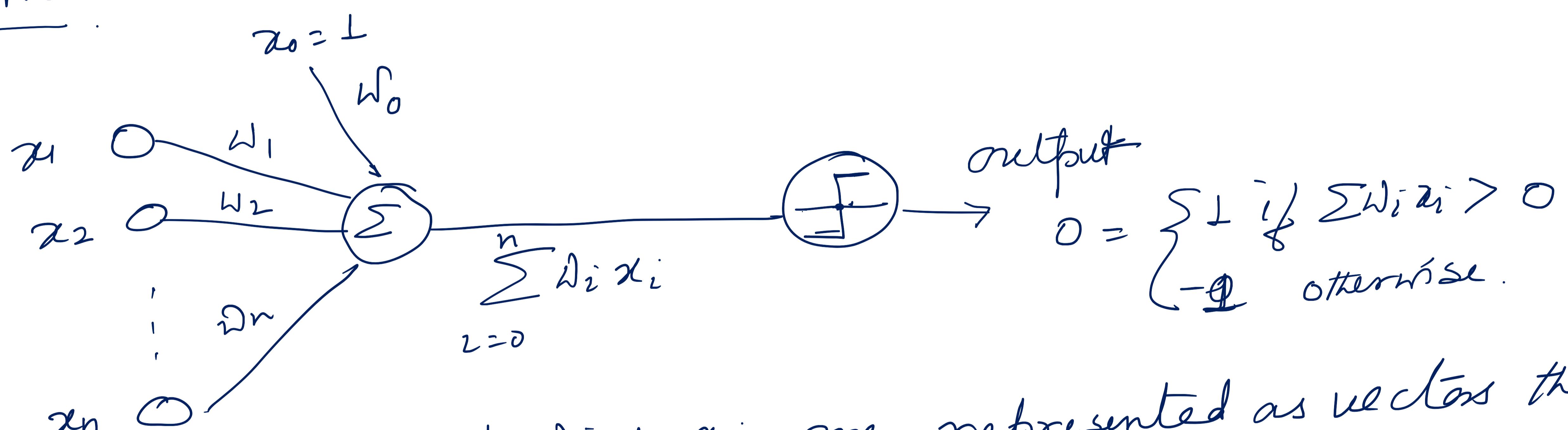


## Node Embeddings

### Neural Networks

→ Creating an automated logic function  
AND, OR, NAND, NOR

#### Perception



So if  $w_i$  &  $x_i$  are represented as vectors then

$$o(\vec{x}) = \text{sgn}(\vec{w} \cdot \vec{x})$$

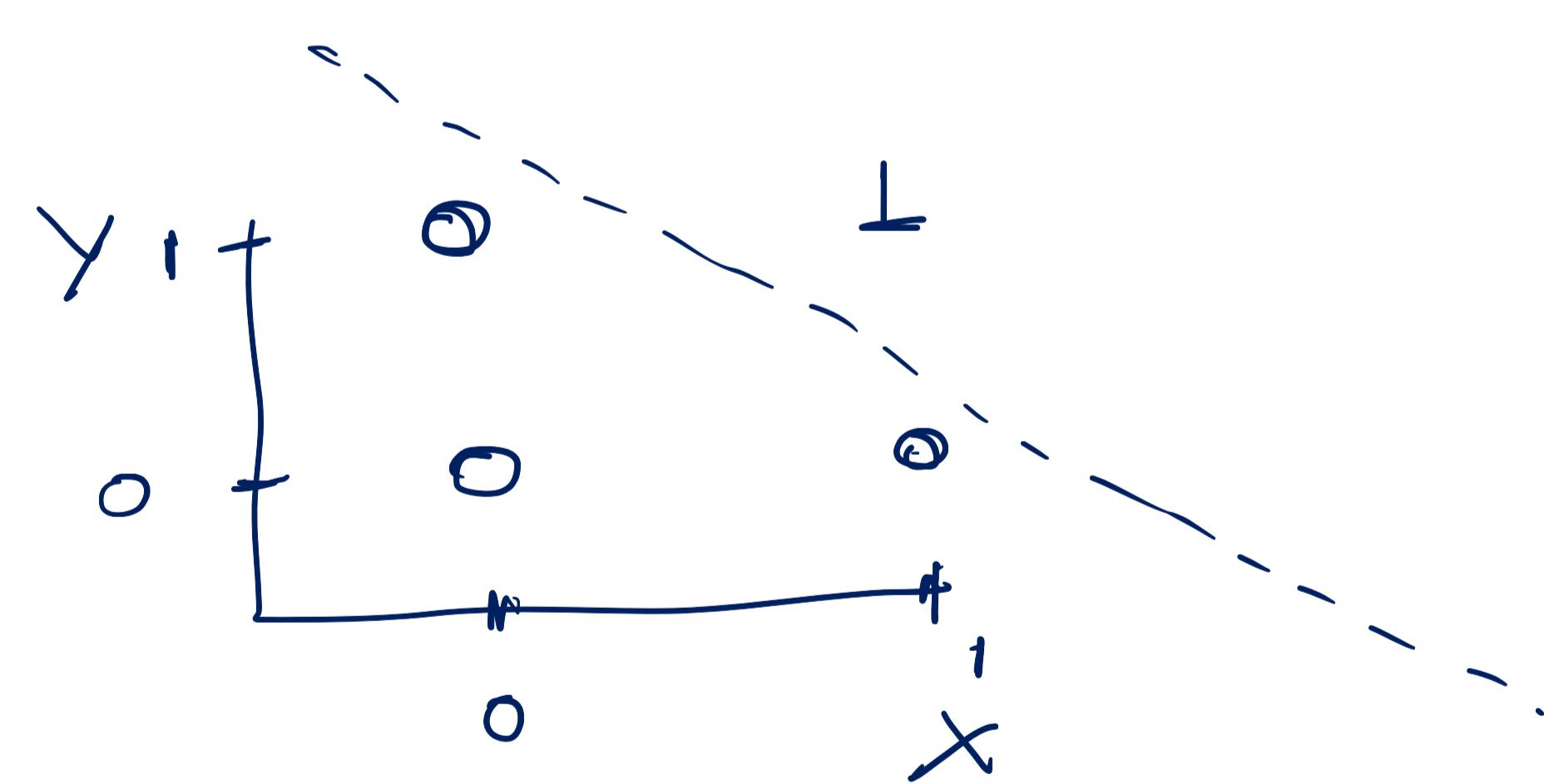
$$\text{sgn}(y) = \begin{cases} 1 & \text{if } y > 0 \\ 0 & \text{otherwise} \end{cases}$$

Learning the perception

- Involves learning  $w_0, w_1, \dots, w_n$

Consider AND

$x$	$y$	target ( $t$ )
0	0	0
0	1	0
1	0	0
1	1	1



How to learn the weight such that for the given  $\theta$  inputs, the outputs ~~result~~ is given as result.

Generally 2 ways — 1) ~~Delta rule~~ Perception training rule  
2) ~~Stochastic~~ Gradient Descent.

