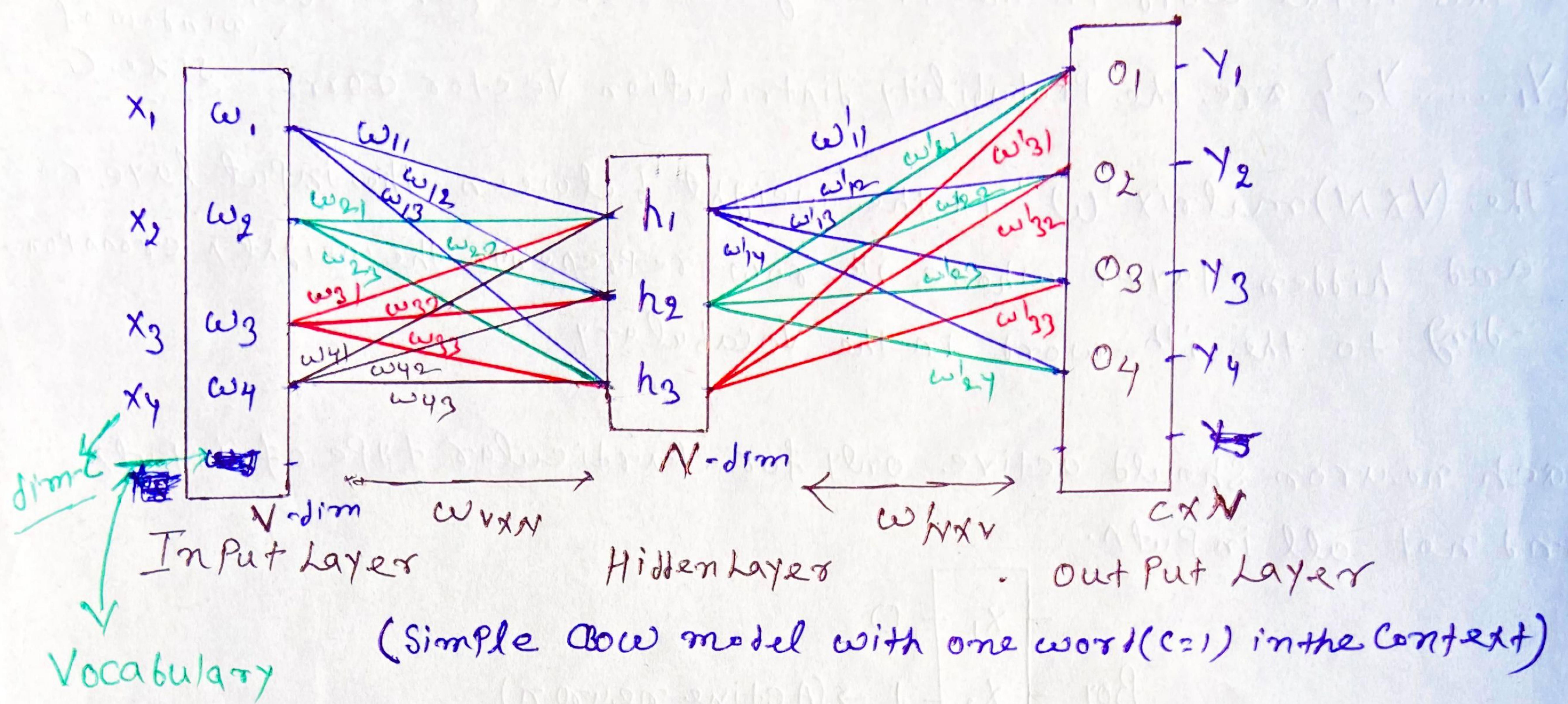


Continuous bag of words (CBOW)



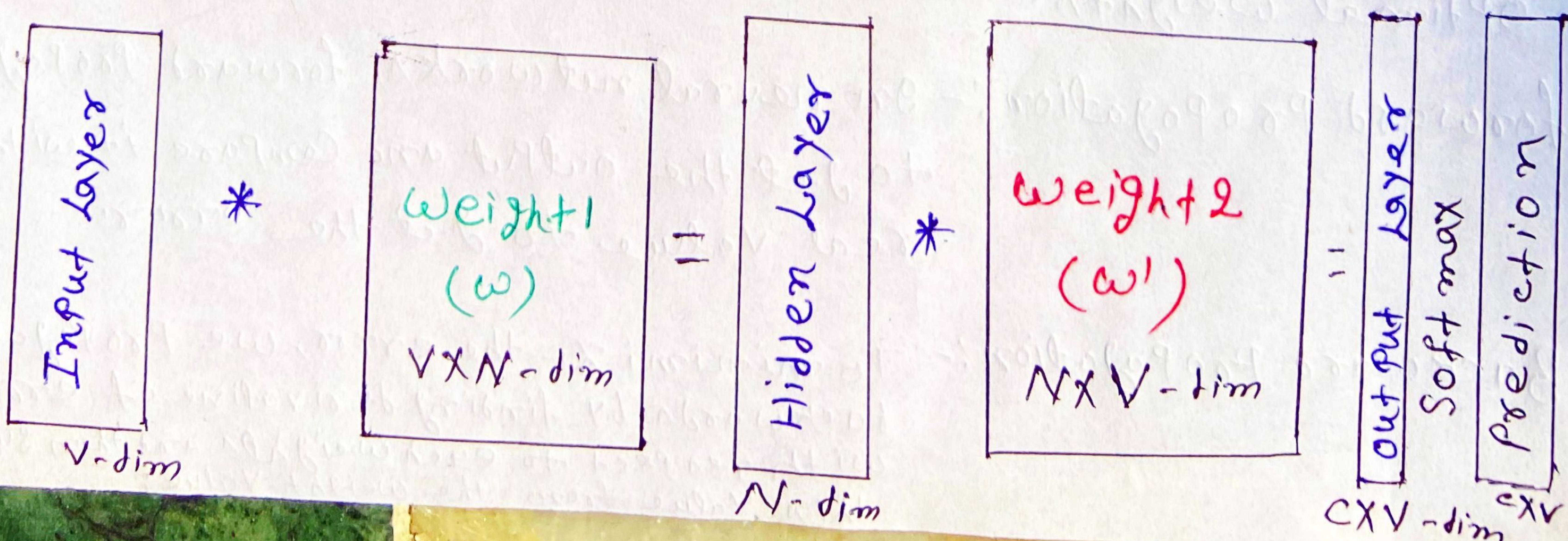
Weight of neuron(w) between Input Layer & hidden layer

$$W = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \\ w_{41} & w_{42} & w_{43} \end{bmatrix}_{V \times N}$$

Weight of neuron(w') between hidden & output Layer.

