Name: Megha Nagabhushan

Class ID: 15

## I. N-Gram

Consider a mini-corpus of three sentences
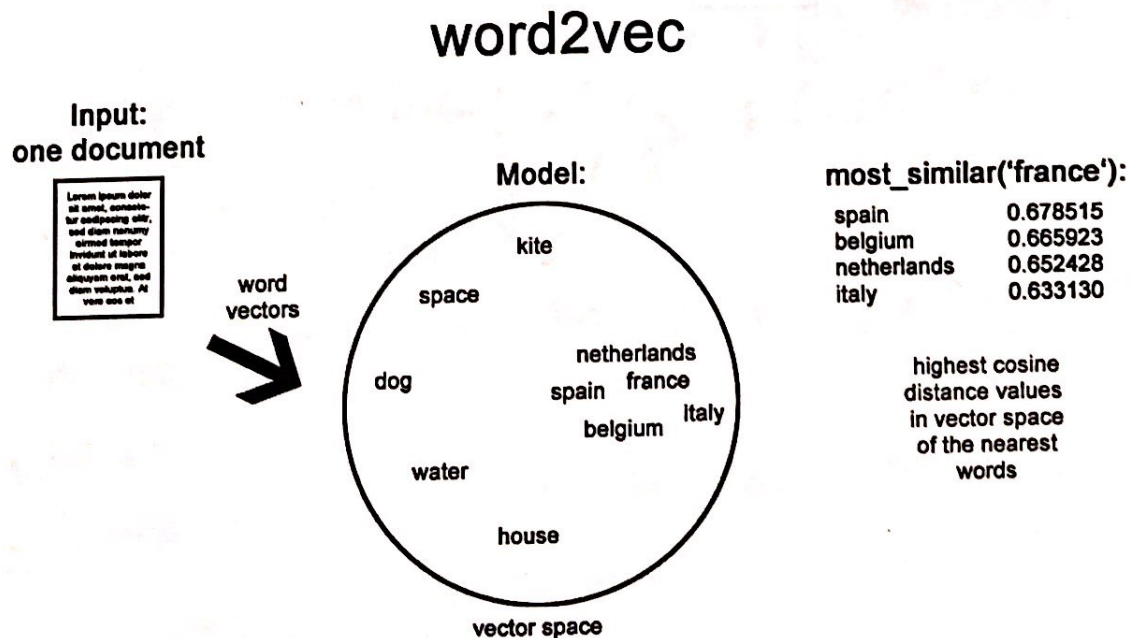
<s> I am Sam </s>

<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.
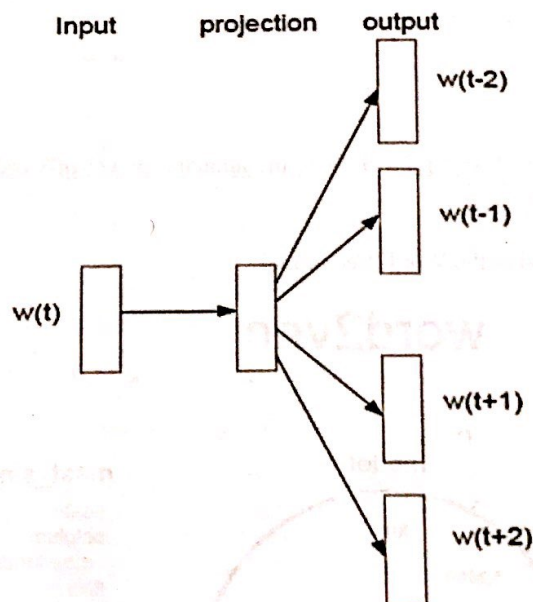


a. Describe the word2vec model

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.

① a) Calculations of bigram probabilities from this corpus.

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

Training corpus:

&lt;s&gt; I am Sam &lt;/s&gt;

&lt;s&gt; Sam I am &lt;/s&gt;

&lt;s&gt; I like green eggs and ham &lt;/s&gt;

$P(I \mid \langle s \rangle) = \frac{2}{3} = 0.67$

$P(am \mid I) = \frac{2}{3} = 0.67$

$P(Sam \mid am) = \frac{1}{2} = 0.5$

$P(\langle /s \rangle \mid Sam) = \frac{1}{2} = 0.5$

$P(\langle /s \rangle \mid am) = \frac{1}{2} = 0.5$

$P(I \mid Sam) = \frac{1}{2} = 0.5$

$P(Sam \mid \langle s \rangle) = \frac{1}{3} = 0.33$

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

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

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

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

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

$P(\langle /s \rangle \mid ham) = \frac{1}{1} = 1$

Bi-gram probability for " I like green eggs and ham"

will be

$P(I \mid \langle s \rangle) \times P(like \mid I) \times P(green \mid like) \times P(eggs \mid green) \times$
$P(and \mid eggs) \times P(and \mid eggs) \times P(ham \mid and) \times$
$P(\langle /s \rangle \mid ham).$

