

Skip-Gram

Input : Large corpus of sentences

Output : \bar{v}_w, \bar{c}_w for each word type w

Hyperparams : word vector dimension d ($\approx 50-300$)
window size K (assume $k=1$)

The film inspired word context
film \rightarrow inspired
film \rightarrow The

Take all neighbors of each word taken up to k positions away

• Skipgram model : Probabilistic model of context given a word

$$P(\text{context} = y \mid \text{word} = x) = \frac{\exp(\bar{v}_x \cdot \bar{c}_y)}{\sum_{y' \in V} \exp(\bar{v}_x \cdot \bar{c}_{y'})}$$

Sum over all vocab

\bar{v}, \bar{c} : model params
 $|V| \times d$ dimensions

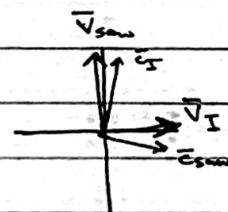
(2 * |V| not params in model)

If \bar{v}_x is similar to \bar{c}_y , y is likely to be in x 's context

Ex.

Corpus = I saw

$$\bar{v}_I = [1, 0] \quad \bar{v}_{\text{saw}} = [0, 1]$$



<u>word</u>	<u>context</u>
I	saw
saw	I

If $\bar{c}_{\text{saw}} = [1, 0]$ and $\bar{c}_I = [0, 1]$, what is $P(\text{context} \mid \text{word} = \text{saw})$?

$$\begin{aligned} \exp(\bar{v}_{\text{saw}} \cdot \bar{c}_I) & \exp(\bar{v}_{\text{saw}} \cdot \bar{c}_{\text{saw}}) \\ \approx 3 & = 1 \text{ (b/c orthogonal)} \end{aligned}$$

$$P(\text{context} \mid \text{word} = \text{saw}) = \frac{3}{4} \quad \text{saw} \mid \text{saw} = \frac{1}{4}$$

- Training Skip-gram

$$\text{Maximize } \sum_{\substack{(x,y) \\ \text{pairs in data}}} \log P(\text{context} = y \mid \text{word} = x)$$

"Impossible Problem"
cannot drive $P \rightarrow 1$

* Initialize parameters randomly (unlike neural nets)

↳ Iterating over data, model pulls similar vectors together over time

* Other Word-Embedding Methods on "seg-19.pdf"