# Skip - Gram - Model

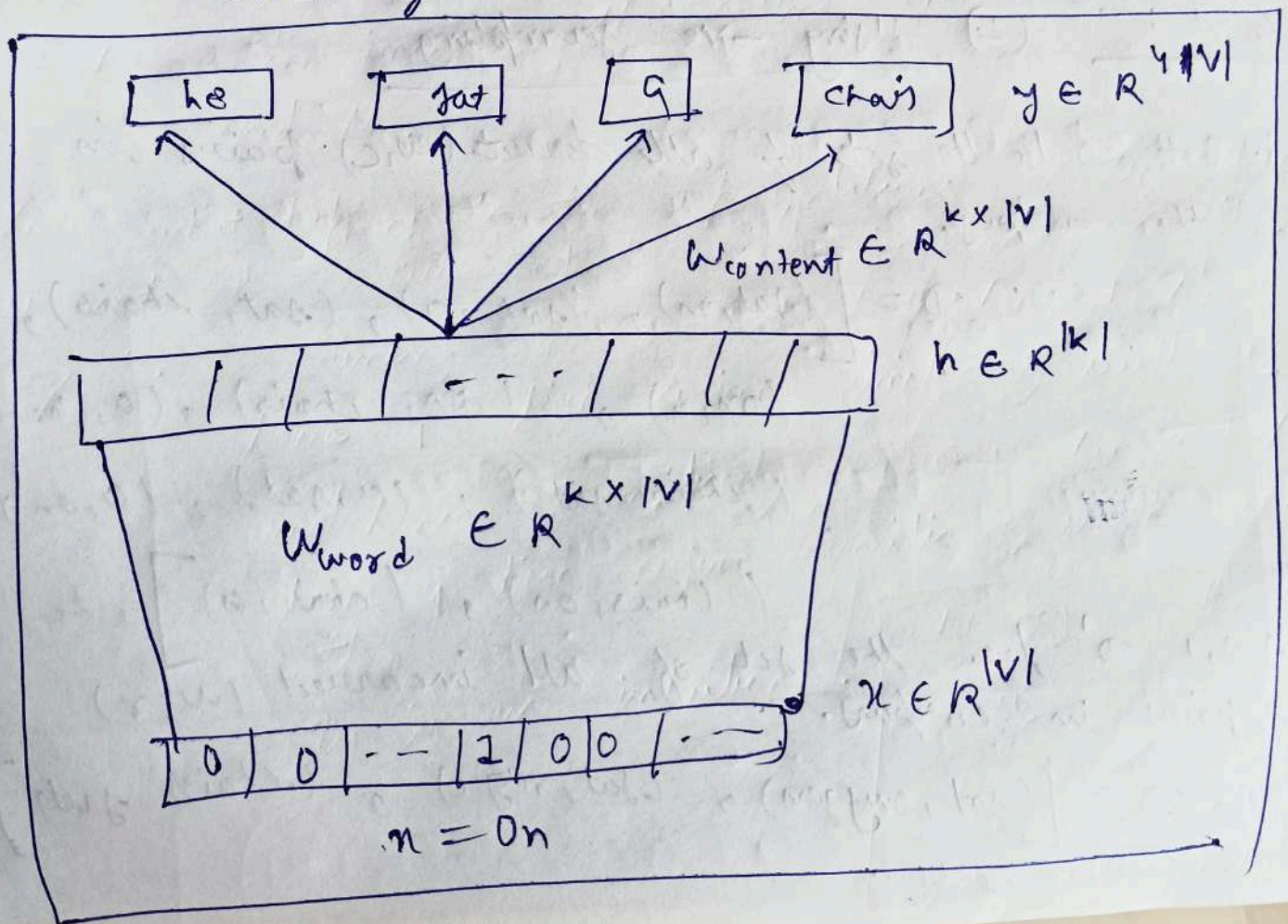→ Recall Previously, CTS bag of words model predicted a word given content.

→ But, skip - gram model will predict content given a word.



he    fat    a    chair    $y \in R^{4 \ast |V|}$

$W_{content} \in R^{k \times |V|}$

$h \in R^{|k|}$

$W_{word} \in R^{k \times |V|}$

| 0 | 0 | -- | 1 | 0 | 0 | -- |

$x \in R^{|V|}$

$n = On$

→ there the loss $\mathcal{L}^n$ will be sum of cross - entropies.

$$\mathscr{L}(\theta) = -\sum_{i=1}^{d-1} \log \hat{y}_{w_i}$$

**Problems :-**

(i.) softmax $f^n$ at the output is is computationally expensive.

