

Doc2vec

발표자: 박채은(DSL 5기)

Doc2vec이란

Doc2vec :

Document embedding with paragraph vectors
document 임베딩 모델



paragraph¹



paragraph²



paragraph³



paragraph⁴

...



paragraph^N

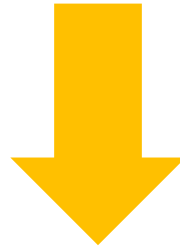
목표: Paragraph vector 구하기

Doc2vec: w2v을 확장해서 만든 모델

Word2vec

: 단어 임베딩 모델

훈련 방식: CBOW, Skip-gram



확장

Doc2vec

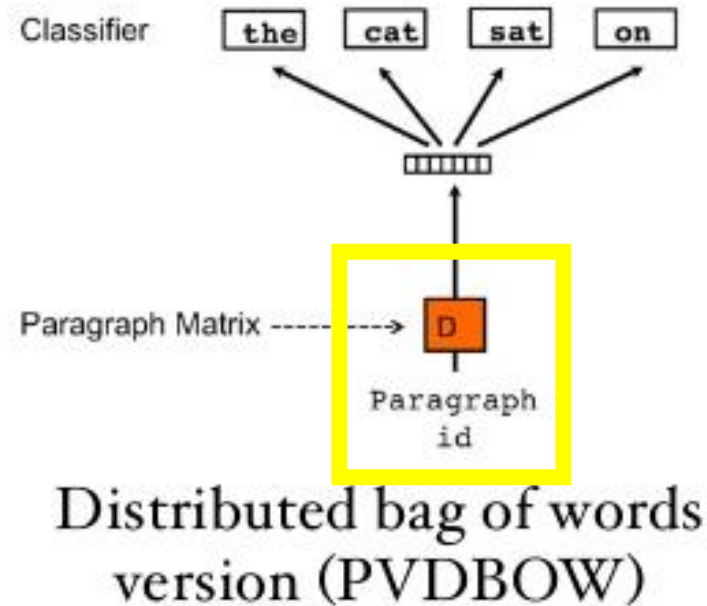
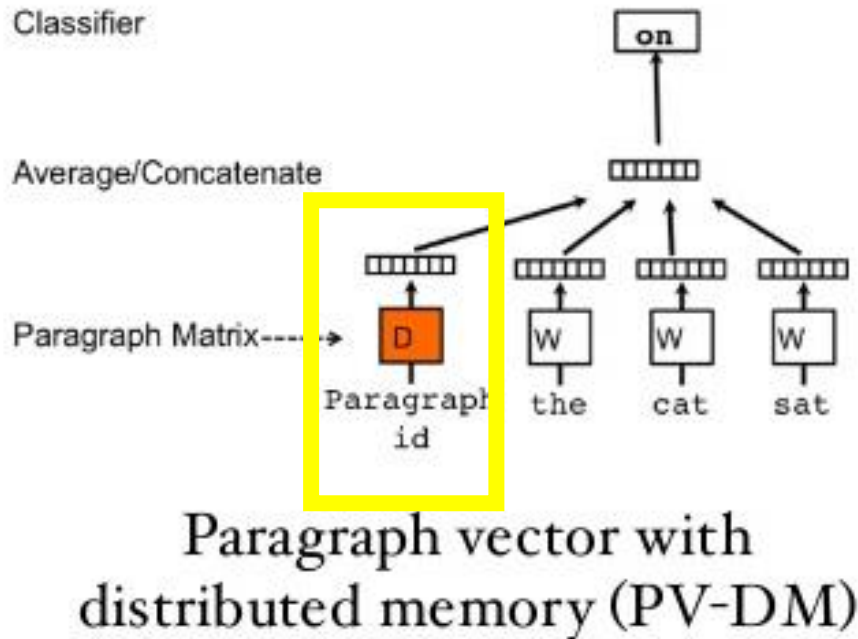
Word2vec의 장점을 그대로 가지게 됨.

