

Convergence - Scaling Axis

Quadratic eq =  $f(x) + (x-x_0) f'(x_0)$

original eq =  $E = \frac{1}{2} \vec{\omega}^T A \vec{\omega} + b^T \vec{\omega} + c$

let  $S = A^{0.5}$ ,  $\vec{\omega} = S\omega$

then  $E = \frac{1}{2} \vec{\omega}^T S^T S \omega + b^T S \omega + c$

equating coefficients,  $\boxed{S^T S = A}$

$$b^T S = b^T$$

$$\therefore \boxed{\hat{b} = A^{-0.5} b}$$

$$E = \frac{1}{2} \vec{\omega}^T \vec{\omega} + \hat{b}^T \vec{\omega} + c$$

$$\vec{\omega} = A^{0.5} \omega$$

$$\omega = A^{-1/2} \vec{\omega}$$

$$\vec{\omega} = S\omega$$

$$\vec{\omega} = A^{1/2} \omega$$

then,  $\frac{\vec{\omega}}{A^{1/2}} = \omega$

$$\boxed{\omega = A^{-1/2} \vec{\omega}}$$

$$A = A^{-1/2} A$$

Using ① in  $E = \frac{1}{2} \vec{\omega}^T A \vec{\omega} + b^T \vec{\omega} + c$

$$E = \frac{1}{2} \vec{\omega}^T \left[ \underbrace{A^{-1/2} A}_{I} \right] A \cdot A^{-1/2} \vec{\omega} + b^T \vec{A}^{-1/2} \vec{\omega} + c$$

$$\boxed{E = \frac{1}{2} \vec{\omega}^T \vec{\omega} + b^T \vec{\omega} + c}$$

$$GD: \hat{w}^{k+1} = \hat{w}^k - \eta \nabla_{\hat{w}} E (\hat{w}^k)^T$$

$$\hat{w} = A^{1/2} w$$

$$\eta^{1/2} \hat{w}^{k+1} = A^{1/2} w^k - \eta A^{1/2} \nabla_w E$$

where,  $\hat{w}^{k+1} = w^k - \eta \nabla_w E$

$$A = \nabla_w^2 E$$

$$E = \frac{1}{2} w^T A w + w^T b + c$$

Newton Method.

Forward Pass

$$y_1 = f_1(w_1 y_0 + b_1), \quad z_1 = w_1 y_0 + b_1$$

$y = \text{activation}(z) = f(z)$

$$y_2 = f_2(z_2)$$

$$y_N = f_N(w_N t_{N-1} + b_N) = f_N(w_N f_{N-1}(w_{N-1} t_{N-2} + b_{N-1}) + \dots + b_N)$$

back prop:-

$$\frac{\partial}{\partial y_i} \text{Div} = \left( \frac{\partial a}{\partial y_i} \right) \text{Dir}(y/a)$$

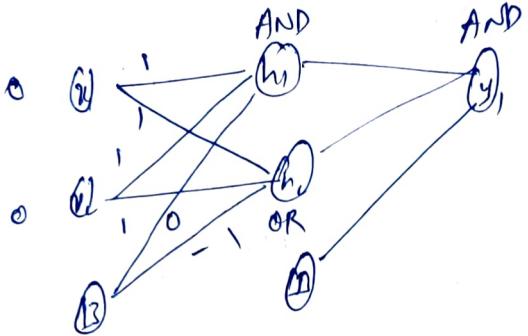
for  $k = N-1, \dots, 0$   
for  $i = 1$ : layer width

$$\frac{\partial}{\partial y_i} \text{Div} = \sum_j w_{ij}^{k+1} \frac{\partial \text{Dirv}}{\partial z_i^{k+1}}$$

$$\frac{\partial}{\partial z_i} \text{Div} = f'(z_i) \cdot \frac{\partial \text{Div}}{\partial y_i}$$

$$\text{and } \frac{\partial \text{Div}}{\partial z_i} = f'_k(z_i) \frac{\partial \text{Div}}{\partial y_i}$$

for all,  $\frac{\partial \text{Div}}{\partial w_{ij}^{k+1}} = y_i^k \frac{\partial \text{Div}}{\partial z_j^{k+1}}$



x	y	<u>XOR</u>
0	0	0
0	1	1
1	0	1
1	1	0

$\text{ReLU}[\omega_1 x_1 + \omega_2 x_2 + b]$  for each neuron

$$h_1 = 0 \cdot 1 + 0 \cdot 1 + 0 = 0$$

$$\text{ReLU}(h_1) = \text{ReLU}(0) = 0$$

Ex 4

Empirical Risk =

$$\text{Loss } \omega = \frac{1}{T} \sum \text{div}(f(x_i, \omega), d_i)$$

emp error - arg min loss  $\omega$

$$\text{div}(f(x_i, \omega), g^{(u)}) = \begin{cases} 0 & \text{if } f(x_i, \omega) = g^{(u)} \\ > 0 & \text{if } f(x_i, \omega) \neq g^{(u)} \end{cases}$$

$$\text{div}(y) = \text{div}(y_i, d) = \frac{1}{2} \|y - d\|^2$$

$$= \frac{1}{2} \sum (y_i - d_i)^2$$

Squared euclidean distance  
between true and desired output

$$\frac{d}{dy_i} \text{div}(y, d) = (y_i - d_i)$$

For binary classification = (as it will be scalar)

$$\frac{d}{dy} \text{Div}(y, d) = \begin{cases} -\frac{1}{y} & \text{if } d=1 \\ \frac{1}{1-y} & \text{if } d=0 \end{cases}$$

$$\text{Div}(y, d) = -d \log y - (1-d) \log(1-y)$$

minimum when  $d=y$

know/see 7

Finding best step size

[34:30] fine story

$$\text{Any quad. eq: } E(\omega) = E(\omega^k) + \underbrace{E'(\omega^k) \cdot (\omega - \omega^k)}_{\text{01}} + \frac{1}{2} E''(\omega^k) \cdot (\omega - \omega^k)^2$$

when derivative: (constant)

$$\frac{d E(\omega)}{d \omega} = 0 + E'(\omega^k) \cdot 1 + \frac{1}{2} E''(\omega^k) (\omega - \omega^k) = 0$$

$$\Rightarrow \boxed{\frac{-E'(\omega^k)}{E''(\omega^k)} = \omega - \omega^k}$$

Compare this with update rule:

$$\omega^{k+1} = \omega^k - \eta \frac{d E(\omega^k)}{d \omega}$$

$$\boxed{\eta_{\text{opt}} = \frac{1}{E''(\omega^k)} = \alpha^{-1}} = \left( \frac{d^2 E(\omega^k)}{d \omega^2} \right)^{-1}$$

if  $\eta < E''(\omega^k)^{-1}$   $\Rightarrow$  goes before the minima.

if  $\eta > E''(\omega^k)^{-1}$   $\rightarrow$  goes above the minima

$\rightarrow \eta > 2\eta_{\text{opt}}$   $\rightarrow$  divergence

Multivariate Quadratic with Diagonal A

$$E = \frac{1}{2} \omega^T A \omega + \omega^T b + c = \frac{1}{2} \sum [a_{ij} \omega_i^2 + b_j \omega_i] + c$$

$$\eta < 2 \min \eta_{i, \text{opt}}$$

- issue 1
- but makes slow learning
  - will oscillate in all dim,  $\eta_{i, \text{opt}} < \eta < 2\eta_{i, \text{opt}}$
  - there estimation are for 2 dimension.  
If there are  $n$  dimension the computing and derivatives are complicated

How to fix these problems?

- Scale the variable one - step thing
- newton method does this by multiplying space by [Hessian]

issue 2 :- learning rate

- for complex models - loss function is not convex

Sol: Start with large  $\eta$ ;  $\eta > 2$

then gradually reduce it with iterations.

Other Sol - train with fixed  $\eta$  until loss is stagnant

iterate { -  $\eta = \alpha \eta$  ;  $\alpha < 1$

## RProp

- this is based on the sign of the derivative
  - . if sign = +ve  $\rightarrow$  move forward
  - . if sign = -ve  $\rightarrow$  move backward
- Simple first order algo.
  - much efficient than gradient descent
  - sometimes better than some good 2nd order algorithms
- only makes minimal assumptions about loss function

## Quickprop

- employs newton update with empirically derived derivative
  - prone to instability for non convex objective fun.
  - still one of the fastest training algo for many problems.
  - keep track of oscillations
- Momentum method.
- maintain running avg.
  - If Smooth - average will have large value
  - If Swings - avg has +ve and -ve values that get ~~cancel~~ cancelled out
  - Update avg.

Unit 8  
Lec 8

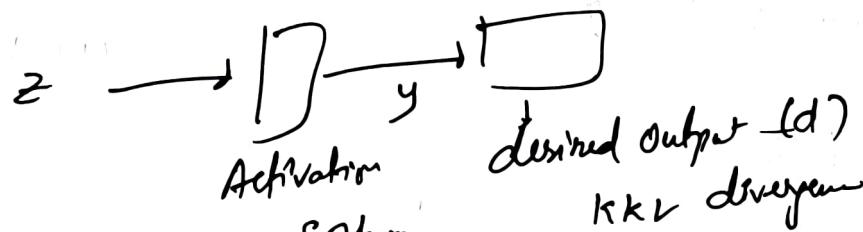
What is a good divergence?

- smooth and less local optima
- low slopes far from optima = bad
- high slopes near optima = bad.



$L_2$  div gradient

$$\nabla \frac{1}{2} \|y - d\|^2 = (y - d)^T$$



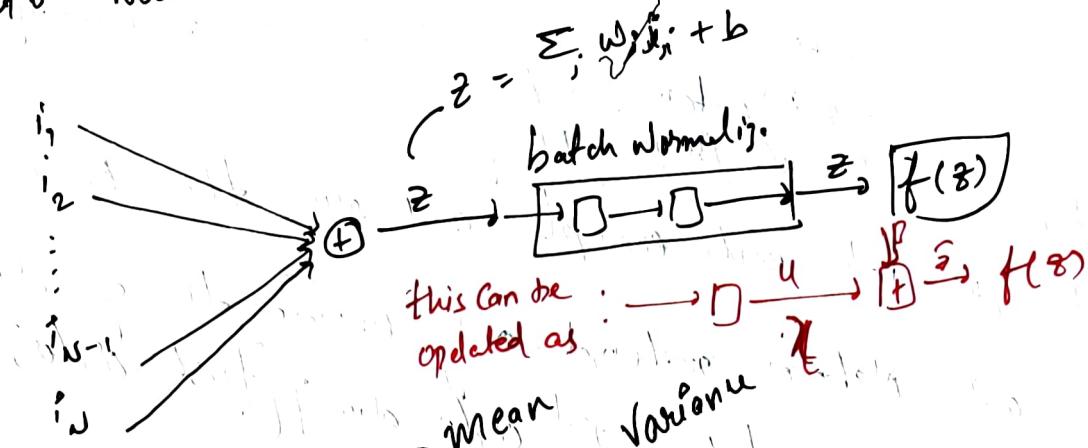
classification : Softmax

Batch Normalization

has 2 steps:

1. move the entire minibatch to origin and normalize  
- has  $\mu$  mean
2. move " " " " " a new location

To move to new location:-



$$\text{Step 1: } \mu_B = \frac{1}{B} \sum_{i=1}^B z_i$$

$$\text{Step 2: } \sigma_B^2 = \frac{1}{B} \sum_{i=1}^B (z_i - \mu_B)^2$$

$$\text{Step 3: } u_i = \frac{z_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

$$\text{Step 4: } z_i = \gamma u_i + \beta$$

$B$  = minibatch size  
 $\mu_B$  = minibatch mean  
 $\sigma_B^2$  = standard deviation  
 $\epsilon$  = small value

Subtract mean from every instance of the minibatch.  
 make sure, collection of minibatches have avg of 0.

new location beta

Scaling to the right/left size and do new location  $\beta$

gamma

→ this [Step 1, 2, 3, 4] will happen for every neuron.  
 [mean, variance, centering and Normalizing]  
 [Scaling + shifting] (also like bias)

$\gamma, \beta$ : Unit specific for every neuron.

Step 1:- It is like having everything at center, at origin when we computing  $\beta$ : it is like Bias of neuron

class 8

21/8/23

Thursday

## Deep Learning

11

### notation

- minibatch  $\rightarrow m_b$
- Back Prop  $\rightarrow BP$

[Sug solns on Monday]

momentum

In BP, we take only examples and even the BP algo.  
momentum and Nesterov method only consider the gradient  
(other also take 2nd order deriv) [first order]

RMS Prop not taking momentum of gradient

compute square of derivative (not 2nd order)  
each individual element of gradient

$$(\partial_w D)^2 = \partial_w^2 D$$

$$= E(\partial_w^2 D)_k = \gamma E(\partial_w^2 D)_{k-1} + (1-\gamma)(\partial_w^2 D)_k$$

$$\Rightarrow w_{k+1} = w_k - \frac{\eta}{\sqrt{E(\partial_w^2 D)_k + \epsilon}} \partial_w D$$

running avg of square of gradient

ADAM : RMS Prop with momentum

Replace gradient with momentum term.

$$m_k = \delta m_{k-1} + (1-\delta)(\partial_w D)_k$$

$$v_k = \gamma v_{k-1} + (1-\gamma)(\partial_w^2 D)_k$$

$$w_{k+1} = w_k - \frac{\eta}{\sqrt{v_k + \epsilon}} m_k$$

$$\hat{v}_k = \frac{v_k}{1-\gamma}$$

$$m_k = \frac{m_k}{1-\delta}$$

In early iteration,  $\delta$  and  $\gamma$   
won't get normalized  
so we do that in last  
layer

optimization done

- Bias-Variance Tradeoff
  - Regularization
  - Dropout, Early stopping
- } Useful but not much are Using

→ model to perform good on unseen data.

- this error can be decomposed as 3 form:-

1. Bias ↑  $B \propto V^2$
2. Variance ↓
3. Noise

### Polynomial (P)

- As P or the function become complex and tries to Overfit the data.
- Test error decrease earlier and as P increases, the test error starts increasing.



Too much complexity of model is called ~~test data~~ Polynomial, P ~~Bias~~.

Variance: In prediction model complexity will be high if model middle part has perfect balance between Bias and Variance.

Underfit: Bias

Overfit: Variance

- High bias & cause high test errors

function we learned is this:- model built from dataset-

$$MSE = E_{D(x,y,\eta)} [(y - f(x:D))^2]$$

with constant  $|D| = N$

effectiveness of model

using mean Squared Error

$E$  = expected value.

