



vocabulary.

Basic representation : 'T', 'h', 'i', 's' ... → characters

Word.

'This', 'is', 'a', 'cat'.
↓
numbers.



word, sentence, paragraph, chapters. → book.

DNN → numbers.

{0-26}.
lower case. few punctuati
etc.

256

Character representation →

Memory consuming → in case of # of cycles required for parsing a sentence.

Word : If we encounter a new type of word in the test sentence, then ^{model} (if) will be unable to predict that word.

Vocabulary → create-

• Byte pair encoding

• Tokenization-

}

Balance
tokenizer

} current
models.

• Sparse representation → One hot vectors.

Let's consider.
→ for now

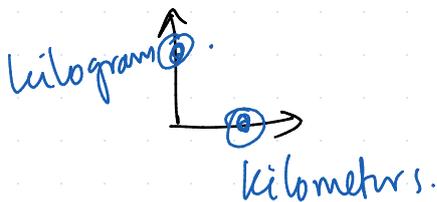
atomic element of a sentence is 'word'-

{ 'kilometers', 'kilograms' }

n x 1
vector.

$$\begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

$$\begin{pmatrix} 0 \\ 1 \end{pmatrix}$$



each of the words are unrelated to each other.

Hence we can have different meanings for each of the individual words.

'queen', 'lady', 'girl' → not completely unrelated.

'king', 'man', 'boy' → not completely unrelated text.

$\begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$ $\begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$ $\begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$

↳ which is not practical / useful.

• since we want similar words to have similar representations.

Hence cosine similarity should be higher for similar words.

Special kinds of tokens :-

<S> , <bos> → start of the sentence.

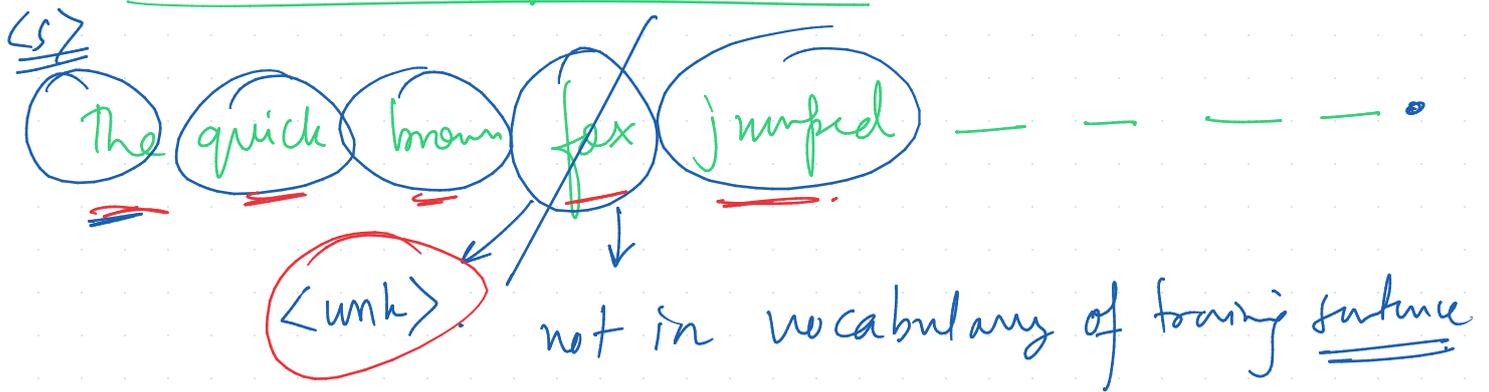
<IS> , <eos> → end of sentence.

<unk> → unknown. : Most of the cases, we will see that a new kind of

<think> → </think> } token is encountered during

reasoning of chatGPT. inference.

Sentence completion task.



Task:

Learn meaningful representations for word

Learn word embeddings. (Token embeddings)

2013

Tomáš Mikolov. Jeff Dean

unsupervised algo - → Google news dataset.

automatically learns meaningful representations.

No annotation required from human side.

