

* Convolutional Neural Network (CNN)

→ 1D CNN

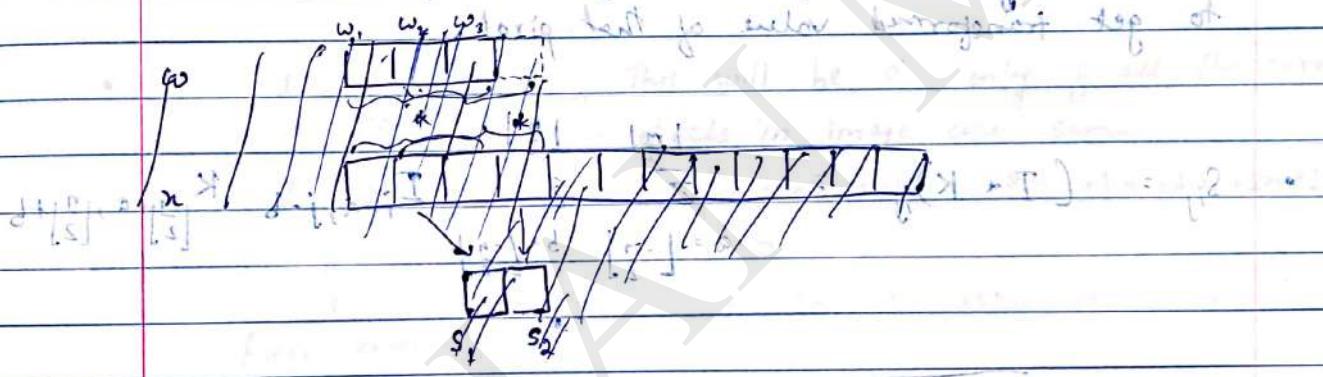
$$s_t = \sum_{a=0}^{n-1} x_{t-a} w_a \quad \text{where } n = (n_1 * w_1) + b$$

size of filter

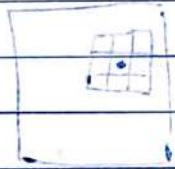
The diagram illustrates the convolution process. An input sequence x is shown as a horizontal row of boxes. A filter w is applied across the input, with weights w_1, w_2, w_3 indicated above it. The result is a single output value s .

→ Let's say $n = 3$

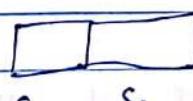
Let's say toxic will be fed in 3x3 matrix problem so



Let's say toxic for value tag

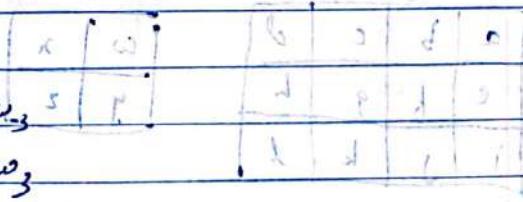


If we consider (x_1, x_2, x_3, x_4) is given with stride 1 then $s_1 = x_1 w_1 + x_2 w_2 + x_3 w_3 + x_4 w_4$



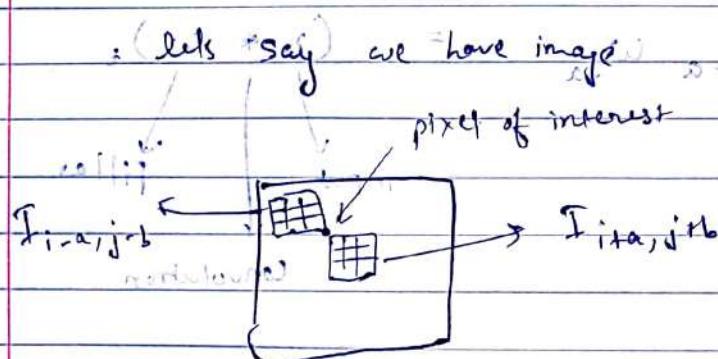
$$s_1 = x_1 w_1 + x_2 w_2 + x_3 w_3 + x_4 w_4$$

$$s_2 = x_2 w_1 + x_3 w_2 + x_4 w_3$$



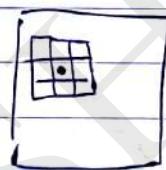
→ 2D Convolutions (With Examples)