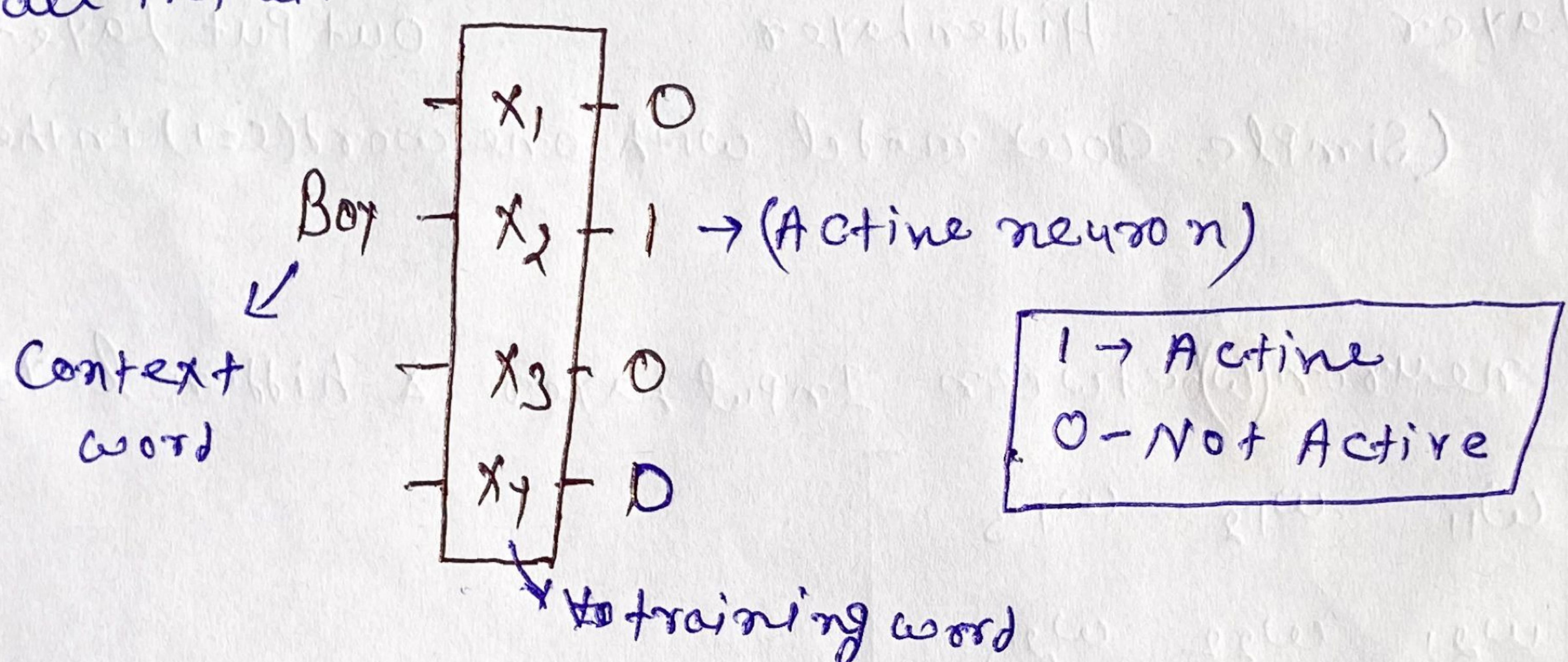
$$W' = \begin{bmatrix} w'_{11} & w'_{12} & w'_{13} & w'_{14} \\ w'_{21} & w'_{22} & w'_{23} & w'_{24} \\ w'_{31} & w'_{32} & w'_{33} & w'_{34} \end{bmatrix}_{N \times V}$$

Model



- In this model ' x ' represent the one-hot encoded vector corresponding to the input word in the training instance. $\{x_1, \dots, x_c\}$ for a word window of size C
- $\{y_1, \dots, y_c\}$ are the probability distribution Vector ~~over~~
- The $(V \times N)$ matrix ' w ' is the weight between the input layer and hidden layer where i th row represents the weights corresponding to the i th word in the vocabulary.

Each neuron should active only for particular type of input and not all inputs.



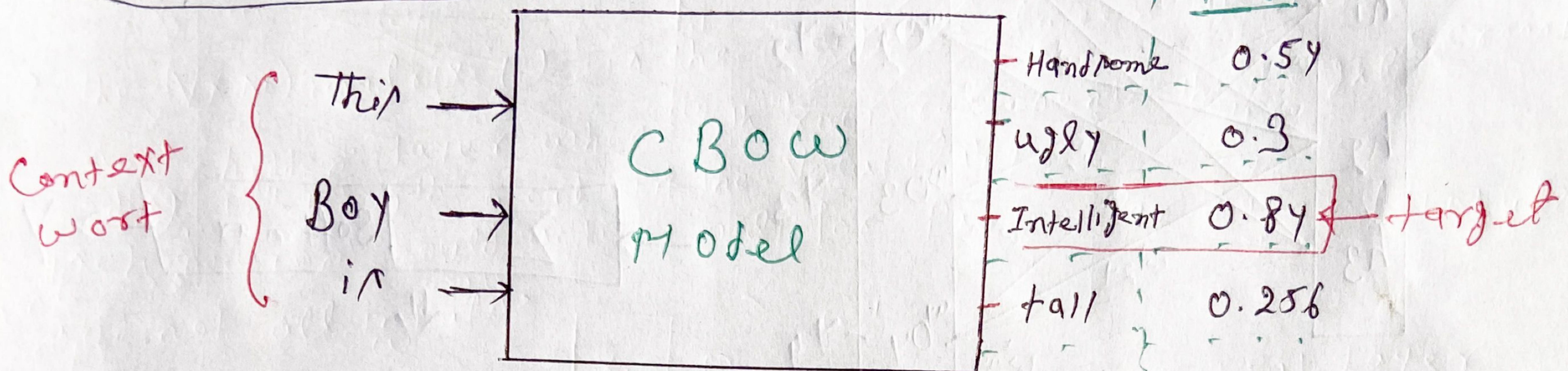
Therefore, by propagation forward we see how well our model is behaving and find the error. If we find out our model has error, we back propagate and use a form of gradient descent to update new values of the weight (w, w'). Then we will again forward propagate to see how well those weights (w, w') are performing and then will backward propagate to update the weight. This will go on until we reach some minimum error value and optimal weights.

forward Propagation:- In neural networks, forward propagate to get the output and compare it with the real value to get the error.

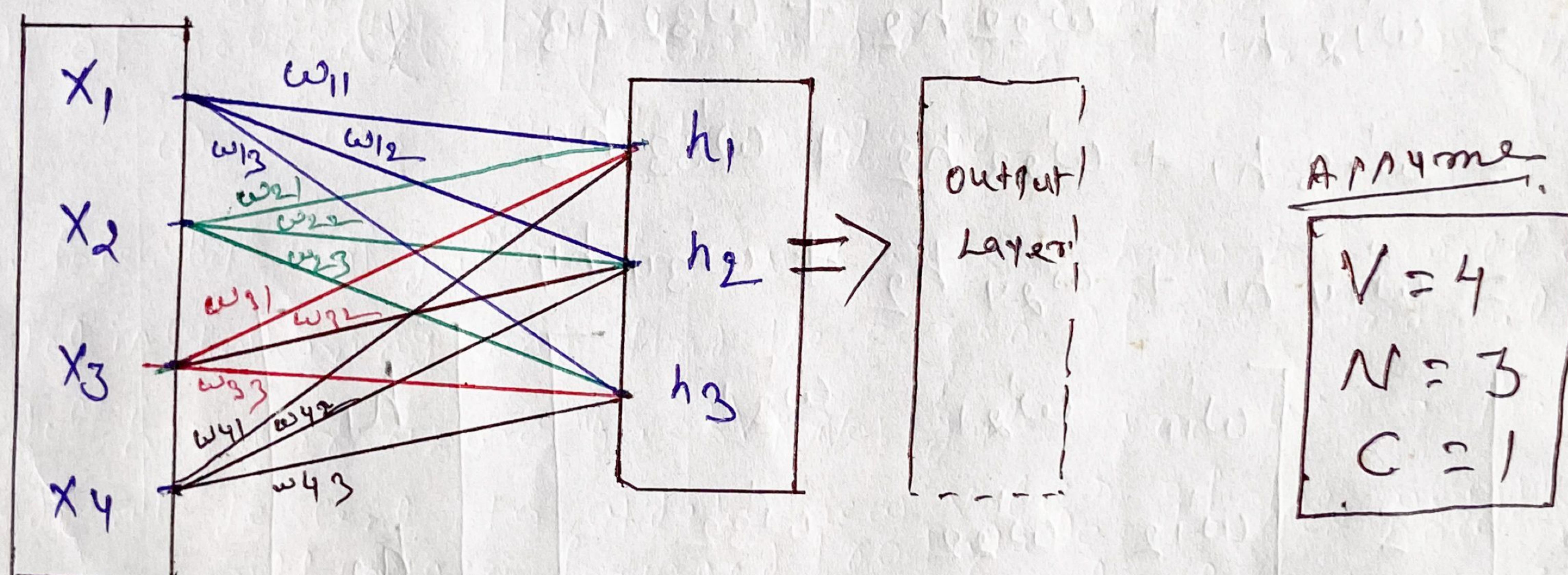
Backward Propagation:- To minimize the error, we propagate backward by finding derivative of error with respect to each weight and then subtracting this value from the weight value.