Delta rule Perceptron training rule.  
Start with some random initial weights  $w_0, w_1, w_2 \dots$

Weight  $w_i$  is revised as

$$w_i = w_i + \Delta w_i \text{ where } \Delta w_i = \eta(t - o)x_i$$

$\eta$  = learning rate

$t$  = target value

$o$  = output value returned by perceptron

$x_i$  = Input

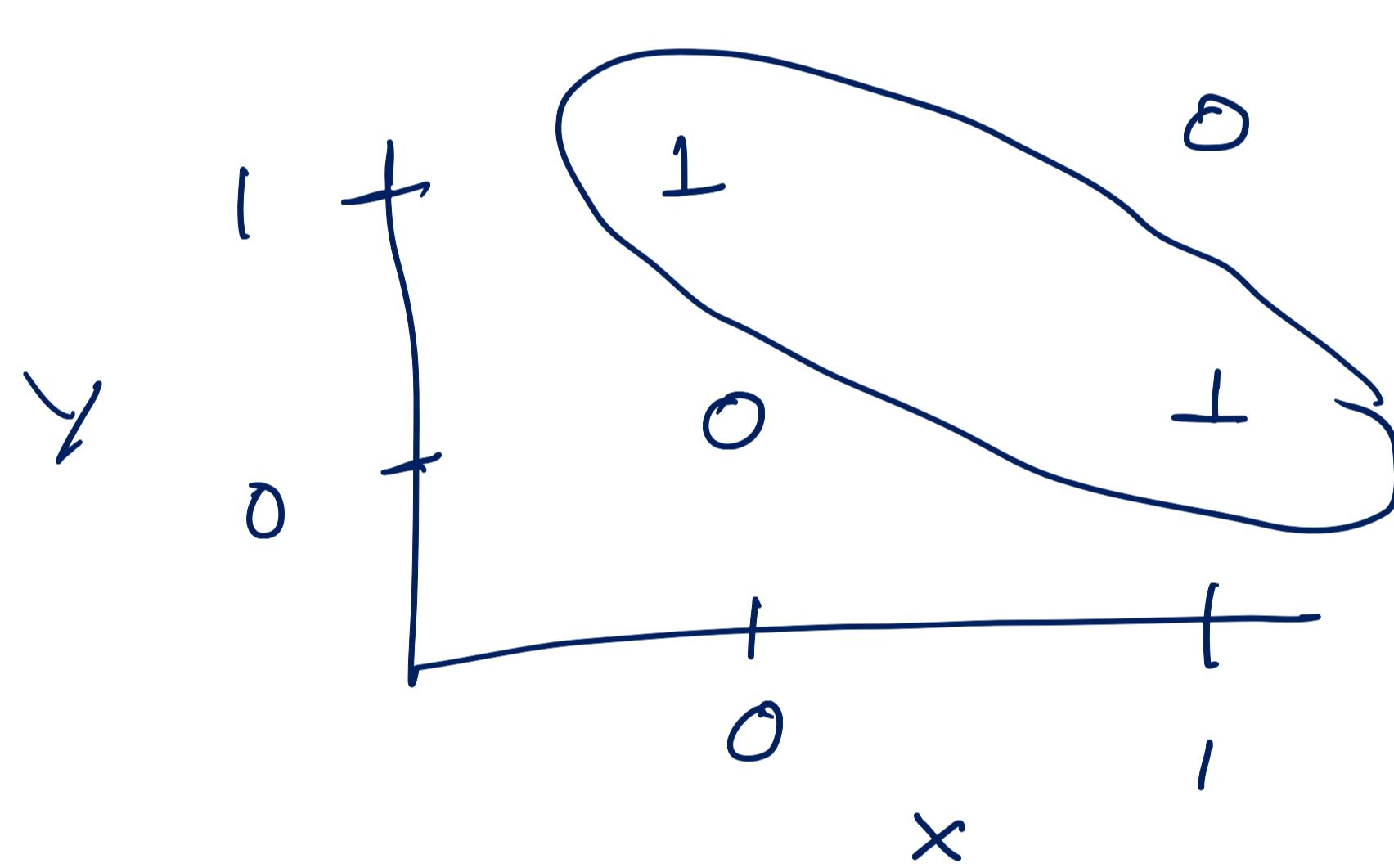
If  $t=0, o=0$  so  $\Delta w_i = 0$

If  $t=1, o=-1$  so  $t-o=2$  i.e.  $w_i$  has to be increased.

So the linear perceptron with delta rule only works for those logics that have a linear separation between the classes.

So it does not work for function like XOR

X	Y	T
0	0	0
0	1	1
1	0	1
1	1	0



For this we handling these cases use gradient descent.

- Need to specify a measure for training error of a hypothesis or a loss function

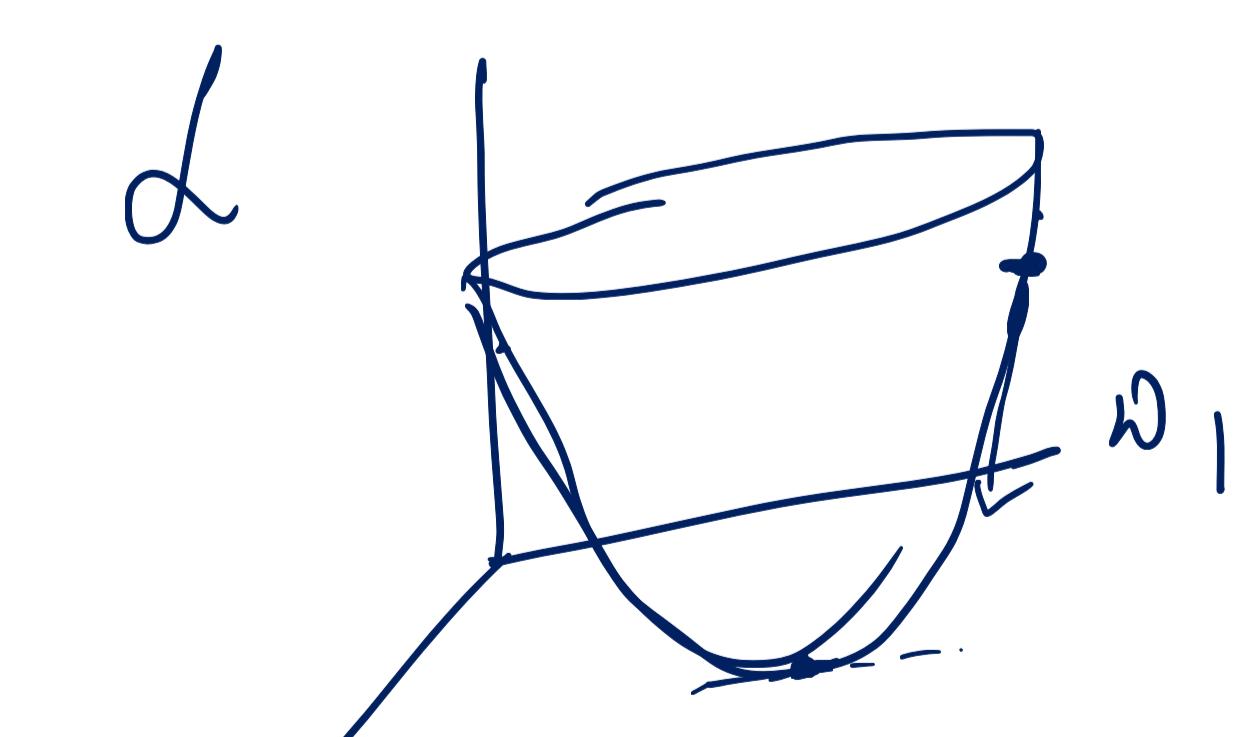
Gradient descent

$$L(\vec{w}) = \frac{1}{2} \sum_d (t_d - o_d)^2$$

$D$  = set of training examples

$o_d$  = output of perceptron for data instance  $d$

$t_d$  = Actual output



Start with some random  $\vec{w}$

Derive Gradient of  $L$ , with respect to  $w_0, w_1, w_2 \dots$   
This would be a vector that gives the direction of the steepest ascent.

