

Text Analytics

Ch5 : Text Representation



서수원

Business Intelligence Lab.
산업경영공학과, 명지대학교

01

NNLM

Word Embedding

- Word Embedding
 - 단어를 특정한 공간의 벡터 스페이스로 맵핑한다.
 - 의미적으로 유사한 단어들이 벡터 공간상 각각에 근처에 위치하게 된다.
 - ✓ 이론적으로 영어의 토큰은 1300만개가 있는데, 이는 1300만 차원이 필요하다는 의미이다.
 - ❖ 이는 비효율이기 때문에, 분산표상을 이용한다.

Word Embedding

- Word vectors : one-hot vector
 - 가장 쉽고 이해하기 쉬운 표현 방법이다.

$$w^{aardvark} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, w^a = \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, w^{at} = \begin{bmatrix} 0 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}, \dots w^{zebra} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}$$

$$\begin{array}{l} \text{motel} [0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0] \\ \text{hotel} [0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0] \end{array} \text{ AND } = 0$$

$$(w^{hotel})^T w^{motel} = (w^{hotel})^T w^{cat} = 0$$

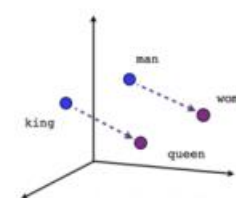
- Word vectors : distributed representation
 - 분산표상, 특정 언어를 특정 차원을 갖는 벡터로 변환하는 매개변수 함수를 의미한다.
 - ✓ $N \ll |V|$ (V : voca.)

$$W : \text{words} \rightarrow \mathbb{R}^n$$

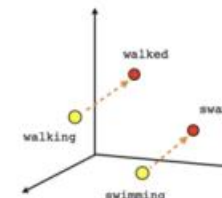
$$W(\text{"cat"}) = (0.2, -0.4, 0.7, \dots)$$

$$W(\text{"mat"}) = (0.0, 0.6, -0.1, \dots)$$

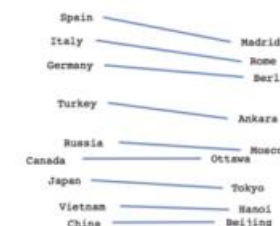
- Word embedding의 특징
 - 단어의 의미적 유사성이 벡터 공간에 보존된다.



Male-Female



Verb tense



Country-Capital

NNLM

- NNLM의 목적
 - 분산표상을 처음으로 사용한 방법론이다.
 - One-Hot vector 의 단점을 분산표상으로 해결함이 목적이다.
 - ✓ One-Hot vector는 희소성의 문제, 저장공간의 문제 등이 있다.
 - 각 단어들을 분산표상에 대한 워드벡터로 표현이 가능하다.
 - 단어의 연속성에 의한 결합확률 분포를 통해 단어가 나올 확률을 알 수 있다.
 - 어떤 벡터가 좋은지에 대한 것과, 단어의 연관성이 높은지에 대한 확률을 동시에 계산한다.

- NNLM

- 분산 표상을 통해, (강아지,고양이),(the,a),(bedroom,room)의 비슷한 의미를 이해 할 수 있다.
 - ✓ 이는 아래와 같은 새로운 문장을 생성할 수 있음을 의미한다.

The cat is walking in the bedroom

to

A dog was running in a room

The cat is running is a room

A dog is walking in a bedroom

The dog was waling in the room

...

- Comparison with Count-based Language Models
 - Count-based Language Models

✓ 체인룰에 의해 모든 분포는 밑과 같이 표현될 수 있다.

