

* Convolutional Neural Network (CNN) (P.S.)

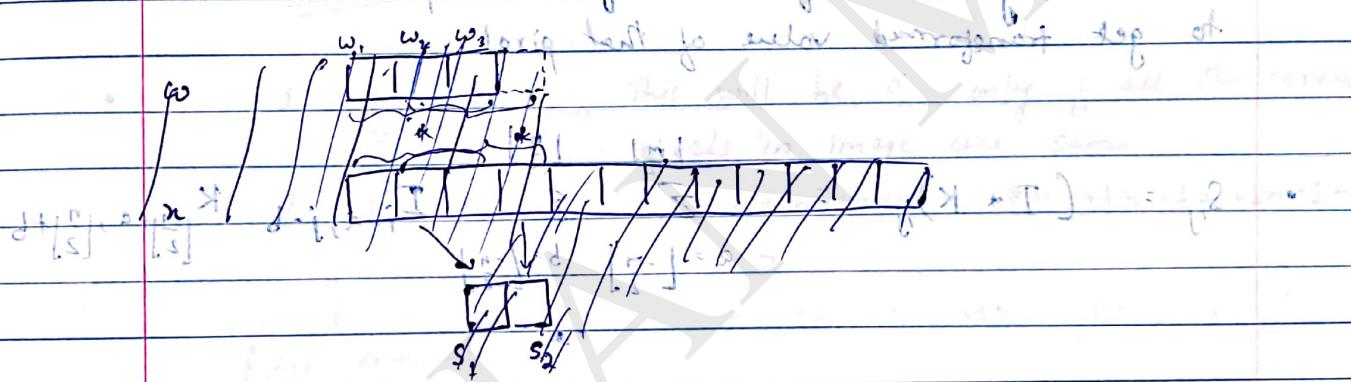
→ 1D CNN Convolution

$s_t = \sum_{a=0}^{n-1} x_{t-a} w_a$ where $w_a = (x_a * w)_a$

size of filter = 3 - 3x3 kernel

→ lets say $n = 3$

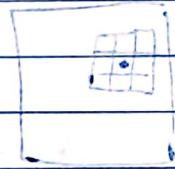
Let's say we have an image of size 3×3 pixels with values x_1, x_2, x_3



→ neighborhood of location for value x_1

→ neighborhood of location for value x_2

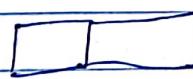
→ 2D CNN Convolution w_1, w_2, w_3



w

Now we have depth of max 1 channel as same

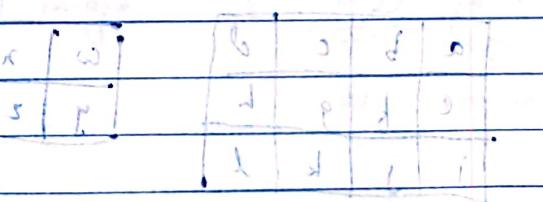
for 3x3x3 input $(x_{i,j})$ if we will multiply $(x_{i,j})$ with w_i then we get s_i



for 3x3x3 input $(x_{i,j})$ if we will multiply $(x_{i,j})$ with w_i then we get s_i

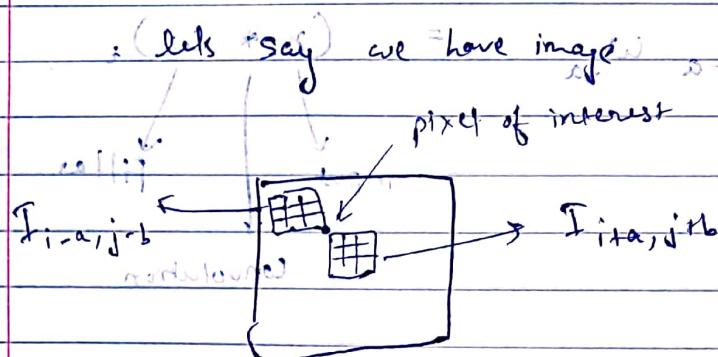
$$s_1 = x_1 w_1 + x_2 w_2 + x_3 w_3$$

$$s_2 = x_4 w_1 + x_5 w_2 + x_6 w_3$$



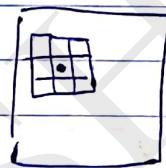
→ 2D Convolutions (With Examples)

$$S_{ij} = (I * K)_{ij} = \sum_{a=0}^{m-1} \sum_{b=0}^{n-1} I_{i+a, j+b} K_{a,b}$$



We multiply matrix after or before the pixel with filter to get transformed value of that pixel.

$$S_{ij} = (I * K)_{ij} = \sum_{a=-\frac{m}{2}}^{\frac{m}{2}} \sum_{b=-\frac{n}{2}}^{\frac{n}{2}} I_{i-a, j-b} K_{\left[\frac{m}{2}\right]+a, \left[\frac{n}{2}\right]+b}$$



get value of pixel by multiplying matrix of pixels around it by filter

Note In above cases filters size is $(m \times n)$ i.e. size of K is $(m \times n)$

a	b	c	d
e	f	g	h
i	j	k	l

w	x
y	z

\Rightarrow

$aw + bx$	$bw + cx$	$cw + dx$
$+ey + fz$	$fy + gz$	$gy + hz$
$ew + fw$	$fw + gx$	$gw + hx$
$+iy + jz$	$+jy + kz$	$ky + lz$

→ Types (of) Filters (size, weight, size-weight) methods *

-  replacing this pixel in image with average of all pixels in window as $\frac{1+1+1+1+1+1+1+1+1}{9}$

ii [↑] Bilir endocrinits over the spine beroefend ist

- $0 \text{ } (-1 \text{ } 4) \text{ } 0, \omega \rightarrow$ central pixel is made 5 times & all other pixels around it are subtracted

Sharpen for size roughly at bottom at $r = 9$ mm
top plane. The wall is covered with almost wavy

- $\begin{matrix} 1 & 1 & 1 \\ 1 & -8 & 1 \\ 1 & 1 & 1 \end{matrix}$ This will be 0 only if all the corresponding pixels in image are same

Edge detector

Page .

$$\text{Current } I = \frac{V - 0.75}{R} = \frac{V - 0.75}{1000} = 10^{-3}V$$

\rightarrow 3D Convolution $\leftarrow + \frac{7-96+16}{?} = 48$

Assume that depth of image & kernel are same - (1)

with assumption ①

~~assignment~~ (Exhibit) * Exhibit

If we apply 3D filter to 3D image a 3D feature map

— map
is obtained

refer /notes/ week 11. / 3D conv. pyg / also ref. part

* Relation (input size, output size, filter size) (logpt)

→ Shows this goes on taking out patches.

→ Let an image have dimensions $W_1 \times H_1 \times D$, filter have dimensions $F \times F$.

The transformed image will have dimensions $W_2 \times H_2 \times D$.

$$\text{Let } W_2 = W_1 + F - 1 \Rightarrow W_2 = W_1 - (F-1)$$

$$\text{Let } H_2 = H_1 + F - 1 \Rightarrow H_2 = H_1 - (F-1)$$

→ In ^{almost} _{all} cases we might have to add padding (^{Transform all pixels by taking them at ~~very less cost~~ _{nearby cost}})
 $P=1$ is added to all four side of input matrix.

The formula then becomes as follows

$$W_2 = W_1 + 2P - F + 1$$

$$H_2 = H_1 + 2P - F + 1$$

→ We may keep a stride as per output dimensions needed

reduces computation

$$W_2 = \frac{W_1 + 2P - F}{S} + 1$$

Main formula

$$H_2 = \frac{H_1 + 2P - F}{S} + 1$$

most common form for stride than $S=1$

(Q)

$$227 \times 227 \times 3 \times (11 \times 11 \times 3) \Rightarrow 9$$

no padding

Input images of size 96 filter will be of size 96 output.

$$W_2 \times H_2 \times D$$

Ans Since filter depth and img depth are same

Thus for each ~~filter~~ ^{1x1x32 before map} 2D feature map is obtained

Since there are 96 filters (in working)

$$D = 96$$

→ given $s=4$, $p=0$

$$W_2 = \frac{227 - 11}{4} + 1 = 55$$

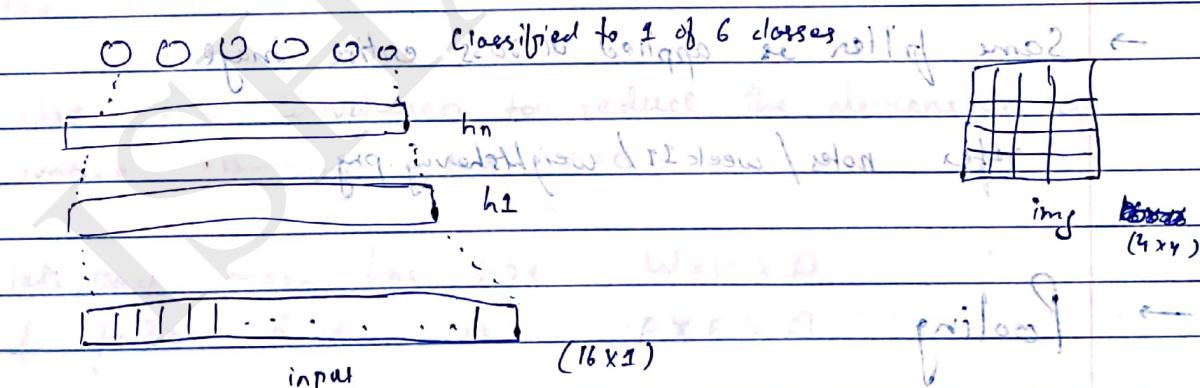
$$H_2 = \frac{227 - 11}{4} + 1 = 55$$

* CNNs

ML approach → Till now we focused on applying predefined static kernels/filters on image. Then we can learn weights to classify.

DL approach → Deep learning focuses on learning these filters. Then we learn weights to classify.

→ For Feed Forward NN



It is a dense connection and there are a lot of parameters since each neuron in previous layer contributes to formation of neuron in next layer.

