



Word2Vec

23.06.01 유하영

Word2Vec

‘비슷한 문맥에서 등장하는 단어들은 비슷한 의미를 가진다’

유사성 있는 단어를 예측

문장 생성은 X -> 다른 모델 필요

원수는 외나무다리에서 만난다
(연어)

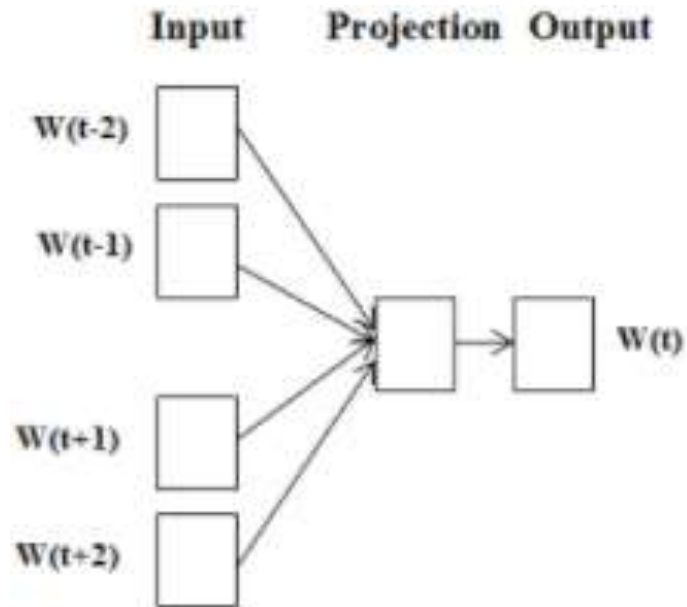
man , woman

word2vec

CBOW

주변단어 \rightarrow 중심단어 예측

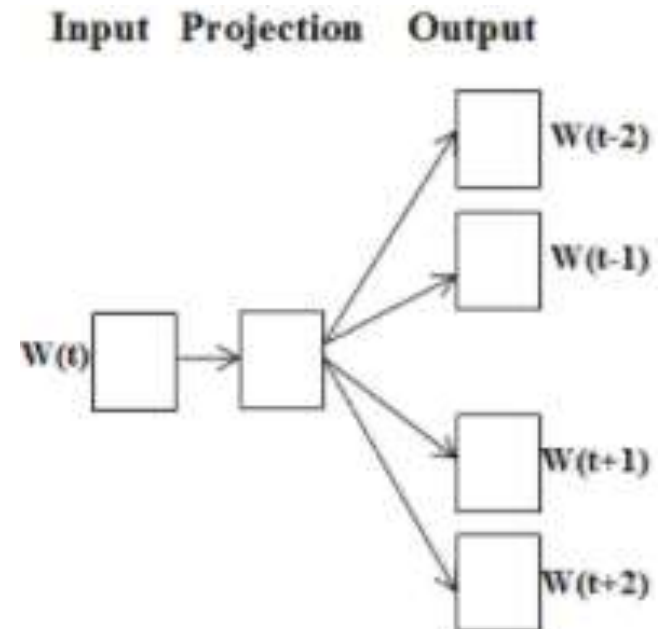
you [?] goodbye and I say hello.



Skip-gram

중심단어 \rightarrow 주변단어 예측

[?] say [?] and I say hello.



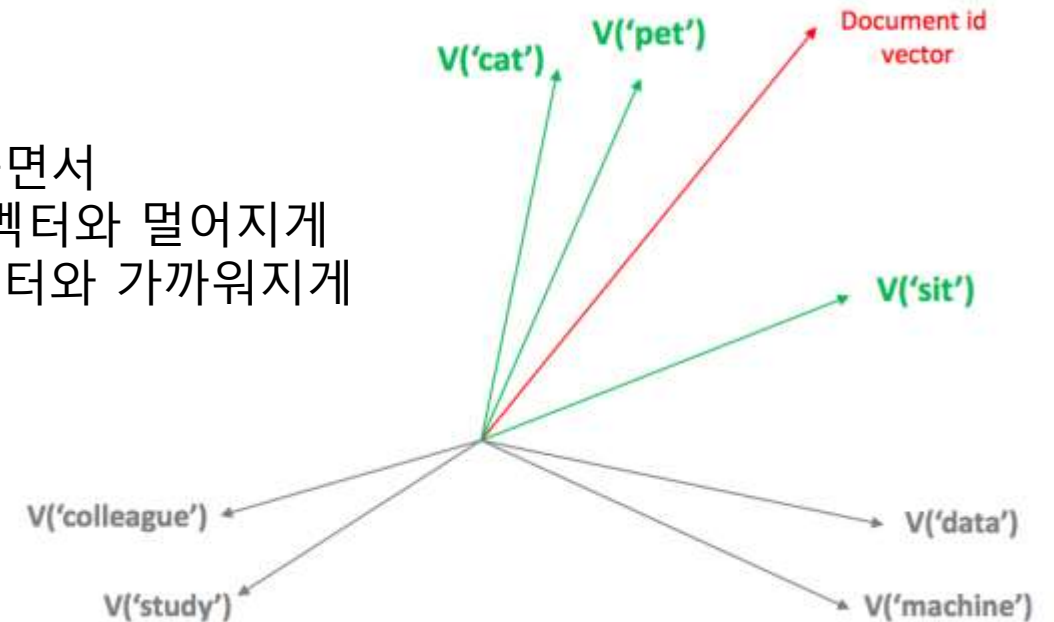
word2vec

CBOW



코퍼스를 슬라이딩 하면서
등장 ↓ 단어의 벡터가 중심단어 벡터와 멀어지게
등장 ↑ 단어의 벡터가 중심단어 벡터와 가까워지게

Skip-gram



학습데이터

The fat cat sat on the mat

중심 단어 주변 단어
↓ ↓
The fat cat sat on the mat

The fat cat sat on the mat

The fat cat sat on the mat

The fat cat sat on the mat

The fat cat sat on the mat

The fat cat sat on the mat

The fat cat sat on the mat

(window size = 2)