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



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

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



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

aw + bw + cw + dw +	bw + cw +	cw + dn +
+ ey + fz	fy + gz	gy + hz
ew + fw +	fw + gx	gw + hx +
+ iy + jz	+ jy + kz	hy + lz

→ Types of Filters (sizing, blur, sharp)

- $\begin{matrix} 1 & 1 & 1 \end{matrix}$ replacing this pixel in image with average
 $\frac{1}{9} \sum_{i=1}^9 M_i$ of all pixels around it

$\begin{matrix} 1 & 1 & 1 \end{matrix}$ 3×3 neighborhood and result

Blur \rightarrow $\frac{1}{9} \sum_{i=1}^9 M_i$ \rightarrow $\frac{1}{9} \sum_{i=1}^9 M_i = \bar{M}$

- $\begin{matrix} 0 & -1 & 0 \\ -1 & 5 & 1 \\ 0 & 1 & 0 \end{matrix}$ central pixel is made 5 times & all other pixels around it are subtracted

$$0 - 1 + 0 = 0$$

Sharpen basic rule of blur is $R=9$
 \rightarrow $\frac{1}{9} \sum_{i=1}^9 M_i = \bar{M}$

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

$$1 + 1 + 1 - 8 + 1 + 1 = -8 + 3 = -5 \neq 0$$

between treatments, we can do this in form of $\frac{1}{9} \sum_{i=1}^9 M_i = \bar{M}$

Edge detector

$$\left[\begin{matrix} 1 & 1 & 1 \\ 1 & -8 & 1 \\ 1 & 1 & 1 \end{matrix} \right] \rightarrow \frac{1}{9} \sum_{i=1}^9 M_i = \bar{M}$$

$$\rightarrow 3D Convolution \left[\begin{matrix} 1 & 1 & 1 \\ 1 & -8 & 1 \\ 1 & 1 & 1 \end{matrix} \right] \rightarrow \frac{1}{9} \sum_{i=1}^9 M_i = \bar{M}$$

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

with assumption (1)

$$\left[\begin{matrix} 1 & 1 & 1 \\ 1 & -8 & 1 \\ 1 & 1 & 1 \end{matrix} \right] \rightarrow \text{EXTRACTOR}$$

if we apply 3D filter to 3D image a 3D feature map is obtained

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

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

→ Shows how output size depends on input size.

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

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

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

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

→ In ^{almost} _{all} cases we might have to add padding. Transform all pixels by taking them at ~~middle~~ ~~center~~ ~~edges~~ after ~~padding~~ ~~done~~ ~~not~~ ~~middle~~ ~~edges~~ ~~padding~~ ~~done~~ ~~not~~ ~~center~~

$P = 1$ 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

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

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

Main formula

(Q)

$$227 \times 227 \times 3 * (11 \times 11 \times 3) \Rightarrow ?$$

② no. of neurons

→ Images going through 96 filters will give 96 output.

$$W_2 \times H_2 \times D$$

Let's do it

Ans Since filter depth and img depth are same

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

Since there are 96 filters (in working)

