# CS371N Lecture 7
## Word Embeddings

## Skip-gram

Mikolov et al. 2013 "word2vec"

Learn <u>2</u> vectors for every word
  word vector
  Context vector
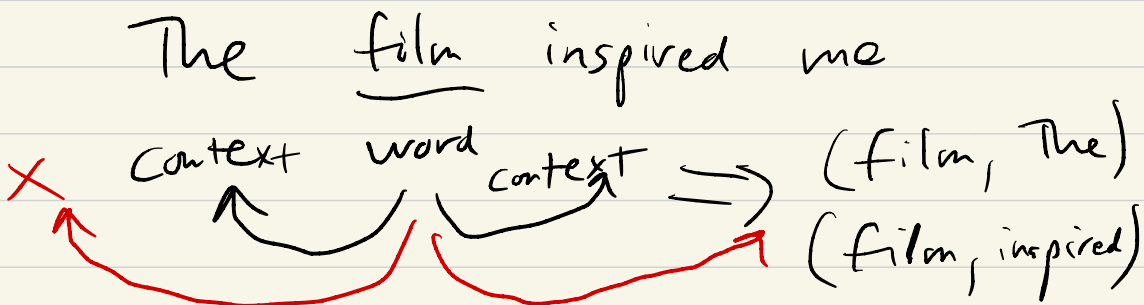
Try to predict context given word

Input: corpus of text

Outputs: $\vec{v}_w , \vec{c}_w$ for each word
  word    context    w in vocab $V$

(In A2: you are given just one vector)

Hyperparameters: $d$ dimension
window size $k$ $(50-300)$

Turn a sentence into (word, context) pairs

The <u>film</u> inspired me
context    word   context   $\implies$   (film, The)
(film, inspired)

$k=2$: Look 2 words away

(film, me)

Loop over words
from offset $\in \{-k, -k+1..-1, 1, ..k\}$
form pair (word, word + offset)

## Model (skip-gram)

$$P(\text{context} = \overset{\bar{c}_y}{y} \mid \text{word} = \overset{\nearrow \bar{v}_x}{x})$$

$$= \frac{e^{\bar{v}_x \cdot \bar{c}_y}}{\displaystyle\sum_{y' \in U} e^{\bar{v}_x \cdot \bar{c}_{y'}}}$$

distribution
over contexts

parameters: word vectors $\bar{v}$    $|U| \times d$

context vectors $\bar{c}$    $|U| \times d$

## randomly initialize

## Training   $(\overset{x}{\text{word}}, \overset{y}{\text{context}})$ examples

minimize $\displaystyle\sum_{(x,y)} -\log P(\text{context} = y \mid \text{word} = x)$
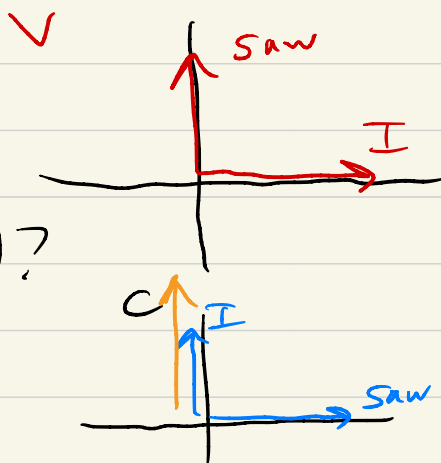
## Ex  Corpus = I saw  k=1

vocab = {I, saw}  d=2

Assume $\overline{V}_I = [1,0]$  $\overline{V}_{saw} = [0,1]$

① Let $\overline{C}_{saw} = [1,0]$

$\overline{C}_I = [0,1]$

what is $P(\text{context} \mid w = saw)$?

2 outcomes $(I, saw)$

$$P(I \mid saw) = \frac{e^{\overline{V}_{saw} \cdot \overline{C}_I}}{e^{\overline{V}_{saw} \cdot \overline{C}_I} + e^{\overline{V}_{saw} \cdot \overline{C}_{saw}}} = \frac{e}{e + 1}$$

$$\approx \frac{3}{4}$$

$$P(saw \mid saw) \approx \frac{1}{4}$$

② How to minimize loss further by changing $\overline{c}$?  $\overline{c}_I = [0, 10] \Rightarrow \frac{e^{10}}{e^{10} + 1}$

③ Why do we need two spaces?
Why $\bar{V} \neq \bar{C}$?

If one space: $P(saw | saw)$ has to
be high! $\overline{V}_{saw} \cdot \overline{V}_{saw}$

## Problems with skip-gram

Suppose we have a 100M word corpus
Vocab size $= 30k$      Vector dim $d = 300$

what's hard here?

$k=1$ : 200M pairs

Each $P(\cdot | \cdot) = O(|V| d)$

$200M \cdot O(|V| d)$

# fixes

① Skip-gram w/ negative sampling        (SGNS)

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

(word, ~ sampled context) as "fake" data

Learn classifier

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

SG: 30k denom.

SGNS: 1 positive + 10 sampled negs. = 11

## ② GloVe

Factorizes a matrix of (word, context)
counts

word

|       | the | I    | saw ... |
|-------|-----|------|---------|
| the   |     | 25   | 12      |
| I     | 25  |      | 1512    |
| saw   | 12  | 1512 |         |

Context

= M

matrix factorization

$$V^T \ C \quad = \ M$$

$(d \times |U|) \ (d \times |U|) \quad (|U| \times |U|)$

Gives the same solution as SG/SGNS