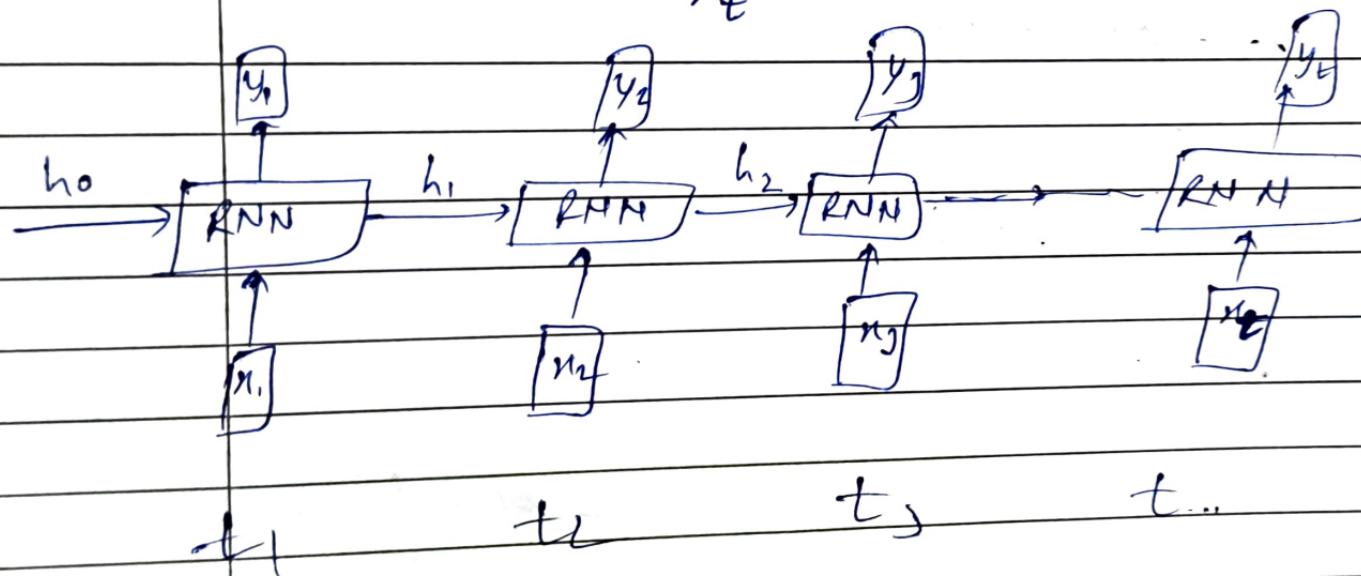
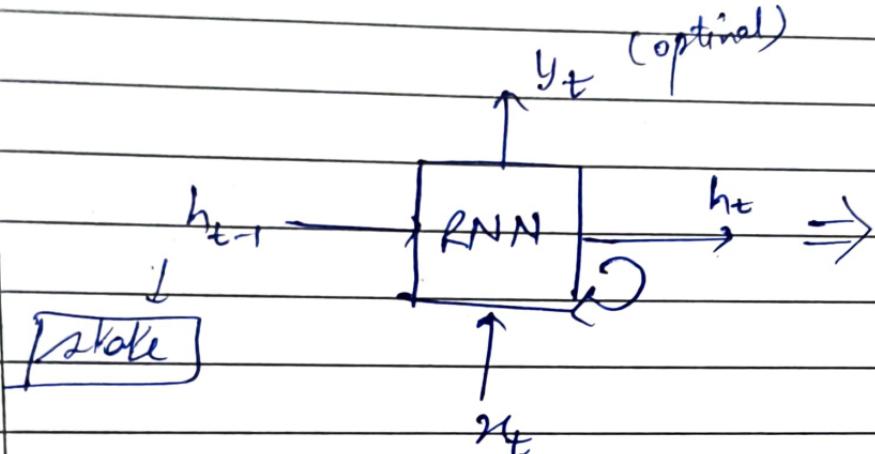
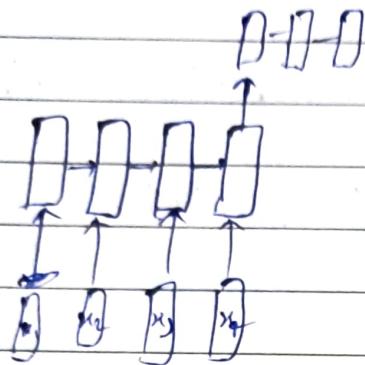


16/9/25

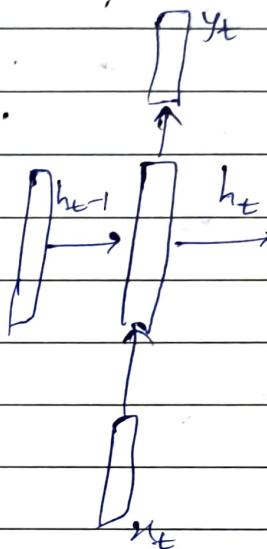


# Machine translation example (english → french)



$$h_t = f_w(h_{t-1}, x_t)$$

$$h_t = \tanh(w_{hh} h_{t-1} + w_{xh} x_t)$$



$$h_t = \tanh \left( w_{hh} h_{t-1} + w_{xh} x_t \right)$$

Below the equation, dimensions are indicated:  $h_t$  is  $4 \times 1$ ,  $w_{hh}$  is  $4 \times 4$ ,  $h_{t-1}$  is  $4 \times 1$ , and  $w_{xh}$  is  $4 \times 1$ .

$$y_t = w_y h_t$$

$$y_t = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Below the equation, dimensions are indicated:  $y_t$  is  $2 \times 1$ ,  $w_y$  is  $4 \times 4$ , and  $h_t$  is  $4 \times 1$ .

①  $h_t$

②  $y_t$

Q1

 $T = 2$  (two time steps)

$$x^{(1)}: x_1 \in \mathbb{R}^2 = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, x_2 \in \mathbb{R}^2 = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

input dim = 2  $\Rightarrow x_t \in \mathbb{R}^2$ hidden state dim = 3  $\Rightarrow h_t \in \mathbb{R}^3$ 

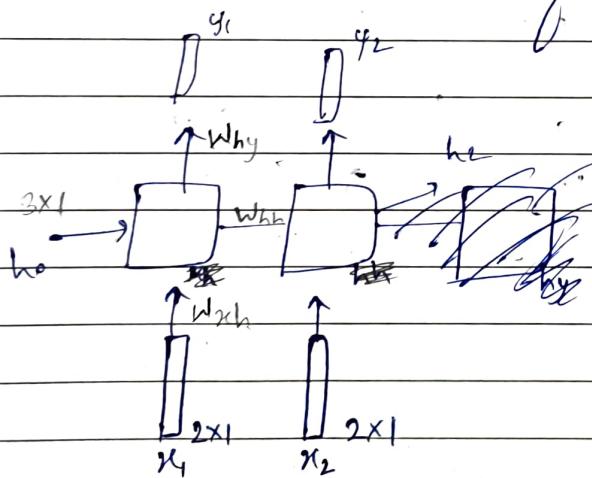
$$h_0 = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, W_{xh} = \begin{bmatrix} 0.5 & -0.3 \\ 0.8 & 0.2 \\ 0.1 & 0.4 \end{bmatrix}, W_{hh} = \begin{bmatrix} 0.1 & 0.4 & 0 \\ -0.2 & 0.3 & 0.1 \\ 0.05 & -0.1 & 0.2 \end{bmatrix}$$

$$y_t \in \mathbb{R}^2$$

$$W_{hy} = \begin{bmatrix} 1 & -1 & 0.5 \\ 0.5 & 0.5 & -0.5 \end{bmatrix}$$

what will be dimensions of all these?

