# CS5560 Knowledge Discovery and Management
## Problem Set 4
### June 26 (T), 2017

Name: SUJITHA PUTHANA

Class ID: 24

## I. N-Gram

Consider a mini-corpus of three sentences
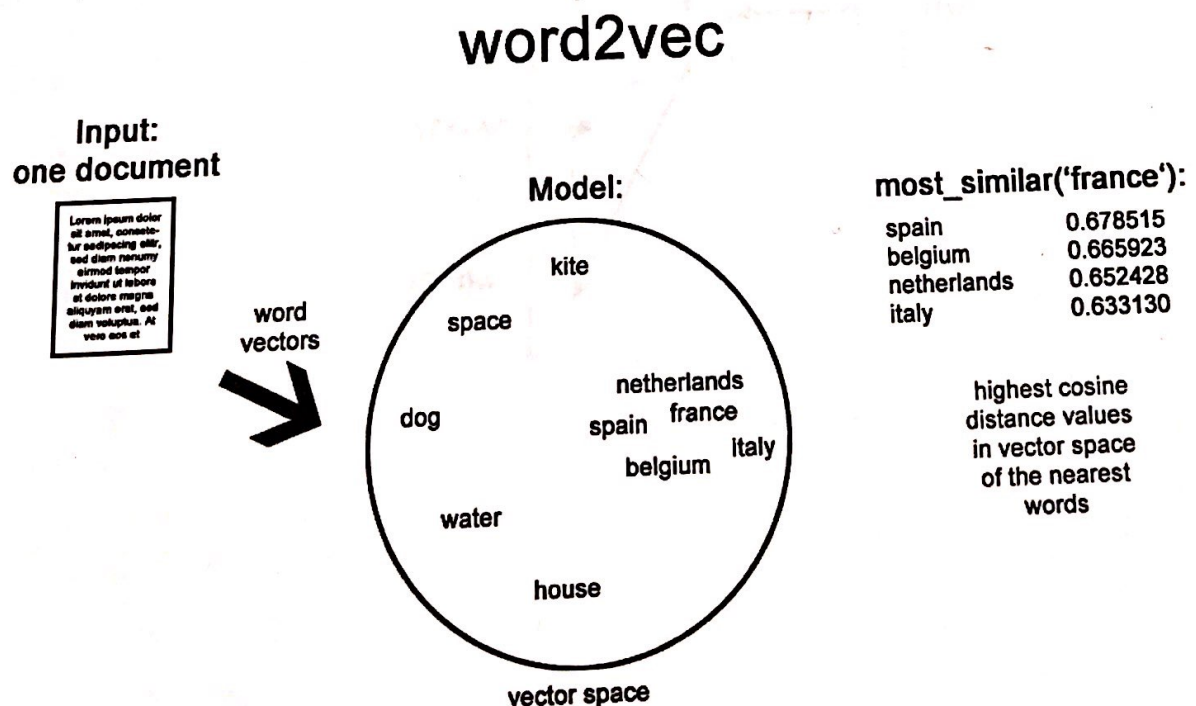
\<s\> I am Sam \</s\>
\<s\> Sam I am \</s\>
\<s\> I like green eggs and ham \</s\>

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.
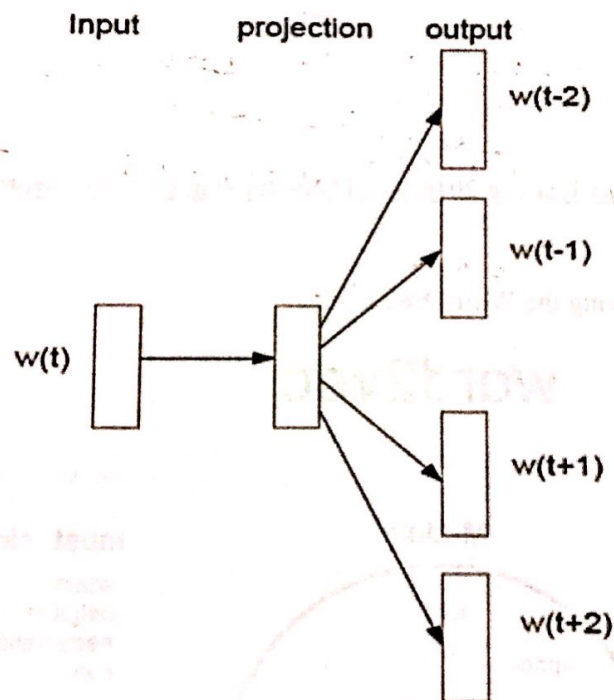


a. Describe the word2vec model

b. Describe How to extend this 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 fog, afternoon light rain,"

- Place the words on the skip-gram Word2Vec model below.
- Draw a CBOW model using the same words.

