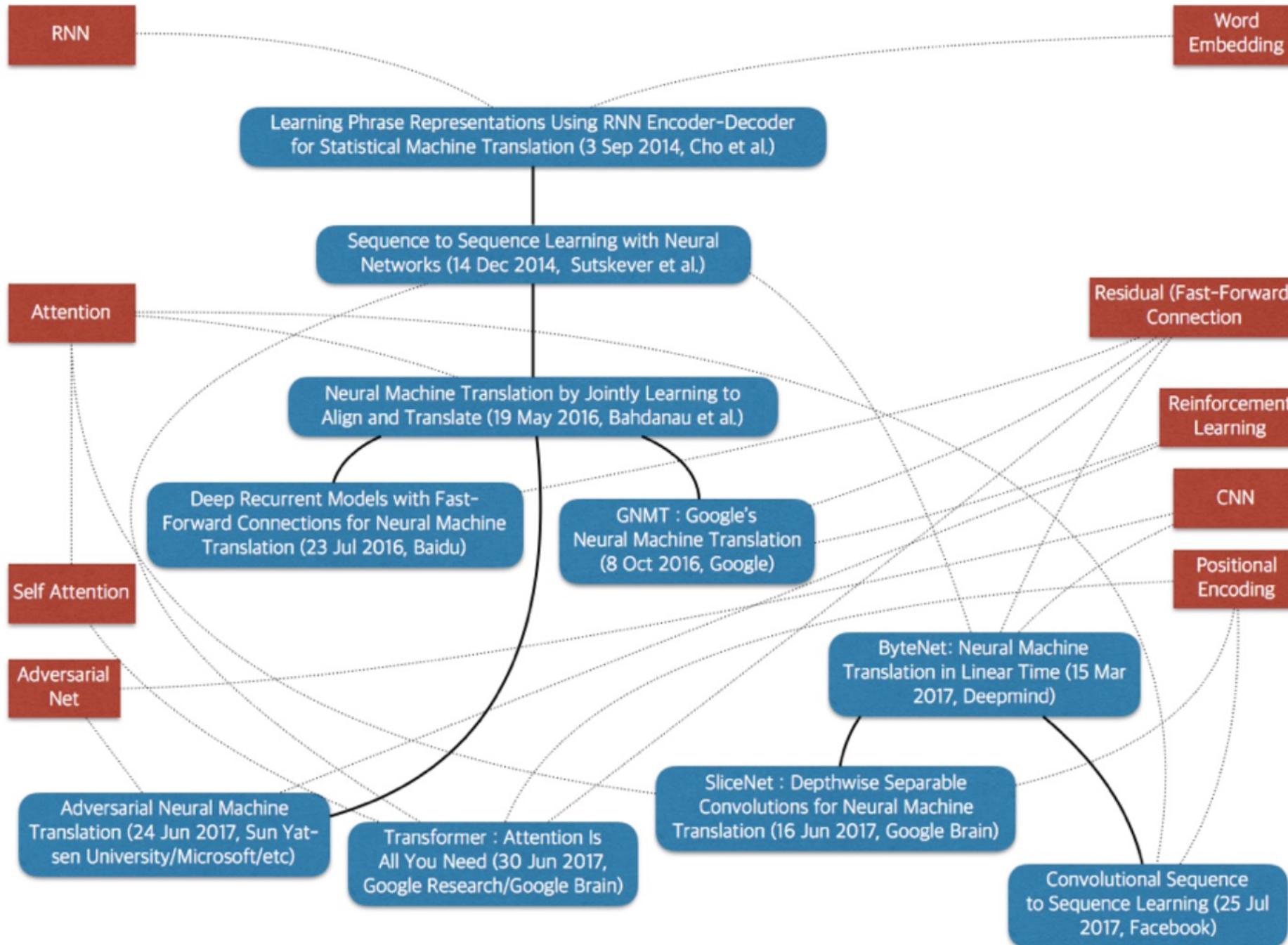


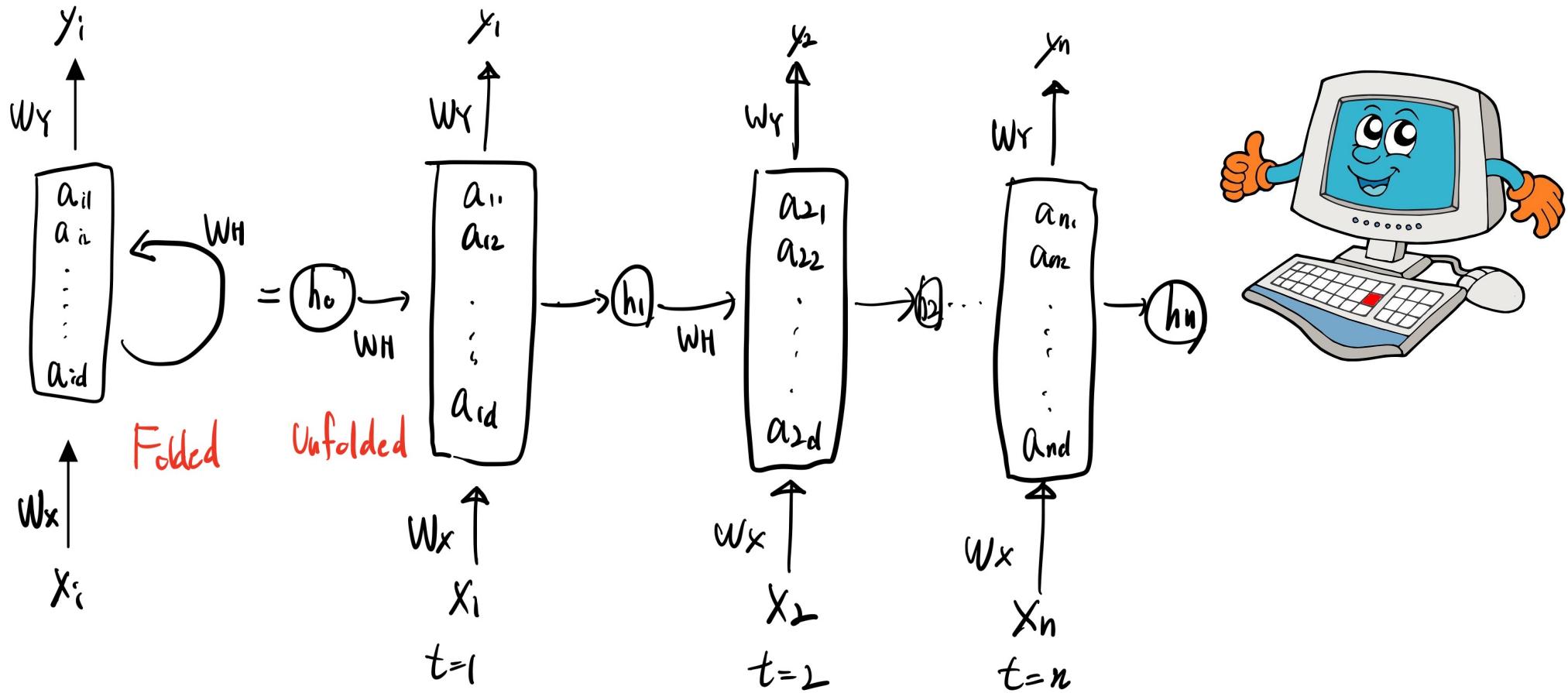
Embedding:

How to Vectorize Text Data





Reminder

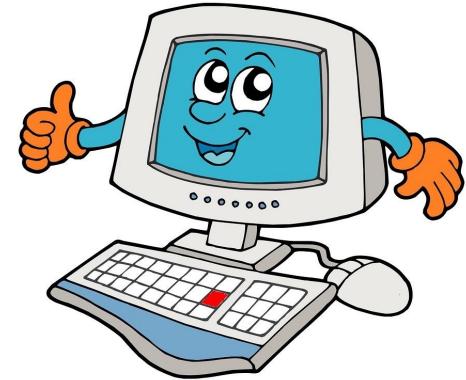
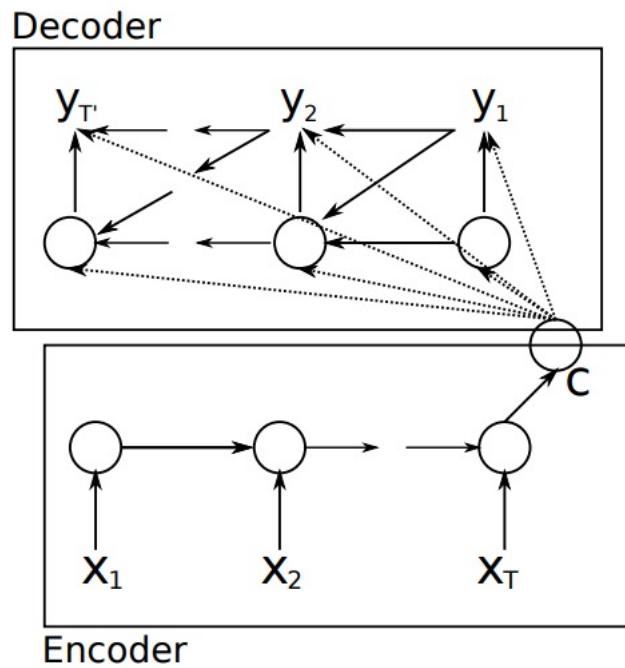


RNN Encoder-Decoder

만나서 반갑습니다



nice to meet you



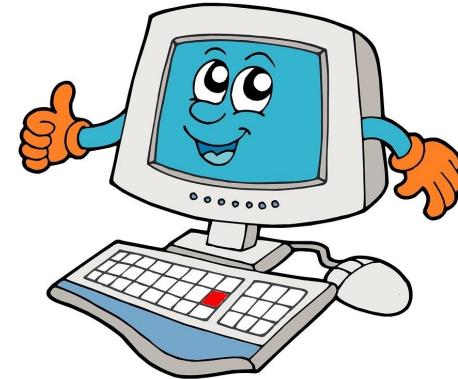
one-hot vector

너구리 : [1, 0, 0]