After getting optional weight of neuron to we predict the target word. output will provide the probability distribution & one correspondence to the word.

"This boy is ___."



Forward Propagation: weight calculation (w): - Input Layer to hidden layer



$$h_1 = w_{11}x_1 + w_{21}x_2 + w_{31}x_3 + w_{41}x_4$$

$$h_2 = w_{12}x_1 + w_{22}x_2 + w_{32}x_3 + w_{42}x_4$$

$$h_3 = w_{13}x_1 + w_{23}x_2 + w_{33}x_3 + w_{43}x_4$$

$$\begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} = \begin{bmatrix} w_{11} & w_{21} & w_{31} & w_{41} \\ w_{12} & w_{22} & w_{32} & w_{42} \\ w_{13} & w_{23} & w_{33} & w_{43} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}$$

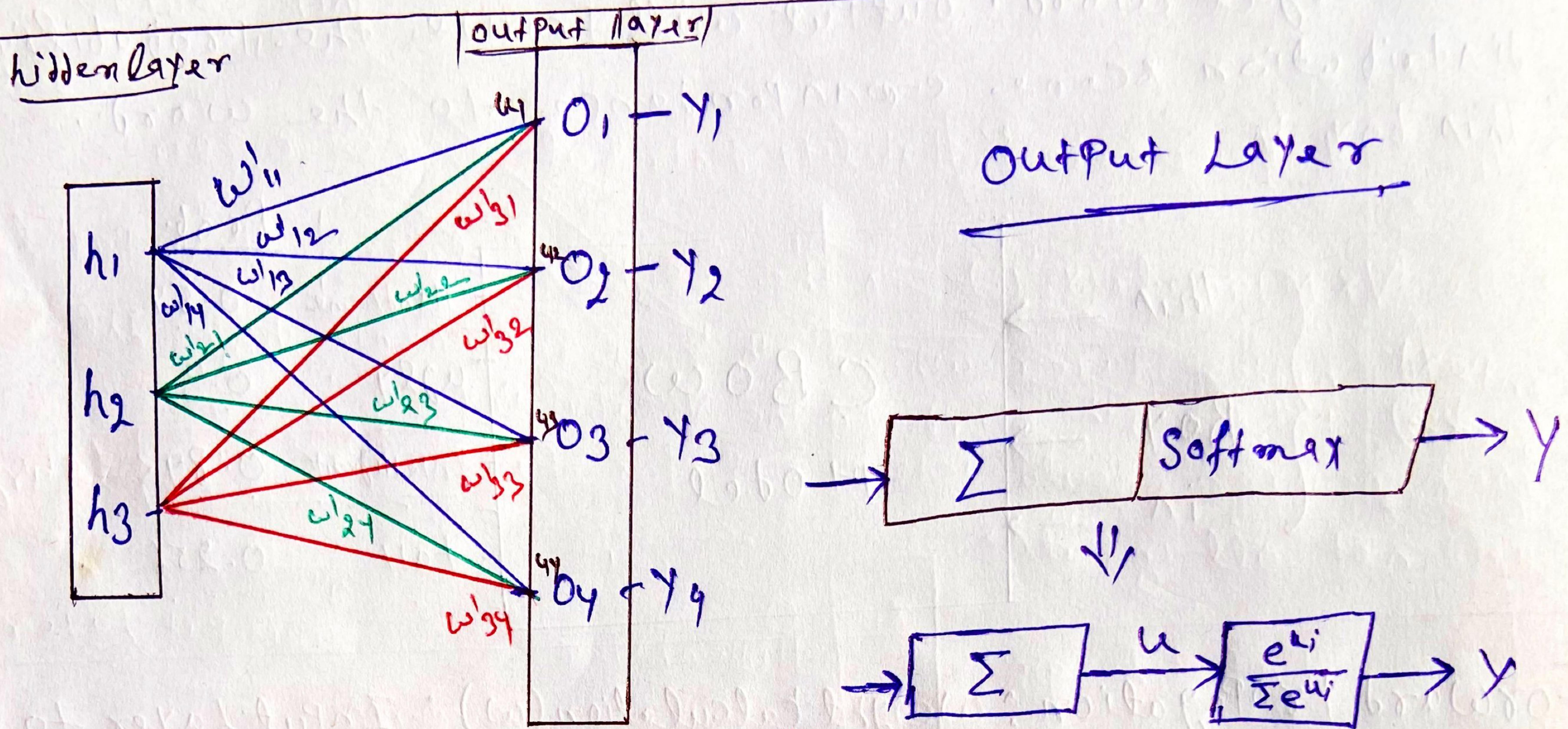
$N=3$

$\frac{V}{N} = 1$

$$h = w^T X$$

(1)

Forward Calculation: Weight Calculation (w'): Hidden to O/P Layer



$$u_1 = w'_{11} h_1 + w'_{21} h_2 + w'_{31} h_3$$

$$u_2 = w'_{12} h_1 + w'_{22} h_2 + w'_{32} h_3$$

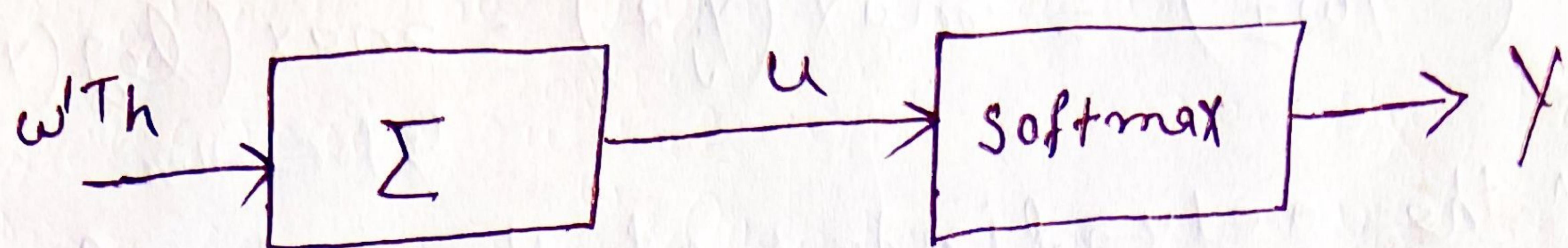
$$u_3 = w'_{13} h_1 + w'_{23} h_2 + w'_{33} h_3$$

$$u_4 = w'_{14} h_1 + w'_{24} h_2 + w'_{34} h_3$$

$$\begin{bmatrix} u_1 \\ u_2 \\ u_3 \\ u_4 \end{bmatrix} = \begin{bmatrix} w'_{11} & w'_{21} & w'_{31} \\ w'_{12} & w'_{22} & w'_{32} \\ w'_{13} & w'_{23} & w'_{33} \\ w'_{14} & w'_{24} & w'_{34} \end{bmatrix} \begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix}$$

$u = w'^T h$

(2)



Softmax function :-

$$f(x_i) = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$$

$$y_1 = \frac{e^{u_1}}{e^{u_1} + e^{u_2} + e^{u_3} + e^{u_4}}$$

$$y_2 = \frac{e^{u_2}}{e^{u_1} + e^{u_2} + e^{u_3} + e^{u_4}}$$

The Softmax output for the word w_j with respect to the given input context word w_I can be given by Conditional Probability, in terms of exponentials, i.e.

$$P(w_j | w_I) = y_j = \frac{e^{u_j}}{\sum_{j=1}^n e^{u_j}} \quad (3)$$

Error Calculation

Assumption:- Let ' w_o ' is the actual output word, w_I is the given input context words and ' v ' is the size of the input context.

