

# CS 378 Lecture 7

## Word Embeddings

### Announcements

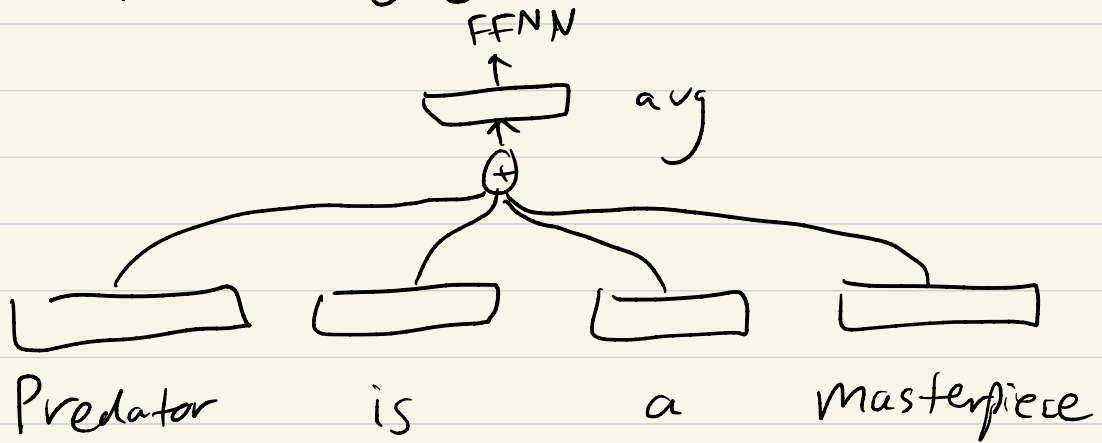
- Fairness response due today
- Friday: talk by Nanyun Peng (UCLA)  
11am 6.302

Recap Last time: FFNNs

$$P(\bar{y} | \bar{x}) = \text{softmax}(Wg(Vf(\bar{x})))$$

Pytorch basics

# Deep Averaging Network (A2)



How does this relate to BoW?

## Word Embeddings

So far: one-hot representations of words

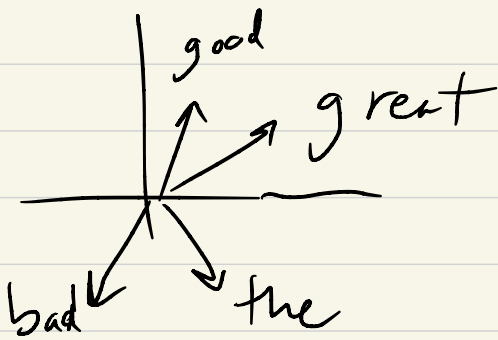
movie was good →

$$\begin{aligned}
 & \left[ \begin{array}{ccc} | & & | \\ \text{movie} & & \text{was} & & \text{good} \end{array} \right] \\
 = & \text{movie} \left[ \begin{array}{cccc} 0 & 0 & 0 & 0 \end{array} \right] + \text{single} \\
 & + \text{was} \left[ \begin{array}{cccc} 0 & 0 & - & - & 1 & 0 & 0 \end{array} \right] \\
 & + \text{good} \left[ \begin{array}{cccc} 0 & 0 & - & - & 1 \end{array} \right]
 \end{aligned}$$

Problem: ① Long vectors  
 ② good vs. great: not similar

film was great | movie is good

Instead of ~10k dims, how about ~100?



$\text{sim}(\text{good}, \text{great})$   
 $>$   $\text{sim}(\text{good}, \text{bad})$

## 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

I took a picture with film

There was a film on the liquid

Polysemous : word has multiple  
senses / meanings

Mikolov et al. 2013 word2vec

Learn 2 vectors for every word

word vec + context vec.

Attempt to predict context  
given word

## Skip-gram

Input: a corpus of text

Output:  $\vec{v}_w$ ,  $\vec{c}_w$  for each  $w$   
word Context in vocab

(for applications: use either  $\vec{v}$   
or  $\vec{c}$  or  $\vec{v} + \vec{c}$ )

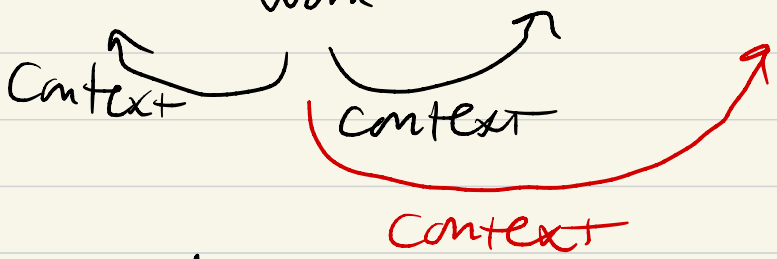
Hyperparameters:  $d$  (50 ~ 300)

window size  $K$

look in both directions

Let  $k=1$       2

The film inspired me  
word



word, context

(film, The)

(film, inspired)

(film, me)

} Training examples

Other pairs (The, film)

(inspired, ...)