\* board

$$NSE = E_{x,y} [(y(x) - t)^2]$$

$$= \int (f(u) - t)^2 P(u, t) du dy$$

$$NSE \rightarrow \int (y(u) - t)^2 P(u, t) du dt$$

need to minimize this then take the gradient, grad it to 0

$P(u, t)$  = joint distribution on  $u$  and  $t$

derivation taken from:- Book

Pattern Recognition and Machine Learning  
by Christopher Bishop

read chapter 4

$t$  = some  $t$  in  $E$  cases  
(Random Variable  $N(0, \dots)$ )

Optimal =  $y^*(x) \rightarrow$  [minimizer that minimises NSE]

$$y^*(x) = \int \frac{t P(x, t)}{P(x)} dt \rightarrow \int t P(t|x) dt$$

conditional prob

$E(t|u)$  = Conditional Expectation

$$y^*(u) = E_{t|u}(t)$$

$$(y(u) - t)^2 = \{ y(u) - E[t|u] + E[t|u] - t^* \}^2$$

$$= [y(u) - E(t|x)]^2 + 2[y(u) - E(t|u)][E(t|u) - t^*]^2 + (E(t|u) - t^*)^2$$

Noise Variance  
\* this is the minimum error, no matter how good model we take. (not in our hands)

$$\text{2nd Term} := E_{u,t} [(y(u) - E(t|u))(E(t|u) - t^*)]$$

$$= E_u \left[ E_{t|u} \left[ \{y(u) - E(t|u)\} \{E(t|u) - t^*\} \right] | x \right]$$

does not depend on  $t$   
So it will come out

$$= E_x \left[ y(u) - E(f|u) \right] \underbrace{E_{f|u} \left[ E(f|u) - f \right] u}_{\text{this will be } 0}$$

$$= E_x \left[ y(u) - E(f|u) \right] \underbrace{\left[ E(f|u) - E(f|u) \right]}_0$$

$$= 0$$

hence, the 2nd term is 0 (vanishing)

hence, MSE will depend on 2 terms.

$$E(L) = \int [y(u) - E(f|u)]^2 P(u) du + \int \text{Var}(f|u) P(u) du$$

$$\text{MSE} = \underbrace{E \left[ (y(u) - E(f|u))^2 \right]}_{\text{1st term}} \quad \text{Per hypothesis}$$

$$= E_u E_D \left[ (y(u:D) - E(f|u))^2 \right]$$

$$(y(u) - E(f|u))^2 = E_D \left[ (y(u:D) - E(f|u))^2 \right]$$

$$\text{MSE} = \underbrace{E_D \left[ y(u:D) - \mu_D \right]^2}_{\text{Bias}} - \underbrace{E_D \left[ (y(u:D) - E_D[y(u:D)])^2 \right]}_{\text{Variance}}$$

(Bias)<sup>2</sup>  
how much your prediction

- differ from the optimal.

[depends on choice of complexity]

y = linear function.

Variance

L<sub>2</sub> regularization

$$L(\omega) = L(\omega) - \frac{1}{2} \alpha \|\omega\|^2$$

MSE

If  $\alpha \gg 1$ , then the algo focuses on 2nd part

If  $\alpha \rightarrow \infty \rightarrow$  high regularization

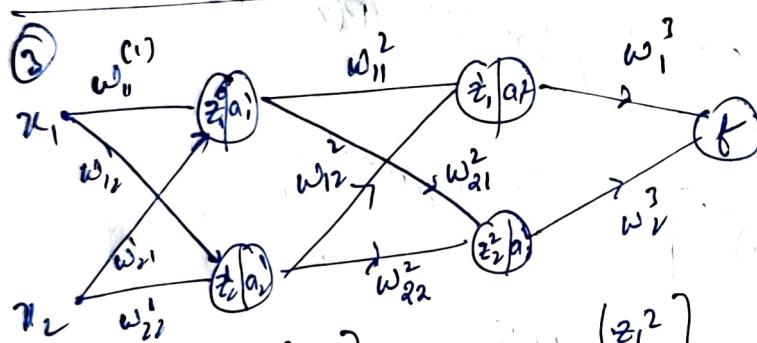
① tanh is odd fn

$$f(x) = \sum_{j=1}^{m_d} a_j \tanh(\omega_j^T x + b_j)$$

$$\bar{\omega}_1 = \begin{pmatrix} \omega_{11} \\ \omega_{12} \\ \vdots \\ \omega_{1n} \end{pmatrix}$$

$$= \sum -a_j \tanh(-\omega_j^T x - b_j)$$

$$= f(x)$$



$$z^1 = \begin{pmatrix} z_1^1 \\ z_2^1 \end{pmatrix}$$

$$z^2 = \begin{pmatrix} z_1^2 \\ z_2^2 \end{pmatrix}$$

$$a^1 = \begin{pmatrix} a_1^1 \\ a_2^1 \end{pmatrix}$$

$$a^2 = \begin{pmatrix} a_1^2 \\ a_2^2 \end{pmatrix}$$

$$f = \omega_1^T a_1^2 + \omega_2^T a_2^2$$

$$f_1 = \frac{\partial f}{\partial z_1^2}$$

activation for

$$f = g(z_1^2) \omega_1^T + g(z_2^2) \omega_2^T$$

$$\frac{\partial f}{\partial z_1^2} = \omega_1^T g'(z_1^2) (1)$$

$$\frac{\partial f}{\partial z_1^2} = \omega_1^T g(z_1^2) (1 - g(z_1^2))$$

$$\delta_2 = \frac{\partial f}{\partial z^2} = \frac{\partial f}{\partial A^2} \cdot \frac{\partial A^2}{\partial z^2}$$

$$\Rightarrow \begin{bmatrix} \frac{\partial f}{\partial z_{12}} & \frac{\partial f}{\partial z_{22}} \end{bmatrix} \rightarrow \text{Row Vector.}$$

$$\delta_2 = \frac{\partial f}{\partial z^2} = \frac{\partial f}{\partial z^2} \cdot \frac{\partial z^2}{\partial z^1} \Rightarrow \frac{\partial f}{\partial z^2} \cdot \frac{\partial z^2}{\partial A^1} \cdot \frac{\partial A^1}{\partial z^1}$$

$$= \delta_2 \cdot \frac{\partial z^2}{\partial z^1} \Rightarrow \delta_2 \cdot \frac{\partial [A^1 \omega^{2T}]}{\partial z^1}$$

$$\text{activation fn} \quad \delta_2 \cdot (\omega^{2T}) \cdot \frac{\partial (A^1)}{\partial z^1}$$

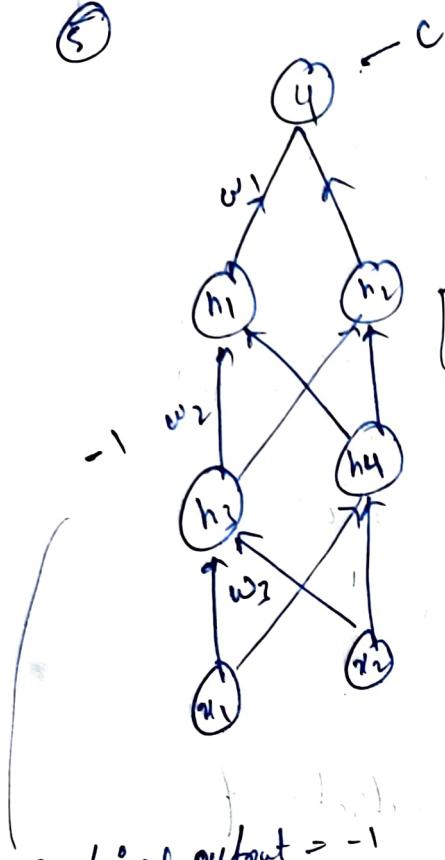
$$\Rightarrow \delta_2 \cdot \omega^{2T} \cdot [g'(z^1)]$$

$$z_1^2 = w_{11}^2 a_1^1 + w_{12}^2 a_2^1$$

$$z_2^2 = w_{21}^2 a_1^1 + w_{22}^2 a_2^1$$

$$\frac{\partial z^2}{\partial A^1} = \begin{bmatrix} \frac{\partial z_1^2}{\partial a_1^1} & \frac{\partial z_1^2}{\partial a_2^1} \\ \frac{\partial z_2^2}{\partial a_1^1} & \frac{\partial z_2^2}{\partial a_2^1} \end{bmatrix}$$

$$= \begin{bmatrix} w_{11}^2 & w_{12}^2 \\ w_{21}^2 & w_{22}^2 \end{bmatrix}$$



$h_i = \max(0, z_i)$   
 $z_i = \text{activation value}$

(a) when of this 2 will be 0

$$(A) \frac{\partial C}{\partial w_1}, \quad (B) \frac{\partial C}{\partial w_2} \\ = 0$$

$$(A) \frac{\partial C}{\partial w_1} = \frac{\partial C}{\partial y} \cdot \frac{\partial y}{\partial w_1}$$

$$(B) \frac{\partial C}{\partial w_2} \neq 0 ? \text{ How } \rightarrow h_2 w_2 + h_4 w_3$$

[need to find the gradient]

combined output = -1

for this  $z_1 = -1$

$$z_1 = -1$$

then  $h_1 = (0, -1)^T \max$

$$h_1 = 0$$

$$z_1 = w_2, h_3 = -1$$

$$\frac{\partial h_1}{\partial z_1} = 0 \quad \text{if } z_1 < 0 \\ \neq 1 \quad \text{if } z_1 > 0$$

① Implement max function using ReLU method. (Had to use ReLU)

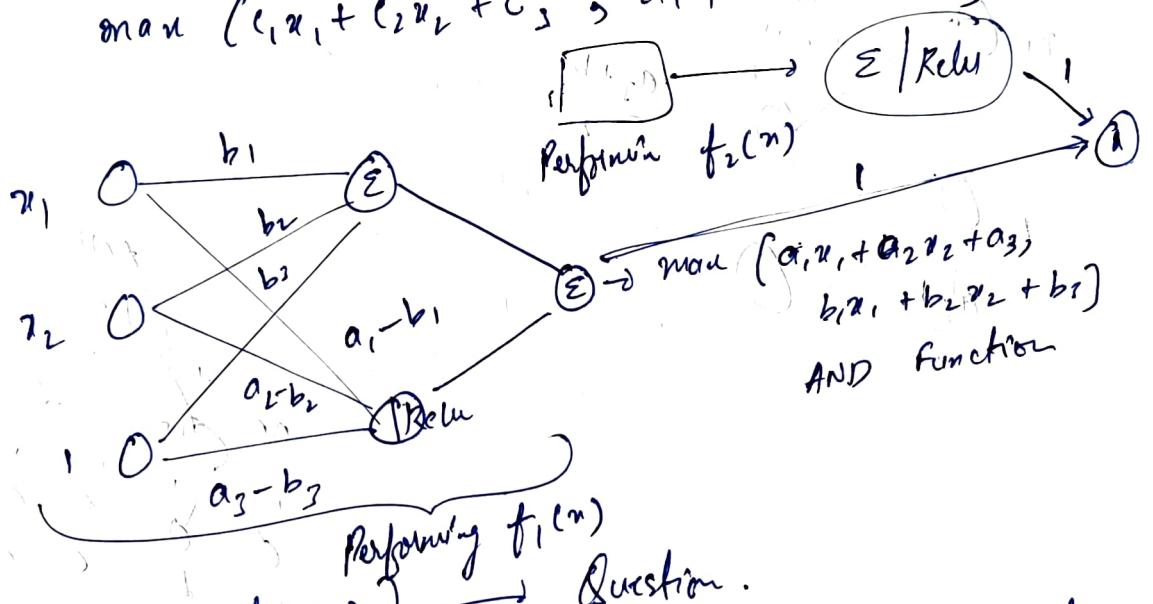
$\max(a, b)$

$$= \max(a, b) + b - b \rightarrow \begin{cases} \text{if } a > b \rightarrow a \\ \text{if } a < b \rightarrow b \end{cases}$$

$$= \max(a - b, 0) + b \rightarrow \begin{cases} a & \text{if } a - b > 0 \\ b & \text{if } a - b \leq 0 \end{cases}$$

$$\min \left[ \max(a_1x_1 + a_2x_2 + a_3, b_1x_1 + b_2x_2 + b_3), \right.$$

$$\max(c_1x_1 + c_2x_2 + c_3, d_1x_1 + d_2x_2 + d_3) \right]$$



$$\min(f_1(x), f_2(x)) \rightarrow \text{Question.}$$

decision boundary - this is the AND function

Q. we need to find out if this  $\min(f_1(x), f_2(x))$  is  $> 0$  or not.

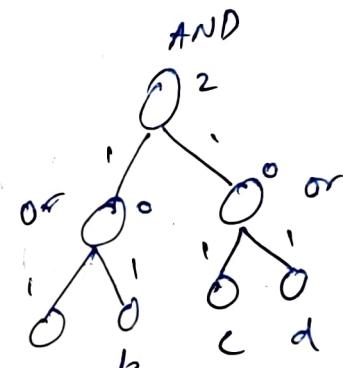
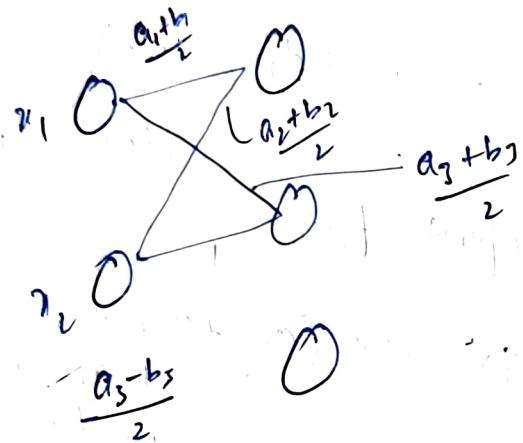
$$\max(a, b) = \frac{a+b+|a-b|}{2}$$

11.

→ Implement max(a<sub>1</sub>, b) Using 2nd method

$$\max(a_1x_1 + a_2x_2 + a_3, b_1x_1 + b_2x_2 + b_3)$$

$$= \frac{1}{2} \left[ (a_1 + b_1)x_1 + (a_2 + b_2)x_2 + \left( \frac{(a_1 + b_1) + (a_2 + b_2)}{a_3 + b_3} \right) + \left| \begin{array}{l} (a_1 - a_2)x_1 + (b_1 - b_2)x_2 \\ + (c_1 - c_2) \end{array} \right| \right]$$



$$\begin{aligned} a &= a_1x_1 + a_2x_2 + a_3 \\ b &= b_1x_1 + b_2x_2 + b_3 \end{aligned}$$

$$\begin{aligned} \text{generally, } \max(a_1, b) &= -\min(-a_1, -b) \\ &= -\max(-a_1, -b) + b - b \\ &= [\max(-a_1, -b) + b] + b \\ &= -\max(b - a_1, 0) + b \end{aligned}$$

## $L_2$ Regularization

$$\tilde{L}(w) = L(w) + \frac{1}{2} \alpha \|w\|^2$$

$$\nabla \tilde{L}(w) = H(w - w^*) \quad w^* = \text{minimum}$$

$H$  is symmetric & the semi definite

$$\nabla \tilde{L}(w) =$$

$$= \nabla L(w) + \alpha w$$

$$= H(w - w^*) + \alpha w$$

$$= Q \Lambda Q^T$$

where  $Q Q^T = Q^T Q = I$

$$\Lambda = \text{Diagonal matrix with eigen values of } H$$

$$\text{If } \nabla \tilde{L}(w) = 0$$

$$H(\bar{w} - \bar{w}^*) + \alpha \bar{w} = 0$$

$$H\bar{w} + \alpha \bar{w} = H\bar{w}^*$$

$$(H + \alpha I) \bar{w} = H\bar{w}^*$$

$$w^* = (H + \alpha I)^{-1} H\bar{w}^*$$

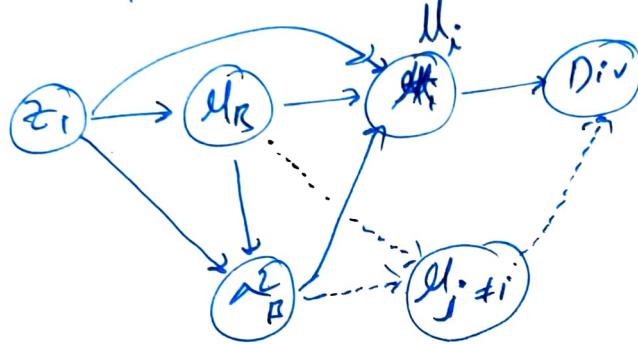
$\alpha = 0 \rightarrow$  have no regularization