훈련방식이 유사

Doc2vec: w2v을 확장해서 만든 모델

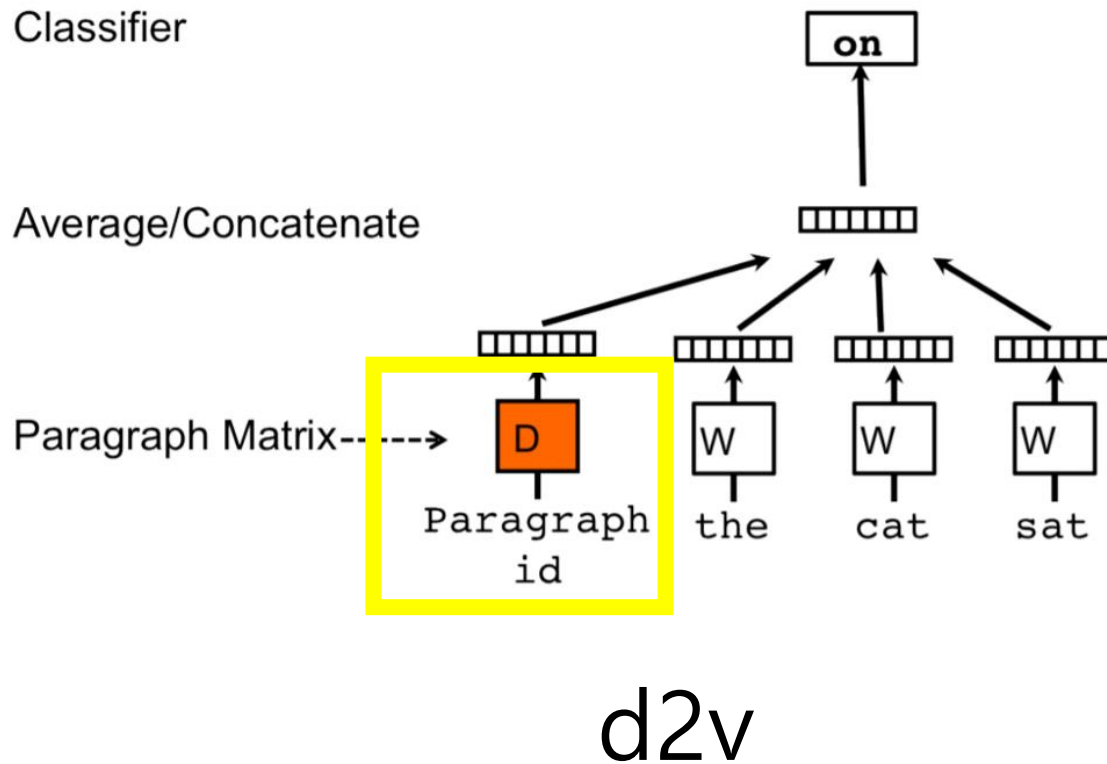
W2v model + 각각의 문서를 나타내는 paragraph_id
-> paragraph vector 얻을 수 있음.

훈련 방식: PV-DM, PV-DBOW

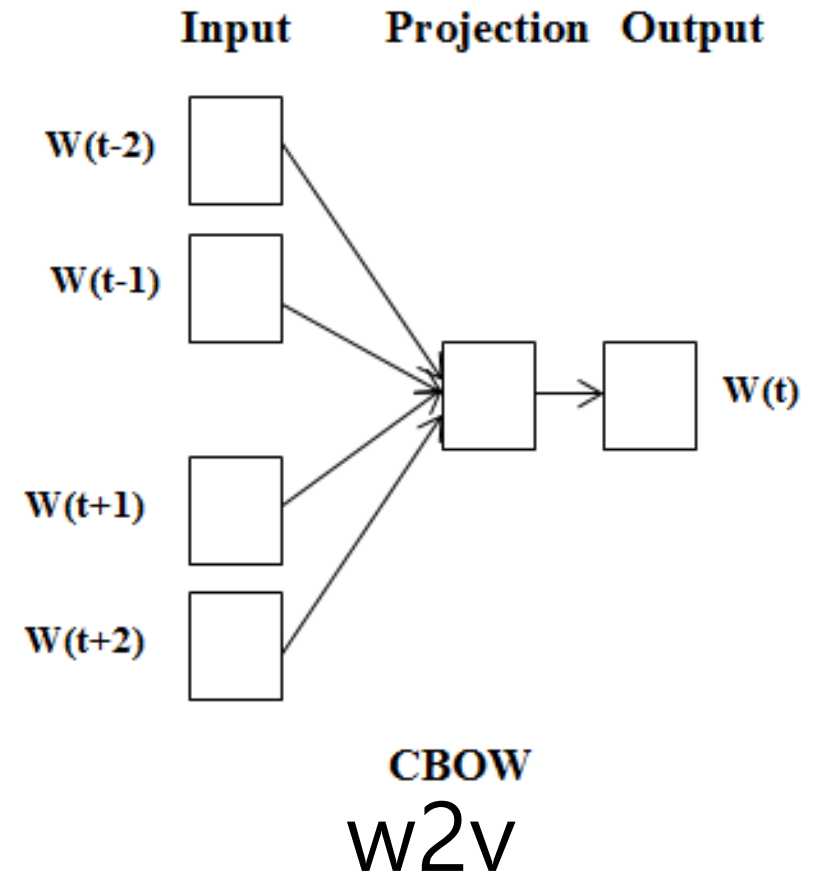


PV-DM

(paragraph vector with distributed memory)