Ans



$$h_t = \tanh(W_{hh} h_{t-1} + W_{xh} x_t)$$

$$y_t = W_{hy} h_t$$

Time step 21, t21

$$h_t = \tanh(W_{hh} h_{t-1} + W_{xh} x_t)$$

$$h_1 = \tanh(W_{hh} h_0 + W_{xh} x_1)$$

$$= \tanh([ ])$$

$$\begin{bmatrix} ? \\ ? \\ ? \end{bmatrix}_{3 \times 1}$$

$$y_t = W_{hy} h_t$$

$$y_1 = W_{hy} h_1$$

Timestep = 2,  $t = 2$

$$h_t = \tanh(W_{hh} h_{t-1} + W_{xh} x_t)$$

At  $t = 1$

$$h_1 = \tanh \left( \begin{bmatrix} 0.1 & 0.4 & 0 \\ -0.2 & 0.3 & 0.2 \\ 0.05 & -0.1 & 0.2 \end{bmatrix}_{3 \times 3} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}_{3 \times 1} + \begin{bmatrix} 0.5 & -0.1 \\ 0.8 & 0.2 \\ 0.1 & 0.4 \end{bmatrix}_{3 \times 2} \begin{bmatrix} 1 \\ 2 \end{bmatrix}_{2 \times 1} \right)$$

$$h_1 = \tanh \left( \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}_{3 \times 1} + \begin{bmatrix} -0.1 \\ 1.2 \\ 0.9 \end{bmatrix}_{3 \times 1} \right)$$

$$h_1 = \tanh \left( \begin{bmatrix} -0.1 \\ 1.2 \\ 0.9 \end{bmatrix} \right) \approx \begin{bmatrix} 0.0995 \\ 0.834 \\ 0.7167 \end{bmatrix} \approx \begin{bmatrix} 0.099 \\ 0.83 \\ 0.716 \end{bmatrix}$$

$$y_1 = W_{hy} h_1$$

$$y_1 = \begin{bmatrix} 1 & -1 & 0.5 \\ 0.5 & 0.5 & -0.5 \end{bmatrix}_{2 \times 3} \begin{bmatrix} 0.0995 \\ 0.834 \\ 0.7167 \end{bmatrix}_{3 \times 1}$$

$$y_1 \approx \begin{bmatrix} -0.37615 \\ 0.1084 \end{bmatrix} \approx$$

0.09975

+ 0.417

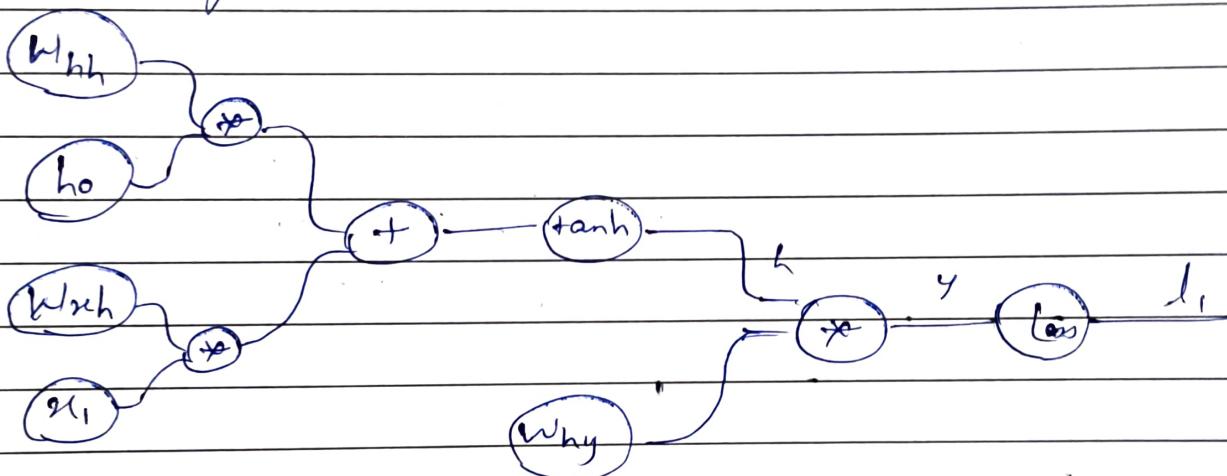
- 0.35835

0.1084

# # Word2Vec

Eg. → Simplified character level Language Model RNN

Computational graph



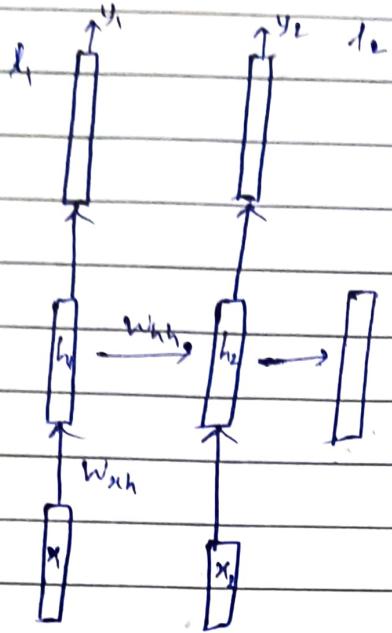
$$\frac{\partial L}{\partial w_{hh}} + \frac{\partial L}{\partial w_{hh}} + \dots \rightarrow \frac{\partial L}{\partial w_{hh}}$$

$$\frac{\partial L}{\partial w_{hh}}$$

$$w_{hh} = w_{hh} - \alpha \frac{\partial L}{\partial w_{hh}}$$

22/9/25

sum the  
whole loss  
and compute  
the gradient



$$L = L_1 + L_2 + L_3 + \dots + L_T$$

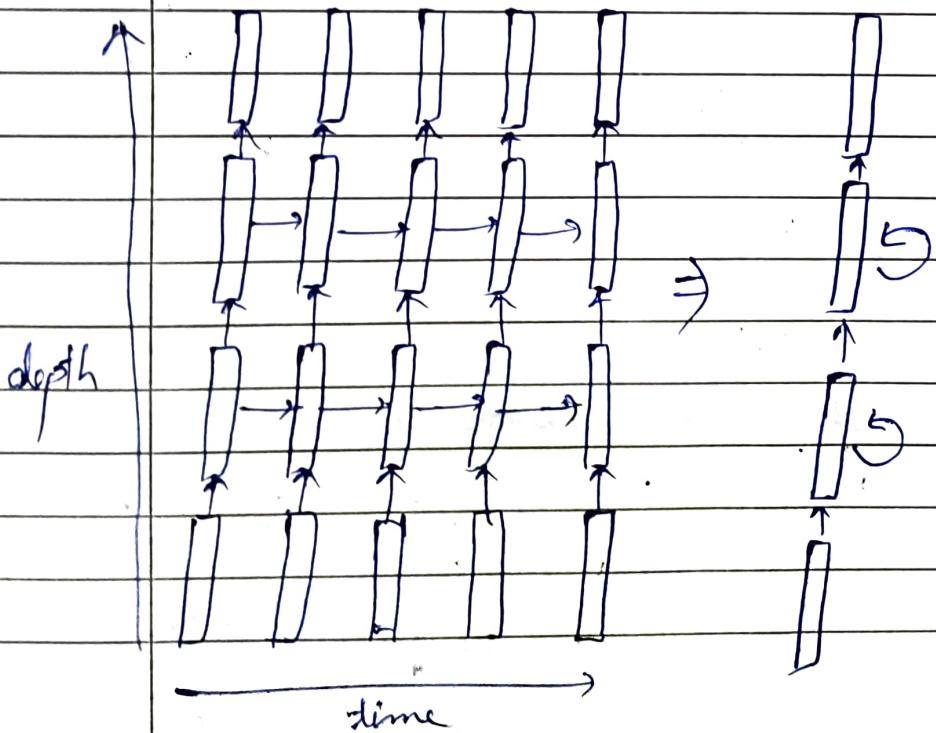
$$\frac{\partial L}{\partial w_{hh}} \quad w_{hh} \leftarrow \frac{\partial L}{\partial w_{hh}}$$

