NAME :- NAGESWARA RAO          CLASS ID - 18

1-①  **Bigram** **probabilities** of the _3rd sentem_

Calculating bigram probability:

$$P(w_i/w_{i-1}) = count(w_{i-1}, w_i) / count(w_{i-1})$$

In english

Probability that word$_{i-1}$ is followed by word$_i$;

$$= \frac{[num\ tiny\ we\ saw\ word_{i-1}\ follows\ by\ word_i]}{[num\ tiny\ we\ saw\ word_{i-1}]}$$

S = begining of sentence

1s = end of sentence

$$P(I/s) = 2/3$$
$$P(like/I) = 1/3$$
$$P(green/like) = 1/1 = 1$$
$$P(eggs/gree) = 1/1 = 1$$
$$P(and/egg>) = 1/1 = 1$$
$$P(1s/ham) = 1/1 = 1$$

② **Trigram** **probability** of the _3rd sentence_

Calculating trigram probability

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

In english :

Probability that we saw word$_{i-1}$ followed by word$_{i-2}$

followed by word$_i$;

$$= \frac{num\ tiny\ we\ saw\ the\ 3\ wong\ moray}{num\ tiny\ we\ saw\ word_{i-1}\ follow\ by\ word_{i-2}}$$

$$P(\text{green}/\text{i like}) = \text{count}(\text{green i like}) / \text{count}(\text{i like})$$

$$= 0/1 = 0$$

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

$$= 0/1 = 0$$

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

$$= 0/1 = 0$$

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

$$= 0/1 = 0$$

(11)    WORD2VEC

(a)   word2vec model.

It is a two-layer neural network that process the text, input is a text corpus. output is a set of vectors. feature vector for word in that corpus. Not a deep neural network. but a numerical form that deep nety can understand. No similarity is expressed as a 98 angle, total similarity OG 1 is a dgree angle

$$\text{similarity} = \cos(\theta) = \frac{AB}{\|A\|\|B\|} = \frac{\sum_{i=1}^{n} A_i B_i}{\sqrt{\sum_{i=1}^{n} A_i^2} \sqrt{\sum_{i=1}^{n} B_i^2}}$$
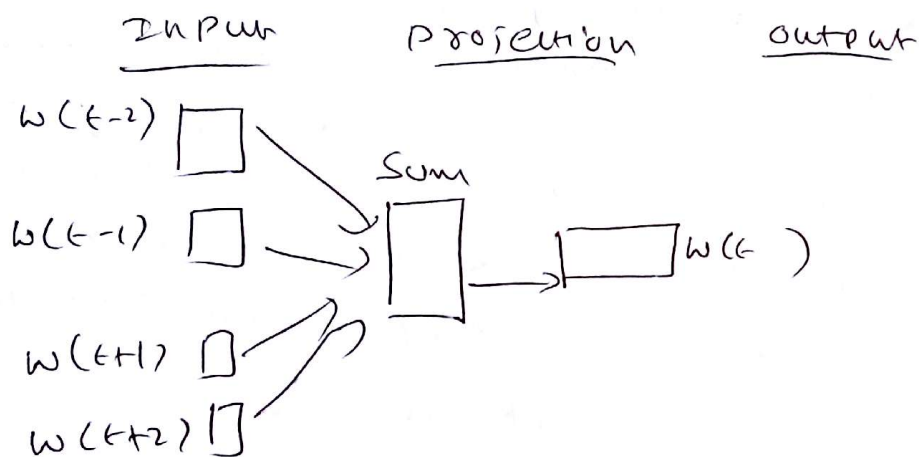
In the model represented in the diagram they have taken a input document and built a word2vec model contuiy word in the document and found the nearest word using cosine similarity

word2vec can utilize either of two model architecture to produce a distributed representation of word:
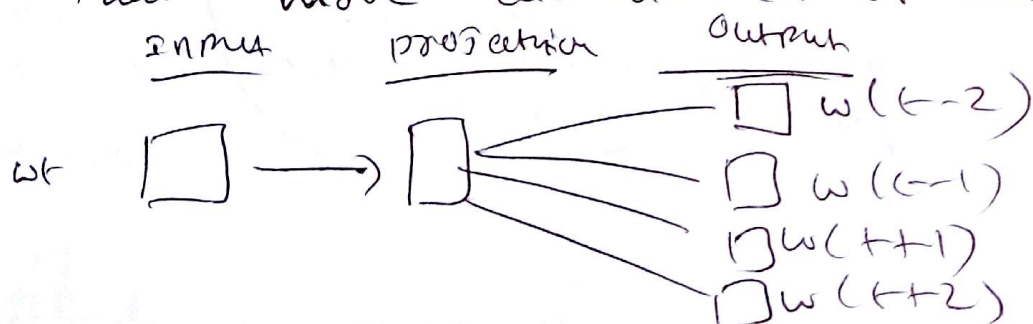
(a) Continous bag of word

(b) Contiouy skip-gram.