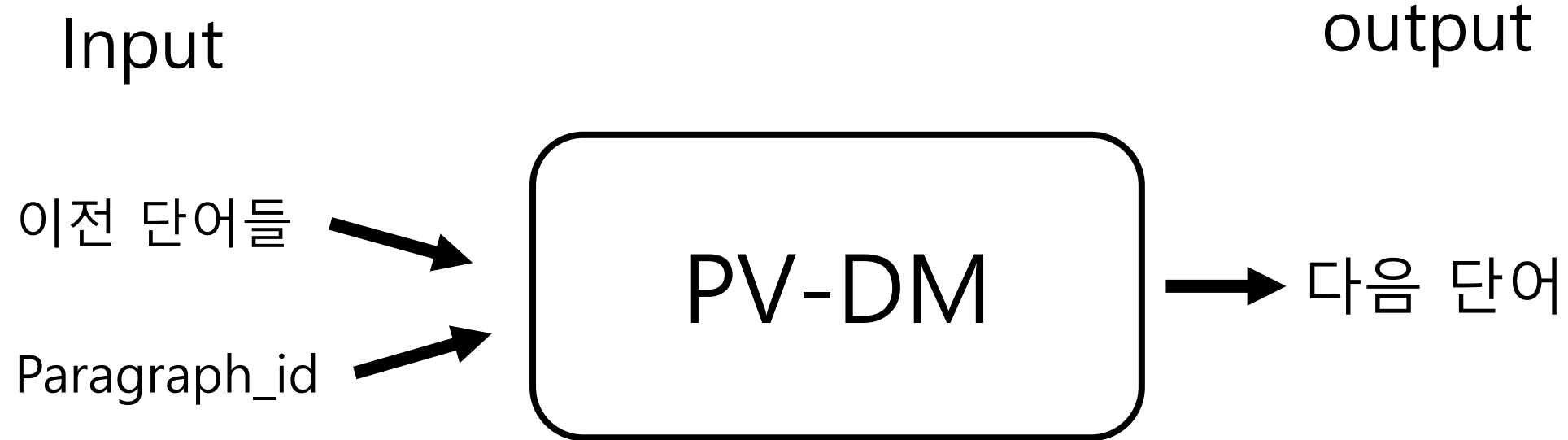


VS



PV-DM

(paragraph vector with distributed memory)