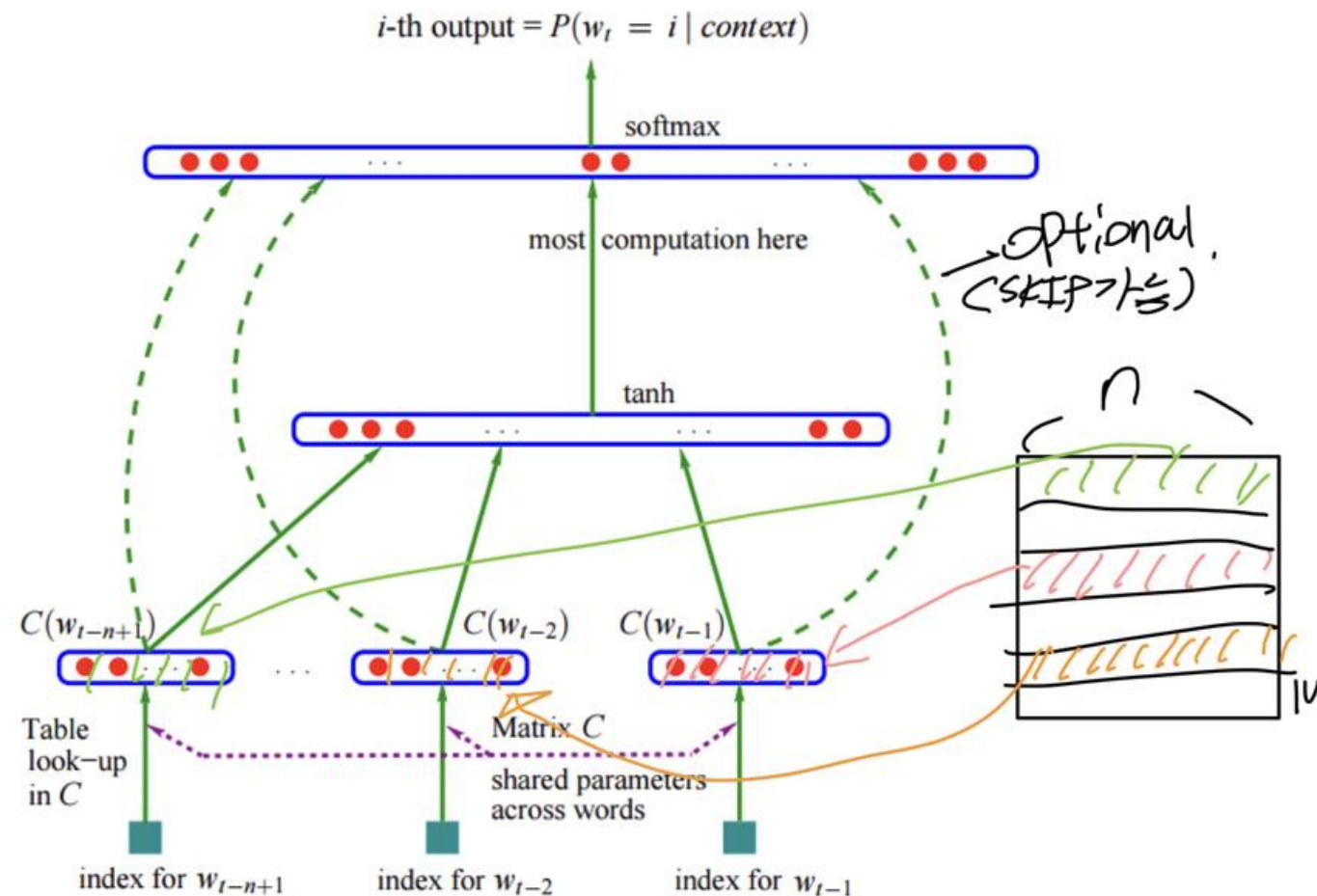
$$p(w_1, \dots, w_T) = \prod_{t=1}^T p(w_t | w_1, \dots, w_{t-1})$$

✓ 카운트 기반 n-gram language models는 마르코프 가정을 이용한다.

$$p(w_t | w_1, \dots, w_t) \approx p(w_t | w_{t-n}, \dots, w_{t-1})$$

NNLM

- NNLM



- NNLM

$$W(\text{"cat"}) = (0.2, -0.4, 0.7, \dots)$$

- 단어를 $R^n (\ln|V|)$ 공간의 차원에 dense vectors로 표현한다..(실수로 되어있는 벡터)
- W_t : One-hot 벡터로 표현된 것을 의미한다.
- 이때 $x_t = XW_t$: Word embedding

$$x_t = \boxed{X} W_t$$

one-hot vector.

look up + able

word embedding

$$P(W_t = j | W_1, \dots, W_{t-1}) = \frac{\exp(p_j \cdot g(x_1, \dots, x_{t-1}) + q_j)}{\sum_{j'} \exp(p_{j'} \cdot g(x_1, \dots, x_{t-1}) + q_{j'})} = \text{softmax}(P_j(x_1, \dots, x_{t-1}) + q_j)$$