중심 단어	주변 단어
[1, 0, 0, 0, 0, 0, 0]	[0, 1, 0, 0, 0, 0, 0], [0, 0, 1, 0, 0, 0, 0]
[0, 1, 0, 0, 0, 0, 0]	[1, 0, 0, 0, 0, 0, 0], [0, 0, 1, 0, 0, 0, 0], [0, 0, 0, 1, 0, 0, 0]
[0, 0, 1, 0, 0, 0, 0]	[1, 0, 0, 0, 0, 0, 0], [0, 1, 0, 0, 0, 0, 0], [0, 0, 0, 1, 0, 0, 0], [0, 0, 0, 0, 1, 0, 0]
[0, 0, 0, 1, 0, 0, 0]	[0, 1, 0, 0, 0, 0, 0], [0, 0, 1, 0, 0, 0, 0], [0, 0, 0, 0, 1, 0, 0], [0, 0, 0, 0, 0, 1, 0]
[0, 0, 0, 0, 1, 0, 0]	[0, 0, 1, 0, 0, 0, 0], [0, 0, 0, 1, 0, 0, 0], [0, 0, 0, 0, 0, 1, 0], [0, 0, 0, 0, 0, 0, 1]
[0, 0, 0, 0, 0, 1, 0]	[0, 0, 0, 1, 0, 0, 0], [0, 0, 0, 0, 1, 0, 0], [0, 0, 0, 0, 0, 0, 1]
[0, 0, 0, 0, 0, 0, 1]	[0, 0, 0, 0, 1, 0, 0], [0, 0, 0, 0, 0, 1, 0]

학습데이터

The fat cat sat on the mat

```
[('fat', 0.7668752074241638), ('mat', 0.6178626418113708), ('the', 0.10980246961116791), ('on', -0.1959061473608017), ('cat', -0.3469330966472626)]
```



대량의 문장으로 이루어진 문서를
전처리하여 토큰화 한 후
모든 단어에 대한 학습



word2vec

```
model.similarity('this', 'is')
```

```
0.407970363878
```

```
model.similarity('post',
```

```
0.057204389197
```

```
7
```

```
'book'
```

```
[ 0.11279297 -0.02612305 -0.04492188  
0.06982422 0.140625 0.03039551 -0.04370117  
0.24511719 0.08740234 -0.05053711 0.23144531 -  
0.07470703 ...
```

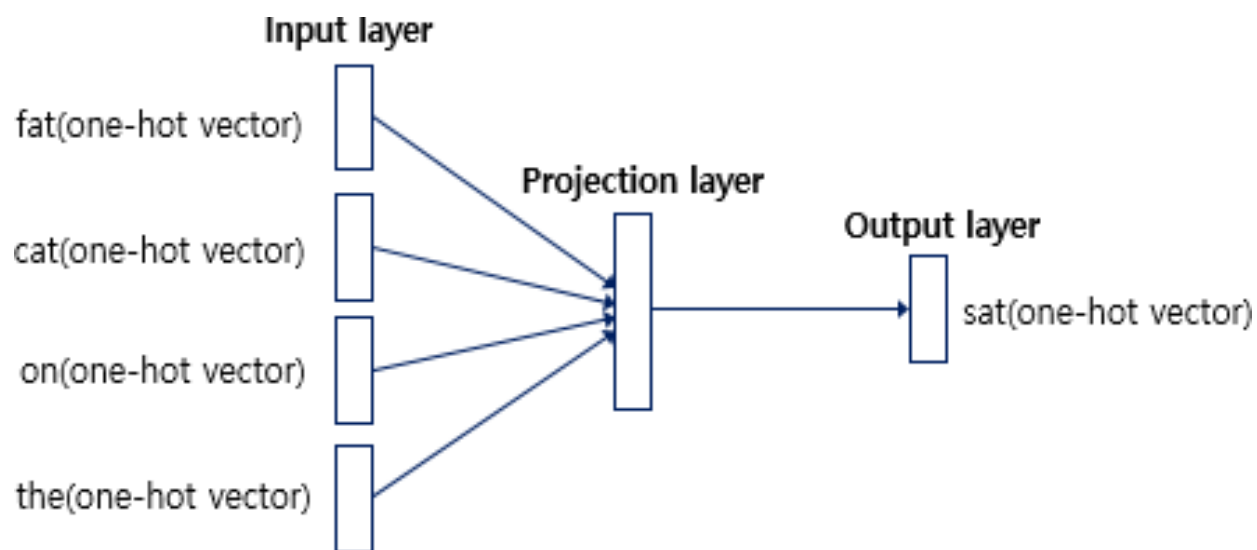
```
중략 ...
```

```
0.03637695 -0.16796875 -0.01483154 0.09667969  
-0.05761719 -0.00515747]
```