$\alpha \neq 0 \rightarrow$   $w$  is rotated back to its original elements

$$\left[ \frac{\lambda_i}{\lambda_i + \alpha} \right] = \begin{cases} \text{close to 1} & \text{for larger eigen values} \\ \text{close to 0} & \text{for smaller eigen values} \end{cases}$$

## Batch Normalization

- Applied before Activation
- input and output are different instances



$$\mu_B = \frac{1}{B} \sum_{i=1}^B z_i$$

$$\sigma_B^2 = \frac{1}{B} \sum (z - \mu_B)^2$$

$$z_i = \frac{z_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

$$\begin{aligned} \frac{d \text{Div}}{d \beta} &= \frac{d \text{Div}}{d \hat{z}_i} \frac{d \hat{z}_i}{d \beta} \\ &\downarrow \\ &\Rightarrow \frac{d \text{Div}}{d \hat{z}_i} \end{aligned}$$

$$\hat{z}_i = \gamma \hat{z}_i + \beta$$

Div = function ( $\hat{z}_i, \mu_B, \sigma_B^2$ )

$$\frac{d \text{Div}}{d \gamma} = \frac{d \text{Div}}{d \hat{z}_i} \underbrace{\frac{d \hat{z}_i}{d \gamma}}_{\mu_i} = \mu_i \frac{d \text{Div}}{d \hat{z}_i}$$

$$\frac{d \text{Div}}{d \mu_i} = \frac{d \text{Div}}{d \hat{z}_i} \cdot \frac{d \hat{z}_i}{d \mu_i} = \frac{d \text{Div}}{d \hat{z}_i} \gamma$$

$$\frac{d \text{Div}}{d z_i} = \frac{d \text{Div}}{d \mu_i} \cdot \frac{d \mu_i}{d z_i} + \frac{d \text{Div}}{d \sigma_B^2} \frac{d \sigma_B^2}{d z_i} + \frac{d \text{Div}}{d \beta} \cdot \frac{d \beta}{d z_i}$$

$$\frac{\partial \text{Div}}{\partial \alpha_B^2} = \frac{-1}{2} (\alpha_B^2 + \epsilon)^{-\frac{1}{2}} \sum_{i=1}^B \frac{\partial \text{Div}}{\partial u_i} (z_i - u_i)$$

$$= \sum \frac{\partial \text{Div}}{\partial u_i} (z_i - u_i) \left[ -\frac{1}{2} \right] (\alpha_B^2 + \epsilon)^{-\frac{1}{2}}$$

$$\boxed{\frac{\partial \text{Div}}{\partial \alpha_B^2} = \frac{-1}{2} [\alpha_B^2 + \epsilon]^{\frac{1}{2}} \sum_{i=1}^B (z_i - u_i) \frac{\partial \text{Div}}{\partial u_i}}$$

$$\frac{\partial \text{Div}}{\partial u_B} = \left[ \sum \frac{\partial \text{Div}}{\partial u_i} \cdot \frac{-1}{\sqrt{\alpha_B^2 + \epsilon}} \right] + \frac{\partial \text{Div}}{\partial \alpha_B^2} \cdot \frac{\sum -2(z_i - u_i)}{B}$$

How?

$$\text{Sol: } \frac{\partial \text{Div}}{\partial u_B} = \sum_i \frac{\partial \text{Div}}{\partial u_i} \cdot \frac{\partial u_i}{\partial u_B} + \frac{\partial \text{Div}}{\partial \alpha_B^2} \cdot \frac{\partial \alpha_B^2}{\partial u_B}$$

$$\frac{\partial \text{Div}}{\partial u_B} = \sum_i \frac{\partial \text{Div}}{\partial u_i} \cdot \frac{-1}{\sqrt{\alpha_B^2 + \epsilon}} + \frac{\partial \text{Div}}{\partial \alpha_B^2} \cdot \frac{\sum -2(z_i - u_i)}{B}$$

This whole term = 0

$$\boxed{\frac{\partial \text{Div}}{\partial u_B} = \sum \frac{\partial \text{Div}}{\partial u_i} \cdot \frac{-1}{\sqrt{\alpha_B^2 + \epsilon}}}$$

→ for each batch, there exists different mean.  
so at the time of inference:-  $B, B-1$  can be ignored

$$\begin{array}{ll} u_B, \alpha_B, & \left\{ u_B = \frac{\sum_i u_i}{k} \right. \\ u_{B_1}, \alpha_{B_1} & \left. \alpha_B^2 = \frac{1}{k} \sum_{i=1}^k \left( \frac{u_i - u_B}{B} \right)^2 \right\} \\ \vdots & \\ u_{B_k}, \alpha_{B_k} & \end{array}$$

$\hat{\mu}$  = unbiased estimator of  $\mu$

$x_1, x_2, \dots, x_n$  are iid rv

$$x_i \sim (\mu, \sigma^2)$$

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

$$E(\hat{\sigma}^2) \neq \sigma^2$$

$$E\left[\frac{n}{n-1} \hat{\sigma}^2\right] = \sigma^2$$

$$\hat{\sigma}^2 = \frac{n}{n-1} \hat{\sigma}^2 = \frac{1}{(n-1)} \sum_{i=1}^n (x_i - \mu)^2$$

$$\begin{aligned} E(\bar{x}) &= E\left[\frac{1}{n} \sum_{i=1}^n x_i\right] \\ &= \frac{1}{n} [\sum E(x_i)] \\ &= \frac{1}{n} \sum \mu \\ E(\bar{x}) &= \mu \end{aligned}$$

- BN is applied to only some layers [or selected neurons only]
- BN improves convergence rate and NN performance
- BN takes care of step size also.

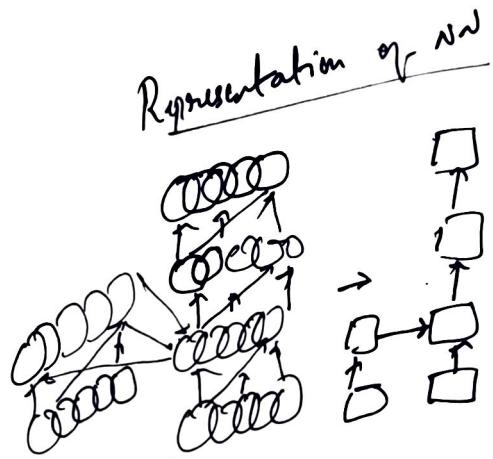
## Modelling Series

RNN

- most models until now consider data was coming from constant series.
- In this data is a non constant like weather data, sequential data, video data, speech data, music.

Input = Sequence of Vector  
Output = Classification Result

Input = Vector  
Each layer has many neurons



→ output of hidden layer at  $t-1$  is the input of the hidden layer at  $t$

finite-response model

→ something that happens today only affects the output of the system for  $n$  days into the future

infinite response systems :-

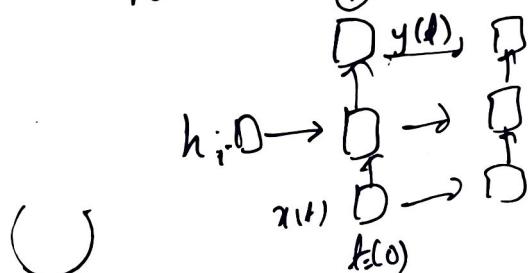
long term effect on output into the future

$$y_t = f(x_t, y_{t-1})$$

that's why Recurrent NN

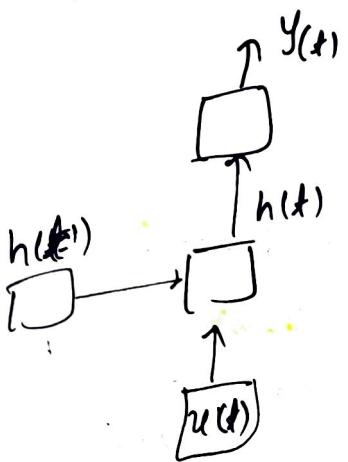
State Space Model

each layer represented by a box have same network replicating each time stamp,  $\textcircled{1} = \textcircled{2}$  with only change in the values.



→ while Backprop in RNN, if at  $t=0$ , gradients are also going to  $(t-1)$ , back in time. This is called backprop <sup>26</sup> through time

given - collection of sequence inputs



Batch Normalization

Convergence can be improved using smoothed updates

$$\underset{w_1, w_2, \dots, w_k}{\text{argmin}} \ L(w_1, w_2, \dots, w_k)$$

Dropout

→ To understand this we need to understand **Bagging**

- turn off some neurons with prob  $(1-\alpha)$
- Pattern of dropped nodes change for each input
- Backprop is only performed over remaining part
- $2^N$  = possible Networks [for network with  $N$  neurons]
- Obtains by choosing different subsets of nodes
- Mechanism to increase pattern diversity
- helps to learn
  - Rich pattern
  - Redundant patterns
  - learns denser patterns + redundancy

Training

- Trains  $2^N$  networks
- Bagged output = Avg over all  $2^N$  networks = statistical expectation of output over all networks

$$\text{Bagged output} = Y = E[\text{network } [y_j^k, j=1 \dots D_k, k=1 \dots k]]$$

Can't compute this Expectation  
So we use approximation

$$E[\text{network } [y_j^k, t_{kj}]] = \text{network } [E[y_j^k] + t_{kj}, j]$$

$$\text{Activation of each neuron} = y_i^k = \text{D}\alpha \left[ \sum w_{ji}^k y_j^{k-1} + b_i^k \right] = \alpha \text{ } \alpha(z_i^k)$$

D = Bernoulli Variable

D = 1 with probability  $\alpha$

thus

$$E[y_i^k] = \alpha \alpha \left[ \sum w_{ji}^k y_j^{k-1} + b_i^k \right]$$

$$w_{\text{test}} = \alpha w_{\text{trained}}$$

NPTEL IIT Guwahati

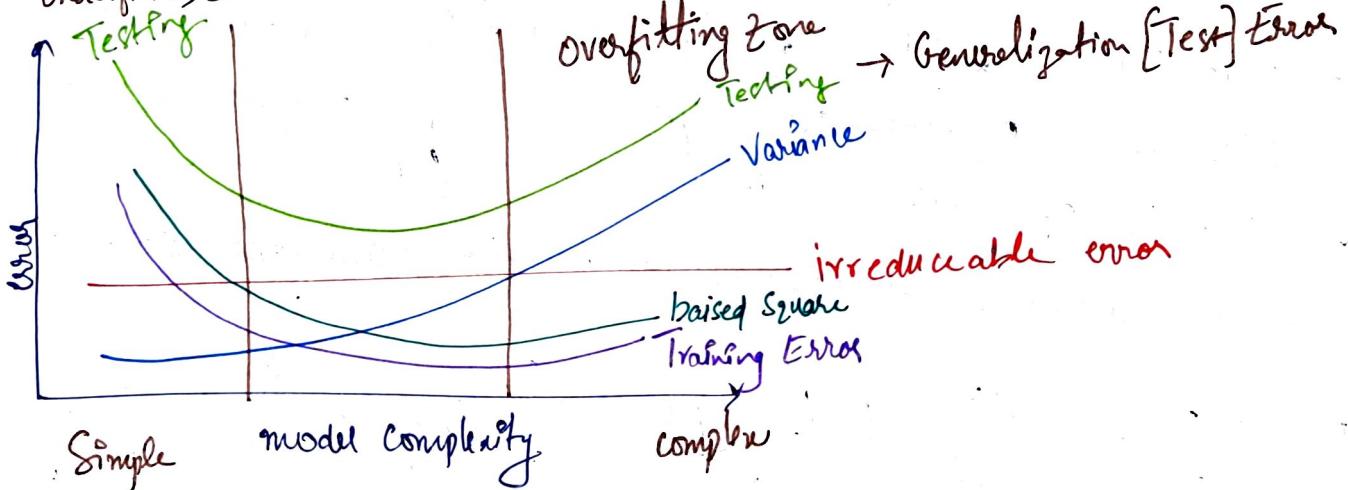
Bias - Variance Trade off

High Variance  $\rightarrow$  Overfitting

High Bias  $\rightarrow$  Underfitting

Test error = Bias Error + Variance + irreducible error

Underfitting zone



Problem Definition

$x$  = Independent Variable

$y$  = Dependent Variable

$y = f(x) + \epsilon \rightarrow \text{Noise}$

$$\text{var}(\epsilon) = E[\epsilon^2] = \sigma_\epsilon^2$$

$\epsilon = \text{Noise, zero mean variance } \sigma_\epsilon^2$   
 $\therefore E[\epsilon] = 0$

$\Rightarrow$  to determine  $\hat{f}$

Q. to find a function  $\hat{f}$  such that it is as close to the true function  $f$  (can learn from training data)

$f$  is learned by minimizing a loss function

Goal  $\rightarrow$  to bring predictions  $\hat{f}(x)$  from training data as close as possible to their observed value.

i.e.,  $y \approx \hat{f}(x)$

MSE  $\rightarrow$  loss function

$$\text{MSE} = E[(y - \hat{f}(n))^2]$$

Squared difference of prediction from the true value  $y$ .

$$\text{Bias} = \text{bias}[\hat{f}(n)] = E[\hat{f}(n)] - f(n) \rightarrow$$

it is the difference of the average value of prediction to the true underlying function  $f(n)$  for a given test point  $x$

$$\text{Variance} = \text{Var}[\hat{f}(n)] = E[(\hat{f}(n) - E[\hat{f}(n)])^2]$$

$\downarrow$   
mean squared deviation of  $\hat{f}(n)$  from its expected value  $E(\hat{f}(n))$  over different realizations of training data.

MSE + bias + variance + irreducible error.  $\rightarrow$  combined formula

$$E[E[y - \hat{f}(n)]^2] = E[\underbrace{\text{bias}[\hat{f}(n)]^2}_{\text{MSE}} + \underbrace{E[\text{Var}[\hat{f}(n)]]}_{\text{Variance}} + \underbrace{\sigma_e^2}_{\text{irreducible error}}]$$

decomposing MSE to bias, variance and irreducible error

minibatches, RMS Prop, Quick Prop ...

[50:00 minutes] SGD uses samples of only one samples gradient  
 SGD converges faster and has large variance. | Provides batch updates  
 So need a better one :-

### ① Minibatch update (MB)

- adjust function at small, randomly chosen subset of points
- Randomly sample (shuffle the data randomly)

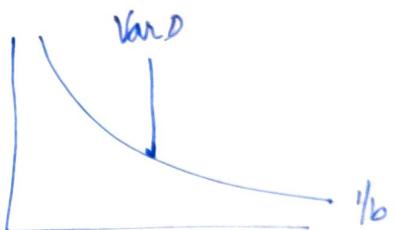
How this works:-

$$\text{MB loss} = \frac{1}{b} \sum \text{div}(f(x_i; \omega), d_i)$$

MB loss = unbiased estimate of expected loss

$$E[\text{MB loss}] = E[\text{div}(f(x; \omega), g(x))]$$

$$\text{Var(batch loss)} = \frac{1}{b}$$



→ Variance is much smaller than SGD

→ Converges faster [overall training loss]

→ loss is computed on entire training set, but then need to process whole data once, hence, it will be complicated / costly

Minibatch size = largest that the hardware supports

Large MB = less variance

= fewer updates per epoch

SGD variance → T than Batch

Online → As the data come in.

MB

- Has low variance than SGD

- Converges worse than SGD

- compensate ↑ by being able to perform batch processing.