Given Sentences

① $\angle s >$ I am Sam $\angle /s>$

② $\angle s>$ Sam I am $\angle /s>$

③ $\angle s>$ I like green eggs and ham $\angle /s>$

Bigram Probability
$$P(w_i|w_{i-1}) = \frac{C(w_{i-1}, w_i)}{C(w_{i-1})}$$

Completes one round calculation of $\angle s>$①

$P(I/\angle s>) = \frac{2}{3} = 0.67$ [Here the $(I/\angle s>)$ exist is ① & ③ sentence and, total available sentence are 3
$I \rightarrow$ ①, ②, ③

$P(am|I) = \frac{2}{3} = 0.67$ [①, ②] $w(I) = 3$

$P(sam|am) = \frac{1}{2} = 0.5$ [①] [$w(am) = 2$]

$P(\angle /s>|sam) = \frac{1}{2} = 0.5$ [①] [$w(sam) = 2$]

$\angle s>$②

$P(sam|\angle s>) = \frac{1}{3} = 0.33$ [②] [$w(\angle s>) = 3$]

$P(I|sam) = \frac{1}{2} = 0.5$

$P(\angle /s>|am) = \frac{1}{2} = 0.5$

$\angle s>$③

$P(like|I) = \frac{1}{1} = 1$

$P(green|like) = \frac{1}{1} = 1$

$P(eggs|green) = \frac{1}{1} = 1$

$P(and|eggs) = \frac{1}{1} = 1$

$P(ham|and) = \frac{1}{1} = 1$

$P(\angle s>|ham) = \frac{1}{1} = 1$

④

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

$$P(I | <S>) * P(like | I) * P(green | like) * P(eggs | green)$$
$$* P(and | eggs) * P(ham | and) * P(</S> | ham)$$

$$= 0.67 * 1 * 1 * 1 * 1 * 1 * 1 = \boxed{0.67}$$

| Bigram Probability = 0.67 |

b) Trigram Probability :-

$$\boxed{P(w_i | w_{i-1}, w_{i-2}) = \frac{C(w_i, w_{i-1}, w_{i-2})}{C(w_{i-1}, w_{i-2})}}$$

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

$$P(like | <S> I) * P(green | I like) * P(eggs | like green)$$
$$* P(and | green eggs) * P(ham | eggs and) * P(</S> | and ham)$$

$$P(like | <S> I) = \frac{1}{2} = 0.5$$

$$P(green | I like) = \frac{1}{1} = 1$$

$$P(eggs | like, green) = \frac{1}{1} = 1$$

$$P(and | green eggs) = \frac{1}{1} = 1$$

$$P(ham | egg and) = \frac{1}{1} = 1$$

$$P(</S> | and ham) = \frac{1}{1} = 1$$

$$\Rightarrow 0.5 * 1 * 1 * 1 * 1 * 1$$

$$= 0.5$$

---

| Trigram Probability = 0.5 |

---

II)

a) <u>Word2Vec Model</u>

→ The input of the word2vec model is a large document and for each word in the document, a vector is built.

→ With all the word vectors, we have vector space which is the model of word2vec.

→ We got the most similar words by calculating the cosine distinance i.e., similarity.
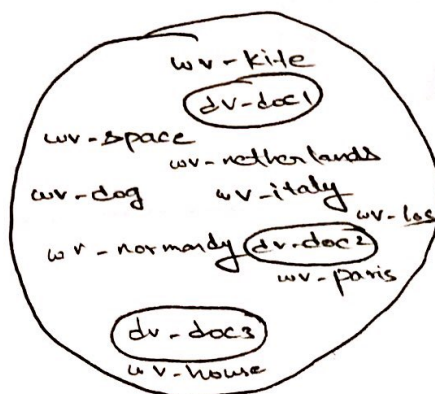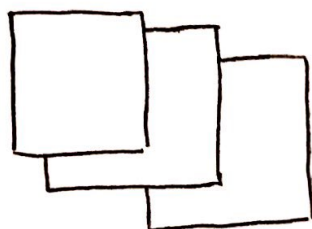
b) <u>Doc2 Vec</u>

→ Doc2Vec, an unsupervised algorithm to generate vectors for sentence, paragraphs, documents.

→ This algorithm is an adaptation of word2vec which can generate vector of words.

→ The vectors generated by doc2 vec can be used for tasks like finding similarity between sentences, paragraphs, documents.

Input
Multiple
documents

wv - kite
dv - doc1
wv - space
wv - netherlands
wv - dog    wv - italy
wv - les
wv - normandy   dv - doc2
wv - paris
dv - doc3
wv - house

Vector space

most_similar ('France')

paris      0.87654

kite       0.7654

dog        0.65

→ Vector space consists of word vectors for each word and additional document vectors.

| Bag of words | Skip gram models |
|---|---|
| ① Input – $w_{i-2}, w_{i-1}, w_{i+1}, w_{i+2}$ <br> Output – $w_i$ | ① Input – $w_i$ <br> Output – $w_{i-1}, w_{i-2}, w_{i+1}, w_{i+2}$ |
| ② Predicts the word given its context | ② Predicts the context given a word |
| ③ Faster to train than skip-gram. Better accuracy for the frequent words. | ③ For even small amount of training data, it gives well even rare words & phrases. |

©

d) **SkipGram** :- | Predicts Content of the given Words |

Output Content

Input word



w(t) | afternoon | → | Projection — afternoon | → morning
→ fog
→ light
→ rain

**Continuous Bag of Words** :- | Predicts Words for given context |

Input Content

| Morning |

| Fog |

| Light |

| Rain |

→ Projection → Output Word

afternoon