## CBOW

In the continouy bag of word architecture, the model predicts the current word from a window of surrounding content words. The order of context words does not influence prediction

Input          Projection          Output



## Continouy skip-gram:

In the Continouy Skin-gram architecture, the model uses the current word to predict the surrounding window of context-words. The skip-gram architecture weighs nearby by content word more heavily then move distant content word
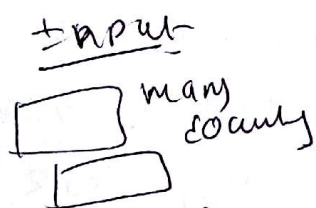
Input          projection          Output

5) Extension of wordzvec for multiple document

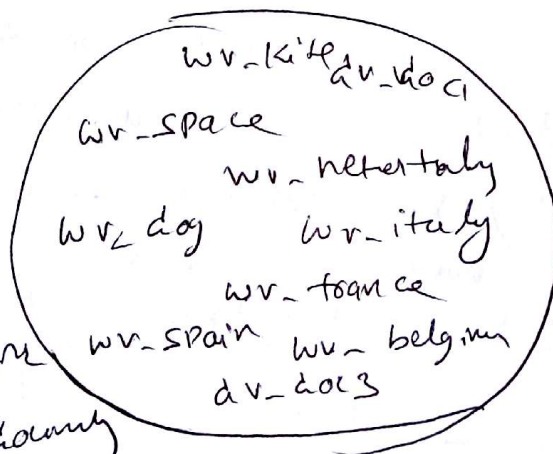An extension of wordzvec to construct Embeddings from entire document is called Paragraph 2vec or docz2vec.

Docz2vec is an unsupervised algorithm to generate vectory for sentence/paragraphs/document. The algorithm is an adaption of wordzvec which can generate vectory for words.

The vectory generated by docz2vec can be used for tasks like finding similarity between sentences/paragraphs/documents doczvec sentence vectory are word order independent. It generate word vectory constructed from character n-grams and then adding up the word vectory to compose a sentence vector. It generate vectory where the vector for a sentence is generated by predicting the adjacent sentency, that are assumed to be semantically related.

DOC2VEC for diagram

Input

many documy

tracining word vector for each word and each documy

wv-kite av-doci
wv-space
wv-netertaly
wv-dog   wv-italy
wv-france
wv-spain  wv-belgim
dv-docz

most-similu ('france')
Paris 0.876543
louvre 0.76543

nostra

highest cosine distance valy in vector space with consideration of the docum

# Difference between CBow and Continous Skip gram

① In CBow we need to think task as "predicting the word given ity content" where as in the skip-gram we think task as "predicting the content given a word".

② Skip-gram works well with small document or few training data, represent well even rare words or phrases.

③ CBow is a several tiny faster to train than the skip-gram. Slightly better accuracy for the frequent words.

Given the sentence is "morning fog, afternoon light rain"

Skip-gram word2vec model for above sentence is

Consider the sentence:-

morning fog, afternoon light rain,
Consider window size is 1

**Input**

morning

fog

afternoon

light

rain

**Training samples**

(morning, fog), (morning, afternoon)

(fog, morning) (fog, afternoon)

(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
(morning, fog, afternoon, light, rain)
consider input is fog then vector representation is
(0, 1, 00, 0)
similarity the vector representation for morning,
afternoon and light are as follow, because
these are in the context of that particular
input word.

morning :- (1, 00, 90)
afternoon :- (0, 0, 1, 0, 0)
light = (0, 0, 0, 1, 0)



Input    Projection    output

w(t)
fog

CBOW model

morning

afternoon

light

Input

hidden    output    fog

w(t-1) morning

w(t+1) afternoon

w(t+2) (light)