$0.67 \times 1 \times 1 \times 1 \times 1 \times 1 \times 1 = \underline{0.67}$

**b)** Trigram probability.

$$P(W_i \mid W_{i-1}, W_{i-2}) = \frac{\text{count}(W_i, W_{i-1}, W_{i-2})}{\text{count}(W_{i-1}, W_{i-2})}$$

$$P(Sam \mid I\ am) = \frac{\text{count}(Sam\ I\ am)}{\text{count}(I\ am)} = \frac{1}{2}$$

$$P(\langle s \rangle \mid am\ Sam) = \frac{1}{1} = 1.$$

$$P(am \mid \langle s \rangle\ I) = \frac{1}{2} \qquad P(am \mid Sam\ I) = \frac{1}{1} = 1 \qquad P(\langle /s \rangle \mid I\ am) = \frac{1}{2}$$

$$P(I \mid \langle s \rangle\ Sam) = \frac{1}{1} = 1$$

$$P(like \mid \langle s \rangle\ I) = \frac{1}{2} \qquad P(green \mid I\ like) = 1 \qquad P(eggs \mid like\ green) = 1$$

$$P(and \mid green\ eggs) = 1 \qquad P(ham \mid eggs\ and) = 1 \qquad P(\langle /s \rangle \mid and\ ham) = 1.$$

Tri-gram probability for

$\langle s \rangle$ " I like green eggs and ham" $\langle /s \rangle$ will be

$$P(like \mid \langle s \rangle\ I) \times P(green \mid I\ like) \times P(eggs \mid like\ green) \times$$

$$P(and \mid green\ eggs) \times P(ham \mid eggs\ and) \times P(\langle /s \rangle \mid and\ ham)$$

$$= \frac{1}{2} \times 1 \times 1 \times 1 \times 1 \times 1 = \frac{1}{2} = \underline{0.5}.$$
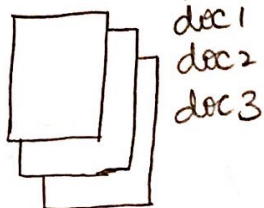
# II.

a) In the Word2Vec model, the input is a large document and for each word in the document, a vector is built. With all the word vectors you have vector space which is the model of word2vec. By calculating the cosine distance (similarity) between word vectors you get the most similar words you looked for a word.
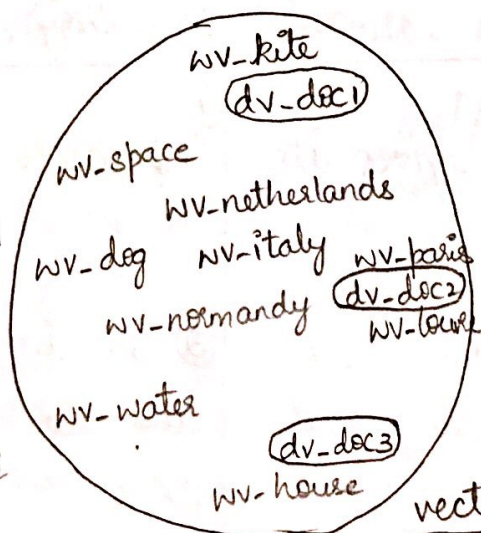
b)

## DOC 2 VEC.

**b)**

Input: many documents

doc1
doc2
doc3

training a word vector for each word and each document gets an ID/tag with a vector while training.

wv-kite
(dv-doc1)

wv-space
wv-netherlands
wv-dog    wv-italy    wv-paris
wv-normandy   (dv-doc2)
wv-louvre
wv-water
(dv-doc3)
wv-house

most_similar ('France'):

| paris | 0.87654: |
| louvre | 0.765432 |
| normandy | 0.65432 |
| ... | |

highest cosine distance values in vector space with consideration of the document vectors.

vector space:

consists of word vectors for each word and additional document vectors.

Word vectors generated by the neural net have nice semantic and syntactic behaviors. We need a clear way to combine them into high quality document vectors.

So here we can use Doc2vec, an unsupervised algorithm to generate vectors for sentence/paragraphs/documents. The algorithm is an adaptation of word2vec which can generate vectors for words.

The vectors generated by doc2vec can be used for tasks like finding similarity between sentences/paragraphs/documents.