So  $-\nabla L(\vec{w})$  gives the direction of the steepest descent.

So then update rule  $w_i = w_i - \eta \nabla L(\vec{w})$

$$L_{(1)} = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

$$\frac{\partial L}{\partial w_i} = \frac{1}{2} \sum_{d \in D} 2(t_d - o_d)(-x_{id})$$

$$= - \cancel{2} \sum_{d \in D} (t_d - o_d) (x_{id})$$

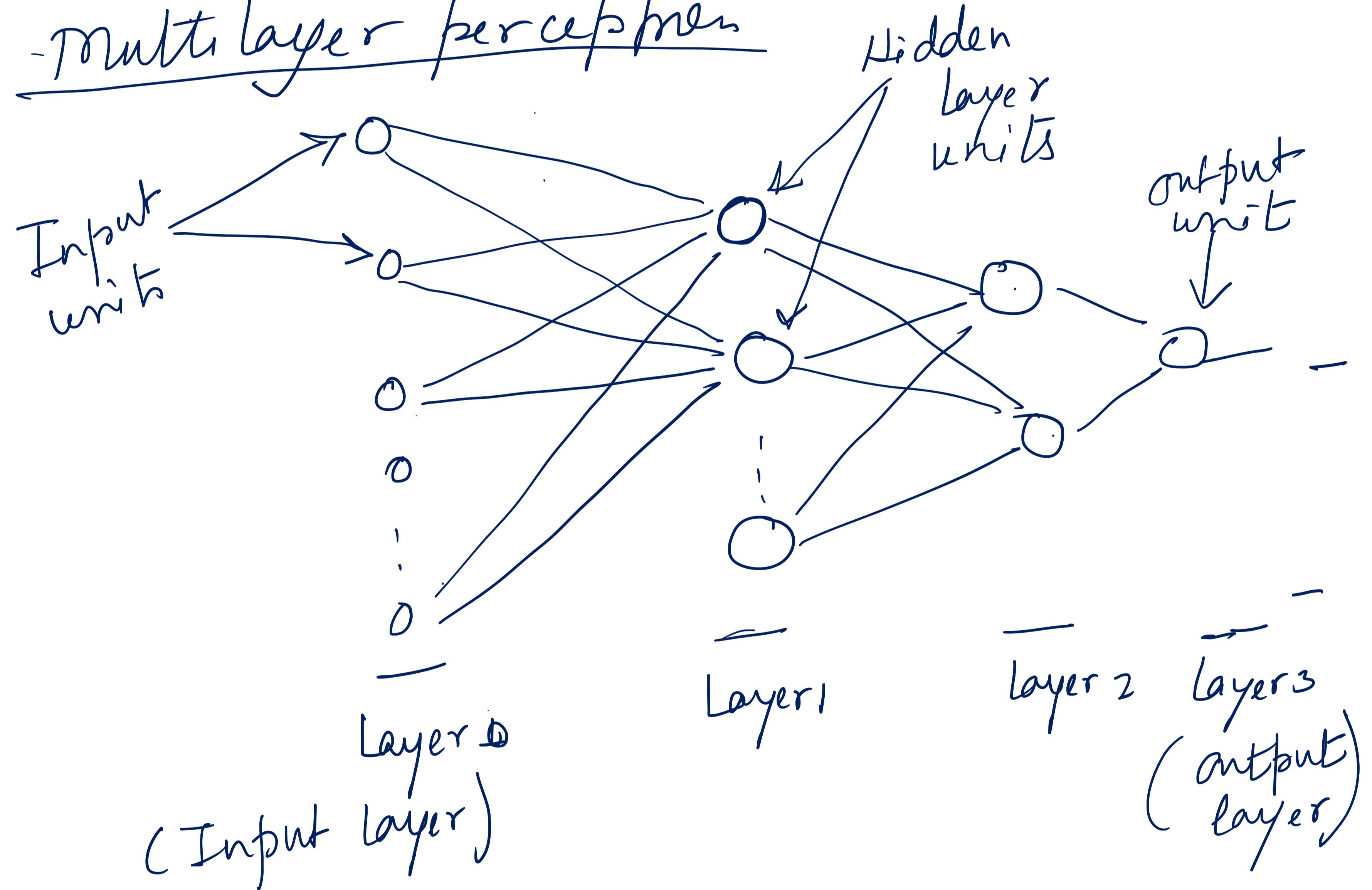
$$\begin{aligned} o_d &= \sum_j w_{ji} x_{ji} \\ \frac{\partial o_d}{\partial w_i} &= x_{ij} \end{aligned}$$

$$\text{So update is } w_i = w_i - \cancel{\alpha} \sum_{d \in D} (t_d - o_d) (x_{id})$$

But the  $\text{sgn}(x)$  function is non differentiable, so replace Sgn with a differentiable function like Sigmoid or tanh. That looks similar to sgn.

For Adding layers to the perceptron.  $\rightarrow$  Enhance the capacity of the model.