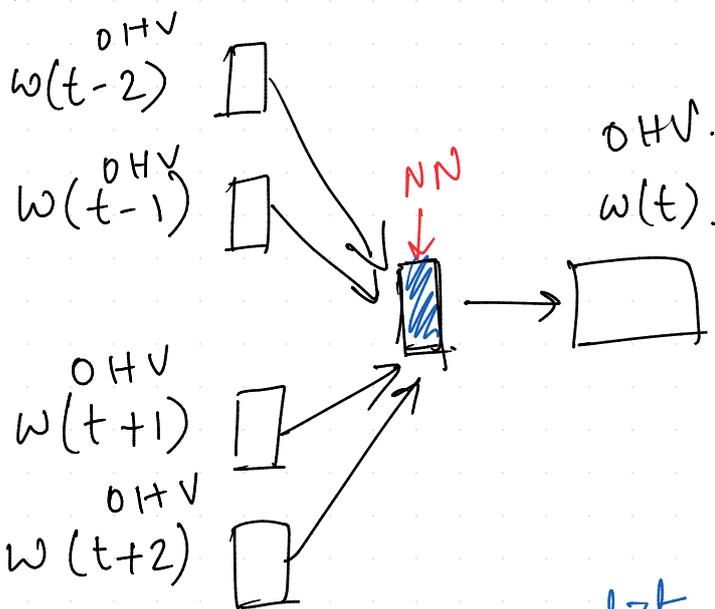
'A word is known by the company it keeps'

Similar words will be closer to each other.

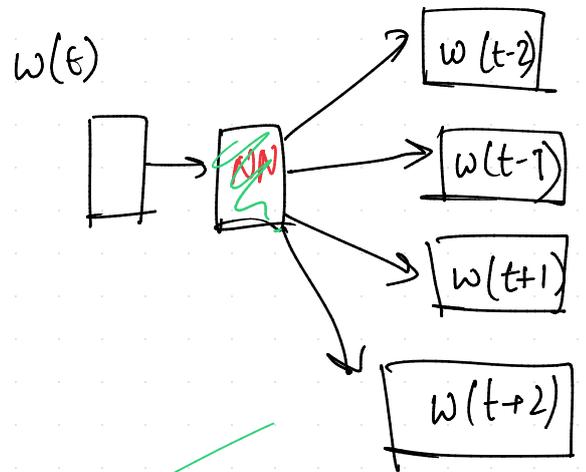
Word2Vec Models.

CBOW

Continuous Bag of Words.



Skip-gram Models.



predict

The cat sat on the mat.

sliding window.

CBOW → predicts a single target based on its surrounding context words.

Skip-gram

Predicts the surrounding context words given a single target word.

Vocabulary \rightarrow $(1, 0, 0, \dots, 0)$ One hot rep.

Categorical Cross Entropy loss \rightarrow Training.

Pytorch ; $\text{OHV} \Rightarrow$ sparse

embedding vector
 \downarrow
function

10,000
 $\begin{pmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & \dots & 1 & \dots & 0 \\ \vdots & & & & \vdots \\ 0 & \dots & \dots & \dots & 0 \end{pmatrix}$ 10000

Sparse rep \rightarrow dense.

vocab size \rightarrow $\mathbb{V} \Rightarrow$ \mathbb{D} dense representation

$\text{OHV} \times$
 $300 \rightarrow$ \mathbb{D} -dim-vector.

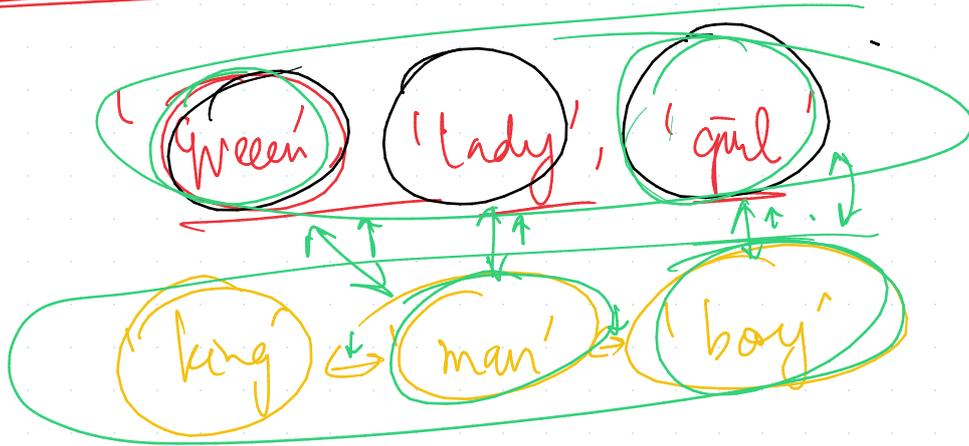
come similarity b/w any two words = 0, if it is sparse.

not sparse
doesn't have much zeros.

