

# Word Embedding

강사 : 백병인

[pi.paek@modulabs.co.kr](mailto:pi.paek@modulabs.co.kr)

모두의연구소 Research Scientist



2019 모두의연구소

# Language Representation to Neural Network

- 컴퓨터는 0과 1의 숫자밖에 모른다.
- 글자에 대응하는 숫자를 부여하면?
  - $A = 1 = 01_{(2)}$
  - $B = 2 = 10_{(2)}$
  - $C = 3 = 11_{(2)}$
- $A + B = C$  의 관계가 있을까?
- A와 B의 거리는 A와 Z의 거리보다 가까울까?



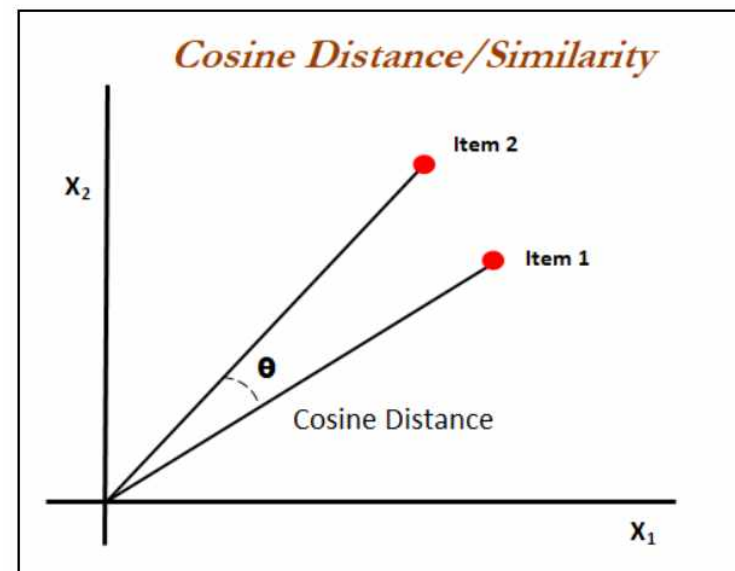
# One-Hot Encoding

- symbol을 vectorize하는 가장 단순한 방법
- 사전의 전체 단어 개수(N) 만큼의 차원을 가지는 vector로 표현
  - $V_{\text{hotel}} = [0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0]$
  - $V_{\text{motel}} = [0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0]$
  - $V_{\text{car}} = [0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0]$
- 문제점
  - 너무 길다
  - 단어간 유사도가 전혀 표현되지 않는다.

# Cosine Similarity

- 좋은 Vector Representation이 가져야 할 조건
  - 만약 language symbol의 vector representation을 만든다면,
  - 단어들 사이의 유사도를 정량적으로 표현할 수 있어야 한다.
- Cosine Similarity
  - 이 값이 클수록, 두 벡터는 유사한 개념을 표현한다고 보는 관점

$$\cos(\theta) = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$



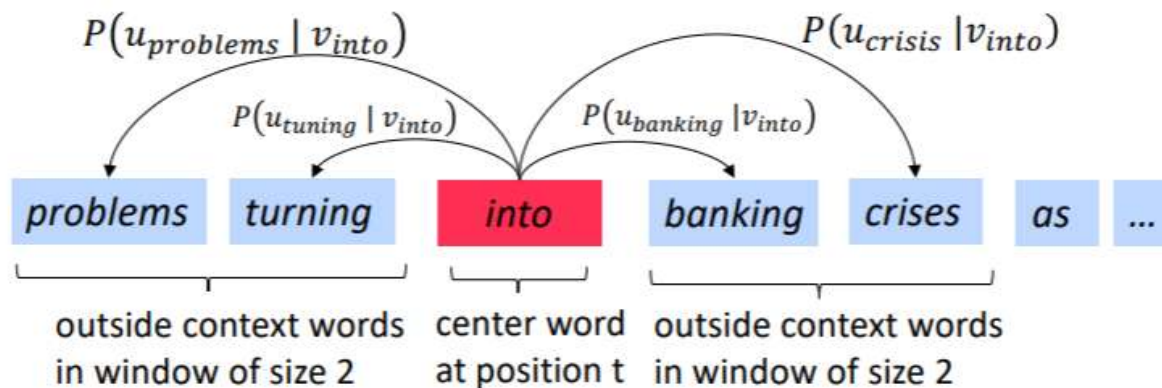
# Word2Vec (Mikolov et.al. 2013)

- 어떻게 word를 의미공간에 vector로 매핑할 수 있을까?
- 과연 word의 의미는 어떻게 발견되는 것일까?
- word(symbol)의 의미는 word 안에 없다.
- word가 사용되는 맥락(context)이 word의 의미를 규정한다.
- Distributional semantics
  - A word's meaning is given by the words that frequently appear close-by.
  - "You shall know a word by the company it keeps" (J. R. Firth 1957: 11)
  - One of the most successful ideas of modern statistical NLP!