Model: (skip-gram)

vocab  $V$


$$P(\text{context} = y \mid \text{word} = x)$$

$$= \frac{e^{\vec{v}_x \cdot \vec{c}_y}}{\sum_{y' \in V} e^{\vec{v}_x \cdot \vec{c}_{y'}}$$

distribution  
over all  
context  
words in  $V$

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

randomly initialized  $\vec{c}$   $|V| \times d$

Training  $(x, y)$   train examples

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

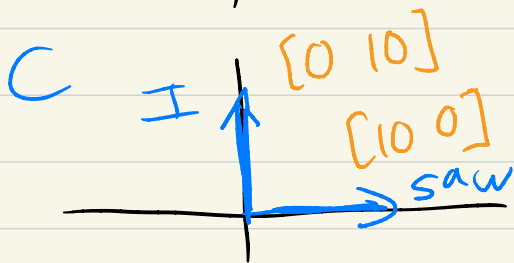


Ex Corpus = I saw  $k=1$   
 vocab = {I, saw}  $d=2$

Assume:  $\vec{v}_I = [1 \ 0]$   $\vec{v}_{saw} = [0 \ 1]$



① Let  $\vec{c}_{saw} = [1 \ 0]$   
 $\vec{c}_I = [0 \ 1]$



What is

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

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

$$= \frac{e^{\vec{v}_{saw} \cdot \vec{c}_{saw}}}{e^{\vec{v}_{saw} \cdot \vec{c}_{saw}} + e^{\vec{v}_{saw} \cdot \vec{c}_I}} = \frac{1}{1 + e^{-1}} \approx \frac{1}{4}$$

$= I = \frac{3}{4}$

② What are the values of  $\bar{c}$  that maximize likelihood?

See prev page

③ Why do we have  $\bar{v} \neq \bar{c}$ ?  
Why two spaces?

dot product of word w/self  
would be high

④ "we saw"  
if we add this, should  
get  $\bar{v}_{we} = \bar{v}_I$

## Other methods

FastText: each word embedding  
= sum of char n-grams

Problem with skip-gram:

For each example, how expensive  
is it to compute  $P(c|w)$ ?

$|V|$  vocab  $d$ -dim vectors

$O(|V|d)$  one evaluation

Corpus:  $|C| \cdot k$   
whole training =  $O(|V|d|C|k)$

Alternative:

Skip-gram with negative sampling  
(SGNS)

Take (word, context) pairs as  
"real" data

Sample "fake" data

Learn a binary classifier

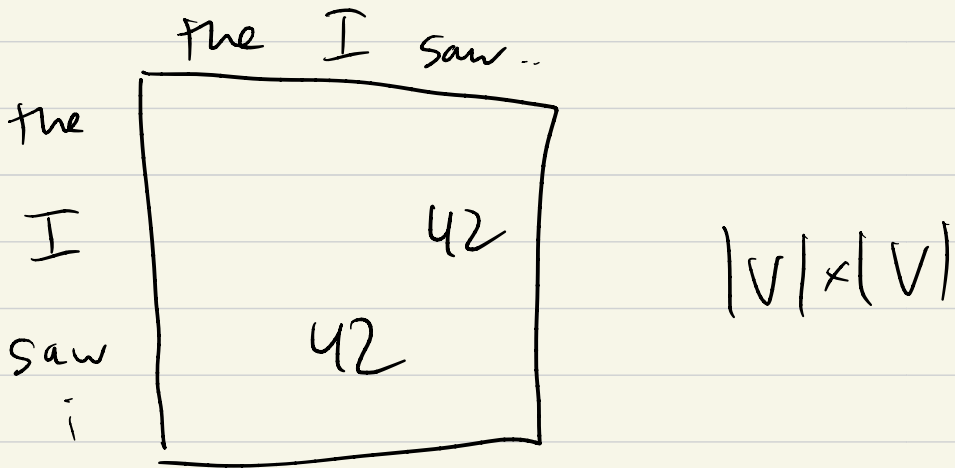
$$P(\text{Real} | x, y) = \frac{e^{\bar{v}_x \cdot \bar{c}_y}}{1 + e^{\bar{v}_x \cdot \bar{c}_y}}$$

(2014)

GloVe

Global Vectors

Factorizes a matrix of  
(word, context) counts



$SG \approx SGNS \approx GloVe$

GloVe has no dependence on  
corpus size