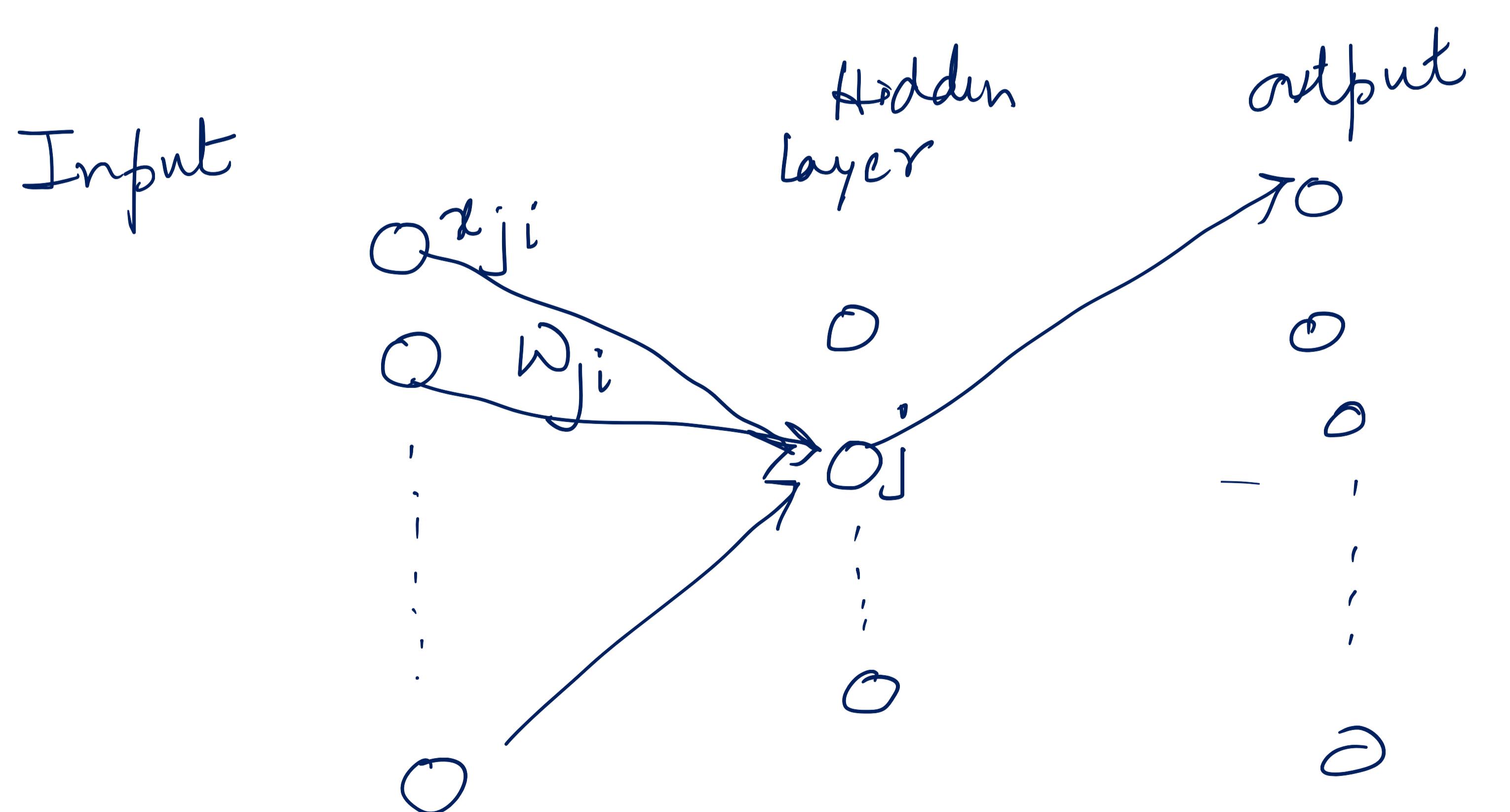
### Multi layer perceptrons



Advantage

- Arbitrary non linear class separation boundaries.

- Higher capacity.



### Symbols

$x_{ji}$  = The  $i^{\text{th}}$  input to the  $j^{\text{th}}$  unit

$w_{ji}$  = The weight associated with the  $i^{\text{th}}$  input to unit  $j$

$\text{net}_j = \sum_i w_{ji} x_{ji}$  (the weighted sum of inputs to unit  $j$ )

$o_j$  = output of the  $j^{\text{th}}$  unit

$t_j$  = target output for unit  $j$

$\sigma$  = sigmoid function

outputs = Set of units in the final layer of the network

$\text{Downstream}(j)$  = Set of units whose immediate inputs include the output of unit  $j$

~~Affine~~ Start with some arbitrary assignment of weights  $\theta_W = [ ]$

for each unit  $j$  in the downstream of the input layer do the following

$$a) \text{net}_j = \sum_i w_{ji} x_i$$

$$b) o_j = \sigma(\text{net}_j) = \frac{1}{1 + e^{-\text{net}_j}}$$

Repeat for further downstream nodes until an output is obtained from the output layer.

→ Feed forward function

↳ calculates the output of the neural network

For updating the weights, use back propagation  
i.e. the weights near to the output layer are updated first.

$$\text{for update} \quad w_{ji} = w_{ji} - \eta \frac{\partial L_d}{\partial w_{ji}}$$

$w_{ji}$  can influence the rest of the network through  $\text{net}_j$

$$\frac{\partial L_d}{\partial w_{ji}} = \frac{\partial L_d}{\partial \text{net}_j} \frac{\partial \text{net}_j}{\partial w_{ji}} = x_i \frac{\partial L_d}{\partial \text{net}_j}$$

Consider 2 cases

- a) where unit  $j$  is an output unit for the network
- b) where unit  $j$  is an internal unit.

Case a) Training rule for output unit weights

→  $w_{ji}$  can influence the rest of the network through  $\text{net}_j$  and  $o_j$   
can influence the rest of the network through  $o_j$

$$\frac{\partial L_d}{\partial \text{net}_j} = \frac{\partial L_d}{\partial o_j} \cdot \frac{\partial o_j}{\partial \text{net}_j} \quad o_j = \sigma(\text{net}_j) = \frac{1}{1 + e^{-\text{net}_j}}$$

$$L_d = \frac{1}{2} \sum_{k \in \text{outputs}} (t_k - o_k)^2$$

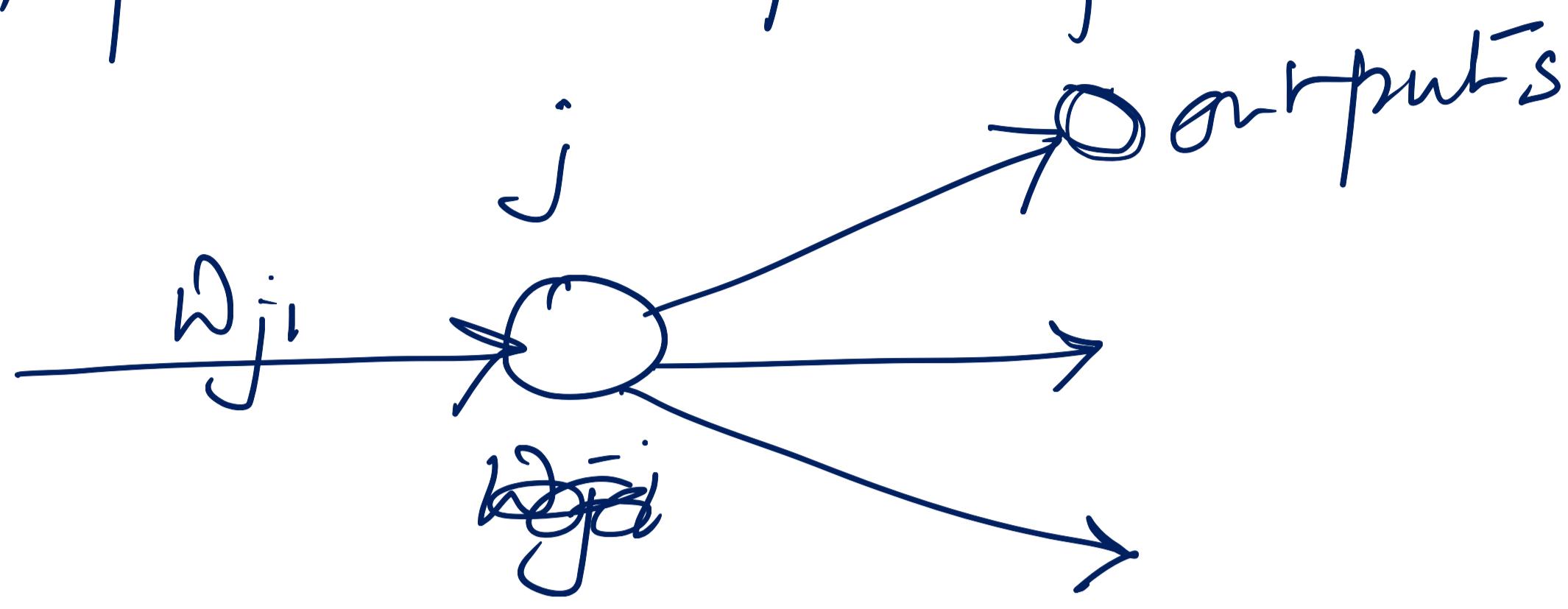
$$\frac{\partial L_d}{\partial o_j} = \frac{1}{2} \times 2(t_j - o_j)(-1) \\ = -(t_j - o_j)$$

$$\frac{\partial o_j}{\partial \text{net}_j} = \frac{e^{-\text{net}_j}}{(1 + e^{-\text{net}_j})^2} = \frac{e^{-\text{net}_j}}{(1 + e^{-\text{net}_j})} \cdot \left( \frac{1}{1 + e^{-\text{net}_j}} \right) = o_j(1 - o_j)$$

$$\frac{\partial L_d}{\partial \text{net}_j} = -(t_j - o_j) o_j(1 - o_j)$$

$$\text{So } \Delta w_{ji} = -\eta \frac{\partial L_d}{\partial w_{ji}} = \eta (t_j - o_j) o_j(1 - o_j) x_i$$

