

CBOW MODEL, ONE WORD CONTEXT $C=1 \rightarrow$ CONTRAST SKIP GRAM

INPUT VECTOR: 1-HOT ENCODING OF LENGTH V
(PUTS ALL prob MASS on the 1 word)

$$\vec{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_k \\ \vdots \\ x_v \end{bmatrix}$$

FOR INPUT WORD k ,
 $x_k = 1$ AND $\forall j, j \neq k, x_j = 0$

$W = V \times N$ (N HIDDEN NODES/UNITS); INPUT \rightarrow HIDDEN LAYER
EACH ROW OF W AN N -DIMENSIONAL VECTOR
REPRESENTATION OF INPUT WORD w , \vec{v}_w

HIDDEN LAYER h

$$h = x^T W = w_k := \vec{v}_{w_I}$$

\vec{v}_{w_I} IS THE VECTOR REPRESENTATION
OF INPUT WORD w_I

$$h = \begin{bmatrix} 0 & \overset{k}{1} & 0 \end{bmatrix} \begin{bmatrix} \vec{v}_1 \\ \vdots \\ \vec{v}_v \end{bmatrix} = \begin{bmatrix} 0 & 1 & \dots \end{bmatrix} 1 \times N = h$$

$x^T: 1 \times V$ $V \times N$ $1 \times N$

k^{th} row of W

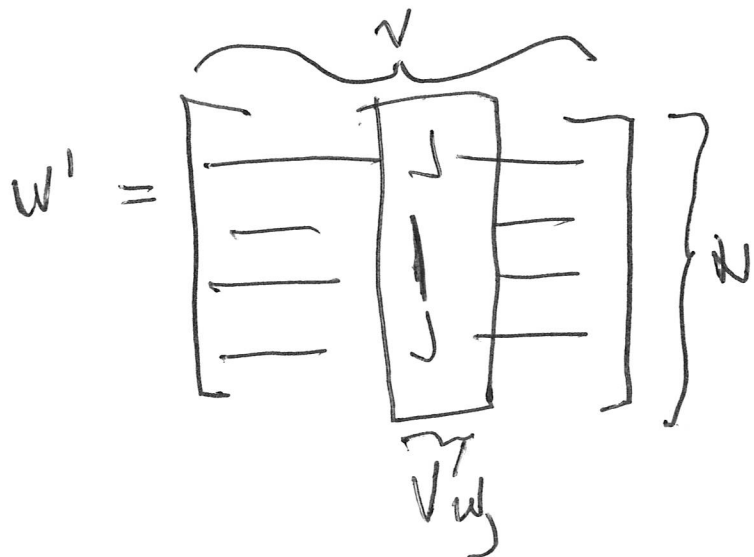
$x^T W = h \Rightarrow$ COPIES THE k^{th}
ROW OF W TO h

HIDDEN LAYER \rightarrow OUTPUT LAYER

$$W' = \{w'_{ij}\} \quad N \times V$$

COMPUTE SCORE u_j

$$u_j = V_{w_j}'^T h \quad V_{w_j}' \text{ IS THE } j^{\text{th}} \text{ COLUMN OF } W'$$



$$V_{w_j}'^T h = \begin{bmatrix} \text{---} & \text{---} & \text{---} \end{bmatrix}_{V_{w_j}'^T} \cdot \begin{bmatrix} \text{---} \\ \text{---} \\ \text{---} \end{bmatrix}_h = \text{SCALAR } u_j$$

$1 \times V$ $1 \times h$

GENERAL SOFTMAX

$$P(w_j | w_I) = y_i = \frac{\exp(u_j)}{\sum_{j=1}^V \exp(u_j)}$$

y_i : output of j^{th} node
in output layer

$$u_j = V_{w_j}'^T h = V_{w_j}'^T X^T w$$

$$p_{\theta}(y|x) = \frac{\exp(\theta^T f(x,y))}{\sum_{y' \in Y(x)} \exp(\theta^T f(x,y'))}$$

$$h = x^T W := v_{w_I} \quad (\text{VECTOR REP of INPUT WORD } w_I)$$

$$u_j = v_{w_j}^T \cdot h \quad \leftarrow x^T W := v_{w_I}$$

SCORE j^{th} column of W' ($N \times V$) $\Rightarrow \begin{bmatrix} w_{1,j} \\ \vdots \\ w_{N,j} \end{bmatrix}$

SOFTMAX POSTERIOR DISTRIBUTION

$$P(w_j | w_I) = y_j = \frac{\exp(u_j)}{\sum_{j'=1}^V \exp(u_{j'})}$$

OUTPUT OF j^{th} NODE IN OUTPUT LAYER

$$P(w_j | w_I) = \frac{\exp(v_{w_j}^T v_{w_I})}{\sum_{j'=1}^V \exp(v_{w_{j'}}^T v_{w_I})}$$

ISSUE: DOT PRODUCT? WRITTEN LIKE OUTER PRODUCT?