Momentum Update 31

- maintains running average of all gradients until current step

$$\Delta w^k = \beta \Delta w^{k-1} - \eta \nabla_w \text{loss}(w^{k-1})^\top \quad \beta = 0.9$$

$$w^k = w^{k-1} + \Delta w^k$$

at any iteration, to get current step :-

- ① compute gradient step at current location
- ② Add scaled previous step
- ③ this is actually running average

To get final step:-

→ Updates are computed in 2 stages:

- ① take step against gradient at current loc.
- ② Add scaled version of previous step

→ this can be more optimal by reversing the order of operations  
→ Newton's method

change the order of steps

- ① extent previous step
- ② compute gradient at resultant position
- ③ add the obtained final step

Momentum is smoothing out the variations

These methods will improve convergence by Normalizing the mean of the derivatives

Other methods also consider their variance

- RMS prop
- Adagrad
- Ada delta
- **ADAM (Popular)**

## RMS Prop

Considers total movement of the gradient.

Variant based on MiniBatch SGD algorithm

[Squares of the derivative]  $\rightarrow (\partial_w D)^2$

$$\text{Running average} = \frac{1}{\sum E[(\partial_w D)_k^2] + \epsilon} \quad \left. \begin{array}{l} \text{Scaling down the running} \\ \text{average.} \end{array} \right\}$$

- Maintain running estimate of MS value of derivatives for each parameter
- Scale update of parameter by square root [mean squared Derivative]

$$E[(\partial_w D_k^2)]_k = \gamma E[(\partial_w D)_k^2] + [1 - \gamma](\partial_w D)_k^2$$

$$w_{k+1} = w_k - \frac{\eta}{\sqrt{E[(\partial_w D)_k^2] + \epsilon}} \partial_w D$$

This is similar to RProp.

Thing that is not happening here are:- Running average of derivatives is not being changed, and Momentum.

## ADAM

Combination of RMS Prop + Momentum  
Utilises Smoothed Version of Momentum-Augmented gradient  
[Considers first and second moments]

$M_k$  = Running Average of derivative

$V_k$  = Squared derivative

$\gamma, \delta$  = To make sure these terms do not dominate in early iterations.

### RMS Prop

$$\eta = 0.001$$

$$\gamma = 0.9$$

### ADAM

$$\eta = 0.001$$

$$\delta = 0.9 \text{ & } \gamma = 0.999$$

} Typical Params

RNN - Modelling Network

- All columns are identical (even incoming edges)

$$h^i(t) = f_i [w^i x(t) + w^{ii} h^i(t-1) + b_i]$$

current weight      Recurrent weights

$$y(t) = f_2 [w^2 h^i(t) + b_k, k=1 \dots n]$$

Notations

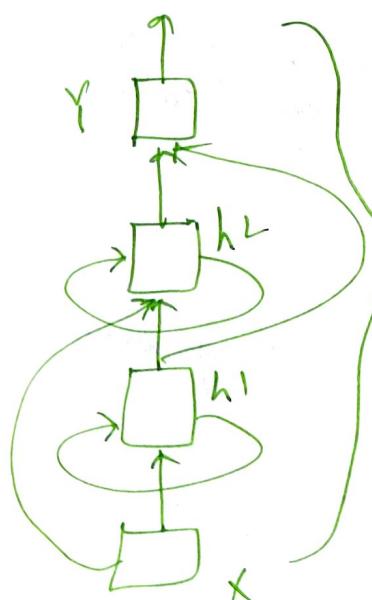
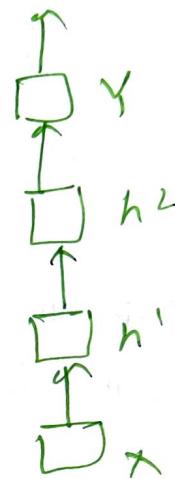
$w^i$  → Weights of layer  $i$

$w^{ii}$  → weights of layer  $i$  to layer  $i$

$w^2$  → weights of layer 2

$$h^2(t) = f_2 [w^2 x(t) + w^{22} h^2(t-1) + b^2]$$

Recurrent State Activation =  $\tanh()$  — typically



$$h^i(t) = f_i [w^{0,i} x(t) + w^{ii} h^i(t-1) + b^i]$$

$$h^2(t) = f_2 [w^{12} h^1(t) + w^{0,2} x(t) + w^{22} h^2(t-1) + b^2]$$

$$y(t) = f_3 [w^3 h^2(t) + w^{13} h^1(t) + b^3]$$

Backprop <sup>through time</sup> focuses on one training instance

Divergence computed is not just the error, it is b/w the sequence of outputs by network and desired seq of outputs

and this is the scalar function of series of vectors

Div is function of all outputs  $y(0) \dots y(T)$

$$\text{global Div} = \sum \text{Div}(t)$$

Backprop : Should start at end

Once we have gradients at  $y(T)$   
going backward accumulating, have  $\frac{dy(T)}{dt}$  what we need

to do to get derivat  
of affine value?

We should do Chain Rule  
 Should do Jacobian of the activation - will give derivative of output  
 Now hidden layer → multiply derivative  $z_2$  to  $\frac{dz_n}{dn}$  (weight matrix)

Class 12

14/9/23

Training the Network

$(x_i, D_i)$  = Inputs

$x_i = x_{i,0}, x_{i,1}, \dots, x_{i,T}$  (given)

$D_i = D_{i,0}, D_{i,1}, \dots, D_{i,T}$  → Targeted Output

$y_i = y_{i,0}, y_{i,1}, \dots, y_{i,T}$  → Output of the Network

$y_i = y_{i,0}, y_{i,1}, \dots, y_{i,T}$  → Predicted Output

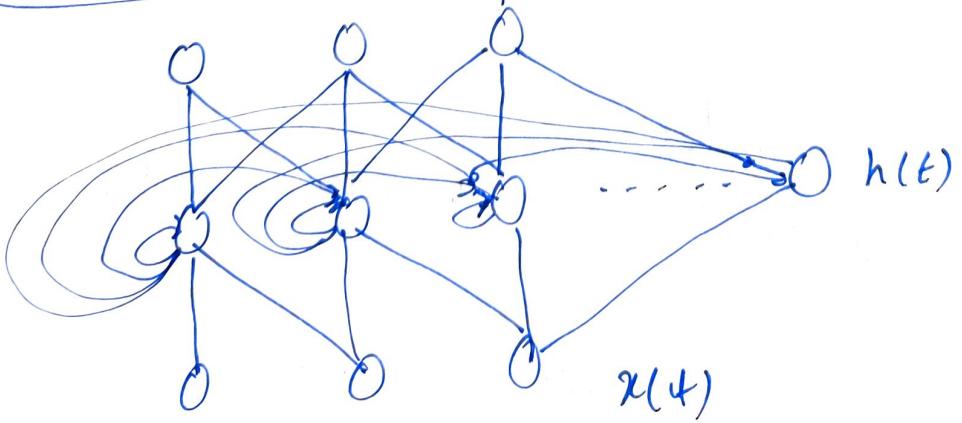
Loss =  $\text{Div} [ [y_{i,0}, y_{i,1}, \dots, y_{i,T}] - [D_{i,0}, D_{i,1}, \dots, D_{i,T}] ]$

$y_i$  compared to  $D_i$  only.

$$= \sum \text{Div} (Y(t) \cdot D(t))$$

$$= \frac{\partial \text{Div} [ [y_{i,0}, y_{i,1}, \dots, y_{i,T}] - [D_{i,0}, D_{i,1}, \dots, D_{i,T}] ]}{\partial Y(t)}$$

$$\boxed{\text{Loss} = \frac{\partial \text{Div} [ Y(t) \cdot D(t) ]}{\partial Y(t)}}$$



$$z_i^2(t) = \sum_{j=1}^{N_i} w_{ji}^2 h_j(t) + b_j^2 \quad i=1, \dots, m$$

$f_i^2(z_i^2(t))$   
activation function

### Background

- will only focus on one training instance
- Using forward pass, we computed  $y(t)$
- Using backprop we need to compute gradient from output layer starting from  $y(t)$

Step 1 - finding gradients wrt  $y(t)$  → of the  $t^{\text{th}}$  output

$$y(t) \in \mathbb{R}^{n_t}$$

$$\frac{\partial \text{Div}}{\partial y(t)}$$

Step 2 : gradients wrt  $z_2^2(t)$

$$\nabla z_2^2(t) \text{ Div} = \nabla y(t) \text{ Div } \nabla z_2^2(t) y(t)$$

activation can be softmax also

Step 3: gradient wrt  $z_i^2(t)$  1st component of  $z_i^2(t)$

Scalar

$$\frac{\partial \text{Div}}{\partial z_i^2(t)} = \frac{\partial \text{Div}}{\partial y_i^T} \cdot \underbrace{\frac{\partial y_i^T}{\partial z_i^2(t)}}_{\text{Vector}}$$

$$\frac{\partial \text{Div}}{\partial z_i^2(t)} = \sum_j \nabla z_j^2(t)$$

derivations → Pg 5 / 26

If we need to find gradient with all components

then

$$(\nabla z^2(t) y(t)) \text{ (Pg 6)}$$

→ each  $z_{iT}$  is formed by applying activation for  $q$   $h_{iT}^{32}$

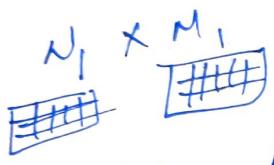
→  $z_{iT}$  has  $n_1$  components.

→  $w_{ij}$  → weight connecting to input layer

$$z_{j,T}^i = \sum w_{ij}^1 w_{1,T} + \sum w_{ij}^{2,2} h_2(T-1) + b_j^i$$

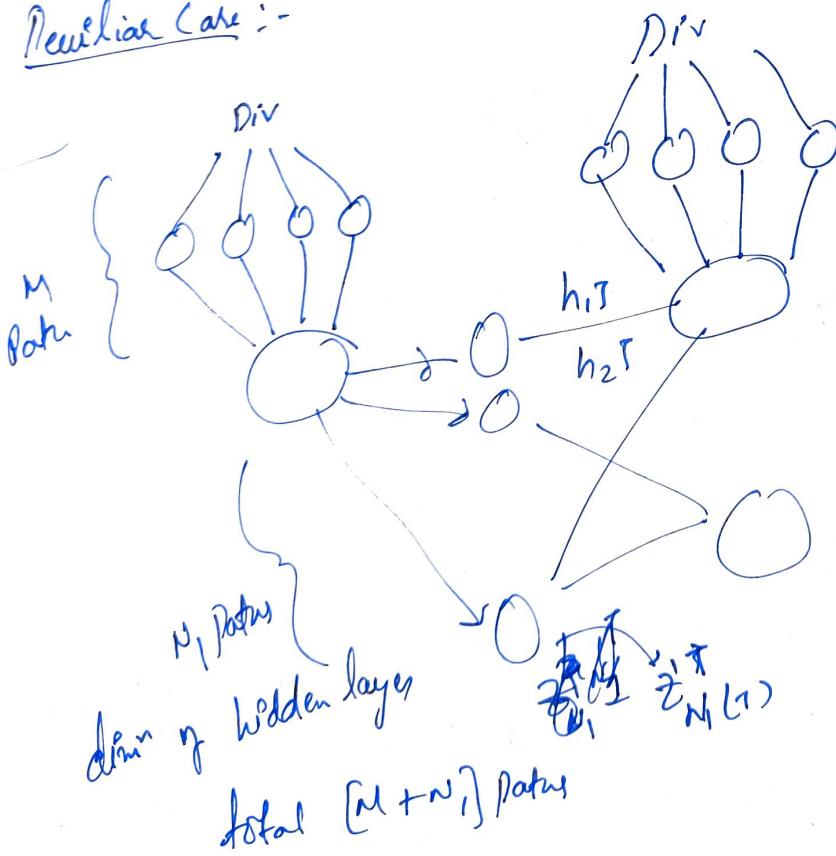
$$\nabla_{w^1} \text{Dir} = x_i \nabla_{z_{1,T}^i} \text{Dir}$$

$$\nabla_{w^n} \text{Dir} = h(T-1) \nabla_{z_{T-1}^i} \text{Dir}$$



Computed one set of gradients at  $T$ , should do the same for time  $(T-1)$

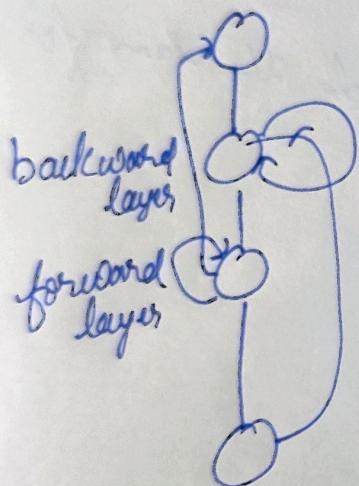
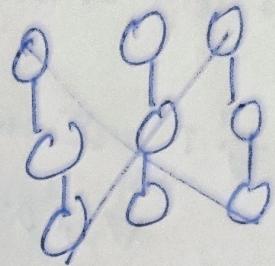
Reular Case :-



$$\frac{\partial \text{Div}}{\partial h_i(t-1)} = \sum_{j=1}^M \frac{\partial \text{Div}}{\partial z_j^2(t-1)} \cdot \underbrace{\frac{\partial z_j^2(t-1)}{\partial h_i(t-1)}}_{w_{ij}^2} + \sum_{j=1}^{N_1} \frac{\partial \text{Div}}{\partial h_j(t)} \cdot \underbrace{\frac{\partial h_j(t)}{\partial h_i(t-1)}}_{w_{ij}^{(1)}}$$

### Bidirectional RNN

has 2 hidden layers  
 $H(t)$  depends on  $U(t-1)$   
 $H(t-1)$  will depend on  $H(t)$



then both are added directly.