# Main Idea of Word2Vec (1)

- **Co-Occurrence** Means Similarity of Meaning!!
- 한 문장 안에 함께 나타난 두 단어 사이에는 의미적 유사도가 있다는 가정



## Word Vector Dictionary

$$\theta = \begin{bmatrix} v_{aardvark} \\ v_a \\ \vdots \\ v_{zebra} \\ u_{aardvark} \\ u_a \\ \vdots \\ u_{zebra} \end{bmatrix} \in \mathbb{R}^{2dV}$$

# Main Idea of Word2Vec (2)

- 그렇다면, 함께 출현한 두 단어 사이의 cosine similarity를 높여 주는 쪽으로 딥러닝을 구현해 보면 어떨까?

Exponentiation makes anything positive

Dot product compares similarity of  $o$  and  $c$ .  
 $u^T v = u \cdot v = \sum_{i=1}^n u_i v_i$   
Larger dot product = larger probability

$$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$

Normalize over entire vocabulary to give probability distribution

- 위  $P(o|c)$ 가 corpus 데이터를 가장 잘 반영하도록 Log Likelihood를 극대화해 보자!
- 그런데, vocabulary 전체에 대해 normalize를 한다면 계산량이 너무 많지 않을까?



# Negative Sampling

- 모든 단어에 대해서 계산하지 말고 실제로 불가능한 예제(네거티브 샘플)에 대해서만 손실함수를 계산하자
  - $p(o|c)$  확률분포를 근사하는 multi-classification 문제를
  - $p(\text{positive or negative}|o,c)$  를 근사하는 binary classification 문제로 바꾼다. ( $o, c$ 는 sampling)

$$J_t(\theta) = \log \sigma(u_o^T v_c) + \sum_{i=1}^k \mathbb{E}_{j \sim P(w)} [\log \sigma(-u_j^T v_c)]$$

negative sample j

$p(\text{positive}|\text{positive sample of } o,c)$

$1-p(\text{positive}|\text{negative sample of } j,c)$

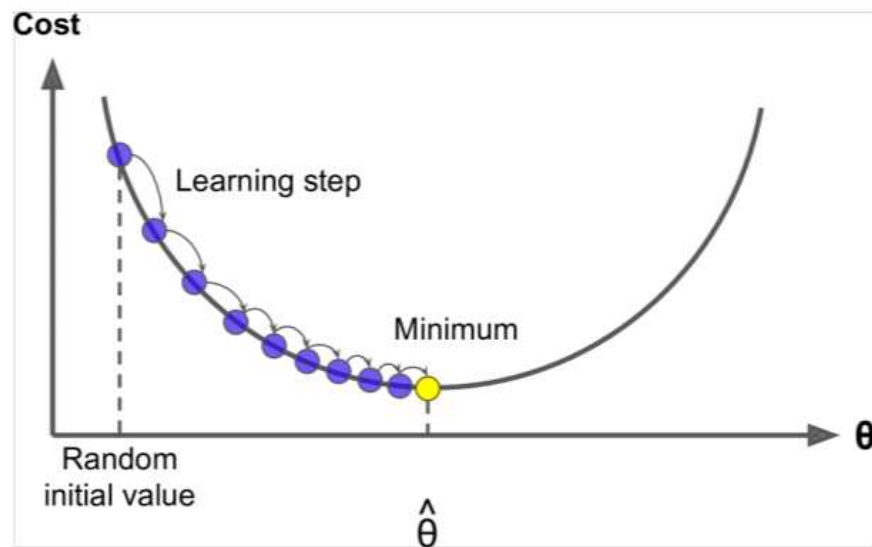
$\text{sigmoid}(-x) = 1 - \text{sigmoid}(x)$





# Optimization with SGD

- 엄청나게 많은 데이터. 그러므로 샘플링을 통한 SGD로 학습.



$$\theta^{new} = \theta^{old} - \alpha \nabla_{\theta} J(\theta)$$

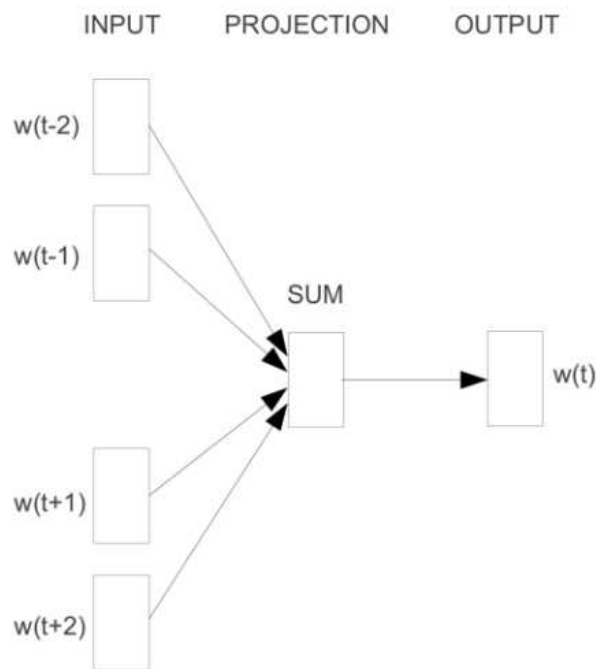
$\alpha$  = *step size* or *learning rate*