normal network

j

t-1개까지의 문장

확률값을 계산한다

- P_j, q_j = Output word embedding, g = NN

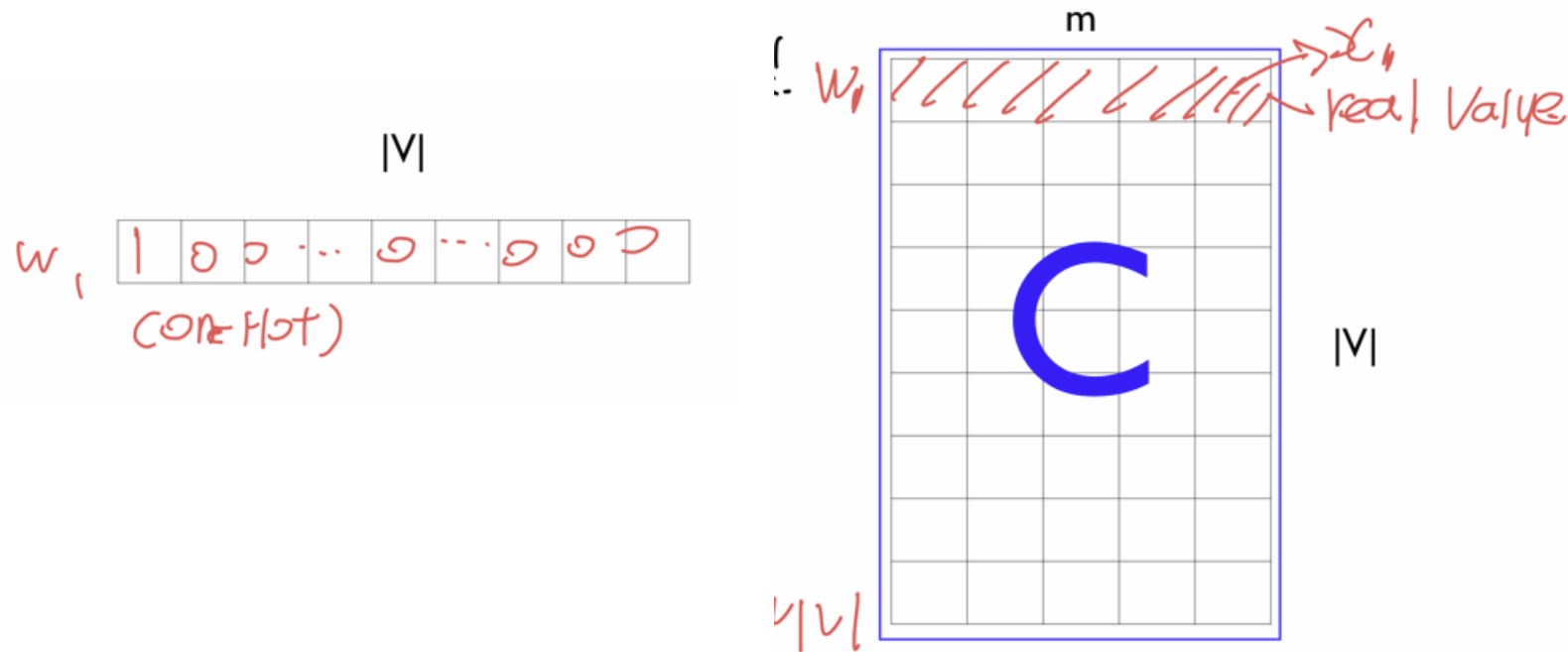
- NNLM

- NNLM의 목적은 높은 우도값을 가지는 $f(w_t, \dots w_{t-n+1}) = \hat{P}(w_t | w_1^{t-1})$, 모델을 만드는 것이다.
 - ✓ W 는 windows를 의미한다.
- 두가지 제약 조건이 존재한다.
 - ✓ $w_1^{t-1}, \sum_{i=1}^{|V|} f(i, w_{t-1} \dots w_{t-n+1}) = 1$ 어떤 조건에서도 이후 단어들이 생성될 확률의 총합은 1이다.
 - ✓ $f \geq 0$ (단어가 생성될 확률은 0보다 크거나 같다.)

NNLM

- NNLM

$$f(w_t, \dots w_{t-n+1}) = \hat{P}(w_t | w_1^{t-1}),$$



- NNLM
 - C는 lookup table에서 꺼내쓰는 벡터이다.
 - $F(w_t \dots)$ 에서 w_t 는 단어 이다.
 - G는 NN부분이다.
 - ✓ I값을 최대로 하는 것이 목표이다.

$$f(w_t, \dots, w_{t-n+1}) = \hat{P}(w_t | w_1^{t-1}), \quad f(i, w_{t-1}, \dots, w_{t-n+1}) = g(C_i, \boxed{C(w_t)}, \dots, C(w_{t-n+1}))$$

(Handwritten diagram and notes below the equation):
 - A red box highlights $C(w_t)$ in the equation.
 - A red arrow points from the box to the text "다음번 Index" (Next Index).
 - A red arrow points from the text "입력 벡터" (Input vector) to the sequence of vectors $C(w_t), \dots, C(w_{t-n+1})$.
 - A red box labeled "hidden" is connected to the input vectors.
 - A red box labeled "output" is connected to the hidden layer.

- NNLM



$g(\text{준다}|\text{너에게, 나의 입술을, 처음으로}) = ?$

$g(\text{지운다}|\text{너에게, 나의 입술을, 처음으로}) = ?$

$g(\text{말킨다}|\text{너에게, 나의 입술을, 처음으로}) = ?$

• NNLM

- Hidden layer을 통할 수도, 그냥 지나칠수도 있다.
 - ✓ Wx 가 optional 이라는 것에서 확인 할 수 있다.
- 바로 밑 식으로 soft-max이고, 합이 1이 되게 한다.
 - ✓ Output lv 의 확률값이 1이 된다.
 - ❖ 제약 조건을 확인하면 된다.