- BIBO : bound Input Bounded Output
- assume activation function as linear
  - Analyse BIBO stability

$$h_k = z_k$$

Linear Revision Polynomial depends on this

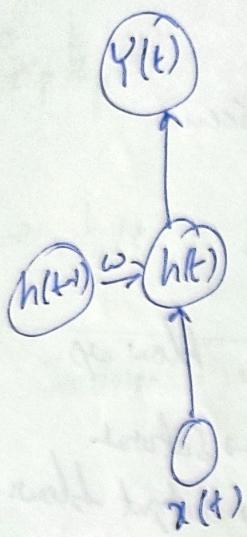
$$h(t) = w(h(t-1)) + c x(t) + b$$

$$h_0(t) = w^t c x(0)$$

only supplying input at  $t=0$  from  $t=1$  there no input

$w$ : matrix of feedback loop

$w < 1 \rightarrow$  Vanishing Gradient  
 $w > 1 \rightarrow$  exploding Gradient



$$\begin{aligned}
 Y(t) &\in R^m \\
 h(t) &= w h(t-1) + c x(t) \\
 &= w [w h(t-2) + c x(t-1)] + c x(t) \\
 &= w^2 h(t-2) + c w x(t-1) + c x(t) \\
 &\vdots \\
 &= w^t h(0) + c \left[ \sum_{k=0}^{t-1} w^k x(t-k) \right] \\
 &= w^t h(0) + c w^t x(0) + \underbrace{c \sum_{k=1}^t w^k x(t-k)}_0
 \end{aligned}$$

$$x(i) = 0, \quad i = 1, \dots, t$$

$$h(t) = \boxed{c w^t x(0)}$$

class 13

18/9/23

39

13/100 stability with linear activation fn. — seen

Now

Non-linear fn's

Sigmoid

Saturates for limited steps, regardless of  $\omega$

Tanh

sensitive to  $\omega$  but eventually saturates

Relu

sensitive to  $\omega$ , can blow up  
all values  $< 1 \rightarrow$  saturate  
at  $> 1 \rightarrow$  output blows

multilayer perceptrons, Deep NN [pg 48]

$$Dv(n) = D \left[ f_n \left( w_{n-1} f_{n-1} \left[ w_{n-2} f_{n-2} \left[ \dots \left[ w_0 u_0 \right] \right] \right] \right) \right]$$

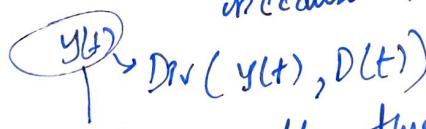
loss of [ output of deep NN ]

[Bas Vanhouw → Cornell]

$Dv$  (activation fn) = Always Bounded

Can think of this as a flowing operator. [pg 52] The Jacobian because of this, there is vanishing gradient

RNN



$Dv(y(t), D(t))$

all of these are contributing to the loss.  
because this is sequential model or infinite

$$\frac{\partial h_t(w)}{\partial \theta} = \frac{\partial h_t(w)}{\partial \theta} \underbrace{h_{t-1}}_{\text{constant wrt } \theta}$$

$$\Rightarrow z_p(t) = \sum_{j=1}^n w_{jp}^{(1)} h_j(t-1) + \sum_{j=1}^n w_{jp} x_j(t) + b_0$$

$$\Rightarrow h_p(t) = f(z_p(t))$$

Jacobian = derivative that comes after applying an activation fn  
 $\frac{d}{du}$  (activation fn) (1 ||) norm = magnitude

$$\frac{\partial h(t)}{\partial z_i(t)} = \text{Jacobian} = f'(z(t))$$

$$|f'(z(t))| \leq \gamma$$

tanh

$$\frac{\partial(h(t))}{\partial h_k} = \prod_{i=k+1}^t w^{(i)} \nabla_{z^{(i)}} h(t)$$

$$\left\| \frac{\partial h(t)}{\partial h_k} \right\| = \prod_{i=k+1}^t \|w^{(i)}\| \left\| \nabla_{z^{(i)}} h(t) \right\|$$

$$\left\| \frac{\partial h(t)}{\partial h_k} \right\| \leq \frac{1}{\gamma} \quad \text{if } \gamma < 1 \quad \text{if then we see Vanishing gradient}$$

multiplication

if greater than will ~~not~~ be exploding gradients

Vanishing | Exploding gradient?  
 will depend on Eigen value (spectral) properties of  $w^{(i)}$

LSTM

(CEF)

Uses concept of Constant Error Flow for RNN to  
 Create constant Error Carousel (CEC), so error gradients  
 don't decay

- State ( $s_t$ ) of RNN records info from all previous time step
- at each time step, old info is replaced by current input.
- This has limited capacity, can't store the entire memory. Some old info has to go down to accommodate new info.
- Similar thing happens when we do backprop in RNN
- if we make wrong prediction in previous step, then the current step will also be wrong. See that this is the backprop.

writing on board - Analogy

writing on board - Analogy for RNN and while mathematically?

(8) How to convert these intuitions for RNN and while mathematically?

$$s_t = f(w s_{t-1} + c z_t + b)$$

this should be  
selective

there exist one output gate that  
selects the value that's important for us to use in  
the next step

$o_{t-1}$  = selector = distributed between 0 and 1  
also called as Output gate

$$o_{t-1} = \sigma(w o_{t-2} + c_o z_{t-1} + b_o)$$

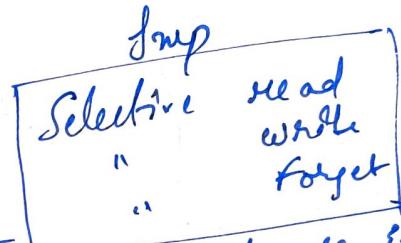
CS7015 DL

$$h_{t-1} = \underbrace{o_{t-1}}_{\text{Vector}} \underbrace{\sigma}_{\text{Vector}} \underbrace{(s_{t-1})}_{\text{Vector}}$$

element wise product

$$\tilde{s}_t = \underbrace{w h_{t-1}}_{\text{selective portion of } s_{t-1}} + (z_t + b)$$

selective portion of  
 $h_{t-1}$



$$j_t = \alpha (\omega h_{t-1} + c_i u_i + b_i)$$

$$j_t \odot \tilde{s}_t$$

Previous State —  $s_{t-1}$

Output gate —  $o_{t-1} =$

Schtr. write —  $h_{t-1} = o_{t-1} \odot a(s_{t-1})$

current (temporary) state —

Input state —

Selective forget —  $s_t = s_{t-1} + l_t \odot \tilde{s}_t$

Forget gate — to forget some part of  $s_{t-1}$

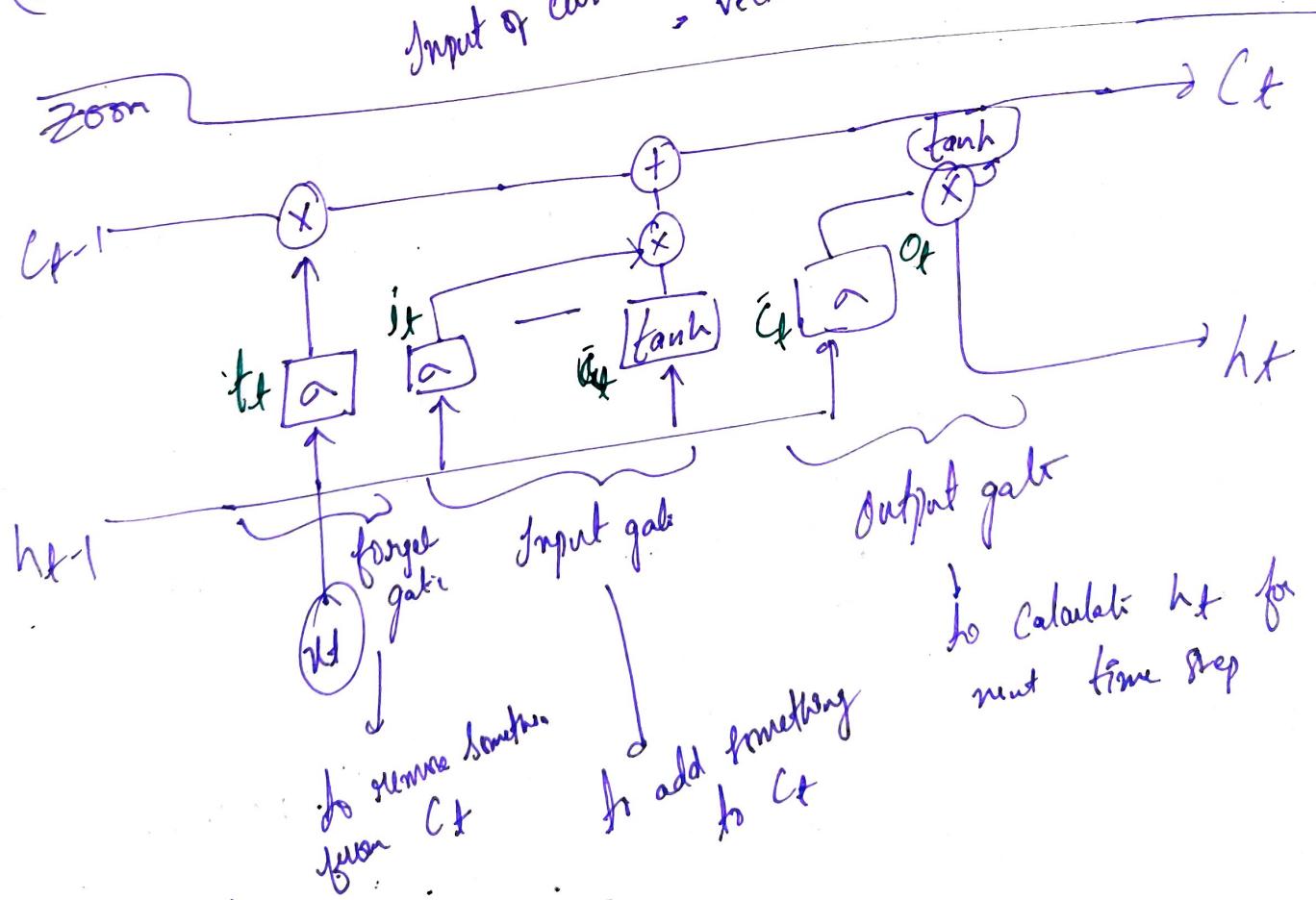
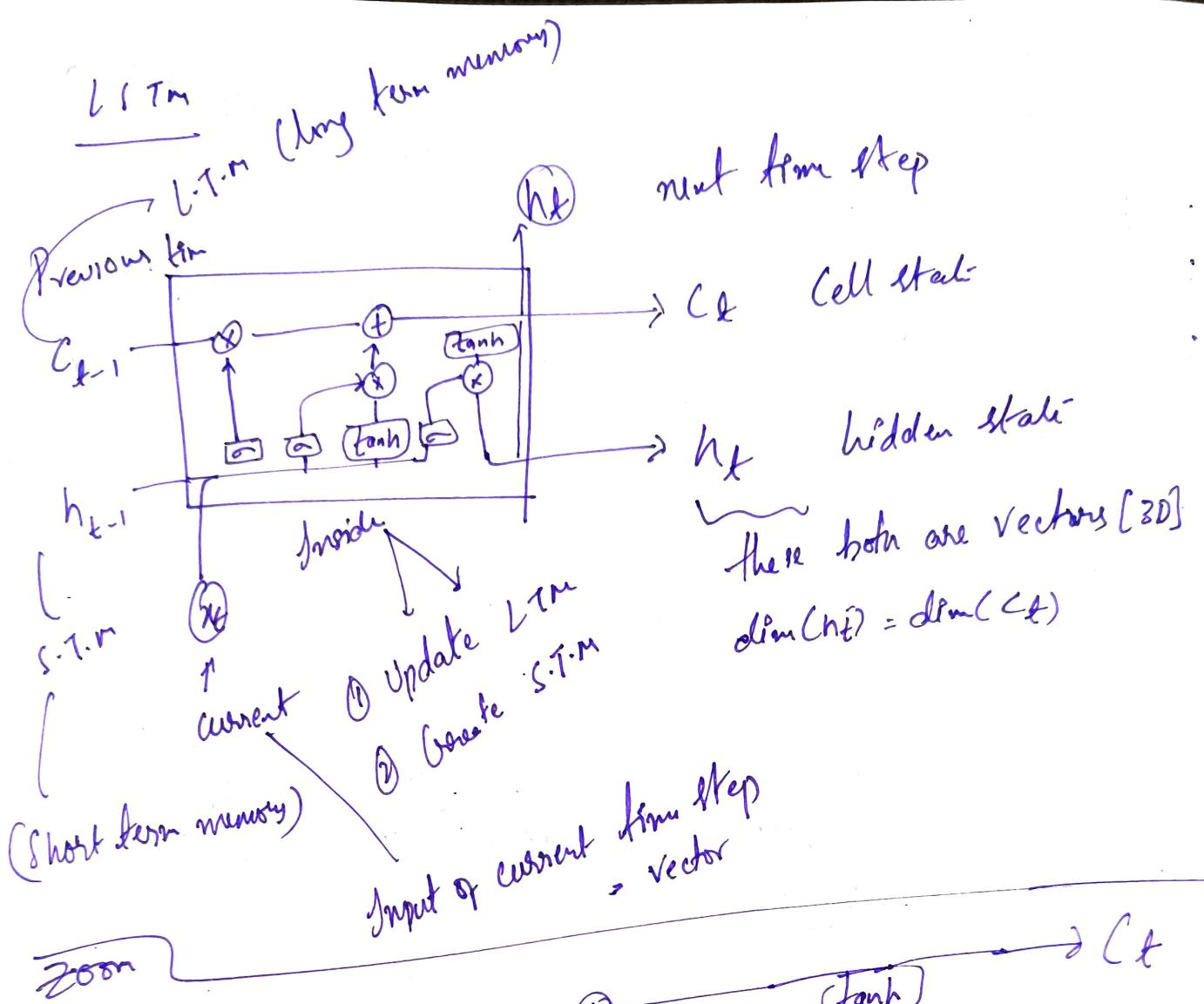
$$f_t = \alpha (\omega_f h_{t-1} + c_f u_t + b_f)$$

$$s_t = f_t \odot s_{t-1} +$$

final output

$$h_t = o_t \odot \alpha(s_t)$$

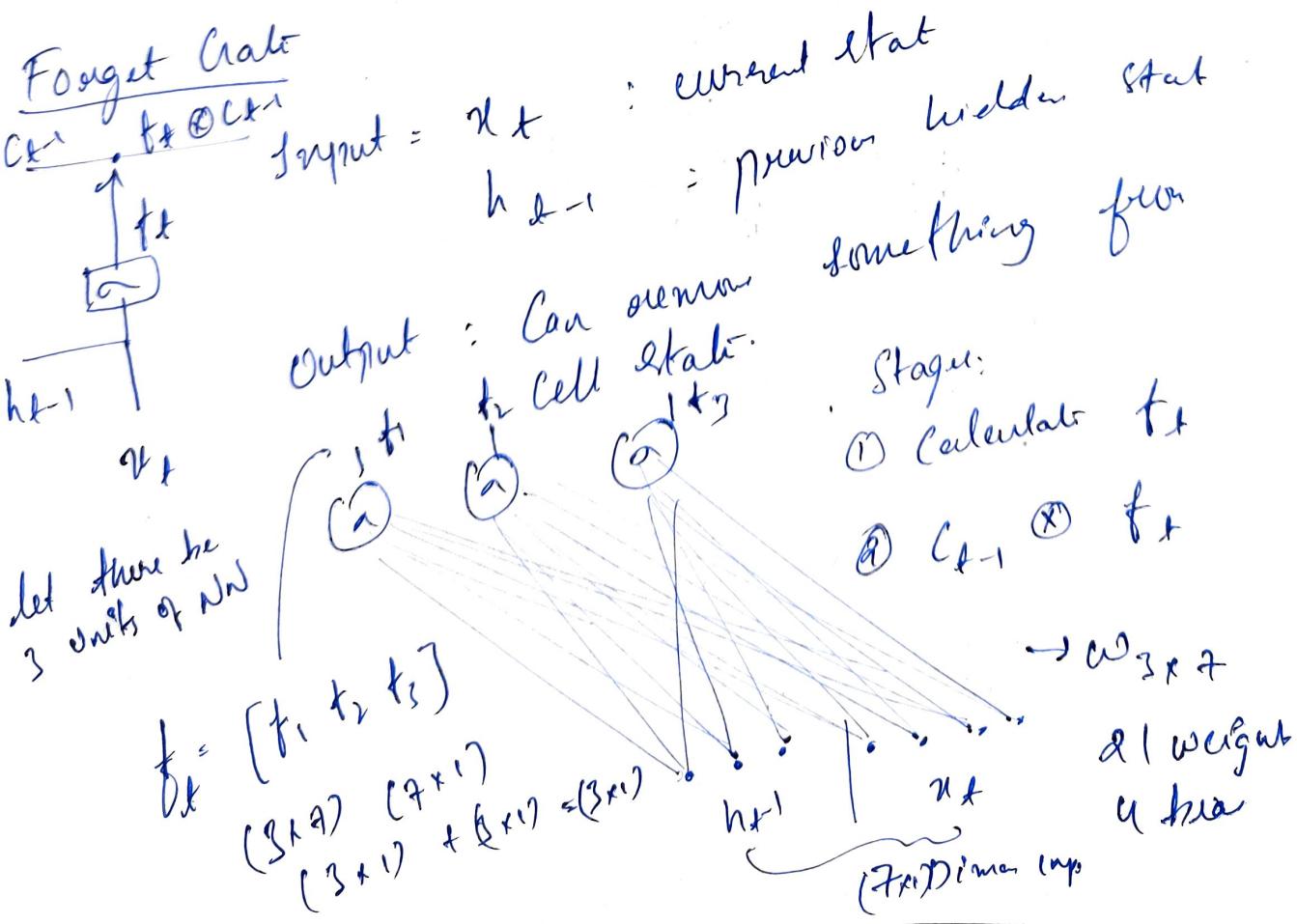
Temporary state  
Input state  
Output state



the gates:- } Forget Gate  
 ② Input  
 ① Forget  
 ③ Output

multihot moded.  
 $f_t$  is  $\tilde{c}_t$  forget input candidate cell state  
 $O_t$  output } vector  
 $\dim(f_t) = \dim(\tilde{c}_t) = \dim(O_t) = \dim(c_t)$   
 $\dim(f_t) = \dim(\tilde{c}_t) = \dim(O_t) = \dim(c_t) = \dim(h_t) = \dim(x_t)$   
 Pointwise operation  $\rightarrow \oplus, \otimes, \text{tanh}$   
 $\hookrightarrow \text{output} = \text{vector}$

