# Word Embeddings

— low dim rep of words capturing similarity
$\underset{\sim 30-500}{c}$

Distributional hypothesis: "You shall know a word by the company it keeps"

Mikolov etal 2013 word2vec

each word w $\longrightarrow$ $\bar{v}_w$ word vector          predict each word's context given that word
$\bar{c}_w$ context vector

# Skip Gram

Input: large corpus of sentences
output: $\bar{v}_w, \bar{c}_w$ for each word w
Hyperparams: word vectors dim d ($\sim 50-300$)
window size k (k > 1)   neighbors of each word taken up to k
the film inspired          film $\to$ inspired                                    positions away
                           film $\to$ the
skip-gram:   context | word   $\longrightarrow$   $P(\text{context}=y|\text{word}=x) = \dfrac{\exp(\bar{v}_x \cdot \bar{c}_y)}{\underset{y' \in V}{\sum} \exp(\bar{v}_x \cdot \bar{c}_{y'})}$
                                                   sum over vocab

if $\bar{v}_x$ is $\sim$ to $\bar{c}_y$, y is likely to     $\bar{V}$ and $\bar{C}$ are model params
be in x's context                                          $2 \cdot |V| \times d$

ex: corpus = I saw
$\bar{v}_I = [1,0]$   $\bar{v}_{saw} = [0,1]$          if $\bar{c}_{saw} = [1,0]$ and $\bar{c}_I = [0,1]$
$\bar{v}_{saw}$ $\bar{c}_I$                                       $P(\text{context}|\text{word}=saw)$?

| word | context |
|------|---------|
| I | saw |
| saw | I |

$= \bar{v}_I$
$\bar{c}_{saw}$

$\underbrace{\exp(v_{saw} \cdot c_I)}_{\approx 3}$  $\underbrace{\exp(v_{saw} \cdot c_{saw})}_{1}$

$P(\text{context}=I | \text{word}=saw) = \dfrac{3}{4}$  $\begin{array}{c} saw|saw \\ = \frac{1}{4} \end{array}$

# Training

max $\underset{(x,y)}{\sum} \log P(\text{context}=y|\text{word}=x)$
     pairs in data          init params randomly

"impossible" to get $P \to 1$