$$D = 96$$

→ given $s=4$, $p=0$

$$w_2 = \frac{227 - 11}{4} + 1 = 55$$

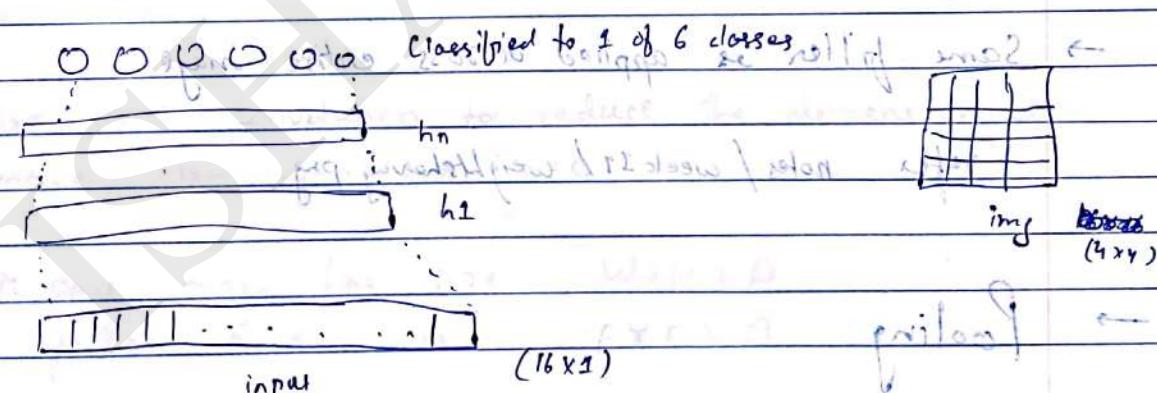
$$\therefore h_2 = \frac{227 - 11}{4} + 1 = 55$$

* 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.

→ Convolution NN

2	4	7
2	5	8
3	6	9

⇒

2	3	4	5	6	7	8	9	10	11	12	13
---	---	---	---	---	---	---	---	----	----	----	----

3x3 img.

3P = 0

9x12 vector

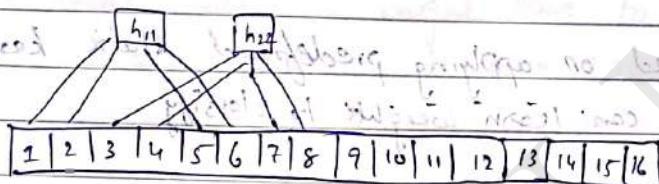
(1)

so this is comparatively sparse

connection

$E + H = 558$

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



(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

using max pool

2×2 (stride=2)

4	9
11	6

feature map obtained after applying filter on img

Similarly ~~Max pooling~~ maxpool with stride 2

$10 \times 7 \times 1 \rightarrow 5 \times 3 \times 1$
 2×2 (stride=2)

$\frac{1+2+3}{4}$	$\frac{5+6+7}{4}$
$\frac{10+11+12}{4}$	$\frac{14+15+16}{4}$

* Image Net

→ ~~AlexNet~~ AlexNet with 5 layers and fc(4096)
resulting in 60 classes

refer notes / week 11 / AlexNet.pptx

→ ZFNet

refer notes / week 11 / ZFNet.pptx

→ VGGNet

refer notes / week 11 / VGGNet.pptx

→ GoogleNet

activation of pre-activations

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

image assuming $p = 0$

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

& filters of size $F \times F \times D$

Total Computations required = $F \times F \times D$

and process by each element of o/p

(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 $D^T D$ computations are reduced
now applying filter size $F \times F \times 9^1$

Total computation for each element of o/p = $F \times F \times D^1$ significantly reduced