CBOW

주변단어 -> 중심단어 예측

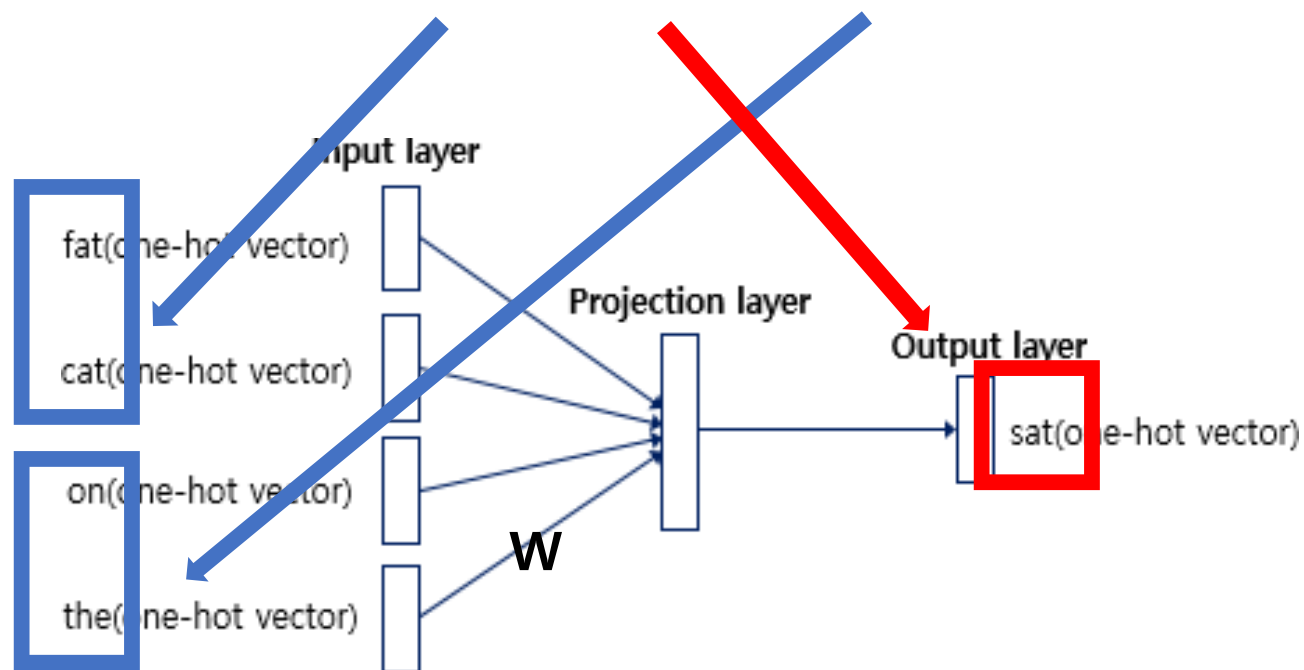
“ The fat cat sat on the mat ”



CBOW

주변단어 -> 중심단어 예측

“ The fat cat sat on the mat ”



예문)

The fat cat sat on the mat.

\ 주변단어 /

중심단어
(예측할 단어)

\ 주변단어 /

중심단어	주변단어
Sat	fat
	cat
	on
	the

→ 원-핫 벡터로 변환



중심단어
(예측할 단어)

sat
[0,0,0,1,0,0,0]

