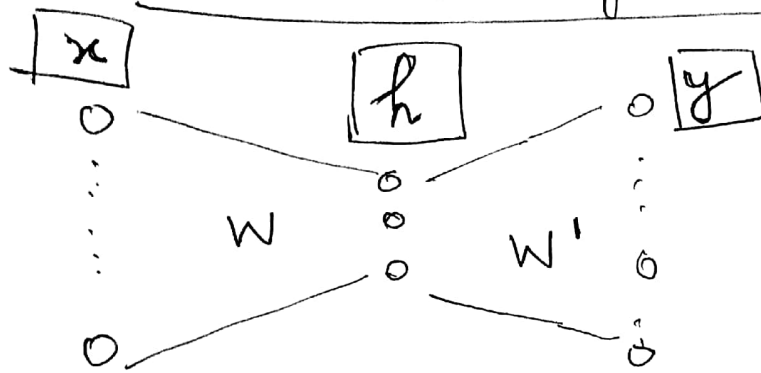
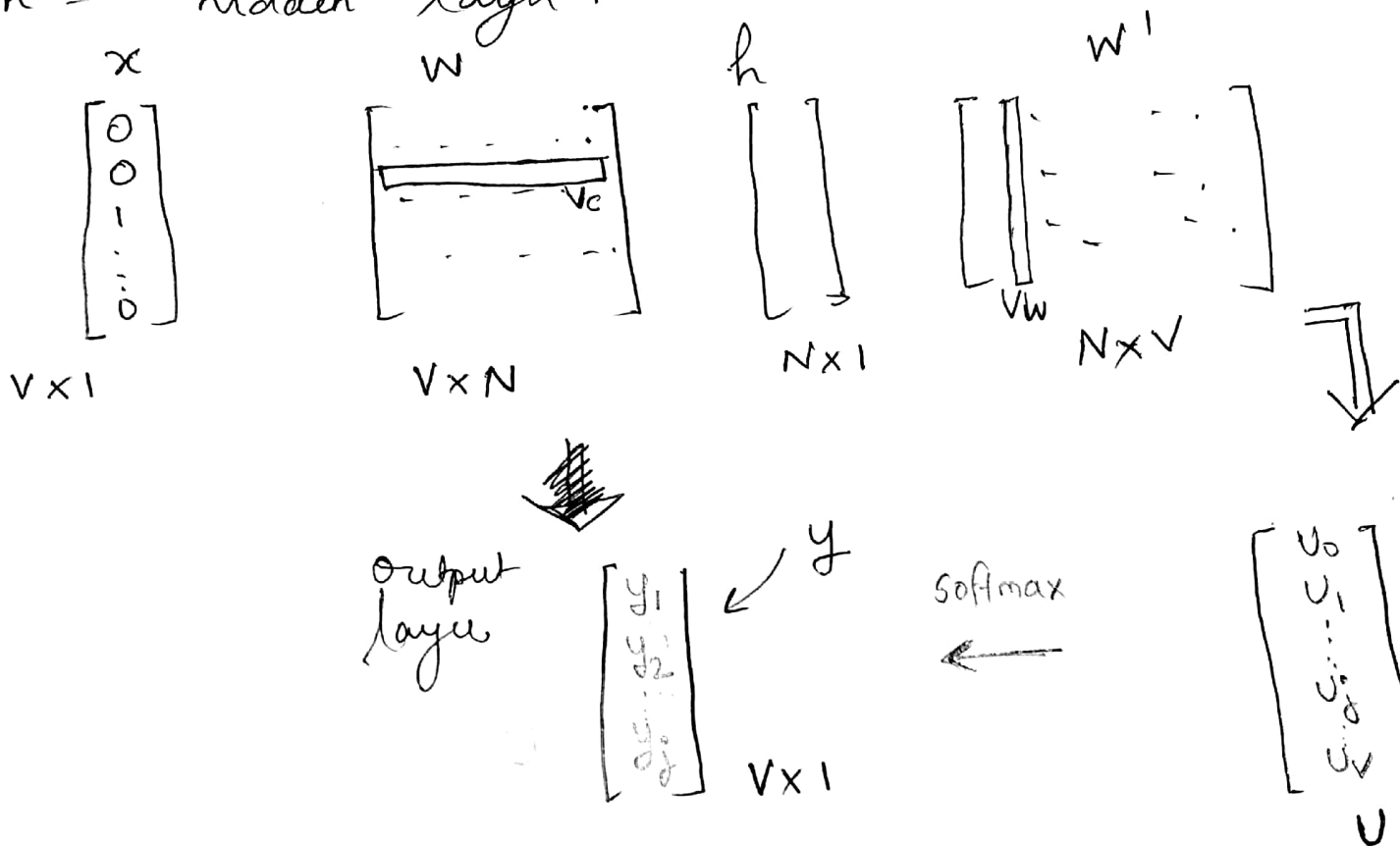