Case (b) Training rule for hidden unit weight  
 When unit  $j$  is a hidden layer unit  
 Then  $\omega_{ji}$  influences the output through net $_j$  of the hidden layer unit, that further influences input the output from the hidden layer



$$\frac{\partial L_d}{\partial \text{net}_j} = \sum_{k \in \text{downstream}(j)} \frac{\partial L_d}{\partial \text{net}_k} \cdot \frac{\partial \text{net}_k}{\partial \text{net}_j}$$

$$L_d = \frac{1}{2} \sum_{k \in \text{outputs}} (t_k - o_k)^2$$

$$\frac{\partial L_d}{\partial \text{net}_k} = \frac{\partial L_d}{\partial o_k} \cdot \frac{\partial o_k}{\partial \text{net}_k}$$

$$\frac{\partial o_k}{\partial \text{net}_k} = o_k(1-o_k)$$

$$\frac{\partial L_d}{\partial o_k} = -(t_k - o_k)$$

$$\frac{\partial L_d}{\partial \text{net}_k} = -(t_k - o_k)o_k(1-o_k)$$

$$\frac{\partial \text{net}_k}{\partial \text{net}_j} = \frac{\partial \text{net}_k}{\partial o_j} \cdot \frac{\partial o_j}{\partial \text{net}_j}$$

$$\frac{\partial o_j}{\partial \text{net}_j} = o_j(1-o_j)$$

$$\frac{\partial \text{net}_k}{\partial o_j} = \omega_{kj}$$

$$\text{so } \frac{\partial \text{net}_k}{\partial \text{net}_j} = \omega_{kj} o_j(1-o_j)$$

$$\text{so } \frac{\partial L_d}{\partial \text{net}_j} = \sum_{k \in \text{downstream}(j)} -(t_k - o_k)o_k(1-o_k)\omega_{kj} o_j(1-o_j)$$

$$= -o_j(1-o_j) \sum_{k \in \text{downstream}(j)} \delta_k \omega_{kj}, = \delta_j$$

$$\text{where } \delta_k = -(t_k - o_k)o_k(1-o_k)$$

$$\text{so } \Delta \omega_{ji} = -\eta \frac{\partial L_d}{\partial \omega_{ji}} = -\eta x_i \delta_j$$

The neural network is used to generate the representation of words appearing in a vocabulary. (Word2Vec)

# Continuous Bag of Word (CBOW) model

Skip gram model.

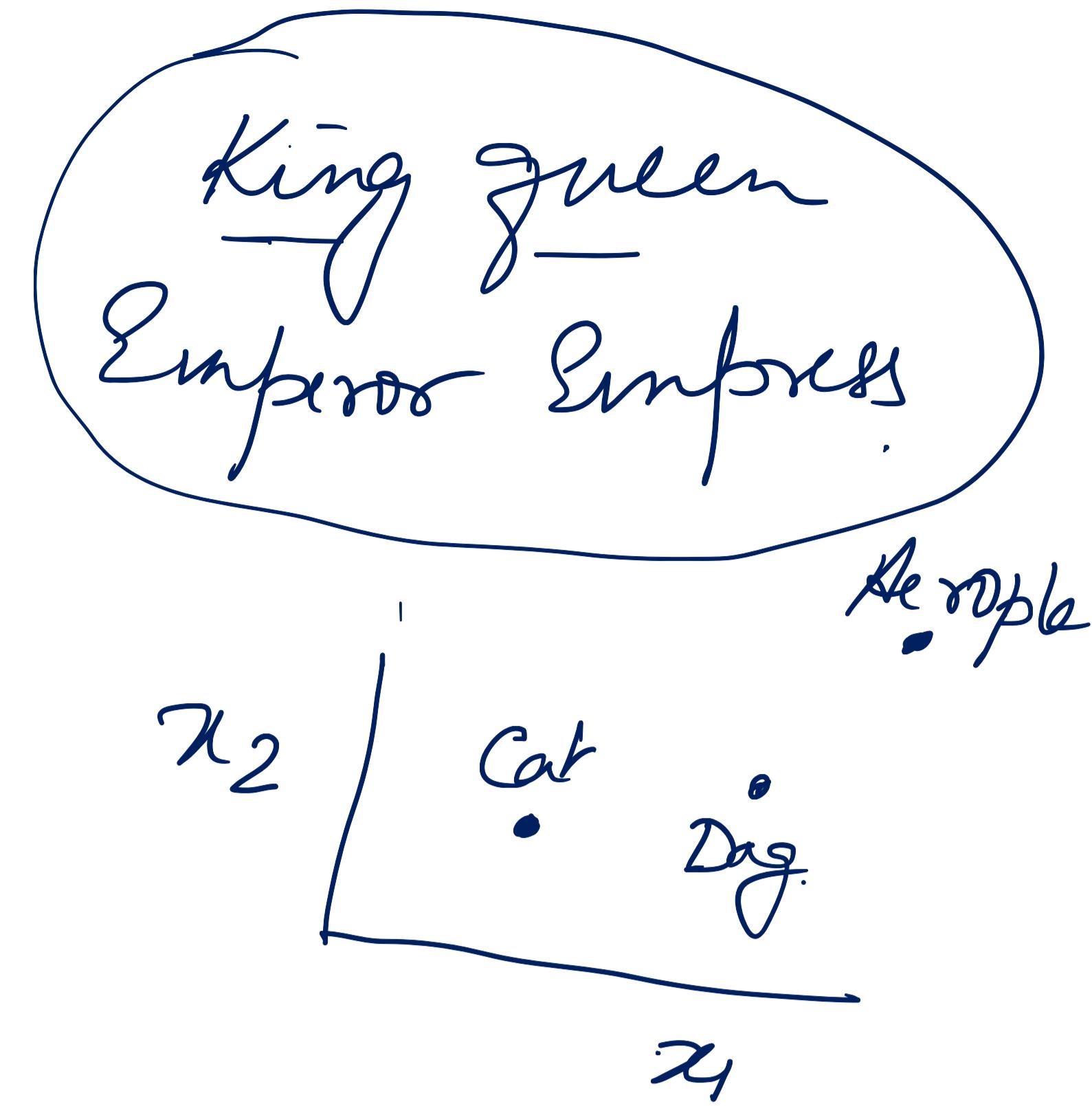
In CBOW model.

I am a good boy.