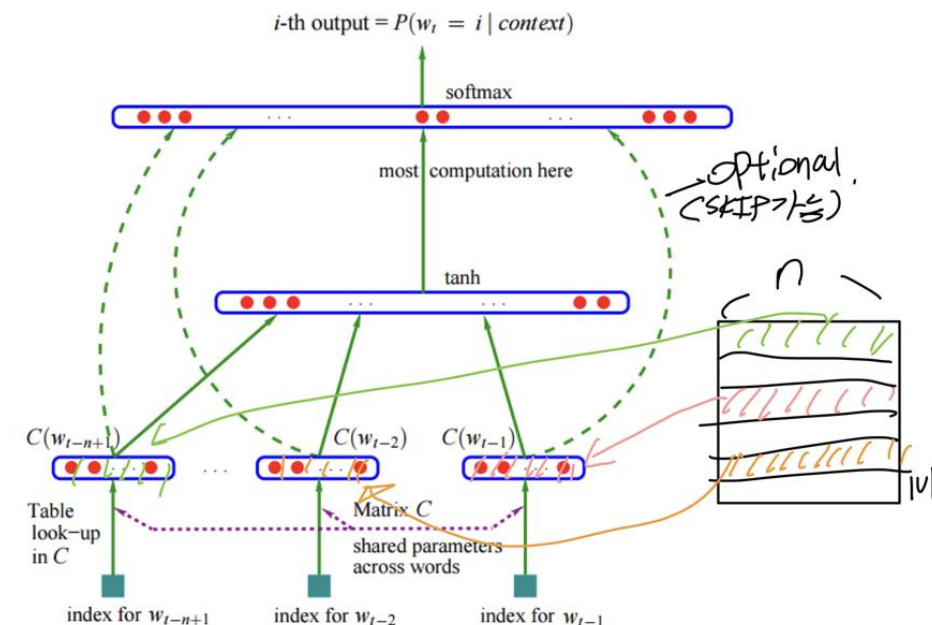
$$\hat{P}(w_t | w_{t-1}, \dots, w_{t-n+1}) = \frac{\exp(y_{w_t})}{\sum_i \exp(y_i)}$$

Handwritten notes and formula:

$$y = b + \boxed{Wx} + U \cdot \tanh(WHx)$$

Annotations:

- \boxed{Wx} : optional
- Wx : 가중치
- $U \cdot \tanh(WHx)$: 히든 레이어의 편향 (상수)
- W : 가중치
- U : 가중치
- WHx : 히든 레이어의 편향 (상수)
- y : 출력
- b : 편향



02

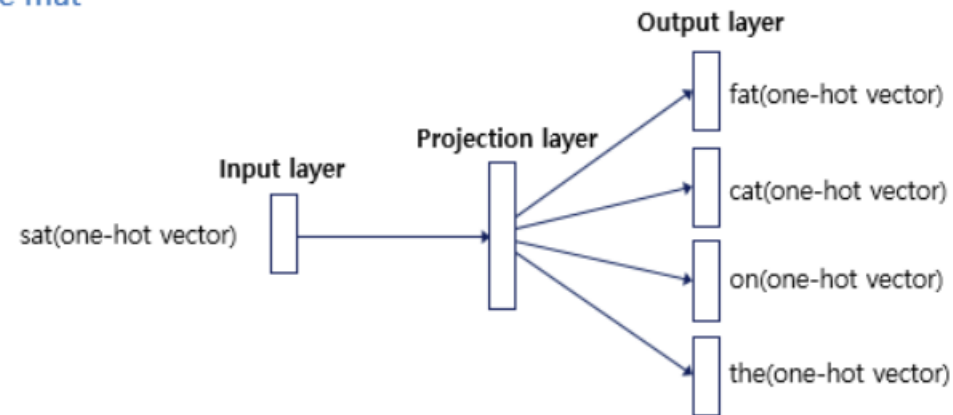
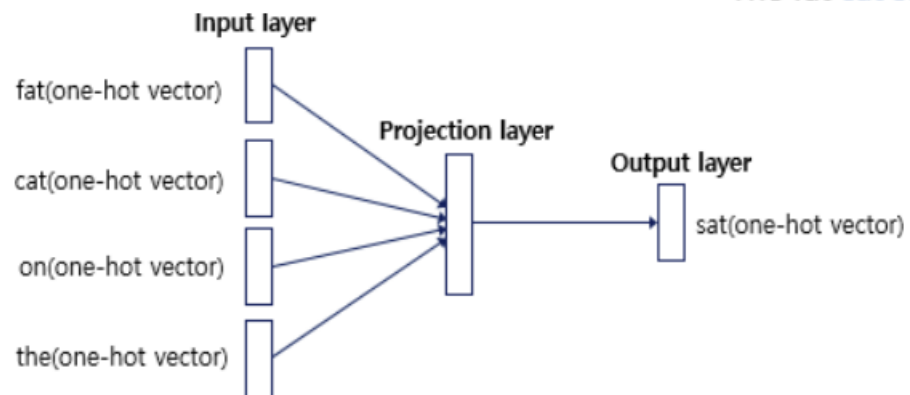
Word2Vec

- Word2Vec : Two Architectures

- CBOW vs Skip-gram

- ✓ CBOW는 주변에 있는 단어들을 입력으로 중간에 있는 단어를 예측하는 방법이다.
 - ❖ 밑 예는 윈도우 크기가 2일때 예시이다.
 - ✓ Skip-gram은 중간에 있는 단어들을 입력으로 주변 주변 단어를 예측하는 방법이다.

