

CS5560 Knowledge Discovery and Management

Problem Set 4

June 26 (T), 2017

Name: SAIJYOTHI AUDIBANDI

Class ID: 7

I. N-Gram

Consider a mini-corpus of three sentences

<s> I am Sam </s>

<s> Sam I am </s>

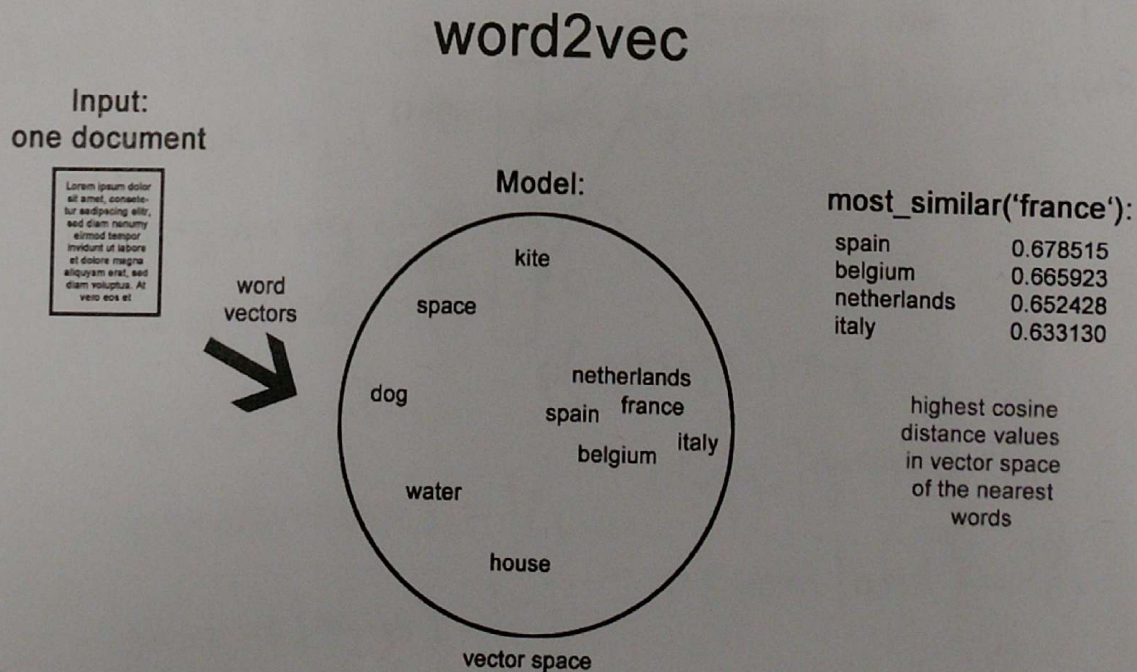
<s> I like green eggs and ham </s>

- 1) Compute the probability of sentence "I like green eggs and ham" using the appropriate bigram probabilities.
- 2) Compute the probability of sentence "I like green eggs and ham" using the appropriate trigram probabilities.

II. Word2Vec

Word2Vec reference: <https://blog.acolyer.org/2016/04/21/the-amazing-power-of-word-vectors/>

Consider the following figure showing the Word2Vec model.



- a. Describe the word2vec model

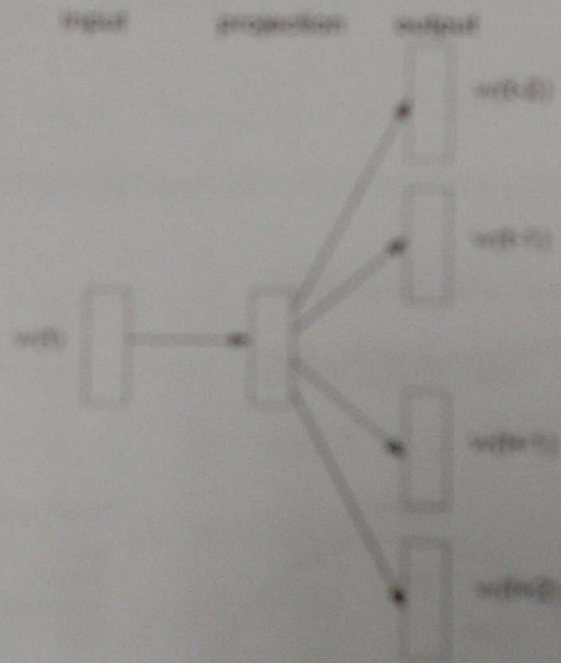
5. Describe how to extend the model for multiple documents. Also draw a similar diagram for the extended model.

