

Word 2 Vec in C ++

documentation

How it works

Step -1 words are converted into one hot vectors

Eg : the - 0, 0, 0, 0, 1, 0, 0.

the size of the one hot vector depends on the number of unique words in the corpus (data set).

Step -2 the one hot vectors are given to a neural network

Step -3

A weight matrix W_o is initialized
with random weights

Step -4

the one hot vector is multiplied
with the weighted matrix

Step -5

the output matrix is run
through a softmax layer and
the output is taken

Step -6

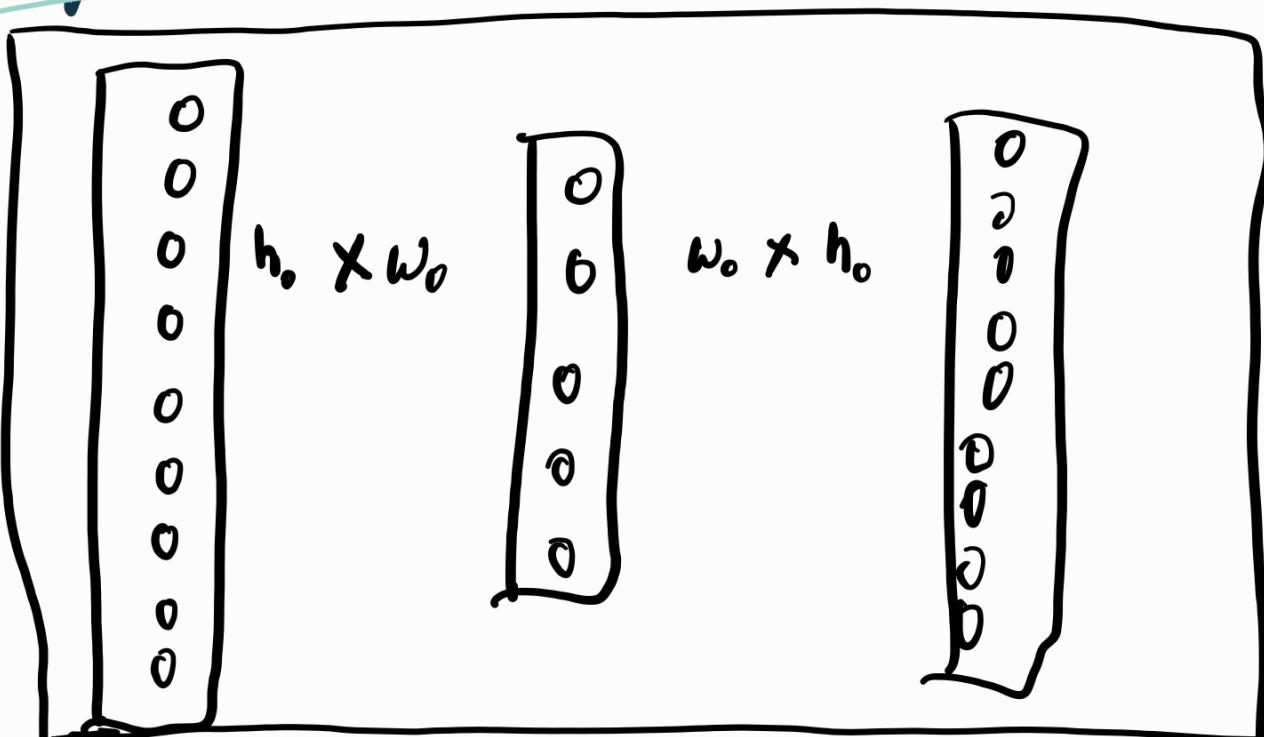
the output is compared to
the required output.

Back propagation is performed
and the weight matrix is
adjusted.

step - 7 the weight matrix is the
embedding value of the
word.

diagram

one hot vector = h_o



in detail code explained

The corpus used here is the book song of ice & fire by george R. R Martin (coolsb.txt). The file is in a .txt format.

i) tokenize Text : this function reads the file and splits it into tokens.

while (file >> word)

splits input by white spaces

Eg: Hello, world is one word

`if(ispunct(c))`

checks if the character is a punctuation

if punctuation the word in currentToken
is added to tokenWords

`else`

it is converted to lowercase
and added to tokenWords

Eg : Hello ! How's NLP going ? Let's tokenize_text

Output : "hello" "how's" "nlp" "going" "lets" "tokenizetext"

why ignore punctuations - how's → hows

else how's → "how", "s"

ii) build Vocabulary

In this function an unordered map called vocab is initialized.

each word is given an index number

Why use an unordered_map - faster search

best case $O(1)$

worst case $O(n)$

iii) create Context Target Pairs Indexed

CBOW algorithm is being used here.

CBOW stands for Continuous Bag Of Words

How CBOW works -

Sentence - My name is Eshvan Balaji

Window size = 2

target word = is

context words - My, name, Eshvan, Balaji

target word = balaji

context words - is, Eshvan

This is what is happening in this code.

Here the context words & the target word are stored as context-target-pairs

Why CBow and not Skip-Gram?

CBow shows better results than skip-gram
on small datasets.

CBow - predicts the target word

Skip-Gram - predicts the context

words given the target word.

iv) initialized weights

rand(time(0)) - random number generator

here $w_1 = \begin{bmatrix} & \text{vocab size} \\ \underbrace{\quad}_{w_1} & \end{bmatrix}$

$\underbrace{\quad}_{w_2}$

embedding dim

random values are assigned to w_1 & w_2

v) softmax

convert all the logit (raw scores) into probability distribution of 0 to 1.

First max score in the vector is found.

$$\text{exp-score}[i] = \exp(score(i) - \text{max-score})$$

here score - max-score

then exponent is done

then the value is added to

sum-exp

then $\frac{\text{val}}{\text{sum-exp}}$ is exp-score

vi) forward Pass

$W_1 = \text{vocab_size} \times \text{embedding_dim}$

$W_2 = \text{embedding_dim} \times \text{vocab_size}$

one hot vector $\times W_1$

then

$W_2 \times \text{one hot vector}$

then

logits are applied to the softmax
to get probability

vii) train Model

epochs are initialized here

learning_rate is initialized here

here the model trained.

i.e the embedding matrix is

built per word.

viii) getWordEmbedding

this function gets embedding vector

for a word from w1

Math

$$\begin{bmatrix} 0 \\ 0 \\ \textcircled{1} \\ 0 \\ 0 \\ 0 \end{bmatrix} \times \begin{bmatrix} 0.2 & \dots & \dots \\ \vdots & \ddots & \vdots \\ 0.6 & \dots & \dots \end{bmatrix} = \begin{bmatrix} \dots & \dots & \dots \\ \dots & \dots & \dots \\ \dots & \dots & \dots \end{bmatrix}$$

one
hot
vector

w_1

$$= \begin{bmatrix} \dots & \dots & \dots \\ \dots & \dots & \dots \\ \dots & \dots & \dots \end{bmatrix}$$

this is the word
embedding

Word2Vec : It is a technique in NLP designed to convert words into dense numerical vectors.

These embeddings capture semantic & contextual relationship b/w words.

semantic relationship : The word king & queen is placed close to each other

contextual relationship : "cat" & "dog" both placed near "pet".

Or - King - man + woman \approx queen

embedding improve performance in tasks like :

Machine translation (understanding word meaning)

Sentiment analysis (capturing context)

Information retrieval (matching related terms)