Algorithm:

```
while True:
    window = sample_window(corpus)
    theta_grad = evaluate_gradient(J, window, theta)
    theta = theta - alpha * theta_grad
```

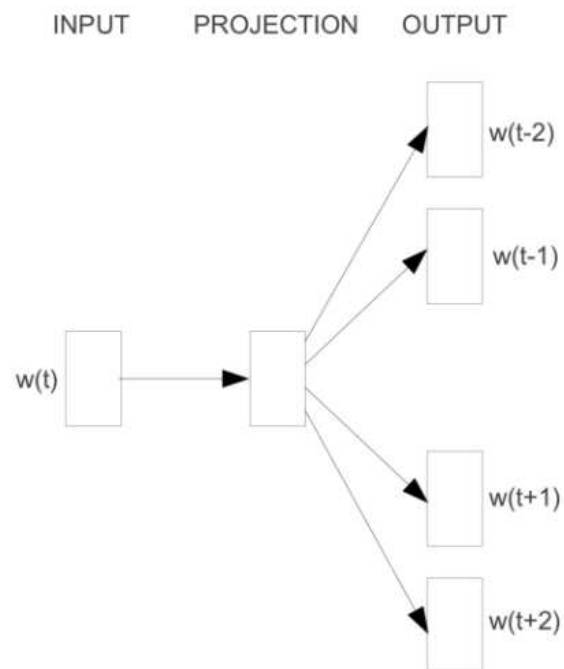


# CBoW & Skip-gram



## CBoW(Continuous Bag-of-words)

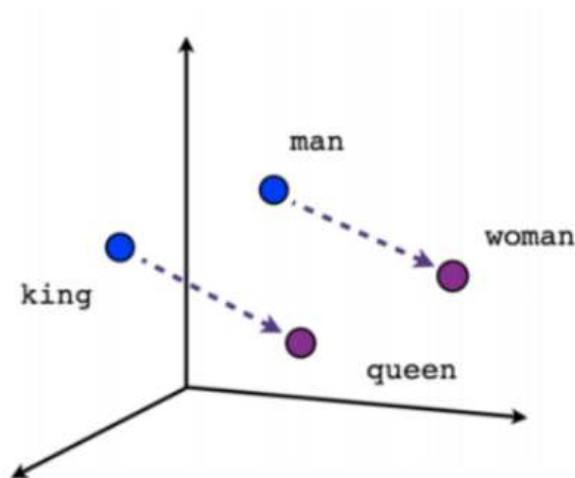
\* 다음 빈칸에 알맞은 말은?  
I go to \_\_\_\_\_ to study.



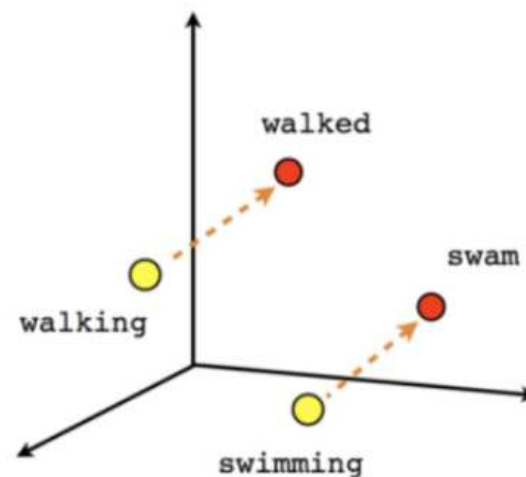
## Skip-Gram

\* 다음 단어로 짧은 글을 지어 보세요.  
\_\_\_\_\_ school \_\_\_\_\_

# Word Vector can do Vector Sum.



Male-Female



Verb tense

Visualize Word Vectors (<https://www.tensorflow.org/images/linear-relationships.png>)

$$V_{\text{king}} - V_{\text{man}} + V_{\text{woman}} = V_{\text{queen}}$$

# Semantic Distribution of Word Vector



Figure 7: Selected words projected along two axes:  $x$  is a projection onto the difference between the embeddings of the words *he* and *she*, and  $y$  is a direction learned in the embedding that captures gender neutrality, with gender neutral words above the line and gender specific words below the line. Our hard debiasing algorithm removes the gender pair associations for gender neutral words. In this figure, the words above the horizontal line would all be collapsed to the vertical line.

Any bias?? Why??

[Bolukbasi et.al. 2016] Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings  
<https://arxiv.org/pdf/1607.06520.pdf>