Suppose good is missing

I am a missing word boy

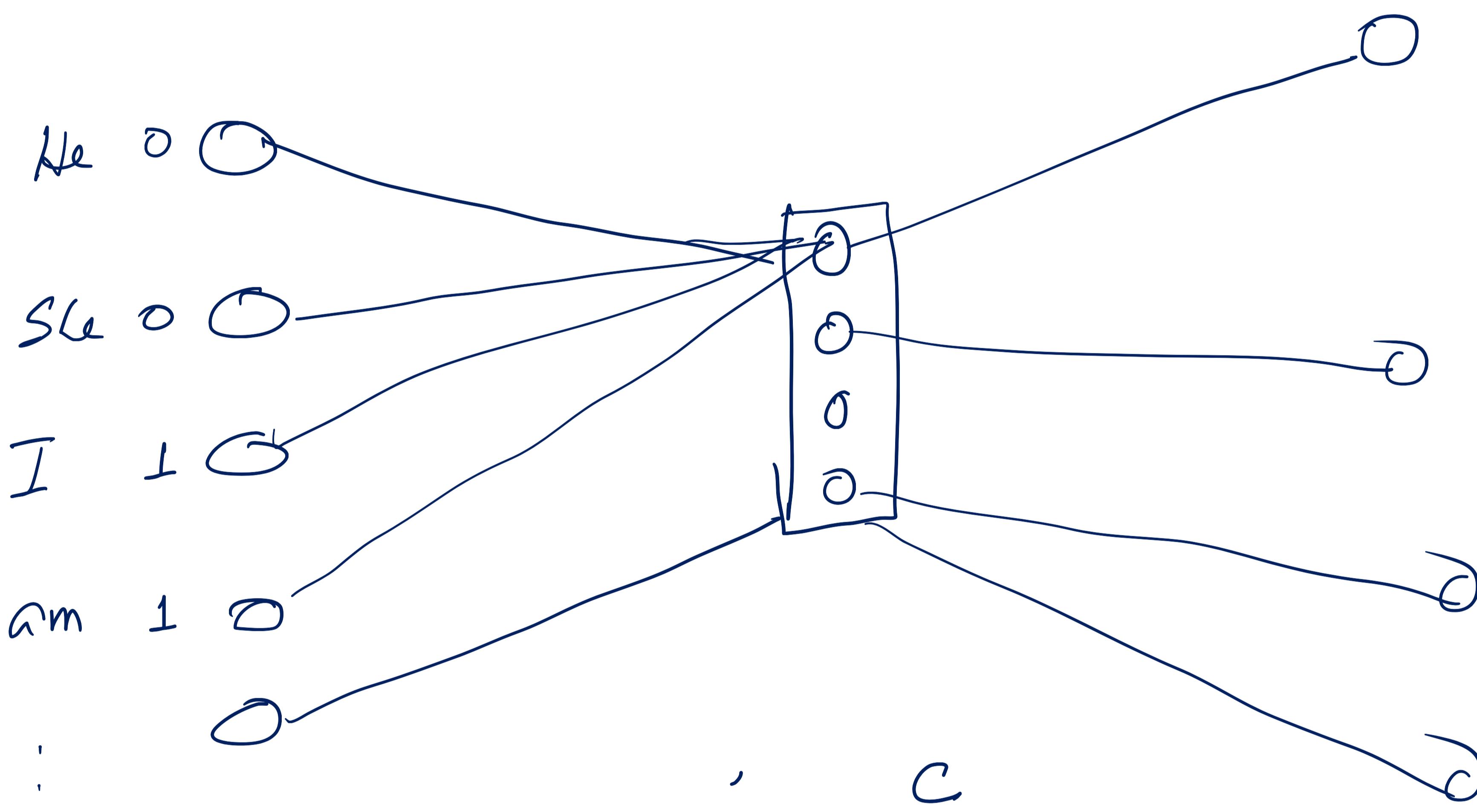
Given a set of context words  $c_1, c_2, c_3, c_4$ , predict the missing word  $w$ .



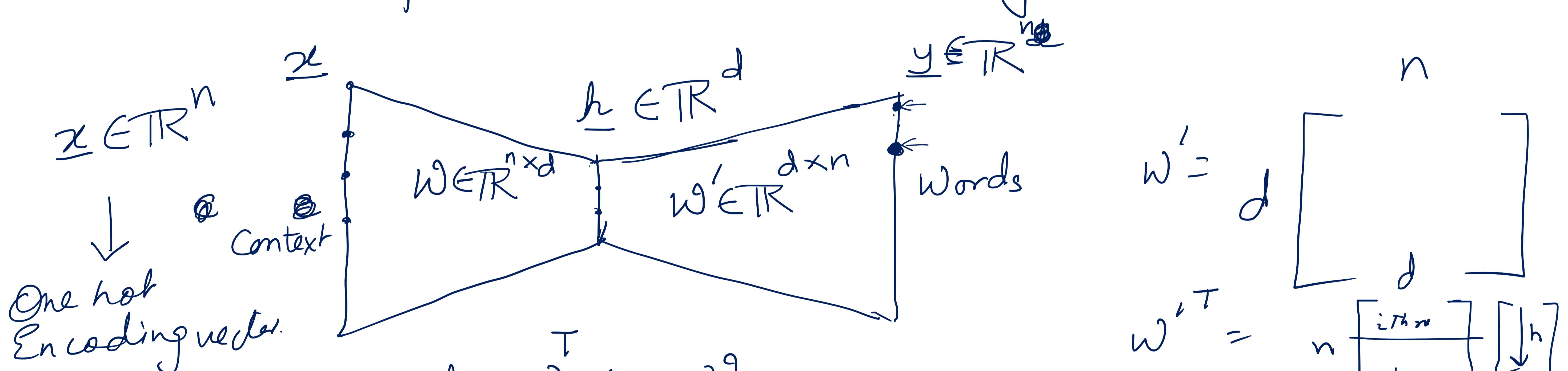
In skipgram model.

Given a word  $w$ , predict the context words.

one hot  
input  
vector



Assuming one context word is given  
Predict the word  $w$  adjacent to it.



$$h = W^T x = \underline{v}_c$$

Let  $\underline{u}$  be the input vector to the output units

So the input to the  $i$ th output unit

$$u_i = \underline{v}_w \cdot \underline{v}_c$$

The output of the  $i$ th output unit is generated using a softmax function

$$y_i = \text{softmax}(u_i) = \frac{e^{u_i}}{\sum_j e^{u_j}}$$

$$= \frac{e^{\underline{v}_w^T \underline{v}_c}}{\sum_{w'} e^{\underline{v}_{w'}^T \underline{v}_c}} = P(w | c) \quad \begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \dots \\ P(w_1|c) P(w_2|c) P(w_3|c)$$

To derive the parameters  $\theta = \{\underline{v}_c, \underline{v}_w\}$ , we need a loss function  
- Suitable loss function is maximum likelihood.