The fat cat sat on the mat



- Word2Vec : Skip-gram

- 원도우의 크기 만큼의 주변 단어들을 예측한다.
- 목적함수는 현재 중앙값으로부터, 주변 단어의 log확률을 최대화 하는 것이다.

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

Handwritten version of the equation with annotations:

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

Annotations in red:

- $J(\theta)$ is circled and labeled "최대화" (Maximize).
- The summation index j is labeled with a red arrow and "대칭" (Symmetric).
- The word w_t in the denominator is labeled with a red arrow and "center".

- Word2Vec : Skip-gram

- 한번에 모든 주변 단어를 다 구해도 되지만, 하나씩 구하더라도 gradient의 값은 동일하다.



knight → The
 knight → mighty
 knight → Lancelot
 knight → fought
 knight → bravely.

- 주어진 단어에 대한 주변 단어의 모델 확률은 밑과 같다.
 - ✓ 분자는 C번째 단어의 내적을 연산한 것과 분모는 모든 단어에 대한 내적을 의미한다.(softmax)

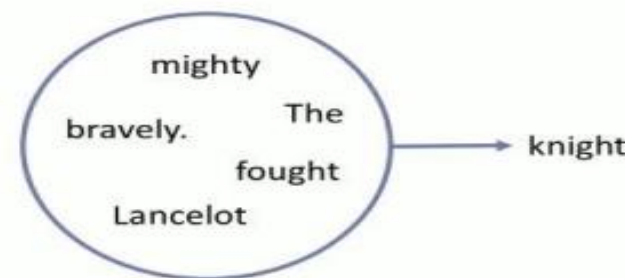
feature for word w : x_w
 classifier for word c : v_c

$$p(c|w) = \frac{e^{x_w^\top v_c}}{\sum_{k=1}^K e^{x_w^\top v_k}}$$

- 워드 벡터의 크기는 다음과 같다.

$$x_w \in \mathbb{R}^d$$

- Word2Vec : CBOW model
 - 주변 단어들을 가지고 단어를 구한다.



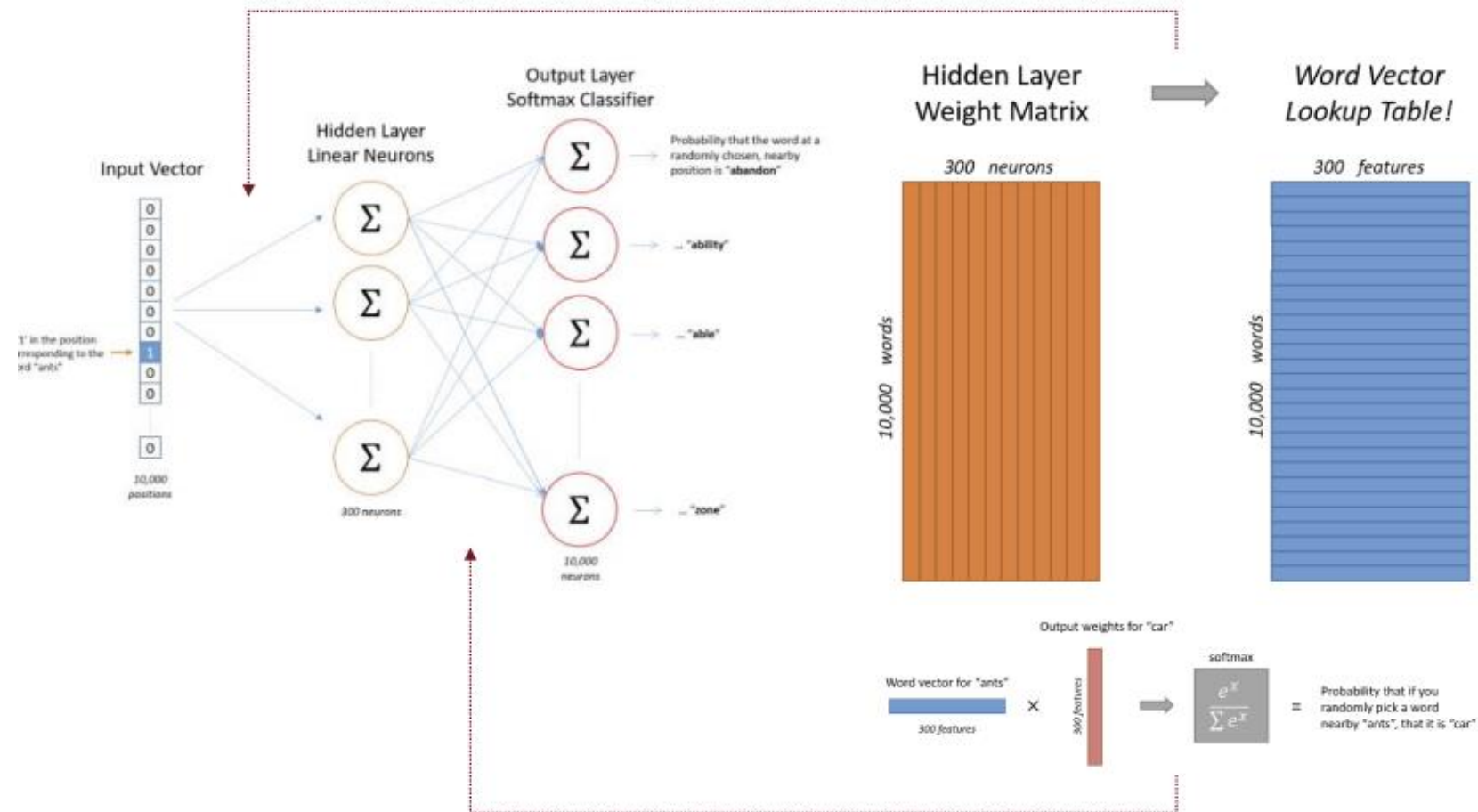
- 주어진 단어에 대한 주변 단어의 모델 확률은 밑과 같다.
 - ✓ 대문자C는 단어의 조합을 의미한다.