Principle of backpropagation

→ Convolution NN ~~using 3x3 and result 3x3~~

2	4	7
2	5	8
3	6	9

⇒

1	2	3	4	5	6	7	8	9
10	11	12	13	14	15	16	17	18

3x3 img.

$$3P = 9$$

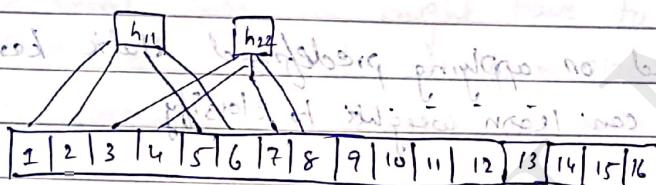
9x1 vector

(1)

so this is comparatively sparse connection

$$P + N = 55$$

1	5	9	13
2	6	10	14
3	7	11	15
4	8	12	16



work well, easily work forward as subset of principal cost

(2)

Weight Sharing

→ Same filter is applied across entire image

refer notes / week 11 / weightsharing.png

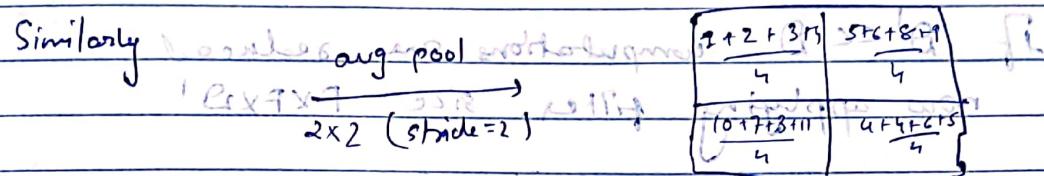
→ Pooling

Used to reduce dimension of image

1	3	5	8
2	4	6	9
10	7	4	6
11	3	5	4

4	9
6	7

feature map obtained after applying filter on img



* Image Net

→ AlexNet (with 5 layers and 3 pool)

refer notes / week 11 / AlexNet.pptx

→ ZFNet

refer notes / week 11 / ZFNet.pptx

→ VGGNet

refer notes / week 11 / VGGNet.pptx

→ GoogleNet

- Use $1 \times 1 \times D$ convolutions to reduce the dimension of

image assuming $P=0$

(i) let say image has size $W \times H \times D$

& filter size is $F \times F \times D$

$$\text{Total Computations required} = F \times F \times D$$

multiplied by each element of o/p

(number of nodes)

(ii) if we use $1 \times 1 \times D$ convolutions & D' such

convolutions then

$$(W \times H \times D) * (1 \times 1 \times D) \Rightarrow W \times H \times D'$$

If $D' \ll D$, computations are reduced.
now applying filter size $F \times F \times D'$

Total computation for each element of o/p = $F \times F \times D'$ significantly reduced

- GoogLeNet has various inception modules having basic layout as follows

refer notes / week 11 / inceptionmodule.png

other inception module structures

refer referenceppt / 11 / slides 423 - 442

- In the last convolution layer

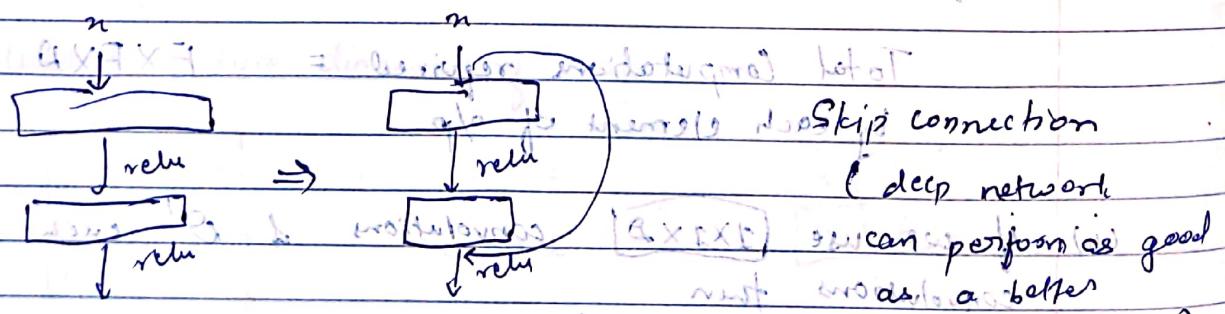
$7 \times 7 \times 1024 \xrightarrow{\text{avg-pool}} 1 \times 1 \times 1024$

equivalent to a vector

all neurons in a layer can be fully connected with all

dense layer with less parameters than earlier

\rightarrow ResNet $\xrightarrow{\text{refer / notes / week 12 / resnet.png}}$



$$H(n) \leftarrow H(n) + (D(n) + S(n))$$

In residual network, in forward pass current layer that has finer features is aggregated with previous layers [2 - 3] layers before the current layer is aggregated]

Generally during backpropagation if network is huge / contains many layers the gradient vanishes during backprop. But if we have residual connection it doesn't allow the gradient to vanish and thus relevant update is made to the weight & the network learns.

$$(x_i w + b_i) \text{ where } w = ?$$

using the same gradient of both old and new weight

$$(d + i w + j z) \text{ where } ?$$

reduced
grad

reduced
grad

gradient learning rate

$$\alpha = 12$$

residual function is 12 weight added at each

18 → 30

work

$$IV \xrightarrow{\frac{12}{12+3}} = (V \times 12) A$$

WEEK 12

Dimensionality reduction *

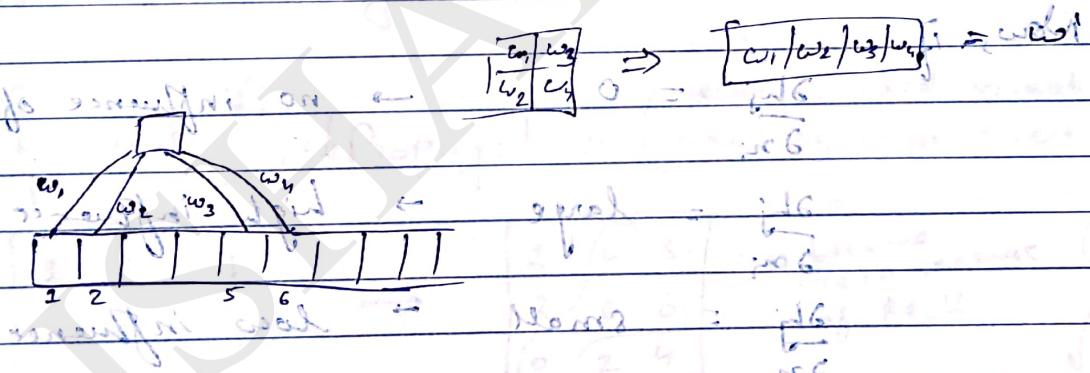
* Visualizing Patches

→ given a neuron in any layer one can always find a patch in the image which is responsible for that neuron to fire.

* Visualizing Filters of a CNN

→ refer autoencoders deep learning (2) how to visualize input using weights from neuron to fire with it.

→ limited to the neurons used certain weights to determine the input & thus determine which patch of input is responsible for a neuron to fire for that.



$$[z_1 \ z_2 \ z_3 \ z_4] = w'$$

thus for all w_1, w_2, w_3, w_4 we get

Thus for these values of input the given neuron will fire.

* Occlusion Experiments

→ Take a patch in image & replace it by black box.
~~now try to classify~~ ~~patch~~ ~~classification~~

If you get correct classification, that patch doesn't affect
class. If it has a very less impact on class, it is doing a

Similarly do it for all patches to understand important
patches for classification. To ~~classification~~

* Finding influence of Input Pixels using Backprop

→ Gradients with respect to input give the influence
lets say h_j is no. neuron of hidden layer j

Now, if $\frac{\partial h_j}{\partial x_i} = 0$ → no influence of x_i on h_j

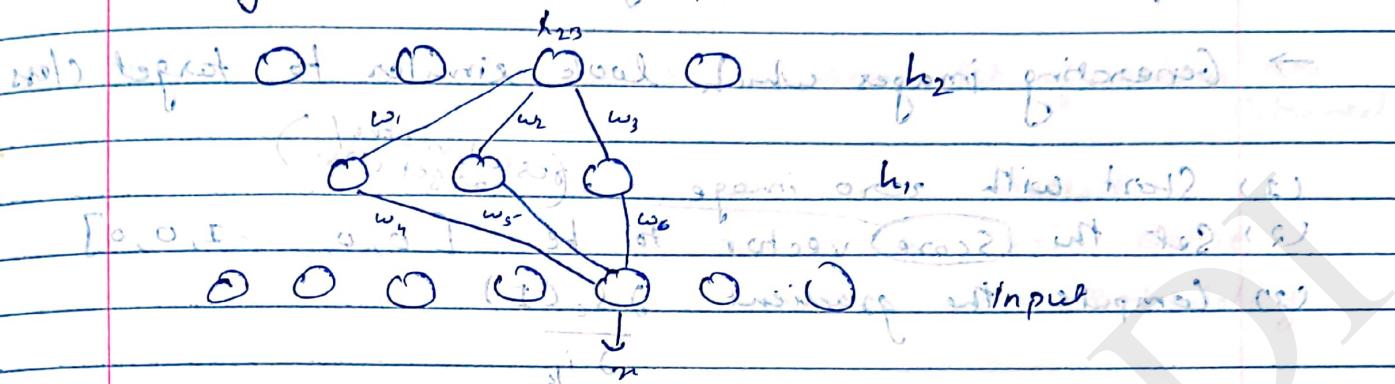
$$\frac{\partial h_j}{\partial x_i} = \text{large} \rightarrow \text{high influence}$$

$$\frac{\partial h_j}{\partial x_i} = \text{small} \rightarrow \text{low influence.}$$

So we represent image as gradients

$\frac{\partial h_j}{\partial x_0}$	$\frac{\partial h_j}{\partial x_1}$	\dots	$\frac{\partial h_j}{\partial x_m}$
.	.	.	.
.	.	.	.
.	.	.	.