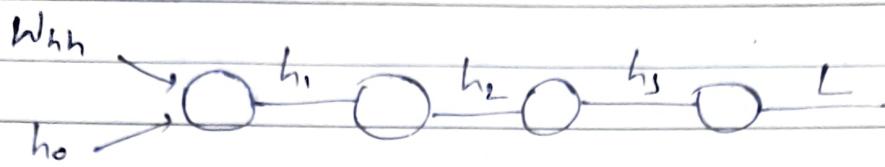
$$\frac{\partial L}{\partial w_{xh}}$$

$$\frac{\partial L}{\partial h_{ny}}$$

BPTT → Backpropagation Through time

# Multilayer RNN





$$\frac{\partial L}{\partial w_{hh}} = \frac{\partial L}{\partial h_3} \cdot \frac{\partial h_3}{\partial h_2} \cdot \frac{\partial h_2}{\partial h_1} \cdot \frac{\partial h_1}{\partial w_{hh}}$$

$$= \frac{\partial L}{\partial h_3} \left( \frac{\partial h_3}{\partial h_2} \cdot \frac{\partial h_2}{\partial h_1} \right) \cdot \frac{\partial h_1}{\partial w_{hh}}$$

$$= \frac{\partial L}{\partial h_3} \left( \prod_{t=2}^T \frac{\partial h_t}{\partial h_{t-1}} \right) \cdot \frac{\partial h_1}{\partial w_{hh}}$$

← generalized form

for T time steps

$$\frac{\partial L}{\partial w_{hh}} = \frac{\partial L_T}{\partial h_T} \cdot \left( \prod_{t=2}^T \frac{\partial h_t}{\partial h_{t-1}} \right) \frac{\partial h_1}{\partial w_{hh}}$$

$$h_t = \tanh(w_{hh} h_{t-1} + w_{xh} x_t)$$

$$\frac{\partial h_t}{\partial h_{t-1}} = \tanh'(w_{hh} h_{t-1} + w_{xh} x_t) w_{hh}$$

$$\frac{\partial L}{\partial w_{hh}} = \frac{\partial L_T}{\partial h_T} \left( \prod_{t=2}^T \tanh'(w_{hh} h_{t-1} + w_{xh} x_t) w_{hh}^{t-1} \right) \frac{\partial h_1}{\partial w_{hh}}$$

If no non-linearity

$$h_t = w_{hh} h_{t-1} + w_{xh} x_t$$

$$\frac{\partial h_t}{\partial h_{t-1}} = w_{hh}$$

\*high vanishing gradient

problem to overcome this, LSTM intro

24/3/25

## LSTM

Q.

$$x_t = [0.5, -0.1]$$

$$w_{hi} = \begin{bmatrix} 0.1 & 0.2 \\ -0.2 & 0.05 \end{bmatrix}$$

$$C_{t-1} = [0.2, -0.2]$$

$$w_{hf} = \begin{bmatrix} -0.4 & 0.2 \\ -0.3 & 0.3 \end{bmatrix}$$

$$h_{t-1} = [0.0, 0.1]$$

$$w_{xi} = \begin{bmatrix} 0.5 & -0.3 \\ 0.4 & 0.1 \end{bmatrix}$$

$$w_{if} = \begin{bmatrix} 0.05 & -0.1 \\ 0.2 & 0.1 \end{bmatrix}$$

$$w_{ho} = \begin{bmatrix} 0.3 & 0.25 \\ -0.2 & 0.2 \end{bmatrix}$$

$$w_{hg} = \begin{bmatrix} -0.5 & 0.4 \\ 0.2 & -0.3 \end{bmatrix}$$

$$w_{xf} = \begin{bmatrix} 0.15 & 0.05 \\ 0.1 & -0.2 \end{bmatrix}$$

$$w_{if} = \begin{bmatrix} 0.2 & 0.1 \\ -0.1 & 0.05 \end{bmatrix}$$

$$h_t = ?$$

$$C_t = ?$$

ans.  $h_t = o_t \odot \tanh(C_t)$

$$C_t = f_t \odot C_{t-1} + i_t \odot g_t$$

$$f_t = \sigma(w_{hf} h_{t-1} + w_{xf} x_t)$$

$$f_t = \sigma \left( \begin{bmatrix} 0.05 & -0.1 \\ 0.2 & 0.1 \end{bmatrix} \begin{bmatrix} 0 \\ 0.1 \end{bmatrix} + \begin{bmatrix} -0.4 & 0.2 \\ 0.3 & 0.1 \end{bmatrix} \begin{bmatrix} 0.5 \\ -0.1 \end{bmatrix} \right)$$

$$= \sigma \left( \begin{bmatrix} 0.05 & -0.1 \\ 0.02 & 0.1 \end{bmatrix} + \begin{bmatrix} -0.22 \\ 0.12 \end{bmatrix} \right).$$

~~$$= \sigma \left( \begin{bmatrix} -0.23 \\ 0.13 \end{bmatrix} \right) = \begin{bmatrix} 0.555 \\ 0.467 \end{bmatrix}$$~~

$$= \sigma \left( \begin{bmatrix} -0.01 \\ 0.01 \end{bmatrix} + \begin{bmatrix} -0.22 \\ 0.12 \end{bmatrix} \right)$$

$$= \sigma \left( \begin{bmatrix} -0.23 \\ 0.13 \end{bmatrix} \right) = \begin{bmatrix} 0.44 \\ 0.532 \end{bmatrix}$$

$$i_t = \sigma \left( \begin{bmatrix} 0.1 & 0.2 \\ -0.2 & 0.05 \end{bmatrix} \begin{bmatrix} 0.0 \\ 0.1 \end{bmatrix} + \begin{bmatrix} 0.5 & -0.3 \\ 0.4 & 0.1 \end{bmatrix} \begin{bmatrix} 0.5 \\ -0.1 \end{bmatrix} \right)$$

$$= \sigma \left( \begin{bmatrix} 0.02 \\ 0.005 \end{bmatrix} + \begin{bmatrix} 0.28 \\ 0.19 \end{bmatrix} \right)$$

$$= \sigma \left( \begin{bmatrix} 0.30 \\ 0.195 \end{bmatrix} \right) = \begin{bmatrix} 0.57 \\ 0.548 \end{bmatrix}$$

$$o_t = \sigma \left( \begin{bmatrix} 0.15 & 0.05 \\ 0.1 & -0.2 \end{bmatrix} \begin{bmatrix} 0.0 \\ 0.1 \end{bmatrix} + \begin{bmatrix} 0.3 & 0.25 \\ -0.2 & 0.2 \end{bmatrix} \begin{bmatrix} 0.5 \\ -0.1 \end{bmatrix} \right)$$