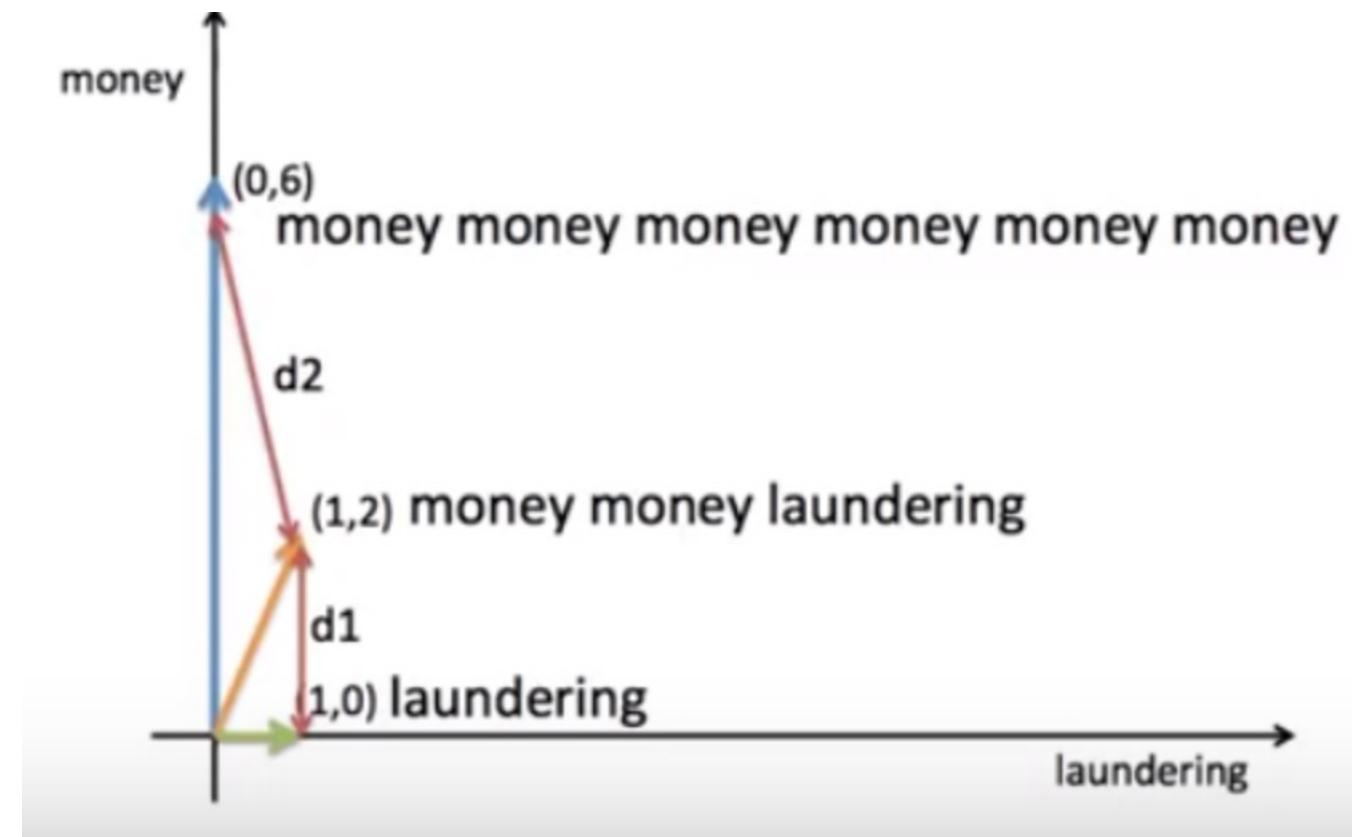
라쿤 : [0, 1, 0]

사람: [0, 0, 1]

1. (데이터가 늘어날수록) 아주 커짐
2. 내적하면 모두 0 (서로 독립)



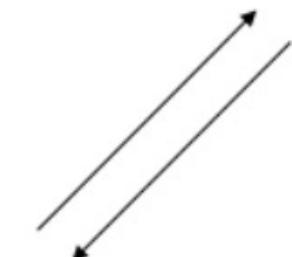
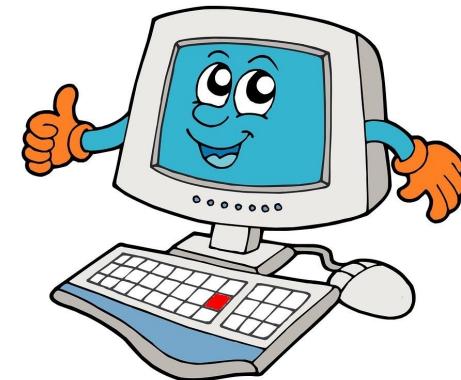
one-hot vector



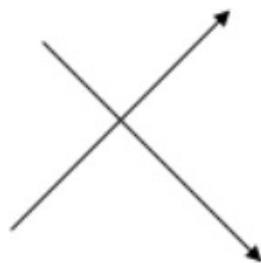
one-hot vector

너구리 : [1, 0, 0]
라쿤 : [0, 1, 0]
사람: [0, 0, 1]

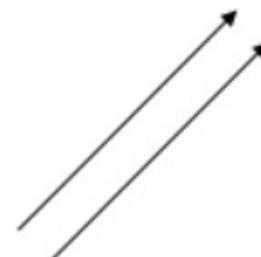
$$\text{similarity} = \cos(\theta) = \frac{A \cdot B}{\|A\|\|B\|} = \frac{\sum_{i=1}^n A_i \times B_i}{\sqrt{\sum_{i=1}^n (A_i)^2} \times \sqrt{\sum_{i=1}^n (B_i)^2}}$$



코사인 유사도 : -1



코사인 유사도 : 0



코사인 유사도 : 1



코드 리뷰

NNLM (Neural Network Language Model)

language model: 주어진 단어로 다음 단어를 예측하는 모델

A neural probabilistic language model. Journal of Machine Learning Research (Bengio et al, 2003)

statistical model of language:

language can be represented by the conditional probability of the next word given all the previous ones.

$$\hat{P}(w_1^T) = \prod_{t=1}^T \hat{P}(w_t | w_1^{t-1}),$$
$$w_i^j = (w_i, w_{i+1}, \dots, w_{j-1}, w_j)$$

A neural probabilistic language model. Journal of Machine Learning Research (Bengio et al, 2003)

fact: temporally closer words in the word sequence are statistically more dependent
-> n-gram models construct tables of conditional probabilities for the next word

$$\hat{P}(w_t | w_1^{t-1}) \approx \hat{P}(w_t | w_{t-n+1}^{t-1}).$$

A neural probabilistic language model. Journal of Machine Learning Research (Bengio et al, 2003)

1. 각각의 단어를 분산된 word feature vector로 본다.
2. 시퀀스의 feature vector를 시퀀스의 joint probability function으로 표현한다.
3. 각각의 word feature vector와 probability function의 파라미터를 동시에 학습시킨다.

그렇다면 확률 함수를 이전 단어들이 주어졌을 때 이후 단어의 conditional probability로 표현 가능함.
이를 학습 데이터들의 Log-likelihood를 최대화하는 식으로 학습 가능.