$$\hat{y}_w = \frac{e^{u_c \cdot v_w}}{\sum_{w' \in V} e^{u_c \cdot v_w}}$$

$u_c = c$-th column $\gamma$ of $W_{content}$

$v_w = w$-th column of $W_{word}$

**Sol$^n$ :-** ① **Using -ve sampling**

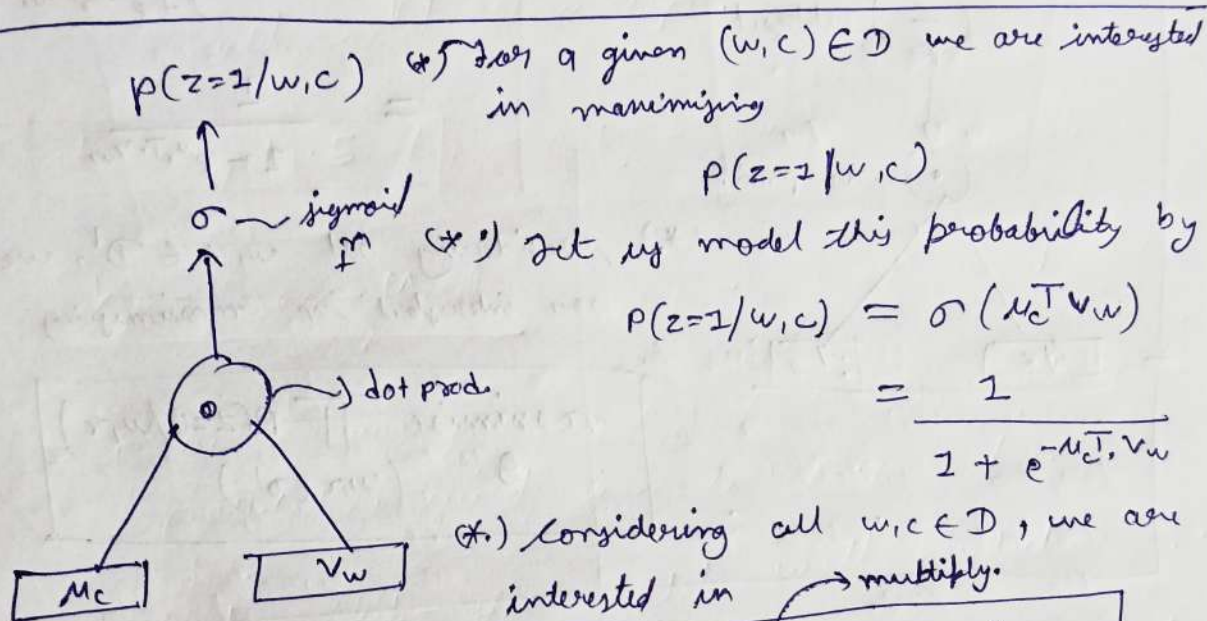(i.) Let $D$ be the set of all correct $(w,c)$ pairs in the corpus.

eg. $D = [\,(sat, on), (sat, a), (sat, chair),$

$(on, a), (on, chair), (a, chair),$

$(chair, sat), (a, sat), (a, on),$

$(chair, on), (chair, a)\,]$ etc.

(ii.) Let $D'$ be the set of all incorrect $(w, r)$ pairs in corpus.

→ $D' = [\,(sat, oxygen), (sat, magic), (chair, sad)\,]$

etc.

(iii.) $D'$ can be constructed by randomly sampling a content word $r$ which has never appeared with $w$ in a creating pair $(w, r)$.

