

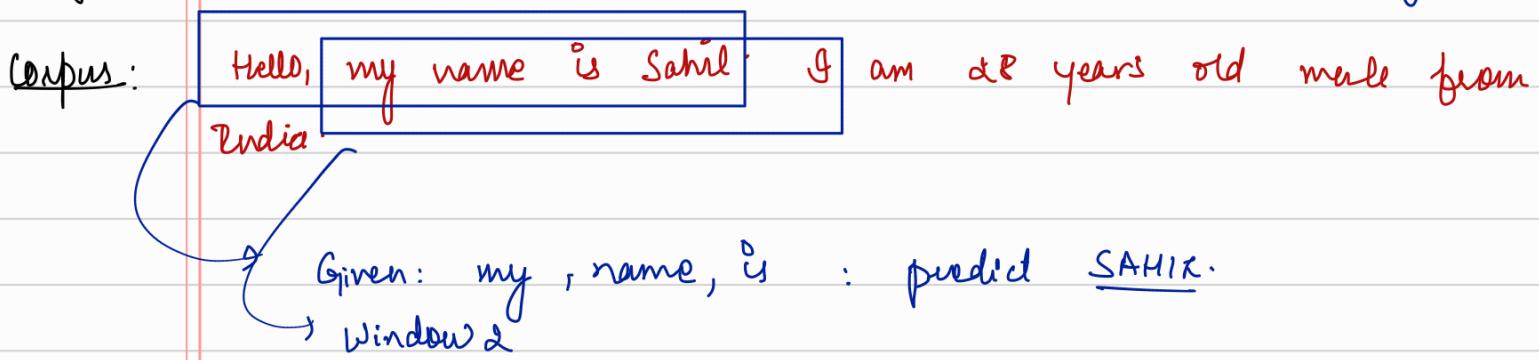
Continuous Bag of Words Model

- Consider the task: Given $(n-1)$ previously observed words, predict the n^{th} word.

Example: He sat on a Chair: Think of this as a classification problem.

y or the output is a one-hot encoded vector of size equal to the size of the vocabulary.

Training Data: From the given corpus, construct a n window training data.



Case when $n=2$

Input: one word: 'BAT': $[0 | 0 | 0 | 1 | 1 | 0 | 0]$
 ← OHE rep of input word

$p(\text{hel})$ $p(\text{chair})$ $p(\text{man})$... $p(\text{on})$

length of OEP = $V = \text{vocab size}$

OEP prob:

$| | | | | \text{BAT} | \dots | | | | |$

$W_{\text{word}} \in \mathbb{R}^{K \times V}$

middle layer: $h \in \mathbb{R}^K$

$W_{\text{context}} \in \mathbb{R}^{K \times V}$

$$\text{SAT: } \boxed{0 \mid 0 \mid 1 \mid 0 \mid \dots \mid 0 \mid 0} \quad \text{OHE: }$$

Also v dimensional

$x \in \mathbb{R}^{(V)}$ { one-hot encoded representation }

$$\text{W context: Dimension} = (K \times V)$$

$$\begin{bmatrix} & \\ K \times V & \end{bmatrix} * \begin{bmatrix} V \times 1 \end{bmatrix} = \begin{bmatrix} K \times 1 \end{bmatrix}$$

W context Input

- Let's explore the $\text{W context} \times 'x'$

$$\begin{bmatrix} \vec{w}_1 & \vec{w}_2 & \dots & \vec{w}_k & \vec{w}_V \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix} \stackrel{\text{K}^{\text{th}} \text{ element} = 1}{=} \vec{w}_k$$

{ Each \vec{w}_i = K dimensional vector }
 { There are V such vectors. }

main takeaway: \rightarrow Since the input is one-hot encoded, the product of the two will be a the K^{th} column n always.

Input one-hot encoding acts as a switch where depending on which element = 1, a column from W context will be selected.

There is 1:1 mapping between number of words in Vocab & the columns of the W context matrix.