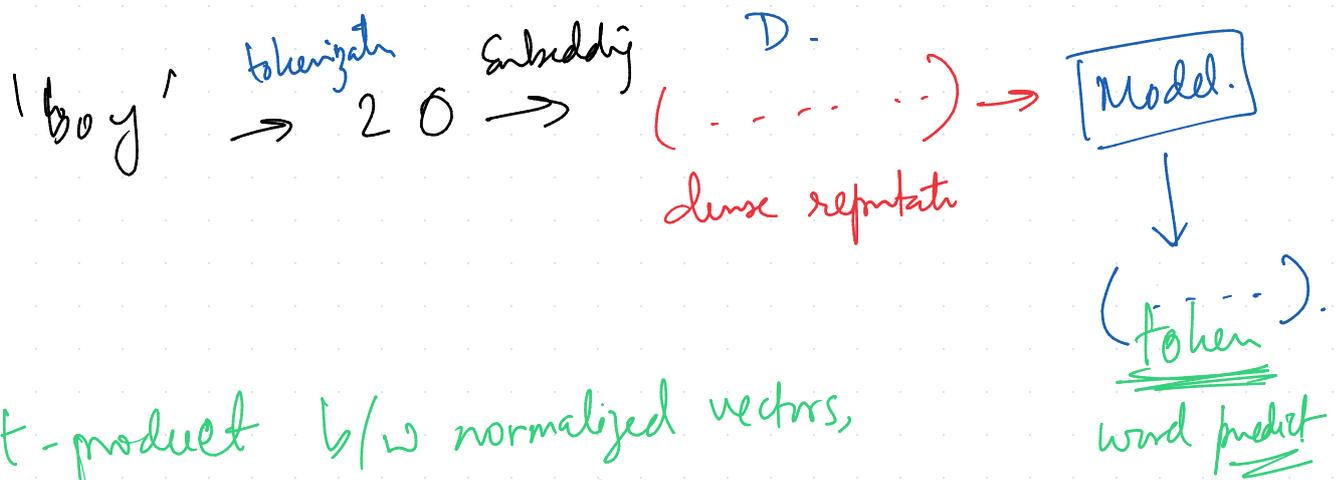
$(2.3, 0.7, 0.6, 1.3, \dots)$
300
 $\sqrt{(2.3)^2 + (0.7)^2 + (0.6)^2 + (1.3)^2 + \dots}$

helpful when computing the cosine similarity b/w word vectors.

Dense representation



meaningful representations are learned.