A neural probabilistic language model. Journal of Machine Learning Research (Bengio et al, 2003)

$w_1 \cdots w_T$ of words $w_t \in V$, training set, where the vocabulary V is a large but finite set.

$f(w_t, \dots, w_{t-n+1}) = \hat{P}(w_t | w_1^{t-1})$. objective: learn a good model

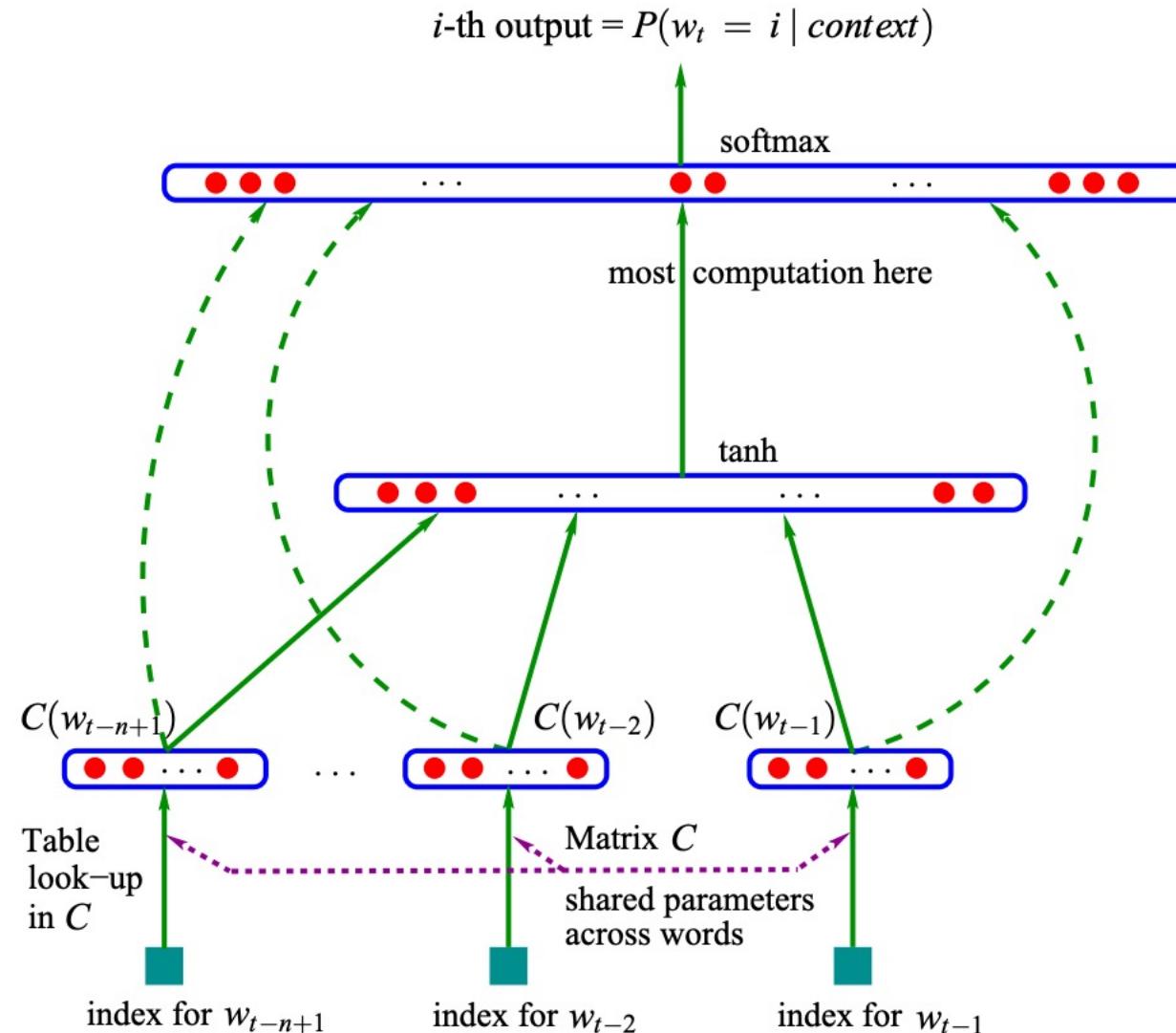
A neural probabilistic language model. Journal of Machine Learning Research (Bengio et al, 2003)

We decompose the function $f(w_t, \dots, w_{t-n+1}) = \hat{P}(w_t | w_1^{t-1})$ in two parts:

1. A mapping C from any element i of V to a real vector $C(i) \in \mathbb{R}^m$. It represents the *distributed feature vectors* associated with each word in the vocabulary. In practice, C is represented by a $|V| \times m$ matrix of free parameters.
2. The probability function over words, expressed with C : a function g maps an input sequence of feature vectors for words in context, $(C(w_{t-n+1}), \dots, C(w_{t-1}))$, to a conditional probability distribution over words in V for the next word w_t . The output of g is a vector whose i -th element estimates the probability $\hat{P}(w_t = i | w_1^{t-1})$ as in Figure 1.

$$f(i, w_{t-1}, \dots, w_{t-n+1}) = g(i, C(w_{t-1}), \dots, C(w_{t-n+1}))$$

A neural probabilistic language model. Journal of Machine Learning Research (Bengio et al, 2003)



A neural probabilistic language model. Journal of Machine Learning Research (Bengio et al, 2003)

$$L = \frac{1}{T} \sum_t \log f(w_t, w_{t-1}, \dots, w_{t-n+1}; \theta) + R(\theta),$$

$$\hat{P}(w_t | w_{t-1}, \dots, w_{t-n+1}) = \frac{e^{y_{w_t}}}{\sum_i e^{y_i}}.$$

A neural probabilistic language model. Journal of Machine Learning Research (Bengio et al, 2003)

$$x_t = C \cdot w_t$$

when $V=3$.