- How do we obtain the required output i.e. $P(\text{On} | \text{Sat})$
 we use the softmax:

$$P(\text{on} | \text{sat}) = e^{\frac{1}{\text{Word} \times K} \sum_i \text{Word} \times K[i]} \quad \begin{array}{l} \text{get the output} \\ \text{and then} \end{array}$$

$$\sum_{i=1}^V e^{w_{\text{word}} h[i]}$$

apply a softmax transformation

hidden layer $h \in \mathbb{R}^k$

$$\text{op layer} = w_{\text{word}}^T h \rightarrow \text{Softmax Transformation}$$

$$\text{op} = W_{\text{word}} \times h$$

(VxK) (Kx1)

*a column of W_{context}
depending on what input is
given*

$$\begin{matrix} & \leftarrow K \rightarrow \\ \begin{matrix} T \\ V \\ \downarrow \end{matrix} & \left[\quad \right] = \left[\quad \right] T \\ & \downarrow \downarrow \downarrow \downarrow \downarrow \end{matrix} \quad \begin{matrix} \text{Each entry} \\ \text{in the op vector} \\ \text{of length } V \end{matrix}$$

this is because hidden layer is created based on OHE vector

y_k : k^{th} entry in the op vector = dot product of k^{th} row of W_{word} and the j^{th} column of W_{context} . j 's value is determined by which element of the OHE vector is equal to '1'.

Thus $P(\text{word} = i | \text{sat})$ thus depends on the i^{th} row of W_{word} .
Thus W_{word} 's i^{th} row also corresponds to a particular word.

so, the model learns two different representations. One when word appears as context word & one when it appears as a target word.

Further Exploration

Let's denote the context word ('sat') by C and the correct output word ('on') by the index ' w '

loss function for the multiclass problem formulation : Cross-entropy

$$L(\theta) = -\log(\hat{y}_w) = -\log P(w|C)$$

it is equivalent to minimizing just the score

of ' w ' or $p(\text{on}|\text{sat})$

$$\text{Cross entropy} = -\sum y_i \log(p_i) \quad y_i = 1 \text{ only for 'on' so it}$$

$h = W_{\text{context}} \times x_c = u_c$] Basically the column of W_{context} , let's say W_j based on what 'j' is in the one-hot rep of the word.

Vectorial representation of the word 'sat'.

$$\hat{y}_w = \frac{\exp(u_c \cdot v_w)}{\sum_{w' \in V} \exp(u_c \cdot v_{w'})}$$

] same as mentioned above
 $v_c \cdot v_w = \text{product of two columns.}$

$$L(\theta) = -\log(\hat{y}_w) = -[\log(\exp(u_c \cdot v_w)) - \log(\kappa)]$$

$$L(\theta) = -u_c \cdot v_w + \log\left(\sum_{w' \in V} \exp(u_c \cdot v_{w'})\right)$$

$$\frac{\partial L}{\partial v_w} = -u_c + \left[\frac{1}{\kappa} \times \exp(u_c \cdot v_w) \times u_c \right]$$

\hookrightarrow Softmax again

\hookrightarrow parameter that we intend to learn.

$$\frac{\partial L}{\partial v_w} = -u_c \cdot (1 - \hat{y}_w)$$

update rule

$$v_w^{\text{new}} = v_w^{\text{old}} - \eta \times \frac{\partial L}{\partial v_w}$$

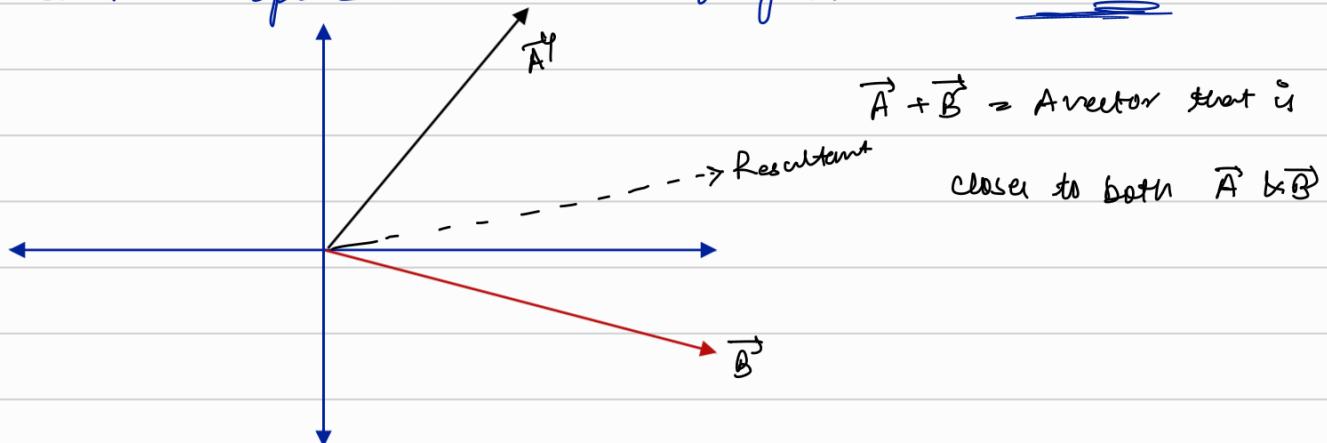
$$v_w^{\text{new}} = v_w + \eta u_c (1 - \hat{y}_w)$$