$$= \sigma \left( \begin{bmatrix} 0.005 \\ -0.02 \end{bmatrix} + \begin{bmatrix} 0.125 \\ -0.12 \end{bmatrix} \right)$$

$$= \sigma \left( \begin{bmatrix} 0.13 \\ -0.14 \end{bmatrix} \right) = \begin{bmatrix} 0.532 \\ 0.465 \end{bmatrix}$$

$$g_t = \tanh \left( \begin{bmatrix} 0.2 & 0.1 \\ -0.1 & 0.05 \end{bmatrix} \begin{bmatrix} 0.9 \\ 0.1 \end{bmatrix} + \begin{bmatrix} -0.5 & 0.4 \\ 0.2 & -0.3 \end{bmatrix} \begin{bmatrix} 0.5 \\ -0.1 \end{bmatrix} \right)$$

$$= \tanh \left( \begin{bmatrix} 0.01 \\ 0.005 \end{bmatrix} + \begin{bmatrix} -0.29 \\ 0.13 \end{bmatrix} \right)$$

$$= \tanh \left( \begin{bmatrix} -0.28 \\ 0.135 \end{bmatrix} \right) = \begin{bmatrix} -0.27 \\ 0.134 \end{bmatrix}$$

$$c_t = \begin{bmatrix} 0.44 \\ 0.532 \end{bmatrix} \odot \begin{bmatrix} 0.2 \\ -0.2 \end{bmatrix} + \begin{bmatrix} 0.57 \\ 0.548 \end{bmatrix} \begin{bmatrix} -0.27 \\ 0.134 \end{bmatrix}$$

$$= \begin{bmatrix} 0.088 \\ -0.1064 \end{bmatrix} + \begin{bmatrix} -0.1529 \\ 0.0734 \end{bmatrix} = \begin{bmatrix} -0.0659 \\ -0.033 \end{bmatrix}$$

Vanishing gradient

$$h_t = o_t \odot \tanh(c_t)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$

will make  $g_t \rightarrow 0$  as no  $\tanh$   
 $\therefore h_t$  will not go to 0

RNN

-/-

$$\frac{\partial h_t}{\partial h_{t-1}} = \tanh' \left( W_{hh} h_{t-1} + W_{wh} x_{t-1} \right)$$

c1

LSTM

$$\frac{\partial h_t}{\partial h_{t-1}} = o_t \tanh'(c_t) \left[ \frac{\partial c_t}{\partial h_{t-1}} + \frac{\partial o_t}{\partial h_{t-1}} \tanh(c_t) \right]$$

$$\frac{\partial c_t}{\partial h_{t-1}} = f_t \cdot \frac{\partial c_{t-1}}{\partial h_{t-1}} + C_{t-1}^{\text{exp}} \frac{\partial f_t}{\partial h_{t-1}} + g_t \frac{\partial i_t}{\partial h_{t-1}}$$

$$+ i_t \frac{\partial g_t}{\partial h_{t-1}}$$

25/9/25

$$\frac{\partial L}{\partial h_{t-1}} = \frac{\partial L}{\partial h_t} \frac{\partial h_t}{\partial h_{t-1}}$$

$$\frac{\partial h_t}{\partial h_{t-1}} = o_t (1 - \tanh^2(c_t)) \frac{\partial c_t}{\partial h_{t-1}} + \frac{\partial o_t}{\partial h_{t-1}} \tanh(c_t)$$

$$\frac{\partial c_t}{\partial h_{t-1}} = \frac{d}{dh_{t-1}} [f_t \odot c_{t-1} + i_t \odot g_t]$$

$$= f_t \frac{\partial c_{t-1}}{\partial h_{t-1}} + C_{t-1} \frac{\partial f_t}{\partial h_{t-1}} + g_t \frac{\partial i_t}{\partial h_{t-1}}$$

