

Neural Models for NLP

Amit Awekar

Language: System, Symbols, Convey meaning

Natural: Human-Human, Machine-Human, Ambiguity

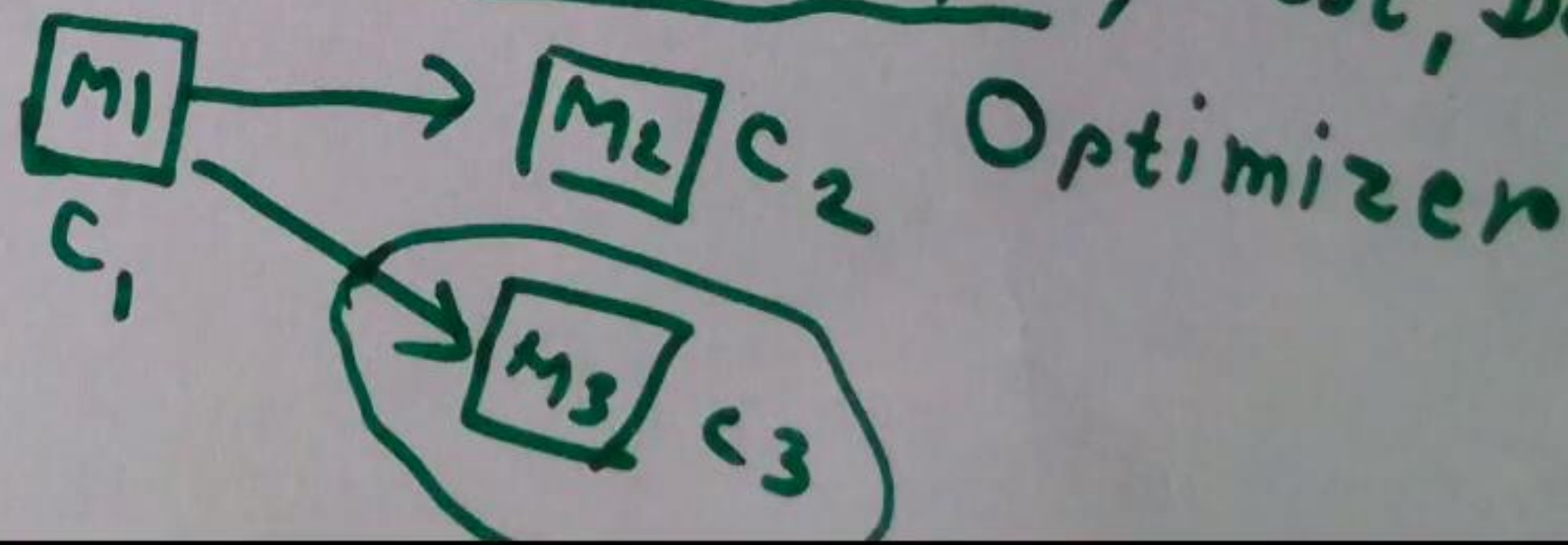
Processing: Goal driven, Task oriented

Models: Machine Learning: Input \rightarrow Output, Cost, Data

Neural:

Brain

Parameters



Input

Output

Model

$$M1: y = x$$

$$M2: y = 2x$$

$$M3: y = x + 1$$

1

2

$$y = \underline{a} \cdot x + \underline{b}$$

2

6

Parameters

3

10

$$a = 1$$

40

38

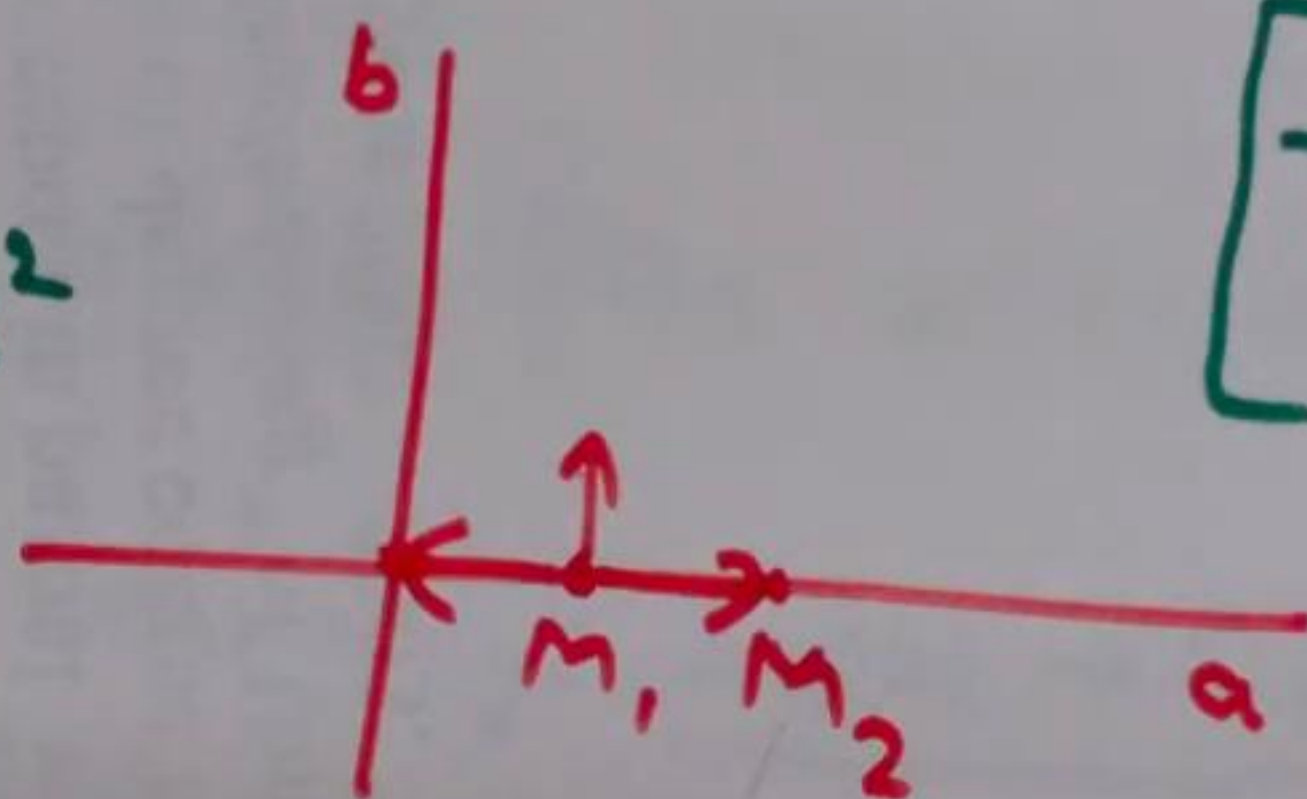
$$b = 1$$

Optimizer

Cost Function

$$C(M_1) = 1^2 + 4^2 + 7^2 + 28^2$$

$$\sum_{i=1}^d (y_i - y_{im})^2$$



$$\left[\frac{-\partial C}{\partial a}, \frac{-\partial C}{\partial b} \right]$$

Step Size
1 unit

$$C(M_2) = 0 + 2^2 + 4^2 + 18^2$$

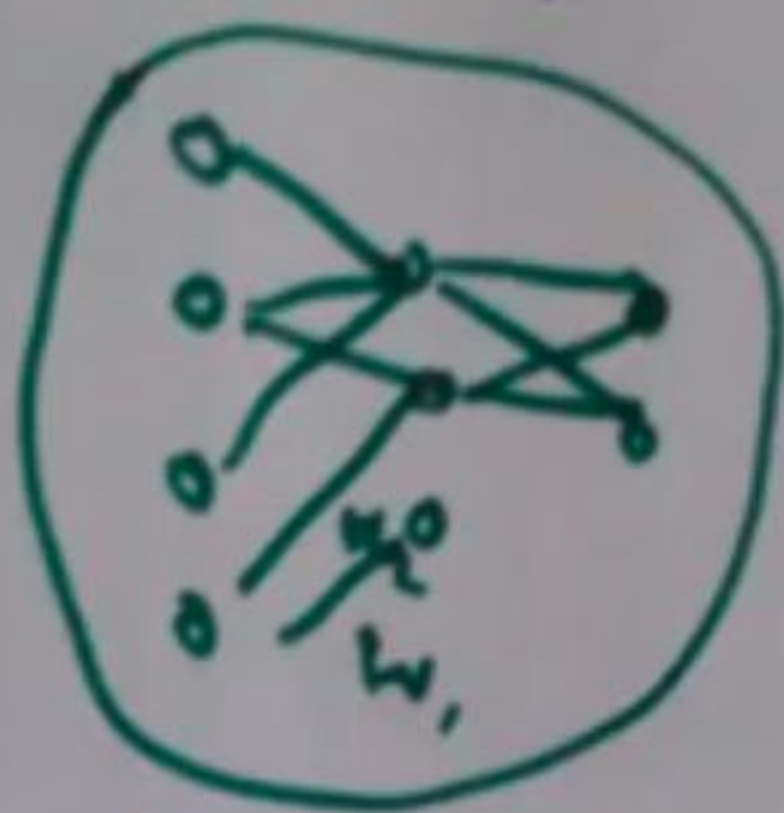
$$C(M_3) = 0 + 3^2 + 6^2 + 27^2$$

$$a = a - \left[\frac{\partial C}{\partial a} \cdot \text{Step Size} \right]$$
$$b = b - \left(\frac{\partial C}{\partial b} \cdot \text{Step Size} \right)$$

Data

Tweet 1 +ve 1
Tweet 2 -ve -1
⋮
Tweet 10^6 +ve

Model



Parameters

$w_1 = 0.1$
 w_2
⋮

w_{1000}

Model Sequence

M_1
 M_2
⋮
 M_{20}

Optimizer

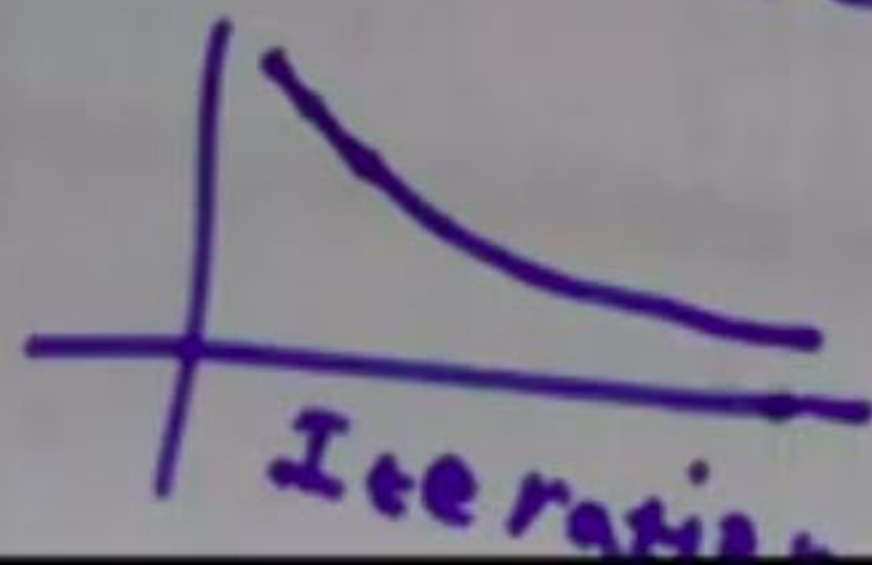
$$w_i = w_i - \frac{\partial C}{\partial w_i} \cdot \text{Step Size}$$

Back propagation

Cost Function
 $C(M_i)$

$$\sum_{i=1}^d (y_i - y_{i,m})^2$$

Step Size



Motivation for ANN

Neurons : 10^{11} : Total

10^4 : Connections per neuron

Excited / Inhibited : States of neuron

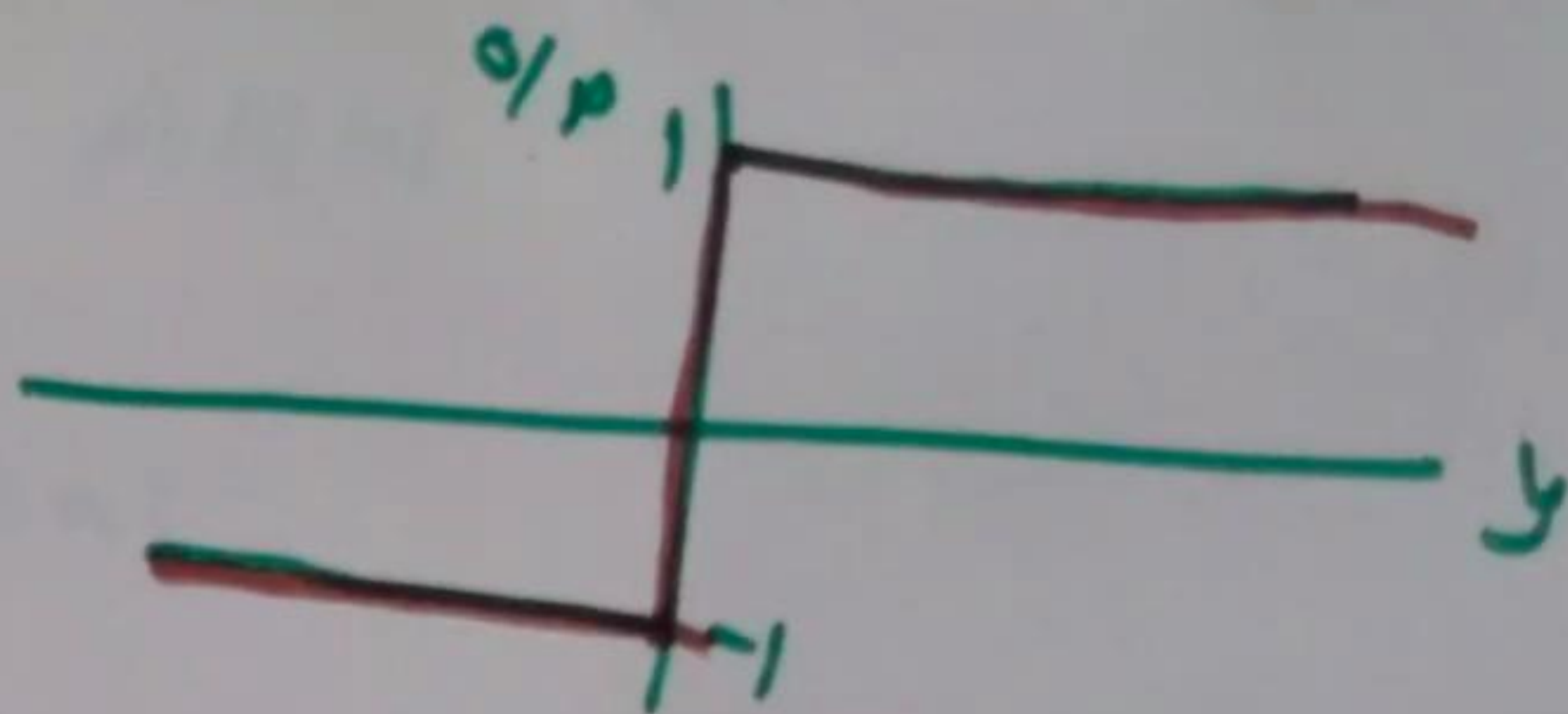
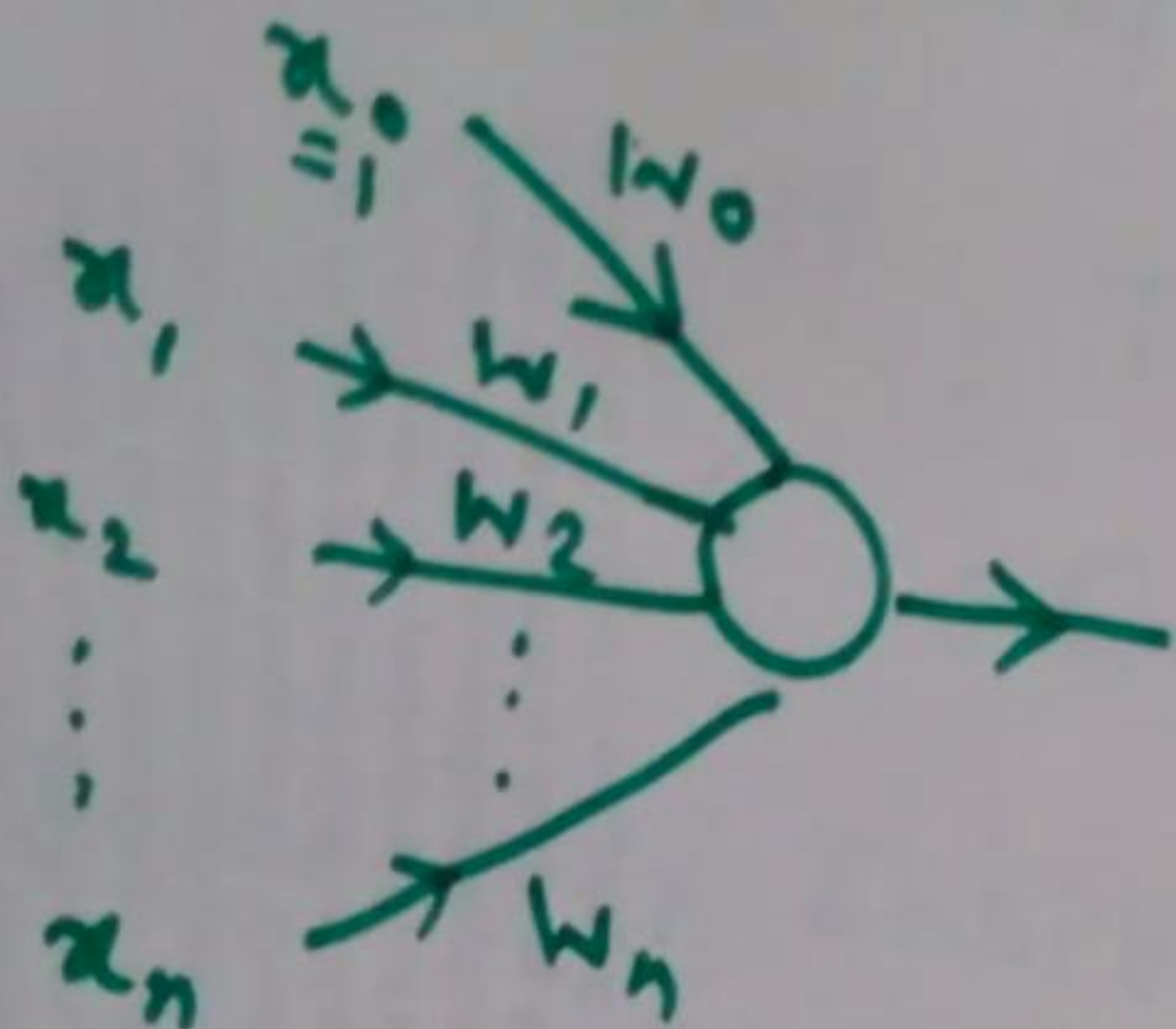
10^{-3} sec : Speed of neuron