Case 1: $\hat{y}_w = 1$] This is the ideal case, meaning that we have learnt a good representation of v_w

Case 2: $\hat{y}_w \approx 0$] This case, the new representation of v_w is

$$v_w^{\text{new}} = v_w + \eta_{\text{dc}}$$

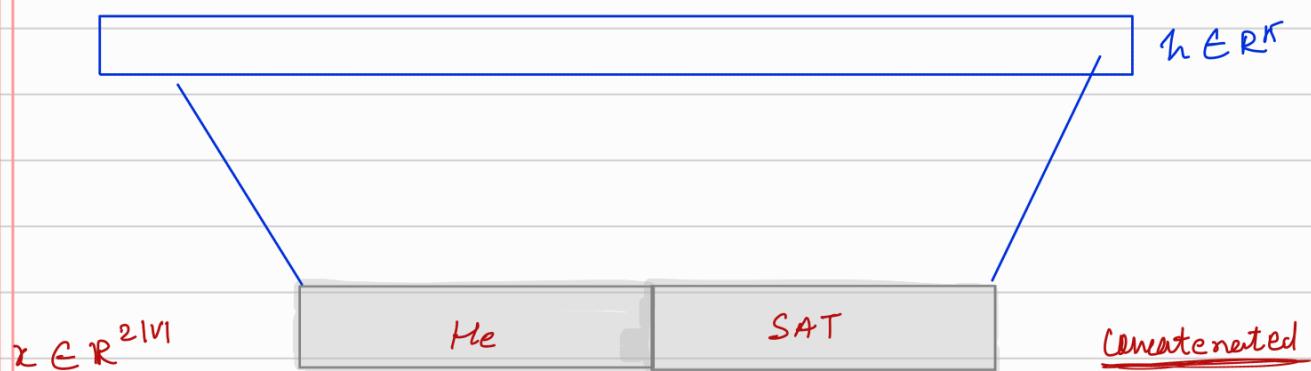
In vector space this means shifting v_w closer to v_{dc} .



Essentially: The algorithm is just trying to make word representation closer to the context representation, thereby ensuring that the cosine similarity between the two vectors is maximized.

- Let's expand on this a little bit: what happens if we change the window to '3'. Meaning, given 2 words, predict the third word.

$f(\text{on}, \{\text{he}, \text{sat}\})$



0	0	1	0	0	...	0	0	he
---	---	---	---	---	-----	---	---	----

0	0	0	0	1	...	0	0	sat
---	---	---	---	---	-----	---	---	-----

$h \in \mathbb{R}^K$ is simply the sum of all the columns that are '1' in the one-hot encoded realm.

For example if he and sat are input, lets say col 5 & col 7. $h = \text{Sum}(C(\text{col } 5 \text{ & col } 7))$ for the W^{context} .

During the back propagation, All columns of W_{word} and a column of the W^{context} will be updated.

Issues with this case

Here, the base issue is with the softmax function.

$$\text{opp} = \mathbb{R}^V = (1 \times K) \times [K \times V] \text{ matrix.}$$

passing to the softmax.

$$\text{Each entry } j, \text{ lets call it } p_j = \frac{e^{(h \times W_{\text{word}}[:, j])}}{\sum_{k=1}^V e^{(h \times W_{\text{word}}[:, k])}}$$