* Finding gradient wrt input w.r.t. maximality *



so we want to find $\frac{\partial h_{23}}{\partial x}$ using topology *

$$\begin{aligned} \frac{\partial h_{23}}{\partial x} &= \frac{\partial h_{23}}{\partial h_{11}} \frac{\partial h_{11}}{\partial x} + \frac{\partial h_{23}}{\partial h_{12}} \frac{\partial h_{12}}{\partial x} + \frac{\partial h_{23}}{\partial h_{13}} \frac{\partial h_{13}}{\partial x} \\ &= w_1 w_3 + w_2 w_5 + w_3 w_6 \end{aligned}$$

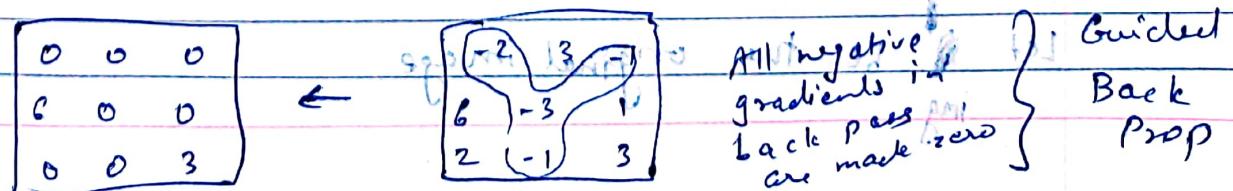
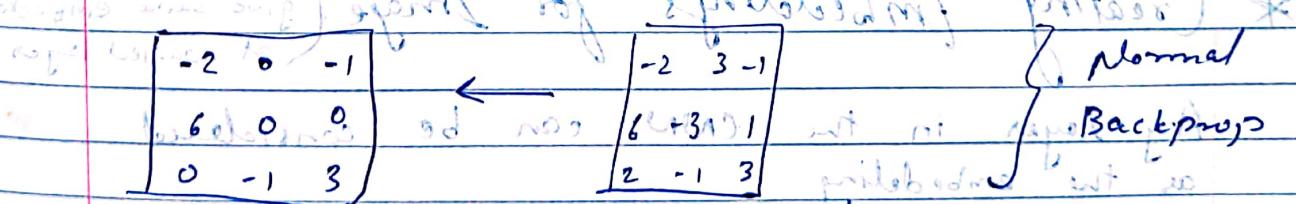
Let's say a test case spans across all nodes ←

Now, plot ~~isolate~~ calculated similarly for all input & plotting

order of grad (1)

what is various backprop technique used visualizing the

* Guided BackPropagation (impact of image pixels on result)



* Optimization Over Images (using gradient)

→ Generating images which look similar to target class

(1) Start with zero image $\xrightarrow{\text{result vector / target class}}$

(2) Set the Score vector to be $[0, 0, \dots, 1, 0, 0]$

(3) Compute the gradient $\frac{\partial S_c(I)}{\partial i_k}$

(4) Update pixel $i_k = i_k + \eta \frac{\partial S_c(I)}{\partial i_k}$

(5) Now again do forward pass

(6) Go to step 2

→ We can also create image such that a particular neuron of any hidden layer fires. (with help of backprop)

(1) Feed img to network

(2) Make activation such that selected neuron is one others are zero (to trigger mitigation of belief)

(3) Backprop to image

(4) $i_k = i_k + \eta \frac{\partial A(j)}{\partial i_k}$, if $A(j)$ is the activation of j^{th} neuron in some layer

* Creating embeddings for Image

Any layer in the CNN can be considered as the embedding

Let img be the original image

Let img_{embed} be the embeded image

0	0	0
0	0	0
0	0	0

We pass this image through the Network or

At the fully connected layer some set of values are obtained



These values are called embedding values and entire set of values is one embedding

$$\text{value} = (33)$$

Now we take a random image and try to reconstruct it such that when it is passed through network the corresponding Fully connected (FC) layer embedding is same.

loss

$$L(c_i) = \| \phi(u) - \phi_o \|_2^2 + \lambda \| \phi(u) \|_1$$

update rule

$$\frac{\partial L}{\partial u} = \eta \frac{\partial \phi(u)}{\partial u}$$

This abstraction of image can be constructed using image \rightarrow image embeddings.

Similarity is possible for convolution layers as well.

* Deep Dream

→ Here we make our loss function & update rule such that ~~more~~ neurons can fire more.

So what happens if a neuron in a convolution layer,

we want to maximize $L(I)$

then we want to maximize $L(I) = h_{ij}$

During backprop. update pixel of image

function of grad loss against gradient of output of well

gradient $\frac{\partial L(I)}{\partial I}$ and $\frac{\partial L(I)}{\partial h_{ij}}$ if h_{ij} has been done in

update rule

$$\| \text{gradient} \| = \lambda \text{ in } \frac{\partial L(I)}{\partial h_{ij}} \text{ where } \lambda = (.)$$

→ Result of this will be as follows:

During training it might have developed relation that if one neuron fires for sky, other neurons might fire for bird, castle, etc. without loss of generality

Thus when we give image of sky and increase the probability of neuron denoting sky to fire. Then it automatically takes knowledge from training data & generates birds, castle as foreground.

thus makes prediction such as sky with a castle being most probable next best

* Deep Art

referred real lost?

→ Making an original image = a cartoon type i.e. styling the image so it looks like cartoon

(1) Content Target

$$L_{\text{content}}(\vec{p}, \vec{n}) = \sum_{ijk} (\vec{p}_{ijk} - \vec{n}_{ijk})^2$$

pixels of embeddings
of actual image pixels of embeddings
of new image.

prob. of pixels + size factor

New image is generated such that all image representation / of assigned ~~are~~ ~~new~~ embeddings generated using new image are similar to embeddings generated by actual image

(2) Style (cartoon)

According to CV study $V \in \mathbb{R}^{64 \times (288 \times 256)}$

computer vision

$V^T V \in \mathbb{R}^{64 \times 64}$ (gives style of image)

$$\therefore L_{\text{style}}(\vec{a}, \vec{n}) = \sum_{l=0}^L w_l E_l \quad \left. \begin{array}{l} \text{or want the style} \\ \text{to match for} \\ \text{all embeddings} \\ \text{/ layers} \end{array} \right\}$$

Here $E_l = \sum_{ij} (G_j^l - A_{ij}^l)^2$

style gram matrix of original image style gram matrix of generated image

[style matrix generated from embedding obtained from original image] [style matrix generated from embedding obtained from new image]

→ Total loss function

total loss

$$\text{envelope } L_{\text{total}}(\vec{p}, \vec{a}, \vec{n}) = \alpha_i L_{\text{content}}(\vec{p}, \vec{n}) + \beta L_{\text{style}}(\vec{a}, \vec{n})$$

refer notes/week 12/depict.py
 $(\vec{p}^* - \vec{p})^2 + (\vec{n}^* - \vec{n})^2$

* Fooling CNNs.

→ Change the target class value & backpropagate & make changes to original image such that you get minimum loss.

(gradient) $\partial \text{loss} / \partial \vec{v}$

$\vec{v} \rightarrow \text{plant vs. gridwall}$

(gradient for alpha exp) $\vec{v} \rightarrow v^T \vec{v}$

alpha exp

more natural

$$(C_{\vec{v}}^* - \vec{v})^2 = \vec{v}^T \vec{v}$$

more natural

more natural

more natural

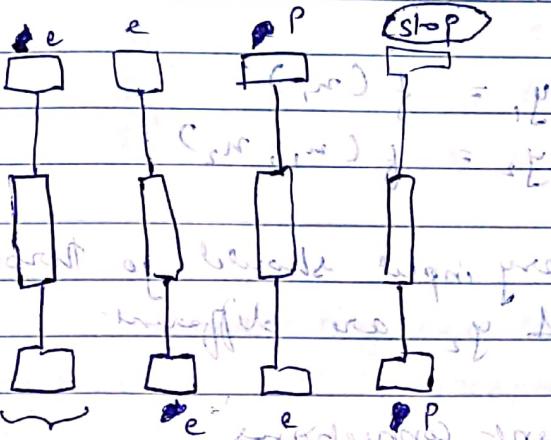
more natural

[most] lossening natural style
most unnatural

[least] lossening natural style
most natural

WEEK 13

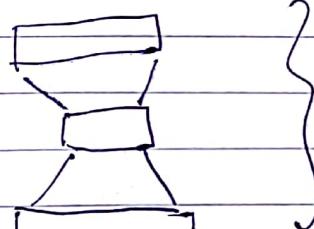
* Sequence Learning Problems



Predicting next character in a word depends on previous characters

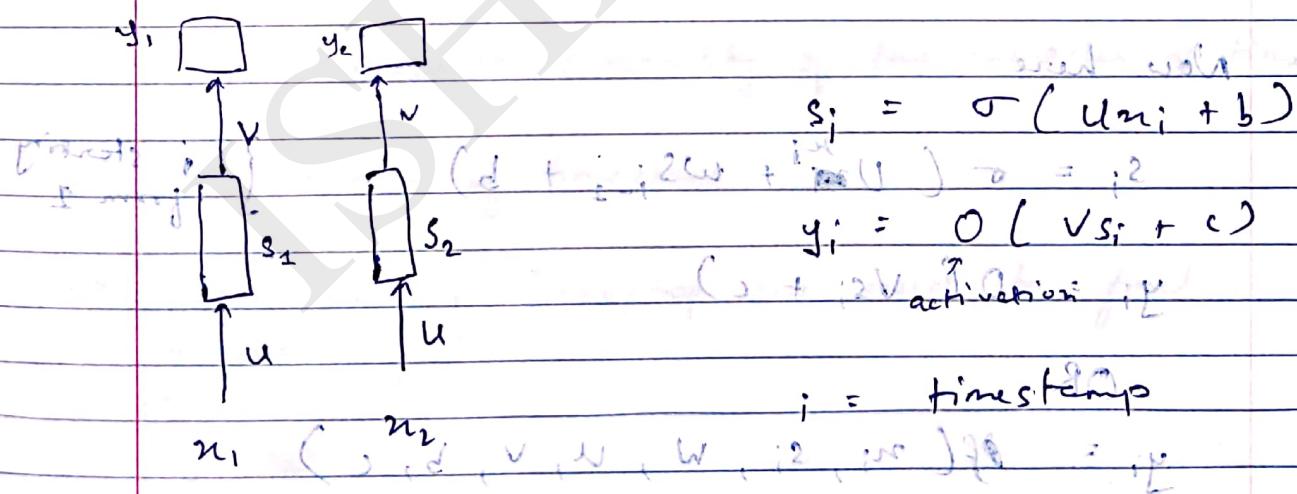
most difficult of these problems is predicting next character in a sequence

difficulty lies in capturing dependencies between elements in sequence



Each one of them is a feed forward NN.