# Word Embeddings - Skip Gram model



$x$  - one hot representation of word  $w$  (the context word)  
 $h$  - hidden layer.



$$h = x^T W = W^T x \quad \{ x^T W \text{ is a row of matrix } W$$

so  $h$  is a row of  $W$  (a row of context words)

$h = v_c = W^T x$

$\{ \because W^T x \text{ is of } N \times 1 = \text{size}$

$$V \times N \quad N \times 1 \Rightarrow \underline{V \times 1}$$

Diagram illustrating matrix multiplication:

Matrix 1 (Input): Dimensions  $N \times T$  (height  $N$ , width  $T$ ). A horizontal slice of width  $V$  is highlighted, resulting in a sub-matrix of dimensions  $V \times N$ .

Matrix 2 (Input): Dimensions  $N \times 1$  (height  $N$ , width  $1$ ). A horizontal slice of width  $V$  is highlighted, resulting in a sub-matrix of dimensions  $V \times 1$ .

Result: The product of the two sub-matrices is a single column vector of dimensions  $V \times 1$ , labeled as the "output layer".

word  $w$  in our corpus.

$$y_i = \sigma(u_i)$$

$$= \frac{e^{v_i}}{\sum_j e^{v_j}} = \frac{e^{v_{\omega}^T v_c}}{\sum_{\omega'} e^{v_{\omega'}^T v_c}}$$

and  $\sum_{w'} w'$  means summation over all words.

words:

$$y_i = \frac{e^{v_c^T v_w}}{\sum_{w'} e^{v_c^T v_{w'}}}$$

$y_i = P(w|c)$  = prob. of a word  $w$  given that a context word  $c$  is given

Our parameters are =  $\{v_c, v_w\}$

$$\theta = \{v_c, v_w\}$$

$$L(\theta) = \prod P(w|c; \theta) \quad \left\{ \begin{array}{l} \text{Product of} \\ \text{probabilities} \end{array} \right\}$$

$$= \prod \frac{e^{v_c^T v_w}}{\sum_{w'} e^{v_c^T v_{w'}}}$$

To maximise this we can take Maximum Log Likelihood

$$l(\theta) = \sum_{w'} v_c^T v_w - \log \sum_{w'} e^{v_c^T v_{w'}}$$

Let find  $\frac{\partial l}{\partial v_w}$  first

$$\frac{\partial l}{\partial v_w} = v_c - \frac{1}{\sum_{w'} e^{v_c^T v_{w'}}} \times e^{v_c^T v_w} \times v_c$$

$$= v_c - P(w|c) v_c$$

$$\frac{\partial l}{\partial v_w} = v_c (1 - P(w|c))$$

Now we can apply gradient descent to find optimal parameters

## Gradient Descent

$$V_w^{\text{new}} \leftarrow V_w^{\text{old}} - \alpha \cdot V_c(1 - P(w|c))$$

↓  
learning rate

1.5      1.1      1.25      1.1      1.2