feature for context \mathcal{C} : $h_{\mathcal{C}}$
 classifier for word w : v_w

$$p(w|\mathcal{C}) = \frac{e^{h_{\mathcal{C}}^T v_w}}{\sum_{k=1}^K e^{h_{\mathcal{C}}^T v_k}}$$

$$h_{\mathcal{C}} = \sum_{c \in \mathcal{C}} x_c$$

- Word2Vec
 - 밑은 다른 방식으로 표현한 것을 나타낸 것이다.



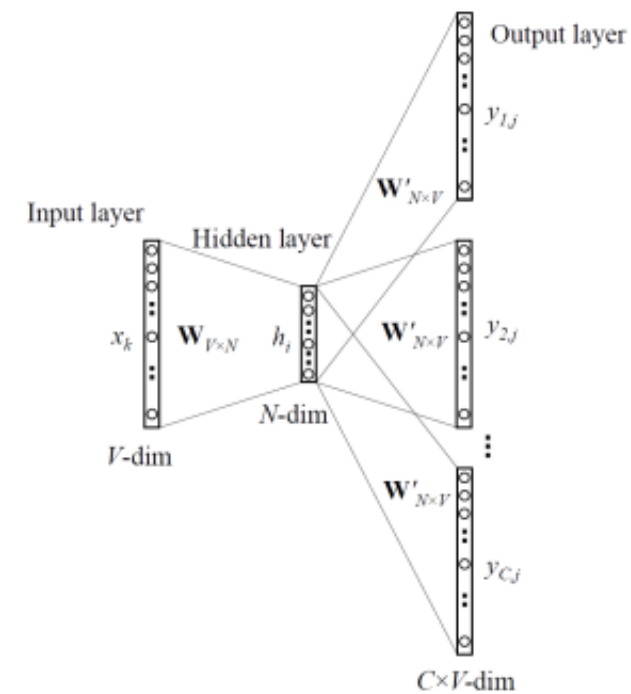
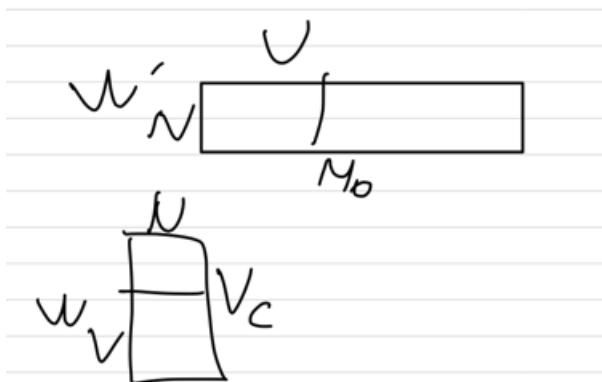
Word2Vec

• Word2Vec

- 간단히 하기위해 $p(w_{t+j}|w_t)$ 대신
 - ✓ O = output, c = center,
 - ✓ W_{t+j} = context , W_t = center

$$p(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w=1}^W \exp(u_w^T v_c)} \quad \text{을 넣어본다.}$$

- 모든 단어는 두개의 벡터를 가진다.
 - ✓ V는 W의 행, U는 W프라임의 열이다.
 - ❖ W프라임과 W는 다르지만 편의상 같다고 가정한다.



- Word2Vec
 - 경사상승법으로 파라미터를 학습한다.