V : 알파벳의 수
 m : x_t 의 차원수

$$\begin{array}{c} C \\ \left(\begin{array}{ccc} a_{11} & a_{12} & a_{13} \\ a_{21} & & \vdots \\ \vdots & \vdots & \\ a_{m1} & \dots & a_{m3} \end{array} \right) \end{array} \quad \begin{array}{c} w_t \\ \left[\begin{array}{ccc} 1 & 0 & a \\ a & 1 & 0 \\ 0 & 0 & 1 \end{array} \right] \end{array} \quad \Rightarrow \quad \begin{array}{c} y_{wt} \\ \left[\begin{array}{c} a \\ a \\ 0 \end{array} \right] \end{array}$$

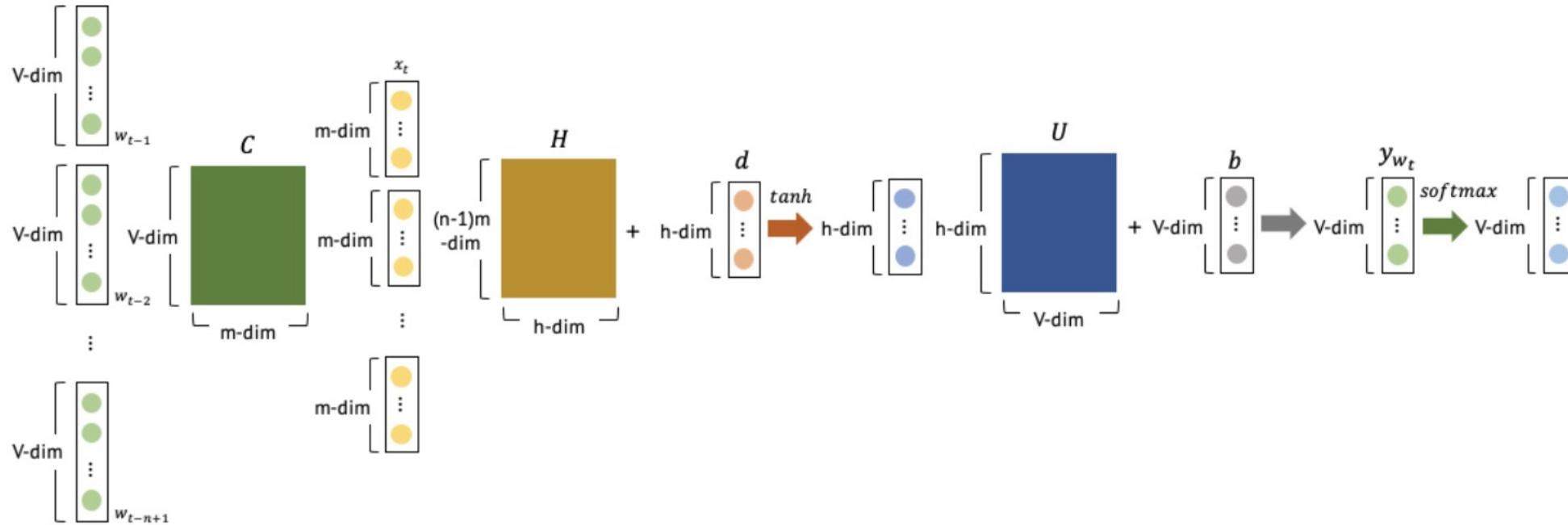
$m \times 3$ $3 \times (V-1)$

A neural probabilistic language model. Journal of Machine Learning Research (Bengio et al, 2003)

$$y_{w_t} = b + U \cdot \tanh(d + Hx_t)$$

$$\begin{aligned} H &\in R^{h \times (n-1)m}, & x_t &\in R^{(n-1) \times m}, & d &\in R^{h \times 1} \\ U &\in R^{|V| \times h}, & b &\in R^{|V|}, & y &\in R^{|V|}, & C &\in R^{m \times |V|} \end{aligned}$$

A neural probabilistic language model. Journal of Machine Learning Research (Bengio et al, 2003)



A neural probabilistic language model. Journal of Machine Learning Research (Bengio et al, 2003)

장점

- 저장 공간 이점
- 단어의 유사도 표현

단점

- 많은 파라미터 training → 과적합 문제, 계산 비용
- n개만의 단어 참고

Efficient Estimation of Word Representations in Vector Space (Mikolov et al. 2013)

computational complexity per each training example

$$Q = N \times D + \underbrace{N \times D \times H}_{\text{projection layer}} + H \times V$$

↑ ↑

projection layer expensive!

N: previous word → one-hot encoding

V : 말뭉치의 단어 수

P: N X D shared projection matrix

we try to minimize computational complexity.

Efficient Estimation of Word Representations in Vector Space (Mikolov et al. 2013)

CBOW (Continuous Bag-of-Words Model)