The objective is to maximize the Conditional Probability of observing the actual output word ' w_o ', given the input context words ' w_I ', thus the loss function can be defined as:-

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

where, j^* is the ~~actual~~ index of the actual output word in the output layer.

Ex

Example if 4th word represent the O/P

$$E(Y_4) = \log P(w_{Y4}/w_I) = \log\left(\frac{e^{u_4}}{e^{u_1} + e^{u_2} + e^{u_3} + e^{u_4}}\right)$$

$$E(Y_4) = u_4 - \log(e^{u_1} + e^{u_2} + e^{u_3} + e^{u_4})$$

As we want to minimize the error, λE .

$$L = -\log P(w_0/w_I)$$

$$E(Y_4) = \log(e^{u_1} + e^{u_2} + e^{u_3} + e^{u_4}) - u_4$$

We can generalize it as:-

$$L = \log \sum_{j=1}^V e^{u_j} - u_{j^*}$$

It is the special case of cross-entropy measurement b/w two probabilistic distribution w_i and u_j

taking derivative L w.r.t. u_4 , we have

$$\frac{dE(Y_4)}{du_4} = \frac{e^{u_4}}{e^{u_1} + e^{u_2} + e^{u_3} + e^{u_4}} - \frac{d(u_4)}{d(u_4)}$$

$$\begin{aligned} \frac{d \log x}{dx} &= \frac{1}{x} \\ \frac{de^x}{dx} &= e^x \end{aligned}$$

$$\frac{dE(Y_4)}{du_4} = Y_4 - 1$$

We can generalize the above derivation.

$$\frac{dE}{dw_j} = \frac{d(\log \sum_{j=1}^V e^{u_j} - u_{j^*})}{du_j} = Y_j - t_j = e_j$$

Note:- $t_j = 1$, if ($t_j = 1$); else $t_j = 0$

Back-Propagation: Output to hidden layer (update w')

In this back-propagation we, update the weight of all neurons from hidden layer to output layer.

Step 1: Taking gradient ' E ' w.r.t; w'_{ii}

$$\frac{\partial E(y_i)}{\partial w'_{ii}} = \frac{\partial E(y_i)}{\partial u_i} \times \frac{\partial u_i}{\partial w'_{ii}}$$

$$\frac{\partial E(y_i)}{\partial u_i} = (y_i - o) = e_i \Rightarrow \text{from eqn (5)}$$

$$\frac{\partial u_i}{\partial w'_{ii}} = \frac{\partial (\omega'_{ii}h_1 + \omega'_{ii}h_2 + \omega'_{ii}h_3)}{\partial w'_{ii}} = h_i \quad \left[\begin{array}{l} u_i = \omega'_{ii}h_1 + \omega'_{ii}h_2 \\ \quad \quad \quad + \omega'_{ii}h_3 \end{array} \right]$$

$$\therefore \frac{\partial E(y_i)}{\partial w'_{ii}} = \frac{\partial E(y_i)}{\partial u_i} \times \frac{\partial u_i}{\partial w'_{ii}}$$

$$\boxed{\frac{\partial E(y_i)}{\partial w'_{ii}} = e_i \cdot h_i} \longrightarrow (6)$$

Step 2: Updating the weight of w'_{ii}

$$(w'_{ii})^{\text{new}} = w'_{ii} - \eta \frac{\partial E(y_i)}{\partial w'_{ii}} = w'_{ii} - \eta (e_i \cdot h_i)$$

Similarly, we can update the weight of $w'_{i2}, w'_{i3}, \dots, w'_{iy}$

$$\boxed{(w'_{ij})^{\text{new}} = (w'_{ij})^{\text{old}} - \eta e_j \cdot h_i} \longrightarrow (7)$$

Where, η = learning rate

$$e_j = y_j - t_j$$

Back-Propagation: Hidden to Input layer (update w)

In this back propagation we, update the weight of all neurons from hidden layer to output layer.

Step 1: Taking gradient 'E' w.r.t. w_{11}

$$\frac{\partial E}{\partial w_{11}} = \frac{\partial E}{\partial h_1} \times \frac{\partial h_1}{\partial w_{11}}$$

$$\frac{\partial E}{\partial h_1} = \left(\frac{\partial E}{\partial u_1} \times \frac{\partial u_1}{\partial h_1} \right) + \left(\frac{\partial E}{\partial u_2} \times \frac{\partial u_2}{\partial h_1} \right) + \left(\frac{\partial E}{\partial u_3} \times \frac{\partial u_3}{\partial h_1} \right) + \left(\frac{\partial E}{\partial u_4} \times \frac{\partial u_4}{\partial h_1} \right)$$

$$\begin{aligned} \frac{\partial E}{\partial h_1} &= (e_1 \cdot w_{11}) + (e_2 w_{12}) + (e_3 w_{13}) + (e_4 w_{14}) \\ &= \sum_{j=1}^4 e_j w_{1j} = EH_i \end{aligned} \quad \text{--- (9)}$$

$$\frac{\partial h_1}{\partial w_{11}} = \frac{\partial (w_{11}x_1 + w_{12}x_2 + w_{13}x_3 + w_{14}x_4)}{\partial w_{11}} = x_1 \quad \text{--- (10)}$$

$$\begin{aligned} \frac{\partial E}{\partial w_{11}} &= \frac{\partial E}{\partial h_1} \times \frac{\partial h_1}{\partial w_{11}} \\ &= \sum_{j=1}^4 e_j w_{1j} \cdot x_1 = EH_i \cdot x_1 \end{aligned}$$

$$\therefore \frac{\partial E}{\partial w_{kj}} = EH_i \cdot x_k$$

— (11)

Step 2: Updating the weight of w_{11}

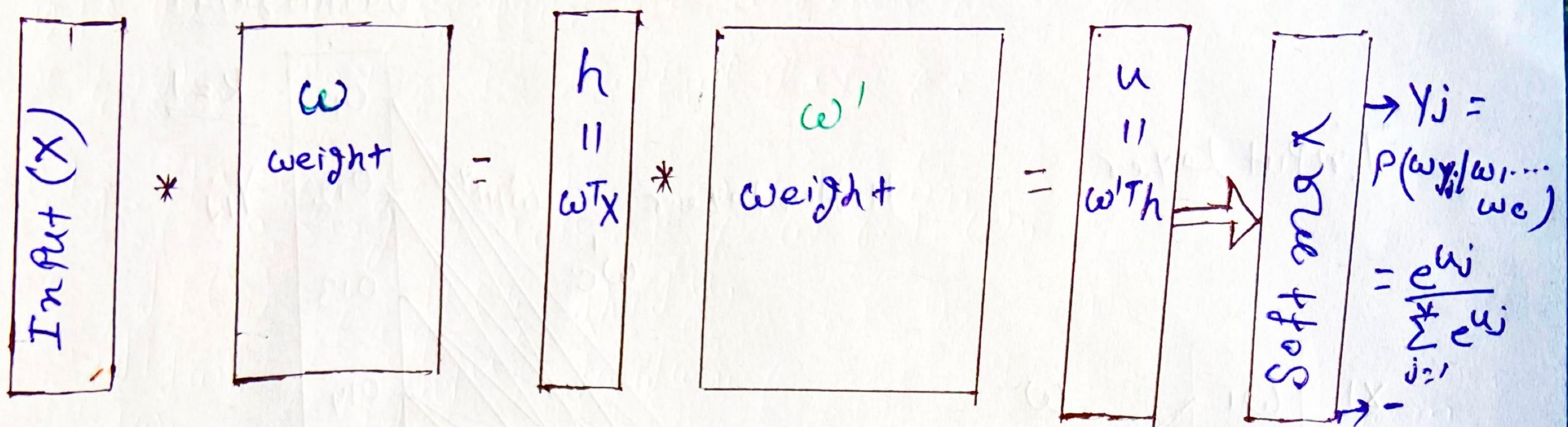
$$(w_{11})^{\text{new}} = (w_{11})^{\text{old}} - \eta \frac{\partial E}{\partial w_{11}} = (w_{11})^{\text{old}} - \eta EH_i \cdot x_1$$