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

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

$$\begin{aligned} \frac{\partial}{\partial v_c} \log p(o|c) &= \frac{\partial}{\partial v_c} \log \frac{\exp(u_o^T v_c)}{\sum_{w=1}^W \exp(u_w^T v_c)} \\ &= \underbrace{\frac{\partial}{\partial v_c} u_o^T v_c}_A - \underbrace{\frac{\partial}{\partial v_c} \log \sum_{w=1}^W \exp(u_w^T v_c)}_B \end{aligned}$$

- Word2Vec
 - 경사상승법으로 파라미터를 학습한다.

✓ For chunk A

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

✓ For chunk B

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

- Word2Vec
 - 경사상승법으로 파라미터를 학습한다.

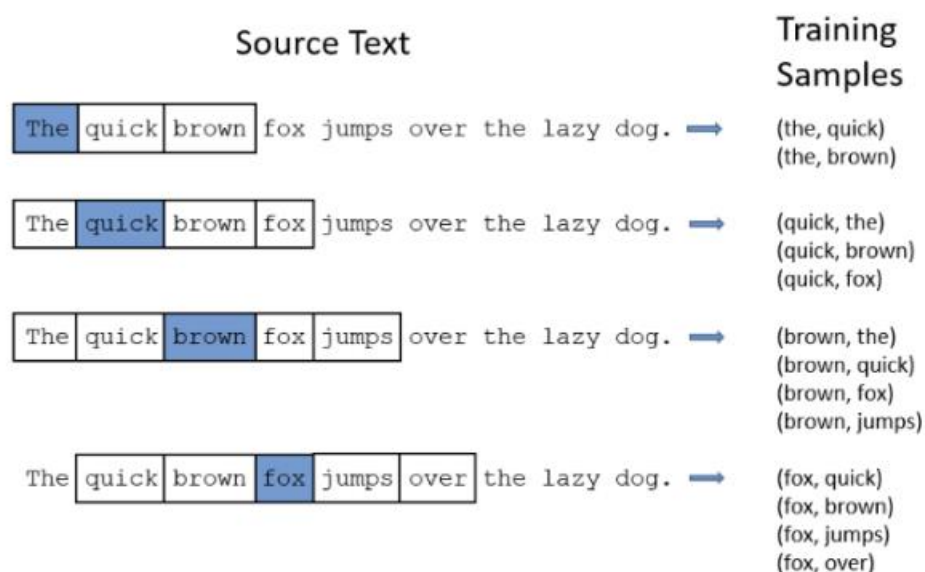
$$\frac{\partial}{\partial v_c} \log p(o|c) = u_o - \sum_{w=1}^W P(w|c) \cdot u_w$$

- 가중치 벡터를 업데이트 하는 식이다.

$$v_c(t+1) = v_c(t) + \alpha \left(u_o - \sum_{w=1}^W P(w|c) \cdot u_w \right)$$

- Word2Vec

- 꼭 한번에 여러 단어를 구할 필요 없이 하나씩 구해도 된다.
 ✓ Gradient는 각각 계산하나, 하나씩 하나 같이 때문이다.



- Word2Vec
 - 가중치가 크다.(2XVXN)
 - ✓ 자주 사용되는 두개의 단어나 구는 하나의 단어로 취급한다.
 - ❖ Ex) machine learning..
 - ✓ 빈번 출현 단어에 대해 Subsampling을 한다.
 - ❖ 코퍼스 등장 확률이 높을수록 지울 확률이 증가 하는 것을 볼 수 있다.

$$P(w_i) = 1 - \sqrt{\frac{t}{f(w_i)}} \quad \begin{array}{l} \longleftarrow \text{Threshold } (10^{-5}) \\ \longleftarrow \text{Term frequency in the corpus} \end{array}$$

