

→ have a mathematical relationship with sparse embedding like PPM1
→ can be seen as optimization of PPM1

Word2Vec

→ SGMIS

(Framework for learning word vectors)

→ embeddings: short dense vectors

→ work better

→ classifier requires to learn far fewer weights

→ smaller parameter space possibly helps with generalization and avoid overfitting

→ dense vector \Rightarrow capture synonyms better

→ fast, easier to train, static embeddings

→ simpler model

→ making binary prediction instead of probability

→ learns one fixed embedding for each word.

→ Instead of counting how often each word occurs, near a word w , train classifier to check how likely word w_1 will show up with w

→ Don't actually care about prediction task. Instead we take the learned classifier weights as the word embeddings

→ Use running text as implicitly supervised learning data (self supervision) → proposed in neural LM

→ Prediction (skip gram)

→ avoid the need of hard labeled supervision again.

1) Treat target and a neighbouring words as true example

2) Sample other words randomly to get -ve sample

3) Use logistic regression

4) Use learned weight as embedding

Sample with $P(w) = U(w)^{3/4} / Z \Rightarrow$ increasing the probability of less frequent words
 \uparrow
 unigram distribution.

$$P(+|w, c), P(-|w, c) = 1 - P(+|w, c)$$

$$\sigma(c \cdot w) \quad \sigma(-c \cdot w)$$

↳ makes the assumption, that all context words are independent

$$\log P(+|w, c, L) = \sum_{i=1}^L \log \sigma(c_i \cdot w) \leq \log \sigma(c \cdot w)$$

↳ store two embeddings for each word

↳ as a target (input embedding)

↳ as a content (output embedding)

↳ Learning skip-gram embeddings

→ Assigning random embedding vector for each of the N vocabulary words

→ iteratively shift the value to be more like embedding of words that occur nearby and less like that don't occur nearby

↳ SGNS uses more -ve examples than +ve's (K being the ratio)

↳ -ve example \Rightarrow noise words \Rightarrow chosen using $P(w)$

$$L = - \left[\log \sigma(c_{pos} \cdot w) + \sum_{i=1}^K \log \sigma(-c_{neg_i} \cdot w) \right]$$

Goal: 1) maximize similarity for (c_{pos}, w)

2) minimize similarity for (w, c_{neg})

t-SNE → projection method ⇒ visualization of embedding

- minimizing loss function using stochastic gradient descent
- tf-idf ⇒ choice of doc length, term frequency affects the performance
- III, word2vec ⇒ context window size affects the performance. tune it on dev set

Other kind of static embedding

↳ FastText ⇒ extension of word2vec ⇒ no global word index
deals with languages with morphology using subword model with unknown words

eg $n=3$ where

wh, whe, her, ere, re >

for each subword a embedding is learned
and the avg word will be represented as the sum of all the embeddings

→ based on vector of probabilities from word co-occurrence matrix while also capturing linear structure used by word2vec

↳ GloVe ⇒ Global vectors

→ It is best to learn multiple embeddings with bootstrap (randomly selecting subset) sampling over documents and average the results as embedding model suffers from inherent variability

Word2Vec: Objective function.

For each position $t=1, \dots, T$, predict context words within a window of fixed size m , given center word w_t

Data likelihood: $L(\theta) = \prod_{t=1}^T \prod_{-m \leq j \leq m, j \neq 0} P(w_{t+j} | w_t; \theta)$

Objective function

$$J(\theta) = -\frac{1}{T} \log L(\theta)$$

$$\begin{bmatrix} w_{t,1} \\ w_{t,2} \\ \vdots \\ w_{t,m} \end{bmatrix} \in \mathbb{R}^{2dm}$$

number of dimension of vectors

So that we don't have to deal with this enormous product

How to calculate $P(w_{t+j} | w_t; \theta)$?

Use two vector per word w :

$v_w \rightarrow$ as a center word

$u_w \rightarrow$ as a context word

$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$ \rightarrow calculating similarity using dot product

context word

center word

$\sum_{w \in V} \exp(u_w^T v_c)$

makes anything true

which is softmax function

because still assign some probability to smaller n_i

because amplifies probability of largest n_i

\rightarrow Different senses of words reside in a linear superposition (weighted sum) in standard word embeddings like word2vec

$$v_w = \alpha_1 v_{w_{s_1}} + \alpha_2 v_{w_{s_2}} + \alpha_3 v_{w_{s_3}} + \dots$$

$$\text{where } \alpha_i = \frac{s_i}{s_1 + s_2 + \dots}$$

Ideally after doing this we can't separate out the individual senses but using "sparse coding" we can do it