$$+ i_t \frac{\partial g_t}{\partial h_{t-1}}$$

put in  $\frac{\partial h_t}{\partial h_{t-1}}$

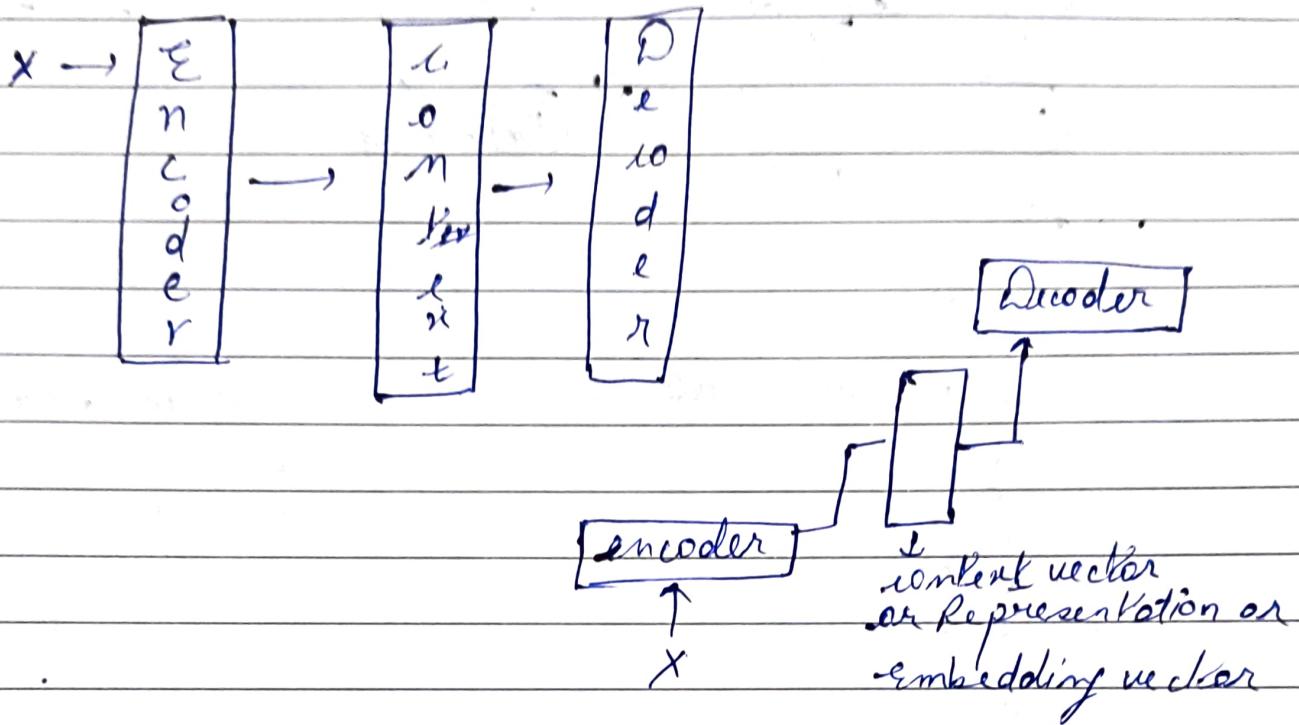
$\frac{\partial h_t}{\partial h_{t-1}}$

we can see there are many additive terms

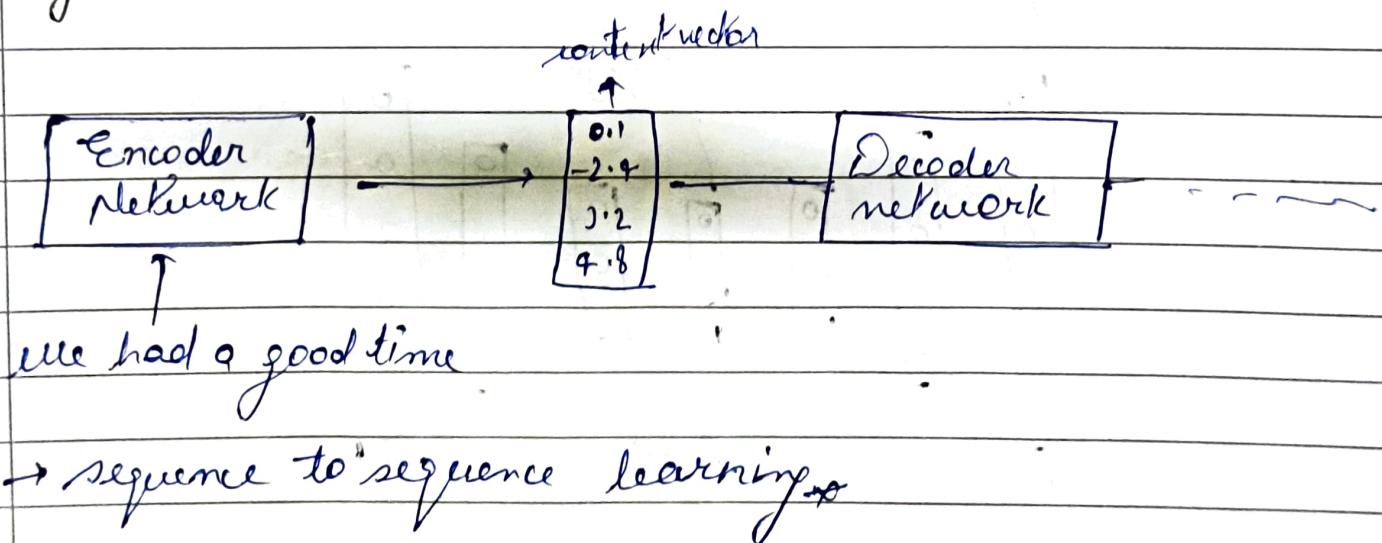
$\therefore$  vanishing gradient is avoided  $\Rightarrow$  theoretically can go to 0

New

## # Encoder Decoder Networks



## Eff Machine Translation Task



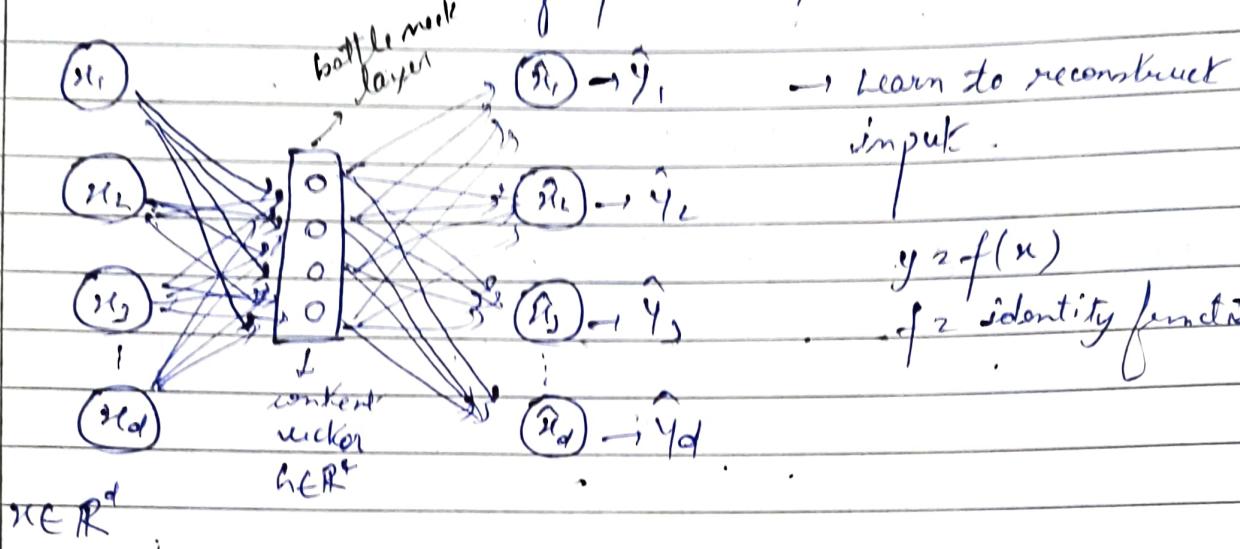
• sweet

• ML to simulate data

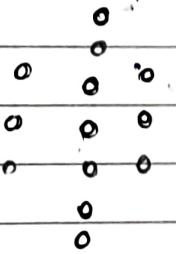
Datasets

- Linear
- PCA
  - t-SNE
- Non-linear

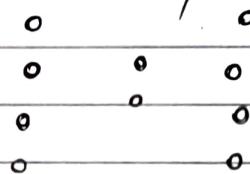
1) Auto-encoders  $\rightarrow$  to better representation of input (main purpose)



over complete autoencoder



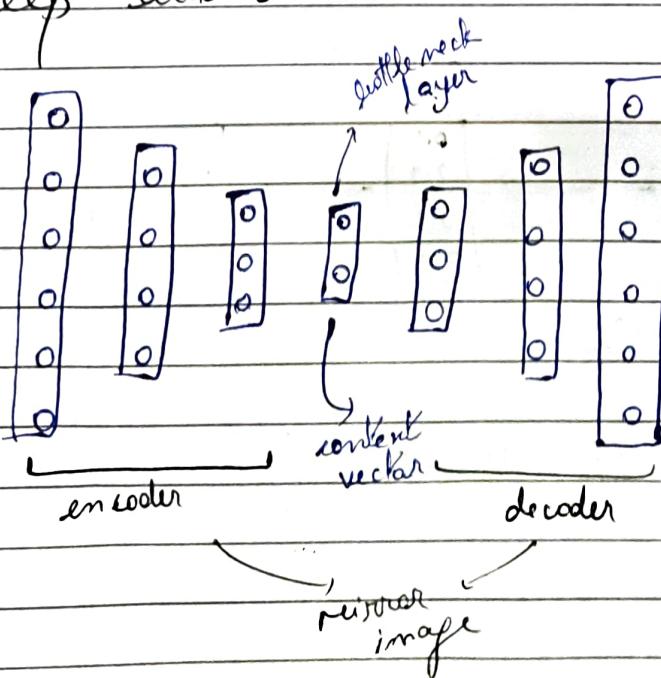
under complete autoencoder



✓

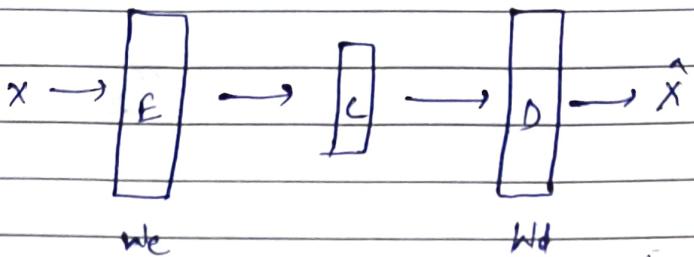
2)