10^{-1} sec : Time required for many human decisions

Parallel

Distributed across multiple neurons

Perceptron



Output = 1 if $y > 0$
 -1 else

$$w_0 + x_1 \cdot w_1 + x_2 \cdot w_2 + \dots + x_n \cdot w_n$$

$$y = \sum_{i=0}^n x_i \cdot w_i$$

$$w_i = w_i + \frac{\Delta w_i}{n(t - 0) x_i}$$



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

$$\frac{1}{2} \sum_{d \in D} (t_d - o_d)^2 \leftarrow E(\vec{w})$$

$$\nabla E(\vec{w}) = \left[\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$$

$$\vec{w} = \vec{w} + \frac{\Delta \vec{w}}{\|\nabla E(\vec{w})\|}$$

$$w_i = w_i + \frac{\Delta w_i}{\|\nabla E(\vec{w})\|}$$

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

$$-\frac{1}{2} \sum_d x_d (t_d - o_d) \frac{\partial}{\partial w_i} (t_d - o_d)$$

$\uparrow x_d \quad \uparrow t_d$

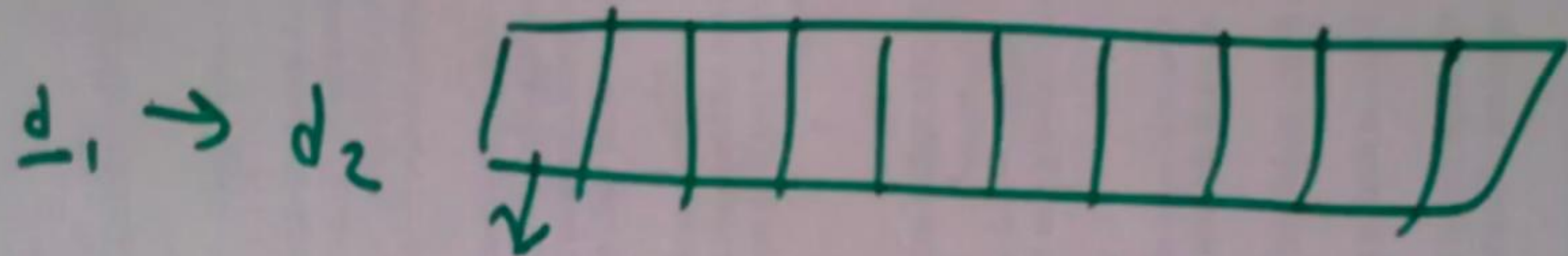
$$\frac{p_i x_i}{w_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n}$$

$$p_i x_i (t_d - o_d) \cdot \sum_d$$

$$w_i = w_i + \boxed{\Delta w_i} \rightarrow -n \cdot \frac{\partial E}{\partial w_i}$$

$$n \cdot \sum_d (t_d - o_d) x_d$$

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

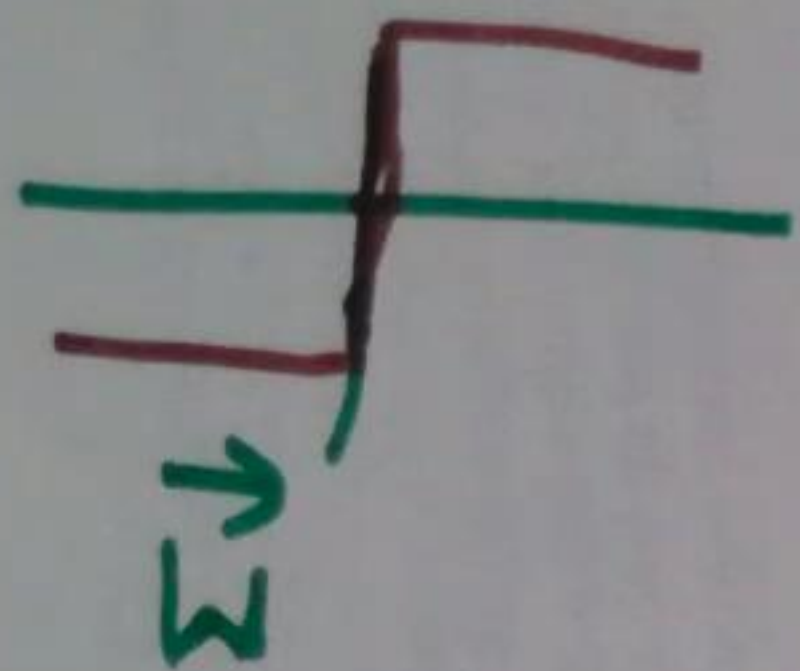


$$E_d(\vec{w}) = \frac{1}{2} (t_d - o_d)^2$$

$$\Delta w_i = n \cdot (t_d - o_d) \cdot x_{i,d}$$

$$E_{d_2}(\vec{w}) = \frac{1}{2} (t_{d_2} - o_{d_2})^2$$

Thresholded



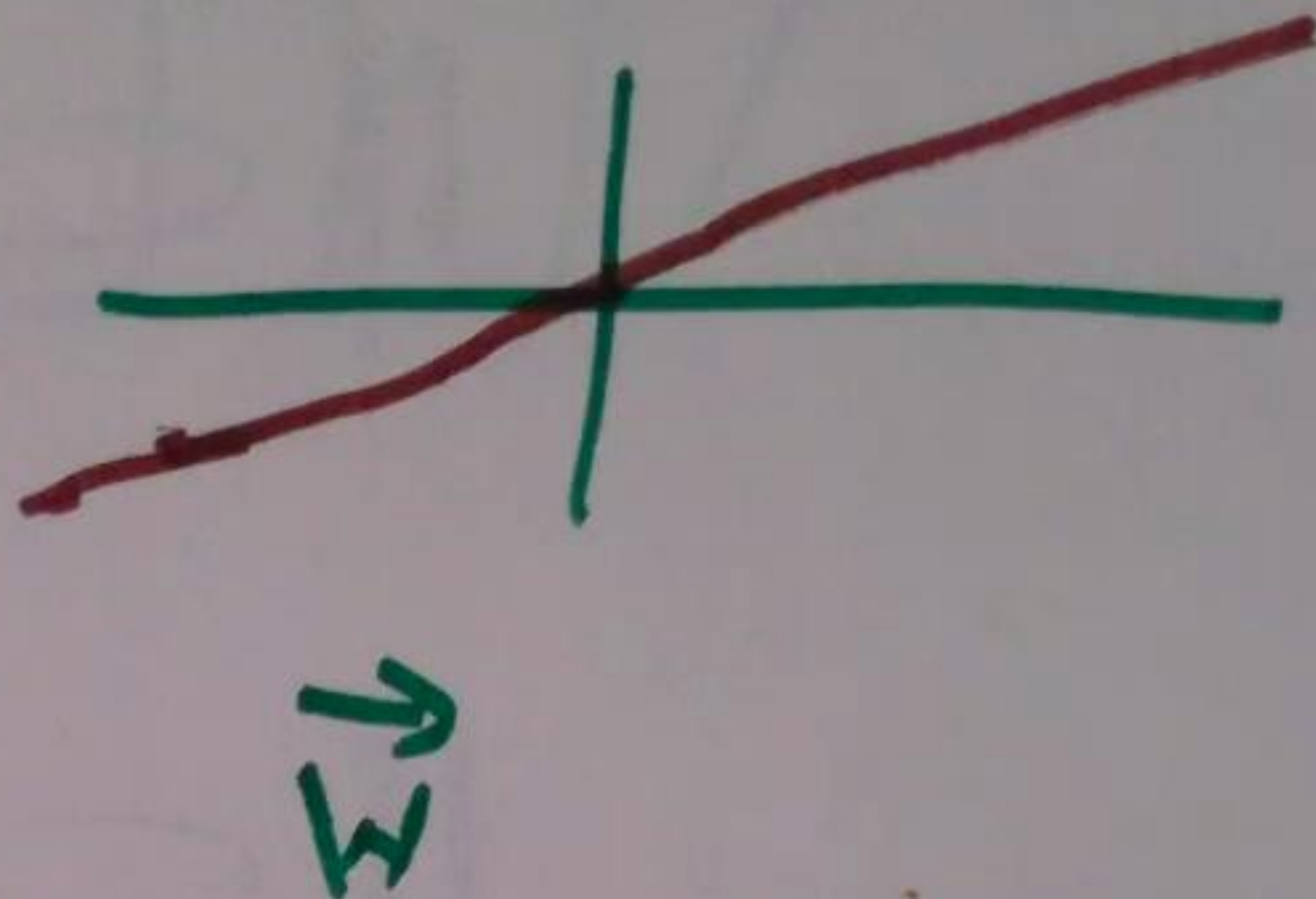
Perceptron rule

$$w_i = w_i + \Delta w_i$$

$$\Delta w_i = \eta \cdot (t_i - o_i) x_i$$

Linearly separable

Universal approximator



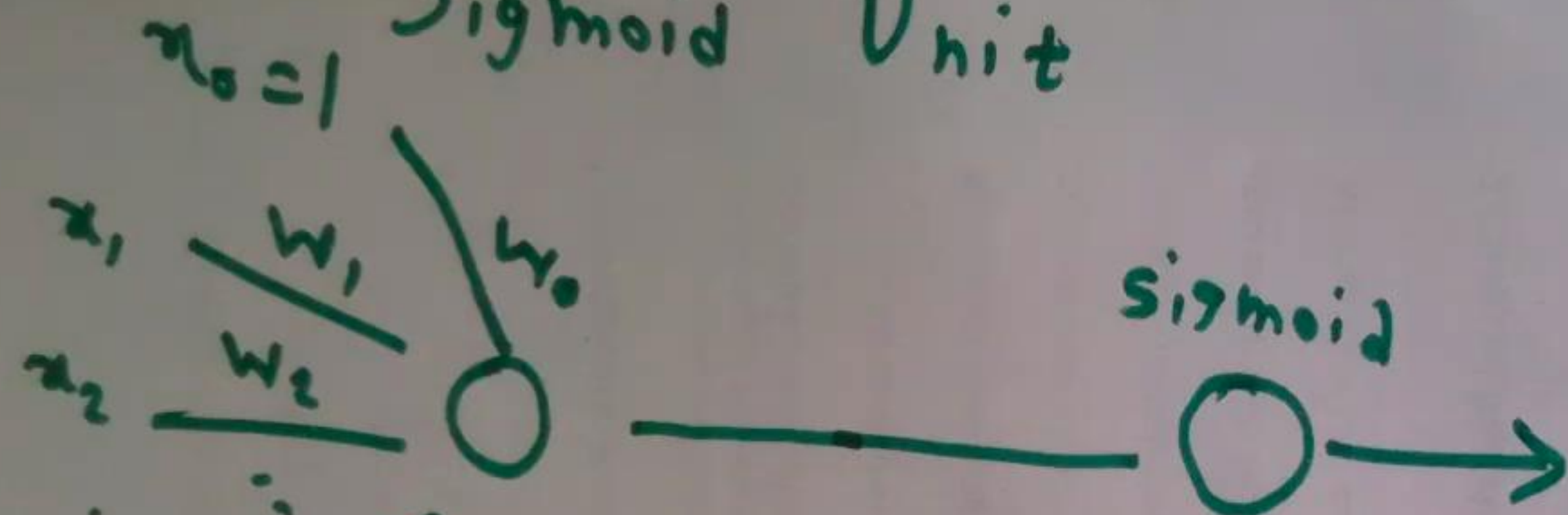
$$\Delta w_i = \eta \cdot \sum_{d \in D} (t_d - o_d) x_i$$