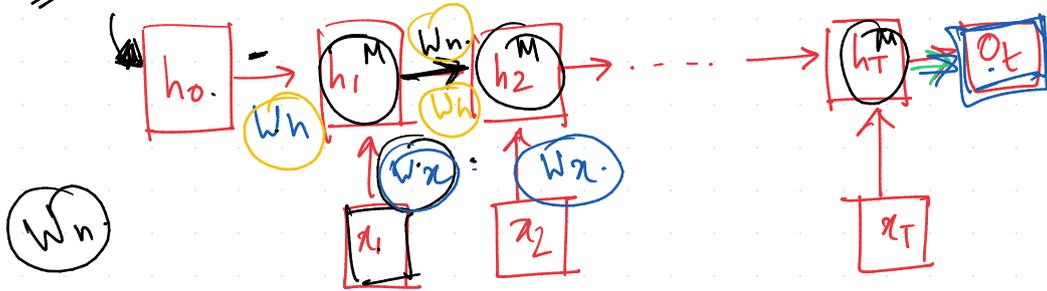


dot-product b/w normalized vectors, then it already calculates cosine similarity

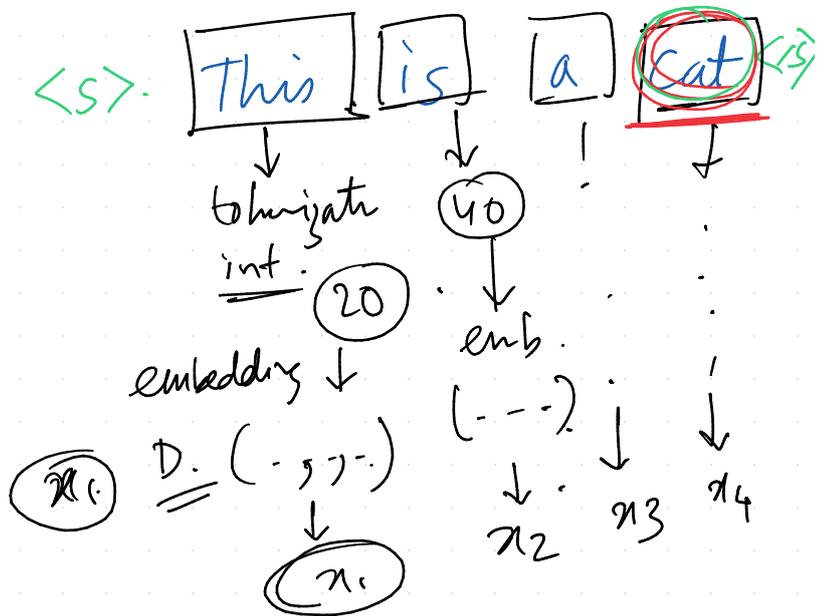
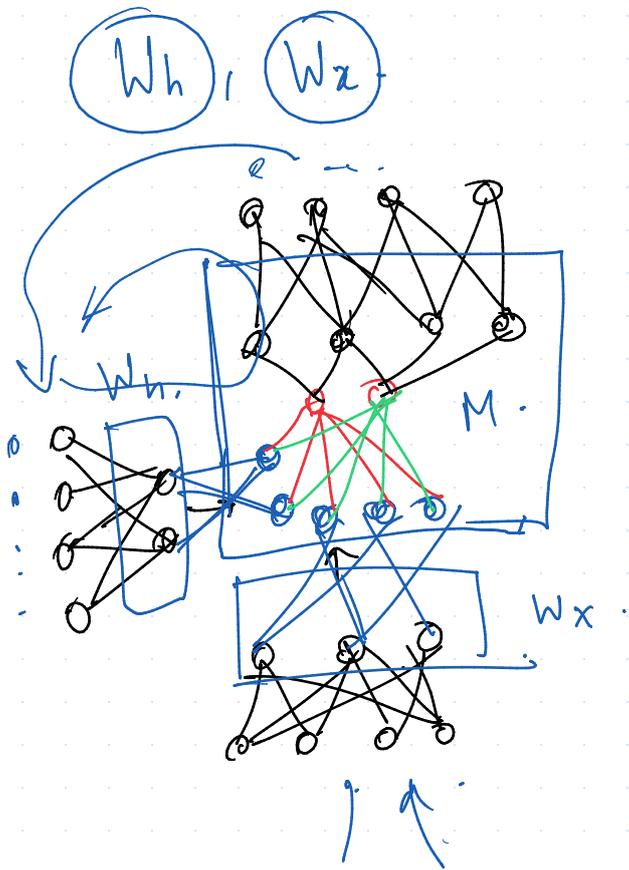
$$\begin{array}{ccc}
 (1, 0, 0) & (0, 1, 0) & (0, 0, 1) \\
 \text{king} & \leftrightarrow \text{man} & \text{boy} \\
 \downarrow & & \downarrow \\
 (1 \cdot 0 + 0 \cdot 1 + 0 \cdot 0) & & \\
 \rightarrow (0.1, 0.9) & (0.9, 0.7) & (0.3, 0.6) \\
 \downarrow & & \\
 0.1 \times 0.2 + 0.9 \times 0.7 & & \\
 0.1 \times 0.9 + 0.9 \times 0.1 & \neq 0 &
 \end{array}$$

Recurrent Neural Networks (RNNs)

(...)



Not computing
o/p for the
time being.



$x_1, x_2, \dots, x_T \rightarrow$ sequence of word vectors.

$g \rightarrow$ non-linear function.

$$h_i = g(W_h \underline{h_{i-1}} + W_x x_i).$$

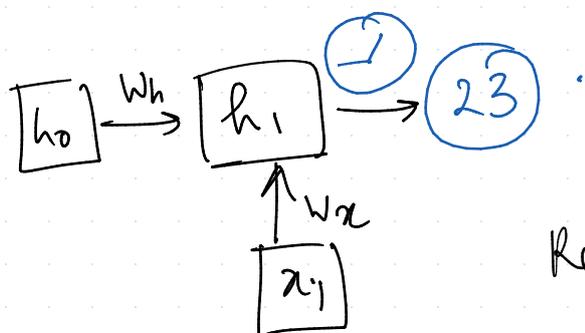
$$y = \text{argmax} \left(\text{Softmax} \left(\underline{\underline{g(\cdot)}} \right) \right)$$

$\langle s \rangle$ hello $\langle /s \rangle$

$$W_x = \begin{bmatrix} 1 & 5 & 3 \end{bmatrix}$$

$$W_h = [7]$$

$$h_0 = [1]$$



$$x_1 \rightarrow \underline{\underline{[3 \ 2 \ 1]}}$$

$$\text{ReLU} \left(\underline{W_h} \cdot \underline{h_0} + \underline{W_x} \underline{\text{embd}(\langle s \rangle)} \right)$$

$$\text{ReLU} \left([7] \times [1] + [1 \ 5 \ 3] \begin{bmatrix} 3 \\ 2 \\ 1 \end{bmatrix} \right)$$

$$\Rightarrow \text{ReLU} \left(7 + 3 + 10 + 3 \right)$$

$$\Rightarrow \text{ReLU}(23) \Rightarrow 23$$

Shortcomings of RNNs.

- Vanishing & exploding gradients.
- History gets forgotten at time.

Temporary solution

- Gradient clipping

Normalize the gradient of a parameter vector when its L2 norm reaches a certain value.