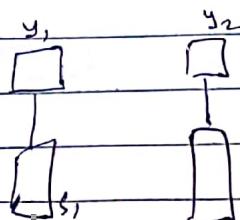
* Recurrent Neural Networks (RNN)



But this doesn't account for previous inputs

refer notes/week 13/wishlist.py for solution

constraint: No extra layers were used. v, w, b optional



can we connect it like this?
The

function is

$y_1 = f(n_1)$

no change here

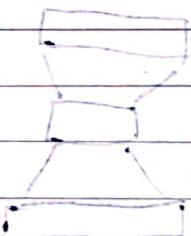
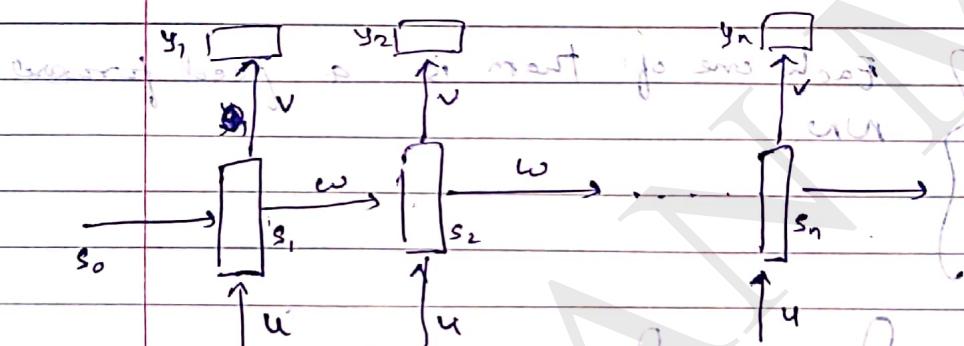
now state s_2

$$y_1 = f(n_1)$$

$$y_2 = f(n_1, n_2)$$

As per wish list every input should go through some function but y_1 & y_2 are different.

→ Thus we make recurrent connections



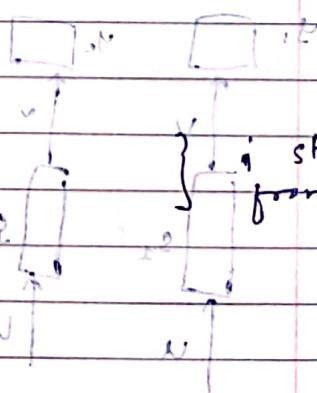
Now here

$$(1 + \text{bias}) = 0$$

$$s_i = \sigma(Vu_i + ws_{i-1} + b)$$

$$(1 + \text{bias}) = 1$$

$$y_i = \text{O}(Vs_i + c)$$



OR write

$$y_i = f(n_i, s_i, w, u, v, b, c)$$

Note

always assume n_i has bias b

s_i is state of network at timestamp i

Weights w , u , v and biases are shared across all timestamps

* Back Propagation Through Time (BPTT) (i)

lets say, loss = $(0), 16 \rightarrow \frac{1}{\sqrt{6}} = (0), 16$

$$z_t \in \mathbb{R}^n$$

$$s_t \in \mathbb{R}^k$$

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

$$U \in \mathbb{R}^{n \times d}$$

$$W \in \mathbb{R}^{d \times k}$$

$$V \in \mathbb{R}^{d \times k}$$

(W here denotes) \Rightarrow

- Example is predicting next character in word using RNN.
refer notes / week 13 / backprop RNN .py.

- Activation at output layer = softmax

$(0), 16 \rightarrow (0), 16 \rightarrow (0), 16 \rightarrow (0), 16 \rightarrow (0), 16$

$$L_t(\theta) = -\log(y_{t,c}) \quad \left. \begin{array}{l} \text{cross entropy at} \\ \text{each time stamp} \end{array} \right\}$$

$$L(\theta) = \sum_{t=1}^T L_t(\theta) \quad \left. \begin{array}{l} \text{total loss} \\ \text{with} \end{array} \right\}$$

$y_{t,c}$ = predicted probability of true character at time-step;

T = number of timesteps

- Here during backprop we need to find gradients for V, W, U (at each time step).

$$V, W, U \in \mathbb{R}^{d \times k} \quad \text{and } \theta \in \mathbb{R}^m$$

now, $\frac{\partial L}{\partial \theta_i} = \frac{\partial L}{\partial y_{t,c}} \cdot \frac{\partial y_{t,c}}{\partial \theta_i}$ with θ

now, $\frac{\partial L}{\partial y_{t,c}} = \frac{\partial}{\partial y_{t,c}} (-\log(y_{t,c})) = \frac{1}{y_{t,c}}$ for $t=1$

(1) Gradients for Viterbi's most probable path *

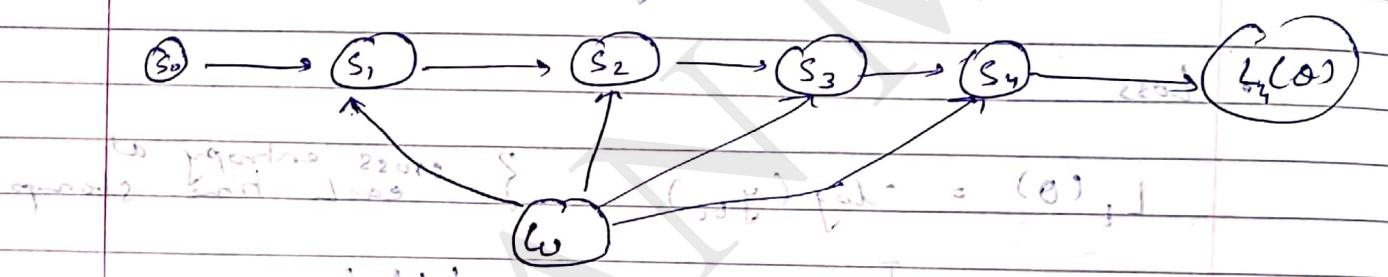
$$\frac{\partial L(\theta)}{\partial v} = \sum_{t=1}^T \frac{\partial L_t(\theta)}{\partial v} \quad \left. \begin{array}{l} \text{add gradient} \\ \text{matrix obtained} \\ \text{at every time stamp} \end{array} \right\}$$

back matrix
of size $N \times N$

(2) Gradient w.r.t W

$$\frac{\partial L(\theta)}{\partial w} = \sum_{t=1}^T \frac{\partial L_t(\theta)}{\partial w} \quad \left. \begin{array}{l} \text{from gradient obtained} \\ \text{from above} \end{array} \right\}$$

For the network = apply layers of the network.



Now,

$$\frac{\partial L(\theta)}{\partial w} = \frac{\partial L_n(\theta)}{\partial s_4} \frac{\partial s_4}{\partial w}$$

refer Backprop derivative to small error first hidden layer for reduction of error first layer required gradient w.r.t w by backprop

$$\text{Now, } s_4 = \sigma(w s_3 + b) \quad \text{in } w, v$$

Here s_4, s_3, \dots, s_1 are functions of w thus we need as follows

Here we have ignored or for simplicity

DELL
PAGE NO.:
DATE:



$$\frac{\partial S_y}{\partial w} = \frac{\partial^+ S_y}{\partial w} + \frac{\partial S_y}{\partial S_3} \frac{\partial S_3}{\partial w}$$

Explicit Implicit

(has direct path to w) (has indirect path to w)

$$= \frac{\partial^+ S_y}{\partial w} + \frac{\partial S_y}{\partial S_3} \left[\frac{\partial^+ S_3}{\partial w} + (\frac{\partial S_3}{\partial S_2} \frac{\partial S_2}{\partial w}) \right]$$

Explicit Implicit

$$\frac{\partial S_y}{\partial w} = \frac{\partial^+ S_y}{\partial w} + \frac{\partial S_y}{\partial S_3} \frac{\partial S_3}{\partial S_2} \frac{\partial S_2}{\partial w} + \frac{\partial S_y}{\partial S_3} \frac{\partial S_3}{\partial S_2} \frac{\partial S_2}{\partial S_1} \frac{\partial S_1}{\partial w}$$

$$= \frac{\partial S_y}{\partial S_4} \frac{\partial^+ S_4}{\partial w} + \frac{\partial S_y}{\partial S_3} \frac{\partial^+ S_3}{\partial w} + \frac{\partial S_y}{\partial S_2} \frac{\partial^+ S_2}{\partial w} + \frac{\partial S_y}{\partial S_1} \frac{\partial^+ S_1}{\partial w}$$

End of end of end

$$\frac{\partial S_y}{\partial w} = \sum_{k=1}^4 \frac{\partial S_y}{\partial S_k} \frac{\partial^+ S_k}{\partial w}$$

$$\therefore \frac{\partial L_y(\theta)}{\partial w} = \frac{\partial L_y(\theta)}{\partial S_y} \left[\sum_{k=1}^4 \frac{\partial S_y}{\partial S_k} \frac{\partial^+ S_k}{\partial w} \right]$$

$$(3) [c_{10} \circ \dots \circ c_{i0} \circ \dots \circ c_{j0}] = ?$$

General

$$\frac{\partial L_f(\theta)}{\partial w} = \frac{\partial L_f(\theta)}{\partial S_i} \left[\sum_{k=1}^4 \frac{\partial S_i}{\partial S_k} \frac{\partial^+ S_k}{\partial w} \right]$$

(3) Gradient cost $\nabla_w L_f(\theta)$
use chain rule through $\nabla_w L_f(\theta)$ and find it



Scanned with OKEN Scanner

* Exploding Δ Vanishing Gradients

→ Now, we find forward pass (forward propagation) (forward pass is feed)

$$\frac{\partial L}{\partial s_j} = \frac{\partial L}{\partial s_k} \frac{\partial s_k}{\partial s_j}$$

$$\frac{\partial s_t}{\partial s_j} = \frac{\partial s_t}{\partial s_{t-1}} \frac{\partial s_{t-1}}{\partial s_{t-2}} \dots \frac{\partial s_i}{\partial s_j} = \prod_{j=k}^{t-1} \frac{\partial s_{i+1}}{\partial s_j}$$

$$\rightarrow \text{Thus } \frac{\partial L}{\partial s_j} = \frac{\partial L}{\partial s_k} \frac{\partial s_k}{\partial s_j} + \frac{\partial L}{\partial s_{k-1}} \frac{\partial s_{k-1}}{\partial s_j} + \dots + \frac{\partial L}{\partial s_1} \frac{\partial s_1}{\partial s_j}$$

