

딥러닝 세미나 Season #6

Word2Vec
2013, Google Inc.

한양대학교
컴퓨터 소프트웨어학과
인공지능 연구실
조건희

Word2Vec

Word2Vec

Word 를 Vector로 만든다!

Word(고차원 공간) \rightarrow Vector(저차원 공간)
투영해준다는 의미

Word(고차원 공간) \rightarrow Vector(저차원 공간)

투영해준다는 의미

Word Embedding이라고 표현하기도 함.

간단히 생각해보자

[아메리카노 카페라떼 카푸치노 마끼아또]

[1	0	0	0]
[아메리카노	카페라떼	카푸치노	마끼아또]

[0	1	0	0]
[아메리카노	카페라떼	카푸치노	마끼아또]

[0	0	1	0]
[아메리카노	카페라떼	카푸치노	마끼아또]

[0	0	0	1]
[아메리카노	카페라떼	카푸치노	마끼아또]

[0 0 0 1]
[아메리카노 카페라떼 카푸치노 마끼아또]

One-Hot Encoding

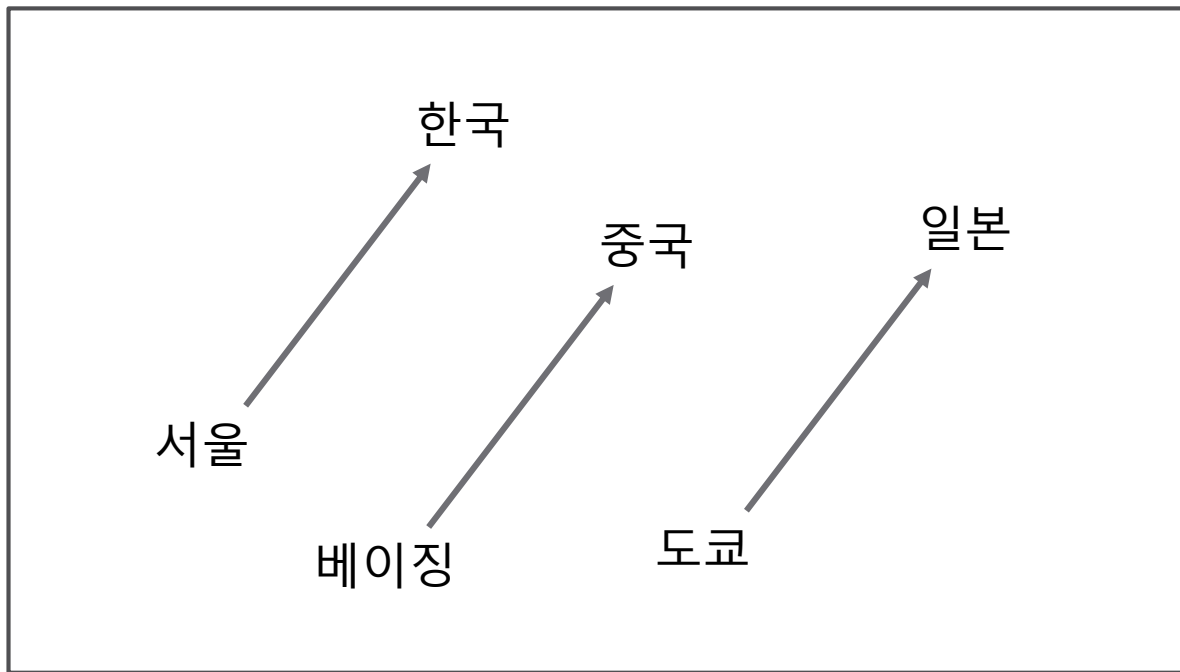
이렇게 간단히 해결될 문제가 아님

One-hot Encoding의 문제점

One-hot Encoding의 문제점
Embedding Vector에
단어의 의미가 드러나지 않음

단어의 의미를 벡터화하면 뭐가 좋을까?

단어의 의미를 벡터화하면 뭐가 좋을까?
벡터 좌표가 단어의 의미를 담고 있으므로
벡터 연산을 통해 추론을 할 수 있음



한국 - 서울 + 도쿄 = ?

