

Skip-gram

Input: large corpus of sentences

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

Hyperparams: word vector dim d ($\sim 50 - 300$)
window size k (assume $k=1$)

The film inspired



word context
film \rightarrow inspired
film \rightarrow The

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

Skip-gram: probabilistic model of context | word

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

\bar{v}, \bar{c} model params

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

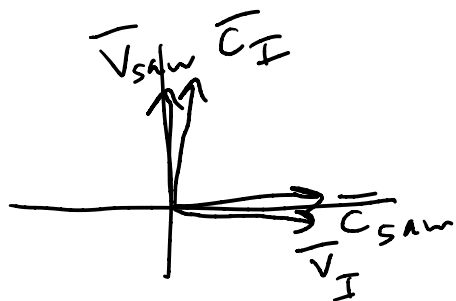
sum over vocab $\rightarrow y' \in V$

$2 \cdot |V| \cdot d$
params in model

Ex Corpus = I saw

$$\bar{v}_I = [1, 0]$$

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



word	context
I	saw
saw	I

If $\bar{c}_{saw} = [1, 0]$ and $\bar{c}_I = [0, 1]$, what is

$$\exp(\bar{v}_{saw} \cdot \bar{c}_I) \approx 3$$

$$\exp(\bar{v}_{saw} \cdot \bar{c}_{saw}) \approx 1$$

$P(\text{context} | \text{word} = \text{saw})?$
 \uparrow vocab

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

Training

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

"Impossible" problem: cannot drive $P \rightarrow 1$

Initialize params randomly