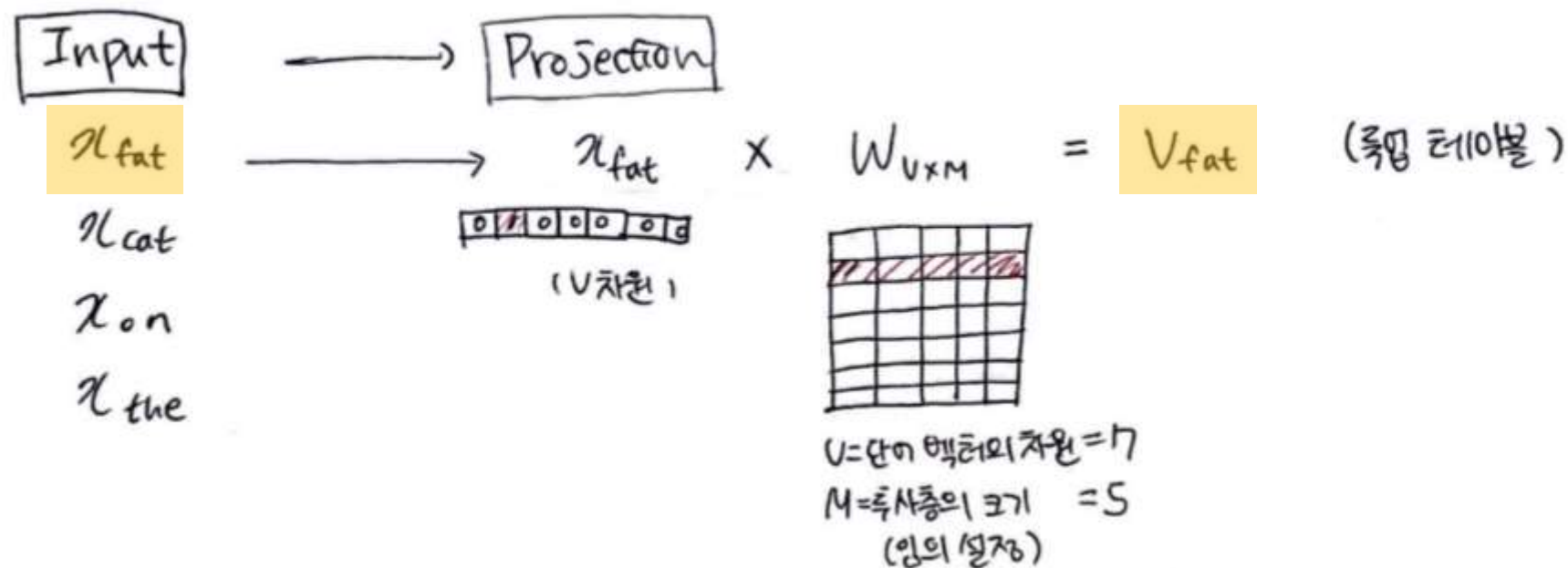
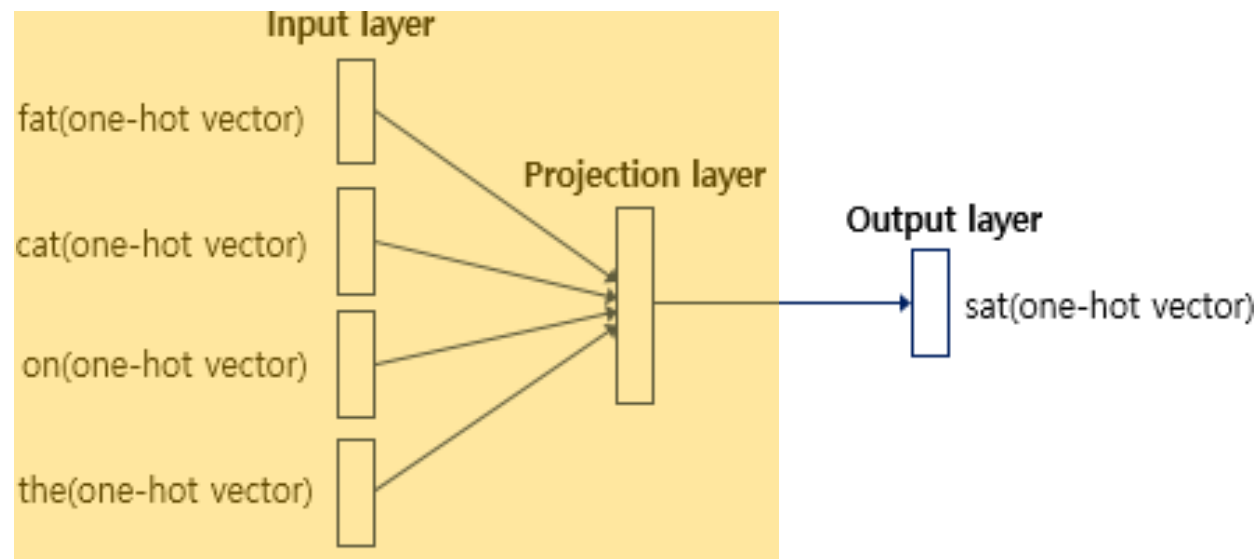
주변단어

fat [0,1,0,0,0,0,0]

cat [0,0,1,0,0,0,0]

on [0,0,0,0,1,0,0]

the [0,0,0,0,0,1,0]



Input

x_{fat}

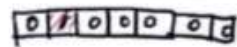
x_{cat}

x_{on}

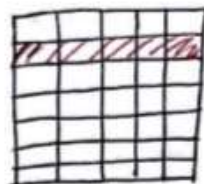
x_{the}

Projection

$$x_{fat} \times W_{U \times M} = V_{fat} \quad (\text{특정 레이어별})$$



(U차원)



U=단어 벡터의 차원=7

M=특성들의 크기=5

(임의의 설정)

