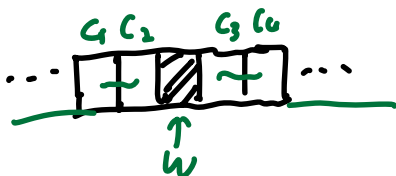


# Lecture 05 & 06 Word Embeddings

★ Skip-Gram model:  $\rightarrow \dots \boxed{c_1} \boxed{c_2} \boxed{w} \boxed{c_3} \boxed{c_4} \dots$



$$c[w] = \{c_1, c_2, c_3, c_4\}$$

Given a target word  $w$ , predict context  $c$ .

Text: "The quick brown fox jumps over the lazy dog."

$\rightarrow \text{Text} = [w_1, w_2, \dots, w_T]$   $T$  tokens.

Goal: for each  $w$ , maximize the corpus probability:

$$\arg \max_{\theta} \prod_{w \in \text{Text}} \left[ \prod_{c \in c[w]} p(c|w; \theta) \right] \quad (1)$$

Define  $D := \{ (w, c) : w \in \text{Text}, c \in c[w] \}$ .

$$\Rightarrow \arg \max_{\theta} \prod_{(w, c) \in D} p(c|w; \theta). \quad (2)$$

How can we properly define  $p(c|w; \theta)$ ?

Using softmax!

$$p(c|w; \theta) = \frac{e^{v_c \cdot v_w}}{\sum_{c' \in C} e^{v_{c'} \cdot v_w}},$$

- $v_c$ : vector representation of word  $c$ .
- $v_w$ : vector representation of word  $w$ .
- $C$ : set of all context words.
- $V$ : set of all target words.

To maximize (2), take log first:

$$\begin{aligned} \arg\max_{\theta} \sum_{(w,c) \in D} \log p(c|w; \theta) \\ := \sum_{(w,c) \in D} (\log e^{v_c \cdot v_w} - \log \sum_{c' \in C} e^{v_{c'} \cdot v_w}) \quad (3) \end{aligned}$$

We hope:

By maximizing (3), we can get  $v_1, v_2, \dots, v_{|V|}$ , which are good embeddings in the sense that similar words will have similar vectors!

But, what is time complexity of maximizing (3)

$|D| = |V| \cdot \text{window} \cdot |C| \approx \text{window} \cdot |V|^2$  : this only one epoch of using SGD.

Reducing time complexity:

①. hierarchical softmax

②. negative sampling ✓

\* negative sampling:

Idea: find parameters  $\theta$  such that the probabilities that all of observations indeed came from the data:

Given pair  $(w, c)$  from  $V \times C$ , what is the probability that  $(w, c)$  is indeed from  $V \times C$ ?

First try:

If I can find  $\theta$  such that  $p(w, c; \theta)$  is really large given  $(w, c) \in D$ , then this  $\theta$  maybe

good ?

To make this concrete, let

$p(D=1|w, c; \theta)$  : prob. that  $(w, c)$  is from  $D$ ,

then we should optimize the following:

$$\arg \max_{\theta} \prod_{(w, c) \in D} p(D=1|w, c; \theta) \quad (4)$$

take log and using  $p(D=1|w, c; \theta) = \frac{1}{1 + e^{-v_c \cdot w_w}}$

$$\Rightarrow \arg \max_{\theta} \sum_{(w, c) \in D} \log \frac{1}{1 + e^{-v_c \cdot w_w}}$$

By setting,  $v_c = v_w$  and  $v_c \cdot v_w = K$ ,  $K \geq 40$ .

then one can optimize (4)!

To prevent all vectors from having same value,  
by disallowing some  $(w, c)$ , to let the model

know :

$\begin{cases} (w, c) \in D \rightarrow p(D=1|w, c; \theta) \text{ should be high!} \\ (w, c) \notin D \rightarrow p(D=1|w, c; \theta) \text{ should be low!} \end{cases}$

We call  $(w, c)$  &  $D$  a negative pair.

How to generate a negative pair?

We have at most window  $\cdot |V|$  "positive" pairs.

What is the probability that randomly generate

$(w, c) \in V \times C$  is a "positive" pair?

$$\bullet \frac{\text{window} \cdot |V|}{|V| \cdot |C|} \approx \frac{\text{window}}{|C|} \text{ quite small!}$$

So,  $D'$  = set of negative samples! then

the goal is

$$\underset{\theta}{\text{argmax}} \prod_{(w, c) \in D} p(D=1 | w, c; \theta) \cdot \prod_{(w, c) \in D'} p(D=0 | w, c; \theta).$$

$1 - p(D=1 \dots)$

Take log and we have:

$$= \underset{\theta}{\text{argmax}} \sum_{(w, c) \in D} \log \frac{1}{1 + e^{-v_c \cdot v_w}} + \sum_{(w, c) \in D'} \log \frac{1}{1 + e^{v_c \cdot v_w}}$$

$$= \underset{\theta}{\text{argmax}} \sum_{(w, c) \in D} \log \sigma(v_c \cdot v_w) + \sum_{(w, c) \in D'} \log \sigma(-v_c \cdot v_w).$$

Remark:

①. For each  $(w, c) \in D$ , we generate  $k$  negatives.

$$p(w, c) \sim p_{\text{word}}(w) \cdot \frac{p_{\text{context}}(c)^{3/4}}{Z}$$

$Z$ : constant.