Describe the differences of the following approaches:

- Continuous bag-of-words model
- Continuous skip-gram model

For the sentence "morning bag, afternoon light rain."

- Draw the words in the skip-gram Word2Vec model below.
- Draw a CBoW model using the same words.



I

N-Gram: A sequential list of the n words, often used in information retrieval and language modeling to encode the likelihood that the phrase will appear in the future.

N-Gram based approaches create applications probabilistic models of n -grams from a given corpus of text and tag new utterances using these models.

Given a min-corpus of 3 sentences

<S> I am sam </S>

<S> sam I am </S>

<S> I like green eggs and ham </S>

1. Calculating the bigram probability of sentence "I like green eggs and ham".

$$p(w_i / w_{i-1}) = \text{count}(w_{i-1}, w_i) / \text{count}(w_{i-1})$$

probability that word _{$i-1$} is followed by word _{i} = [Num times we saw word _{$i-1$} followed by word _{i}] / Num times we saw word _{$i-1$} .

s - beginning of sentence

/s - end of sentence.

$$p(I/s) = 2/3 \quad p(\text{like}/I) = \frac{1}{3}$$

$$p(\text{green}/\text{like}) = 1/1 \quad p(\text{eggs}/\text{green}) = 1/1$$

$$p(\text{and}/\text{eggs}) = 1/1 \quad p(\text{ham}/\text{and}) = 1/1$$

2. calculating the probability of sentence "I like green eggs and ham" using trigram probabilities.

$$P(w_i | w_{i-1} w_{i-2}) = \text{count}(w_i, w_{i-1}, w_{i-2}) / \text{count}(w_{i-1}, w_{i-2})$$

Probability that we saw word_{i-1} followed by word_{i-2} followed by word_i = [no. of times we saw the 3 word in order] / [no. of times we saw word_{i-1} followed by word_{i-2}]

$$P(\text{green} | \text{I like}) = \text{count}(\text{green I like}) / \text{count}(\text{I like}) = 0 / 1 = 0.$$

$$P(\text{eggs} | \text{like green}) = \text{count}(\text{eggs like green}) / \text{count}(\text{like green}) = 0$$

$$P(\text{and} | \text{green eggs}) = \text{count}(\text{and green eggs}) / \text{count}(\text{green eggs}) = 0$$

$$P(\text{ham} | \text{eggs and}) = \text{count}(\text{ham eggs and}) / \text{count}(\text{eggs and}) = 0.$$

II

(a) word 2vec model

A two layer neural net that process text.

- Input is a text corpus.

- output is a set of vectors.

- Not a deep neural net, but numerical form that deep net can understand.

measuring cosine similarity

- No similarity is expressed as a 90 degree angle.

- Total similarity of 1 is a 0 degree angle.

$$\text{similarity} = \cos(\theta) = \frac{A \cdot B}{|A| |B|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$