$\frac{\partial s_j}{\partial s_{j-1}}$ is to be found

let,

$$a_j = w s_{j-1} + b + u_{n,j}$$

$$s_j = \sigma(a_j)$$

$$a_j = [a_{j1}, a_{j2}, a_{j3}]^T \quad s_j = \sigma(a_j)$$

$$s_j = [\sigma(a_{j1}), \sigma(a_{j2}), \dots, \sigma(a_{jd})] - ②$$

→ Now

$$\frac{\partial s_j}{\partial s_{j-1}} = \frac{\partial s_j}{\partial a_j} \frac{\partial a_j}{\partial s_{j-1}}$$

will be diagonal matrix \times (forward pass is feed)

w big loss forward pass is feed

from ① & ②

$$\frac{\partial \hat{S}_j}{\partial a_j} = \begin{bmatrix} \frac{\partial \hat{S}_{j1}}{\partial a_{j1}} & \frac{\partial \hat{S}_{j2}}{\partial a_{j2}} & \dots \\ \vdots & \ddots & \vdots \\ 0 & 0 & 0 \end{bmatrix} \times w$$

multiplied

$$= \begin{bmatrix} \sigma'(a_{j1}) & 0 & \dots \\ 0 & \sigma'(a_{j2}) & \dots \\ \vdots & \vdots & \ddots & \sigma'(a_{jd}) \end{bmatrix} \times w$$

$$= \text{diag}(\sigma'(a_j)) \times w$$

Taking magnitude

$$\left\| \frac{\partial \hat{S}_j}{\partial S_{j-1}} \right\| = \left\| \text{diag}(\sigma'(a_j)) w \right\|$$

by property of norm

$$\leq \left\| \text{diag}(\sigma'(a_j)) \right\| \|w\|$$

(if $\|ab\| \leq \|a\| \|b\|$)

$$\sigma'(a_j) \leq \frac{1}{4} = \gamma \quad (\text{if } \sigma \text{ is logistic})$$

$$\sigma(a_j)(1 - \sigma(a_j)) \leq \frac{1}{4}$$

\downarrow
max value
taken by derivative

$$\text{is } \frac{1}{4}$$

at max of sigmoid and vector length

$$\sigma'(a_j) \leq 1 = \gamma \quad (\text{if } \sigma \text{ is tanh})$$

(2) If $\gamma \lambda > 1$

$$\leq \gamma \|\omega\|$$

$$\begin{array}{c} \text{Q} \\ \times \\ \text{L} \\ \text{gradient} \end{array} \quad \begin{array}{c} \text{S} \\ \times \\ \gamma \lambda \end{array}$$

Now,

$$\left\| \frac{\partial S_t}{\partial S_k} \right\| = \left\| \prod_{j=k+1}^t \frac{\partial S_j}{\partial (S_j^{(0)})^{(0)}} \right\|$$

$$\left\| \frac{\partial S_t}{\partial S_k} \right\| \leq \prod_{j=k+1}^t \gamma \lambda$$

$$\leq (\gamma \lambda)^{t-k}$$

 $\rightarrow \gamma \lambda < 1$ vanishing gradient

distant from origin

 $\rightarrow \gamma \lambda > 1$ exploding gradient

* Important

$$(\text{check in } i) \quad \gamma \geq (c_{\text{opt}})^{(0)}$$

$$\frac{\partial L_t(\theta)}{\partial w} = \frac{\partial L_t(\theta)}{\partial S_2} \sum_{k=1}^t \frac{\partial S_k}{\partial S_{k-1}} \frac{\partial S_{k-1}}{\partial w}$$

dxd

✓

dxd

✓

dxd

gradient

dxdxd

froced values are known to us

$$(\text{check in } i) \quad \gamma \geq t \geq (c_{\text{opt}})^{(0)}$$

→ Let's compute any one element of tensor $d \times d \times d$

$$\frac{\partial^+ s_{kp}}{\partial w_{qr}}$$

(p, q, r) element of tensor.

Note $\Delta p = 1$ if $i = j$
diff wrt w are 0

$$a_k = w s_{k-1} + (b + U a_k)$$

$$s_k = \sigma(a_k)$$

$$\therefore a_{lc} = w s_{k-1}$$

$$\begin{bmatrix} a_{k1} \\ a_{k2} \\ \vdots \\ a_{kd} \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1d} \\ w_{p1} & w_{p2} & \dots & w_{pd} \\ \vdots & \vdots & \ddots & \vdots \\ w_{d1} & w_{d2} & \dots & w_{dd} \end{bmatrix} \begin{bmatrix} s_{k-1,1} \\ s_{k-1,2} \\ \vdots \\ s_{k-1,d} \end{bmatrix}$$

$$\therefore a_{kp} = \sum_{i=1}^d w_{pi} s_{k-1,i} \quad \textcircled{1}$$

$$s_{kp} = \sigma(a_{kp})$$

$$\frac{\partial s_{kp}}{\partial w_{qr}} = \frac{\partial s_{kp}}{\partial a_{kp}} \frac{\partial a_{kp}}{\partial w_{qr}}$$

$$= \sigma'(a_{kp}) \frac{\partial a_{kp}}{\partial w_{qr}}$$

~~Initial vector for hidden state has two steps in it~~

$$\frac{\partial \text{step} p}{\partial w_{qr}} = \sum_i \frac{\partial p(s_{t-1}, i)}{\partial w_{qr}}$$

$$\begin{matrix} q+2 & 6 \\ q+3 & 6 \end{matrix}$$

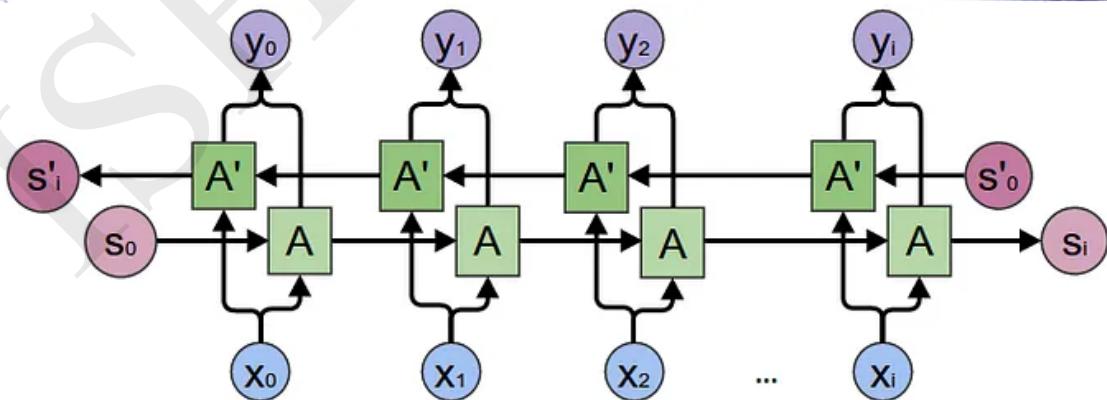
$$= s_{t-1, i} \quad \text{if } p = q \quad \text{else}$$

$$= 0 \quad \text{otherwise}$$

$$\therefore \frac{\partial S_{kp}}{\partial w_{qr}} = \sigma'(a_{kp}) s_{t-1, r} \quad \text{if } p = q \quad \text{else}$$

$$\begin{cases} 1 & \text{if } p = q \\ 0 & \text{otherwise} \end{cases}$$

* Bi-Directional RNNs



WEEK 14

* Selective Read, Write, Forget.

→ In a RNN the problem of vanishing and exploding gradient occurs during backprop.

But, during forward prop also we keep on updating the state vector. This vector is also of finite dimension so by reaching at end of RNN it will also forget info from 1st state.

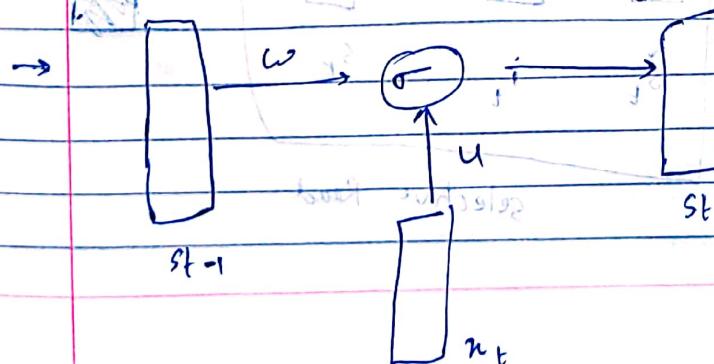
During forward pass & Back pass info is tend to be forgotten.

→ Selective read, write, forget using white-board analogy

* LSTM (long-short term memory) (W, V, U)

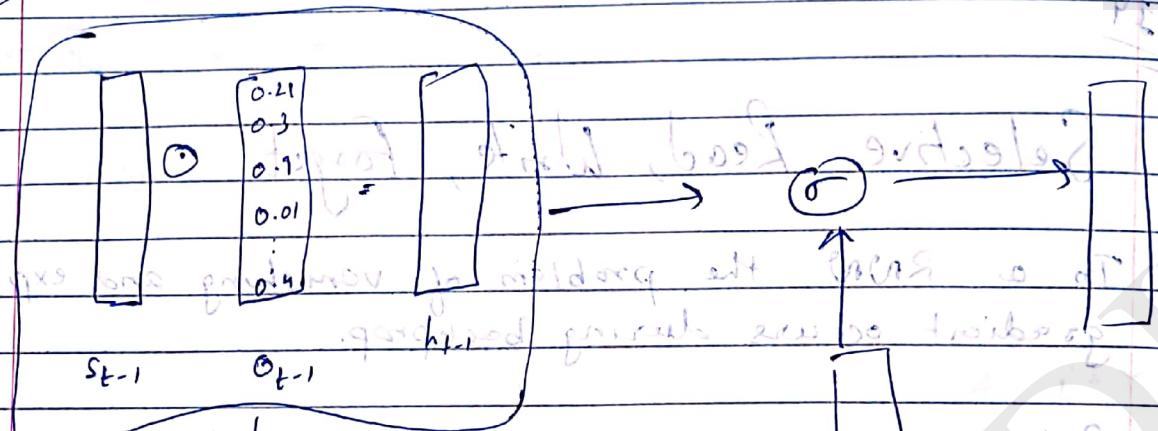
→ lets consider task of sentiment analysis.

- Forget - the information like stop words
- Selective Read - the info added by previous sentiment bearing words
- Selective Write - new information from current word to state.



