

NOTES

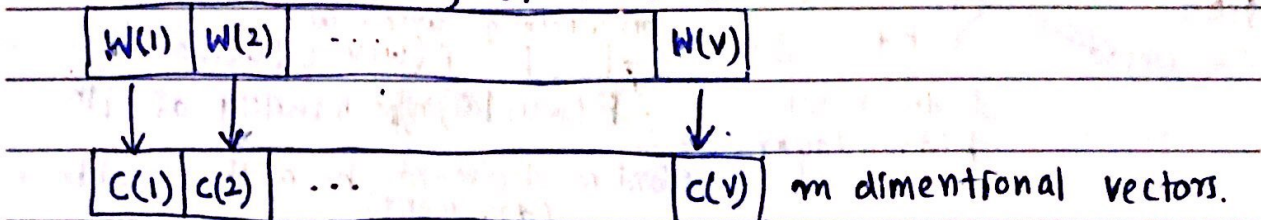
Natural language processing - Take natural texts and try to make predictions based on them.

Words in a vocabulary are mapped into multidimensional vectors (word embedding / word to vector).

Each vector represents a point in that multidimensional map.

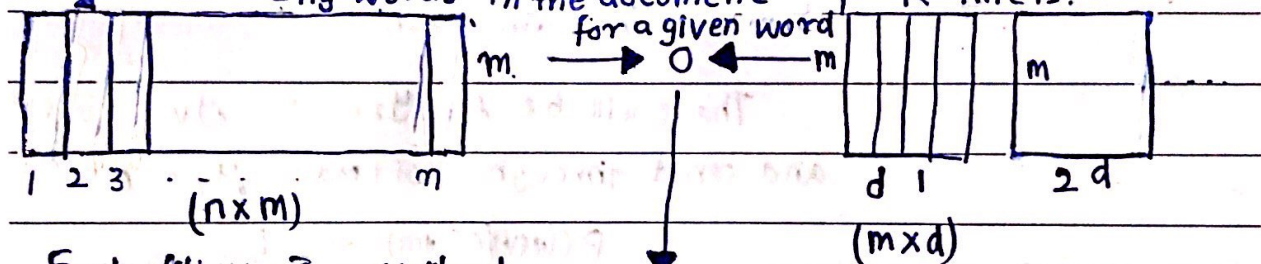
Same category words have nearby points. (feature space)

Vocabulary of V no of words.



Vocabulary codebook. (This will also be learned as a part of parameters)

Vectorized text document. A good word-vec system can predict surrounding words in the document. One filter has d no of words.

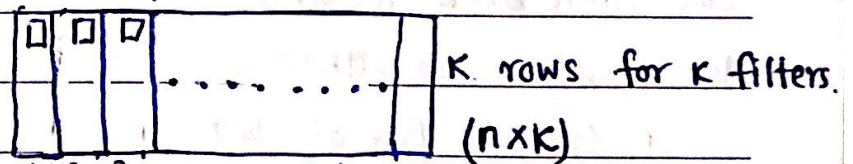


Each filter is convolved

and check how much the each words matches with the

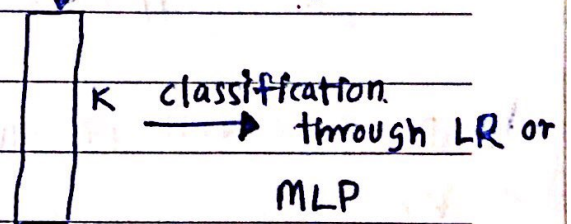
words in the filter.

A correlation map is created.



max pooling

Max value from each row will be taken. (Highest match with the filter)



For parameter learning in NLP we don't use a true label.

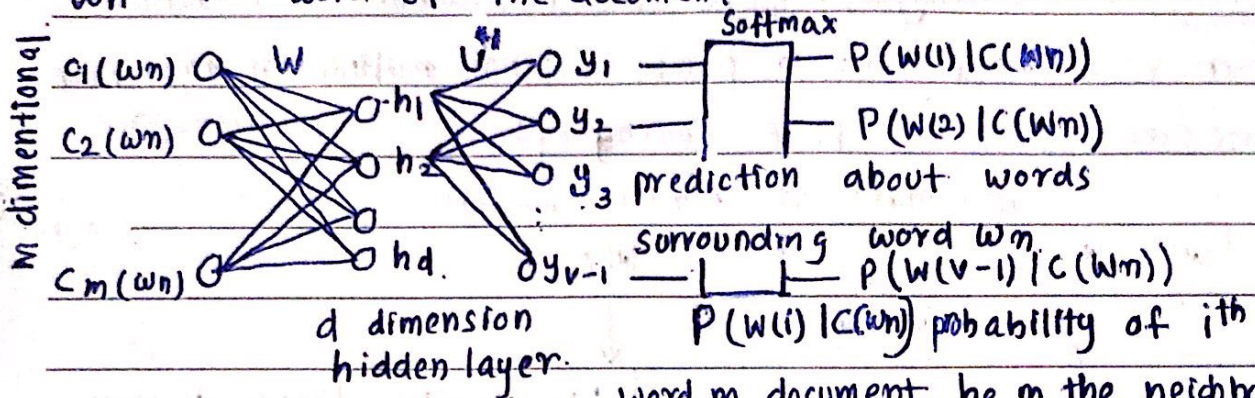
∴ In order to create a true label for the text a human will have to read the whole document. ∴ costly.

NOTES

Vocabulary has V no of words.

w_n - n^{th} word of the document.

Neural text model



$$h = W \odot C + b$$

word m document be in the neighbour hood of w_n . n^{th} word in document

$$y = U \odot h + b_1$$

↓ Softmax

To take the probability h is send through another set of weights. U with $v-1$ dimension.

There will be y_1, y_2, \dots, y_{v-1} outputs

and send through softmax. $y_v = 0$

Since there are no label

$$P(w(v)|C(w_n)) = \frac{1}{1 + \sum_{i=1}^{v-1} \exp(y_i)}$$

data we can use either

1. Continuous Bag of words.

= 1 - sum of $(v-1)$ prob

$$C(w_n) = \text{Avg.} \sum_{\substack{i=-a \\ i \neq 0}}^a C(w_{n-i}) \quad (a \in 1, 2, 3, \dots)$$

Sum of neighbouring words' vectors except w_n .

$P(w(i)|C(w_n))$ will give the probability of i^{th} word being w_n

or 2. Skip-Gram model

Input is $C(w_n)$. The highest probabilities in the output indicates words in vocabulary that are in the document neighbourhood of w_n

to learn the parameters W, U , biases and word vectors C .

NOTES

How to get most suitable parameters? (cost function should be highest)

$$\text{Cost func} = - \sum_{i=1}^M \log \text{probability for } i^{\text{th}} \text{ input word.}$$

function of out parameters.
training.

M - depends on no of documents we have

Take the negative value of the
cost function value and perform SGD