- Googlenet has various inception modules having basic layout as follows

refer notes / week 11 / inceptionmodule.py

other inception modules structures

refer referenceppt / 11 / slides 423 - 442

- In the last convolution layer

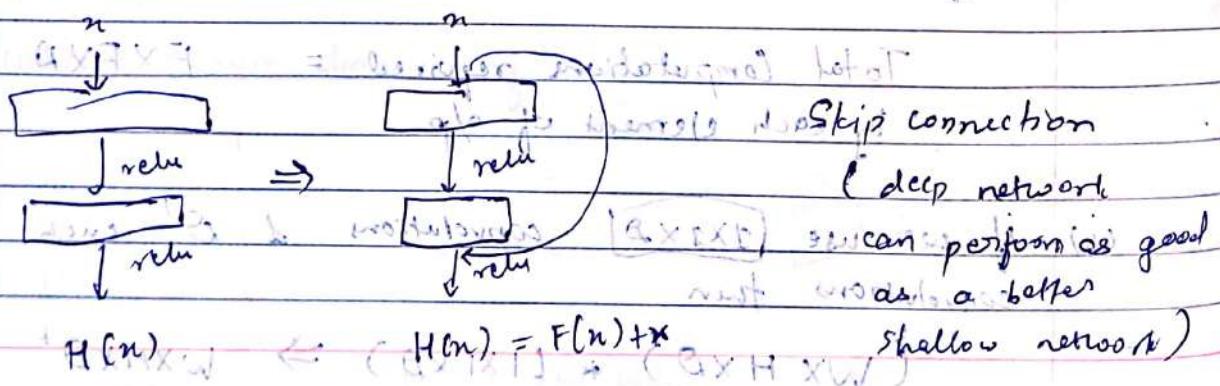
$7 \times 7 \times 1024 \xrightarrow{\text{avg pool}} \underline{1 \times 1 \times 1024}$

equivalent to a vector

in network with lesser can be fully connected with it.

dense layer with less parameters than earlier

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



29 MARCH 2019

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.

$$(d + \delta w + \Delta w) \text{ where } \delta w = \dots$$

more working left of bottom is not working well

$$(d + \delta w + \Delta w) \text{ where } \delta w = \dots$$

reduces noise

and this is working well

$$\delta w = \dots$$

working fine now left in 2 except works not in left

left to right

works

$$V \xrightarrow{\begin{smallmatrix} 1 \\ 2 \\ 3 \\ 4 \end{smallmatrix}} = (V \times A) A$$

WEEK 12

Dimensionality reduction

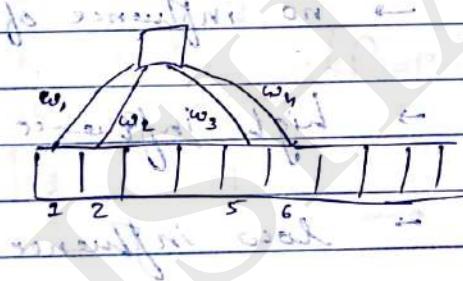
* Visualizing Patches

→ 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 for a neuron to fire without loss.

→ limited to the specific learned weights to determine the input λ thus determine which patch of input is responsible for a neuron to fire.



$$\begin{bmatrix} w_1 & w_2 \\ w_4 & w_5 \end{bmatrix} \Rightarrow [w_1, w_2, w_3, w_4] = \text{weight}$$

$$n_1 | n_2 | n_3 | n_4 = w'$$

thus if $w' = 16$ then $n_1, n_2, n_3, n_4 = 16$

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

* Occlusion Experiments

→ Take a patch in image & replace it by black lo.
~~now try to classify~~ ~~and let's go~~ ~~position~~

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

Similarly do it for all patches to understand important
patches for classification. To ~~and let's go~~ ~~position~~

* Finding influence of Input Pixels using Backprop

Gradients with respect to input give the influence
lets say h_j is a neuron of hidden layer j

Now, if

$$\frac{\partial h_j}{\partial x_i} = 0 \rightarrow \text{no influence of } x_i \text{ on } h_j$$