Stochastic Gradient Descent

$$E(\vec{w}) \rightarrow \begin{matrix} E_1(\vec{w}) \\ E_2(\vec{w}) \end{matrix}$$

$$\Delta w_i = \eta \cdot (t - o) x_i$$

Sigmoid Unit



$f = \text{activation function}$

$$\vec{x} \cdot \vec{w} \Rightarrow \text{net} = \sum_{i=0}^n x_i \cdot w_i$$

$$\text{output} = f(\text{net})$$

$$\text{sigmoid}(\text{net}) = \frac{1}{1 + e^{-\text{net}}}$$



$$\text{sigmoid}(x) = \frac{1}{1+e^{-x}} = y$$

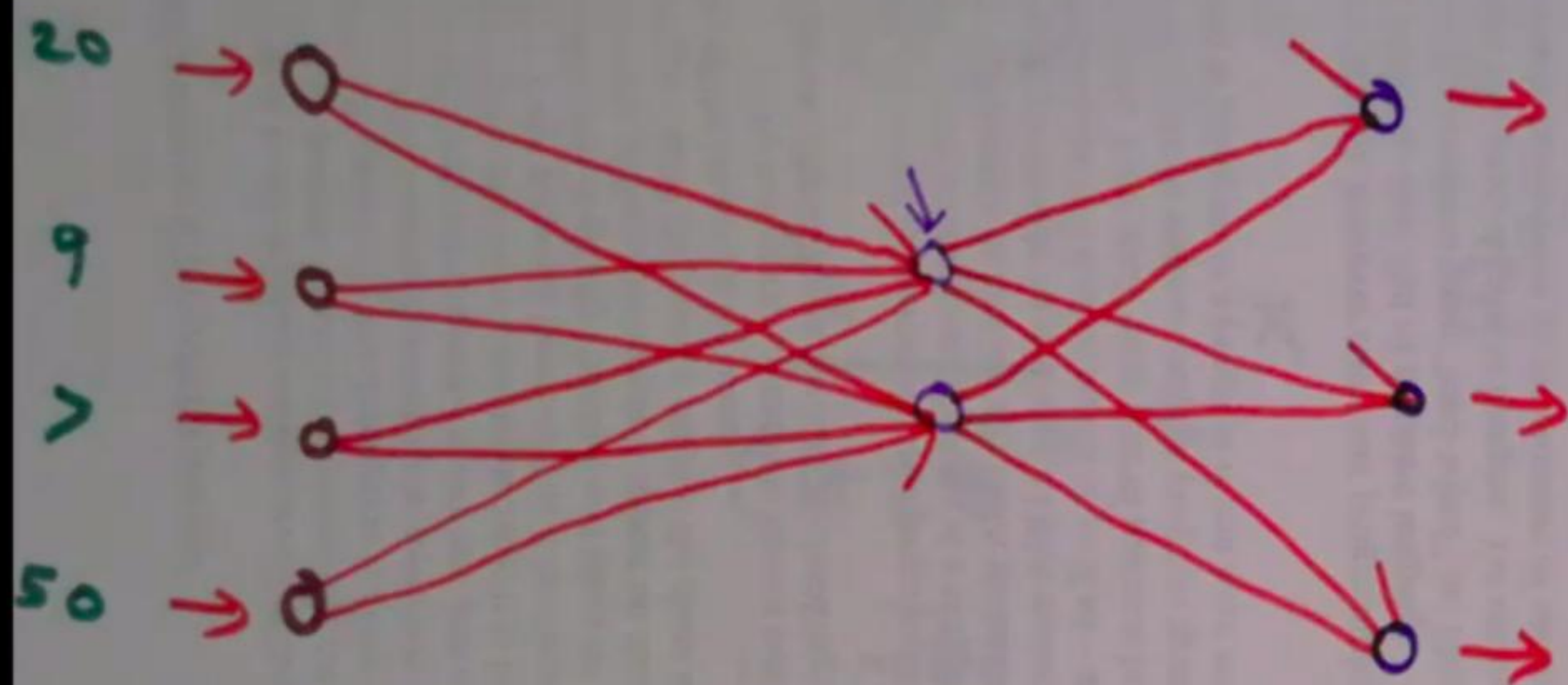
$$\frac{d}{dx} \left(\frac{1}{1+e^{-x}} \right) = \frac{-1}{(1+e^{-x})^2} \frac{d}{dx} (1+e^{-x})$$

$$= \frac{-1}{(1+e^{-x})^2} \times e^{-x}$$

$$= \frac{1 \cdot e^{-x} + 1 \cdot (-1)}{(1+e^{-x})^2}$$

$$= \frac{e^{-x} + 1}{(1+e^{-x})^2} - \frac{1}{(1+e^{-x})^2}$$

$$y - y^2$$



$$W_{ji} = W_{ji} + \underbrace{\Delta W_{ji}}_{n \cdot \delta_j \cdot x_{ji}}$$

Output layer: k

$$\delta_k = o_k (1 - o_k) (t_k - o_k)$$

Hidden layer: h

$$\delta_h = o_h (1 - o_h) \sum_{k \in \text{Output}} \delta_k \cdot W_{hk}$$

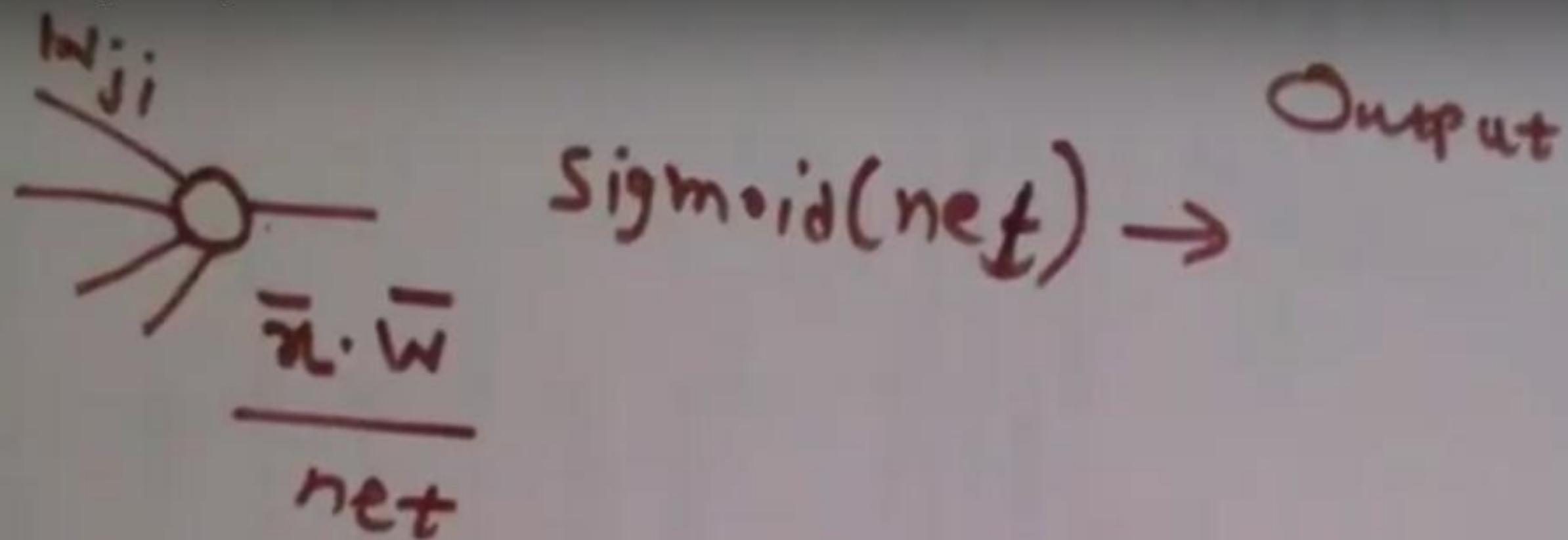
W_{ji}

x_{ji}

x, t

o_k

t_k



$$w_{ji} = w_{ji} + \Delta w_{ji}$$

$$-n \cdot \frac{\partial E}{\partial w_{ji}} \cdot \frac{\partial E}{\partial net_j} \cdot \frac{\partial net_j}{\partial w_{ji}} \rightarrow \frac{\partial w_{ji} \cdot x_{ji}}{\partial w_{ji}} \rightarrow x_{ji}$$

$$E_d(\bar{w}) = \frac{1}{2} \sum_{k \in output} (t_k - o_k)^2$$

$$\delta_j = - \frac{\partial E}{\partial net_j}$$

$$o_k$$

$$\frac{\partial E}{\partial net_k}$$

$$= \frac{1}{net_k} \left[\frac{\partial E}{\partial net_k} \right]$$

$$o_k$$

$$(t_k - o_k)^2$$

$$-\frac{\partial o_k}{\partial net_k} \cdot \boxed{\frac{\partial E}{\partial o_k}}$$

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

$$- o_k (1 - o_k) \cdot \frac{1}{2} \frac{\partial}{\partial o_k} (t_k - o_k)^2$$

$$\frac{1}{2} \cdot 2 (t_k - o_k) \cdot \frac{\partial}{\partial o_k} (t_k - o_k)$$

$$o_k (1 - o_k) \cdot (t_k - o_k)$$

$$\frac{\partial net_i}{\partial net_j}$$

$$\begin{aligned}
 & \frac{\partial E}{\partial \text{net}_k} = - \frac{\partial O_k}{\partial \text{net}_k} \cdot \boxed{\frac{\partial E}{\partial O_k}} \\
 & \frac{1}{2} \sum_{k \in \text{output}} (t_k - O_k)^2 \\
 & - O_k(1 - O_k) \cdot \frac{1}{2} \frac{\partial}{\partial O_k} (t_k - O_k)^2 \\
 & \frac{1}{2} \cdot 2(t_k - O_k) \cdot \frac{\partial}{\partial O_k} (t_k - O_k) \\
 & O_k(1 - O_k) \cdot (t_k - O_k)
 \end{aligned}$$

$$W_{ji} = W_{ji} + \Delta W_{ji}$$

$$-n \cdot \frac{\partial E}{\partial W_{ji}}$$

$$\frac{x_{ji} \cdot n \cdot o_j \cdot (1 - o_j) \sum \delta_k \cdot w_{kj}}{\delta_j}$$

$$-n \cdot \frac{\partial E}{\partial net_j} \cdot \boxed{\frac{\partial net_j}{\partial W_{ji}}} \cdot x_{ji}$$

$$\sum_{k \in \text{Downstream}(j)} \left(\frac{\partial E}{\partial net_k} \cdot \frac{\partial net_k}{\partial net_j} \right)$$

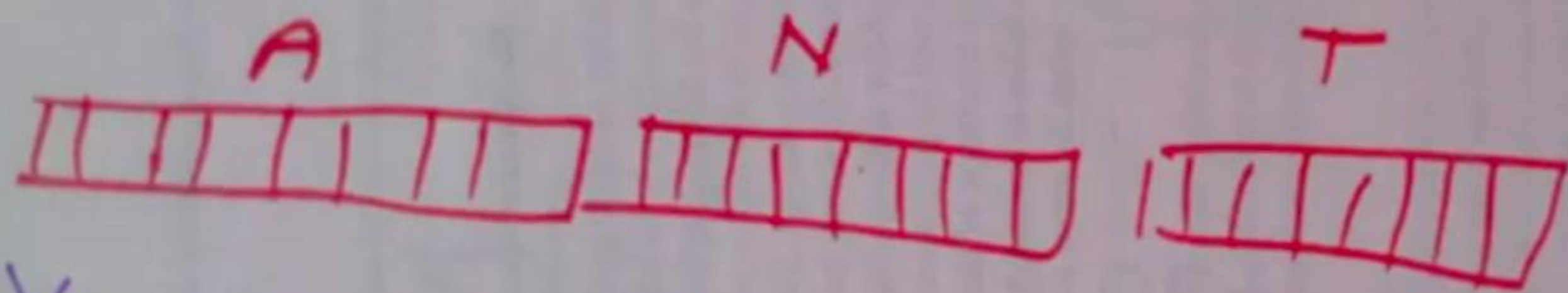
$$x_{ji} \cdot n \cdot \sum \delta_k \cdot \left[\frac{\partial o_j}{\partial net_j} \cdot \frac{\partial net_k}{\partial net_j} \right]$$

Data Representation: Vector

Characters, Words, Phrases, Sentences, Paragraph,
→ Document, Section, Book / Collection, Corpus
Audio

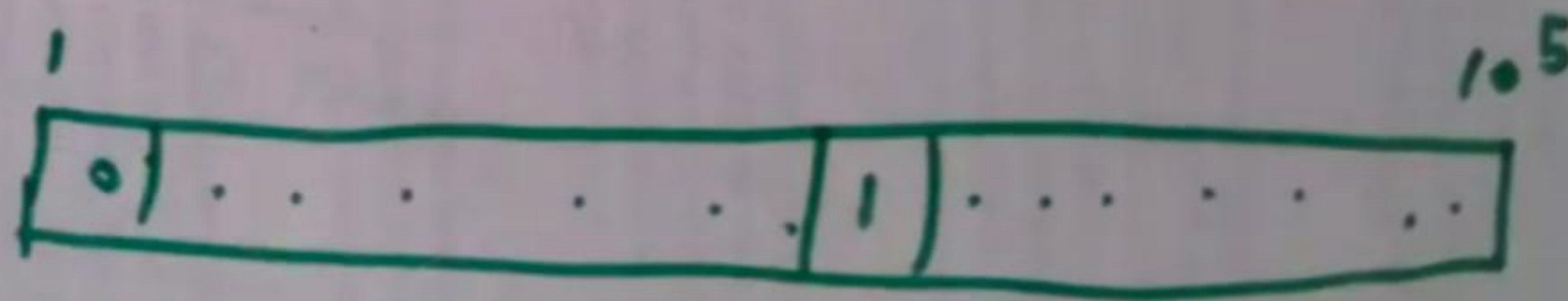
Characters: ASCII
8 bits

UNICODE
16 bits



Vocab size: 10^5

One hot encoding:



CBOW

Given: Context

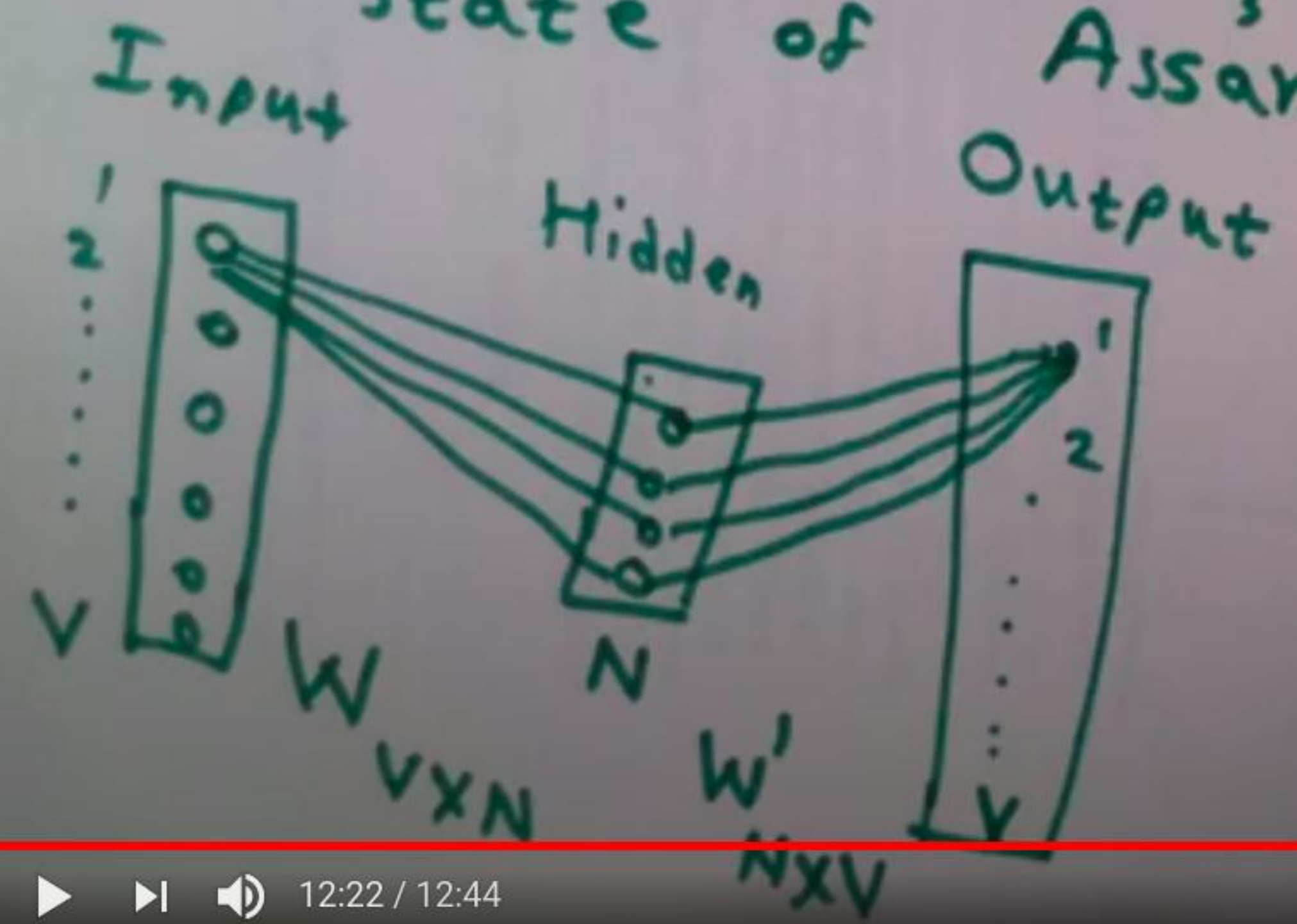
Predict: Word

Skip Gram

Given: Word

Predict: Context

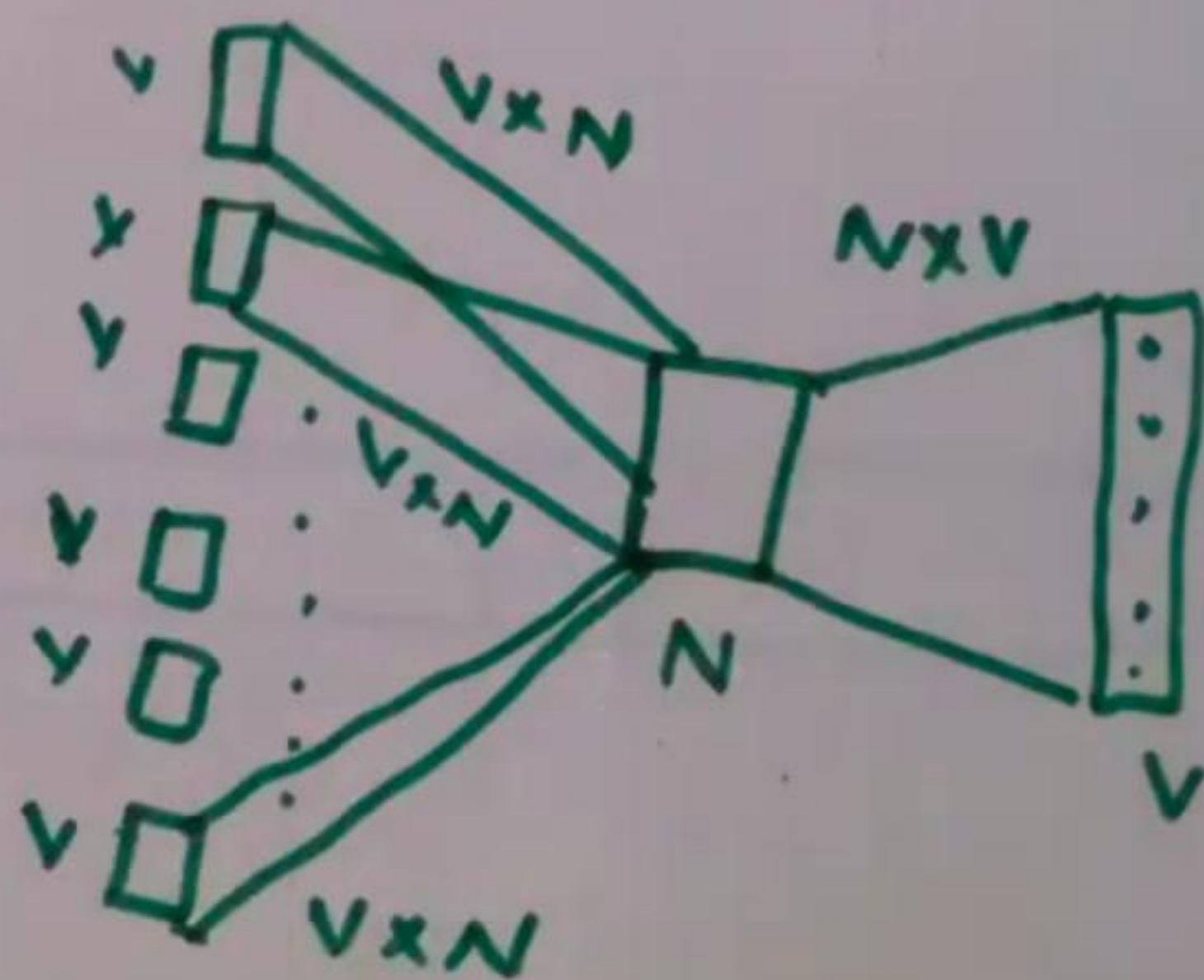
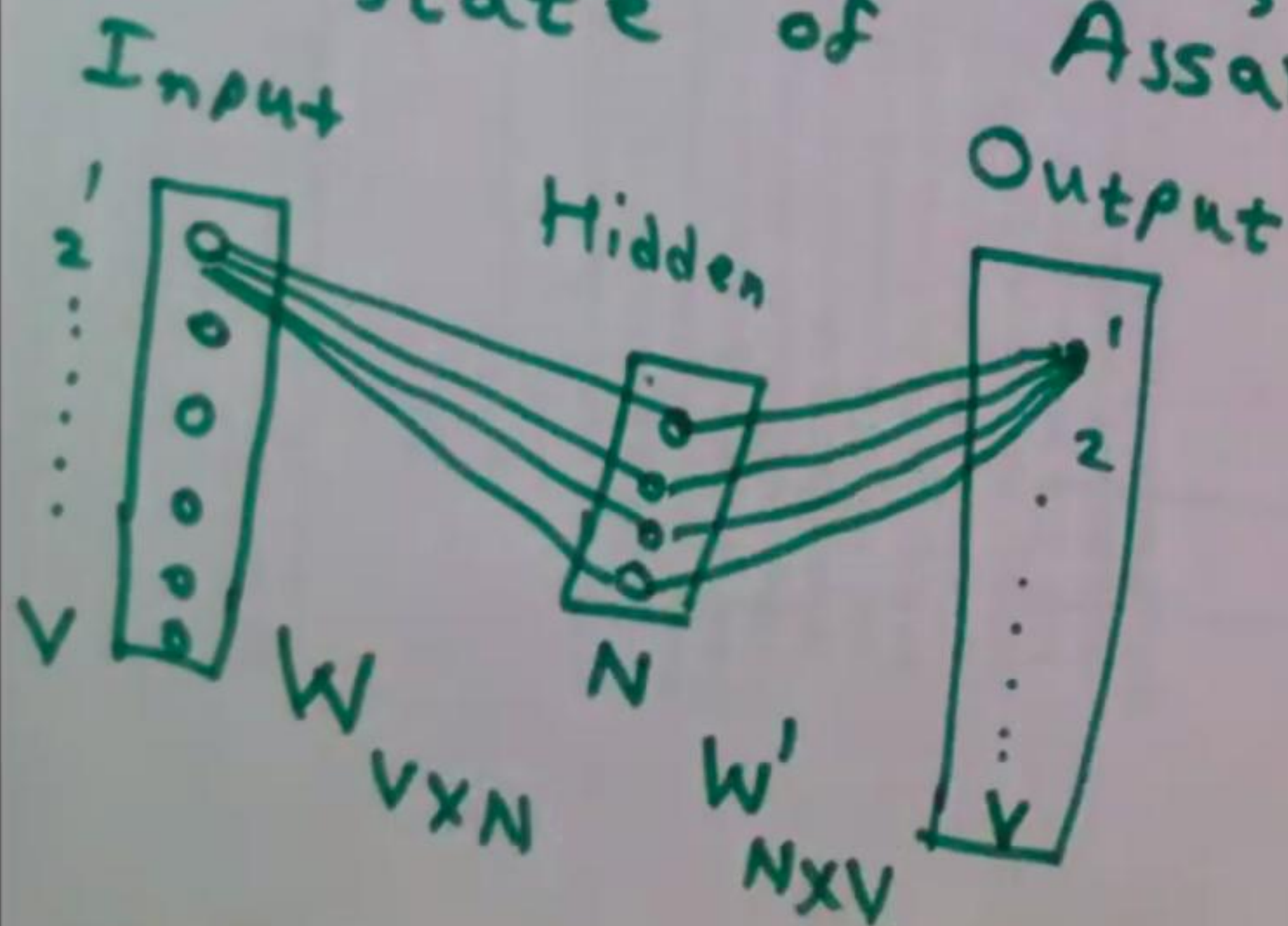
$\frac{I}{1}$ $\frac{I}{2}$ $\frac{T}{3}$ Guwahati $\frac{is}{3}$ located $\frac{in}{4}$ $\frac{north}{5}$ $\frac{eastern}{6}$
State of Assam.

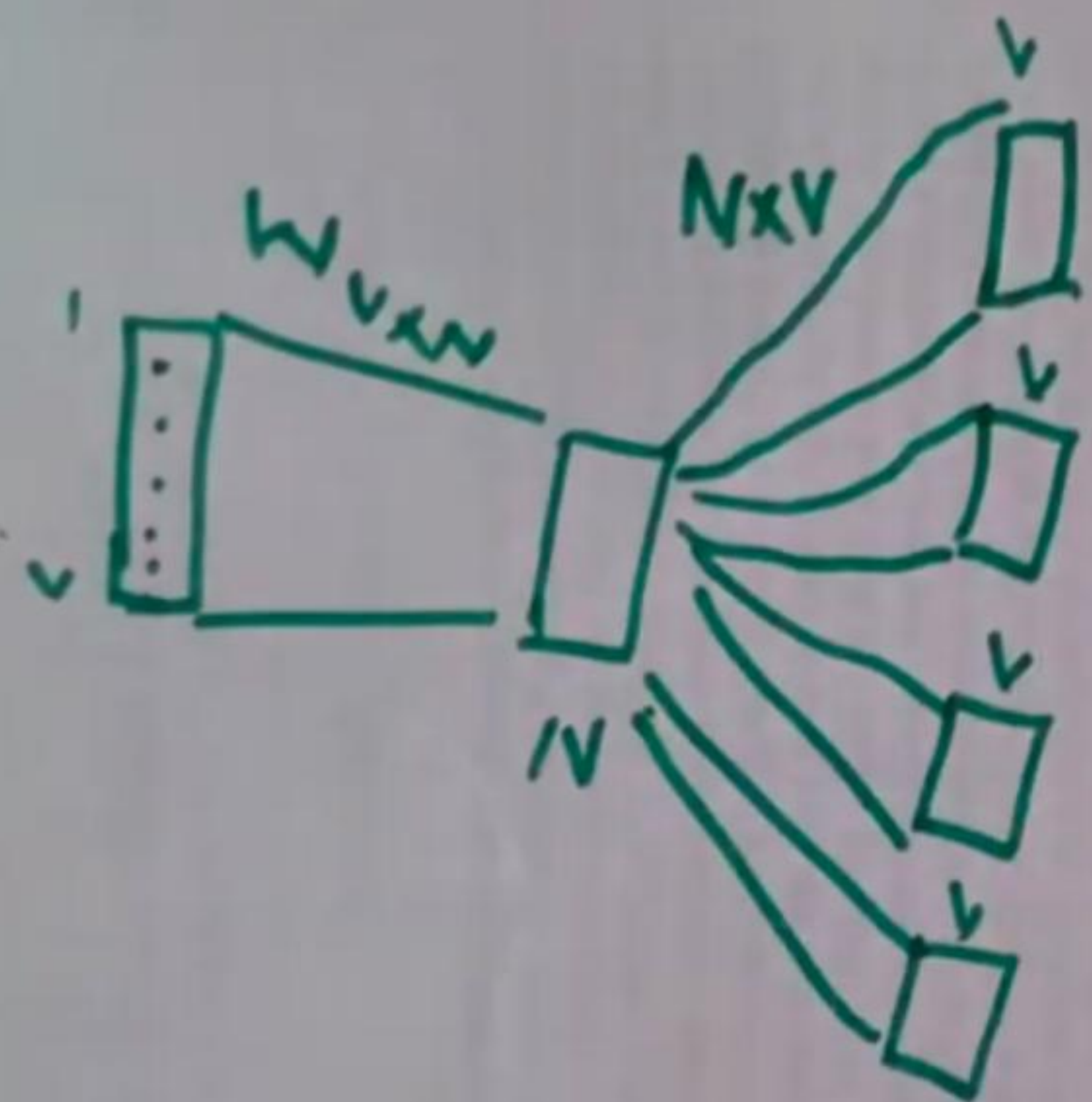


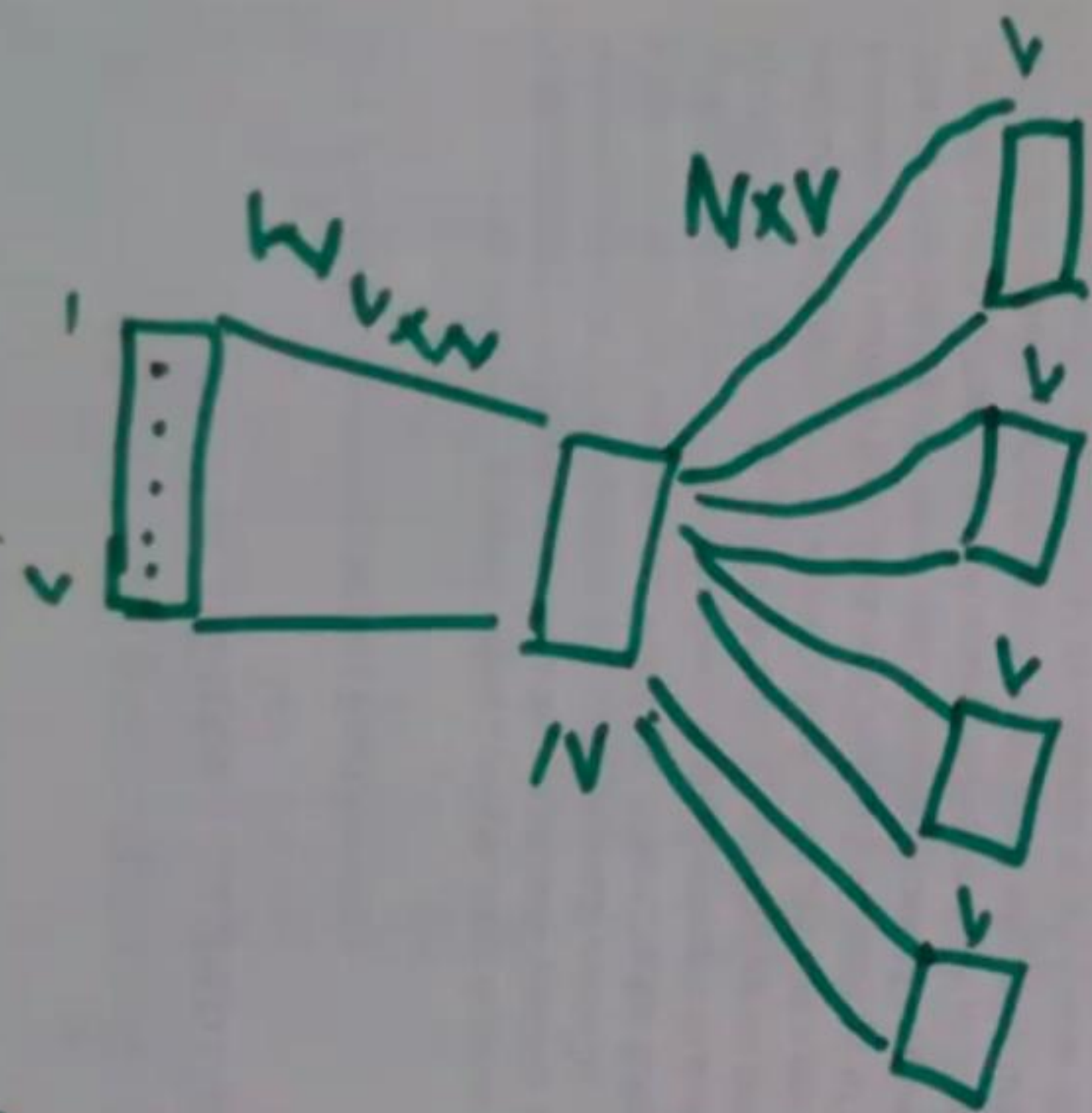
Context
Predict: Word

Given: Word
Predict: Context

1 2 3 4 5 6
1 2 3 4 5 6
IIT Guwahati is located in north eastern
State of Assam.