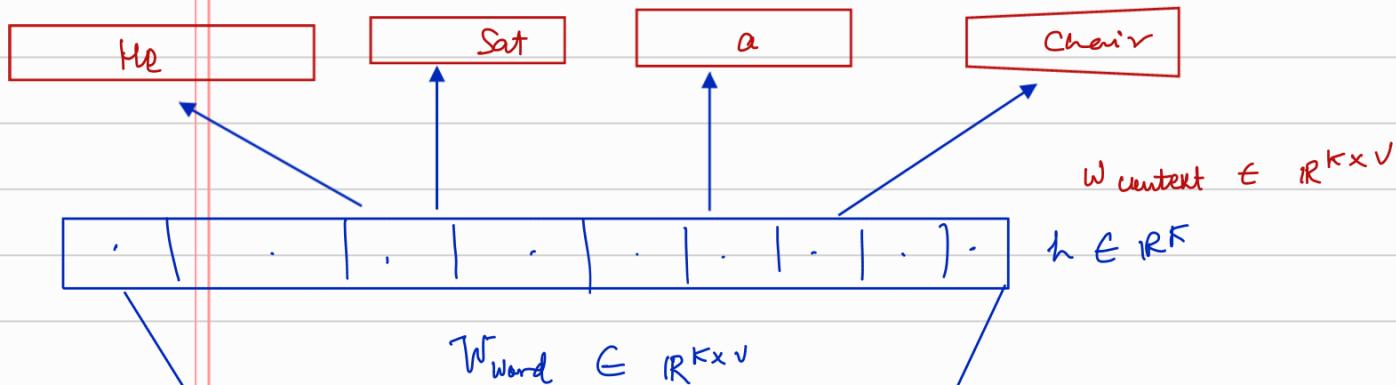
Can be of the order of millions.

Computationally complex

Skip-Gram Model

This model is kinda reverse of the previous CBOW model.

objective: Given a word, predict a set of context word.



0	0	...	0	1	0	0	0
---	---	-----	---	---	---	---	---

$$x \in \mathbb{R}^{1 \times 1}$$

\leftarrow \rightarrow on

Set of input word \rightarrow obviously it will be a v -dimensional vector which is one-hot encoded.

Loss function: Set of 4 categorical cross entropies

Also, we have flipped W_{word} and W_{context} in this case.

$$\text{Loss function} = L(\theta) = - \sum_{i=1}^{d-1} \log (\hat{y}_{w_i}) \rightarrow \text{probability}$$

d = 5 in our case, which is the window.

Loss function = sum of 4 categorical cross entropies

some problems: Same issue as bag of words. Instead of one softmax, now we have four softmax which adds a lot to the computational issue.

Solution 1: use negative sampling:

Sol 2: Contrastive estimation

Sol 3: hierarchical softmax.

\rightarrow Negative Sampling: let D be a set of corpus $D = \text{Set } (w, c)$ pairs
 $D' = (w, s) \rightarrow$ Synthetically created pairs that don't belong together.

$$p(z=1 | w, c)$$



v_w = Rep. of word

v_c = Rep. of context word

model of the problem

Given (w, c) , I want to maximize the prob. $p(z=1 | w, c)$

v_c

v_w

$$p(z=1 | w, c) = \sigma(v_c^T v_w)$$

Consider $(w, c) \in D$, I want to maximize the $p(z=1 | w, c)$

$$\text{Loss}(u, v) = - \prod_{(w, c) \in D} p(z=1 | w, c)$$

$$\log(\text{loss}) = - \sum_{(w, c) \in D} p(z=1 | w, c)$$

simultaneously $\# w, r \in D'$, I want to max. $p(z=0 | (w, r))$

$$\text{objective} = \max \prod_{(w, r) \in D'} p(z=0 | (w, r))$$

objective: For every pair in 'D' $\max \prod p(z=1 | (w, c))$ AND for every pair in D' , $\max \prod p(z=0 | (w, r))$

$$\text{obj} = \max \left[\underbrace{\prod_{(w, c) \in D} p(z=1 | (w, c))}_A \times \underbrace{\prod_{(w, r) \in D'} p(z=0 | (w, r))}_B \right]$$

prob is modeled as a Sigmoid function

$$p(z=1 | (w, c)) = \sigma(u_c^T v_w) \quad \begin{matrix} \rightarrow \text{Vector rep. of} \\ \text{context} \end{matrix}$$

$$\log(\text{objective}) = \sum_{(w, c) \in D} \log(\sigma(u_c^T v_w)) \quad \begin{matrix} \rightarrow \text{Vector rep. of} \\ \text{word} \end{matrix}$$