$$\text{if } f(w_i) = 10^{-4}, P(w_i) = 1 - \sqrt{\frac{1}{10}} = 0.6838$$

$$\text{if } f(w_i) = 10^{-2}, P(w_i) = 1 - \sqrt{\frac{1}{1000}} = 0.9684$$

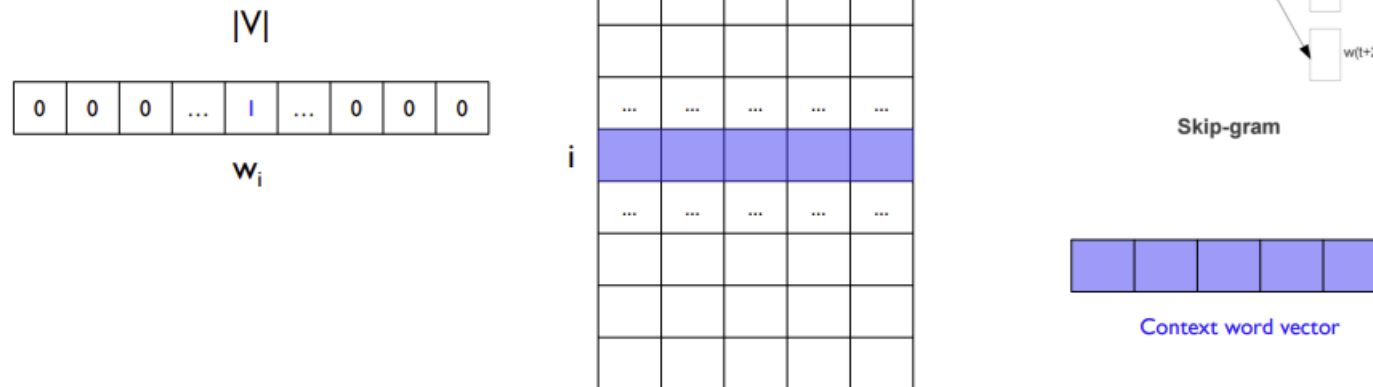
- Word2Vec
 - 가중치가 크다.(2XVXN)
 - ✓ Negative sampling을 활용한다.
 - ❖ 가중치 업데이트에 모든 단어를 사용하지 않고 몇 개의 예만 사용한다.

$$J(\theta) = \frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log p(w_{t+j} | w_t) = \frac{1}{T} \sum_{t=1}^T J_t(\theta)$$

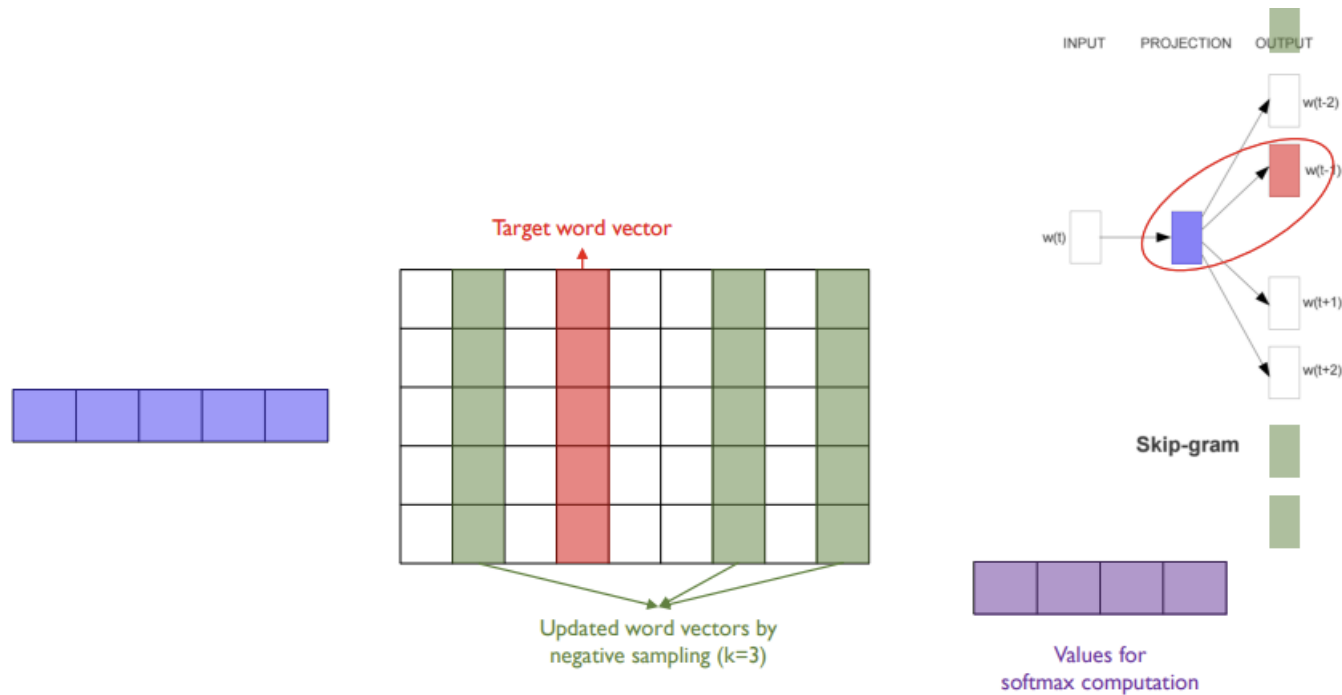
$$J_t(\theta) = \log \sigma(u_o^T v_c) + \sum_{i=1}^k E_{i \sim P(w)} [\log \sigma(-u_i^T v_c)]$$

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

- Word2Vec
 - 가중치가 크다.(2XVXN)
 - ✓ Negative sampling을 활용한다.
 - ❖ 가중치 업데이트에 모든 단어를 사용하지 않고 몇 개의 예만 사용한다.



- Word2Vec
 - 가중치가 크다.($2 \times V \times N$)
 - ✓ Negative sampling을 활용한다.
 - ❖ 가중치 업데이트에 모든 단어를 사용하지 않고 몇 개의 예만 사용한다.



- Word2Vec