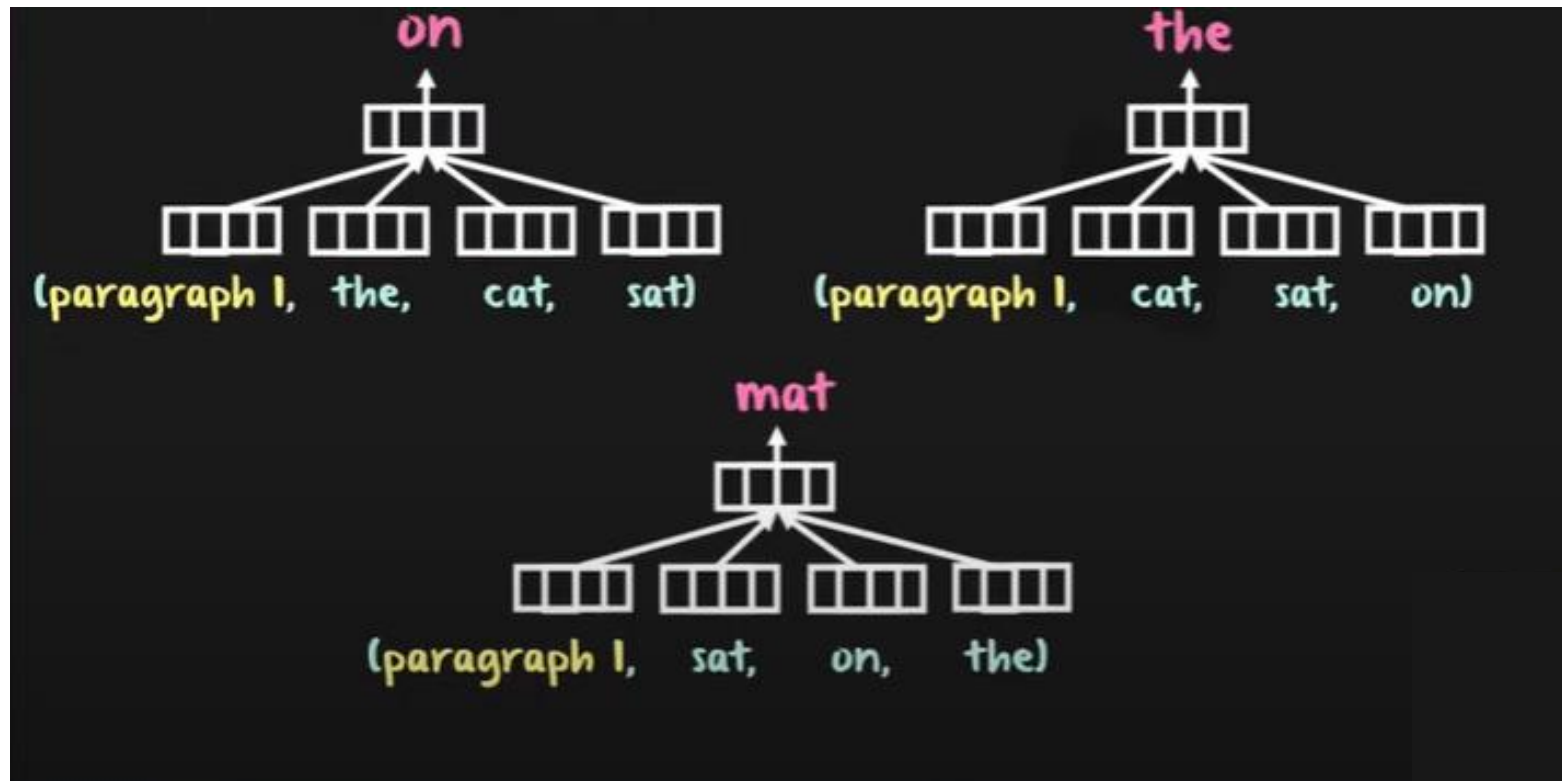


PV-DM

(paragraph vector with distributed memory)

- 예) Paragraph_1 : 'the cat sat on the mat.'

윈도우 사이즈: 3

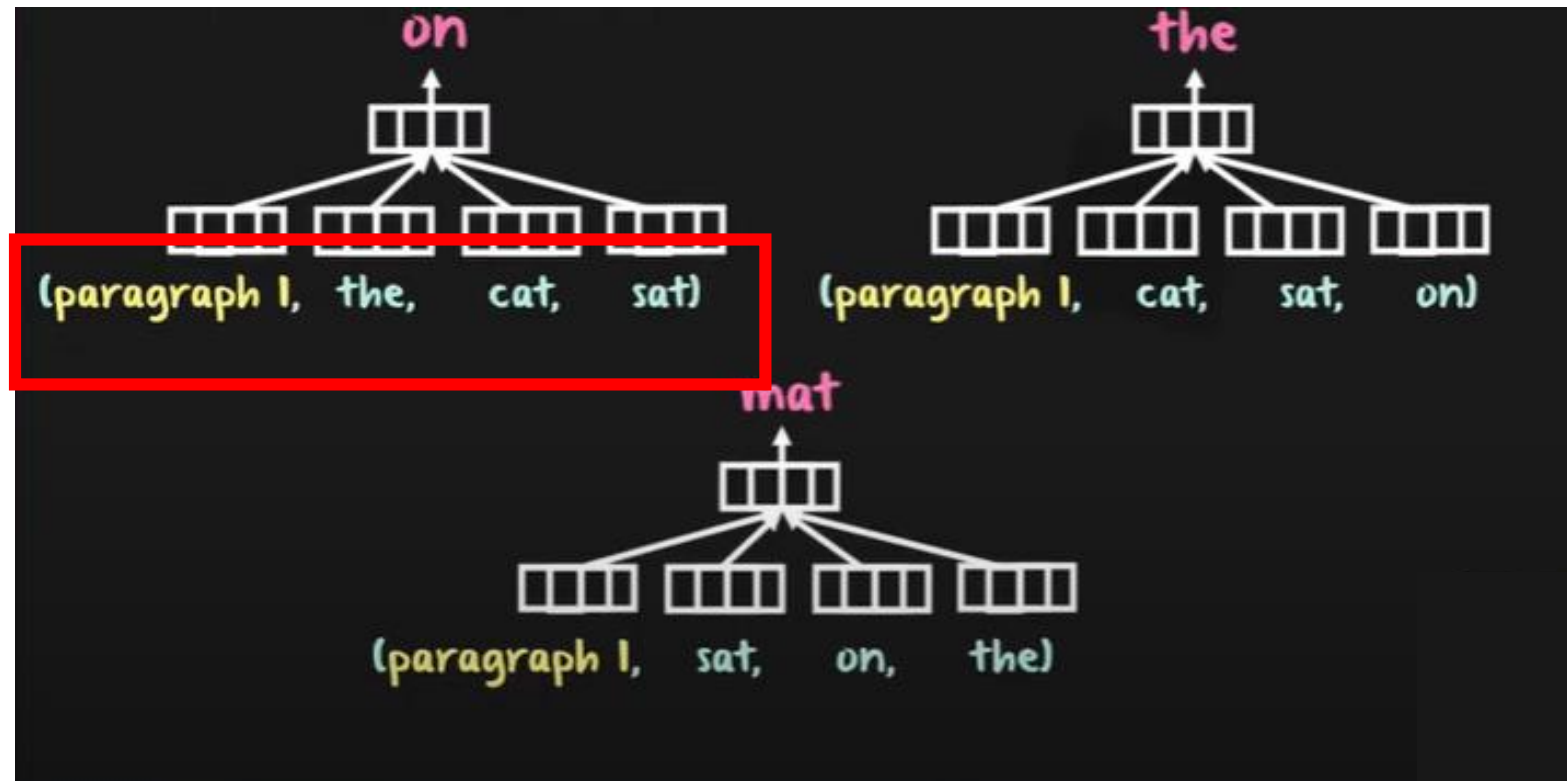


PV-DM

(paragraph vector with distributed memory)

- 예) Paragraph_1 : 'the cat sat on the mat.'