$\alpha, \text{tanh}$  : neural networks collection of multiple nodes with sigmoid and tanh activations.  
 Sigmoid:

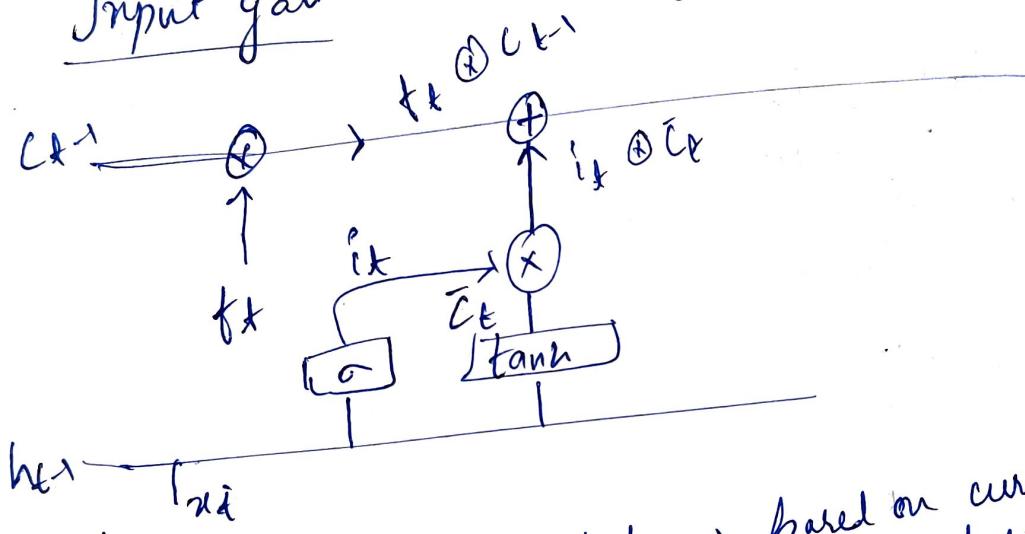


$$f_t = \sigma \left[ w_f \underbrace{[h_{t-1}, x_t]}_{\substack{1 \times 7 \\ 7 \times 1}} + b_f \right] \rightarrow (3 \times 1)$$

$[h_{t-1}, x_t] \rightarrow$  concatenation of 2 vectors

$$f_t \odot C_{t-1} \rightarrow 3 \times 1$$

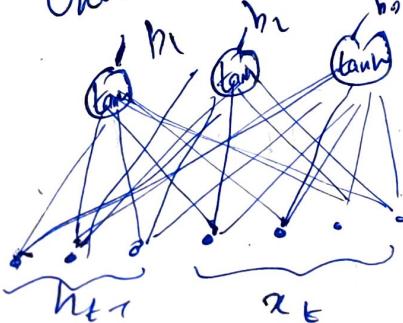
Input gate adds new info to the data



Stages:

- ①  $\tilde{C}_t$ : Candidate cell state  $\rightarrow$  based on current input and previous cell state  $\rightarrow$  we add new info that's needed to be added.
- ②  $i_t$ : decide what item needed to be added to final cell state.
- ③  $C_t$ : calculate the final cell state.

Units from NN  $\Rightarrow$   $\tilde{C}_t$



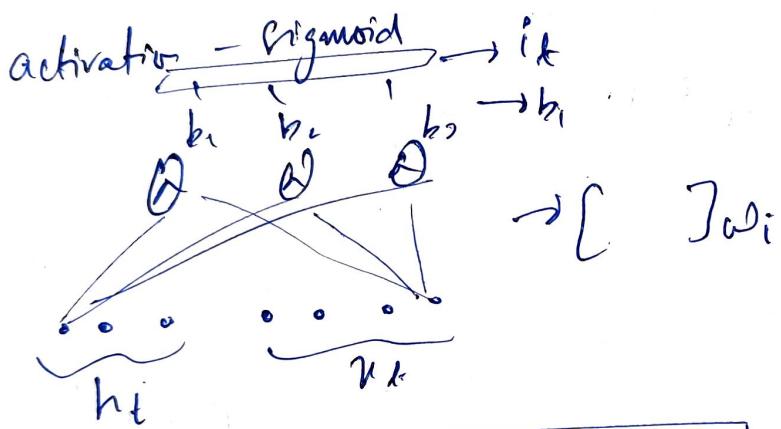
$$\rightarrow w_R (3 \times 3)$$

$$\tanh [w_c (h_{t-1}, x_t) + b_c]$$

$$\tilde{c}_t = \tanh \left[ w_c \underbrace{\left( h_{t-1}, x_t \right)}_{(7 \times 1)} + b_c \right]$$

(3x7) (3x1)  
(3x1) (3x1)

Potential import info =  $\bar{C}_t$   
 is filter from  $\bar{C}_t$ , next step, calculating it

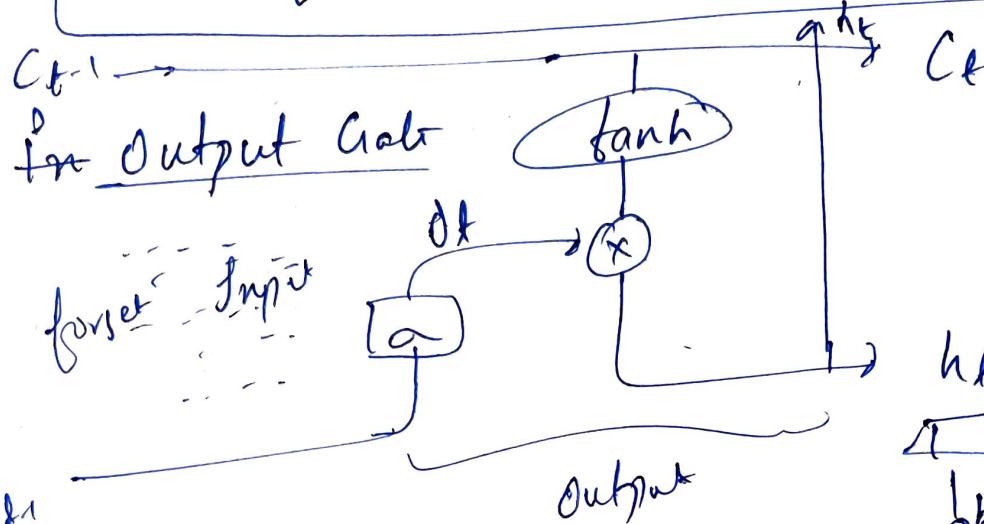


$$l_t = \alpha \underbrace{\left( w_i [h_{t-1}, x_t] + b_i \right)}_{\text{Prediction}} + \beta \underbrace{\left( \sum_{j=1}^{t-1} \left[ h_j, x_t \right] \right)}_{\text{Past Inputs}}$$

$$\boxed{i_t \otimes C_t \rightarrow \bar{C}_t^*} \rightarrow \text{filtered Candidate Cell Stab}$$

$$\text{Now, } \tilde{C}_t = \begin{pmatrix} f_t \otimes C_{t-1} \\ i_t \otimes \tilde{C}_t \end{pmatrix}$$

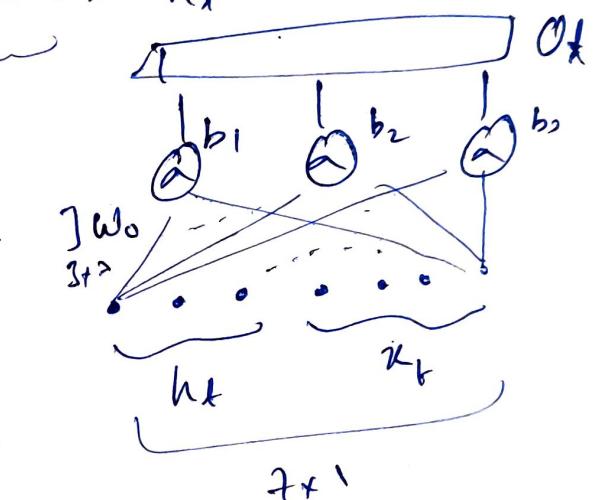
$$C_t = [f_t \odot C_{t-1}] \oplus [P_t \odot \tilde{C}_t]$$



$$O_t \odot \tanh(C_t)$$

$$O_t = \sigma [W_0 [h_{t-1}, u_t] + b_0]$$

$3 \times 7$        $7 \times 1$        $7 \times 1$   
 $3 \times 1$   
 $\sigma$   $(3 \times 1) \times (3 \times 1)$



$$h_t = O_t \odot \tanh(C_t) \rightarrow 3 \times 1$$

$3 \times 1$   
 $3 \times 1$

Class 14  
25/9/23

mid term

- (Q1) Cross entropy is not a very stable loss to use it  
 $H(P, Q) = -\sum p(u) \log q(u)$   
 ↓  
 True probability      Predicted probability
- (Q4) RNN only forward pass  
 [Solutions on Monday]

Module 15-3

LSTM can't handle exploding gradients. But can handle vanishing gradients.

(Q) How they control vanishing gradients?

Ans)

$$\begin{aligned} s_t &= i_t \odot \tilde{s}_t + f_t \odot s_{t-1} \\ h_t &= s_t \odot o_t \end{aligned}$$

$\odot$  : element wise Product

backprop is very complicated to show, this blocks the gradient  
 \* If  $s_{t-1}$  is not contributing much then it means  $f_t$  might be 0, so in this case  $s_t$  is not responsible for  $s_t$ 's vanishing gradients

→ V.G happen in RNN when gradients flow in all paths is 0, cause we are multiplying with power of w

Gradients will explode depending on the power of  $\omega$ .

In LSTM there exist one path too where gradients won't  
vanish or explode.

(Q) Is  $\frac{\partial L}{\partial s_n}$  alive or Vanished?

Ans) There are many paths to  $L_t$  from  $s_n$ , so should calculate the perfect path where gradients doesn't vanish.

$$h_t = \sigma(s_t) \odot o_t$$

$$h_t = a(s_t) \odot o_t$$

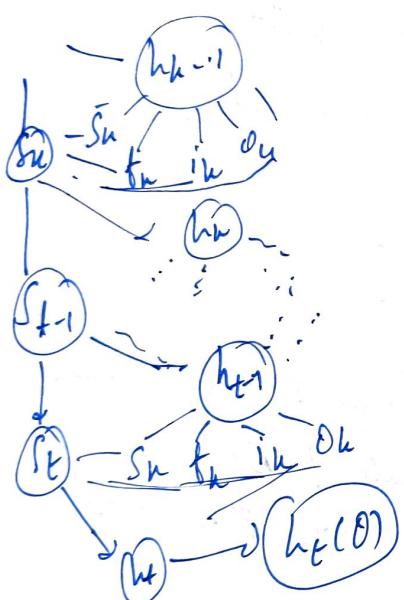
$$\frac{\partial h_t}{\partial s_{t+1}} = a'(s_{t+1}) \times (o_t)$$

This will be a Jacobian Matrix

$$\frac{\partial h_k}{\partial s_{t+1}} = \begin{pmatrix} \frac{\partial h_1}{\partial s_{t+1}} & 0 & \cdots & 0 \\ 0 & \frac{\partial h_2}{\partial s_{t+1}} & \cdots & 0 \\ \vdots & \vdots & \ddots & 0 \\ 0 & \cdots & \cdots & \frac{\partial h_{t+2}}{\partial s_{t+1}} \end{pmatrix}$$

$$= D \begin{bmatrix} o_t \odot a'(s_t) \end{bmatrix}$$

Diagonal



$$\frac{\partial s_t}{\partial s_{t-1}} = ?$$

$$s_t = f_t \circ s_{t-1} + i_t \circ \tilde{s}_t$$

$$\tilde{s}_t = [W h_{t-1} + V x_t + b]$$

$$h_{t-1} = a(s_{t-1}) \circ o_{t-1}$$

Encoder / Decoder Model, Attention Mechanism

reading entire sentence — encoder part  
generating entire sentence in another lang - decoder  
needs 2 RNN's.

Image Captioning — needs 2 blocks.

attention mechanism :- when encoding info in one language,  
do we need entire info to translate? No we can  
translate as an ongoing process. A.m does this work  
give priority for which part to give

class 15  
5/10/23

Mid Sem Paper discussion

(8) multilabel problem classification  
multihot vector representation.

$n$  classes

if class  $1, 3 \rightarrow$  true for some, then its multihot  
vector representation =  $[1, 0, 1, 0]$

$l_{\text{cross ent}} \geq l \log l$

S.T.: not a good loss fn.

$k$ , total no. of classes

$L$  = example labels are true

$l \leq k$

$l_{\text{ce}}(y, \hat{y}) \geq l \log l$   
 $y = [y_1, \dots, y_l] \quad y_i \in \{0, 1\}$

$l_{\text{ce}} = - \sum_{i,y=1}^l \log t_{y_i}$

$l_{\text{ce}} = - \frac{l}{k} \sum_{i,y=1}^k \frac{1}{l} \log t_{y_i}$   
convex fn

$$\sum_{y_i=1} \frac{1}{L} \log f_{y_i} \leq \log \left( \sum_{y_i=1} \frac{1}{L} f_{y_i} \right)$$

$$\sum \lambda_i f(u_i) \leq f(\pi \in \Delta; u_i)$$

$$\geq -L \log \left[ \sum_{y_i=1} \frac{1}{L} f_{y_i} \right]$$

$$= -L \log \frac{1}{L} \sum_{y_i=1} f_{y_i}$$

$$\geq -L \log \frac{1}{L}$$

$$\geq L \log L$$

$$\textcircled{2} \quad \omega^{(n)} = \omega^{(n-1)} - \eta \nabla_{\omega} J(\omega^{(n-1)})$$

Weight decay - multiply each parameter  $\omega$  by  $(1-\epsilon)$   
 Should show this is same as  $L_2$  Reg

$$\omega'^{(n)} \xrightarrow{\rightarrow} (1-\epsilon) \omega^{(n)}$$

$$\omega'^{(n)} = \mapsto (1-\epsilon) \left( \omega^{(n-1)} - \eta \nabla_{\omega} J(\omega^{(n-1)}) \right)$$

$$\omega'^{(n)} = \omega^{(n-1)} - \eta \left[ (1-\epsilon) \nabla_{\omega} J(\omega^{(n-1)}) + \frac{\epsilon}{n} \omega^{(n-1)} \right]$$

gradient of some  
 (new) objective function

To get object fn we need to integrate it

$$J_{new}(\omega) = (1-\epsilon) J(\omega) + \frac{\epsilon}{2n} \|\omega\|^2$$

{Can we take anal order of Taylor series expansion of  $\omega$ ? No  
because  $\omega$  is parameter}

