

빠르게 해보는 Word Embedding
- word2vec, doc2vec with gensim

신윤식

Word Embedding의 개념

Word Embedding의 역사 및 주요 방법

Word Embedding으로 풀 수 있는 문제의 유형

주요 방법들의 동작원리와 이론

적용 가능한 API 등 소개

실제 예제적용 및 결과

•

•

15분

모든 것을 할 수는 없다.

한번쯤 거쳐가는 방법

보면서 쉽게 이해할 수 있는 주제

쉽고 빠른 구현이 가능한 API 존재



방법 - word2vec, doc2vec

예시 - 문서분류(Document Classification)

Syntactic Similarity

밥을 먹었다.

밥이 먹고 싶다.

밥에 콩을 넣었다.



“구문”의 차이

Semantic Similarity

성질이 많이 죽었다.

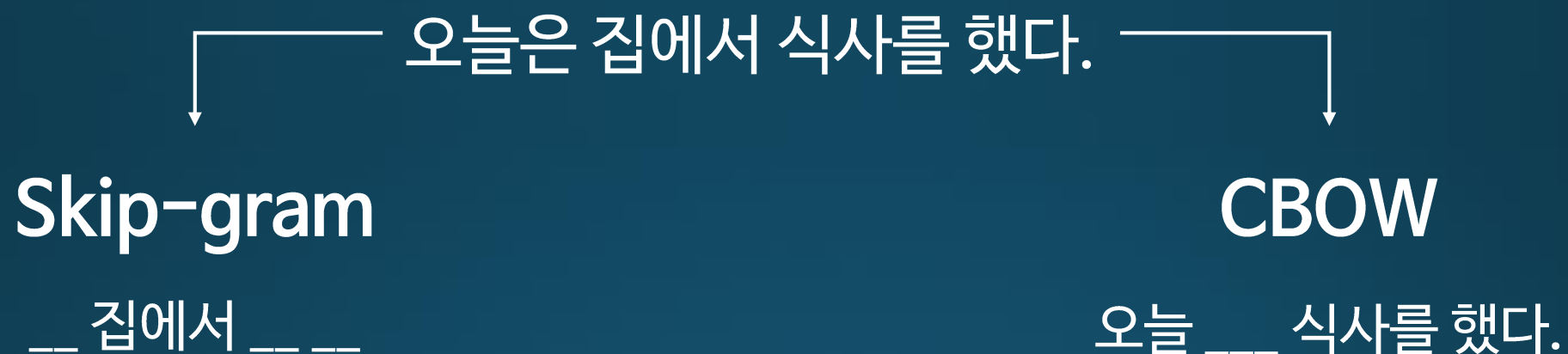
솜이불 솜이 죽었다.

기르던 화초가 죽었다.



“의미”의 차이

word2vec



word2vec

중심단어

주변단어

오늘은 집에서 식사를 했다.

window size = 1,

오늘은 → 집에서

집에서 → 오늘은

집에서 → 식사를

식사를 → 집에서

식사를 → 했다.

했다. → 식사를

word2vec

주변단어 중심단어 주변단어
오늘은 집에서 식사를 했다.

window size = 1,

오늘은 → 집에서

집에서 → 오늘은

집에서 → 식사를

식사를 → 집에서

식사를 → 했다.

했다. → 식사를

word2vec

오늘은 ^{주변단어} 집에서 ^{중심단어} 식사를 ^{주변단어} 했다.

window size = 1,

오늘은 → 집에서

집에서 → 오늘은

집에서 → 식사를

식사를 → 집에서

식사를 → 했다.

했다. → 식사를

word2vec

오늘은 집에서 주변단어 식사를 중심단어 했다.