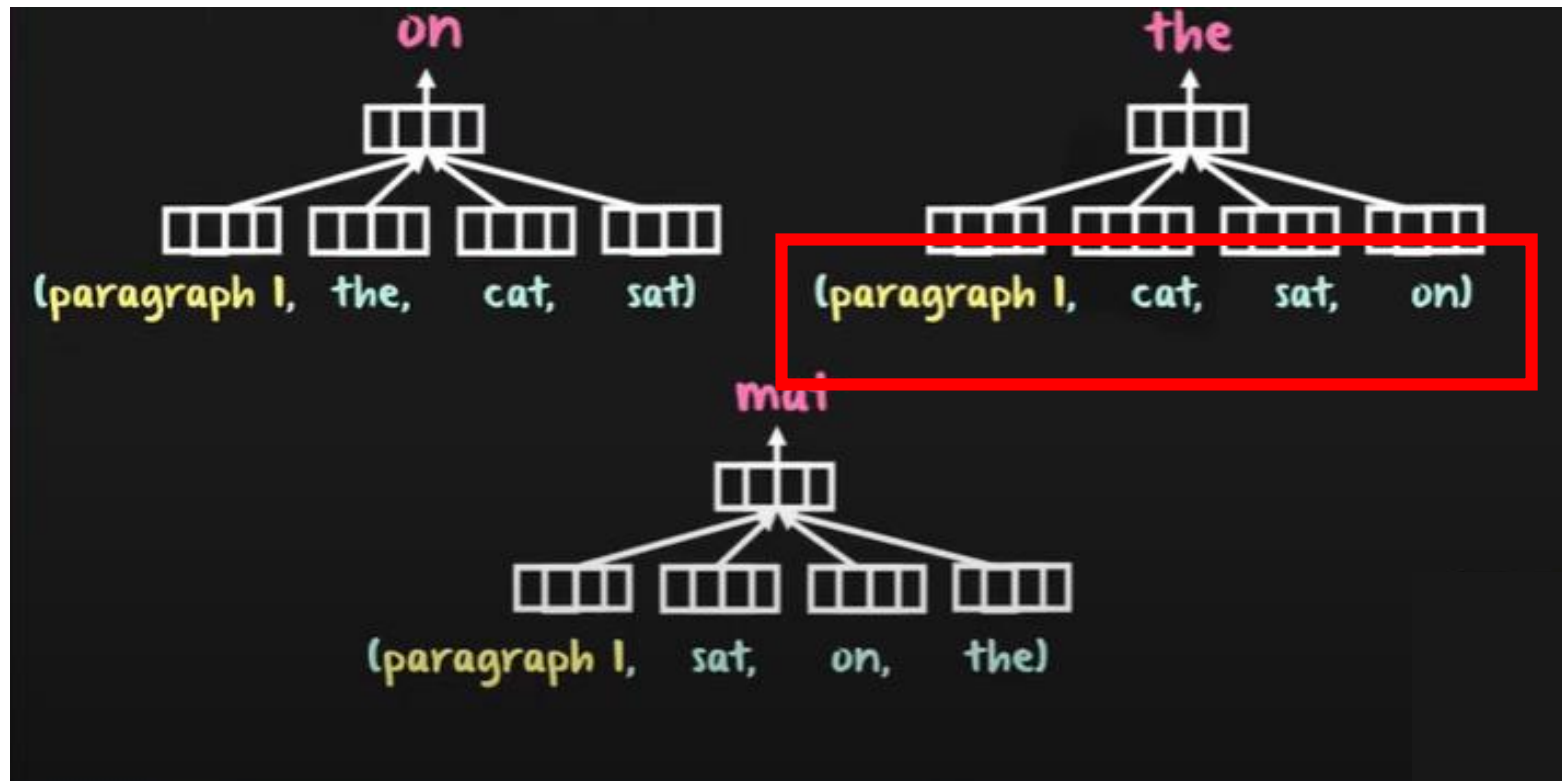
윈도우 사이즈: 3



PV-DM

(paragraph vector with distributed memory)