$$(w_{ij})^{\text{new}} = (w_{ij})^{\text{old}} - \eta EH_i \cdot x_i \quad \text{--- (12)}$$

where,

η = learning rate
 $EH_i = \sum_{j=1}^4 e_j w_{ij}$

CBOW Model



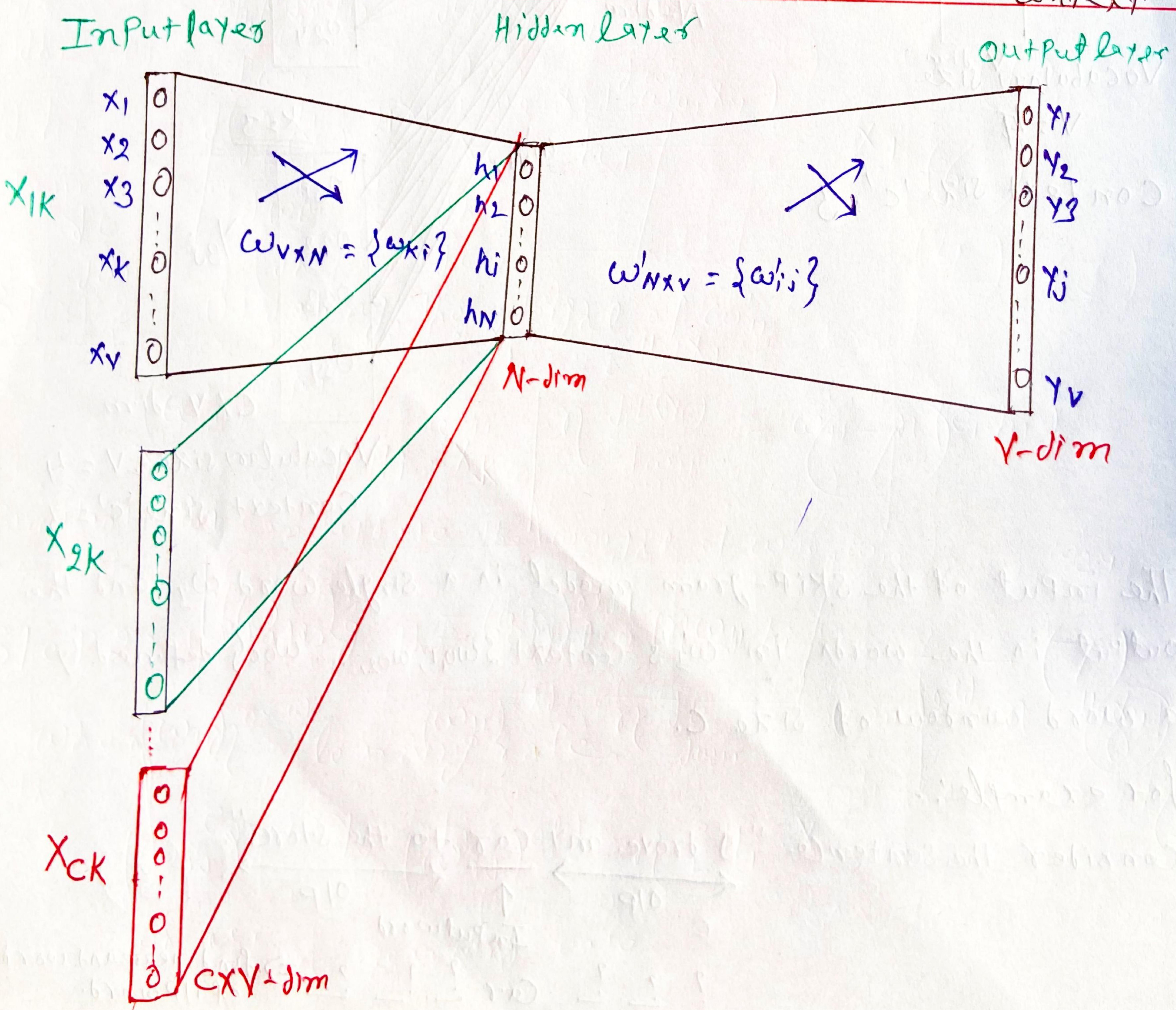
$$(w)^{\text{new}} = (w)^{\text{old}} - \eta_{EH} x$$

$$(w')^{\text{new}} = (w')^{\text{old}} - \eta_{eh}$$

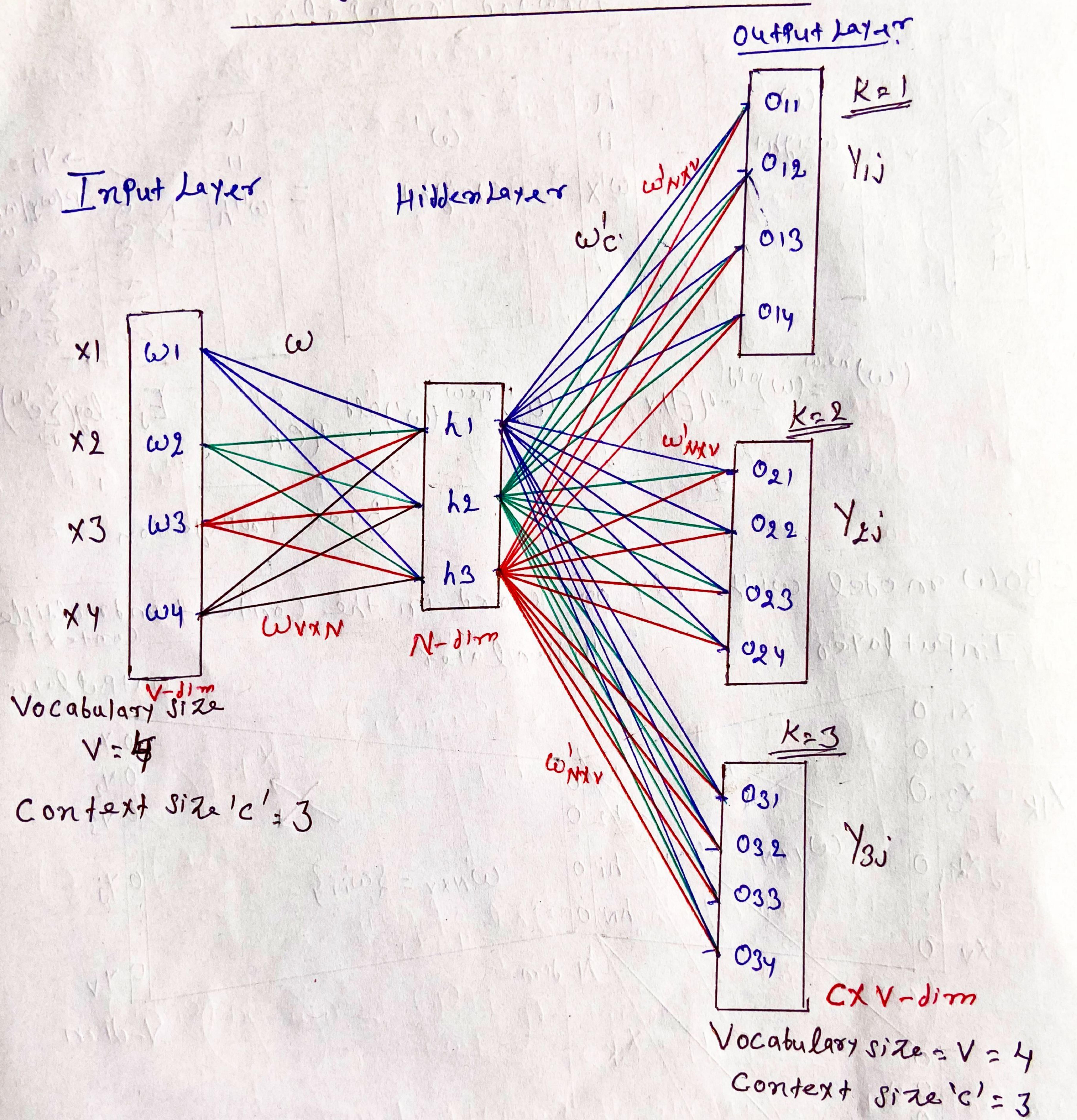
$$E_j = \log(\sum_{i=1}^V e^{u_i}) - q_j^*$$