window size = 1,

오늘은 → 집에서

집에서 → 오늘은

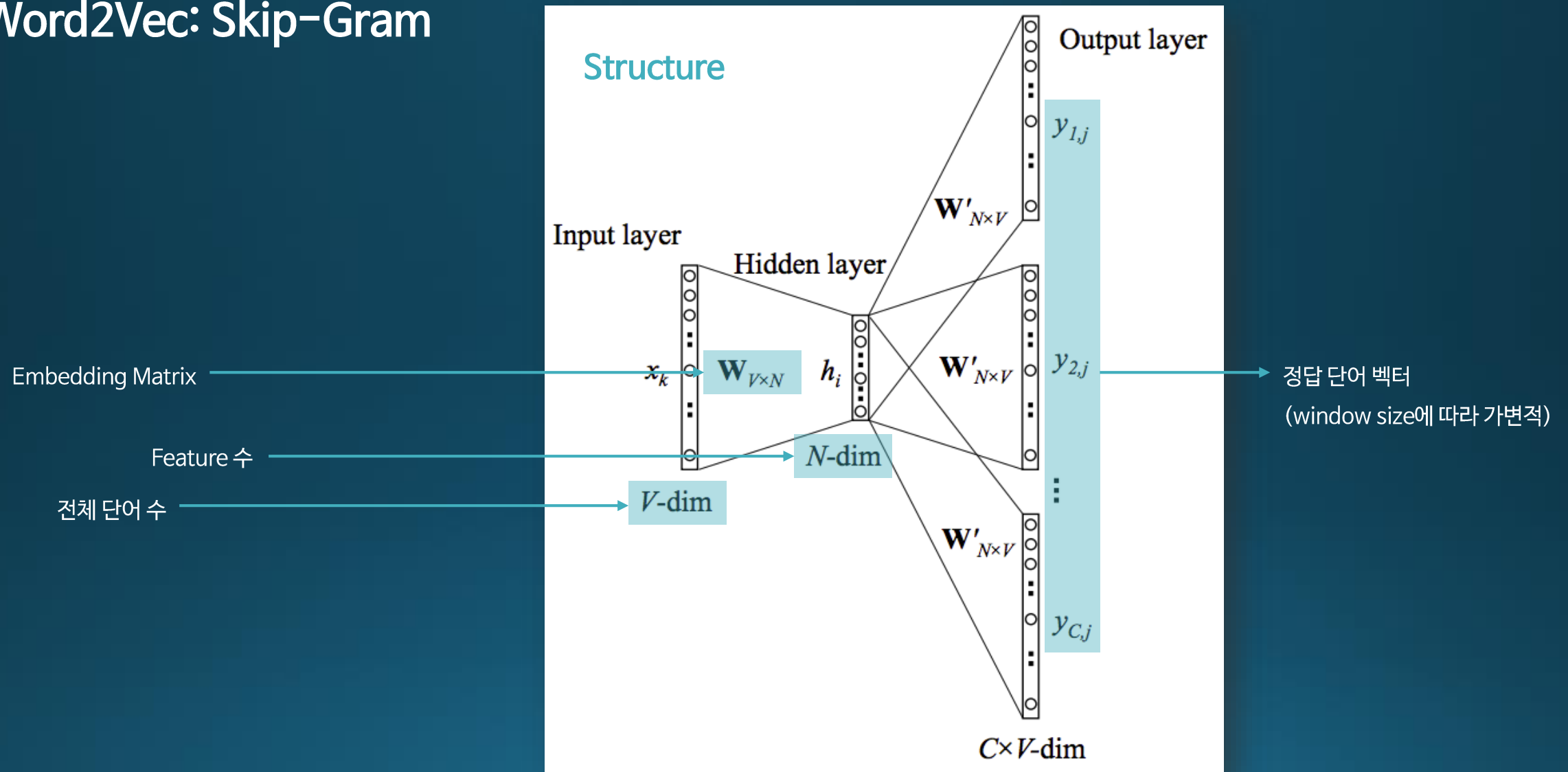
집에서 → 식사를

식사를 → 집에서

식사를 → 했다.

했다. → 식사를

Word2Vec: Skip-Gram



Rong, Xin. "word2vec parameter learning explained." *arXiv preprint arXiv:1411.2738* (2014).

Word2Vec: Skip-Gram

$$\text{similarity} = \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\|_2 \|\mathbf{B}\|_2}$$