Word2Vec Models

Word2Vec Models

CBOW & Skip-gram

CBOW

Context(주변 단어) → Target(목표 단어) 예측

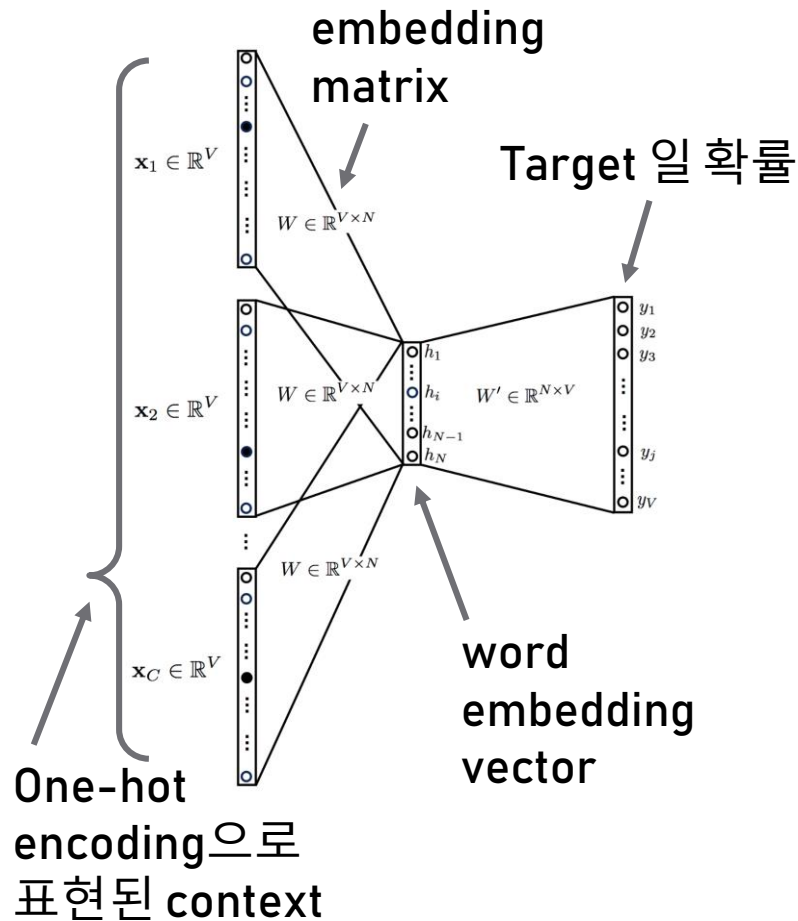


Skip-gram

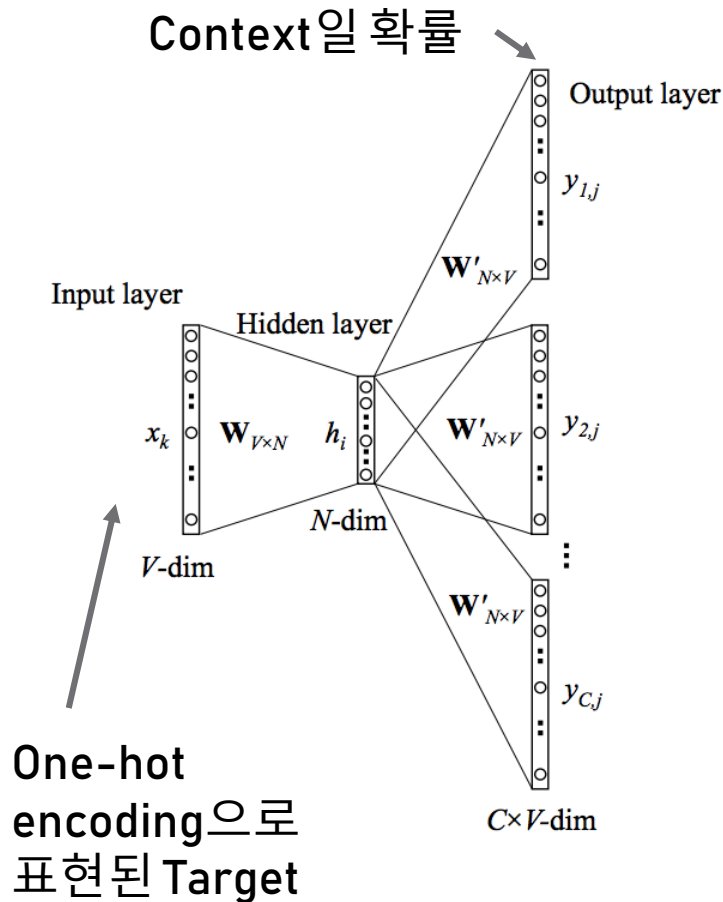
Target(목표 단어) → Context(주변 단어) 예측



CBOW



Skip-gram



Skip-gram

x : One-hot encoded input word (V 차원)