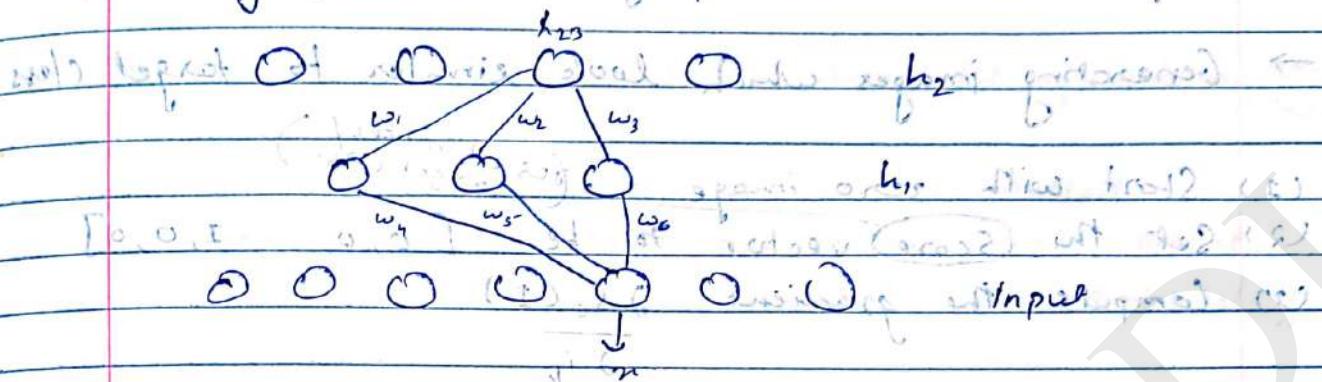
$$\frac{\partial h_j}{\partial x_i} = \text{large} \rightarrow \text{high influence}$$

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

So we represent image as gradients

$\frac{\partial h_j}{\partial x_0}$	$\frac{\partial h_j}{\partial x_1}$	\dots	$\frac{\partial h_j}{\partial x_m}$
1	1	1	1
1	1	1	1
1	1	1	1
1	1	1	1

* Finding gradient wrt inputs w.r.t. maximization

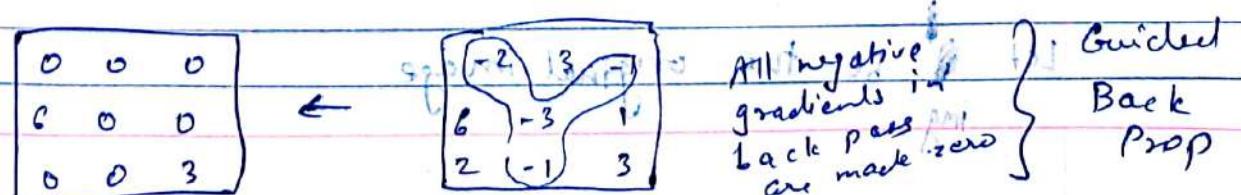
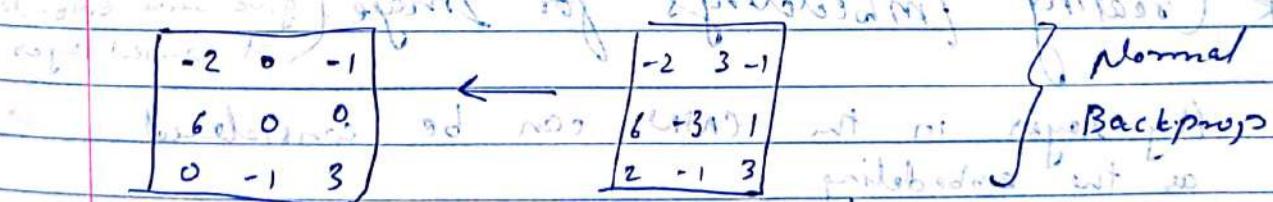
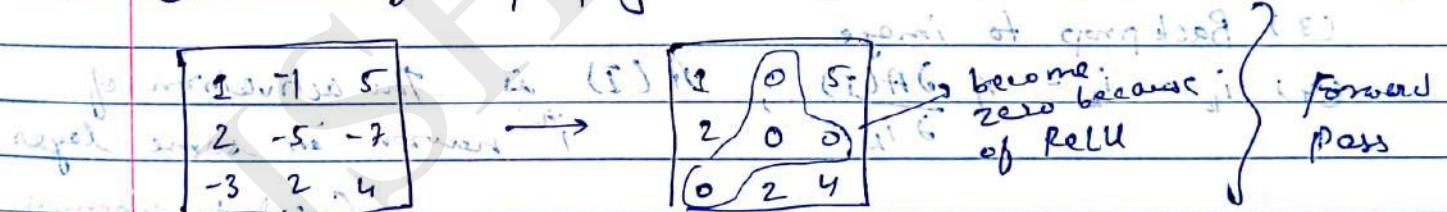