$$x_{fat} \rightarrow x_{fat} \times W_{U \times M} = V_{fat}$$



$$x_{cat} \rightarrow x_{cat} \times W_{U \times M} = V_{cat}$$



$$x_{on} \rightarrow x_{on} \times W_{U \times M} = V_{on}$$



$$x_{the} \rightarrow x_{the} \times W_{U \times M} = V_{the}$$



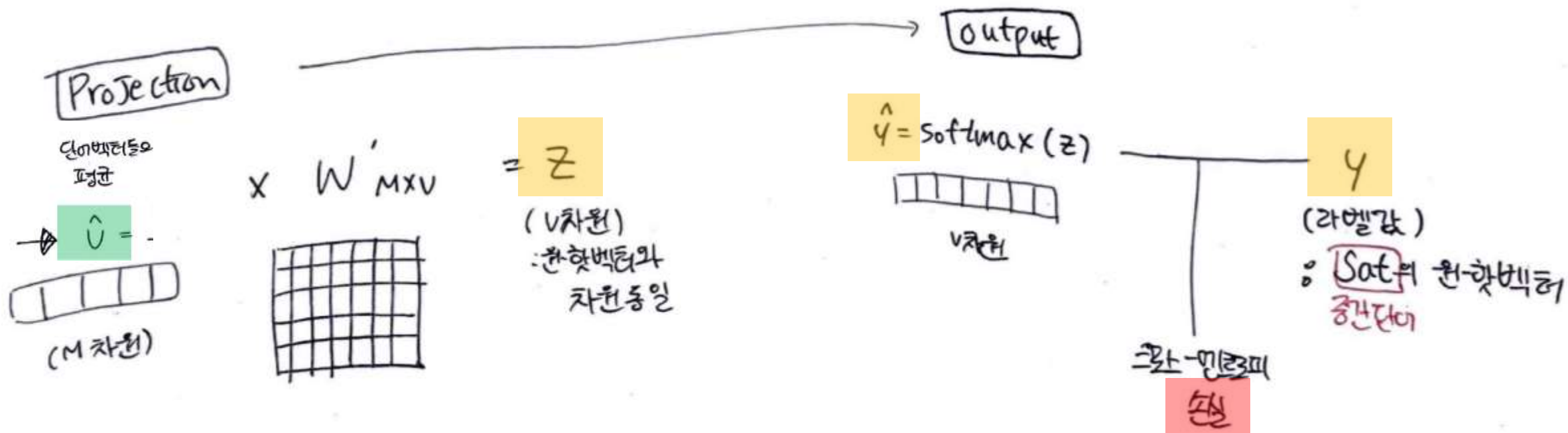
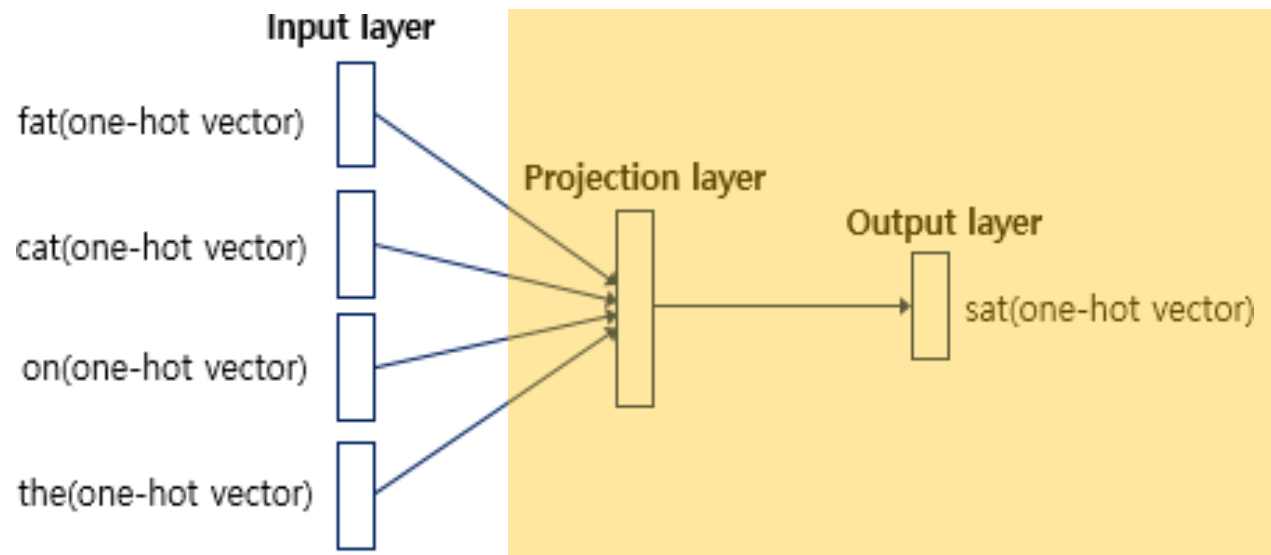
(M차원)

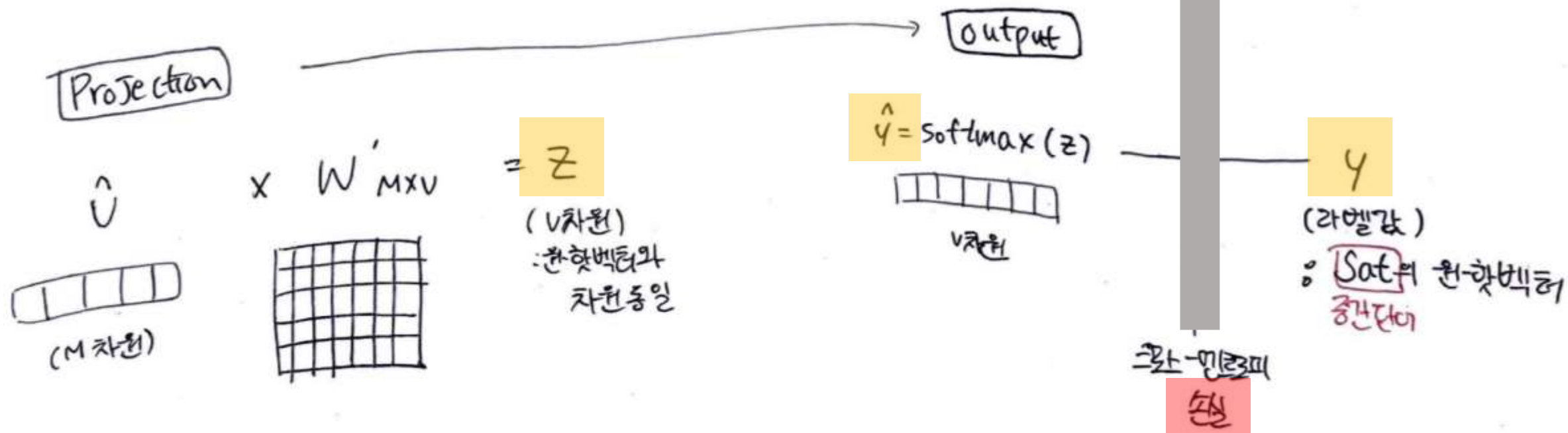
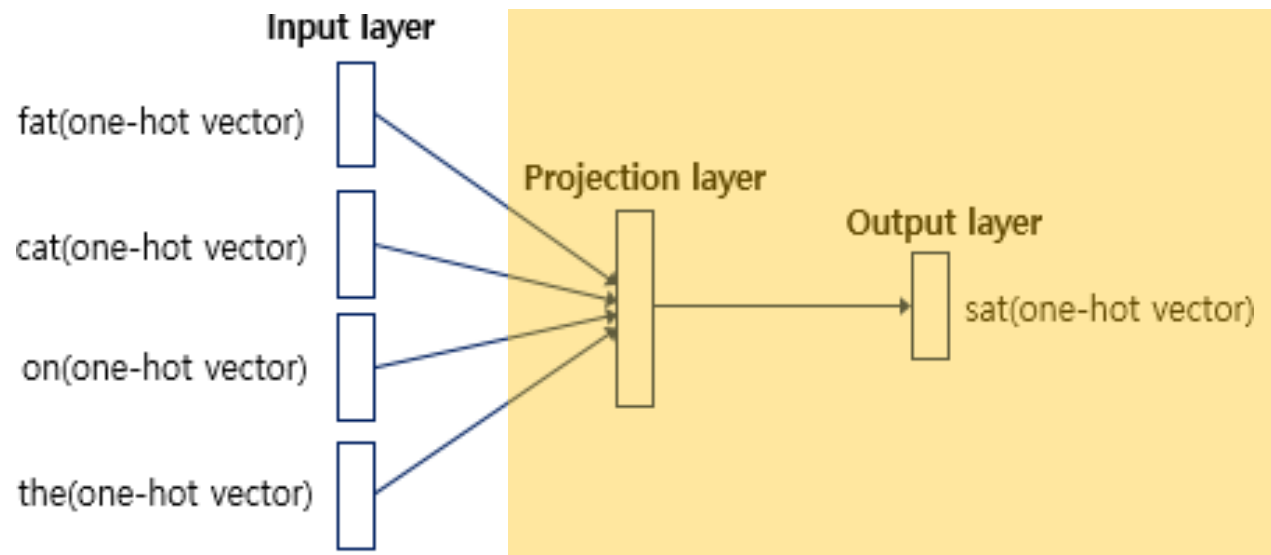
단어 벡터들의