NOTE: v_w AND $v_{w'}$ ARE 2 REPRESENTATIONS OF INPUT WORD w
 v_w COMES FROM W , THE INPUT \rightarrow HIDDEN WEIGHT MATRIX
 the rows of

$v_{w'}$ - columns of W' , hidden \rightarrow output

TRAINING

IDEA IS TO MAXIMIZE

$$P(w_0 | w_I) = \frac{\exp(v_{w_0}'^T v_{w_I})}{\sum_{j'=1}^V \exp(v_{w_{j'}}'^T v_{w_I})}$$

$$\Rightarrow \max (P(w_0 | w_I)) = \max y_{j^*} \\ = \max \log y_{j^*}$$

$$= \max \log \left(\frac{\exp(u_{j^*})}{\sum_{j'=1}^V \exp(u_{j'})} \right)$$

$$\begin{aligned} &= v_{w_0}'^T h + \log \sum_{j'=1}^V \exp(v_{w_{j'}}'^T h) \\ &\quad \leftarrow = \max \log(\exp(u_{j^*})) - \log \sum_{j'=1}^V \exp(u_{j'}) \\ &= u_{j^*} - \log \sum_{j'=1}^V \exp(u_{j'}) \\ &:= -E \quad \left(\log \right) \end{aligned}$$

so $E = -\log P(w_0 | w_I)$ is the loss function
(want to minimize this), j^* is the index
of the actual output word in the output layer

BTW, $E = -\log p(w_0 | w_T)$ IS A SPECIAL CASE OF THE CROSS ENTROPY ($y=1$)

UPDATE HIDDEN \rightarrow OUTPUT

NEED THE DERIVATE OF E WRT THE j TH NODES IN POT, NAMELY $u_j \Rightarrow$

$$\frac{\partial E}{\partial u_j} = \partial (u_{j^*} - \log(\sum_{j'=1}^V \exp(u_{j'})))$$

$$= \frac{\partial u_{j^*}}{\partial u_j} - \frac{\partial \log \sum_{j'=1}^V \exp(u_{j'})}{\partial u_j}$$

$$= t_j - \frac{1}{\sum_{j'=1}^V \exp(u_{j'})} \quad \text{where}$$

$t_j = 1 (j=j^*)$ ie 1 if $j=j^*$, 0 otherwise

$$\frac{1}{\sum_{j'=1}^V \exp(u_{j'})} \sim y_j \quad \left(\text{ISSUE - REALLY NEED } \exp(u_j) \text{ IN THE NUMERATOR} \right)$$

$$\Rightarrow \frac{\partial E}{\partial u_j} = y_j - t_j := e_j \leftarrow \begin{array}{l} \text{PREDICTION} \\ \text{ERROR } e_j \text{ in the output} \\ \text{LAYER} \end{array}$$

NEXT: DERIVATIVE OF w'_{ij} TO GET

The gradient of the hidden \rightarrow output weights

$$\frac{\partial E}{\partial w'_{ij}} = \frac{\partial E}{\partial u_j} \cdot \frac{\partial u_j}{\partial w'_{ij}} \quad (\text{chain rule})$$

$$= e_j \cdot \frac{\partial u_j}{\partial w'_{ij}}$$

$$= e_j \cdot \frac{\partial (V_{w'_j}^T \cdot h)}{\partial w'_{ij}}$$

$$V_{w'_j} = \begin{bmatrix} w'_{1j} \\ w'_{2j} \\ w'_{3j} \end{bmatrix}$$

$$= e_j \cdot h_i$$

$$V_{w'_j}^T = [w'_{1j} \ w'_{2j} \ w'_{3j}]$$

UPDATE INPUT \rightarrow HIDDEN

$$\frac{\partial E}{\partial h_i} = \sum_j \frac{\partial E}{\partial u_j} \cdot \frac{\partial u_j}{\partial h_i} \quad (\text{chain rule})$$

$$= \sum_{j=1}^J e_j \cdot w'_{ij} = E h_i$$

WHAT EXACTLY IS w'_{ij} ?

$$\frac{\partial V_{w'_j}^T}{\partial w_{ij}} = [0 \dots 1 \dots 0]$$

$$\frac{\partial h}{\partial w'_{ij}} = \vec{h}$$

$$\frac{\partial (V_{w'_j}^T \cdot \vec{h})}{\partial w_{ij}} = h_i$$

SELECTS the i th
elt of h