Deep auto encoders



$$L = \frac{1}{2} \|x - \hat{x}\|_2^2$$

30/9/25



$\bullet x \in \mathbb{R}^2$ ,  $x = \begin{bmatrix} 2 \\ 0 \end{bmatrix}$

Autoencoder - Loss  $\rightarrow$  MSE

| = PCA if = Linear (Encoder)

| = Linear (Decoder)

$$h \in \mathbb{R}^1$$

content  
vector

1x2

2x1

$$\cdot \begin{bmatrix} 0.5 & -1.0 \end{bmatrix}_{1 \times 2}$$

$$= \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix}_{2 \times 1}$$

One forward pass and loss = ?

$$h_{\cancel{\text{out}}} = \begin{bmatrix} 2 \\ 0 \end{bmatrix}_{2 \times 1} \begin{bmatrix} 0.5 & -1.0 \end{bmatrix}_{1 \times 2} \overset{2}{=} \begin{bmatrix} 1 & -2 \\ 0 & 0 \end{bmatrix}_{2 \times 2}$$

$$\hat{x}_2$$

$$h_2 = \begin{bmatrix} 0.5 & -1.0 \end{bmatrix} \begin{bmatrix} 2 \\ 0 \end{bmatrix} = [1]$$

$$\hat{x}_2 = h \text{ Wd} = [1] \times \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix} = \begin{bmatrix} 1 \\ 0.5 \end{bmatrix}$$

$$\text{Loss} \Rightarrow \frac{1}{2} ((1-2)^2 + (0.5-0)^2)$$

$$= \frac{1}{2} (1 + 0.25) = \frac{1.25}{2}$$

$$= 0.625$$

7.4  
2.7

11

E.g.  $h_1 = [1, 0, 1]$ ,  $h_2 = [0, 1, 1]$ ,  $h_3 = [1, 1, 0]$   
 $s_{t-1} = [1, 0, 1]$ , score function = dot product.  
 $c_t = ?$

$$c_t = \sum_{j=1}^3 \alpha_{t,j} h_j$$

$$c_t = \alpha_{t,1} h_1 + \alpha_{t,2} h_2 + \alpha_{t,3} h_3$$

$$\alpha_{t,1} = \frac{\exp(\text{score}(s_{t-1}, h_1))}{\exp(\text{score}(s_{t-1}, h_1)) + \exp(\text{score}(s_{t-1}, h_2)) + \exp(\text{score}(s_{t-1}, h_3))}$$
$$= \frac{7.4}{12.8} = 0.58$$

$$\alpha_{t,2} = \frac{2.7}{12.8} = 0.2$$

$$\alpha_{t,3} = \frac{2.7}{12.8} = 0.2$$

$$c_t = 0.58 \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} + 0.2 \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix} + 0.2 \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} 0.58 \\ 0 \\ 0.58 \end{bmatrix} + \begin{bmatrix} 0 \\ 0.2 \\ 0.2 \end{bmatrix} + \begin{bmatrix} 0.2 \\ 0.2 \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} 0.78 \\ 0.4 \\ 0.78 \end{bmatrix}$$

9/10/25

Thinking

Thinking

$$\textcircled{1} \quad \begin{aligned} p_1 &= \text{Thinking} \\ k_1 &= \text{Thinking} \\ v_1 &= \text{Thinking} \end{aligned}$$

$$\textcircled{2} \quad \begin{aligned} p_1 &= \text{Thinking} \\ k_2 &= \text{Machines} \\ v &= \text{Thinking} \end{aligned}$$

Machines

$$\begin{aligned} q_1 &= \text{Machines} \\ k_1 &= \text{Thinking} \\ v &= \text{Machines} \end{aligned}$$

$$\begin{aligned} q_1 &= \text{Machines} \\ k_2 &= \text{Machines} \\ v &= \text{Machines} \end{aligned}$$

~~Thinking~~

$$i/p = \text{Playing Outside} \Rightarrow z_1, z_2$$

$$q_1 = [0.212 \quad 0.04 \quad 0.63 \quad 0.36]^T$$

$$k_1 = [0.31 \quad 0.84 \quad 0.963 \quad 0.57]^T$$

$$v_1 = [0.36 \quad 0.6] \quad 0.1 \quad 0.98]^T$$

$$q_2 = [0.1 \quad 0.14 \quad 0.86 \quad 0.77]^T$$

$$k_2 = [0.45 \quad 0.94 \quad 0.77 \quad 0.58]^T$$

$$v_2 = [0.31 \quad 0.36 \quad 0.19 \quad 0.72]^T$$

PlayingOutside

$$Z_1 = [0.1114 \quad 0.588 \quad 0.196 \quad 0.559]$$

$$[0.1116 \quad 0.6046 \quad 0.1431 \quad 0.5930]$$

11

Some

$$\frac{1}{2} + k_1$$

$$= 0.06572 + 0.0336 \\ + 0.60669 + 0.2052 \\ = 0.91121$$

$$\frac{1}{2} + k_2$$

$$= 0.0954 + 0.0376 \\ + 0.4599 + 0.2008$$

$$= \cancel{+0.51} \quad 0.8017$$

$$\text{Divided by } \sqrt{2} = \frac{0.91121}{2} \\ = \frac{\sqrt{9}}{\sqrt{2}} = 0.455605 \\ \approx 0.45$$

$$\frac{\cancel{0.51}}{2} \quad \frac{0.8017}{2}$$

$$\cancel{0.6255} = 0.90085 \\ \approx 0.9$$

$$\text{Softmax} \quad \frac{1.57}{1.57+1.49} = \frac{1.57}{3.06} \\ = 0.513$$

$$\frac{1.49}{1.57+1.49} = \frac{1.49}{3.06}$$

$$0.487$$

$$\text{Softmax} \times \text{Value} = 0.513 \cdot \begin{bmatrix} 0.31 \\ 0.32 \\ 0.19 \\ 0.72 \end{bmatrix} + 0.487 \cdot \begin{bmatrix} 0.31 \\ 0.32 \\ 0.19 \\ 0.72 \end{bmatrix}$$

$$= \begin{bmatrix} 0.18968 \\ 0.42579 \\ 0.0513 \\ 0.19994 \end{bmatrix} + \begin{bmatrix} 0.15037 \\ 0.17592 \\ 0.0925 \\ 0.35068 \end{bmatrix}$$

$$\begin{bmatrix} 0.33565 \\ 0.60111 \\ 0.14383 \\ 0.59558 \end{bmatrix}$$

# Flux Balance

— / —

Diffusion model

$$x_0 \in \mathbb{R}^d \sim p_{x_0}$$

ground truth dist.