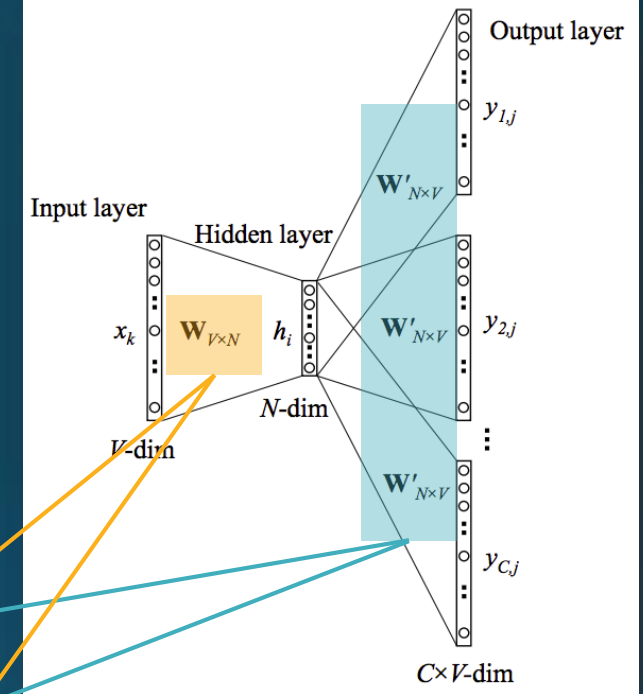
Cosine similarity From Wikipedia

maximize $v(\text{중심단어}) \times v(\text{주변단어})$
 minimize $\sum v(\text{중심단어}) \times v(\text{모든단어})$

Objective Function

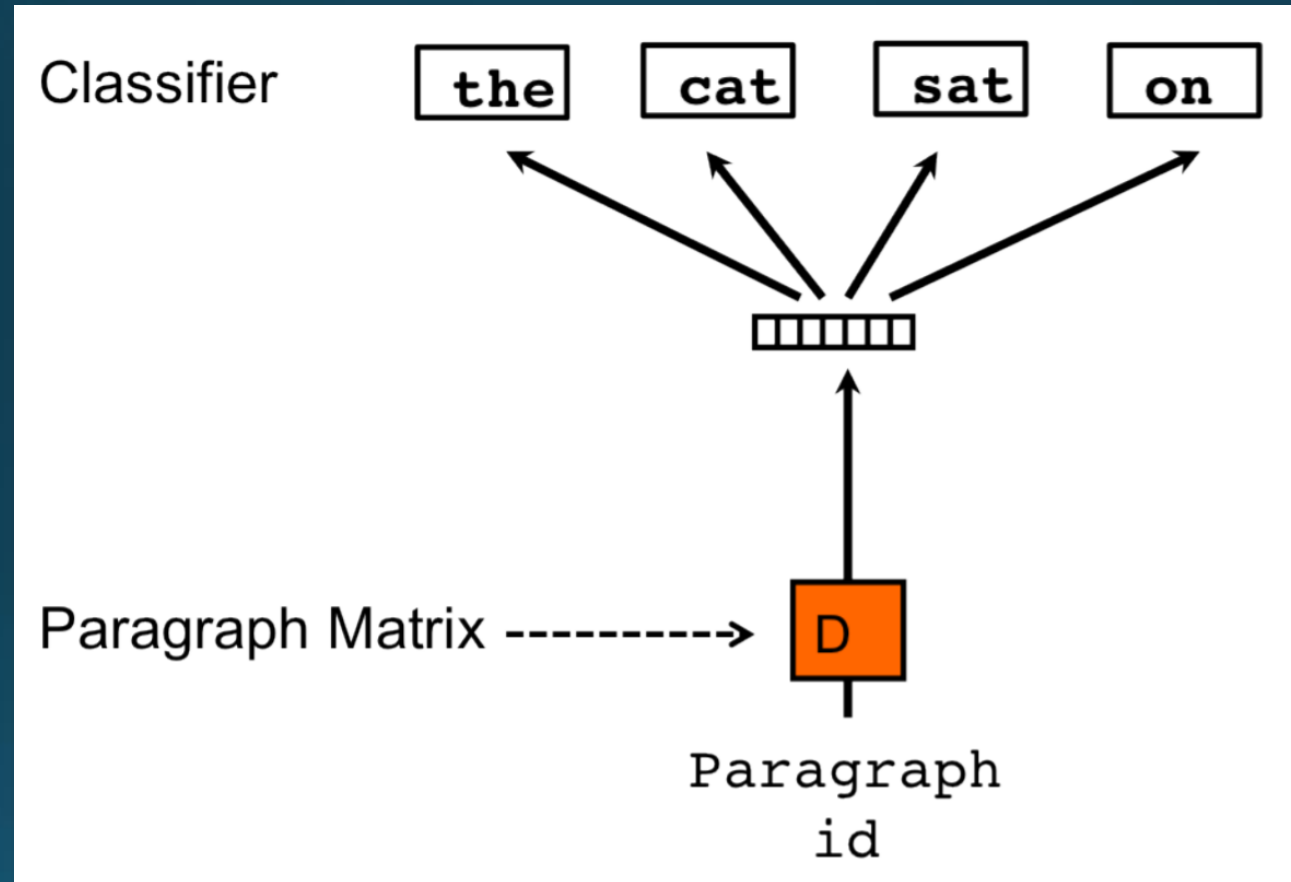
$$p(w_O | w_I) = \frac{\exp \left(v'_{w_O}{}^T v_{w_I} \right)}{\sum_{w=1}^W \exp \left(v'_w{}^T v_{w_I} \right)}$$

Mikolov, Tomas, et al. "Distributed representations of words and phrases and their compositionality." *Advances in neural information processing systems*. 2013.