Here if we want to selectively write values of s_{t-1}

Selective Write



Selective
write

including fractions like 1/2, 1/3, etc. can also be stored in the state cell. The output from S_{t-1} is used to select which cell to make some values in S_{t-1} close to zero & some close to 1 & so on as per importance.

at first in first step itself it was known which cell is important and

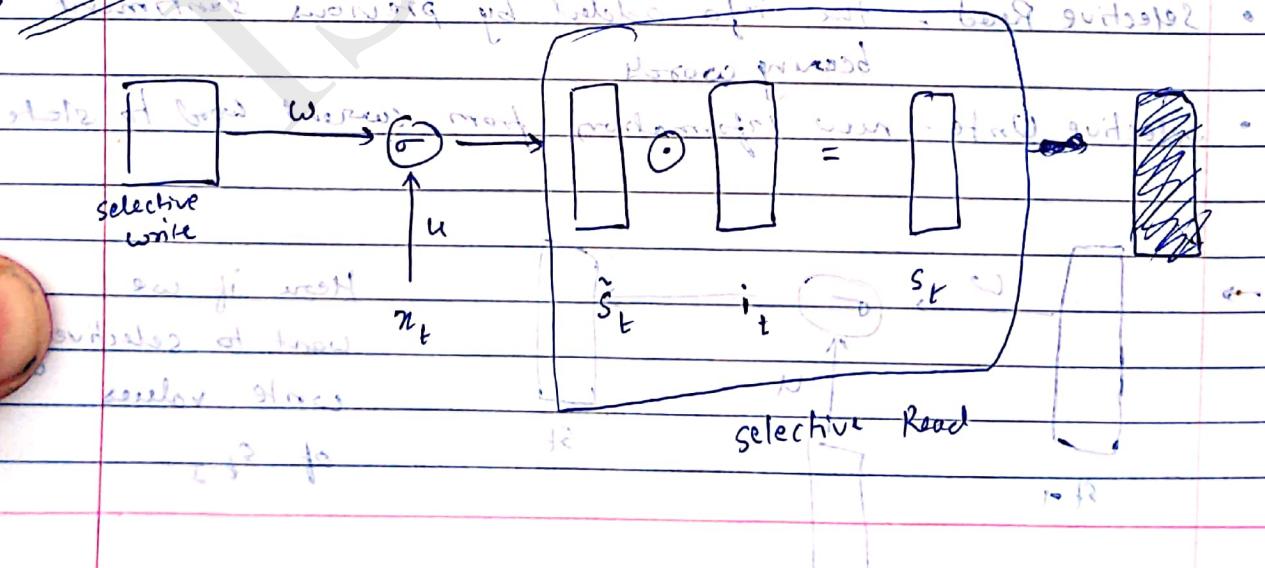
Now, we need to learn O_{t-1} (Output Gate)

$$O_{t-1} = \sigma(W_o h_{t-1} + V_o x_{t-1} + b_o)$$

$\rightarrow (W, V, W)$ are parameters of RNN (O_{t-1}) [TTS]

b_o, W_o are parameters of RNN (O_{t-1}) [TTS]

Selective Read



$$\tilde{s}_t = \sigma(W h_{t-1} + U x_t + b) \quad \text{Eq. 1}$$

$$s_t = \tilde{s}_t \odot i_t$$

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i)$$

(tiny memory leak) [18]

Note

→ w, u, b & w_i, u_i, b_i are different parameters

~~Till now~~

- Previous state $\rightarrow s_{t-1} = W h_{t-1} + U x_{t-1} + b$

~~del~~

$$(j_d + j_m + j_w) \rightarrow i_t$$

- Output gate $\rightarrow o_{t-1} = \sigma(W_o h_{t-1} + U_o x_{t-1} + b_o)$

- Selective Write $\rightarrow h_{t-1} = o_{t-1} \odot \sigma(s_{t-1})$

$$(j_d + j_m + (j_w \odot j_o)) \rightarrow i_t$$

- Current (temp) state $\rightarrow \tilde{s}_t = \sigma(W h_{t-1} + U x_t + b)$

$$(j_d + j_m + (j_w \odot j_o)) \rightarrow i_t$$

- Input gate $\rightarrow i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i)$

~~forget~~ | ~~read~~ Selective Read $\rightarrow j_i \wedge O_2 \tilde{s}_t$ | ~~forget~~ layer

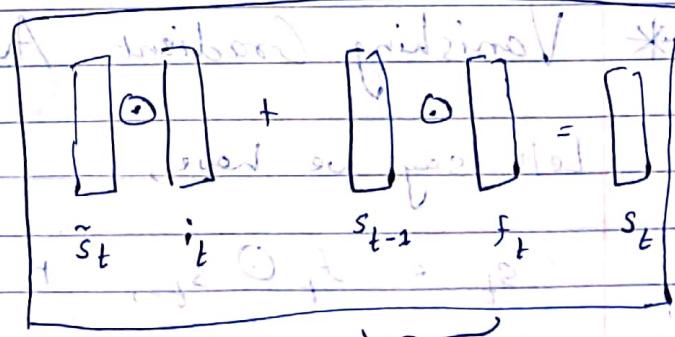
~~forget~~ | ~~read~~ $\rightarrow j_i \wedge O_2 \tilde{s}_t$ | ~~forget~~ layer

Selective Forget



~~forget~~ | ~~read~~ | ~~forget~~

$$n_t \odot i_t$$



~~forget~~

~~read~~

~~forget~~

~~forget~~

$$f_t = \sigma(C_d w_i h_{t-1} + U_d v_t + b_d)$$

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

* GRU (Gated Recurrent Unit)

→ Gates

$$\bullet o_t = \sigma(C_o s_{t-1} + U_o v_t + b_o)$$

$$\bullet i_t = \sigma(C_i s_{t-1} + U_i v_t + b_i)$$

→ States

$$\bullet \tilde{s}_t = \sigma(C_s o_t \odot s_{t-1} + U_s v_t + b_s)$$

$$\bullet s_t = (1 - i_t) \odot s_{t-1} + i_t \odot \tilde{s}_t$$

Note refer Notes / week 14 / LSTM.png & Notes / week 14 / lstm.png
 notes / week 14 / GRU.png & Notes / week 14 / gru.png

* Vanishing Gradient Avoiding LSTM

Let's say, we have,

$$s_t = f_t \odot \tilde{s}_{t-1} + i_t \odot \tilde{s}_t$$

In the worst case we say that \tilde{s}_t tends to 0 or vanishes

Thus we have $(\alpha)_t$ for backprop all the way up to s_t . Now we just need to do forward pass which starts with $s_t = f_t \circ s_{t-1}$.

$$s_t = f_t \circ s_{t-1} \quad \text{Initial, } B \text{ of } (\alpha)_t$$

$\frac{\partial E}{\partial s_t} = \alpha_t f_t$ for $(\alpha)_t$ for backprop up to s_t . We need to propagate it.

as we go back pass $f_t \times f_{t-1} \times f_{t-2} \times \dots \times f_1$

$$= (f_t)^T$$

(β_t, α_t) $\rightarrow i$

(β_t, α_t) Here also during back pass f_t vanishes

(β_t, α_t) But during forward pass

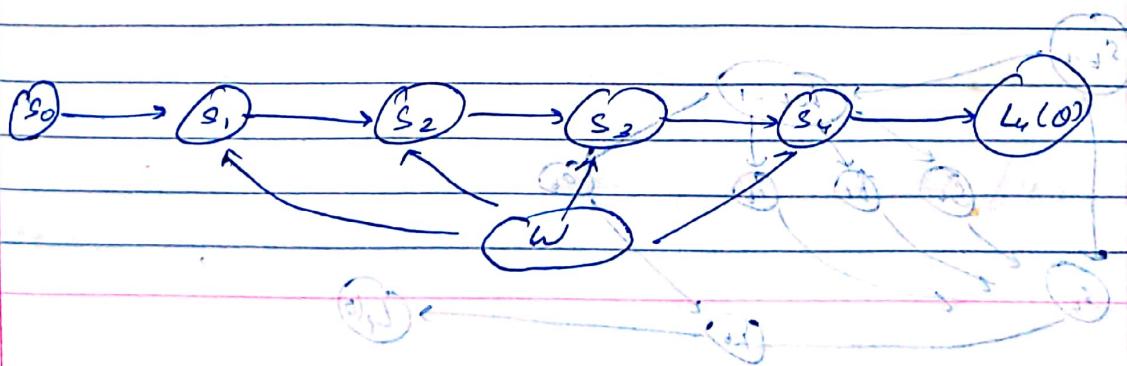
f_t information goes from i state to next