V : 단어사전 크기

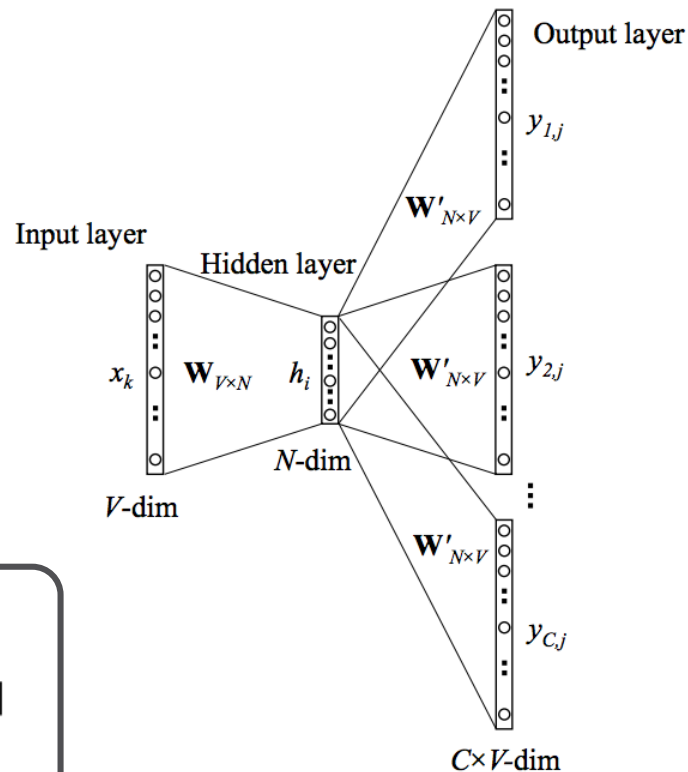
W : 임베딩 매트릭스($V \times N$)

N : 임베딩 사이즈

h : word embedding 벡터 [$h = x * W$] (N 차원)

W' : $N \times V$ 매트릭스

y : output [$y = \text{softmax}(h * W')$] (V 차원)



$$\begin{array}{c}
 [0 \quad 0 \quad 0 \quad 1 \quad 0] \\
 \mathbf{x}
 \end{array}
 \times
 \begin{array}{ccc}
 17 & 24 & 1 \\
 23 & 5 & 7 \\
 4 & 6 & 13 \\
 10 & 12 & 19 \\
 11 & 18 & 25
 \end{array}
 = [10 \quad 12 \quad 19]$$

\mathbf{W} \mathbf{h}

Skip-gram

말뭉치

Window 크기: 2
학습 데이터 샘플
(input, output)

The quick brown fox jumped over the lazy dog. →

(the, quick)
(the, brown)

The quick brown fox jumped over the lazy dog. →

(quick, the)
(quick, brown)
(quick, fox)