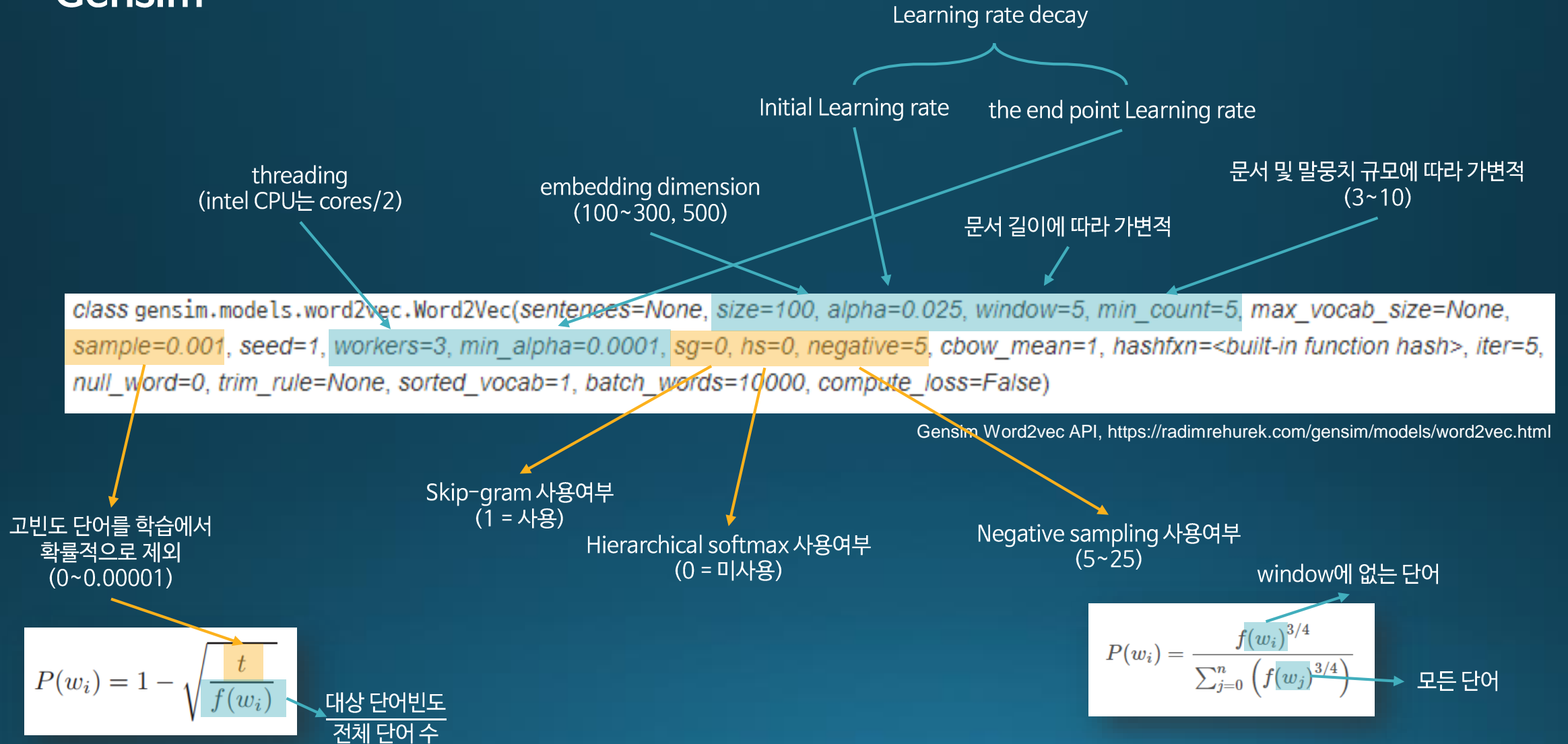
Rong, Xin. "word2vec parameter learning explained." *arXiv preprint arXiv:1411.2738* (2014).

Doc2Vec: DBOW



Le, Quoc, and Tomas Mikolov. "Distributed representations of sentences and documents." *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*. 2014.

Gensim



Q&A

E.O.D