$$\therefore (f_t)^T \text{ for forward.}$$

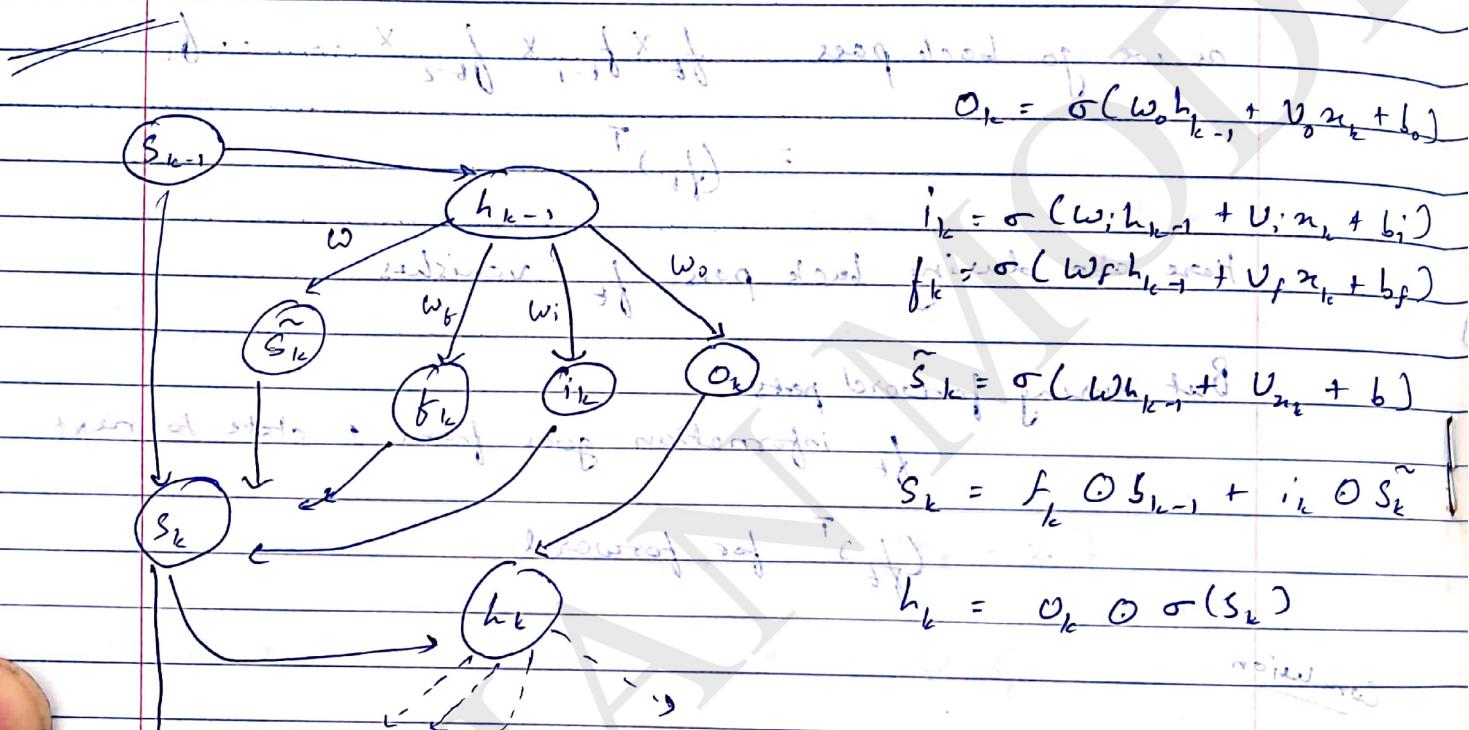
Conclusion

- If the state at time $t-1$ did not contribute to state at time t (ie, if $\|f_t\| \approx 0$ and $\|\alpha_{t-1}\| \approx 0$) then during backprop the gradients flowing into s_{t-1} will vanish.

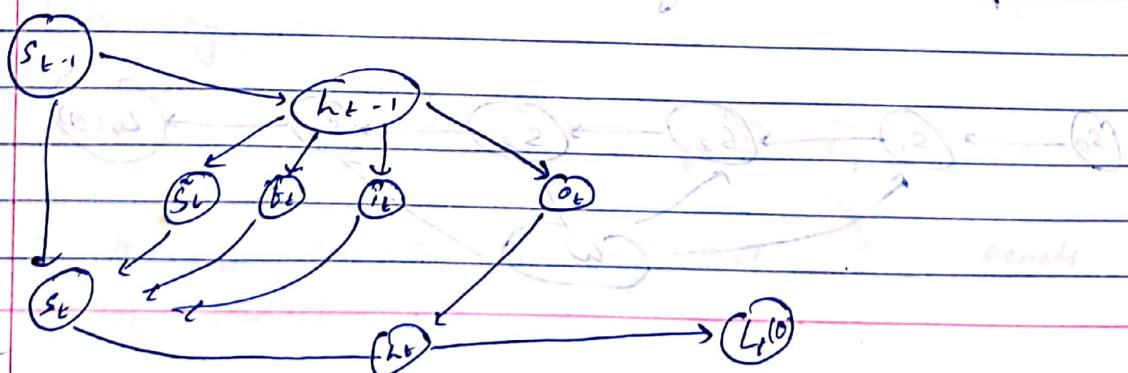
* Proof (Vanishing doesn't occur)



- In general, the gradient of $L_t(\theta)$ w.r.t θ_i vanishes when the gradients flowing through each and every path from $L_t(\theta)$ to θ_i vanishes.
- The gradient of $L_t(\theta)$ w.r.t θ_i explodes when the gradient flowing through at least one path explodes.



Note: At standard time step t we have state s_t and hidden state h_t .
 At time $t+1$ we will have all hidden states h_1, h_2, \dots, h_t and new hidden state h_{t+1} .
 We will also have new states s_1, s_2, \dots, s_t and new state s_{t+1} .



→ Now lets say w_j causes error in output & we want to backprop gradient to (w_j) .

Then we need to prove that there exist at least one path from $L(0)$ to S_k that allows greatest flow so that it does not vanish in later part.

$$\frac{\partial L_t(\theta)}{\partial S_{1:t}} = \frac{\partial L_t(\theta)}{\partial h_t} \frac{\partial h_t}{\partial S_t} \left(\frac{\partial S_{t+1}}{\partial S_t}, \dots, \frac{\partial S_{L+1}}{\partial S_t} \right)$$

(~~longest~~
can be
found) \rightarrow quickstep \downarrow (2) fast

$$\left(\begin{array}{cc} 06 & 16 \\ 16 & 06 \end{array} \right) = \cancel{116} \left(\begin{array}{cc} 06 & 16 \\ 16 & 06 \end{array} \right) \overset{(6)}{\cancel{16}} =$$

$$\textcircled{1} \quad \frac{\partial h_t}{\partial s_1} = \sigma(O_t O \sigma(s_1)) = \text{Diagonal matrix } \begin{pmatrix} O_t O \sigma'(s_1) \\ \vdots \\ O_t O \sigma'(s_1) \end{pmatrix}$$

$$(2) \quad \frac{\partial s_t}{\partial s_{t-1}} = \frac{\partial f_t \circ s_{t-1} + (i_t \circ s_L)}{\partial s_{t-1}} \xrightarrow{\text{ignoring worst case}}$$

$$\|f_t\| = \sqrt{\int_{\Omega} |f_t(x)|^2 dx} = \sqrt{\int_{\Omega} (c_t(x))^2 dx} = \sqrt{\int_{\Omega} c_t^2(x) dx} = \|c_t\|$$

skew-symmetric diagonal matrix will be zero on the principal

$\left(\frac{\partial S_t}{\partial S_{t+1}} \text{ weight } - \frac{\partial S_{t+1}}{\partial S_t} \text{ weight } \right) \text{ (if } S_t \text{ is present in front of } S_{t+1})$

$\therefore \Delta L_t \rightarrow \Delta L_{t-1} \rightarrow \dots \rightarrow \Delta L_0 \rightarrow \Delta L_{t-k+1}$

$$= L'_t(h_t) \text{ of } D(\sigma(s_t) \circ o'_t) \text{ at } h_t$$

Need to find out if $D(\sigma^t_{b=k+1} \circ \dots \circ \sigma^t_0)$ will vanish towards last (2) or (3) , may stop here.
 Thus gradient vanish during backprop only if it has vanished in forward prop. If vanish only when it needs to this proves (3) .

* Proof (Exploding can be stopped)

path $L_t(o) \rightarrow h_t \rightarrow o_t \rightarrow h_{t-1} \dots h_k \rightarrow o_k \rightarrow h_{k-1}$

$$= \frac{\partial L_t(o)}{\partial h_t} \left(\frac{\partial h_t}{\partial o_t} \frac{\partial o_t}{\partial h_{t-1}} \right) \dots \left(\frac{\partial h_k}{\partial o_k} \frac{\partial o_k}{\partial h_{k-1}} \right)$$

$$\stackrel{(1) \rightarrow 0, 0}{\text{using}} = (2) \sigma(s_0) \circ o_0 \quad (1)$$

$$\stackrel{\text{using } L'_t(h_t) \cdot (D(\sigma(s_t) \circ o'_t) \cdot w_0)}{=} \dots \quad (2)$$

$$\text{similarly } (2) \sigma(s_1) \circ o_1 \quad (2)$$

$$\| \tau \| \leq \| L'_t(h_t) \| (\| k \| \| w_0 \|)^{t-k+1}$$

Depending on norm of $\| w_0 \|$ gradients may explode

Thus use Gradient Clipping if norm exceeds certain value bring it down to a threshold.

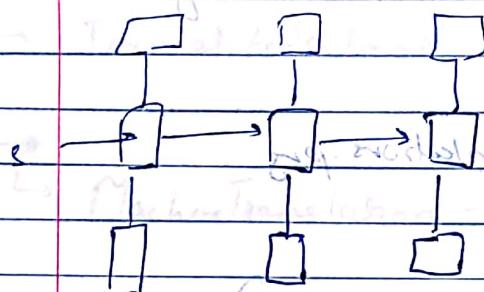
~~WEEK 15~~ What will be the next token in given sequence based on previous tokens?

* Introduction to Encoder-Decoder networks

Input sequence with final output tokens are being predicted.

→ Predicting next word in sequence

$$y^* = \operatorname{argmax} P(y_t | y_1, y_2, \dots, y_{t-1})$$



At each stage in RNN, it will predict a probability distribution & word with max prob. will be chosen.

Using RNN we can compute

$$P(y_t=j | y_1^{t-1}) = \text{softmax}(v_s + c)_j$$

$y_t = j$ is vocabulary element of v . v is vocabulary.

(j is vector of y_t)

$$P(y_t=j | y_1^{t-1}) = P(y_t=j | s_{t-1})$$

All the information of $(y_1, y_2, \dots, y_{t-1})$ is contained in s_{t-1}

Thus model is

refer notes/week15/lec1/model.png

Links of file will be soon uploaded.

Thank you for watching.

During training we would have all the input along with sequence

But, during test we will be given only a word and we will have to build entire sequence using it

refer notes/week 15/lec 1/testcase.py

→ We will now use following representations for RNN, GRU, LSTM

This refers to notes/week 15/lec 1/representations.py

and this refers to notes/week 15/lec 1/testcase.py

→ Image Captioning (Encoder - Decoder)

$$P(y_t | y_{t-1}) = P(y_t | s_t)$$

$$i(s + \epsilon v) \rightarrow P(y_t | s_t) = P(y_t | i = s_t)$$

Here we have image as well.

$$P(y_t | s_t, v, I)$$

cannot take entire image
we need representation of image

$$P(y_t | s_t, f_{\text{enc}}(I))$$

fully connected layer
in VGG.

benefits of using VGG

in hidden unit

This becomes s_0

pre-trained (1000) 21 layers refer

$\therefore s_0 = f_{\text{enc}}(I)$ now we feed it to RNN
for caption generation

