

Skip-gram

Mikolov et al. 2013 "word2vec"

Learn 2 vectors for every word

word vector

Context vector

Try to predict context given word

==

Inputs: a corpus of text

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

(for AZ: use \vec{v} , or $\vec{v} + \vec{c}$)

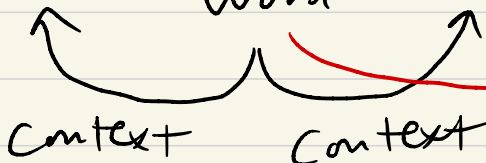
Hyperparameters: d (50 ~ 300)
window size k

Let $k=1$

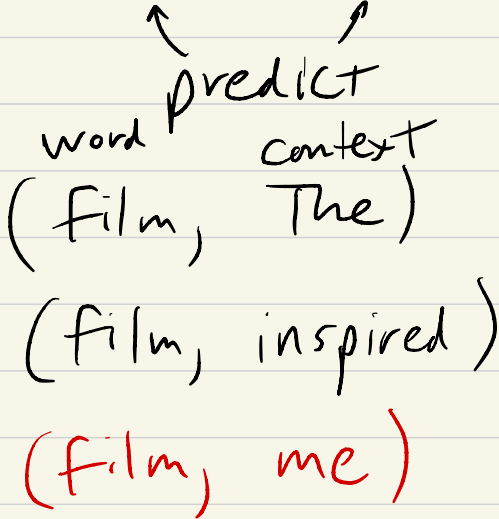
$k=2$

(skip-gram)
includes
 $k=1$

The film inspired me



context skip bigram



Training examples

Model (skip-gram) \bar{v}_x

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

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

distribution
over context
words in V

parameters: vectors \bar{V} $|V| \times d$

context vecs \bar{C} $|V| \times d$

randomly initialized

Training (x, y) train exs

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

Ex Corpus = I saw $k=1$

vocab = $\{I, \text{saw}\}$ $d=2$

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

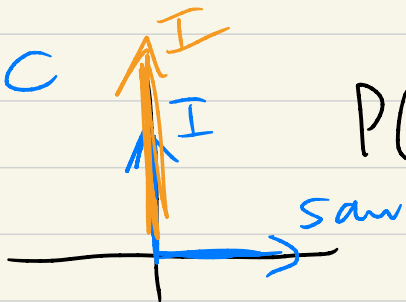


① Let $\bar{c}_{\text{saw}} = [1, 0]$

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

what is

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



2 outcomes (I, saw)

$$e^{\bar{v}_{\text{saw}} \cdot \bar{c}_I} = e \quad e^{\bar{v}_{\text{saw}} \cdot \bar{c}_{\text{saw}}} = 1$$

$$P(I | \text{saw}) = \frac{e}{e+1} \approx \frac{3}{4} \quad P(\text{saw} | \text{saw}) = \frac{1}{e+1} \approx \frac{1}{4}$$

(2) How to minimize loss further by changing \bar{c} ?

$$C_{\pm} = [0 \ 2] \Rightarrow C_{\pm} = [0 \ 10] \\ \frac{e^{10}}{e^{10} + 1} \approx 0.999$$

(3) Why is $\bar{v} \neq \bar{c}$? Why two spaces?

(saw, saw) always be high!

word vector selects for words that are near it

noun \rightarrow verbs

noun \nrightarrow nouns

Problems with skip-gram

If we ran this training over 100M word corpus with $V = 30K$, what's going to be hard?

- polysemy: different word senses
 - different vector per sense?
 - train on a homogeneous corpus
 - context dependent vectors (BERT, GPT)

- Computation: $|V|, d \quad 50 \sim 300$

$$P(y|x) = O(|V|^d)$$

For training: do that $\times 100M$

Two fixes

Skip-gram w/negative sampling

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

(word, \sim sampled-context) as fake
data

Learn a classifier

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

(film, buy) is this fake?

GloVe

Factorizes a matrix of
(word, context) counts ($k=1$)

the I saw ...

the	25	12	
I	25	1512	
saw	12	1512	
⋮			

$|V| \times |V|$

matrix factorization

$$(|V| \times v) \times (v \times |V|) \approx |V| \times |V|$$

Same as SG + SGNS