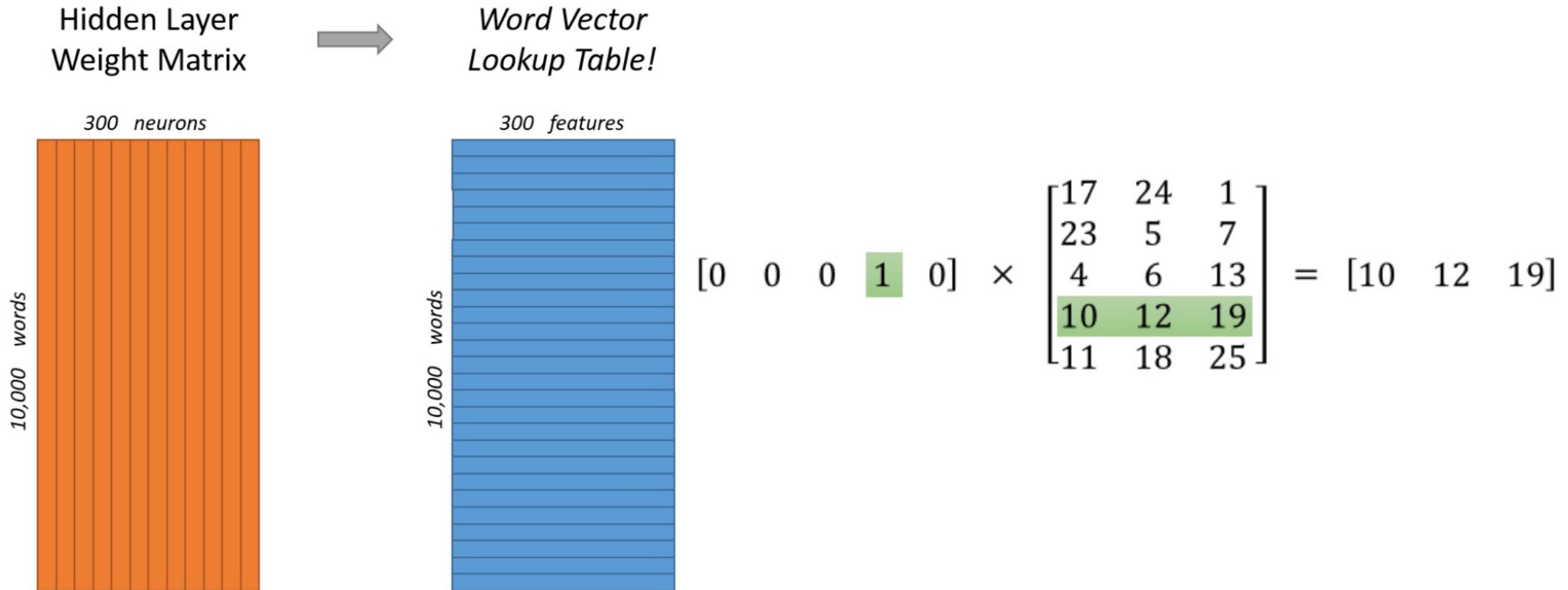
1. remove non-linear hidden layer
2. projection layer is shared for all words (order of words does not influence)
3. use words from the future

$$Q = N \times D + D \times \log_2(V).$$

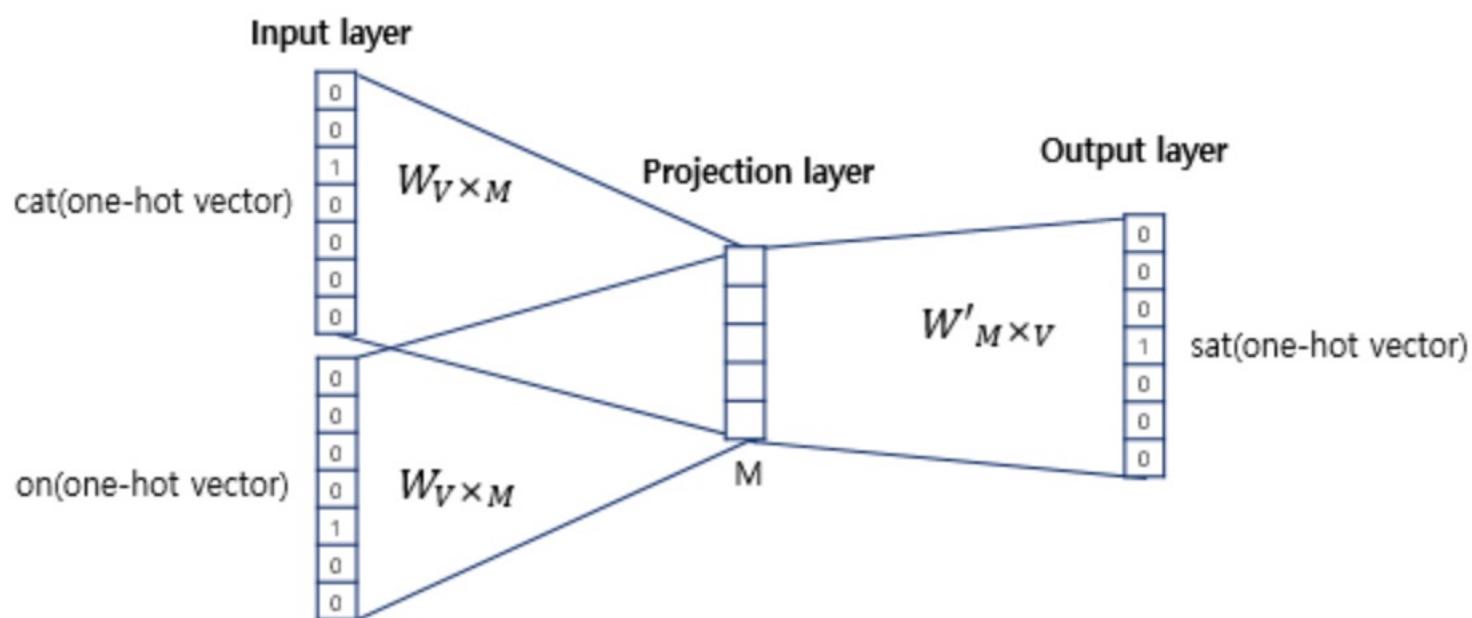
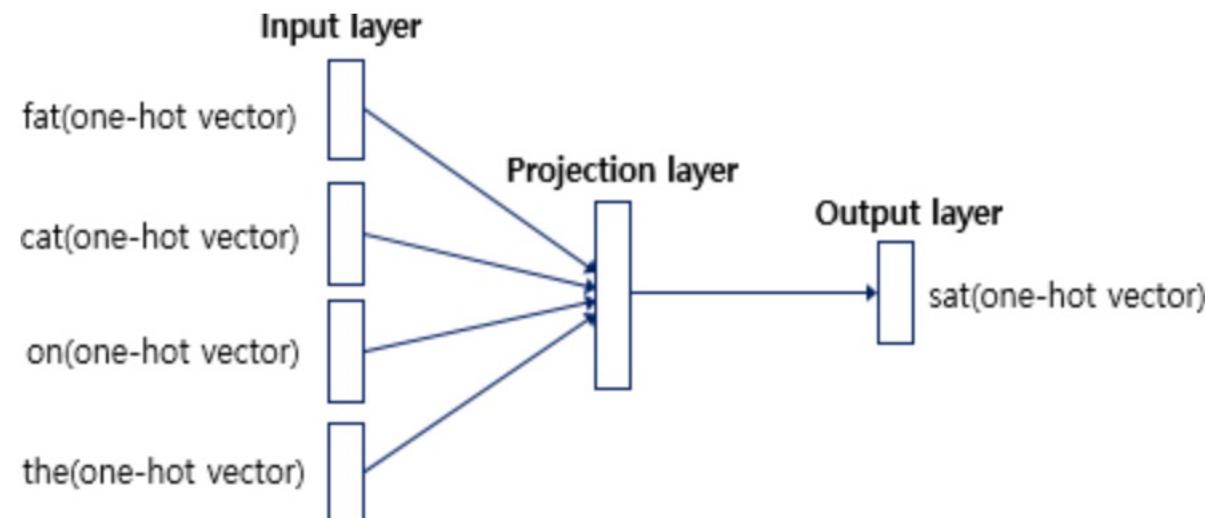
Efficient Estimation of Word Representations in Vector Space (Mikolov et al. 2013)