Back-propagation

CBOW model with only one word in the context and multiple contexts



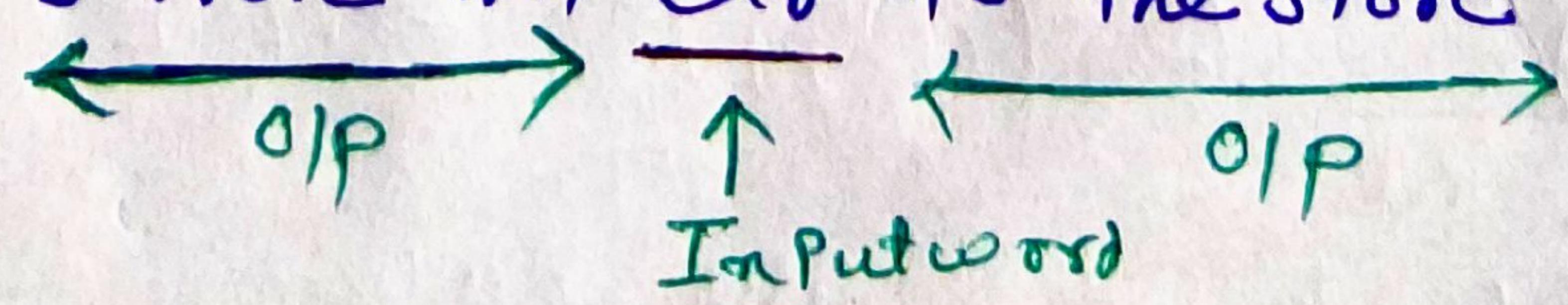
Skip-Gram Model



The input of the Skip-Gram model is a single word w_I and the output is the words in w_I 's context $\{w_{01}, w_{02}, \dots, w_{0c}\}$ defined by a word window of size c .

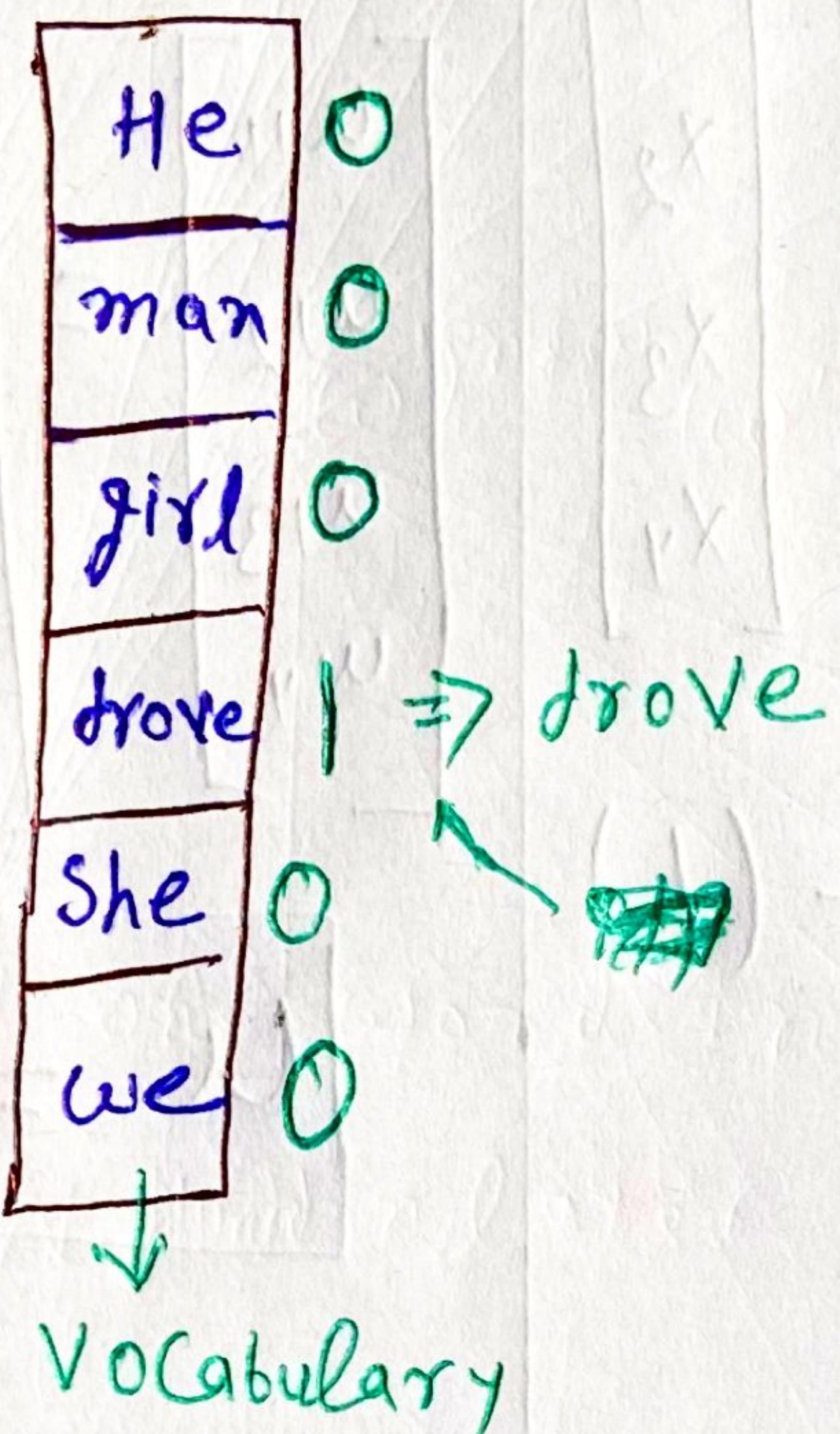
for example:-

Consider the sentence "I drove my car to the store".



? ? ? car ? ! ? \Rightarrow find nearest word of i/p word.

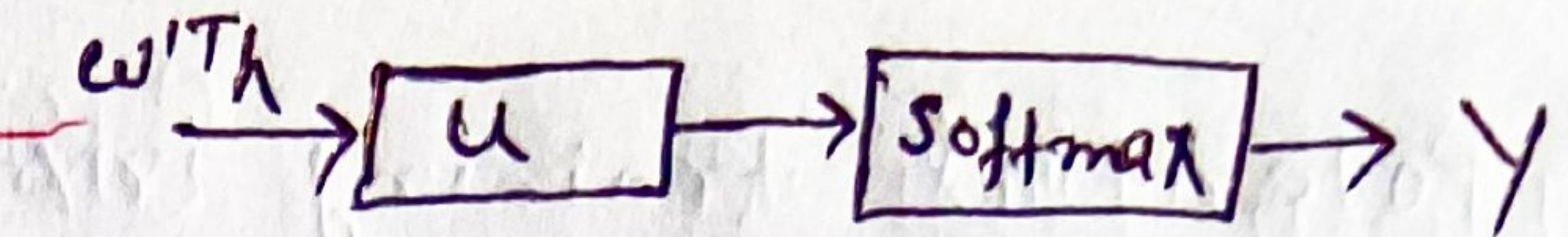
A potential training instance could be the word "Car" as an input and the words $\{\text{"I"}, \text{"drove"}, \text{"my"}, \text{"to"}, \text{"the"}, \text{"store}\}$ as output. All of these words are one-hot encoded meaning they are vector of length V (the size of the vocabulary) with a value 1 at the index ~~corresponding~~ corresponding to the word and zero or 0 in all other indexes.