so we want to find $\frac{\partial h_{23}}{\partial x}$ using steepest

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

Now, plot image calculated similarly for all input & plotting.

* Guided BackPropagation (impact of image pixels on result)



* Optimization Over Images

→ 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 layer fires. with

(1) Feed img to network

(2) Make activation such that selected neuron is one other are zero

(3) Backprop to image

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

* Creating embeddings for Image (used to construct an image that would give same embedding at wanted layer)

Any layer in the CNN can be considered as the embedding

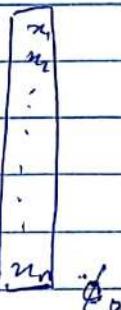
Let img be the original image

Let img be the original image

0	0	0
0	0	0
0	0	0

We pass this image through the Network.

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.

$$\text{embedding} = (\text{values})$$

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

$$L(\phi_i) = \|\phi(u) - \phi_i\|^2 + \lambda \|\phi(u)\|^2$$

obtained embedding required embedding

update rule

$$\text{image} = \eta \frac{\partial L}{\partial u}$$

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

Similarity is possible for convolution layers as well.

Similarity is not good if we want to compare two images.

Deep Dream

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

So we if this is a neuron in convolution layer,

we want to maximize $L(I)$

$$L(I) = h_{ij}^2$$

During backprop. update pixel of image

functionality of web based systems maintained by various web hosts
hosted on $\partial L(I)$ and $\partial \bar{L}(I)$ is highly local and does not
overlap with $\text{dim}(\mathbb{Z})$ which is often quite large and
updateable.

$$\| \text{softmax} \| = \lambda \text{ max} + \| \Delta L(I) \circ \sigma' \| = (1, 1)$$

→ 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. available fire report on opening

Thus when we give image of sky and increase the probability of neuron denoting sky to 1. Then it automatically loses knowledge from training data & generates birds, castle as foreground.

* Deep Art

referred book: DCCP

→ 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.

New image is generated such that all image representations / of original image embeddings generated using new image are similar to embeddings generated by actual image

(2) Style (cartoon)

According to CV study $V \in \mathbb{R}^{64 \times (232 \times 256)}$

↑ computer vision ↓ dimension of embedding

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

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

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

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

refer notes/week 12/depict.py

* Fooling CNNs

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

(gradient) $\delta p / \delta v$

$\nabla L / \nabla v \rightarrow \nabla p / \nabla v$ plus $v^T v$ or $v^T v$

(gradient of $\delta p / \delta v$) $\nabla v^T v = v^T v$

$$(C_{\text{style}} - v^T v)^2 = \beta \text{ loss}$$

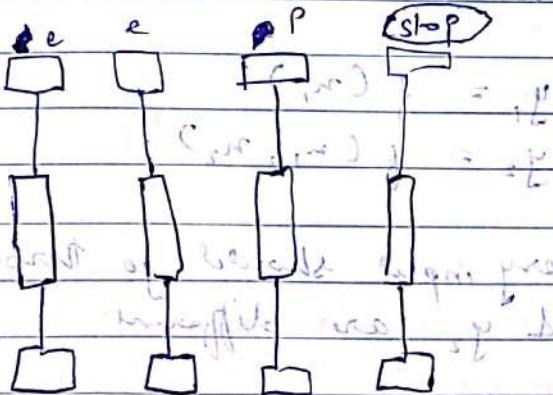
$\nabla C_{\text{style}} / \nabla v^T v$ more robust