Problem : sample from  $p_{x_0}$  (unknown)

Suppose  $p_\theta$  is the "model" distribution

$$D_{KL}(p_{x_0} || p_\theta) = \int_{x_0} p_{x_0} \log \frac{p_{x_0}}{p_\theta} dx_0$$

$$\theta^* = \operatorname{argmin}_\theta D_{KL}[p_{x_0} || p_\theta]$$

$$= \operatorname{argmin}_\theta \left[ \int p_{x_0} \log p_{x_0} - \int p_{x_0} \log p_\theta \right]$$

$$= \operatorname{argmax}_\theta \int_{x_0} p_{x_0} (\log p_\theta) dx_0$$

$$= \operatorname{argmax}_\theta \frac{1}{p_{x_0}} \int f(x) p_\theta(x) dx$$

$$x \sim p_x \int f(x) p_\theta(x) dx$$

$$= \frac{\mathbb{E}_{p_x} (f(x))}{p_x}$$

$$\frac{\mathbb{E}_{p_x} f(x)}{p_x} \approx \frac{1}{N} \sum_{i=1}^N f(x_i)$$

$$x_i \sim p_x$$

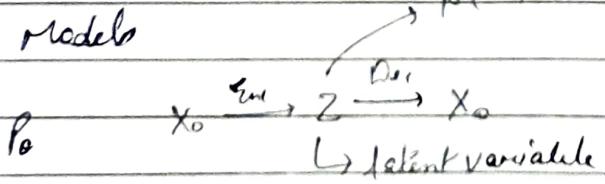
$$\ln \frac{\mathbb{E}_{p_x} (f(x))}{p_x}$$

$$\approx \operatorname{argmax}_\theta \frac{1}{N} \sum_{i=1}^N \log p_\theta(x_i)$$

$$x_i \sim p_{x_\theta}$$

MLE

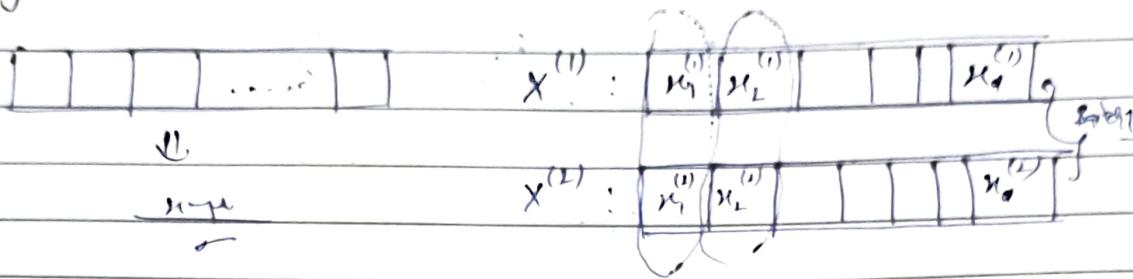
Diffusion models



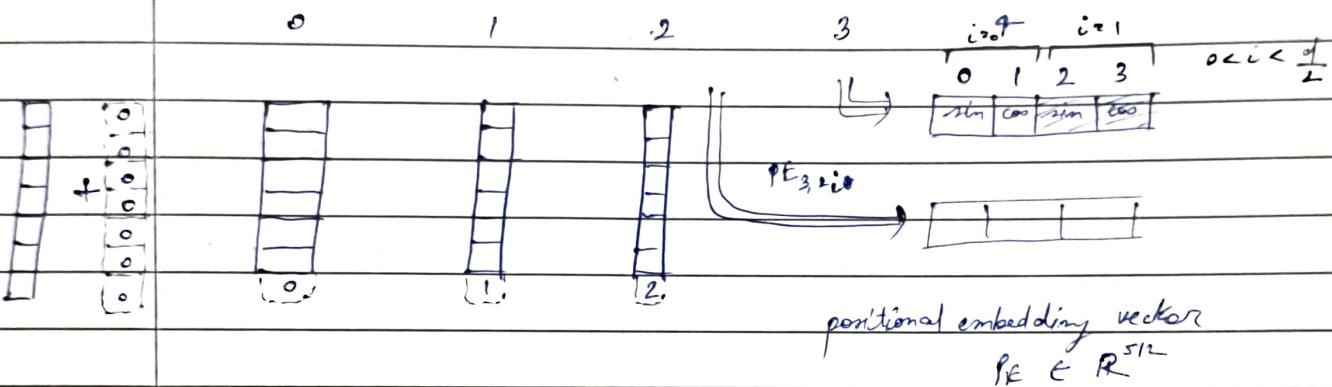
Latent model

13/10/25

# Layer Normalization  $\Rightarrow$  Batch Normalization



Key components of Transformers architecture



# Positional Encoding

$$PE_{pos,2i} = \sin\left(\frac{pos}{10000^{(2i)/emb}}\right)$$

$$PE_{pos,2i+1} = \cos\left(\frac{pos}{10000^{(2i)/emb}}\right)$$

e.g.

I am a robot

$PE: d_{24} \rightarrow \mathbb{R}^d$	1	→	0	→	$P_{00} P_{01} P_{02} P_{03}$
am	→	1	→	$P_{10} P_{11} P_{12} P_{13}$	
a	→	2			
robot	→	3	→	$P_{30} P_{31} P_{32} P_{33}$	

— / —

$$P_{00} = \sin\left(\frac{0}{100 \frac{\pi}{180}}\right)$$

$$\therefore \sin(0) = 0$$

$$P_{01} = \cos\left(\frac{0}{100 \frac{\pi}{180}}\right) = \cos(0)$$

$$\therefore \cos(0) = 1$$

$$P_{02} = \sin\left(\frac{0}{100 \frac{\pi}{180}}\right) = \sin(0) = 0$$

$$P_{03} = \cos(0) = 1$$

$$P_{10} = \sin\left(\frac{1}{100 \frac{\pi}{180}}\right) = \sin(1) \\ \approx 0.84$$

$$P_{11} = \cos\left(\frac{1}{100 \frac{\pi}{180}}\right) = \cos(1) \\ \approx 0.54$$

$$P_{12} = \sin\left(\frac{1}{100 \frac{\pi}{180}}\right) = \sin\left(\frac{1}{10}\right) = \sin\left(\frac{1}{10}\right) \\ \approx 0.0998$$

$$P_{13} = \cos\left(\frac{1}{100 \frac{\pi}{180}}\right) = \cos\left(\frac{1}{10}\right) = 0.995$$

similarly for rest

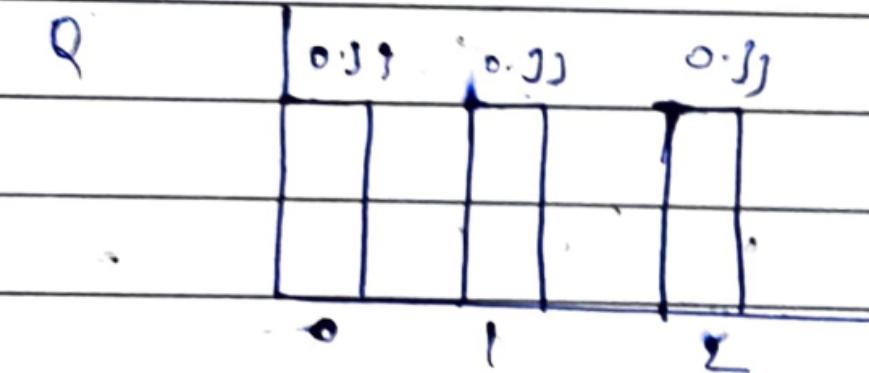
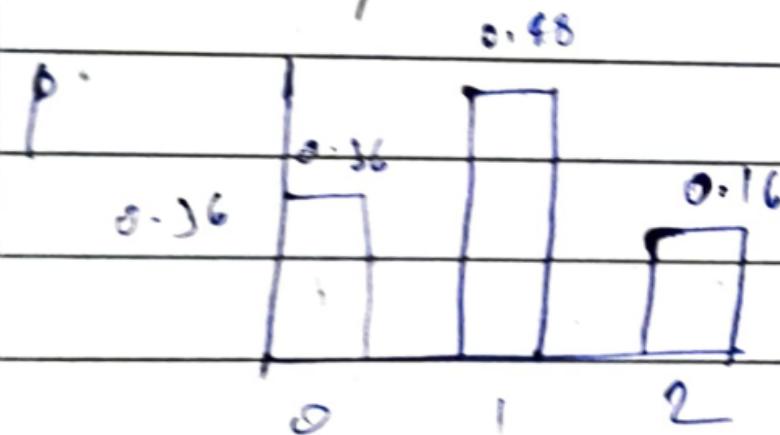
#  
P<sub>000</sub> of del product =