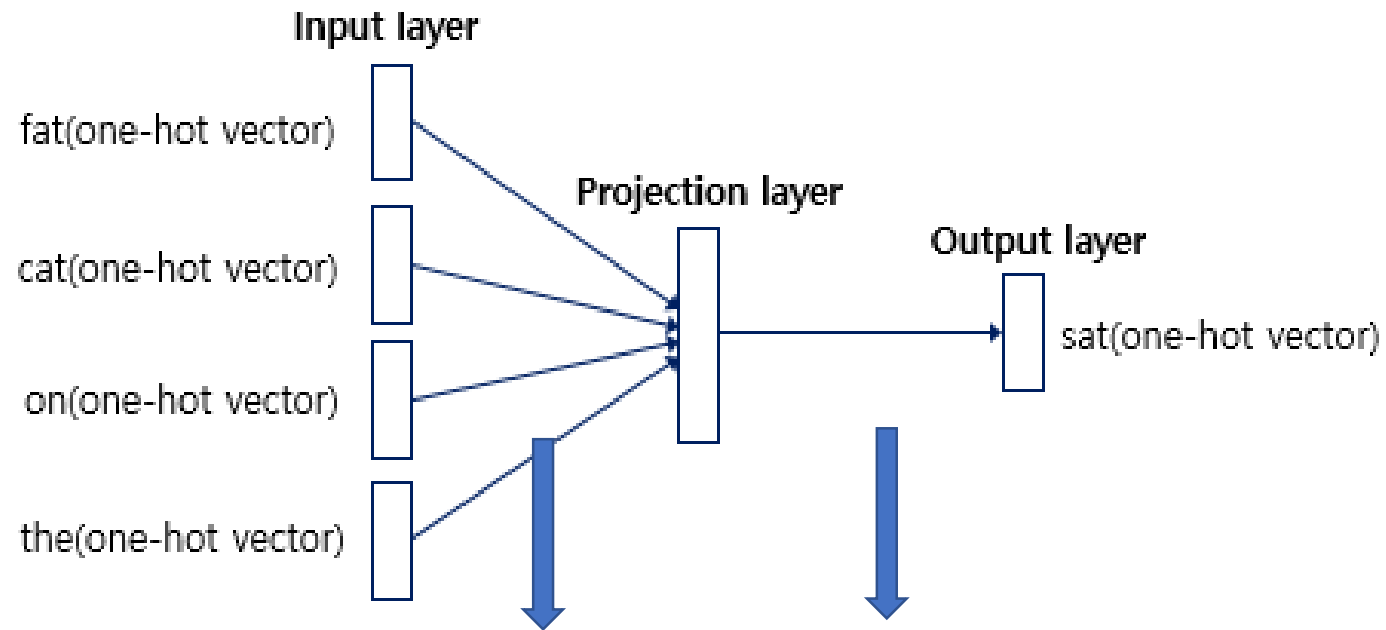
평균

$$\hat{U} = \frac{V_{fat} + V_{cat} + V_{on} + V_{the}}{2n} :$$

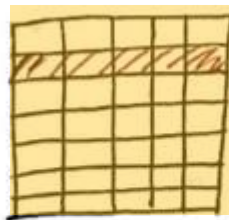
(M차원의 임베딩 벡터)





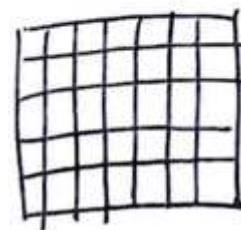


$W_{U \times M}$



U =단어 벡터의 차원=7
 M =특성들의 크기=5
 (임의의 설정)