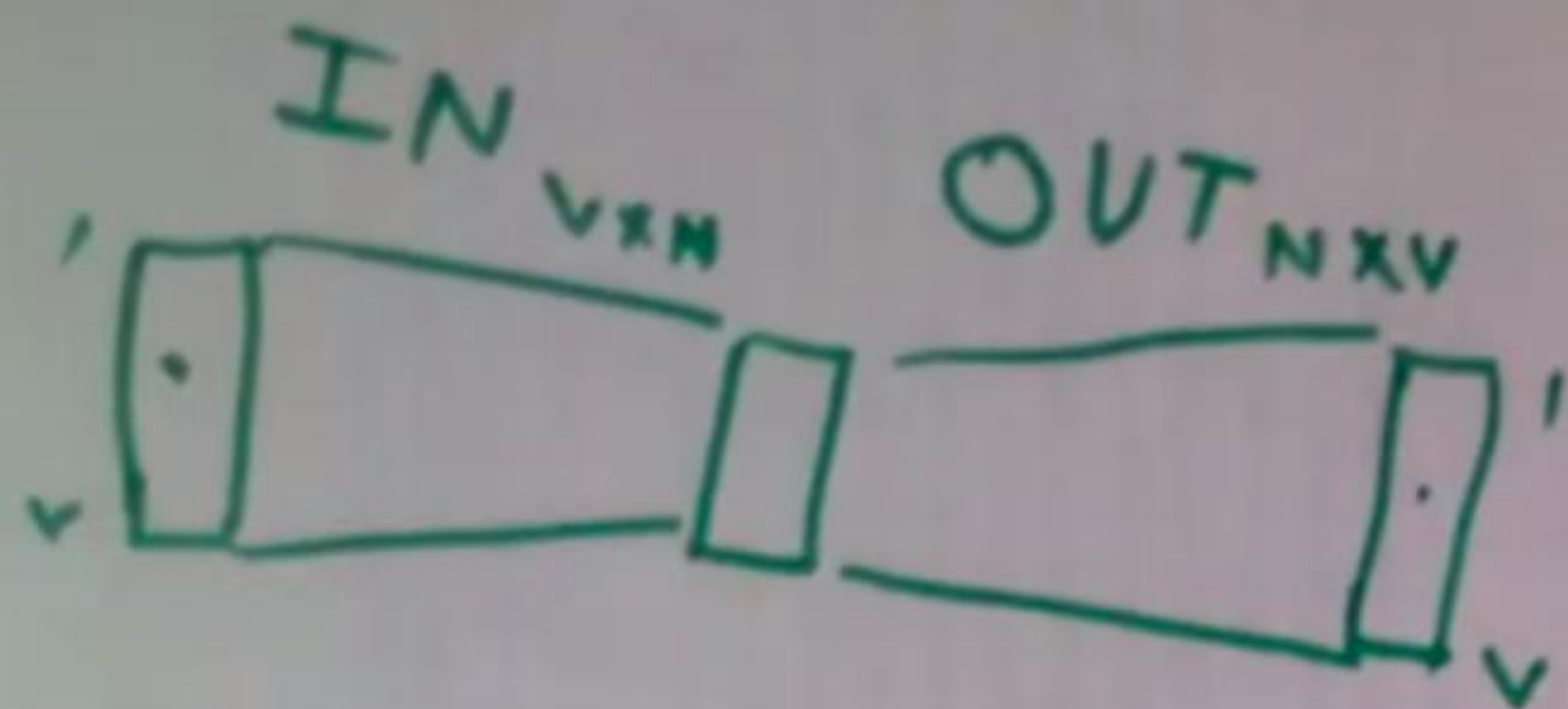




$$\begin{aligned}
 &P(w_1 | w_j) \\
 &P(w_2 | w_j) \\
 &P(w_3 | w_j) \\
 &\prod_{j=1}^T \prod_{i=1}^K P(w_i | w_j) = L(\Theta)
 \end{aligned}$$

$$F(\Theta) = -\frac{1}{T} \log(L(\Theta)) = -\frac{1}{T} \sum_{j=1}^T \sum_{i=1}^K P(w_i | w_j, \Theta)$$

$$P(w_j | \underline{w_i}, \theta)$$



Input
located
located
:

v_i

Output
Guwahati
in
:

u_j

Output Value

$$w_j : v_i \cdot u_j : e^{v_i \cdot u_j}$$

$$w_{j+1} : v_i \cdot u_{j+1} : e^{v_i \cdot u_{j+1}}$$

Soft Max

$$w_j : \frac{e^{v_i \cdot u_j}}{\sum_{k=1}^V e^{v_i \cdot u_k}}$$

$$p(w_{j^*} | w_i, \theta) = \frac{\exp(u_{j^*})}{\sum_{j'=1}^V \exp(u_{j'})} = y_{j^*} \quad E = -\log(y_{j^*})$$

$$-t_j + y_j \quad -u_{j^*} + \log \sum_{j'=1}^V \exp(u_{j'})$$

$$\frac{\partial E}{\partial w'_{ij}} =$$

$$\frac{\partial E}{\partial u_j}$$

$$\frac{\partial u_j}{\partial w'_{ij}}$$

 h_i

$$\frac{\partial u_{j^*}}{\partial u_j} = \begin{cases} 0 & \text{if } j \neq j^* \\ 1 & \text{if } j = j^* \end{cases}$$

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

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

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

$$w'_{ij} = w'_{ij} + \underbrace{\Delta w'_{ij}}_{-n \cdot h_i (y_j - t_j) e_j}$$
$$w'_{ij} - n \cdot e_j \cdot h_i$$

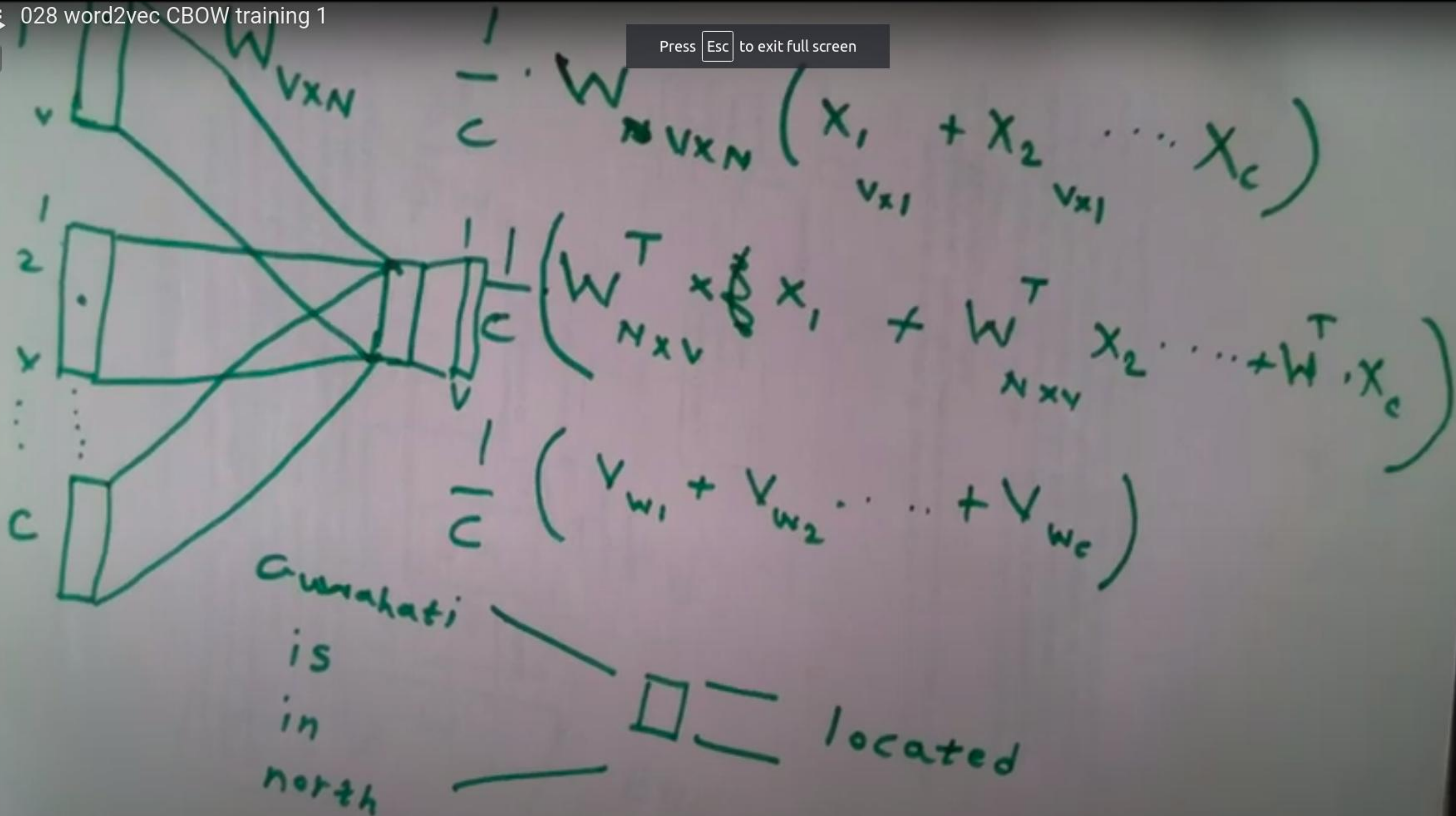
$$W_{ki} = W_{ki} + \Delta W_{ki}$$

$$-n \frac{\partial E}{\partial W_{ki}}$$

$$-n \cdot x_k \cdot E_{H_i}$$

$$\frac{\partial E}{\partial W_{ki}} = \left[\frac{\partial E}{\partial h_i} \right] \left[\frac{\partial h_i}{\partial W_{ki}} \right] x_k$$

$$\rightarrow \sum_{j=1}^n x_j \left[\frac{\partial E}{\partial u_j} \right] \left[\frac{\partial u_j}{\partial W_{ki}} \right] w'_{ji}$$



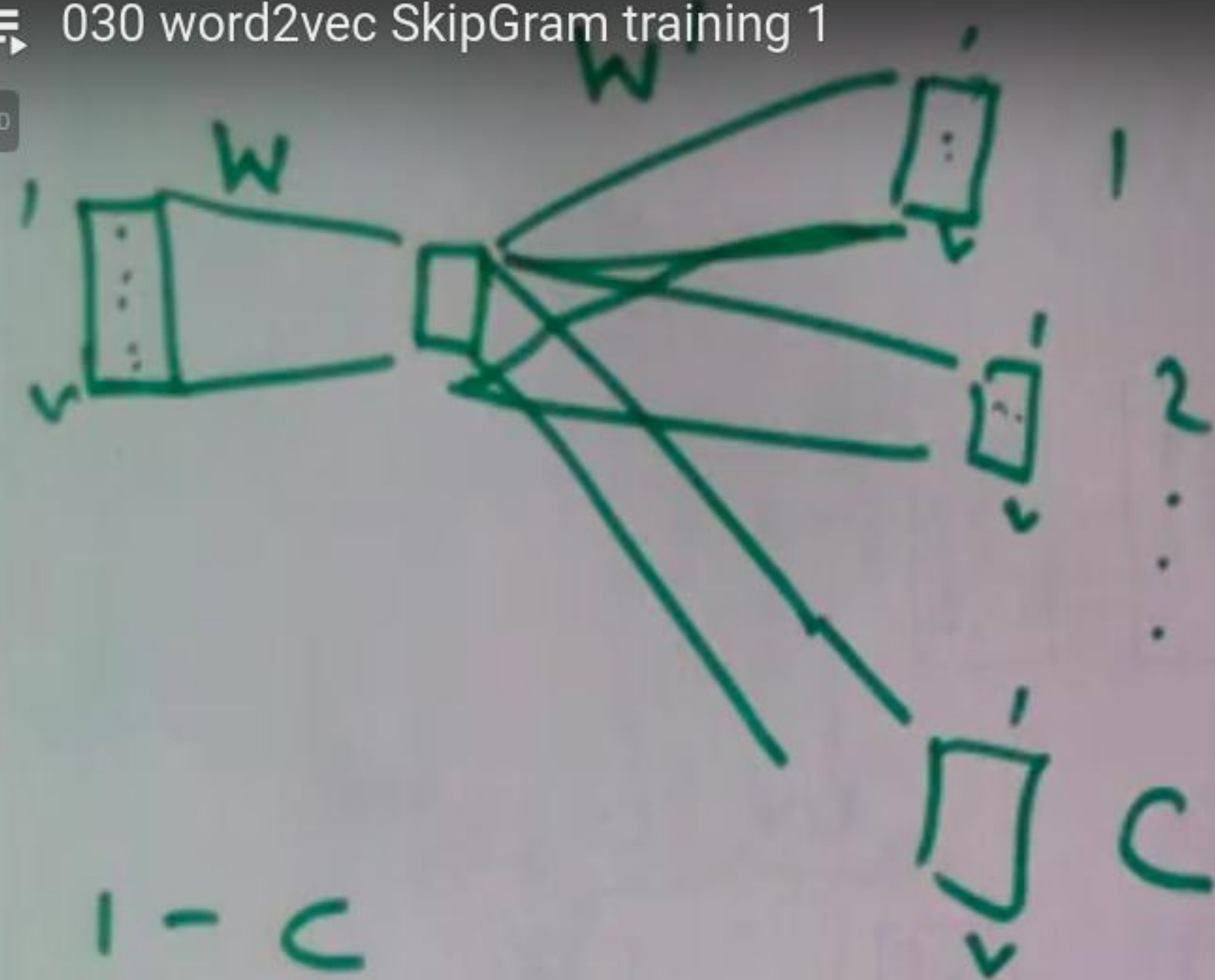
$$E = -\log(p(w_{j*} | w_1, w_2, \dots, w_c))$$

$$= -u_{j*} + \log \sum_{j'=1}^v \exp(u_{j'})$$

~~the~~

$$w_{ki} = w_{ki} + \frac{\Delta w_{ki}}{-\eta \cdot \frac{\partial E}{\partial w_{ki}}}$$

$$-\eta \cdot \alpha_k \cdot E H_i \cdot \frac{1}{c}$$

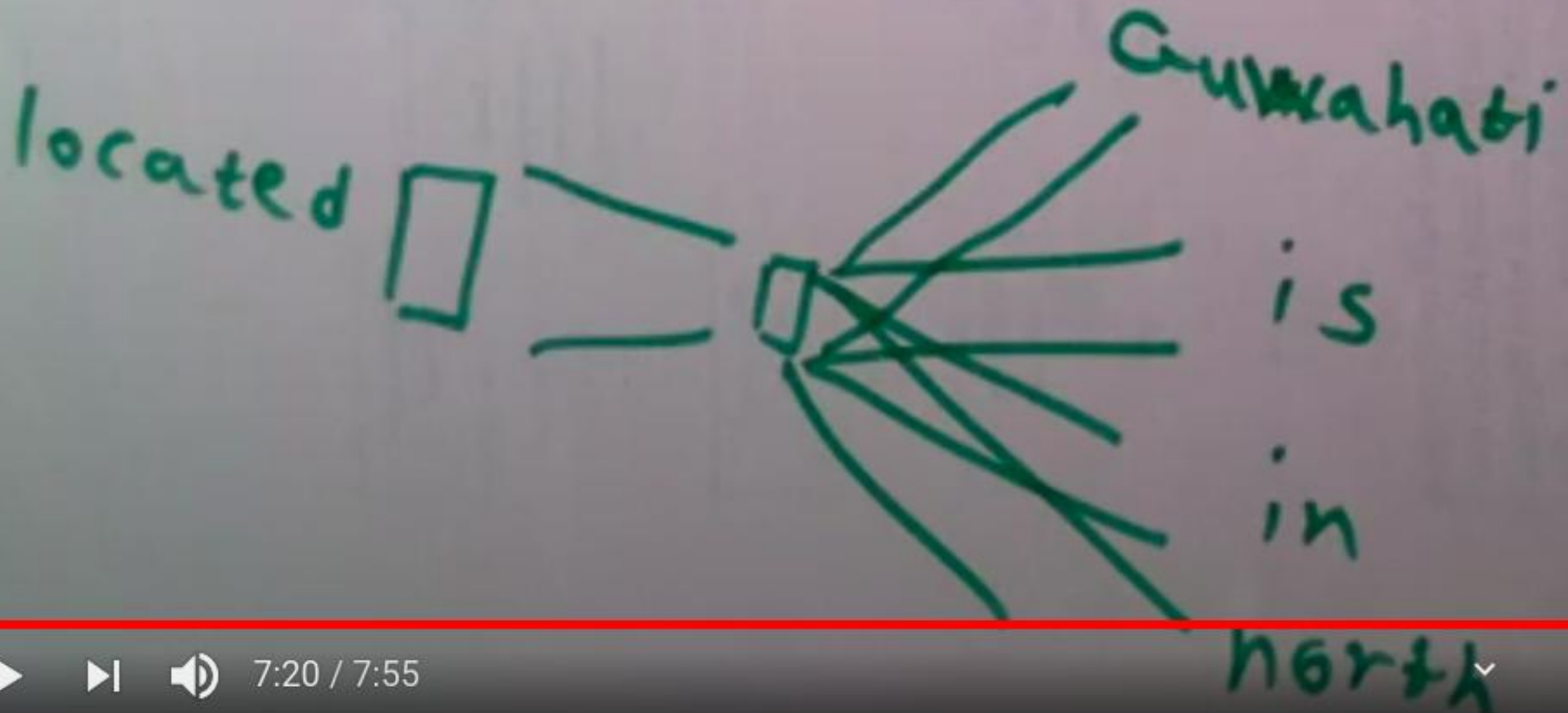


$$W_{V \times N} \cdot X_{V \times 1}$$

$$W^T_{N \times V} \cdot X_{V \times 1} \quad h_{N \times 1}$$

$$P(w_{c,j} | w_i) = y_{c,j} = \frac{\exp(u_{c,j})}{\sum_{j'=1}^V \exp(u_{j'})}$$

