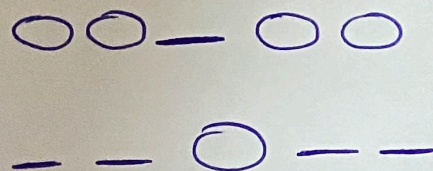


word 2 vec  $\rightarrow$  skip-Gram  
 $\rightarrow$  CBOW

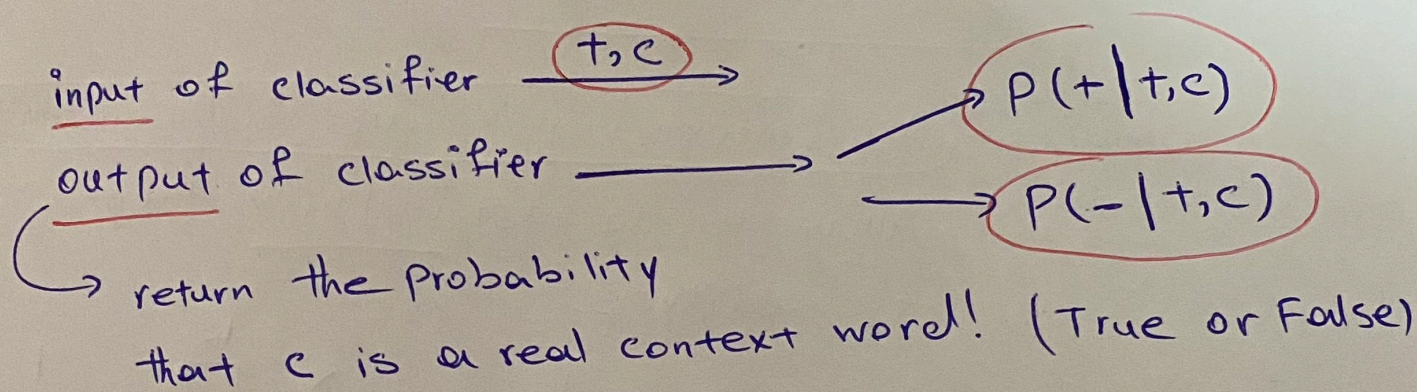


1

unsupervised classifier

skip Gram with negative sample (SGNS)

$P(+|t, c)$



compute the probability  $P$ ? intuition of skip-Gram is to base this probability on similarity

(A word is likely to occur near the target if its embedding is similar to target embedding)

$\rightarrow$  similarity in vectors (two vectors) are similar if they have a high dot product

so  $\Rightarrow$  similarity  $(t, c) \approx t \cdot c$

problem  $\rightarrow$  similarity of  $\vec{t}$  and  $\vec{c}$  is not Probability  
 $\rightarrow -\infty$  to  $+\infty$

solution  $\rightarrow$  sigmoid  $\sigma(x) = \frac{1}{1+e^{-x}}$



② So probability that word c is a real context

$$\text{True probability} \Rightarrow P(+ | t, c) = \frac{1}{1 + e^{-t \cdot c}}$$

note  $\Rightarrow$  this equation gives us the probability for one word. But we need to shake it over all context window.

\* Assumption: All context words are independent.  
Allow us to just multiply their probabilities.

$$P(+ | t, c) \longrightarrow P(+ | t, c_{1:k}) = \prod_{i=1}^k \frac{1}{1 + e^{-t \cdot c_i}}$$

$$* k = 2 * C$$

\*  $C$  = context window size

## Learning skip-Gram

initial embedding process with some pos and neg samples.

... lemon, a [tablespoon] of apricot jam, a] pinch...

+

t	c
apricot	of
apricot	tablespoon
apricot	jam
apricot	a

-

t	c	t	c
apricot	my	apricot	seven
⋮	⋮	⋮	⋮



③ \* training a Binary classifier with pos and neg samples.

\* skip-gram uses more neg examples than pos

How many  $\Rightarrow k \rightarrow$  in our example  $= \underline{2}$

so  $\rightarrow$  for each pos example  $t, c$ , we have 2 neg

→ Add noise  $\rightarrow$  should we add every words in our neg samples? NO

if we sample based on unweighted freq,  $P(w)$   
so  $\rightarrow P(\text{"the"})$ ,  $P(\text{"aardvark"})$  are same. which is actually not a right manner.

we need a weight  $\rightarrow \alpha = 0.75$

$$P^{\frac{3}{4}}(w) \Rightarrow \underline{P_{\alpha}(w) > P(w)}$$

→ Because it gives rare words

Finally  $\rightarrow$  the Goal of classifier

$$L(\theta) = \sum_{(t,c) \in +} \log(+|t,c) + \sum_{(t,c) \in -} \log(-|t,c)$$

objective

maximize

minimize



We can use SGD to train this objective.

↳ iteratively modifying the parameters to max obj  
(embedding for each target word  $t$  and each context word)

### embedding matrices

target embedding  $t \Rightarrow 1 \times d$

context embedding  $c \Rightarrow d \times 1$