$W'_{M \times U}$



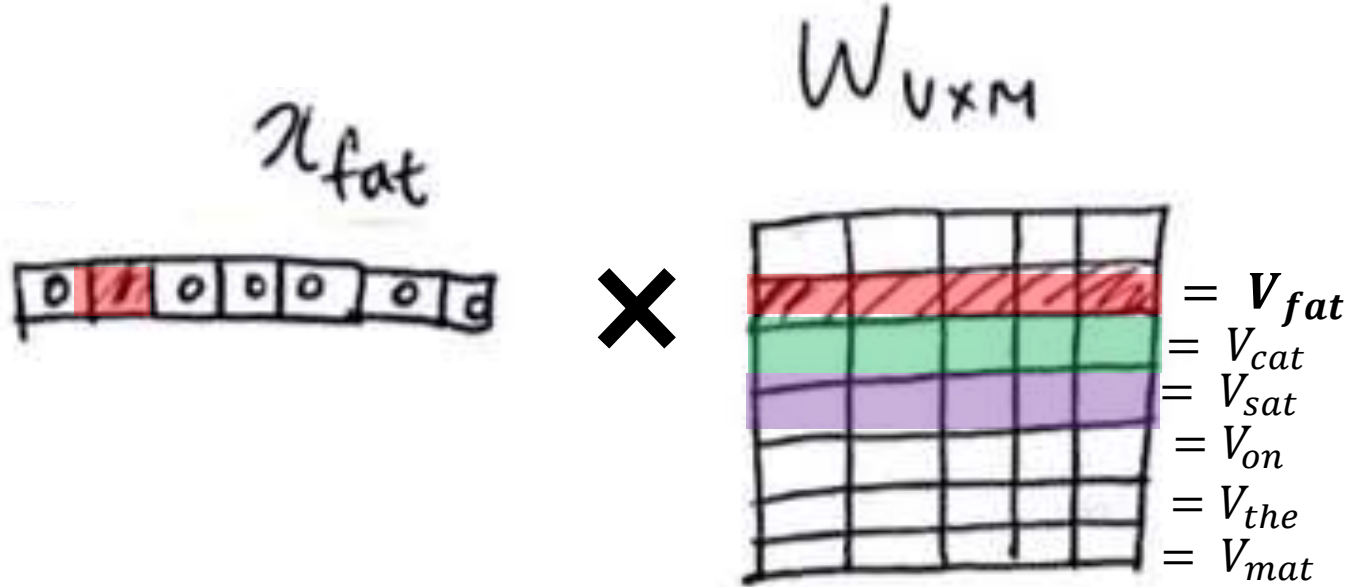
가중치 행렬 w 와 w'

- 차원은 전치한 것과 같다.
- 완전히 동일한 행렬은 아니다.
 (때에 따라서 두개를 하나의 행렬로 취급해서 학습할 수도 있다.)

- 중심단어를 맞추기 위해 W 와 W' 를 조금씩 업데이트 하며 학습하는 것
- W 는 가중치 행렬이자 word2vec의 최종 결과물인 임베딩 단어벡터의 모음이다.

가중치 선택

lookup table

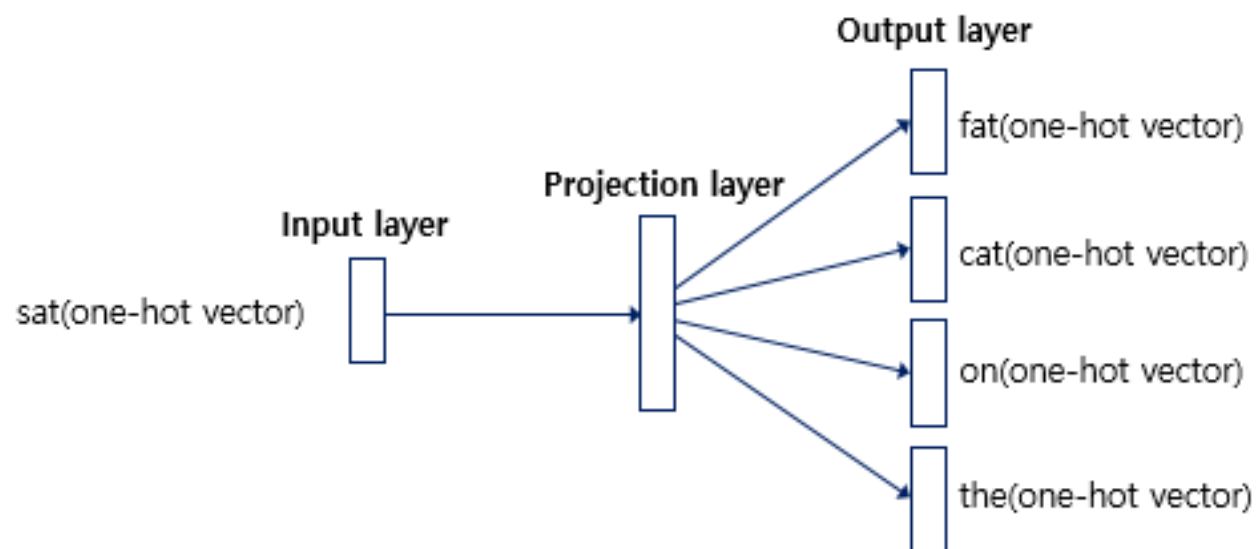


- A. 가중치행렬 **W**를 최종단어의 분산 표현으로 이용한다.
- B. 가중치행렬 **W'**를 최종단어의 분산 표현으로 이용한다.
- C. 양쪽 가중치 모두 최종단어의 분산표현으로 이용한다.