$\nabla C_{\text{style}} / \nabla v^T v$ more robust

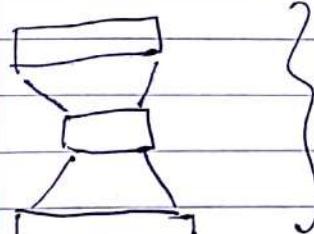
[more lossless update rule] [less noisy update rule]

WEEK 13

* Sequence Learning Problems

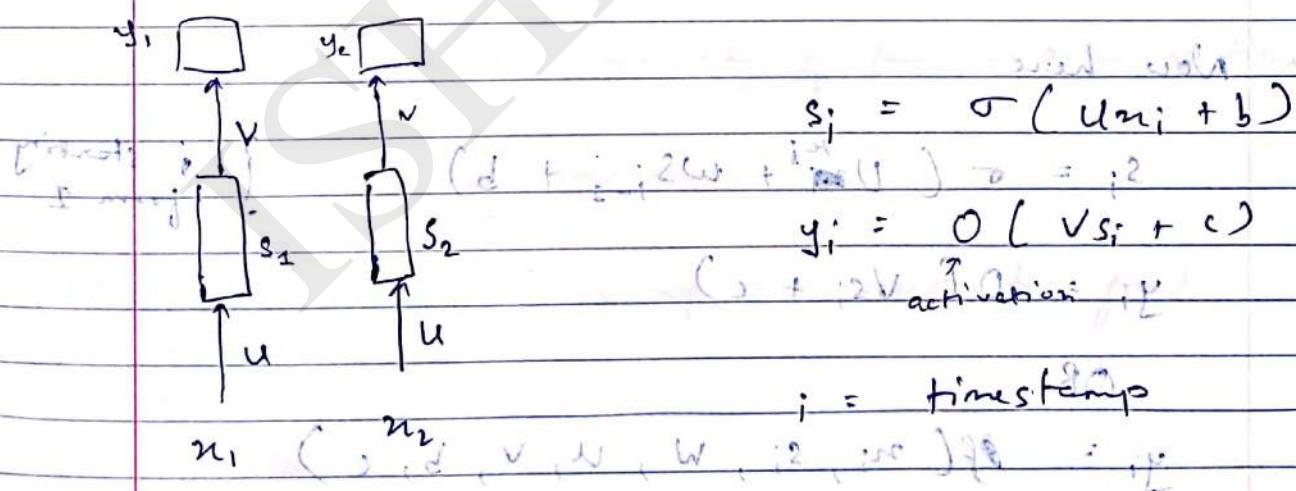


Predicting next character in a word depends on previous characters



Each one of them is a feed forward NN.

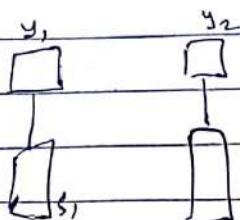
* Recurrent Neural Networks (RNN)



But this doesn't account for previous inputs

refer notes / week 13 / wishlist.py for state in

state in the loop the current state is v, w, b updated



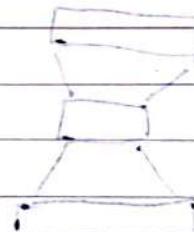
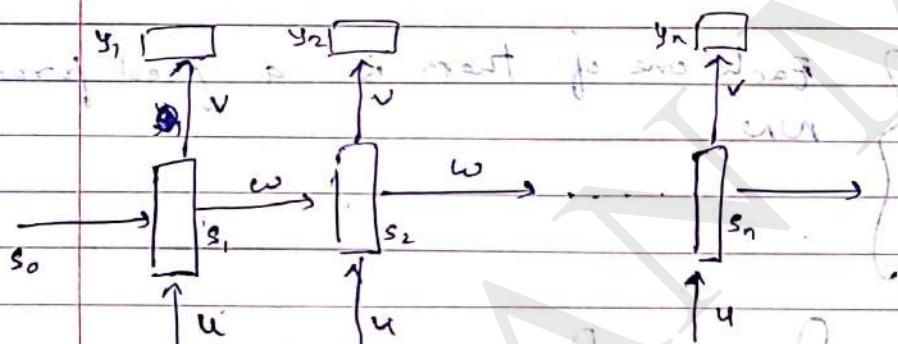
if we connect it like this?
then