---
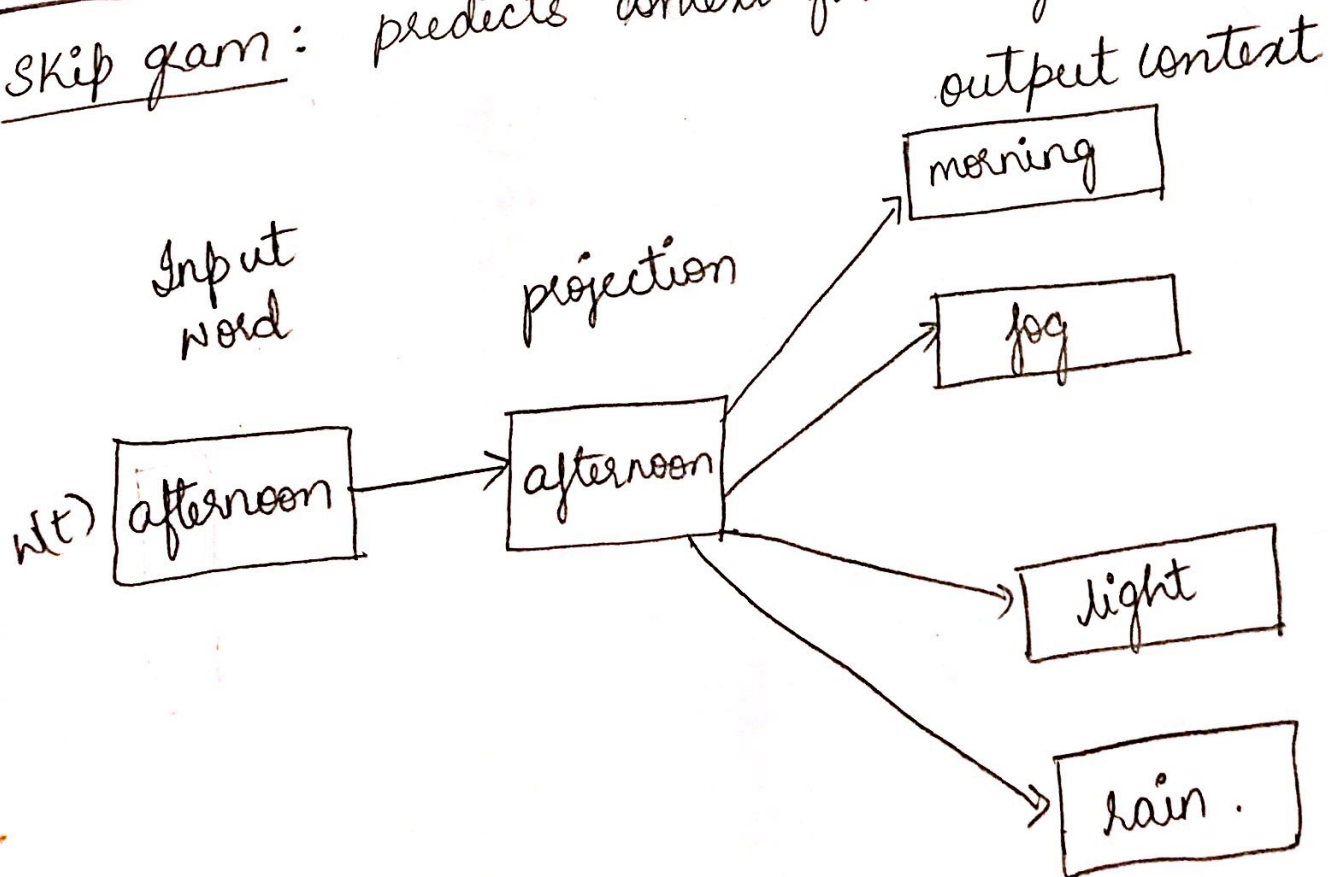
## Bag of Words model vs Skip gram model.

| Bag of Words model | Skip gram model |
|---|---|
| ① "predicts the word given its context" | ① "predicts the context given a word". |
| ② The input could be $W_{i-2}$, $W_{i-1}$, $W_{i+1}$, $W_{i+2}$, the preceding and following words of the current word we are at. The output of the neural network will be $W_i$. | ② The input to the model is $W_i$ and the output could be $W_{i-1}$, $W_{i-2}$, $W_{i+1}$, $W_{i+2}$. |
| ③ ~~several times~~ | |

③ several times faster to train than the skip-gram and better accuracy for the frequent words.

③ words well with small amount of training data, represents well even rare words or phrases.

④ Ex: "Hi fred, how was the pizza?"

Bag: of words - 3 grams.

"hi fred how"  "fred how was"

"how was the" - - - - - -

Skip - gram - 3 grams

"Hi fred how", "Hi fred was"

"fred how was", "fred how t

Here, "Hi fred was" skips

"how".

---

Skip gram : predicts context for the given word.



output context

# Continuous Bag of Words : predicts word for a given context.

Input
Context

output word
for the
context

```
┌──────────┐
│ morning  │─────┐
└──────────┘     │
┌──────────┐     │
│   fog    │─────┤        ┌──────────────┐              ┌────────────┐
└──────────┘     ├───────▶│  projection  │─────────────▶│ afternoon. │
┌──────────┐     │        └──────────────┘              └────────────┘
│          │─────┤
└──────────┘     │
┌──────────┐     │
│  light   │─────┤
└──────────┘     │
┌──────────┐     │
│   rain   │─────┘
└──────────┘
```