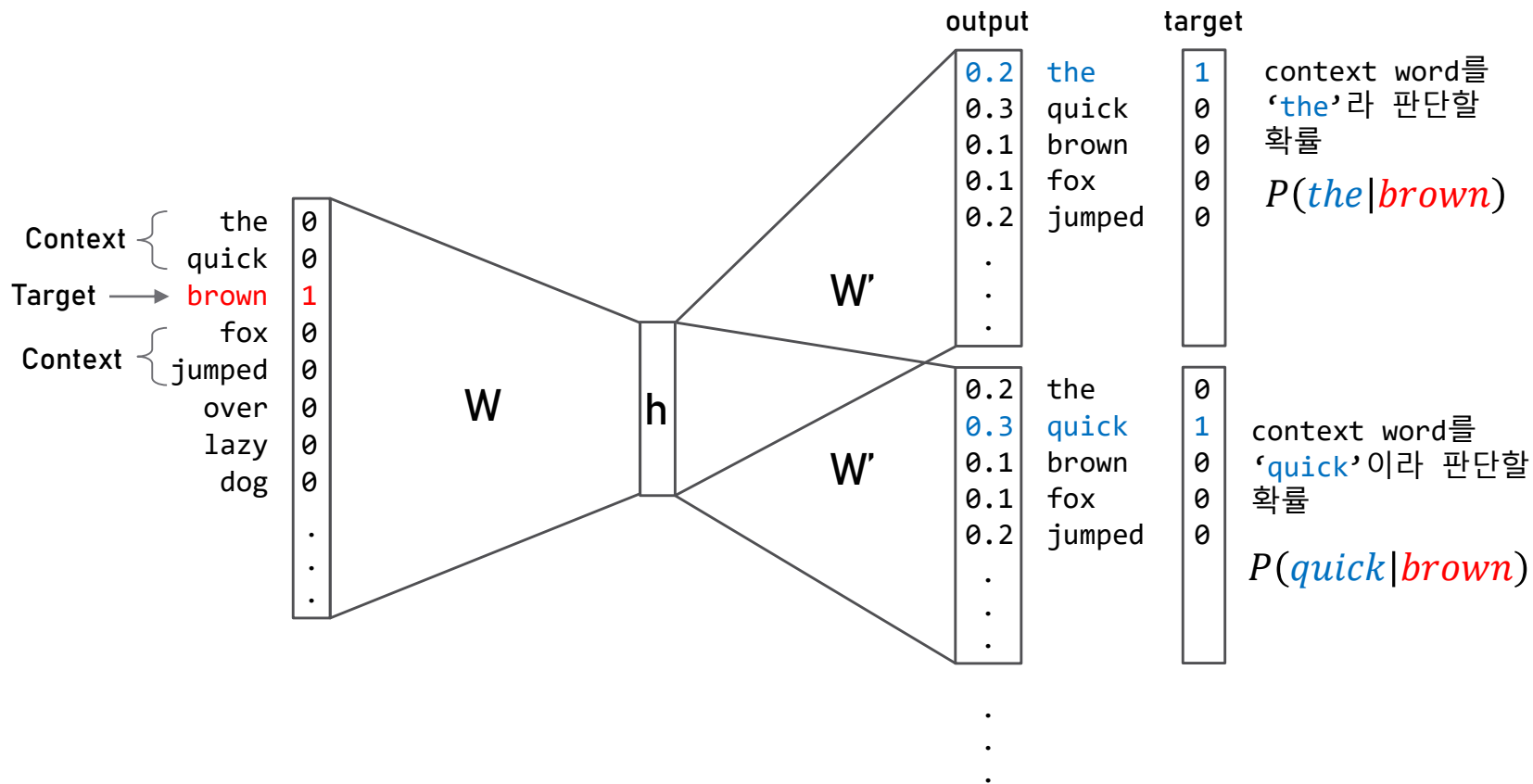
The quick brown fox jumped over the lazy dog. →

(brown, the)
(brown, quick)
(brown, fox)
(brown, jumped)

말뭉치 내에 존재하는 모든 단어를 슬라이딩하며 학습!

Skip-gram

The quick brown fox jumped over the lazy dog.



Skip-gram

The quick brown fox jumped over the lazy dog.

$$P(\text{quick}|\text{brown}) = \frac{\exp(u_{\text{quick}}^T v_{\text{brown}})}{\sum_{w=1}^W \exp(u_w^T v_{\text{brown}})}$$

target 단어를 입력으로 받았을 때
이 단어가 target의 context 단어일 확률

실제 context일 경우 → 최대화
실제 context가 아닐 경우 → 최소화

Word2Vec의 학습기법

Word2Vec의 학습기법

Subsampling frequent words

Negative sampling

Subsampling frequent words

보통 말뭉치의 차원수는 10만개(V) 안팎
임베딩 차원수가 100차원(N)이라고 해도

$$W \rightarrow 10\text{만} \times 100 \text{ (V} \times \text{N)}$$

$$W' \rightarrow 100 \times 10\text{만} \text{ (N} \times \text{V)}$$

총 2000만개의 파라미터

연산량이 너무 많음!

Subsampling frequent words

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

$P(w_i)$: w_i (i 번째 단어)를 학습에서 제외할 확률

$f(w_i)$: w_i (i 번째 단어)가 말뭉치에 등장한 비율

t : *threshold*

Negative sampling

softmax를 계산할 때
전체 단어에 대하여 계산하면
연산량이 너무 많아짐.

Negative sampling

softmax를 계산할 때
전체 단어에 대하여 계산하면
연산량이 너무 많아짐.

일부만 sampling해서 softmax를 계산해보자!

Negative sampling

The quick brown fox jumped over the lazy dog.

Positive sample

(the, quick)
(the, brown)
(quick, the)
(quick, brown)
(quick, fox)
(brown, the)
(brown, quick)
(brown, fox)
(brown, jumped)
(...)

Negative sample

(the, fox)
(quick, dog)
(quick, lazy)
(brown, over)
(...)

Negative sampling

윈도우 사이즈가 5일 경우

Positive sample 5개와

Negative sample 20개만 가지고 softmax 계산

윈도우에 등장하지 않은 단어가 negative sample로 등장할 확률

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

감사합니다