

CS 378 Lecture 7: Word Embeddings

[add pronouns to zoom]

Announcements

- AZ out
- Readings updated (bolding)
- Mid-semester survey
- Download AZ code for this lecture

Today

- Intro to word embeddings
- Explore embeddings
- Skip-gram: model + training
- Revisit DANs

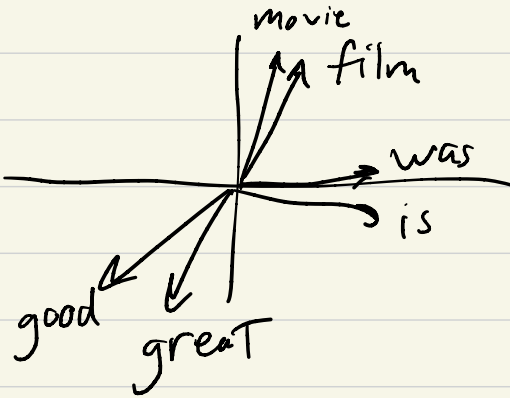
Word Embeddings

$$\begin{aligned} \text{Movie was good} &\rightarrow \left[\begin{array}{c} | \\ \text{movie} \end{array} \quad \begin{array}{c} | \\ \text{good} \end{array} \quad \begin{array}{c} | \\ \text{was} \end{array} \quad \dots \right] \\ &= \left[\begin{array}{c} | \\ \text{movie} \end{array} \right] + \left[\begin{array}{c} - \quad | \quad - \\ \text{good} \end{array} \right] \\ &\quad + \left[\begin{array}{c} - \quad | \quad - \\ \text{was} \end{array} \right] \end{aligned}$$

Each word is a $|V|$ -len vector
w/ a single 1

film is great
movie was good \Rightarrow dot prod = 0

Word embs: low-dimensional representations
(50-300)
that capture similarity



Distributional Hypothesis

JR Firth 1957: "you shall know a word by the company it keeps"

I watched the movie

I watched the film

The film inspired me

The movie inspired me

movie and film can show up in similar contexts

Are movie + film always substitutable?
polysemy: one word has multiple senses

Mikolov 2013: word2vec

Learn word + context vectors for each word

Attempt to predict context given word

Embedding properties

$$\text{sim}(\text{good}, \text{bad}) \approx 0.8$$

Skip-gram

Input: corpus of sentences

Output: \bar{v}_w, \bar{c}_w for each word w in the vocab

(what people use: \bar{v}_w OR \bar{c}_w
OR $\bar{v}_w + \bar{c}_w$)

Hyperparameters: dimension d
window size k

Let $k=1$

Form (word,
context pairs)

The film inspired

(film, The)

look k words in each
direction

(film, inspired)

$V = \text{vocab}$

Skip-gram model

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

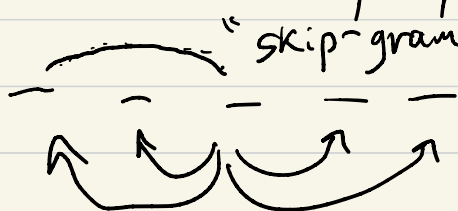
params: vectors \vec{v} $|V| \times d$
 \vec{c} $|V| \times d$

Training Take our corpus
Get (x, y) pairs
word context

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

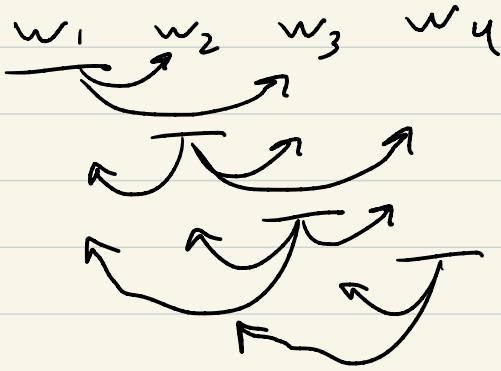
Randomly initialize \vec{v}, \vec{c} , use SGD

$k=2$



$k=3$: go 3 out, etc

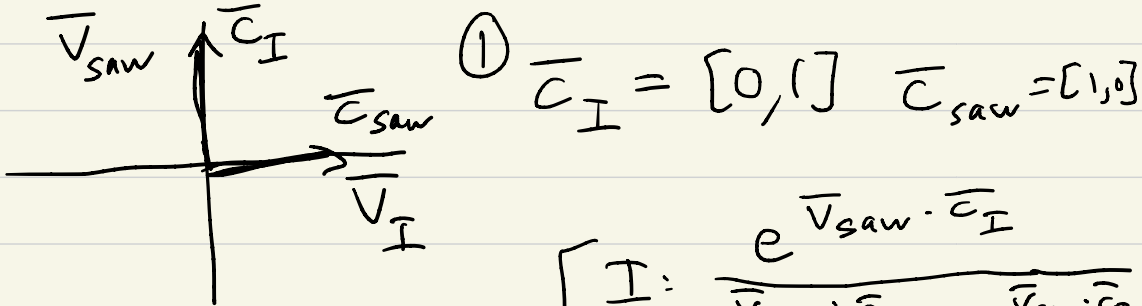
$k=2$



10 pairs as training data

Ex Corpus = I saw $k=1$
vocab = {I, saw}

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



$$P(\text{context} \mid \text{word} = \text{saw}) = \begin{cases} \text{I: } \frac{e^{\bar{v}_{saw} \cdot \bar{c}_I}}{e^{\bar{v}_{saw} \cdot \bar{c}_{saw}} + e^{\bar{v}_{saw} \cdot \bar{c}_I}} \\ \text{saw: } 1/4 \end{cases}$$

$3/1+3 = 3/4$

② (word = I, context = saw) ★
saw I

$$\bar{c}_{\text{saw}} = [100, 0] \quad \bar{c}_I = [0, 100]$$

$$\frac{e^{100}}{1 + e^{100}} = 0.99999 \dots$$

Maximizing likelihood is "impossible"!

(saw, I) can't assign prob 1 to
(saw, you) each

skip-gram is slow

for each example: $O(|V| \cdot d)$
multiplies