Source Text	Training Samples
The quick brown fox jumps over the lazy dog. →	(the, quick) (the, brown)
The quick brown fox jumps over the lazy dog. →	(quick, the) (quick, brown) (quick, fox)
The quick brown fox jumps over the lazy dog. →	(brown, the) (brown, quick) (brown, fox) (brown, jumps)
The quick brown fox jumps over the lazy dog. →	(fox, quick) (fox, brown) (fox, jumps) (fox, over)

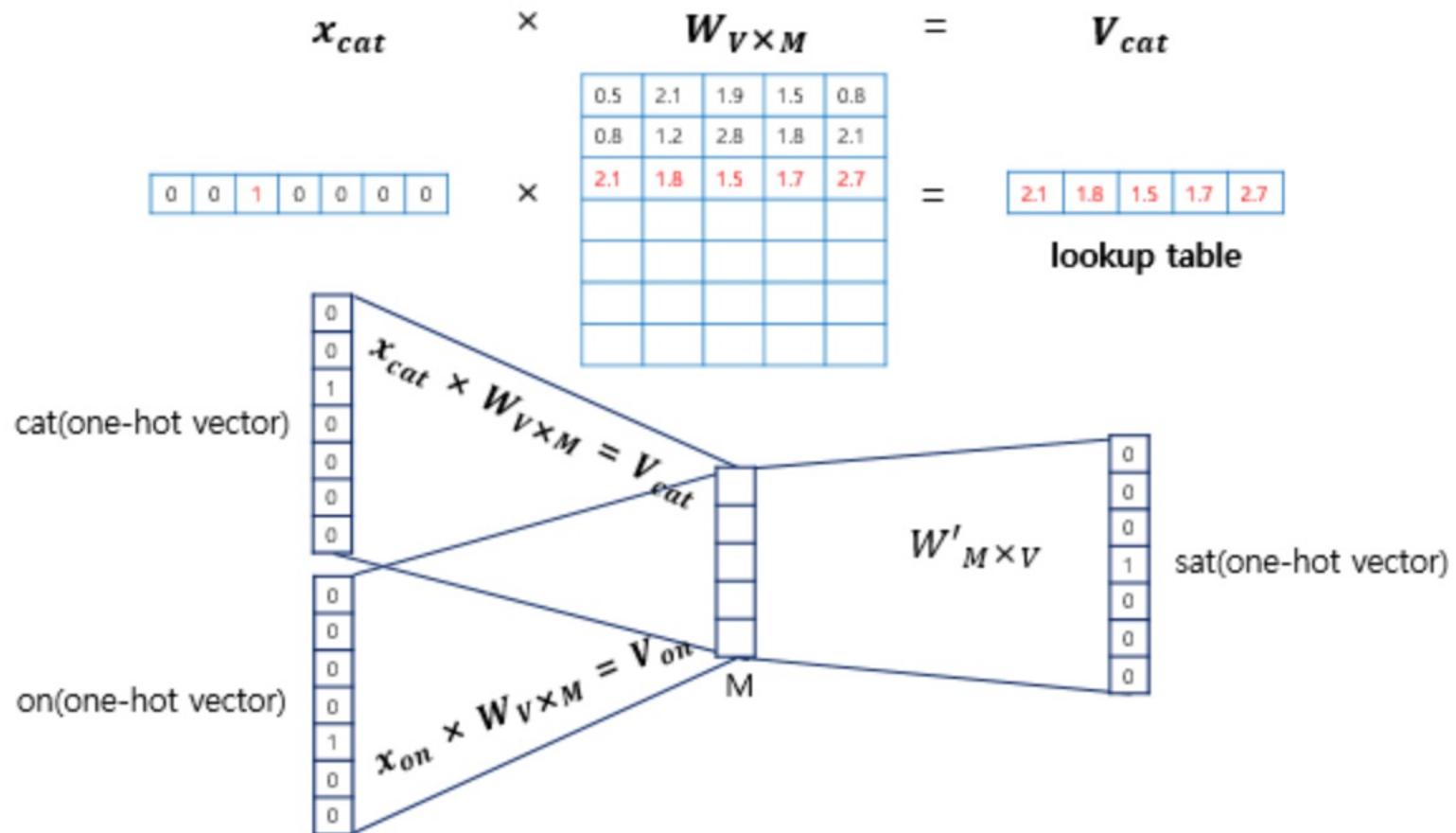
Efficient Estimation of Word Representations in Vector Space (Mikolov et al. 2013)