$$\text{The likelihood } L(\theta) = \prod_w P(w | c; \theta)$$

$$= \prod_w \frac{e^{\underline{v}_w^T \underline{v}_c}}{\sum_{w'} e^{\underline{v}_{w'}^T \underline{v}_c}}$$

$$- \log L(\theta) = -l(\theta) = -\sum_w \log \left[ \frac{e^{\underline{v}_w^T \underline{v}_c}}{\sum_{w'} e^{\underline{v}_{w'}^T \underline{v}_c}} \right]$$

$$= - \left[ \sum_w \underline{v}_w^T \underline{v}_c - \sum_w \log \left( \sum_{w'} e^{\underline{v}_{w'}^T \underline{v}_c} \right) \right]$$

$$\frac{\partial l(\theta)}{\partial \underline{v}_w} = \left[ \underline{v}_c - \frac{1}{\sum_{w'} e^{\underline{v}_{w'}^T \underline{v}_c}} \cdot e^{\underline{v}_w^T \underline{v}_c} \cdot \underline{v}_c \right]$$

$$= \left[ \underline{v}_c - \underline{v}_c \cdot \left[ \frac{e^{\underline{v}_w^T \underline{v}_c}}{\sum_{w'} e^{\underline{v}_{w'}^T \underline{v}_c}} \right] \right]$$

$$= \left[ \underline{v}_c - \underline{v}_c \cdot P(w | c) \right]$$

$$= \left[ \underline{v}_c [1 - P(w | c)] \right]$$

So update  $\underline{v}_w$  as  $\underline{v}_w^{new} = \underline{v}_w^{old} + \eta \underline{v}_c [1 - P(w | c)]$

$$y = f(x)$$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$\frac{\partial y}{\partial x} = \left[ \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right]$$

The bottleneck of this iterative method is that during every update we have to take the sum  $\sum_{w'} e^{\underline{v}_{w'}^T \underline{v}_c}$

involving all the words in the vocabulary.

The number of words in the vocabulary may run up to million. Hence this method does not work practically. So use a different approach

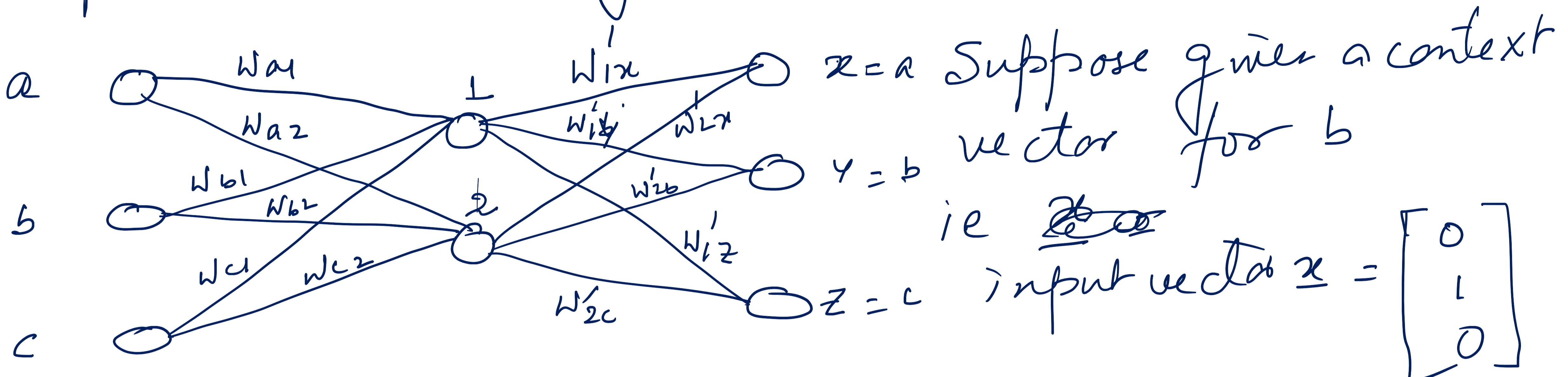
## Negative Sampling

## A working Example

Assume there are only 3 words in the vocabulary.

6 Try to work on an example of CBOW model.

We are given a context word and we need to predict the missing word.



Let  $w$  be the weights to the hidden layer units

$$w = \begin{bmatrix} w_{a1} & w_{a2} \\ w_{b1} & w_{b2} \\ w_{c1} & w_{c2} \end{bmatrix} \quad \text{Then } \underline{h} = w^T x$$

$$= \begin{bmatrix} w_{a1} & w_{b1} & w_{c1} \\ w_{a2} & w_{b2} & w_{c2} \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} w_{b1} \\ w_{b2} \end{bmatrix} = \underline{y}_c = \underline{h}_b$$

$$w' = \begin{bmatrix} w'_{1x} & w'_{1y} & w'_{1z} \\ w'_{2x} & w'_{2y} & w'_{2z} \end{bmatrix} \quad \text{Vector representation for the word } z = \underline{y}_w$$

So the output vector is  $\underline{u} = w'^T \underline{h}$

$$= \begin{bmatrix} w'_{1x} & w'_{2x} \\ w'_{1y} & w'_{2y} \\ w'_{1z} & w'_{2z} \end{bmatrix} \begin{bmatrix} w_{b1} \\ w_{b2} \end{bmatrix} = \begin{bmatrix} w'_{1x} \cdot w_{b1} + w'_{2x} \cdot w_{b2} \\ w'_{1y} \cdot w_{b1} + w'_{2y} \cdot w_{b2} \\ w'_{1z} \cdot w_{b1} + w'_{2z} \cdot w_{b2} \end{bmatrix}$$

$$= \begin{bmatrix} v_w^{(a)} \cdot v_c^{(b)} \\ v_w^{(b)} \cdot v_c^{(b)} \\ v_w^{(c)} \cdot v_c^{(b)} \end{bmatrix}$$

Convert to softmax score

$$y = \text{softmax}(\underline{u}) \quad \underline{v}_c^T \underline{v}_w$$

$$y_i = \text{softmax}(u_i) = \frac{e^{u_i}}{\sum_j e^{u_j}}$$

$$y_a = P(w=a|c=b); \quad y_b = P(w=b|c=b); \quad y_c = P(w=c|c=b)$$

Loss function to derive the Parameters  $\theta = \{v_c, v_w\}$

Maximum likelihood

Consider the dataset

$$D = (c_1, w_1) (c_2, w_2) (c_1, w_3) \dots (c_2, w_3), (c_2, w_5) \dots$$

↓  
Context  
 $\{a, b, c\}$       word  
 $\{a, b, c\}$

- Attempt to learn the word distribution for each context word.

Suppose consider the context word  $c_1 \rightarrow$  find the prob of appearance of each word adjacent to it.

Suppose in the actual data the following context & word pair appears

$$(a, b) (a, c), (b, a) (c, a) (a, b) (a, c) (b, a) (a, c)$$

~~Given Context word a ; b & c appears as word.~~

~~Given Context word b ; the prob of a & c appearing as words.~~

$$= P(a|b) \cdot P(a|b)$$

$$L = \cancel{\text{Ans}} \frac{e^{v_c^{(b)^T} v_w^{(a)}}}{\sum_{w'} e^{v_c^{(b)^T} v_w^{(w')}}} \times \frac{e^{v_c^{(b)^T} v_w^{(c)}}}{\sum_{w'} e^{v_c^{(b)^T} v_w^{(w')}}}$$

Now Log L and try to maximize Log L with respect to  ~~$v_c$  &  $v_w$~~

So to avoid calculating the denominator repeatedly, we use negative sampling.

Negative Sampling  $\leftarrow$  Distributed representation of words & phrases and their compositionality (mikolov et al)

for skipGram

Convert this problem to a supervised problem.

Consider 2 datasets:  $D$ : set of all pairs  $(c, w)$  that are in text

$D'$ : set of all pairs  $(c, w)$  that are NOT in text

Let  ~~$Z = 1$  if~~ label  $Z = 1$  for all pairs in  $D$   
 $Z = 0$  for all pairs in  $D'$

Turned the problem to a classification prob.

