

Background

- 원-핫 벡터의 한계
 - 단어 벡터 간 유의미한 유사도를 계산할 수 없음
- 단어 벡터 간 유사도를 반영할 수 있도록 단어 의미를 수치화할 수 있는 방법 필요 => 임베딩
- 희소표현은 2차원에 각 차원이 분리된 표현방법이고, 분산표현은 저차원에 단어의 의미를 여러 차원에다 분산하는 표현방법 => 분산표현을 사용해서

• W2V 으로 무엇을 할 수 있을까?

(한국 - 서울 + 도쿄 = 일본
 박찬호 - 야구 + 축구 = 손나래)

단어 벡터 간 유의미한 유사도 계산
 ↳ 대표적 방법이 word2vec

→ 위와 같은 연산이 가능한 이유는 각 단어 벡터가 단어 벡터 간 유사도를 반영한 값을 가지고 있기 때문

W2V 학습방식

• CBOW

→ 주변에 있는 단어들을 입력으로 중간에 있는 단어들을 예측하는 방법

→ 메커니즘

" The (fat cat sat on the) mat "

주변 단어 중심 단어

→ 윈도우 크기 = n 실제 주변 단어와 개수 = 2n
 → 중심 단어를 예측하기 위해 앞, 뒤로 몇개와 단어를 볼지에 대한 범위

• 슬라이딩 윈도우 (sliding window)

중심 단어 주변 단어

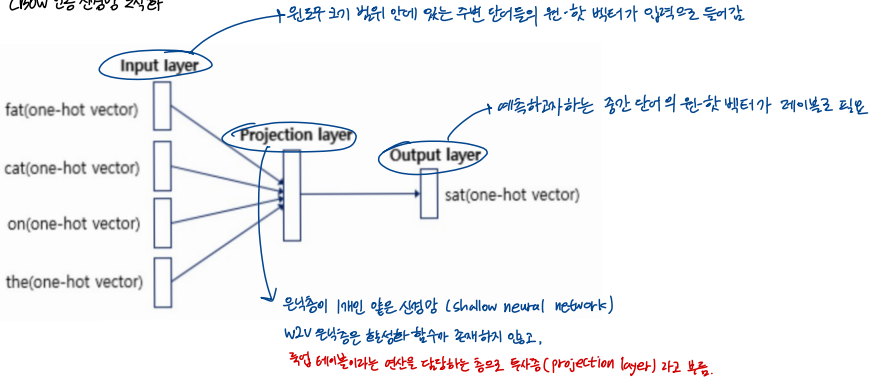
↓ ↓

중심 단어	주변 단어
The fat cat sat on the mat	[1, 0, 0, 0, 0, 0, 0, 0] [0, 1, 0, 0, 0, 0, 0, 0] [0, 0, 1, 0, 0, 0, 0, 0]
The fat cat sat on the mat	[0, 1, 0, 0, 0, 0, 0, 0] [0, 0, 1, 0, 0, 0, 0, 0] [0, 0, 0, 1, 0, 0, 0, 0]
The fat cat sat on the mat	[0, 0, 1, 0, 0, 0, 0, 0] [0, 1, 0, 0, 0, 0, 0, 0] [0, 0, 1, 0, 0, 0, 0, 0]
The fat cat sat on the mat	[0, 0, 0, 1, 0, 0, 0, 0] [0, 1, 0, 0, 0, 0, 0, 0] [0, 0, 1, 0, 0, 0, 0, 0]
The fat cat sat on the mat	[0, 0, 0, 0, 1, 0, 0, 0] [0, 1, 0, 0, 0, 0, 0, 0] [0, 0, 1, 0, 0, 0, 0, 0]
The fat cat sat on the mat	[0, 0, 0, 0, 0, 1, 0, 0] [0, 1, 0, 0, 0, 0, 0, 0] [0, 0, 1, 0, 0, 0, 0, 0]
The fat cat sat on the mat	[0, 0, 0, 0, 0, 0, 1, 0] [0, 1, 0, 0, 0, 0, 0, 0] [0, 0, 1, 0, 0, 0, 0, 0]
The fat cat sat on the mat	[0, 0, 0, 0, 0, 0, 0, 1] [0, 1, 0, 0, 0, 0, 0, 0] [0, 0, 1, 0, 0, 0, 0, 0]