- 예) Paragraph_1 : 'the cat sat on the mat.'
윈도우 사이즈: 3

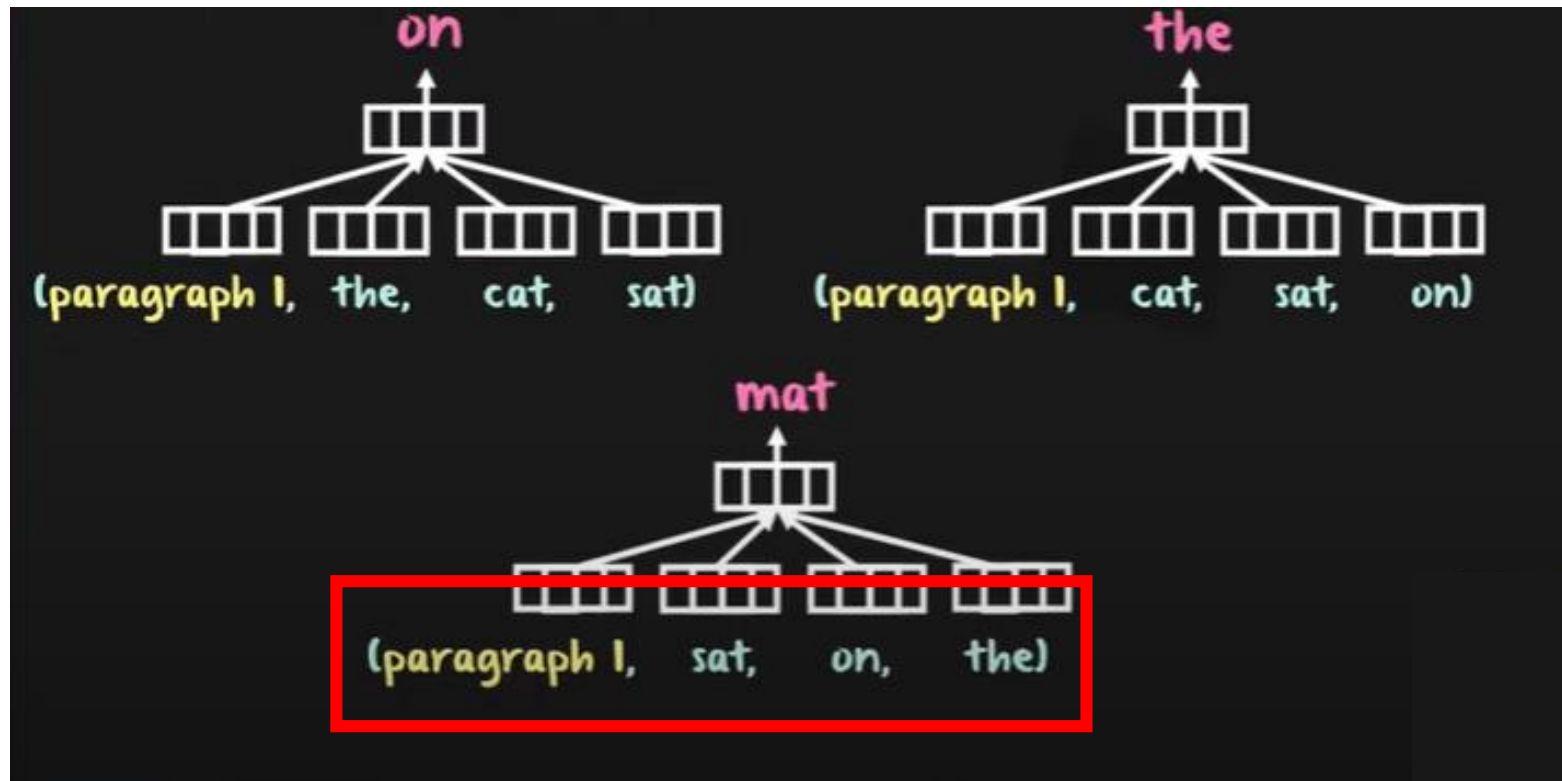


PV-DM

(paragraph vector with distributed memory)

- 예) Paragraph_1 : 'the cat sat on the mat.'