(Q6) curve fitting problem

$(x_i, y_i) \rightarrow$  given

$x_i, y_i \in \mathbb{R}$

linear mapping  $f: x \rightarrow y$

$$y = w_2 + w_1 x$$

without noise :-

objective fn :-

$$L = \frac{1}{N} \sum_{i=1}^N (y_i - w_0 - w_1 x_i)^2$$

$$\tilde{x}_i = x_i + \epsilon_i$$

$$\epsilon_i \sim N(0, \sigma^2)$$

$$L = \frac{1}{N} \sum_{i=1}^N (y_i - w_0 - w_1 \bar{x}_i)^2$$

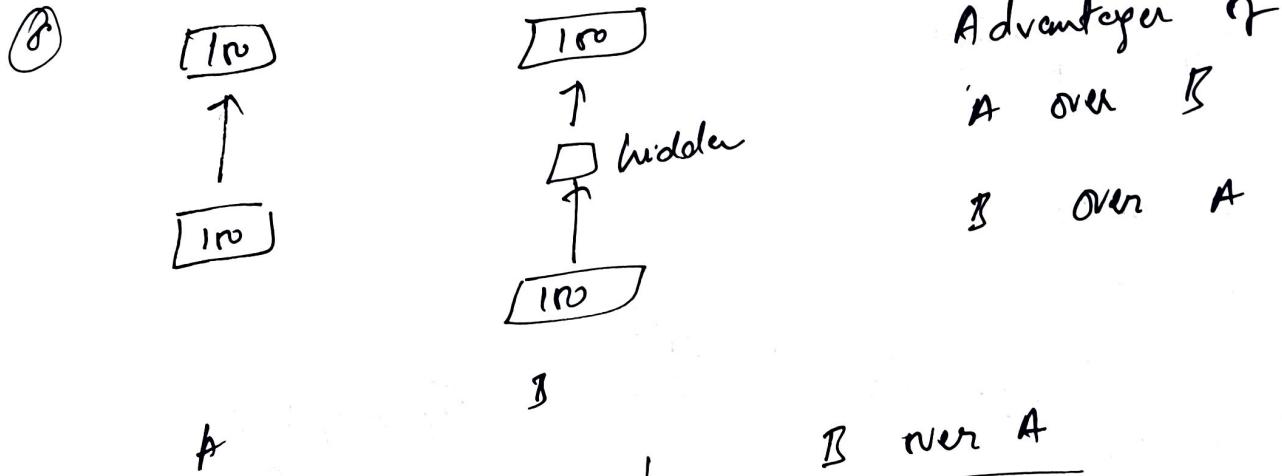
$$= \frac{1}{N} \sum [y_i - w_0 - w_1 x_i - w_1 \epsilon_i]^2$$

$$= \frac{1}{N} \sum [y_i -$$

$$= \frac{1}{N} \sum (y_i - w_0 - w_1 x_i)^2 + \frac{1}{N} w_1^2 \sum \epsilon_i^2 - \frac{2}{N} w_1 \sum (y_i - w_0 - w_1 x_i)$$

$$\tilde{L} = L + \frac{1}{N} \sum_i \omega_i^2 \mathbb{E}[\epsilon_i^2] - \frac{2}{N} \sum_i \mathbb{E}[\epsilon_i] (y_i - w_0 - \omega_i z_i)$$

$$\begin{aligned} \mathbb{E}[\tilde{L}] &= L + \frac{1}{N} \sum_i \omega_i^2 \mathbb{E}[\epsilon_i^2] - \frac{2}{N} \sum_i \mathbb{E}[\epsilon_i] (y_i - w_0 - \omega_i z_i) \\ &= L + \frac{\omega_1^2}{N} \sum_{i=1}^N \alpha_i^2 \\ &\rightarrow L + \omega_1^2 \alpha^2 \quad \rightarrow L_2 \text{ regularization} \end{aligned}$$



- A over B
- ① number of weights are more [disadvantage]  
 $\rightarrow 15000$
  - ② nodes are lesser

B over A

- ① less number of weights.  
 $1000 + 1000$

(5)



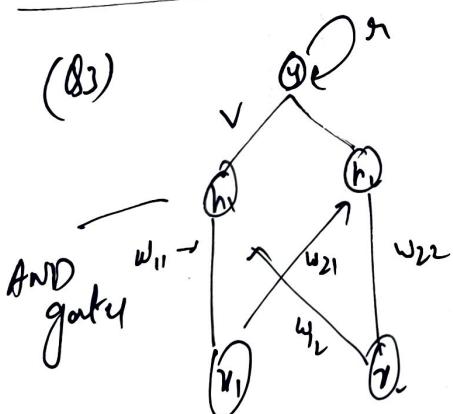
$$\begin{aligned}
 h_t &= u_t - h_{t-1} \\
 &= x_t - (x_{t-1} - h_{t-2}) \\
 &\Rightarrow u_t = x_{t-1} + h_{t-2} \\
 &= u_t - u_{t-1} + x_{t-2} - h_{t-3}
 \end{aligned}$$

What does it do?

$$h_t = \begin{cases} \sum_{j=1}^{\frac{d}{2}} x_{2j} - \sum_{j=1}^{\frac{d}{2}} x_{2j-1} & \text{if } t = \text{even} \\ \sum_{j=1}^{\frac{(t+1)}{2}} x_{2j-1} - \sum_{j=1}^{\frac{(t+1)}{2}} x_{2j} & \text{if } t = \text{odd} \end{cases}$$

$$y_t = \alpha(\log h_t)$$

(Q3)



find parameters of KNN that  
depends if the 2 inputs are  
same or not.

$$i_1 = x_1(1) \quad x_2(1) \quad \dots$$

$$i_2 = x_1(2) \quad x_2(2) \quad \dots$$

$$\text{are } i_1 = i_2$$

done using RNN?

will since they are same if  $y(t-1) = 1$   
and their current bits are same.

$$h^{(t)} = \phi(w u^{-(t)} + b) \quad t > 1$$

$$y^t = \begin{cases} \phi(v^{T-h(t)} + ry^{t-1} + c) \\ \phi(v^{T-h(t)} + c_0) \end{cases} \quad t=1$$

find  $v, r, c, w, b$

$$h_1(t) = u_1(t) \cdot u_2(t)$$

$$h_2(t) = \bar{u}_1(t) \cdot \bar{u}_2(t)$$

$$w = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}$$

$$b = \begin{bmatrix} 2 \\ 0 \end{bmatrix}$$

$$\text{at } t > 1 \quad h_1(t) + h_2(t)$$

$$y(t) = h_1(t) + h_2(t)$$

$$r = \begin{bmatrix} 1 & 1 \end{bmatrix}$$

$$v = \text{threshold} = 1$$

$$c_0 = \text{threshold} = 1$$

$$\text{at } t > 0 \quad y(t) = y(t-1) (h_1(t) + h_2(t))$$

$$r = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$r > 1$$

$$c = -2$$

(Q7) design RNN architecture  
start emitting  $\otimes 0$  once it sees the first 0

$$\text{Input : } \begin{array}{l} 1110101 \\ 1110000 \end{array}$$

Output :  $\otimes$  0 or not

$(h_t)$   $\rightarrow$  Should record if it received 0 or not

Truth table:-

$h_{t-1}$	$x_t$	$h_t$
0	0	1
0	1	0
1	0	1
1	1	1

todays class  
 - encoder  
 - decoder  
 - attention mechanism

601

12/10/23

## Variational Autoencoders

$$\hat{P}_t = \frac{1}{m} \sum_{i=1}^m h(x_i; \theta) \quad l = \text{neuron} \rightarrow \text{spike}$$

$$R(\theta) = \sum P \log \frac{P}{\hat{P}} +$$

$L(\theta)$ : squared error diff between  $x_i$  and output

$$= L(\theta) + R(\theta)$$

$r \rightarrow$  vector

$$\frac{\partial \hat{P}}{\partial w} = \left[ \frac{\partial f_1}{\partial w} \quad \frac{\partial f_2}{\partial w} \quad \dots \quad \frac{\partial f_n}{\partial w} \right]$$

$$\frac{\partial \hat{P}_i}{\partial w_{it}} = \frac{\partial}{\partial w} \left[ \frac{1}{m} \sum_{i=1}^m g(w^T x_i + b_t) \right]$$

$$= \frac{1}{m} \sum \frac{\partial (g(w^T x_i + b_t))}{\partial w_n} = \frac{1}{m} \sum g'(w_i^T x_i + b_t) x_i$$

$$\frac{\partial \hat{P}_i}{\partial w} = x_i [g'(w^T x_i + b)]^T$$

# \* 281 deep learning Stanford - Justin Johnson

19 May 2 202-

Sequence to sequence with RNNs

Self attention

Input  $\rightarrow$  seq:  $x_1, \dots, x_T$

Output  $\rightarrow$  seq:  $y_1, \dots, y_T$

$$\text{encoder: } h_t = f_w [x_t, h_{t-1}]$$

Initial decoder state  $\rightarrow s_0$

content vector  $\rightarrow c$  [often  $c = h_1$ ]

$\rightarrow$  Can't achieve parallelism using LSTM, due to feedback going back to the model.

$\rightarrow c_{tj}$  capturing contribution of  $j^{\text{th}}$  input in predicting  $t^{\text{th}}$  output

$h_i \rightarrow$  encoder state. (hidden state)

$s_t \rightarrow$  hidden state of decoder

$$\text{content vector} = c_t = \sum_{i=1}^T \alpha_{t-i} h_i$$

HK

TK

TIKH

MATH

IR

IR

class 18

16/10/23

Paper's

→ Attention is all you need

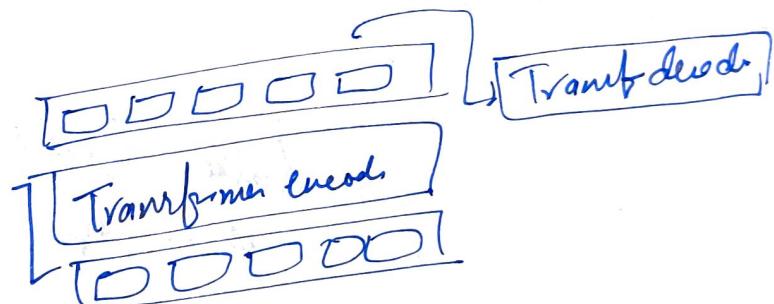
[by google Research] 2017

→ Opportunity and Risk of Foundation model.  
Survey model.

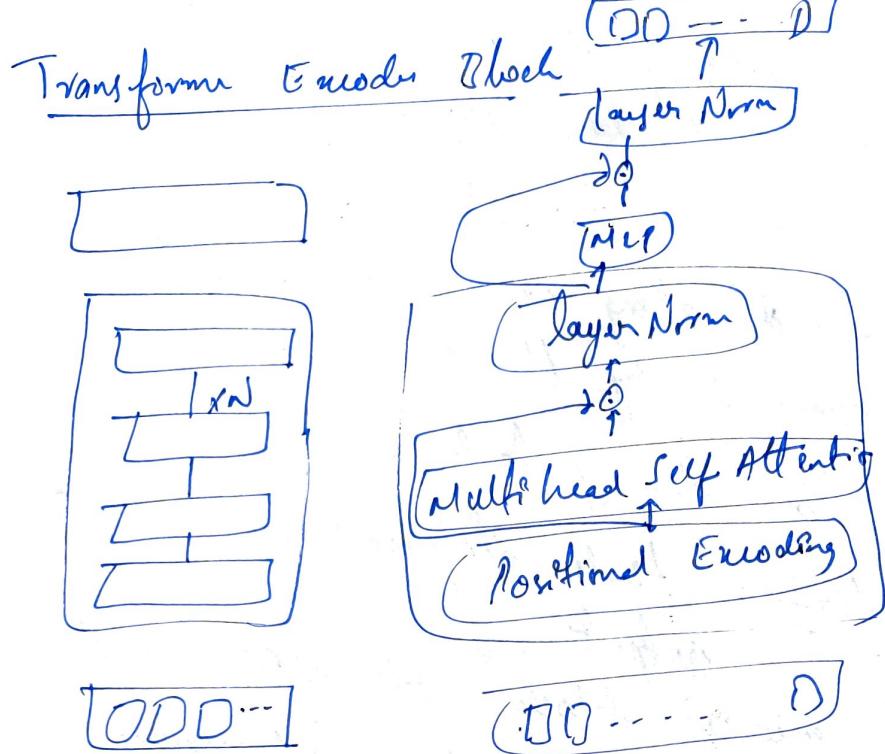
"Image Captioning using Transformer"

Input → Image  
Output → words seq

Image → CNN → features



decoder →  $y_t = \tau$



→ batch norm → helps in invariant that

→ layer norm → helps in normalizing output of all layers, take mean, std hidden layer

$$\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} - a_{i(m)} \quad \mu = \frac{1}{n} \sum_{i=1}^n a_{i(m)}$$

$$a^2 = \frac{1}{n} \sum [a_{i(m)} - \mu]^2$$

no separate parameters like in batch norm here, output normalization is directly done and the forward work is done

Main idea is self attention here.

even though we are doing position encoding, we still do self attention and layer norm

### Decoder block

alignment is found along with each output.  
Violet box is produced by the encoder.

$x_1, x_2, x_3, \dots$  are the outputs.

Multi-head attention → finds alignment between the output

kinda similar to encoder block with violet addition

This works better than ResNet (match performance)

Vision Transformer → clipped ff

GPT's are the decoders of these transformers.

### Latent Semantic Indexing (LSI) / (DFA)

Ideas for giving vector representation for words

words = frequency of occurrence

$$\text{words} = \begin{bmatrix} d_1 & d_2 & \dots & d_n \end{bmatrix}$$

D = document

Document Frequency Matrix

do singular Value decomposition → we get singular value →  
using these singular values we get matrix/vector  
of words

(8) what is one hot representation? — Search on GfYT

SVD doesn't capture relation between words well  
computational cost is high -  $O(nm^2)$

### Word2Vec

associates words to points in the space  
NN — single hidden layer network  
supervised model.

Skip graph — another architecture  
Auto encoder — denoising

### Remaining

CNN variants  
Res Nets.