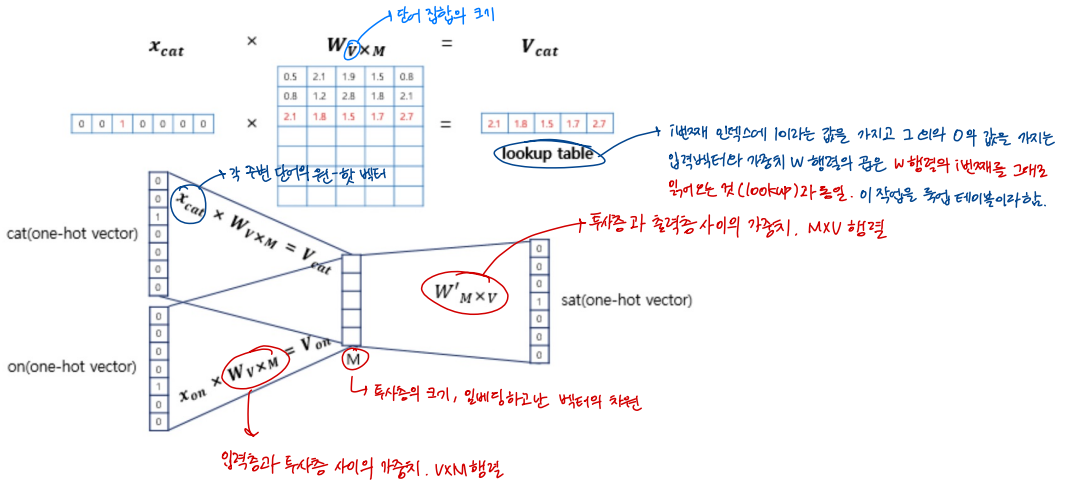
↑ 입력한 one-hot vector

→ 윈도우를 옆으로 움직여서 주변 단어와 중심단어의 선택을 변경하며 학습 데이터 셋을 만드는 방법

• CBOW 인공 신경망 구조화



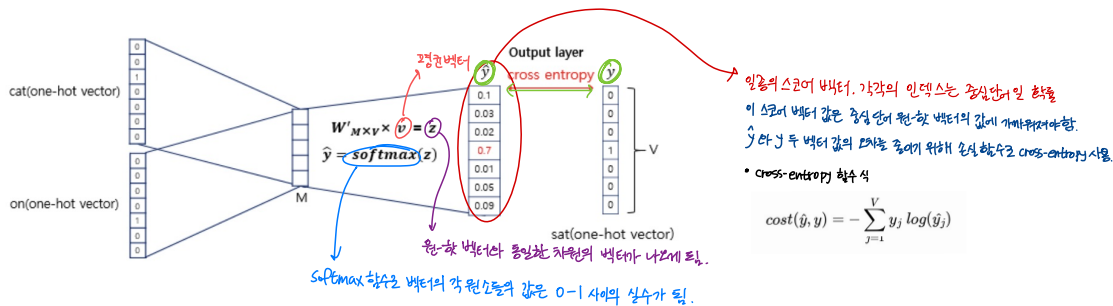
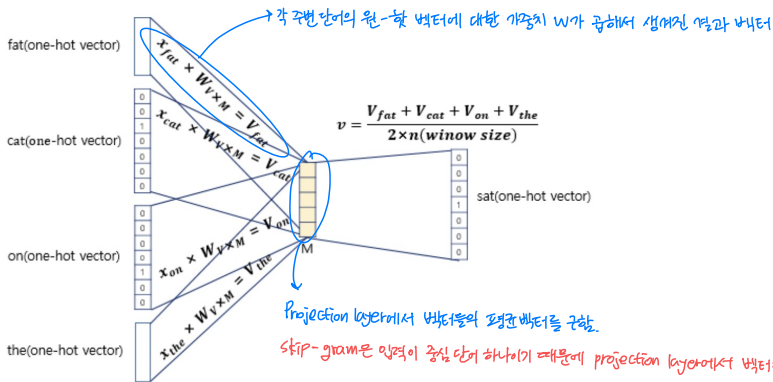
• 인공 신경망 확대 및 동작 예제니준



→ 인공 신경망 훈련 전, 가중치 행렬 W 와 W' 는 랜덤 값을 가지게 되고,

주변단어와 중심 단어는 더 정확하게 맞추기 위해 계속해서 W 와 W' 를 학습해가는 구조.

→ lookup 해은 W 의 각 행 벡터가 W2V 학습 후, 각 단어의 M 차원의 임베딩 벡터로 간주됨.



→ 역전파를 수행하면 W와 W'가 학습되고
M차원의 크기를 갖는 W의 행렬의 행은 각 단어의 임베딩 벡터로 사용하거나
W와 W' 행렬 두 가지 모두를 가지고 임베딩 벡터를 사용할 수 있음

• skip-gram

→ 중심 단어에서 주변 단어를 예측함

→ 중심 단어 하나에서 주변 단어를 예측하므로 투시층에서 벡터의 평균을 구하지 않음

중심 단어	주변 단어
cat	The
cat	Fat
cat	sat
cat	on
sat	fat
sat	cat
sat	on
sat	the

중심 단어

주변 단어

The fat cat sat on the mat

The fat cat sat on the mat