$p_{\text{word}}(w)$ : unigram distribution of words

$p_{\text{context}}(c)$ : unigram distribution of contexts.

$$\Rightarrow k \cdot |D| \approx |D'|.$$

Then, we can get the loss

$$- \sum_{t=1}^T \left[ \log \sigma(c_{\text{pos}} \cdot w_t) + \sum_{j=1}^k \log \sigma(-c_{\text{neg}_j} \cdot w_t) \right]$$

$J(c_{\text{pos}}, w_t, c_{\text{neg}_{1:k}}; \theta)$ :  $k+2$  vectors  
for each.

$$\text{SGD: } \theta^{t+1} = \theta^t - \eta_t \cdot \nabla J(c_{\text{pos}}, w_t, c_{\text{neg}_{1:k}}; \theta)$$

## ★ Training skip-gram of negative sampling.

Setting:

- $V$ : set of target words
- $C$ : set of context words (usually,  $V = C$ ).
- $\theta = [W_{in}, W_{out}]$ , where
- $W_{in} = W_{input}$ : hidden layer, projection layer
- $W_{out} = W_{output}$ : output layer.

(Recall our goal is given  $w_{in}$ , we predict context).

Objective:

$$\min_{\theta = [W_{in}, W_{out}]} - \sum_{t=1}^T \left[ \log \sigma(w_t^T \cdot c_{pos}) + \sum_{i=1}^k \log \sigma(-w_t^T \cdot c_{neg_i}) \right].$$

$w_t$ : current target word ( $t$ ).

$c_{neg_i}$ :  $i$ th negative word.  $i=1, 2, \dots, k$ .

$c_{pos}$ : the context word.

$$J(c_{pos}, c_{neg_{1:k}}, w) = - \left[ \log \sigma(w^T \cdot c_{pos}) + \sum_{i=1}^k \log \sigma(-w^T \cdot c_{neg_i}) \right]$$

SGD:

$$\theta^{t+1} = \theta^t - \eta_t \cdot \underline{DJ}(\theta^t).$$

given each pair  $(w, c_{pos}, c_{neg}; k)$ , the gradient:

$$Q. \frac{\partial J}{\partial c_{pos}} = - \frac{1}{\sigma(w^T \cdot c_{pos})} \cdot \sigma'(w^T \cdot c_{pos}) \cdot w$$

$$\text{Recall } \sigma'(x) = \sigma(x)(1 - \sigma(x))$$

$$= [\sigma(w^T \cdot c_{pos}) - 1] \cdot w.$$

$$(2) \frac{\partial J}{\partial c_{neg_i}} = ? \text{ How.}$$

Imagining  $w^T \cdot c_{pos} \Leftrightarrow (-w^T) \cdot c_{neg_i}$ ,

so, replacing  $w$  by  $-w$  on  $\frac{\partial J}{\partial c_{pos}}$ , you

will get:  $\frac{\partial J}{\partial c_{neg_i}}$ , we have:

$$\begin{aligned} \frac{\partial J}{\partial c_{neg_i}} &= [\sigma(-w^T c_{neg_i}) - 1] (-w) \\ &= [1 - \sigma(-w^T c_{neg_i})] \cdot w \end{aligned}$$



$$= \sigma(W^T \cdot C_{neg_i}) \cdot W, \quad i=1, 2, \dots, k.$$

Recall that:  $\sigma(-x) = 1 - \sigma(x)$ .

③,  $\frac{\partial J}{\partial W}$ , pretty easy,

$W^T \cdot C_{pos} = C_{pos}^T \cdot W$ : two variables have

no difference, the same:

$W^T \cdot C_{neg_i} = C_{neg_i}^T \cdot W$ . We have:

$$\frac{\partial J}{\partial W} = [\sigma(W^T \cdot C_{pos}) - 1] \cdot C_{pos} + \sum_{i=1}^k [\sigma(C_{neg_i} \cdot W)] \cdot C_{neg_i}.$$

Then, we have the updates of

$W, C_{pos}, C_{neg_{i:k}}$  at time  $t$ :

$$C_{pos}^{t+1} = C_{pos}^t - \eta_t \cdot [\sigma(C_{pos}^t \cdot W^t) - 1] \cdot W^t$$

SGD:  $C_{neg_i}^{t+1} = C_{neg_i}^t - \eta_t \cdot [\sigma(C_{neg_i}^t \cdot W^t)] \cdot W^t$

$$W^{t+1} = W^t - \eta_t \cdot \left[ [\sigma(C_{pos}^t \cdot W^t) - 1] \cdot C_{pos} + \right.$$

$$\left. \sum_{i=1}^k [\sigma(C_{neg_i}^t \cdot W^t)] \cdot C_{neg_i} \right]$$