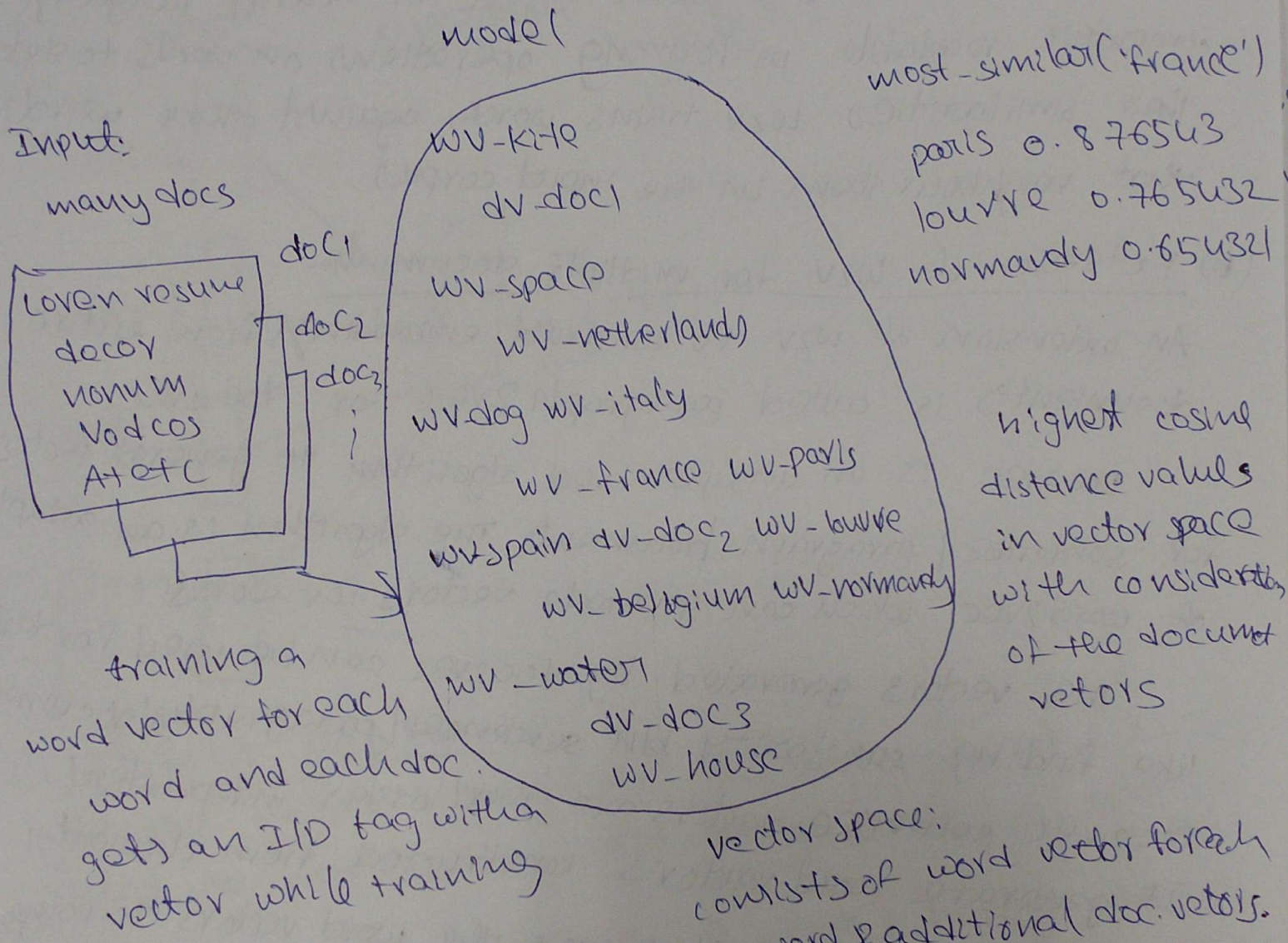
word2vector "vectorizes" about words for neural language computer readable performing operations on words to detect their similarities. w2v trains words against other words that neighbour them in the input corpus.

(b) Extension of w2v for multiple documents

An extension of w2v to construct embeddings from entire documents is called paragraph2vec or doc2vec.

Doc2vec is an unsupervised algorithm to generate vectors for sentence / paragraph documents. The algorithm is an adaption of word2vec which can generate vectors for words.

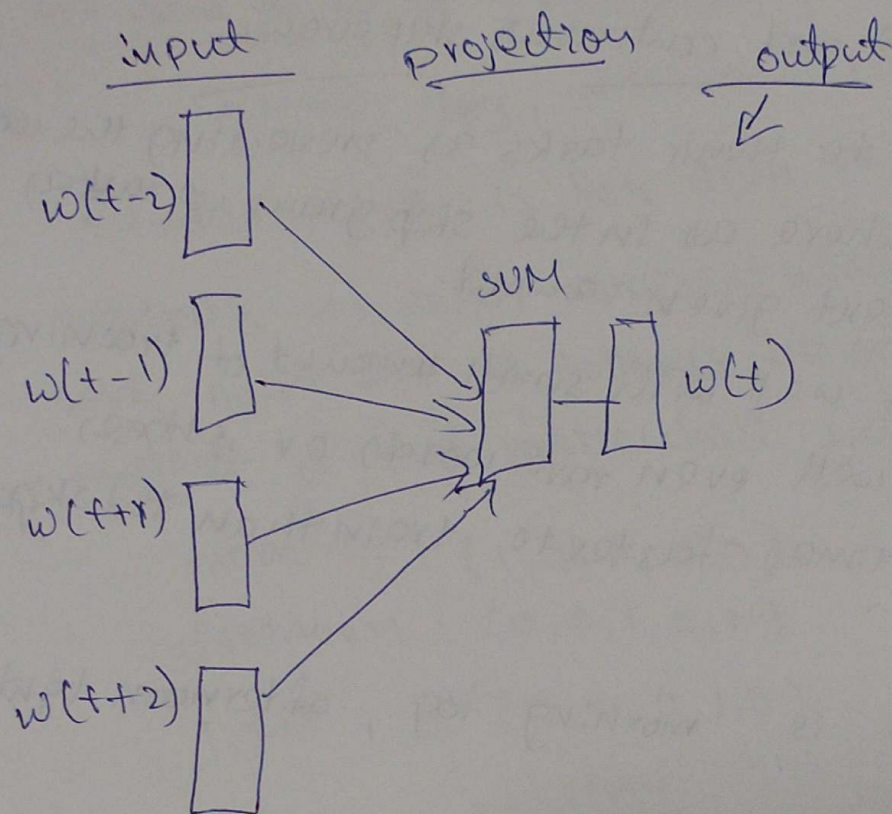
The vectors generated by doc2vec can be used for tasks like finding similarity b/w sentences / paragraphs / documents. Doc2vec sentence vectors are word order independent. It generate word vectors constructed from character ngrams and then adding up to the word vectors to compute a sentence vector. It generated by predicting the adjacent sentences, that are assumed to be semantically related.