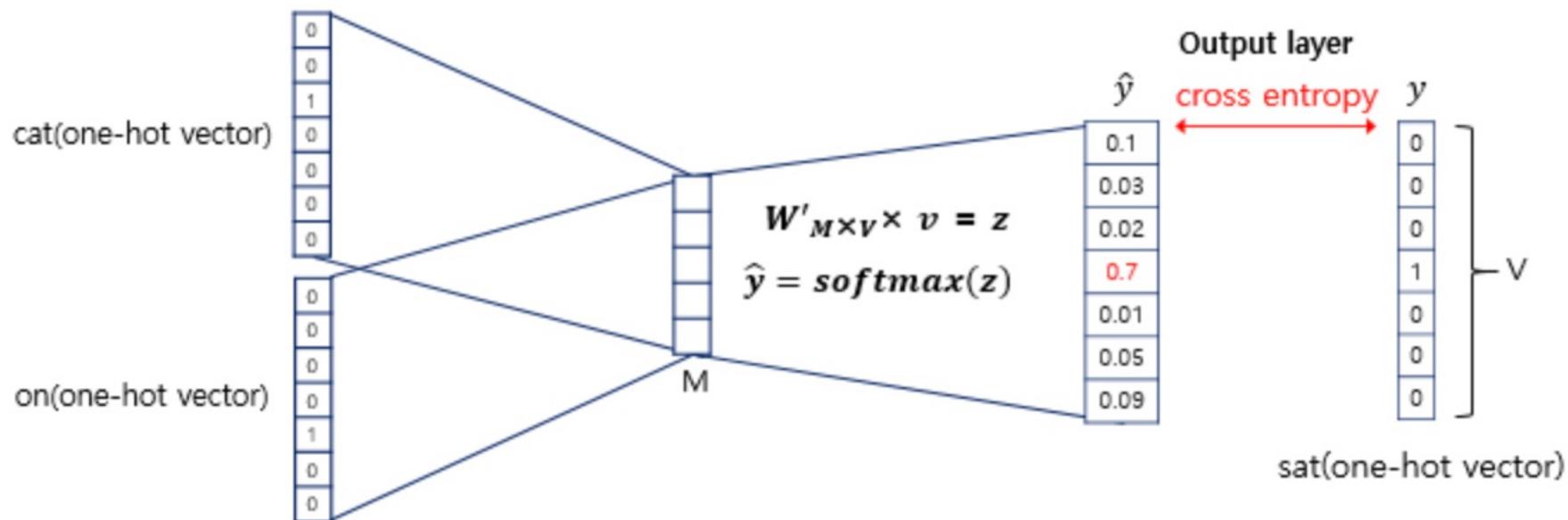
Efficient Estimation of Word Representations in Vector Space (Mikolov et al. 2013)



Efficient Estimation of Word Representations in Vector Space (Mikolov et al. 2013)



Efficient Estimation of Word Representations in Vector Space (Mikolov et al. 2013)



$$H(\hat{y}, y) = - \sum_{j=1}^{|V|} y_j \log(\hat{y}_j)$$

Efficient Estimation of Word Representations in Vector Space (Mikolov et al. 2013)

$$\begin{aligned} \text{minimize } J &= -\log P(w_c | w_{c-m}, \dots, w_{c+m}) \\ &= -\log P(u_c | v) \\ &= -\log \frac{\exp(u_c^\top \hat{v})}{\sum_{j=1}^{|V|} \exp(u_j^\top \hat{v})} \\ &= -u_c^{intercal} \hat{v} + \log \sum_{j=1}^{|V|} \exp(u_j^\top \hat{v}) \end{aligned}$$

Efficient Estimation of Word Representations in Vector Space (Mikolov et al. 2013)

$$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w=1}^W \exp(u_w^T v_c)}$$

$$\begin{aligned}\frac{\partial}{\partial v_c} \ln P(o|c) &= \frac{\partial}{\partial v_c} \ln \frac{\exp(u_o^T v_c)}{\sum_{w=1}^W \exp(u_w^T v_c)} \\&= \frac{\partial}{\partial v_c} u_o^T v_c - \frac{\partial}{\partial v_c} \ln \sum_{w=1}^W \exp(u_w^T v_c) \\&= u_o^T - \frac{1}{\sum_{w=1}^W \exp(u_w^T v_c)} \left(\sum_{w=1}^W \exp(u_w^T v_c) \cdot u_w \right) \\&= u_o^T - \sum_{w=1}^W \frac{\exp(u_w^T v_c)}{\sum_{w=1}^W \exp(u_w^T v_c)} \cdot u_w \\&= u_o^T - \sum_{w=1}^W P(w|c) \cdot u_w \\v_c^{t+1} &= v_c^t + \alpha(u_o^T - \sum_{w=1}^W P(w|c) \cdot u_w)\end{aligned}$$

Efficient Estimation of Word Representations in Vector Space (Mikolov et al. 2013)

subsampling : 1 번째 단어를 학습에서 제외 (계산량 감소)

$$P(w_i) = 1 - \sqrt{\frac{t}{f(w_i)}}$$

Efficient Estimation of Word Representations in Vector Space (Mikolov et al. 2013)

negative sampling : 전체 단어를 구하지 않고, 일부 단어를 계산

$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=0}^n f(w_j)^{3/4}}$$

superhero dataset에서 superman과 history가 가장 비슷한 히어로 10명 찾기

<https://www.kaggle.com/jonathanbesomi/superheroes-nlp-dataset>