윈도우 사이즈: 3

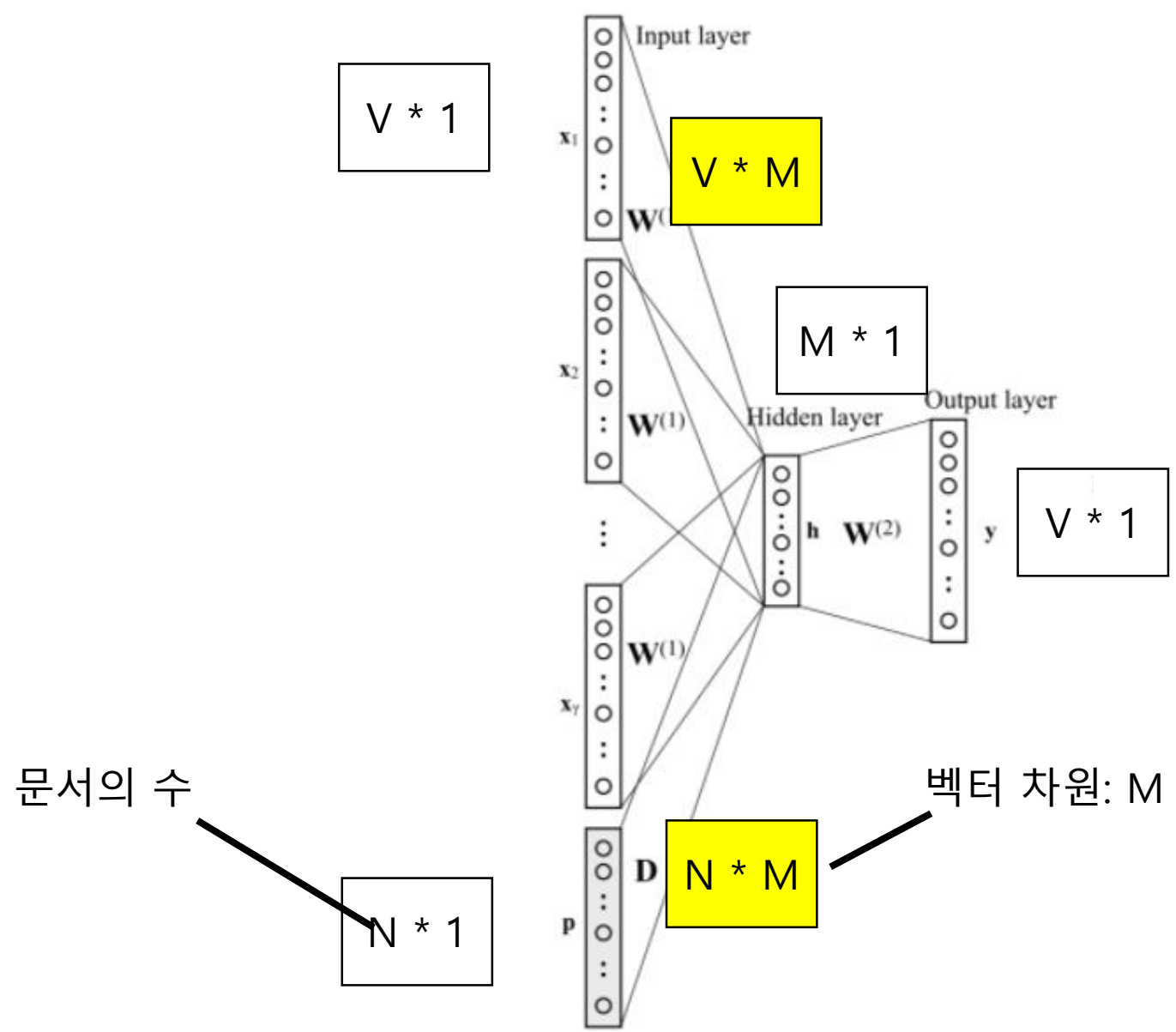


PV-DM

(paragraph vector with distributed memory)

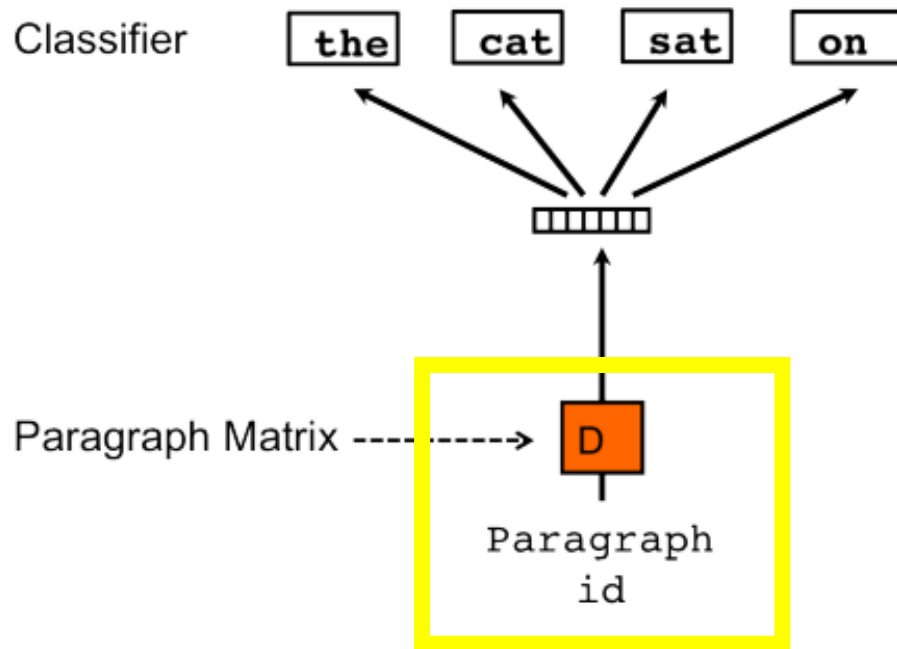
PV-DM이라는 이름이 붙은 이유는?