$$u_{c,j} = u_j = h \cdot v'_j$$



$$E = -\log \prod_{c=1}^C p(w_{oc} | w_i) \rightarrow$$

$$\sum_{c=1}^C \log p(w_{oc} | w_i)$$

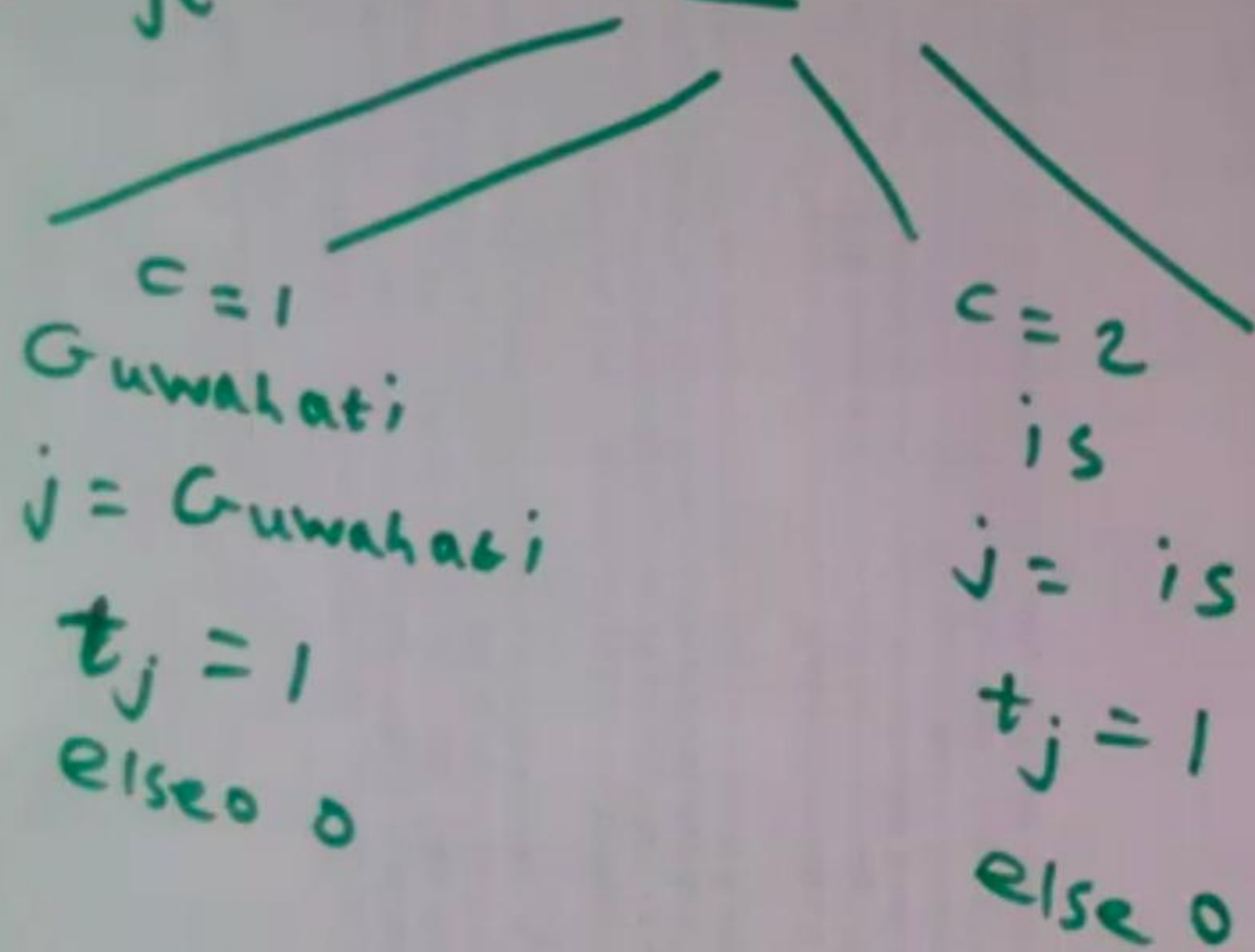
$$\boxed{\exp(u_{cj_c^*})}$$

$$\boxed{\sum_{j'=1}^V \exp(u_{j'})}$$

$$= -\sum_{c=1}^C u_{j_c^*} + C \cdot \log \sum_{j'=1}^V \exp(u_{j'})$$

$$\Delta w'_{ij} = -\eta \cdot \frac{\partial E}{\partial w'_{ij}} = -\eta \cdot \sum_{c=1}^C \frac{\partial E}{\partial u_{cj_c^*}}$$

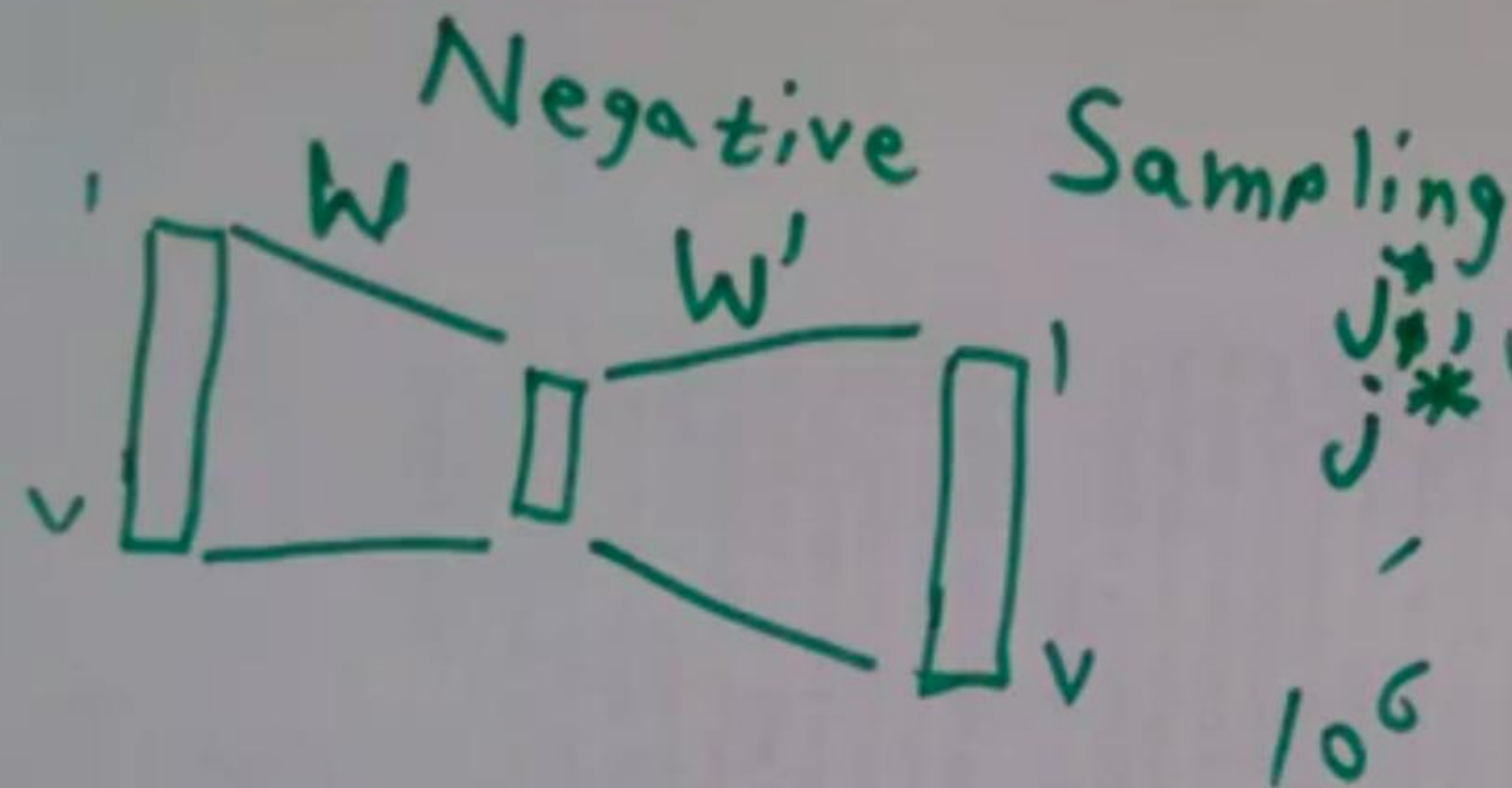
$$\frac{\partial E}{\partial u_{jc}} = y_{cj} - t_{cj} = e_{cj}$$



$$EI_j = \sum_{c=1}^c e_{cj}$$

$$\Delta W'_{ij} = -n \cdot \left(\sum_{c=1}^c e_{cj} \right) \cdot h_i$$

$$-n \cdot h_i \cdot EI_j$$



$j_1^*, j_2^*, \dots, j_c^*$

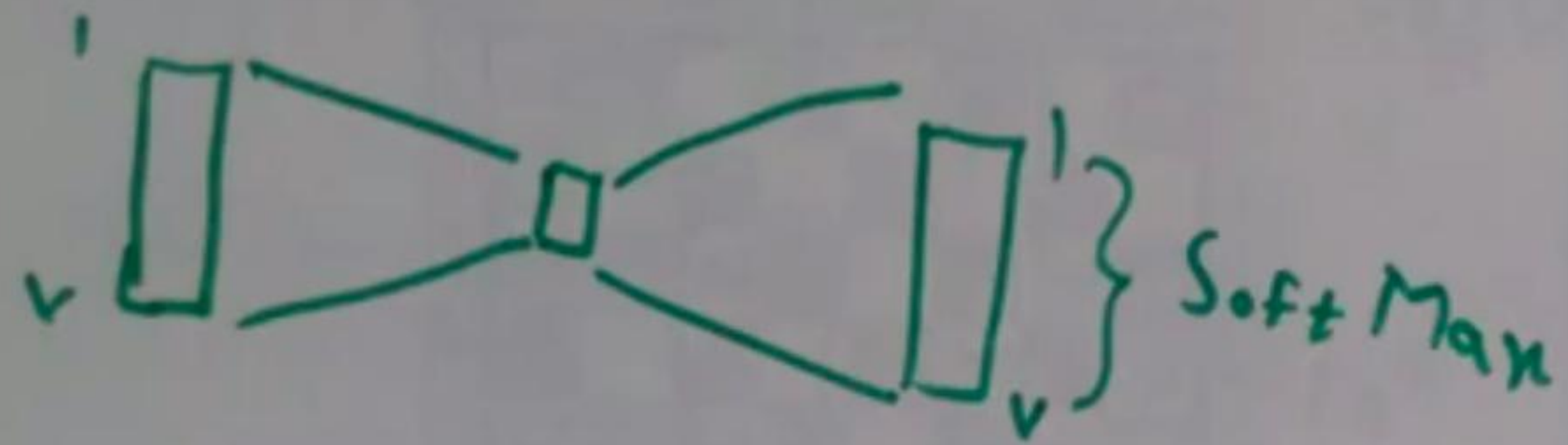
5

5

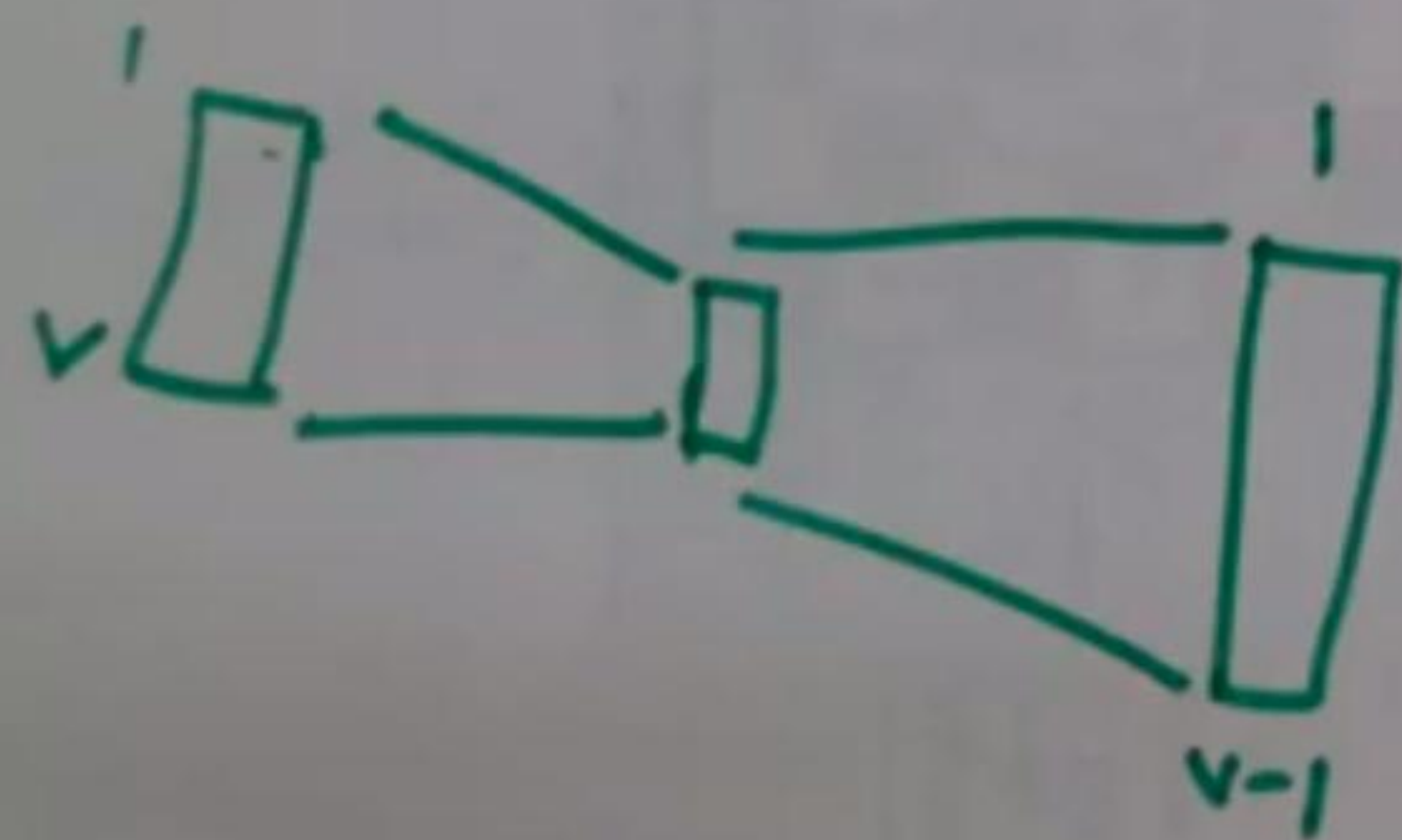
$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^c f(w_j)^{3/4}}$$