w2v can utilize either of two model architectures to produce a distributed representation of words.

(a) continuous bag of words:- In the continuous bag of words architecture, the model predicts the current word from a window of surrounding context words. The order of context words does not influence prediction.

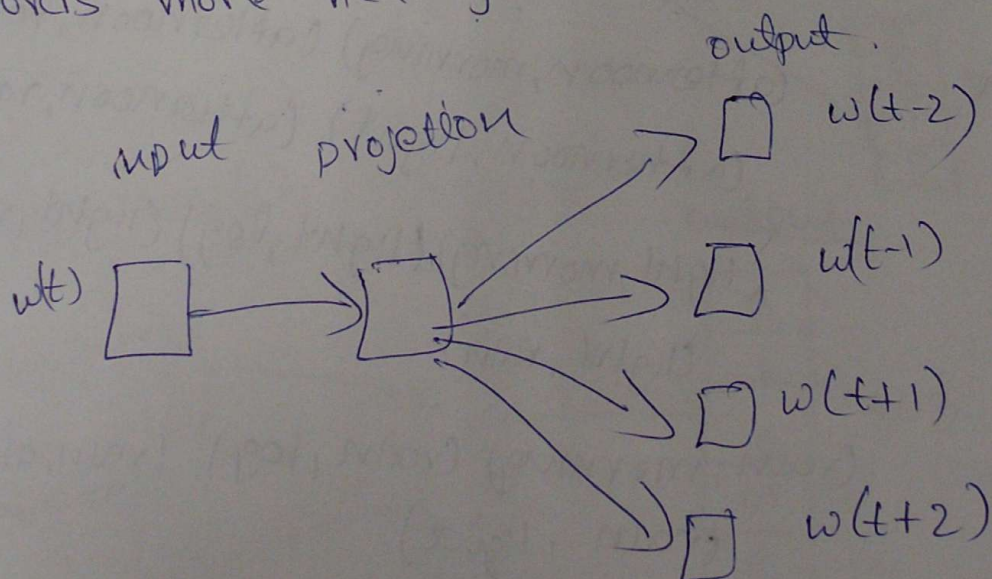
(b)



(b) continuous - skip-gram

In the continuous skip-gram, the model uses the current word to predict the surrounding window of context words.

The skip-gram architecture weighs nearby context words more heavily than more distant words.



Differences b/w CBOW and continuous skipgram

1. In CBOW we need to think tasks as 'predicting the word given its context' where as in the skip gram we task as predicting the context given a word.
2. skip-gram works well with small amount of training data, represents well even rare words or phrases.
3. CBOW is several times faster to train than the skip gram.

→ Given the sentence is 'morning fog, afternoon light rain'

skip-gram:-

consider window size as '1'.

input

morning

fog

afternoon

light

rain

Training samples

(morning, fog), (morning, afternoon)

(fog, morning) (fog, afternoon) (fog, light)

(afternoon, morning) (afternoon, fog)

(afternoon, light) (afternoon, rain)

(light, morning) (light, fog) (light, afternoon)

(light, rain)

(rain, morning) (rain, fog) (rain, afternoon)

(rain, light)

We need to build a vocabulary of words:-

(words), fog, afternoon, light, rain)

consider, input is fog then vector representation is (0,1,0,0,0)
 similarly the vector representation for morning, afternoon and light are as follows because there are in the context of that particular input word.

morning : (1,0,0,0,0)

afternoon : (0,0,1,0,0)

light : (0,0,0,1,0)