(iv.) As before let $V_w$ be the representation of the word $w$ and $M_c$ be the representation of the content word $c$.
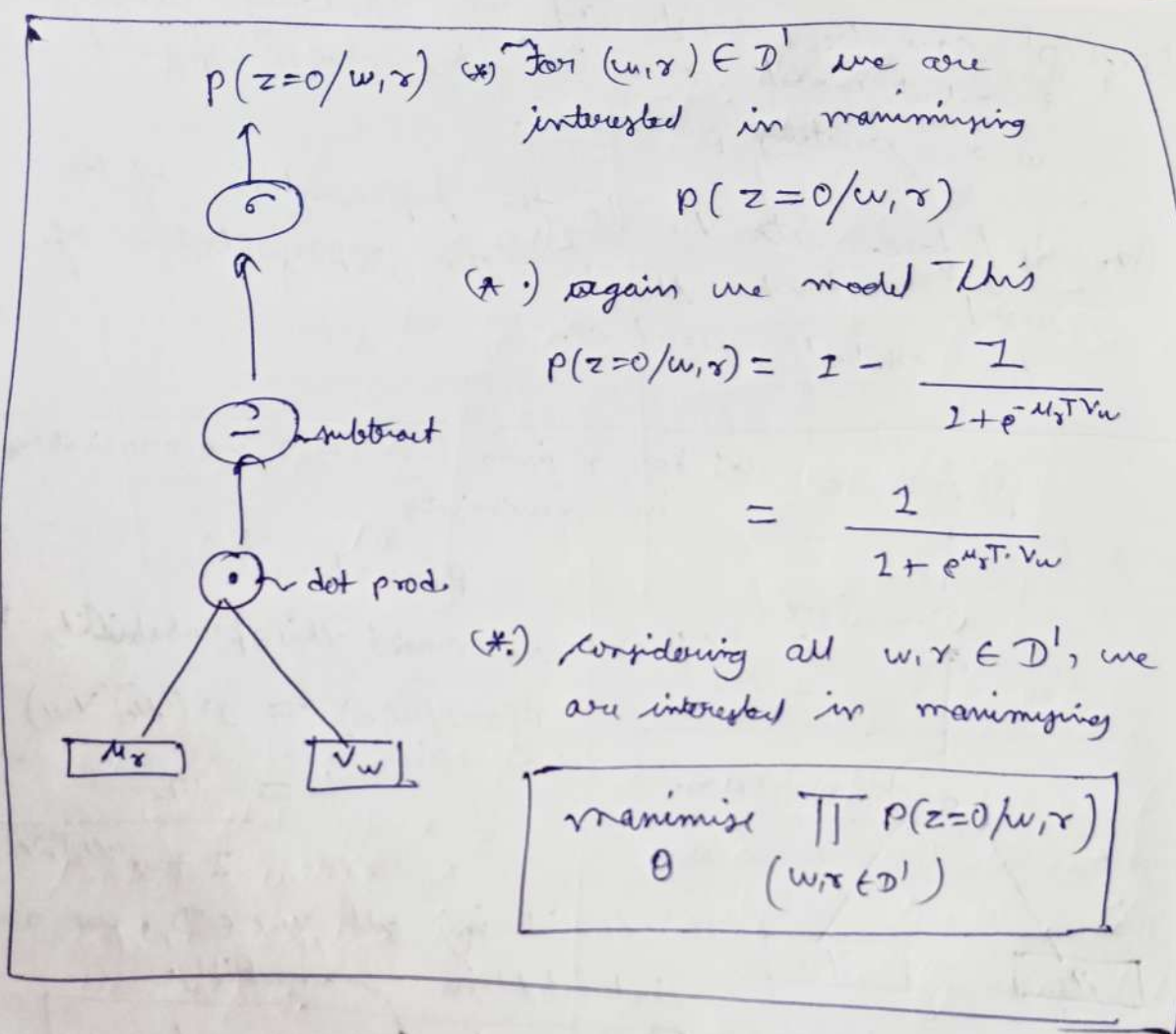


$p(z=1/w,c)$ ← $\sigma$ ~ sigmoid

*) For a given $(w,c) \in D$ we are interested in maximizing

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

**) Let us model this probability by

$$p(z=1/w,c) = \sigma(M_c^T V_w)$$

$$= \frac{1}{1 + e^{-M_c^T \cdot V_w}}$$

dot prod.

$M_c$     $V_w$

*) Considering all $w,c \in D$, we are interested in ⟶ multiply.

$$\text{maximise} \prod_{(w,c) \in D} p(z=1/w,c)$$
$$\theta$$

$\theta$: is the word representation $(V_w)$ and content representation $(M_c)$ for all words in our corpus.

we want to maximise the $p=1$ (probability) for each word in our corpus.

∴ multiplication (and).

$p(z=0/w,r)$

$\uparrow$

$\sigma$

$\uparrow$

$\overset{-}{\ominus}$ subtract

$\odot$ dot prod.

$\boxed{\mu_r}$ $\boxed{v_w}$

(*) For $(w,r) \in D'$ we are interested in maximising

$$p(z=0/w,r)$$

(A.) again we model this

$$p(z=0/w,r) = 1 - \frac{1}{1+e^{-\mu_r^T v_w}}$$

$$= \frac{1}{1+e^{\mu_r^T \cdot v_w}}$$

(*.) considering all $w,r \in D'$, we are interested in maximising

$$\boxed{\underset{\theta}{\text{maximise}} \prod_{(w,r \,\in D')} P(z=0/w,r)}$$

Note that we want $p(z=1/w,c)$ for each correct word $\in D$ and we also want $P(z=0/w,r)$ for each incorrect word $\in D'$. Both simultaneously i.e. AND.

Hence our final goal is

$$\underset{\theta}{\text{maximise}} \prod_{(w,c) \in D} P(z=1/w,c) \prod_{(w,r) \in D'} \big(P(z=0/w,r)\big)$$

By simplifying and taking log, we get.

$$\text{maximise}_{\theta} \sum_{w,c \in D} \log\left(\sigma\left(u_c^T v_w\right)\right) + \sum_{(w,r)\in D'} \log\left(\sigma\left(-u_r^T v_w\right)\right)$$

$$\sigma(x) = \frac{1}{1+e^{-x}}$$

Note, in the original research paper.

(i) Size of $D'$ is $k$ times the size of $D$.

(ii) random content word is drawn from a modified unigram distribution

$$\gamma \sim \left(P(\gamma)\right)^{3/4}$$

⊛ where, $P(\gamma) = \frac{count(\gamma)}{N}$

$N = $ total no. of words in corpus.