$\frac{6}{108}$

Hierarchical Soft Man

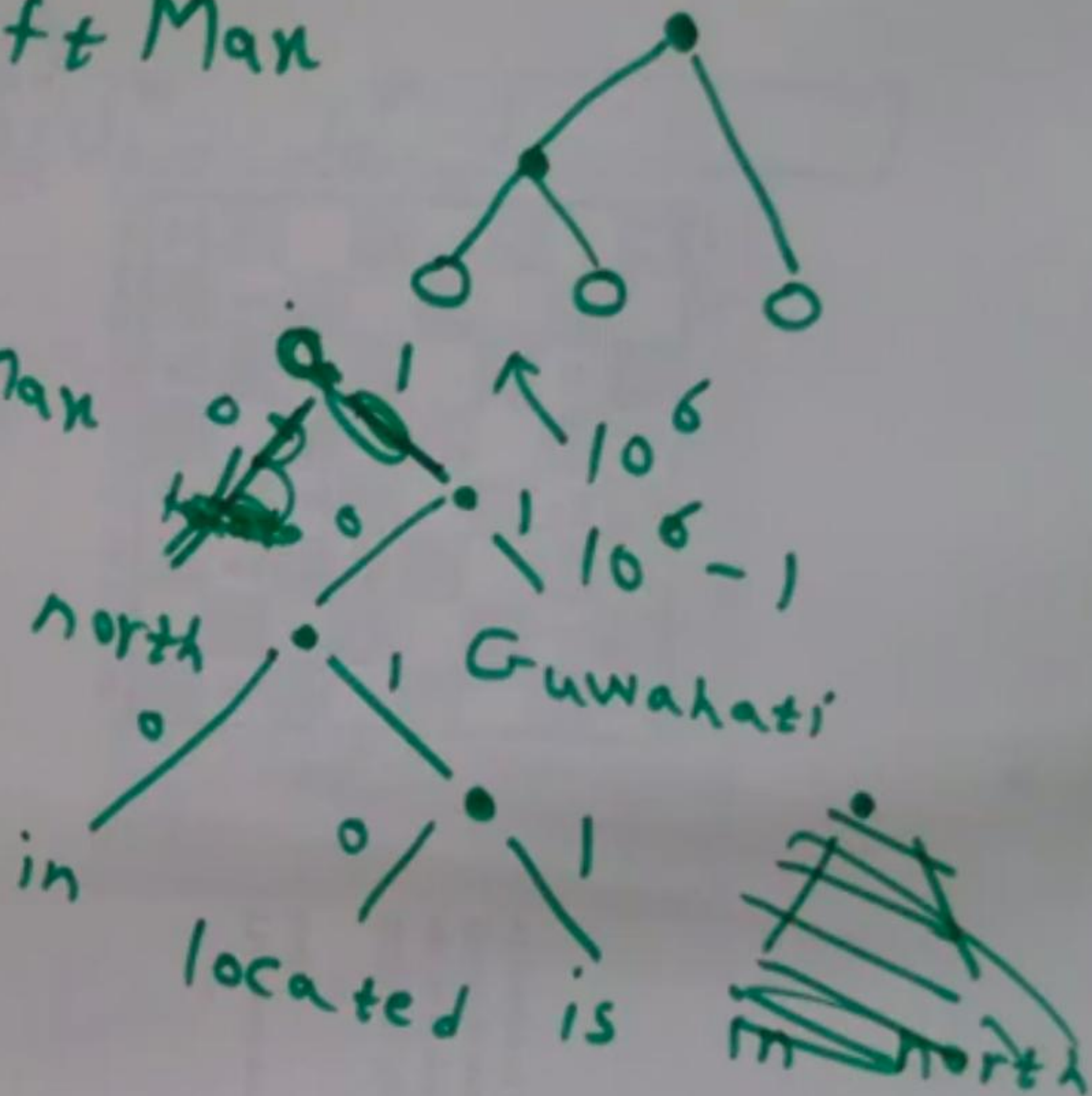


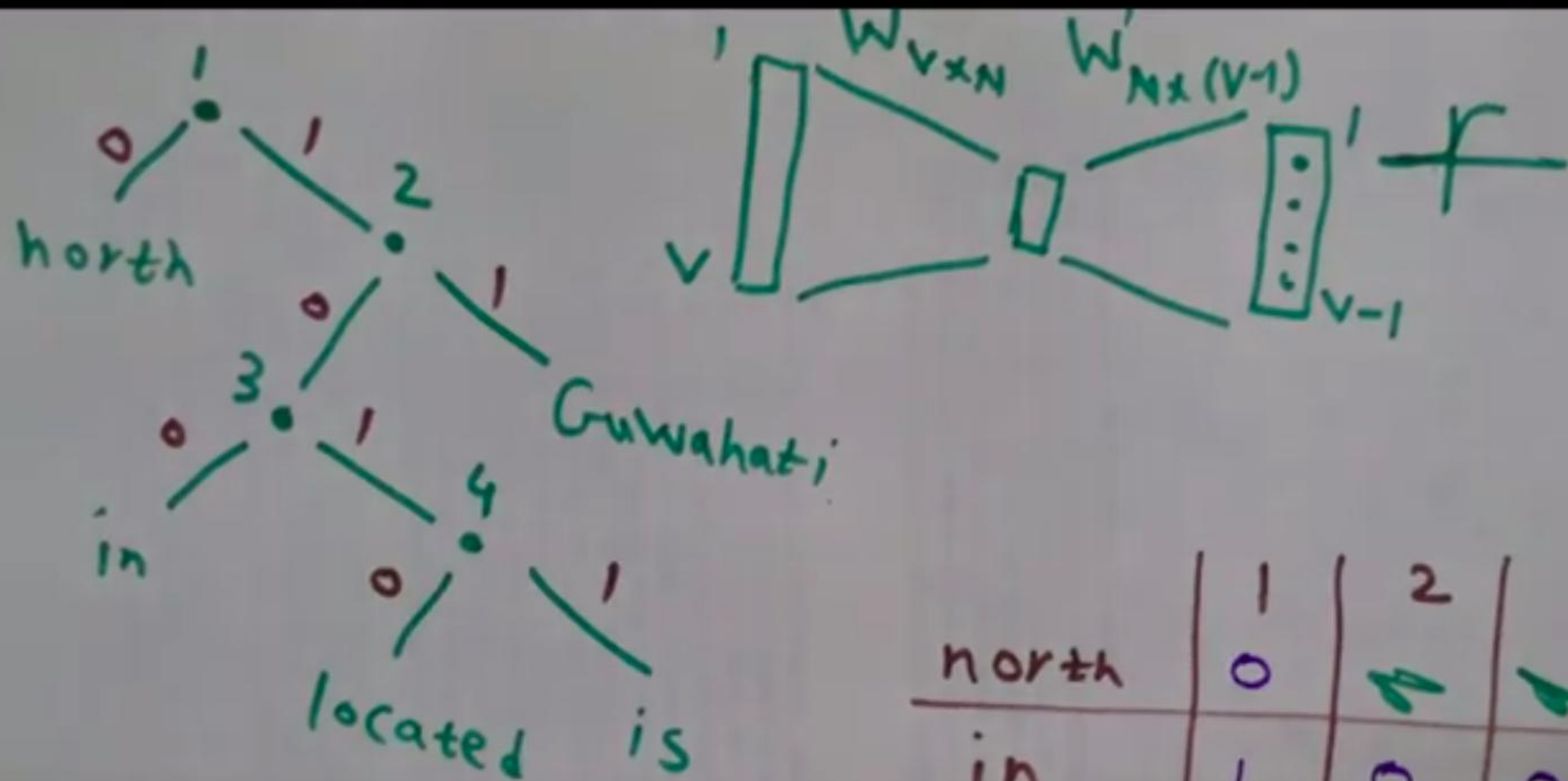
Soft Man



north: 0
in: 100

Guwahati: 11





north: 0
in : 100

	1	2	3	4
north	0	1	1	1
in	1	0	0	1
located	1	0	1	0
Guwahati	1	1	1	1
is	1	0	1	1

Phrases

→ Times of India

Indian Express

$$\text{score}(w_i, w_j) = \frac{\text{count}(w_i, w_j) - \delta}{\text{count}(w_i) \times \text{count}(w_j)}$$

Evaluation of Word Vectors

Objective

- Training speed
- Memory requirement
- Large training data
- Downstream task

Subjective

- Word analogies

$$\begin{aligned} &\text{vector}(\text{King}) \\ &- \text{vector}(\text{Male}) \\ &+ \text{vector}(\text{Female}) \end{aligned}$$

vector(X)

bigger

- big
+ small

smaller

Country: Capital

India: New Delhi

USA: ?

Bigger: Big :: Queen

Smaller: Small

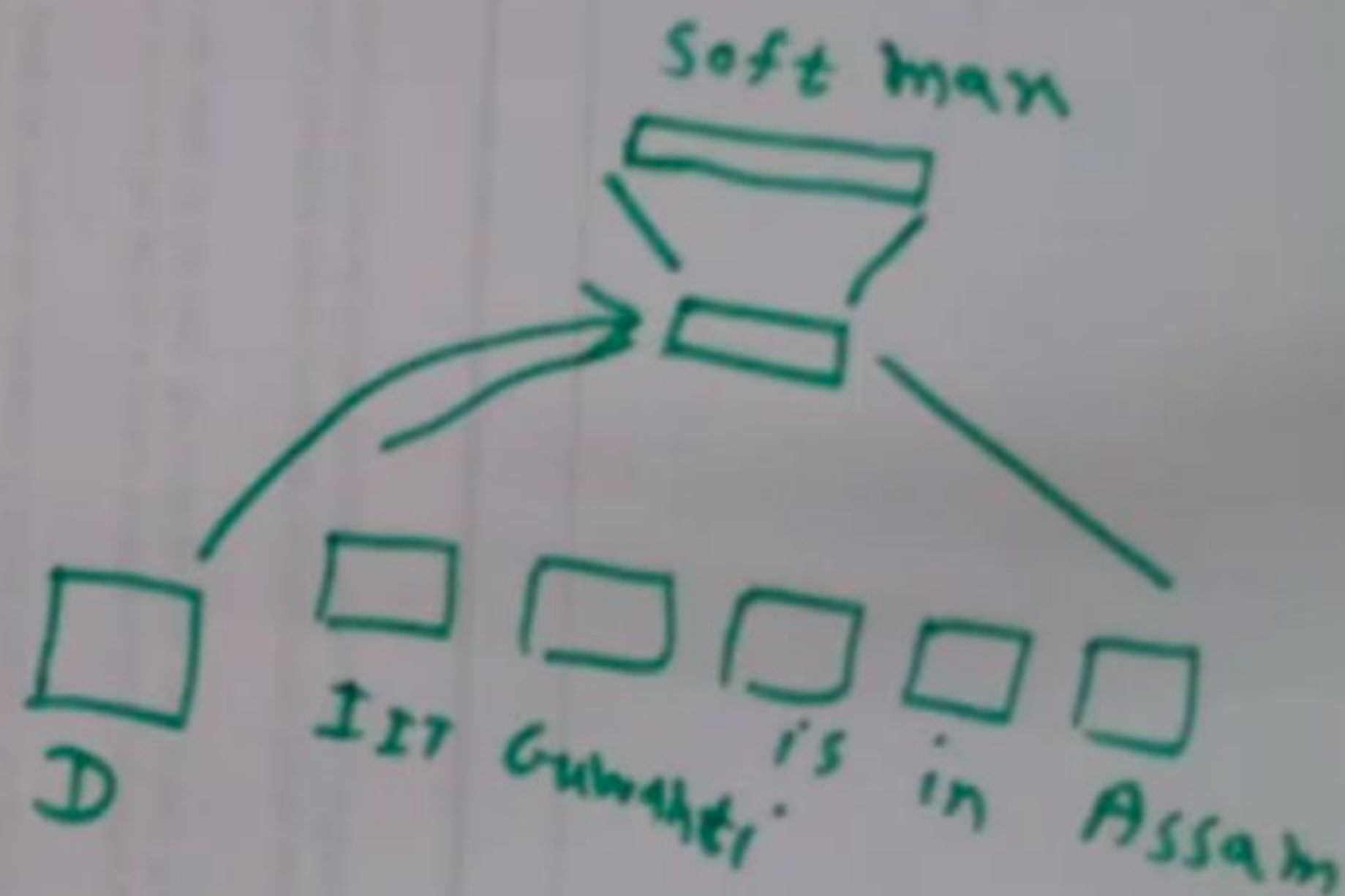
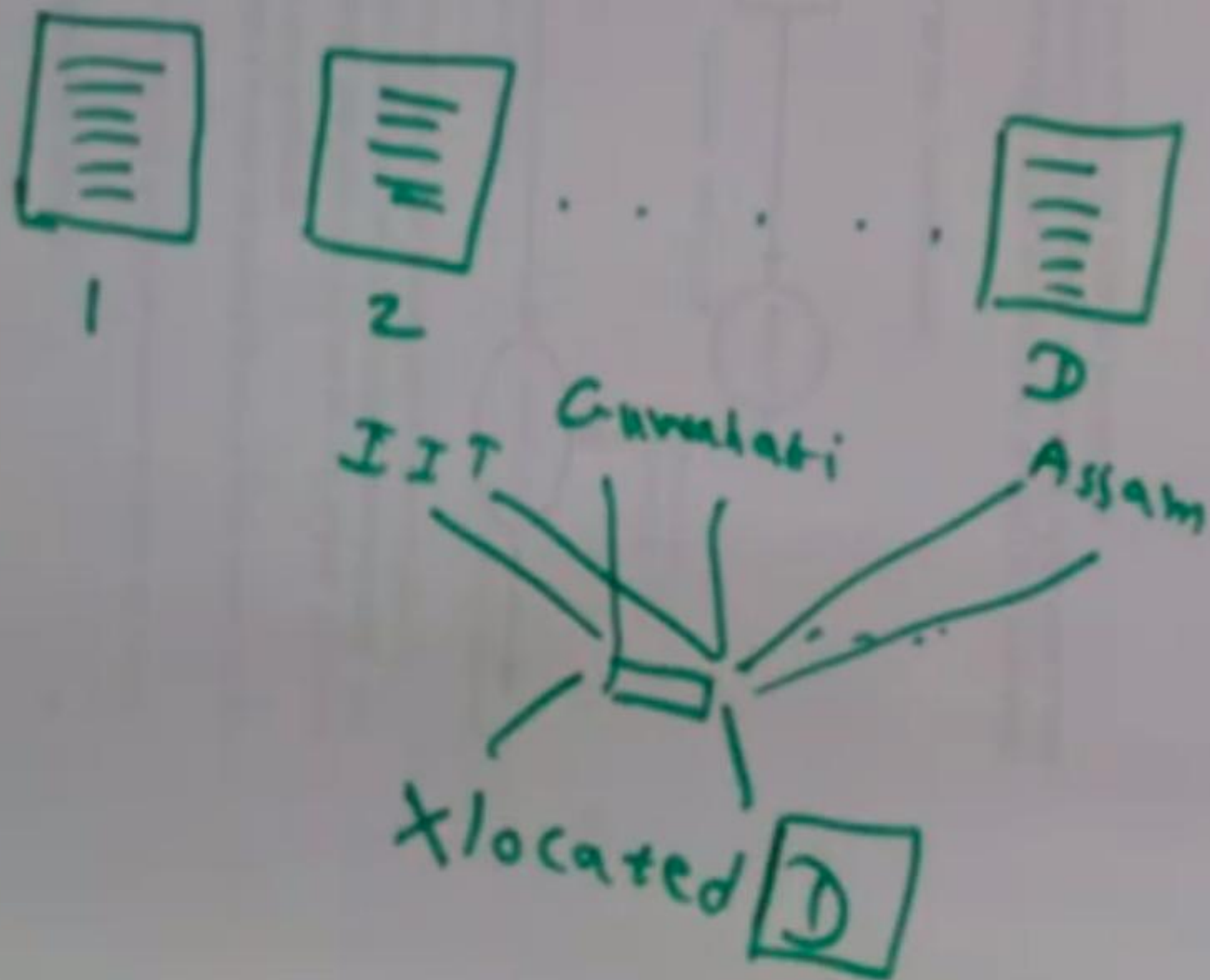
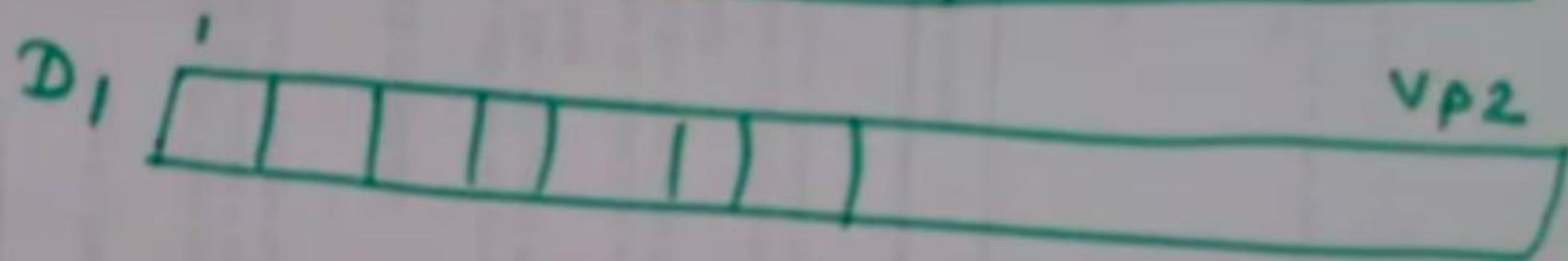
King: Male :: Queen: Female