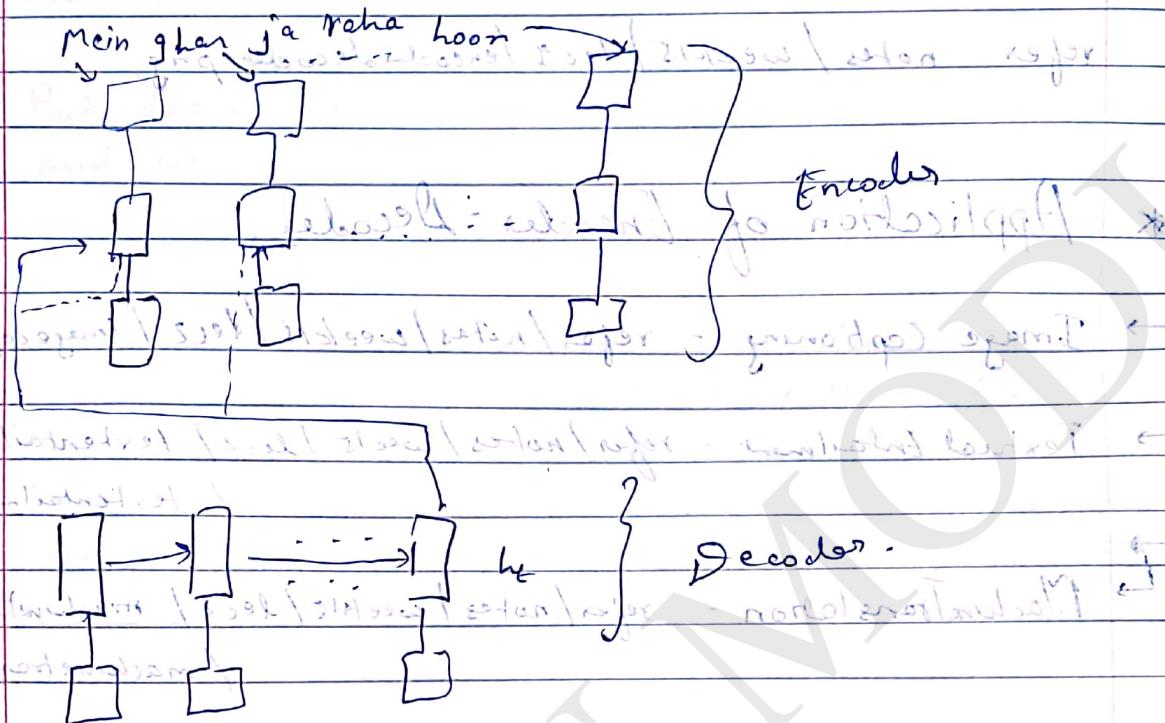
We can also feed in fig (ii) to all states.

refer notes / week15 / lec 2 / encoderdecoder.py

* Application of Encoder - Decoder

- Image Captioning - refer / notes / week15 / lec 2 / imagecaptioning.py
- Textual Entailment - refer / notes / week15 / lec 2 / textentailment.py
/ textentailment2.py
- Machine Translation - refer / notes / week15 / lec 2 / machinetranslation.py
/ machinetranslation2.py
- Transliteration - refer / notes / week15 / lec 2 / transliteration2.py
/ transliteration2.py
- Image QA - refer / notes / week15 / lec 2 / imageqa.py
- Document summarization (read page no 1) - refer / notes / week15 / lec 2 / docsummary.py
- Video Captioning - refer / notes / week15 / lec 2 / videocaption.py
- Video Classification - refer / notes / week15 / lec 2 / videoclassification.py
- Dialog - refer / notes / week15 / lec 2 / dialog.py

* Attentional Mechanism

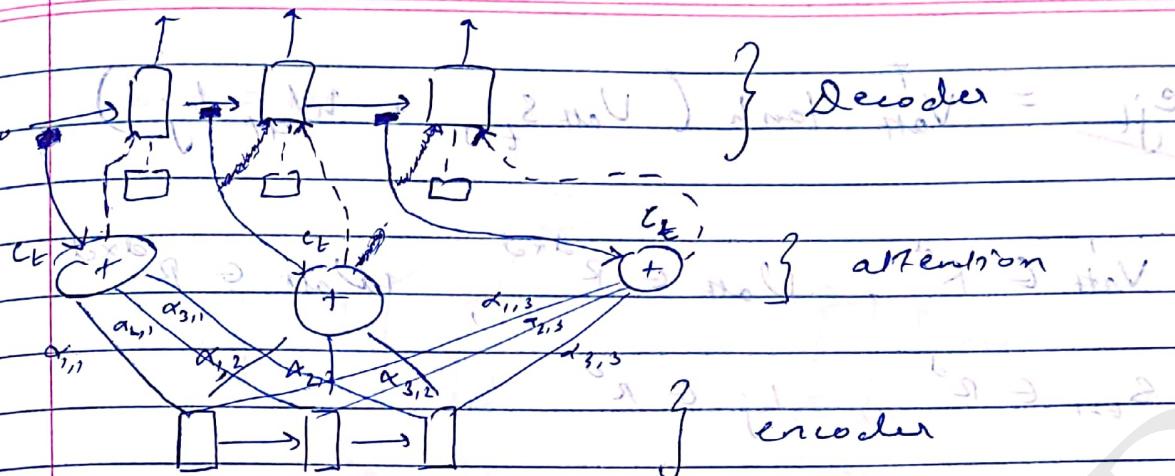


So here we can see the passing of h_t to the network. This h_t contains all the information ie (I am going home) that we can observe that you can see the word after

home → ghar These are words that are highly amies & closely related.

going → ja raha

So the intuition is why not pass this relationship instead of h_t



small weight matrix α is to be learned between s_t and d_i .

Here $\alpha_{1,1}, \alpha_{1,2}, \alpha_{1,3} | \alpha_{2,1}, \alpha_{2,2}, \alpha_{2,3} | \alpha_{3,1}, \alpha_{3,2}, \alpha_{3,3}$

are all parameters which are to be learnt.

→ Now, ~~decoder state~~ ~~encoder state~~ ~~decoder state~~ ~~encoder state~~

$$e_{j,t} = f_{ATT}(s_{t-1}, h_j)$$

against now position j

j^{th} word \rightarrow t^{th} time stamp in decoder of j^{th} word. $e_{j,t}$ will be t^{th} time stamp in encoder of j^{th} word. $e_{j,t}$ is called j^{th} word's t^{th} time stamp.

$$\alpha_{j,t} = \frac{\exp(e_{j,t})}{\sum_{j=1}^M \exp(e_{j,t})}$$

mark $\alpha_{j,t}$ for j^{th} word

Thus $\alpha_{j,t}$ is probability

j^{th} word with respect to t^{th} time stamp.

$$e_{jt} = \text{Vatt}^T \tanh (\text{Vatt} s_{t-1} + \text{Watt} h_j)$$

scales $\text{Vatt} \in \mathbb{R}^d$, $\text{Vatt} \in \mathbb{R}^{d \times d}$, $\text{Watt} \in \mathbb{R}^{d \times d}$

$s_{t-1} \in \mathbb{R}^d$, $h_j \in \mathbb{R}^d$

This thing works because it is better modeling choice

Also, e_{jt} is used for encoder.

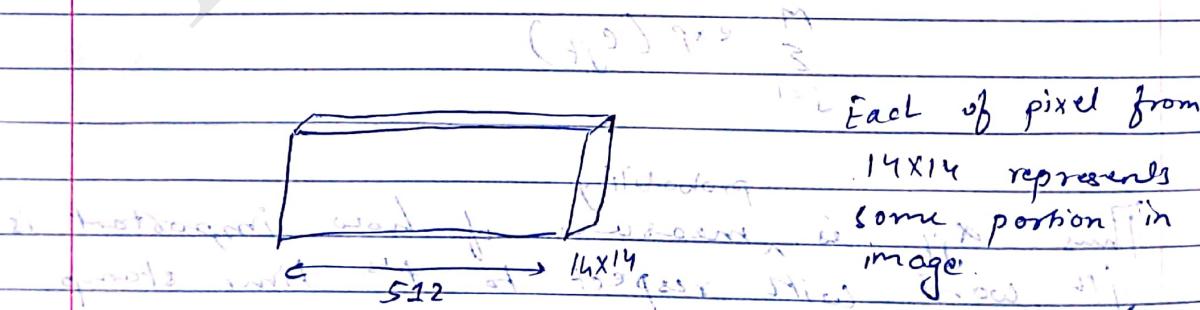
$$e_{jt} = \sum_{j=1}^T \alpha_{jt} h_j$$

refer notes / week 15 / dec 3 / attention mechanism.pdf

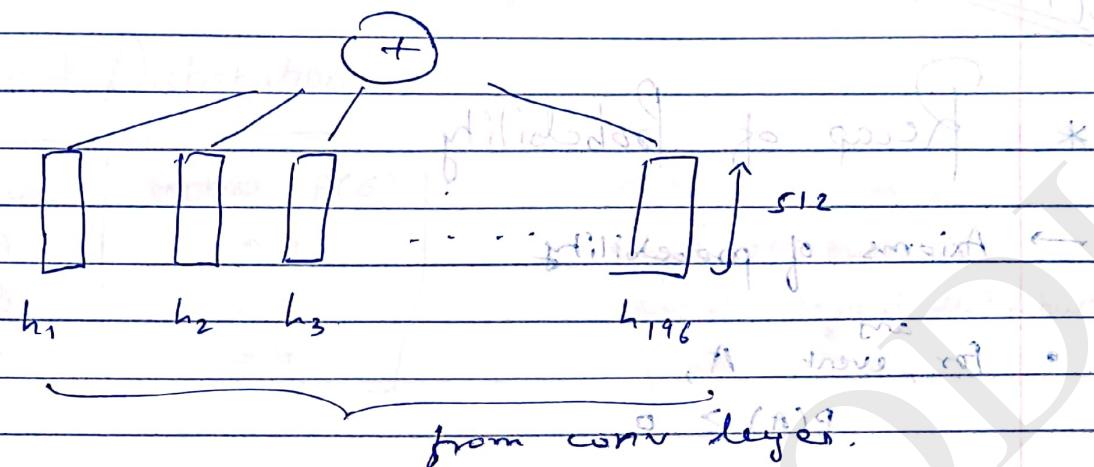
* Attention Over Images

Unlike RNN, where each state gives word to conv for image given in FC, which doesn't give us which entry is for which pixel in image

Thus we use the convolution layers before FC,



$$14 \times 14 = 196$$



$$\alpha_{tj} = f_{att}(s_{t-1}, h_j)$$

$$f_{att} = v^T \tanh(u_s^T s_{t-1} + w_h h_j + b)$$

* Hierarchical Attention

→ Understanding hierarchical RNN structure

refer notes/week15/lec5/Dialog hierarchy RNN.py
/Document hierarchy spng.

→ Now we use attention mechanism in these hierarchical models.

• refer notes/week15/lec5/att mech hierarchy encoder.py
/att mech hierarchy decoder.py.

higher level -> higher level information
lower level -> lower level information