P<sub>001</sub>

P<sub>0000</sub> of del product =   
P<sub>1000</sub>

input:  [i]  
[j] + [P<sub>00</sub>]  
P<sub>00</sub>:  (in the paper)

# KL example



Binomial distante

Uniform dist

Distribution<sup>x</sup>

$x$	0	1	2
$P(x)$	$\frac{9}{25}$	$\frac{12}{25}$	$\frac{4}{25}$
$Q(x)$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$

$$\begin{aligned}
 D_{KL}(P||Q) &= \sum_{x \in X} P(x=x) \ln \frac{P(x=x)}{Q(x=x)} \\
 &= \frac{9}{25} \ln \left( \frac{9/25}{1/3} \right) + \frac{12}{25} \ln \left( \frac{12/25}{1/3} \right) \\
 &\quad + \frac{4}{25} \ln \left( \frac{4/25}{1/3} \right)
 \end{aligned}$$

$$\approx 0.0852936$$

$$D_K(Q||P) = \frac{1}{3} \ln \left( \frac{1/3}{9/25} \right) + \frac{1}{3} \ln \left( \frac{1/3}{12/25} \right)$$

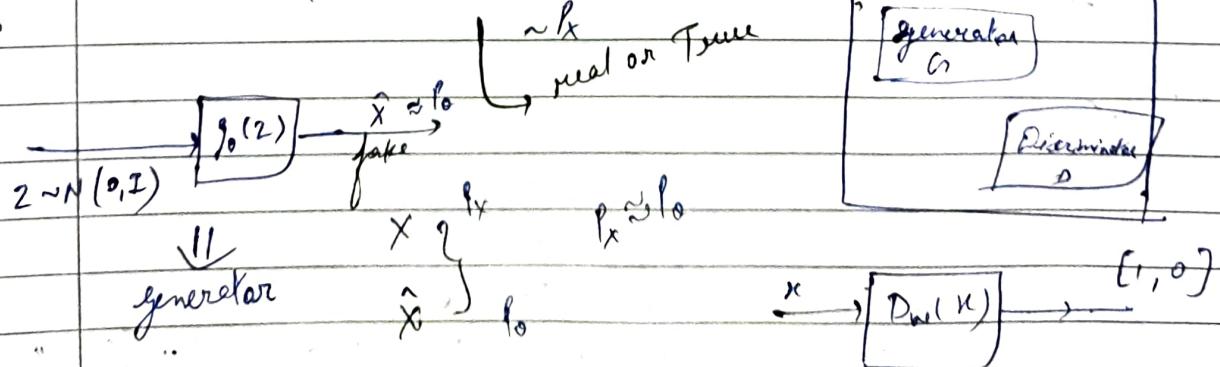
$$+ \frac{1}{3} \ln \left( \frac{1/3}{4/25} \right)$$

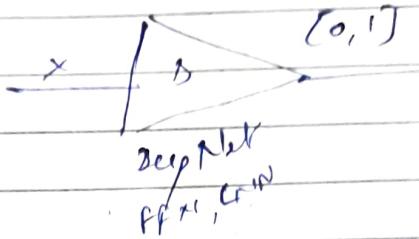
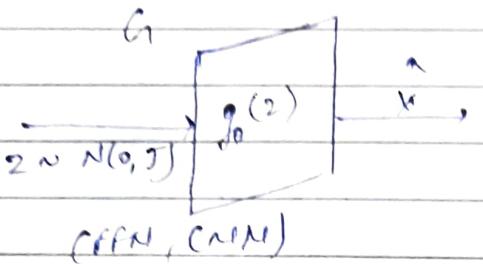
$$\approx 0.039755$$

# Variational Divergence minimization  
 $f(u) = u \log u - (u+1) \log(u+1)$

Generative Adversarial Network (GAN)

$$X = \{x_1, x_2, x_3\} \sim \text{iid}$$





$$J_{GAN}(\theta, w) = \mathbb{E}_{x \sim P_r} \log D_w(x) + \mathbb{E}_{\hat{x} \sim P_\theta} \log (1 - D_w(\hat{x}))$$

$$\hat{x} = f_\theta(z)$$

Input:  $\{x_1, x_2, \dots, x_n\} \sim \text{iid } \sim P_x$

$$B_1 : x_1, x_2, \dots, x_{B_1} \sim P_x \quad \min_{\theta} \max_{w} \left[ \frac{1}{B_1} \sum_{i=1}^{B_1} \log D_w(x_i) + \frac{1}{B_2} \sum_{j=1}^{B_2} \log (1 - D_w(f_\theta(z_j))) \right]$$

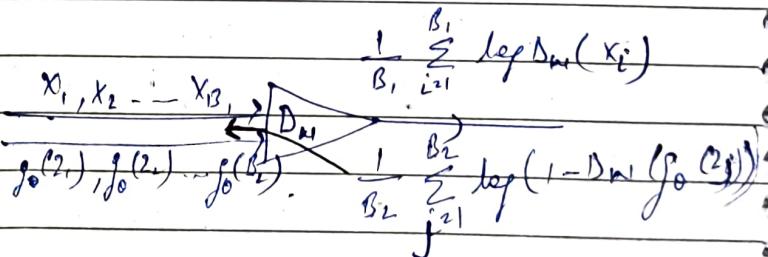
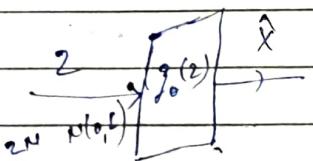
$$B_2 : \hat{x}_1, \hat{x}_2, \dots, \hat{x}_{B_2} \sim P_\theta$$

$$J_{GAN}(\theta, w) = \arg\max_w \left[ \frac{1}{B_1} \sum_{i=1}^{B_1} \log D_w(x_i) + \frac{1}{B_2} \sum_{j=1}^{B_2} \log (1 - D_w(f_\theta(z_j))) \right]$$

*edit*

To Discriminator

Keep  $\theta$  as constant



$$\nabla_w J = w^T w + \lambda D_w$$

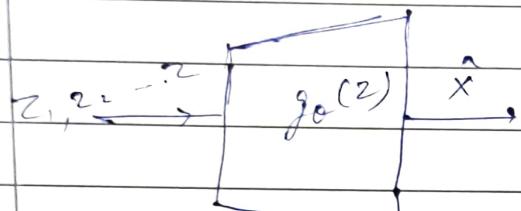
$$\theta^* = \arg\min_{\theta} J_{GAN}(\theta, w)$$

$$\approx \arg\min_{\theta} \left[ \frac{1}{B_1} \sum_{i=1}^{B_1} \log D_w(x_i) + \frac{1}{B_2} \sum_{j=1}^{B_2} \log (1 - D_w(f_\theta(z_j))) \right]$$

$$\theta^* \approx \arg \min_{\theta} \left[ \frac{1}{B_2} \sum_{j=1}^{B_2} \log (1 - D_w(g_\theta(z_j))) \right]$$

$$\theta : \theta \leftarrow \alpha_L \nabla_{\theta} J(\theta, w)$$

To train generator



$$J_{GAN} = \left[ \frac{1}{B_2} \sum_{j=1}^{B_2} \log (1 - D_w(g_\theta(z_j))) \right]$$

$$\theta^{t+1} = \theta^t - \alpha_L \nabla_{\theta} J$$