In this model X represents the one-hot encoded vector corresponding to the input word in the training instance and $\{y_1, y_2, \dots, y_c\}$ are the one-hot encoded vectors corresponding to the output words in the training instance. The $V \times N$ matrix W is the weight matrix between the input layer and hidden layer where j th row represents the weights corresponding to the j th word in the vocabulary. Each output word vector also has an associated $N \times V$ out matrix W' . There is also a hidden layer consisting of N nodes.

$$W = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \\ w_{41} & w_{42} & w_{43} \end{bmatrix} \quad \text{Input \& hidden layer} \quad W' = \begin{bmatrix} w'_{11} & w'_{12} & w'_{13} & w'_{14} \\ w'_{21} & w'_{22} & w'_{23} & w'_{24} \\ w'_{31} & w'_{32} & w'_{33} & w'_{34} \end{bmatrix} \quad \text{Weight b/w hidden \& O/P layer.}$$

Forward Propagation: weight calculation



Input layer to hidden layer

$$h_1 = w_{11}x_1 + w_{21}x_2 + w_{31}x_3 + w_{41}x_4$$

$$h_2 = w_{12}x_1 + w_{22}x_2 + w_{32}x_3 + w_{42}x_4$$

$$h_3 = w_{13}x_1 + w_{23}x_2 + w_{33}x_3 + w_{43}x_4$$

$$\begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} = \begin{bmatrix} w_{11} & w_{21} & w_{31} & w_{41} \\ w_{12} & w_{22} & w_{32} & w_{42} \\ w_{13} & w_{23} & w_{33} & w_{43} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}$$

$$h = w^T x \quad \text{--- (1)}$$

Hidden layer to output layer

$$u_{11} = w'_{11}h_1 + w'_{21}h_2 + w'_{31}h_3$$

$$u_{12} = w'_{12}h_1 + w'_{22}h_2 + w'_{32}h_3$$

$$u_{13} = w'_{13}h_1 + w'_{23}h_2 + w'_{33}h_3$$

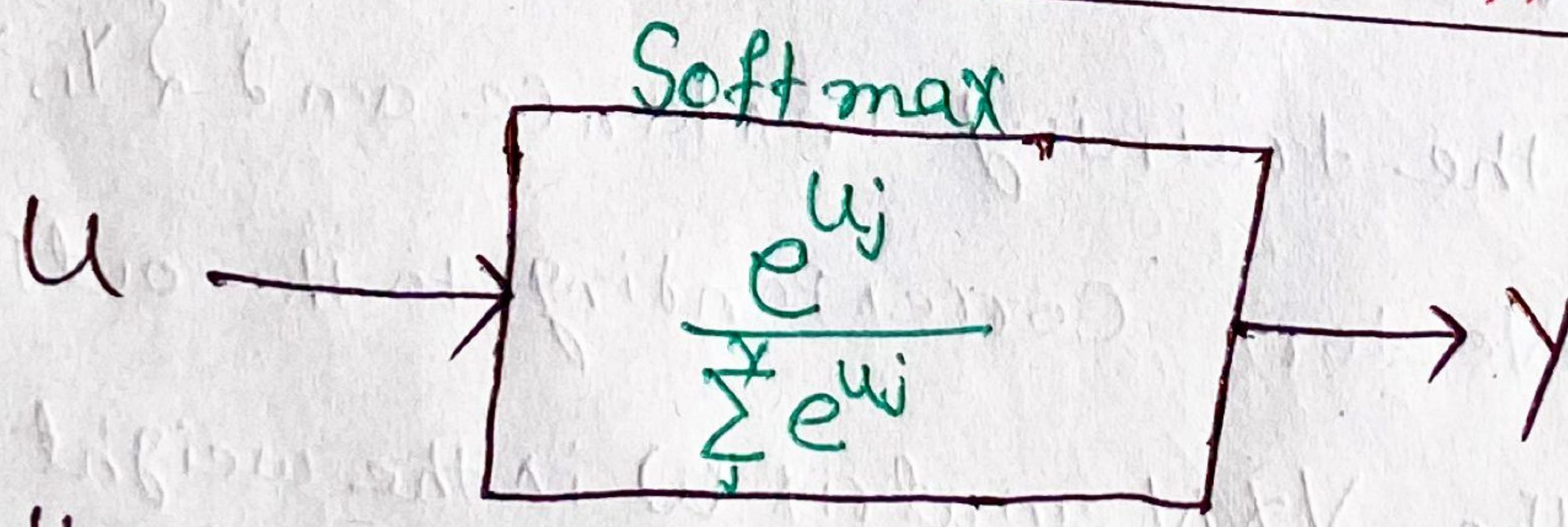
$$u_{14} = w'_{14}h_1 + w'_{24}h_2 + w'_{34}h_3$$

$$\begin{bmatrix} u_{11} \\ u_{12} \\ u_{13} \\ u_{14} \end{bmatrix} = \begin{bmatrix} w'_{11} & w'_{21} & w'_{31} \\ w'_{12} & w'_{22} & w'_{32} \\ w'_{13} & w'_{23} & w'_{33} \\ w'_{14} & w'_{24} & w'_{34} \end{bmatrix} \begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix}$$

$$u = w'^T h \quad \text{--- (2)}$$

$$w'^T \cdot h \rightarrow \sum \rightarrow u$$

Forward Propagation: Calculating Softmax output



$$y_{11} = \frac{e^{u_{11}}}{e^{u_{11}} + e^{u_{12}} + e^{u_{13}} + e^{u_{14}}}$$

$$y_{12} = \frac{e^{u_{12}}}{e^{u_{11}} + e^{u_{12}} + e^{u_{13}} + e^{u_{14}}}$$

The Softmax output for the word w_j with respect to the given output ~~word~~ word w_I can be given by Conditional Probability, i.e.

$$P(w_j/w_I) = y_j = \frac{e^{u_j}}{\sum_{j=1}^V e^{u_j}} \quad \text{--- (3)}$$

Error Calculation

Assumption: Let, $\{w_{1j}, w_{2j}, w_{3j} \dots w_{Cj}\}$ are actual output words, for the given single input word w_I . where, w_{ij} is the j th word of 1^{st} context window, ~~and~~

The Training objective is to maximize the conditional probability of observing the actual output words, $\{w_{1j}, w_{2j}, w_{3j} \dots w_{Cj}\}$ given the single input context word w_I . Additionally, we want to minimize the error, so the loss function can be defined by Conditional Probability, in term of exponentials, i.e.

$$E = -\log P(w_{1j}, w_{2j}, \dots, w_{Cj} | w_I) = -\log \prod_{c=1}^C \left[\frac{e^{u_{cj}}}{e^{u_{1j}} + e^{u_{2j}} + \dots + e^{u_{Cj}}} \right]$$

$$\therefore E = -\log \prod_{c=1}^C \left[\frac{e^{u_{cj}}}{\sum_{j=1}^J e^{u_{1j}}} \right] \quad (4)$$

$$E = C \cdot \log \sum_{j=1}^J e^{u_j} - \sum_{c=1}^C u_{cj} \quad (4)$$

where,

C = total Context window

J = Vocabulary size

Now, taking derivative, w.r.t,

$$\frac{\partial E}{\partial u_{cj}} = \frac{e^{u_{cj}}}{\sum_{j=1}^J e^{u_j}} - Z_{cj} = Y_{cj} - Z_{cj} \quad (5)$$