PV_DM을 통해 만들어진 벡터가
문서의 주제를 잡아주는 저장소의 역할을 하기 때문에



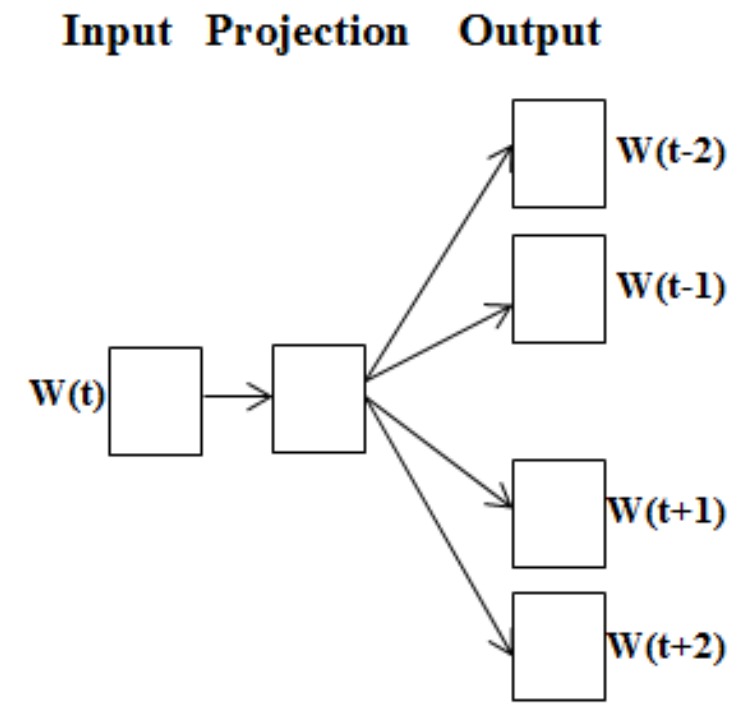
PV-DBOW

(Paragraph Vector with Distributed Bag Of Words)



d2v

VS

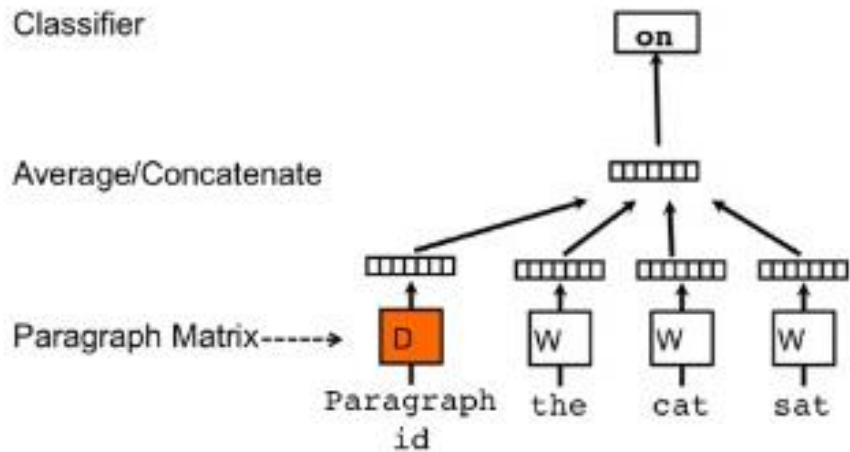


Skip-gram

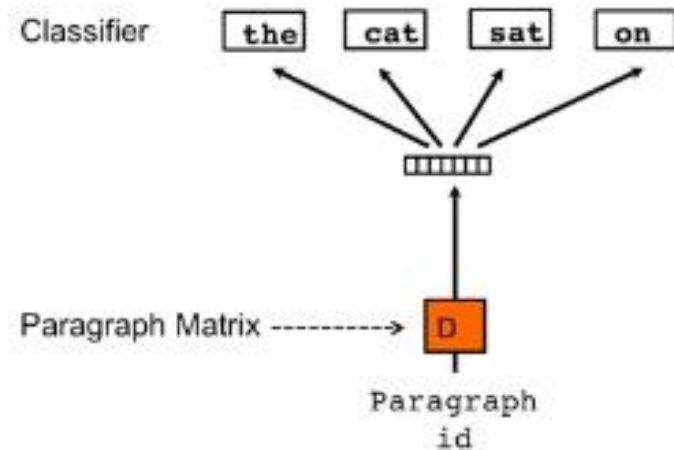
w2v

PV-DM과 PV-DBOW의 차이점

- word vector 생성 여부



Paragraph vector with distributed memory (PV-DM)



Distributed bag of words version (PVDBOW)



Thank
You!

A bright yellow rectangular sticky note is shown at a slight angle. It has a piece of translucent grey tape at the top-left corner. The words "Thank" and "You!" are written in a blue, hand-drawn, sketchy font. The note has a soft drop shadow on the white background.