$P(Z=1 | (c, w))$  is the prob that  $(c, w)$  appears in  $D$

$P(Z=0 | (c, w))$  " " " " "  $(c, w)$  does not appear in  $D$ .

$$P(Z|X) = \frac{P(X|Z)P(Z)}{P(X)}$$

Posterior                          class conditional  
    Prior  
    marginal

So to directly find the posterior without calculating the class conditionals or priors, we can ~~use~~ assume a function that ~~returns~~ returns the prob of  $Z$  given  $X$ .

They assume  $P(Z|X)$  follows a logistic Regression function

$$\text{i.e } P(Z=1 | (c, w)) = \frac{1}{1 + e^{\underline{v}_c^T \underline{v}_w}}$$

$$P(Z=0 | (c, w)) = 1 - \frac{1}{1 + e^{\underline{v}_c^T \underline{v}_w}} = \frac{e^{\underline{v}_c^T \underline{v}_w}}{1 + e^{\underline{v}_c^T \underline{v}_w}} = \frac{1}{1 + e^{-\underline{v}_c^T \underline{v}_w}}$$

$$\text{The } P(Z | (c, w)) = \left( \frac{1}{1 + e^{\underline{v}_c^T \underline{v}_w}} \right)^Z \left( \frac{1}{1 + e^{-\underline{v}_c^T \underline{v}_w}} \right)^{1-Z}$$

so  $P(Z | (c, w))$  is also parameterized by  $\theta = \{\underline{v}_c \underline{v}_w\}$

To derive these parameters we use maximum likelihood.

$$\text{let } D = \{(c_1, w_1), (c_2, w_2), (c_3, w_3), \dots\} \leftarrow \text{Positive pairs}$$

$$D' = \{(c_4, w_4), (c_5, w_5), \dots\} \leftarrow \text{Negative pairs}$$

$$L(\theta) = \prod_{(c, w) \in D \cup D'} P(Z | (c, w))$$

$$= \text{Log L}(\theta) = l = \sum_{(c, w) \in D \cup D'} \log P(z | (c, w)) = \sum_{(c, w) \in D \cup D'} \log \left[ \frac{z}{1 + e^{v_c^T v_w}} \right] \left( \frac{1-z}{1 + e^{-v_c^T v_w}} \right)$$

Minimizing  $-\text{Log}(L(\theta)) = -l$

$$= - \sum_{(c, w) \in D \cup D'} z \log \left( \frac{1}{1 + e^{v_c^T v_w}} \right) + (1-z) \log \left( \frac{1}{1 + e^{-v_c^T v_w}} \right)$$

$$= - \left[ \sum_{(c, w) \in D} \log \left( \frac{1}{1 + e^{v_c^T v_w}} \right) + \sum_{(c, w) \in D'} \log \left( \frac{1}{1 + e^{-v_c^T v_w}} \right) \right]$$

The skipgram paper considers one positive pair and a small no. of negative pairs (around 5) for each gradient update

How are (-)ve pairs identified.

- They consider derive  $P_n(w) = \frac{\#w}{\sum \#w'}$   
For each word they consider randomly sample 2 words based on this probability.
- And they assume that these 2 words will not appear as (word, context) pair in the corpus

Then they derive the gradient

$$\frac{\partial l}{\partial v_c}, \frac{\partial l}{\partial v_w}$$

$\begin{matrix} 1 \\ 0.2 \\ 1.2 \\ 0.25 \end{matrix}$	$\begin{matrix} 0.25 \\ -0.50 \end{matrix}$	$\dots$
the	by	
a	$(a, \text{the}) \rightarrow \in D'$	

Choose arandomly between 0 & 1

Suppose 1st  $\rightarrow 0.3$   
2nd  $\rightarrow 0.1$  (the)

$$P_1(\text{the}) = \frac{20}{100} = 0.2$$

$$P_1(\text{boy}) = \frac{5}{100} = 0.05$$

$$P_1(a) = \frac{25}{100} = 0.25$$

$$v_c^{\text{new}} = v_c^{\text{old}} - \eta \frac{\partial l}{\partial v_c}$$

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

→ Set of running text documents.

→ From the set of running text document & create (word & context) pairs  
Let that be marked as  $D'$

→ Minimize the (-)ve loglikeit

- Create the negative sample set  $D'$
- Create a batch of one positive and 5 negative samples
- Minimize the (-)ve log likelihood of generating the batch for using the model.

Consider the text I am a very good boy.

$$D = \begin{matrix} (\text{I}, \text{am}) \\ \text{word} \end{matrix} \quad \begin{matrix} (\text{am } \text{a}) \\ \text{word cont} \end{matrix} \quad \begin{matrix} (\text{am }, \text{I} \\ \text{word}, \text{cont}) \end{matrix} \quad \begin{matrix} (\text{a}, \text{very}) \\ \text{word} \end{matrix} \quad \begin{matrix} (\text{a}, \text{am}), \\ \text{cont} \end{matrix} \\ (\text{very}, \text{good}), (\text{very}, \text{a}), (\text{good boy}), (\text{good very}) \\ (\text{boy}, \text{good})$$

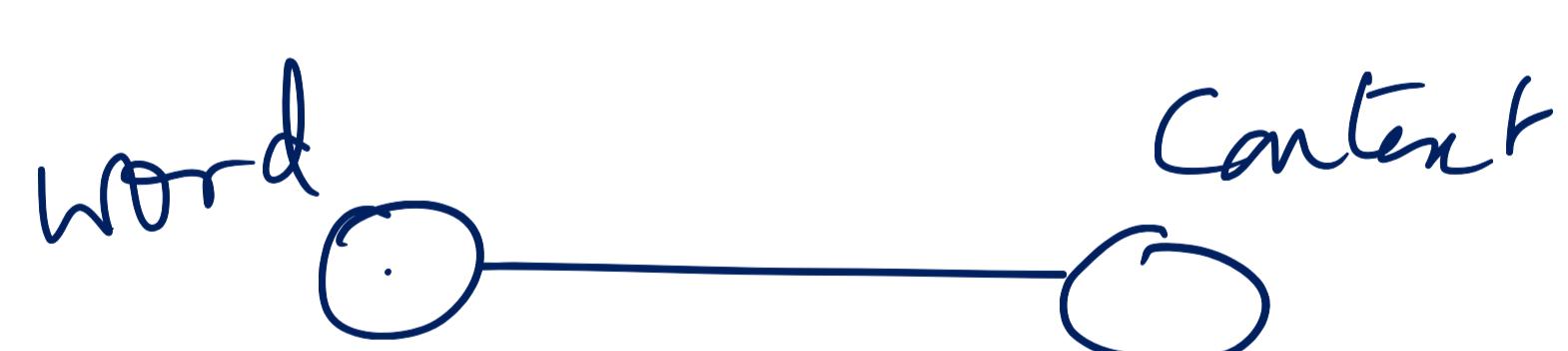
For each start word  $w$  initialize with a random embedding  $v_w$

For each content  $c$   $v_c$   $\in \mathbb{R}^{n \times n}$  random  $n \in V_c$

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

$$v_c^{\text{new}} = v_c^{\text{old}} - \eta \frac{\partial l}{\partial v_c}$$

What should be ~~the same~~  $D$  for a graph



### Assignment

for a given  $u$ , the random walk should be biased towards nodes with degree similar to  $u$ .