$$+ \sum_{(w, r) \in D'} \log(1 - \sigma(u_r^T v_w)) \quad \begin{matrix} \rightarrow \text{Vector rep. of} \\ \text{randomly selected} \\ \text{negative word} \end{matrix}$$

$$1 - \sigma(u_r^T v_w) = 1 - \frac{1}{1 + e^{-\sigma \cdot u_r^T v_w}}$$

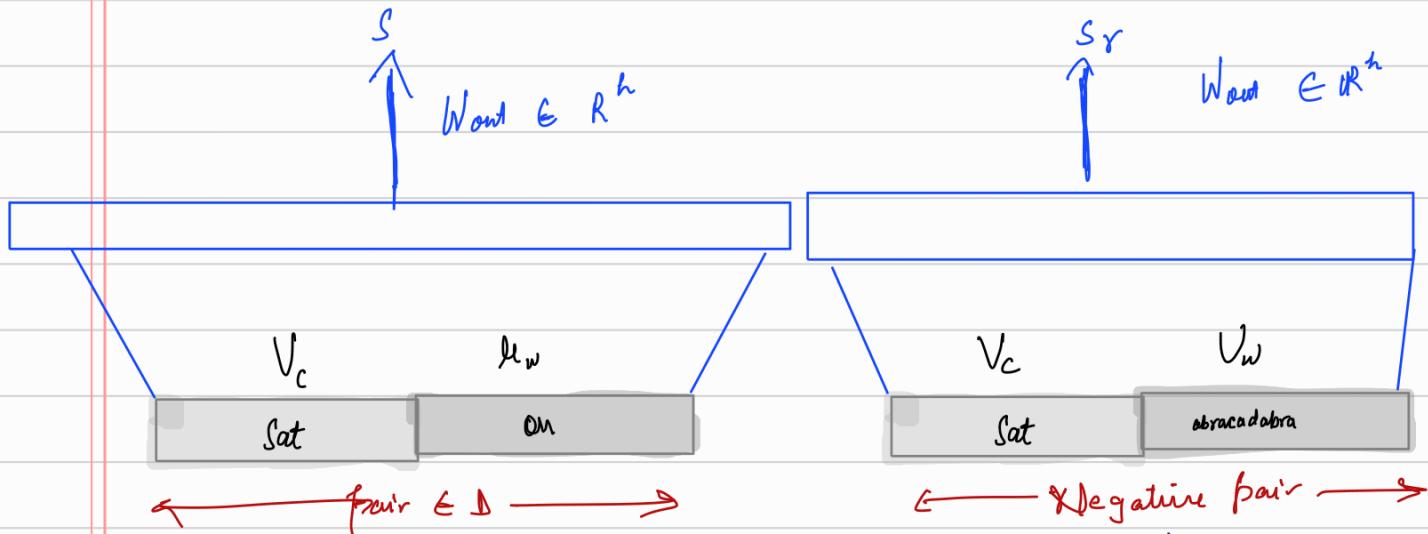
$$= \frac{1 - e^{-\sigma \cdot u_r^T v_w}}{1 + e^{-\sigma \cdot u_r^T v_w}} \quad \begin{matrix} -x \\ \cancel{-x} \end{matrix} \quad = \frac{1}{1 + e^{\sigma \cdot u_r^T v_w}}$$

$$\log(\text{Objective}) = \sum_{w, c \in D} \log(\sigma(\underbrace{\mu_c^T V w}_{\text{dot product}})) + \sum_{w, r \in D'} \log(\sigma(\underbrace{-\mu_r^T V w}_{\text{dot product}}))$$

Key Takeaway :- In addition to bringing close words together (part A) of above 3gn, we are also moving dissimilar words away from each other (part B) of 3gn.

Contrastive Estimation

Instead of formulating the problem as a probability problem like above, we can formulate it as following:



Instead of calculating the $p(z=1 | V_c, V_w)$ we can train it as a regression problem where one has to calculate the output score 'S'.

We would like for S to be much greater than S_r . Not only that, I want S to be greater than S_r by a margin.

$S - S_r$ should be at least 'm'] m will be controlled by me.

Thus, we are interested in $S - (Sr + m)$

loss \neq if $S > Sr + m ; \lambda = 0$

$S < Sr + m ;$ calculate loss and
adjust your
embeddings.

$$\text{loss} = \max(0, S - (Sr + m))$$

→ no softmax
use here.