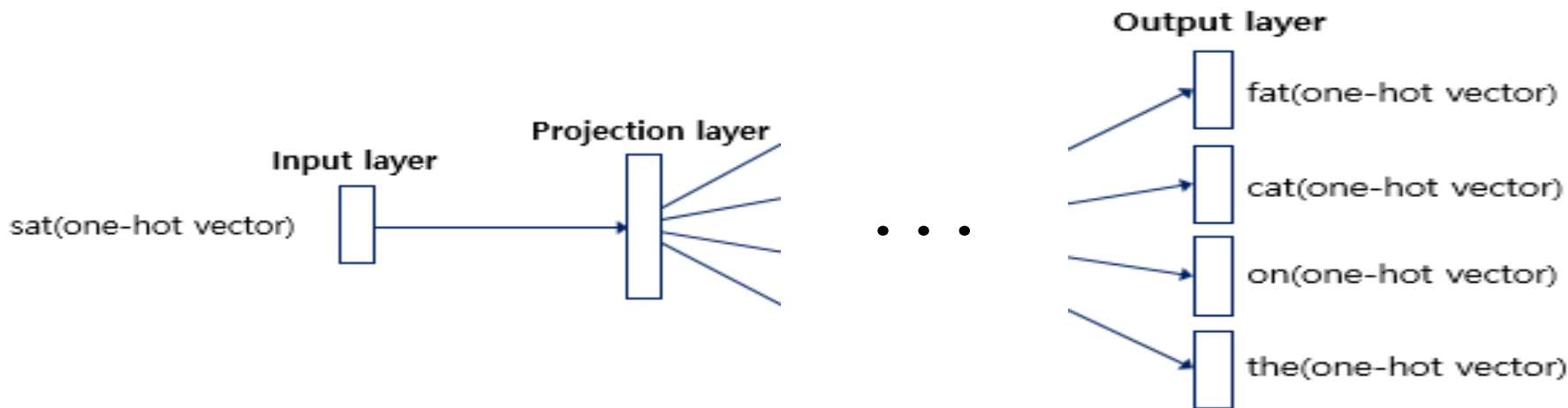
Skip-gram

중심단어 -> 주변단어 예측

The fat cat sat on the mat



중심단어	주변단어
sat	fat
sat	cat
sat	on
sat	the



output

Input

projection

중심단어.
(one-hot vector)
"sat"
= x_{sat}

중심단어

$$x_{sat} \times W_{V \times M} = V_{sat} \text{ (M차원)}$$

이단어
주변단어

① $V_{sat} \times W'_{M \times V} = Z \rightarrow \text{softmax}(Z) = \hat{y}$

② $V_{sat} \times W'_{M \times V} = Z \rightarrow \text{softmax}(Z) = \hat{y}$

③

④

= 총 4번의 가중치 행렬 update (back-propagation).

Handwritten notes and labels include: "output (보여주기)", "실제값 (y)", "예측값 (y-hat)", "2차", "update", "up-", and "Sat의 임베딩 행렬 update".

you got me looking for attention.

CBOW

주변단어

중심단어

학습횟수

got → you = 1

you
me } \bigcirc got = 1

got
looking } me = 1

me
for } looking = 1

looking
attention } for = 1

for — attention. = 1

Skip-gram

중심단어

주변단어

학습횟수

you — got = 1

got — { you → up-
me → up- = 2

me — = 2

looking = 2

for = 2

attention. = 1

SGNS

skip-gram negative sampling

skip-gram

전체 단어 집합

SGNS

일부 단어 집합
에 집중

다중 분류
류



**negative
sampling**

이중분

I love you

I, you 주어졌을때
가운데 나올단어가 무엇인가?

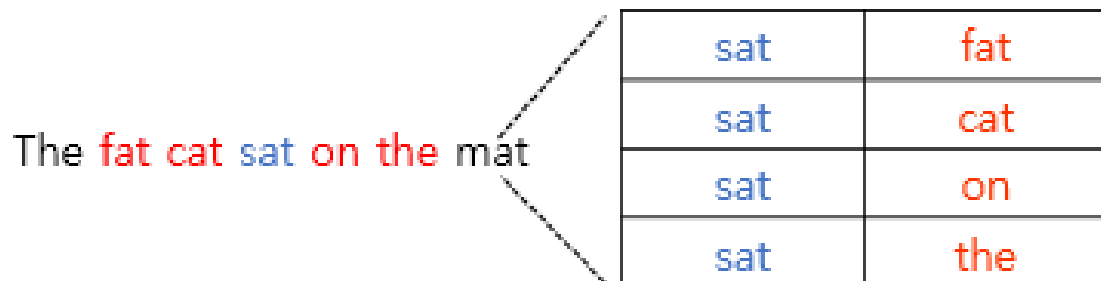
I, you란 단어가 주어졌을때
가운데 나올단어가 love인가?

positive

실제로 관측된 단어쌍의 점수 최대화

negative

무작위로 선택된 부정샘플의 집합점수 최소화



skip-gram

입력과 레이블의 변화

입력1	입력2	레이블
cat	The	1
cat	fat	1
cat	sat	1
cat	on	1

입력1	입력2	레이블
cat	The	1
cat	fat	1
cat	pizza	0
cat	computer	0
cat	sat	1
cat	on	1

단어 집합에서 랜덤으로
선택된 단어들을
레이블 0의 샘플로 추가.

negative sampling