Where,

Y_{cj} = out score of j th word of c^{th} context window

$Z_{cj} = 1$, if the j th word of c^{th} context window is the actual output word, otherwise, we put $Z_{cj} = 0$

Back-propagation: Output to Hidden layer (update w')

In this back-propagation phase we, update the weight of all neurons from hidden layer to output layer. (w'_{11}, w'_{12}, \dots)

Step1: Taking gradient of E' w.r.t, w'_{11}

$$\frac{dE(Y_1)}{dw'_{11}} = \sum_{c=1}^C \frac{\cancel{dE(Y_1)}}{dU_{c1}} \times \frac{dU_{c1}}{dw'_{11}} \quad \xrightarrow{\text{from eqn (5)}}$$

$$\sum_{c=1}^C \frac{dE(Y_1)}{dU_{c1}} = \cancel{\frac{dE(Y_1)}{dU_{c1}}} \sum_{c=1}^C (Y_{ci} - Z_{ci}) = e_1$$

$$\sum_{j=1}^c \frac{dU_{c1}}{dw'_{11}} = \sum_{j=1}^c \frac{d(w'_{11}h_1 + w'_{12}h_2 + w'_{13}h_3)}{dw'_{11}} = h_1$$

$$\therefore \frac{dE(Y_1)}{dw'_{11}} = \frac{dE(Y_1)}{du_1} \times \frac{du_1}{dw'_{11}} = e_1 \cdot h_1 \quad \boxed{6}$$

Step2: Updating the weight of w'_{11}

$$(w'_{11})^{\text{new}} = (w'_{11})^{\text{old}} - \eta \frac{dE(Y_1)}{dw'_{11}} \quad \left[\begin{aligned} \frac{dE}{dw'_{11}} &= \sum_{c=1}^C \frac{dE}{dU_{c1}} \times \frac{dU_{c1}}{dw'_{11}} \\ &= \sum_{c=1}^C (Y_{ci} - Z_{ci}) \cdot h_i \end{aligned} \right]$$

$$(w'_{11})^{\text{new}} = (w'_{11})^{\text{old}} - \eta (e_1 \cdot h_1) \quad \boxed{7}$$

$$(w_{ij})^{\text{new}} = (w_{ij})^{\text{old}} - \eta \sum_{c=1}^C (Y_{ci} - Z_{ci}) \cdot h_i \quad \boxed{7}$$

Back-propagation: Hidden to Input Layer (update w)

In this back-propagation phase we update the weight of all neurons from hidden layer to output layer.

Step 1: Taking gradient of E' w.r.t. h_i

$$\frac{\partial E}{\partial h_i} = \sum_{j=1}^J \frac{\partial E}{\partial u_{cj}} \times \frac{\partial u_{cj}}{\partial h_i} = \sum_{j=1}^J \sum_{c=1}^C (y_{cj} - z_{cj}) \cdot w'_{ij}$$

Step 2: Calculating gradient of E' , w.r.t. w_{ii}

$$\frac{\partial E}{\partial w_{ii}} = \frac{\partial E}{\partial h_i} \times \frac{\partial h_i}{\partial w_{ii}} = \sum_{j=1}^J \sum_{c=1}^C (y_{cj} - z_{cj}) \cdot w'_{ij} \times (x_i)$$

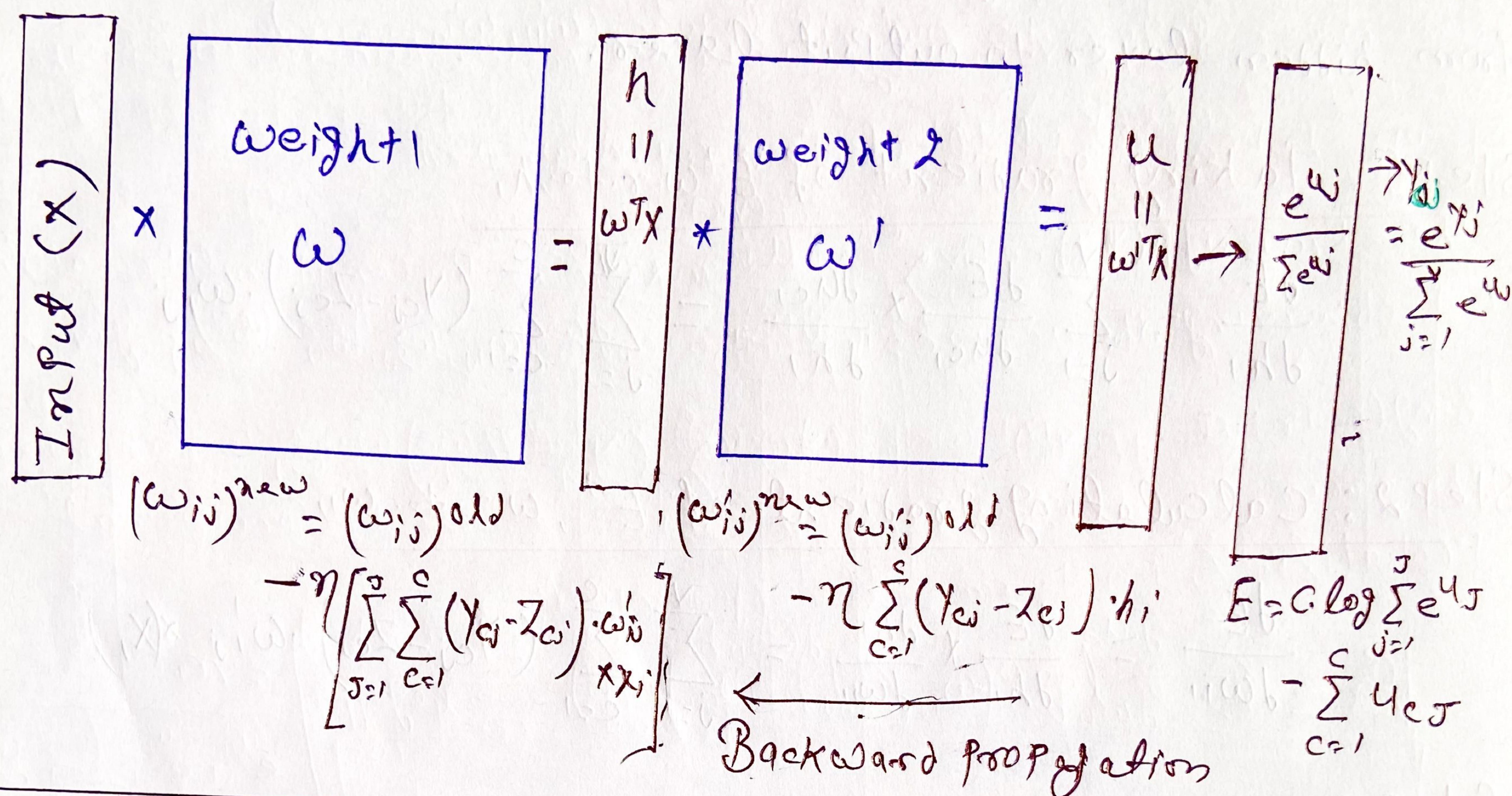
Step 3: updating the weight of w_{ii}

$$\begin{aligned} (w_{ii})^{new} &= w_{ii} - \eta \frac{\partial E}{\partial w_{ii}} \\ &= w_{ii} - \eta \left(\sum_{j=1}^J \sum_{c=1}^C (y_{cj} - z_{cj}) \cdot w'_{ij} \times (x_i) \right) \end{aligned}$$

$$(w_{ij})^{new} = (w_{ij})^{old} - \eta \left(\sum_{j=1}^J \sum_{c=1}^C (y_{cj} - z_{cj}) \cdot w'_{ij} \times (x_i) \right) \quad (B)$$

Model

forward Propagation



Skip-gram

Input layer