Documents

Bag of words: tf idf

Bag of n-grams:

doc2vec

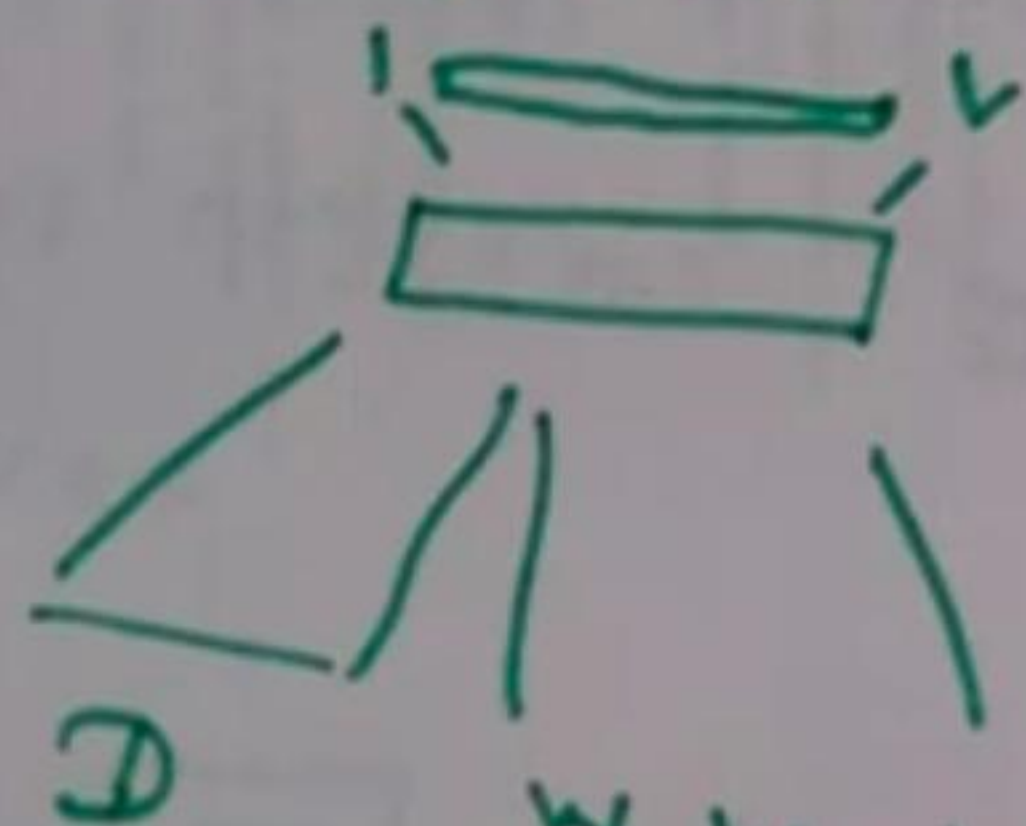


PV-DM: Paragraph Vector Distributed Memory

N: Paragraphs / Documents
V: Vocabulary size

p

$$q \quad (N \times p) + (V \times q)$$



	D_1	D_2	D_3
1	0.7	0.5	0.1
2	0.9	0.6	0.2

$$\begin{aligned} 0.7 &\rightarrow 0.7001 \\ 0.9 &\rightarrow 0.89999 \end{aligned}$$

Finance

bank
money
lending
interest
player
is
cost

w_1, w_2, \dots, w_c
Sports

bat
ball
money
is
this

computer science

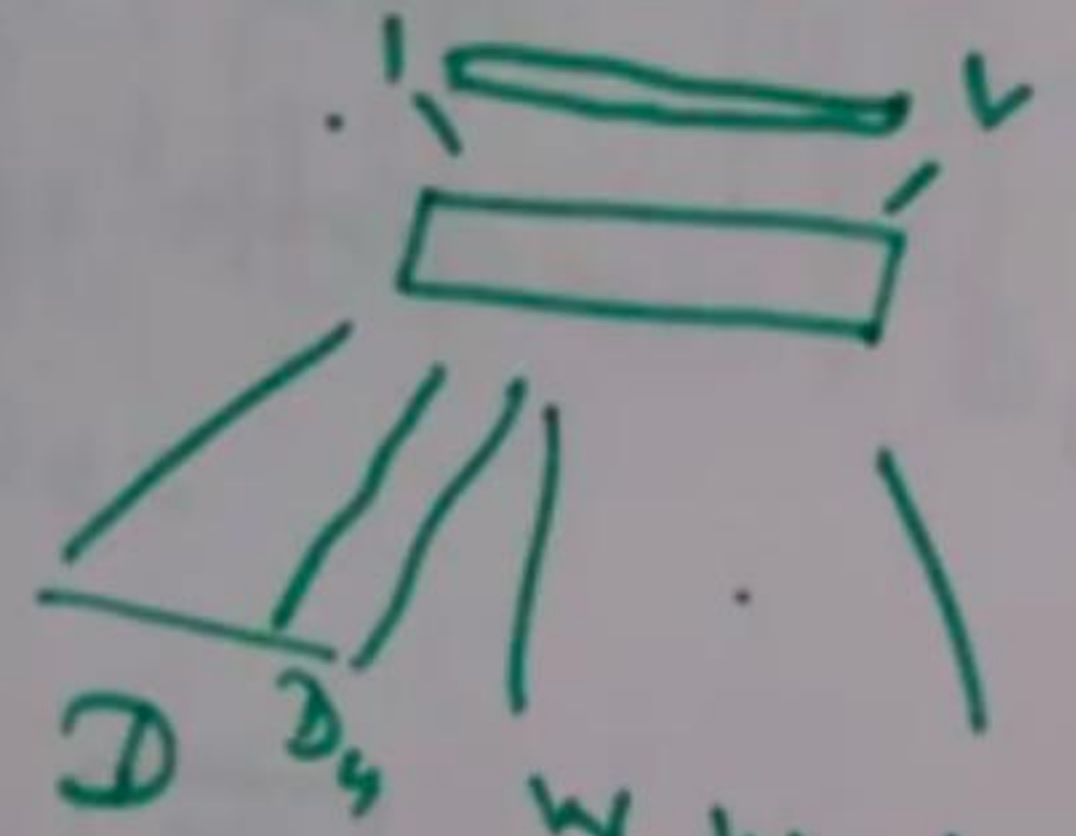
algorithm
cost
is
this

Paragraph Vector Distributed Memory

N: Paragraphs / Documents
 V: Vocabulary size

: p

$$: q \quad (N \times p) + (V \times q)$$



	D1	D2	D3	D4	
1	0.7	0.5	0.1	0.7	0.7
2	0.9	0.6	0.2	0.2	0.9

0.7 → 0.7001
 0.9 → 0.89999

Finance

bank
 money
 lending
 interest
 player
 is
 cost

Sports
 w1, w2, ..., wc

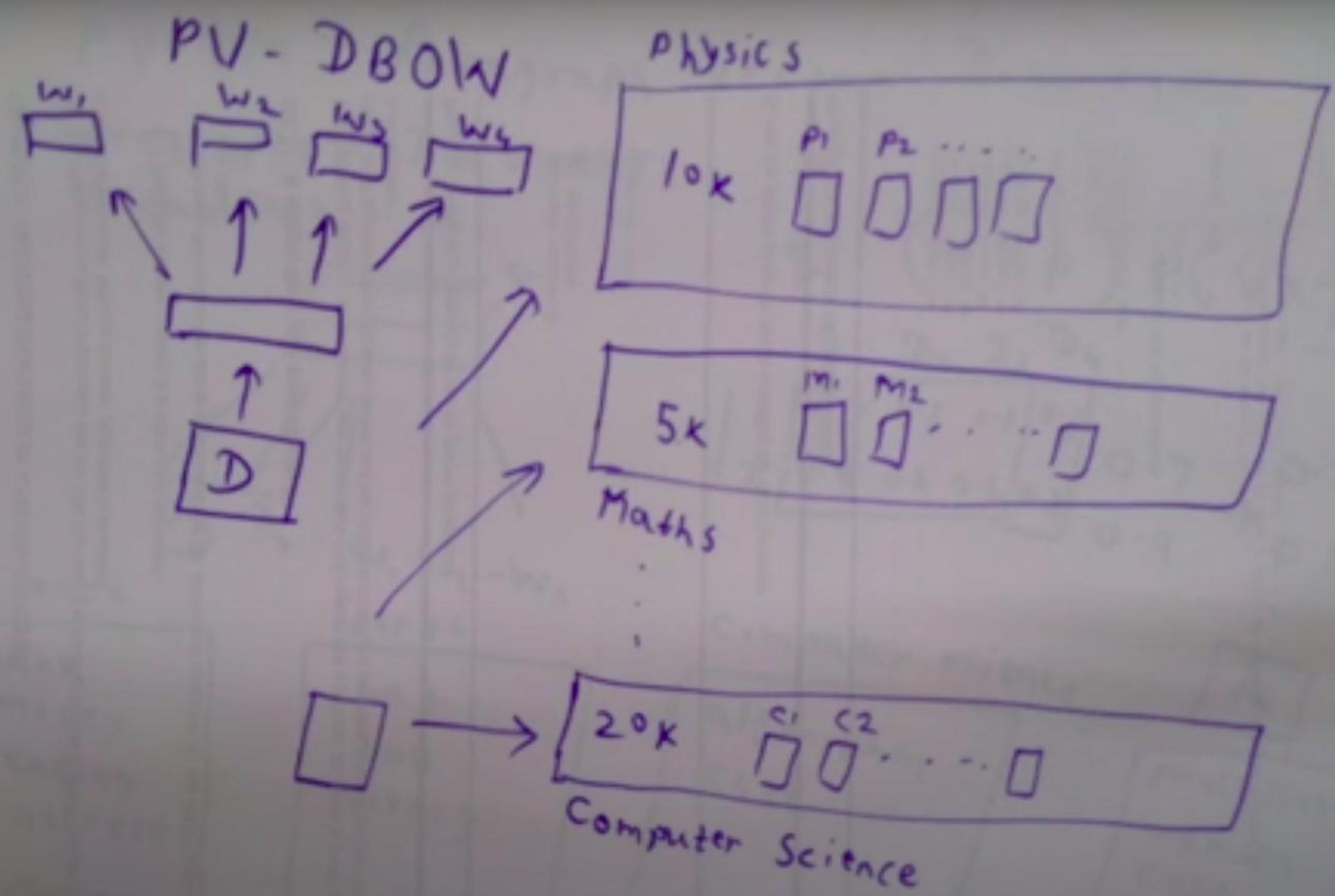
bat
 ball
 money
 ...
 is
 this

Computer science

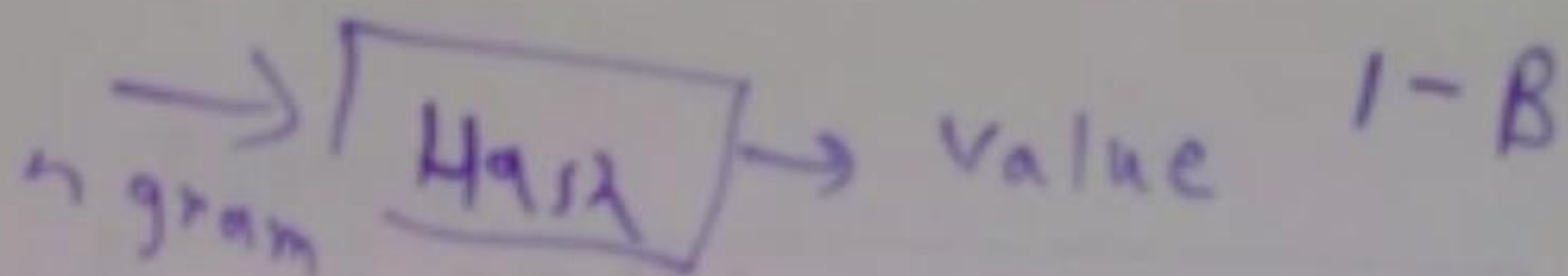
algorithm
 cost
 ...
 is
 this

D4

processor
 memory
 cost
 network
 is
 this



Fast Text



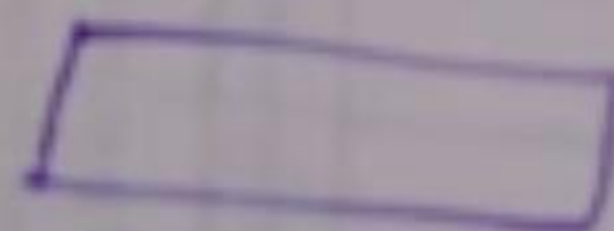
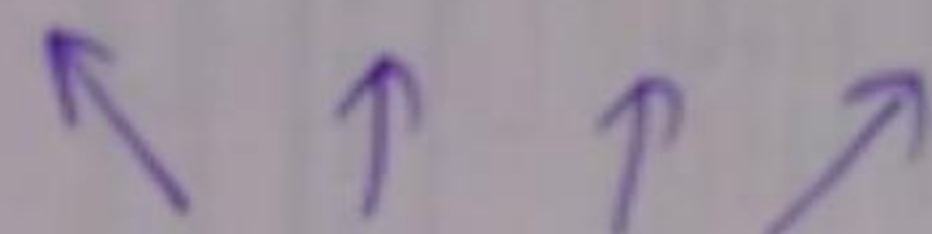
India

Indian

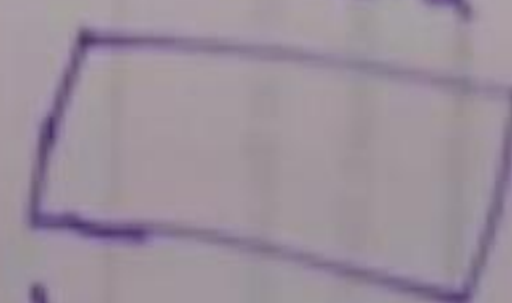
3 gram

1 < In
2 Ind
3 ndi
4 dia
5 ia >

< In
Ind
ndi
dia
ian
an >

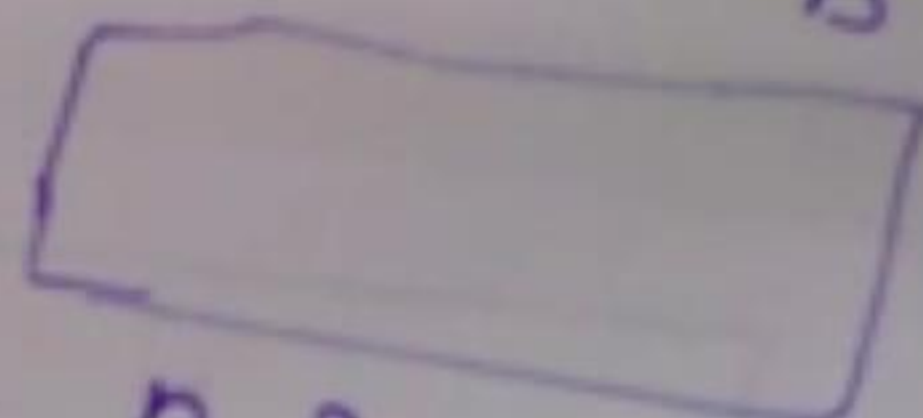


India



Word
vectors

1 — B



n gram
vectors

India \rightarrow Asia Population Software Education