GANs  
Diffusion models

### CNN

- learning filters that detect different things.
- CNN finds these automatically and finds out and learns output.
- when a neuron minimum fires, when input that passed to highly correlated with the output neuron.

$$x = (x_1, \dots, x_w)$$

Kernel:  $u = (u_1, \dots, u_w)$

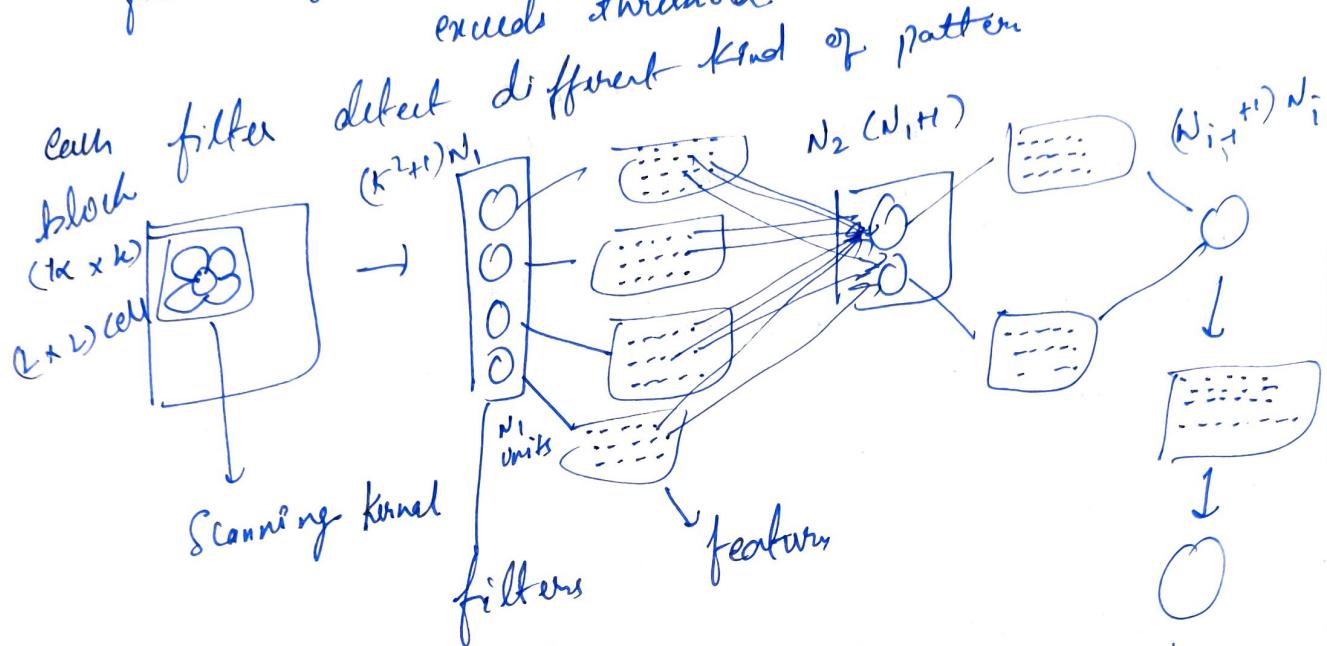
$$(x \otimes u)_i = \sum_{j=1}^w x_{i-1+j} u_j = [x_1, \dots, x_{i+w-1}] u$$

$$(1, 2, 3, 4) * [3, 2] = [3+4, 6+6, 9+8] = [7, 12, 17]$$

26/10/23

CNN's

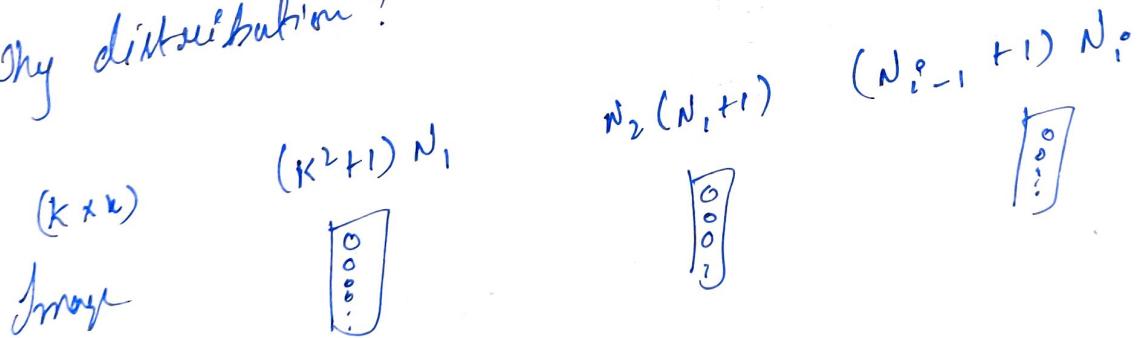
- Neuron fires when a certain pattern is detected
- Correlation operation is done
- fire → yes if correlation between weight and input exceeds threshold



(Q) Why multiple CNN layers?

(A) To decrease calculations as it will have lesser number of parameters. To detect same task

(Q) Why distribution?



Weights  
1st layer =  $N_1(K^2+1)$  weights.

2nd layer =  $\left[\left(\frac{K}{2}\right)^2 N_1 + 1\right] N_2$  weights

Total parameters =  $O\left[1^2 N_1 + \left(\frac{K}{2}\right)^2 N_1 N_2 + N_2 N_3 + \dots\right]$

receptive field = the input that makes neuron fire

2/1/23

## Diffusion Models

generative model

trying to learn distribution of some data

- has a component - Encoder (Decoder)
- takes sample data and map it to series of intermediate latent variables

$$x \rightarrow z_1 \dots z_T$$

sample data

latent variable  $\rightarrow$  noise

Adds noise to image and reconstruct original image

Encoder :-

$$z_1 = \sqrt{1-\beta_1} u + \sqrt{\beta_1} \epsilon_1$$

noise standard normal distribution

$$z_t = \sqrt{1-\beta_t} z_{t-1} + \sqrt{\beta_t} \epsilon_t$$

If  $\beta_t > 1$  then  $\epsilon_t = 0$

$\epsilon_t = 0 \Rightarrow$  mean is 0  
mean of  $z_t$  is more closer to 0, as we are decreasing  
B the mean is getting closer to 0.

Mean of  $z_1 = ? =$  if  $u$  is fixed then distribution  
is gaussian

$$(\sqrt{1-\beta_1}) u = \text{mean}$$

$$\text{Covariance} = (\beta_1, \text{Identity Matrix})$$

Forward eq:-

$N \rightarrow$  Normal distribution

$$q(z_t | x) = N(\sqrt{1-\beta_t} u, \beta_t I)$$

$$q(z_t | z_{t-1}) = N[\sqrt{1-\beta_t} z_{t-1}, \beta_t I]$$

This forms markov chain because  $P(z_t)$  depends only on value of immediate preceding variable  $z_{t-1}$  with sufficient steps in  $T$ ;  $q(z_T | x) = q(z_T)$

Become Standard Normal distribution

Joint distribution of  $z_1 - z_T$

$$q(z_1, z_2, \dots, z_T | x) = q[$$

$$z_1 = \sqrt{1-\beta_1} u + \sqrt{\beta_1} \epsilon_1 \quad \text{--- ①}$$

$$z_2 = \sqrt{1-\beta_2} z_1 + \sqrt{\beta_2} \epsilon_2 \quad \text{--- ②}$$

= Replace ① in ②

$$z_2 = \sqrt{1-\beta_2} [\sqrt{1-\beta_1} u + \sqrt{\beta_1} \epsilon_1] + \sqrt{\beta_2} \epsilon_2$$

$$z_2 = \sqrt{(1-\beta_2)(1-\beta_1)} u + \sqrt{\beta_1(1-\beta_2)} \epsilon_1 + \sqrt{\beta_2} \epsilon_2$$

$$\boxed{E[z_2] = \sqrt{(1-\beta_2)(1-\beta_1)} u = \sqrt{\alpha_2} u = E(z_2)}$$

$$E[z_2^2] = \sqrt{(1-\beta_2)(1-\beta_1)} u^2 + \beta_1 \beta_2 + \beta_2^2 = \beta_1 \beta_2 + \beta_2^2$$

$$\text{Var}(z_2) = \beta_1 \beta_2 + \beta_2^2 - (\beta_1 \beta_2 + \beta_2^2) = \beta_2(1-\beta_1)$$

$$\text{Var}(z_2) = \beta_1 + \beta_2 - \beta_1 \beta_2$$

$$\boxed{\text{Var}(z_2) = 1 - \alpha_2}$$

## Marginal and conditional Distribution

$$q(x_t) = \int q(\theta_t | \alpha) q(\alpha) d\alpha \quad : \text{Marginal}$$

$$q(\theta_{t-1} | \theta_t) : q(\theta_{t-1} | z_t) = \frac{q(\theta_t | \theta_{t-1}) q(\theta_{t-1})}{q(\theta_t)}$$

conditional → Intractable  
done by sampling as can't compute marginal distribution of  $q(\theta_{t-1})$

useful when training decoder.

couldn't find reference for page 12  $\rightarrow$  13.

Page 13, ① conditional distribution is used in marginal.

Pg 12  
Simplification of log term:-

$$\log \left[ \frac{P(x, \theta_1, \dots, \theta_T | \phi_1, \dots, \phi_T)}{q(\theta_1, \dots, \theta_T | x)} \right]$$

$$= \log \left[ \frac{P(x | \theta_1, \phi_1) \prod P(\theta_{t+1} | \theta_t, \phi_t) P(\theta_T)}{q(\theta_1 | x) \prod_{t=2}^T q(\theta_t | \theta_{t-1})} \right]$$

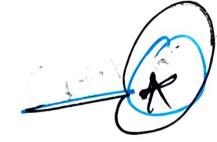
$$= \log \left[ \frac{P(x | \theta_1, \phi_1)}{q(\theta_1 | x)} \right] + \log \left[ \frac{\prod P(\theta_{t+1} | \theta_t, \phi_t)}{\prod q(\theta_t | \theta_{t-1})} \right] + \log P(\theta_T)$$

↓  
Reverse Probability.

we observe that

$$q(z_t | z_{t-1}) = q(z_t | z_{t-1}, u)$$

$$q(z_t | z_{t-1}) = \frac{q(z_{t-1} | z_t, u) q(z_t | u)}{q(z_{t-1} | u)}$$



Using ① in ① :

$$\log \left[ \frac{p(x, z_{1:T}) / \theta_{1:T}}{q(z_{1:T} | u)} \right] = \log \left[ \frac{p(u | z_1, \phi_1)}{q(u | u)} \right] + \log \left[ \frac{\pi p(z_{t-1} | z_t, \phi_t)}{\pi p(z_{t-1} | u)} \right]$$
$$+ \left[ \frac{q(z_{t-1} | u)}{q(z_t | u)} \right] \log \frac{p(z_t | z_{t-1}, \phi_t)}{p(z_t)}$$
$$= \log p(x | z_1, \phi_1) + \log \frac{\pi p(z_{t-1} | z_t, \phi_t)}{\pi q(z_{t-1} | u)} + \log \frac{p(z_t)}{q(z_t | u)}$$

Keep this out of  
optimization  
as its in  
normal diff.

$$\text{ELBO}(\theta_{1:T}) = \int q(z_{1:T} | u) \log \frac{p(x, z_{1:T} | \theta_{1:T})}{q(z_{1:T} | u)} dz$$
$$= \int q(z_{1:T} | u) \left\{ \underbrace{\log p(x | z_1, \phi_1)}_{\text{(this is conditional fn of } q(z_1 | z_1))} + \underbrace{\log \frac{\pi p(z_{t-1} | z_t, \phi_t)}{\pi q(z_{t-1} | z_t, \phi_t)}}_{dz} \right\}$$
$$= \int q(z_{1:T} | u) \left[ \log p(x | z_1, \phi_1) \right] + \int q(z_{1:T} | u) \sum \log \frac{p(z_{t-1} | z_t, \phi_t)}{q(z_{t-1} | z_t, \phi_t)} dz$$

$$= E_{q(z_1|z_T)} \left[ \log P(x|z_1, \theta) \right] + \sum \int q(z_1 \dots z_T | u)$$

$$\log \frac{P(z_{t-1}|z_t, \theta_b)}{q(z_{t-1}|z_t, u)} dz_1 \dots z_T$$

$$= E_{q(z_1|u)} \left[ \log \frac{P(x|z_1, \theta_b)}{q(z_1|z_T, u)} \right] - \sum \int q(z_1 \dots z_T | u)$$

$$\log \frac{q(z_{t-1}|z_t, x)}{P(z_{t-1}|z_t, \theta)} dz_1 \dots z_T$$

$$= E_{q(z_1|u)} \left[ \log \frac{P(x|z_1, \theta_b)}{q(z_1|z_T, u)} \right] + \sum \int q(z_{t-1}|z_t, u)$$

$$\log \frac{q(z_{t-1}|z_t, u)}{P(z_{t-1}|z_t, \theta)} dz_1 \dots z_T$$

thus is KL divergence between 2 elements

$$= E_{q(z_1|u)} \log P(x|z_1, \theta_b) - \sum \text{KL} (q(z_{t-1}|z_t, u) \| P(z_{t-1}|z_t, \theta_b))$$

16/11/23

V.A-E

Auto encoder

Encoder

Decoder

$$x \rightarrow z$$

$$z \rightarrow x$$

Variational Auto encoder

$$x \rightarrow P[z|x]$$

$$z \rightarrow P(x|z)$$

"Skipping relationships pg 12 Lecture 21 pdf"

$\hat{u}$  [reconstruction]  
↑  
Encoder

Decoder  $P_\theta(x|z)$

Latent Variable are distributed  
as Standard Normal  
distribution.

Encoder  $q_\phi(z|x)$

$x$  (data)

$$z \sim N[\mu(\theta, x), \Sigma(\theta, x)]$$

Parameters of  
multivariate Normal  
distribution.

try to enforce the  
be same as standard  
Normal distribution by  
minimizing KL divergence.

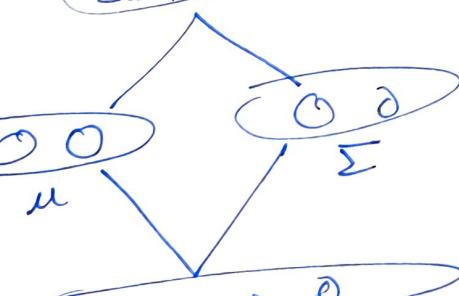
$0 0 0 0$   $P_\theta(x|z)$

Sample

Sample  $z$

$\theta$  are the parameters  
of encoder network

$\phi$  parameters of  
decoder network



$q_\phi(z|x)$

$x_t$

encoder and decoder  
both are neural  
networks.

object fn of decoder  $\rightarrow$  maximize  $P(x_i|z)$   $\rightarrow$  likelihood Marginalize

$$P(x_i) = \int P(z) P(x_i|z) dz$$

$$= + \mathbb{E}_{z \sim Q_\phi(z|x_i)} [\log P_\theta(x_i|z)]$$

\* Find the divergence between 2 distributions Pg 29 Lee 21  
(Q) Find the divergence between 2 distributions (in exam)

VAE Approximate it using a single sample.  
using  $\frac{1}{N} \sum_{n=1}^N f(n)$   $\rightarrow$  generally

$$\mathbb{E}[f(n)] \approx f(n) \rightarrow \text{In VAE}$$

finding expectation is not tractable.  
distribution that's outputted by encoder.  
decoder network predicts mean with  
Identity covariance matrix.

find log of it  
no need to find expectation as we are  
Sampling with sample  $f(n) \rightarrow f(n)$

$$\log [P(x_i|z)] = C - \frac{1}{2} \|x_i - \mu(z)\|^2$$

(gaussian distribution)

output of  
decoder network  
(mean of distribution)  
(can be ignored)

Overall obj of VAE = pg 31 lec 21

Can update parameters using backprop.  
(wrt  $\theta, \phi$ )

encoder      decoder.

Problem is not end to end differentiable  
output of decoder is not  
difficult to implement backprop. as chain rule is  
hard in this case.

hence  $f_\phi^{(z)}$  is non deterministic.  
 $\xrightarrow{\text{Random for } z}$  Random for  $x$   
not continuous for  $x$

de parameterization trick

$\phi$  is also multivariate gaussian distribution :- If we  
Introducing new param or gaussian  
Want to generate any sample  
[multiply by Standard dev]  $\times$  [ ]  
(Shift)

$$\leftarrow \rightarrow N(\mu, \sigma^2)$$

$$x \rightarrow N(\mu, \sigma^2)$$

$$u = \mu + \sigma^2$$

Inputs  $\rightarrow$  data points.

Sampling Step  $\rightarrow$   $\epsilon \in$  Sampled  $\mathcal{Z}$  that can be mapped to sketch

$$[\mu(\mathbf{x}) + \Sigma(\mathbf{x})\epsilon]$$

"This can be done at Sampling".

$$\epsilon \sim N(0, I)$$

$$\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\Rightarrow \mathbf{x} = \boldsymbol{\mu} + \boldsymbol{\Sigma}^{1/2}\epsilon$$

$\rightarrow$  after updating Sampled  $\mathcal{Z}$ , the backprop will be done easily.

$\rightarrow$  Randomness is now associated to  $\epsilon$  not  $\mathbf{x}$ .  
This has the problem is taken care of.

Classes done