$$n_1 \rightarrow n_2$$

as shared bias
not good enough.

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

→ Thus we make recurrent connections

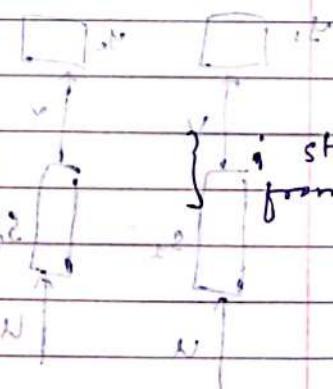


Now here

$$(1 + \text{exp})^{-1}$$

$$s_i = \sigma (U_{\text{in}}^{n_i} + w s_{i-1} + b) \quad \left. \begin{array}{l} \text{i starting} \\ \text{from 1} \end{array} \right\}$$

$$y_i = \text{f}(s_i)$$



OR

$$y_i = \text{f}(n_i, s_i, w, u, v, b, c)$$

Note

s_i is state of network at timestamp i
Weights w, u, v and biases b are shared across all timestamps

* Back Propagation Through Time (BPTT) (i)

lets say, $x_1 \rightarrow (0), 1, 6 \rightarrow (0), 1, 6$
 input with x_1, x_2, x_3

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

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

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

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

$$w \in \mathbb{R}^{d \times d}$$

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

(i) for s_1, s_2, s_3

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

- Activation at output layer = softmax

Loss

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

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

y_{tc} = predicted probability of true character at time-step;

T = number of timesteps

- Here during backprop we need to find gradients for

$$V, W, U \quad (1 + \epsilon) \rightarrow 1 + \epsilon$$

... $\rightarrow 1 + \epsilon$ with θ

... $\rightarrow 1 + \epsilon$ for gradient

(1) Gradients for Viterbi's mitigation tool.

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

test matrix

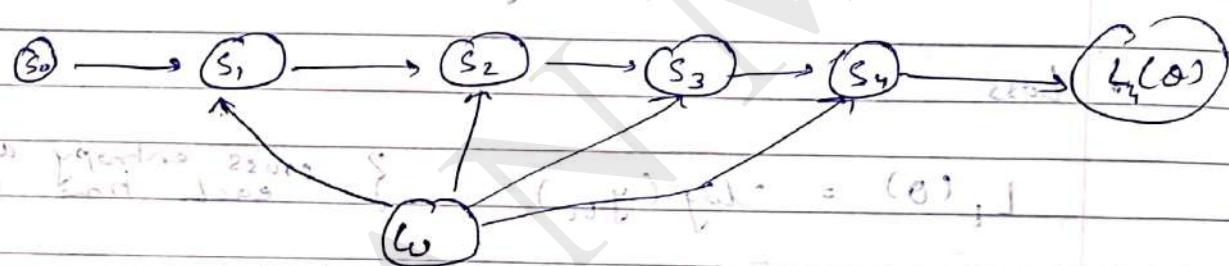
size 4×4

(2) Gradient w.r.t W

$$\frac{\partial L(\theta)}{\partial W} = \sum_{t=1}^T \frac{\partial L_t(\theta)}{\partial W} \quad \text{for gradient computation} \rightarrow$$

gradient of loss function w.r.t W

For the network = output layers of the architecture.



Now,

$$\frac{\partial L_4(\theta)}{\partial w} = \frac{\partial L_4(\theta)}{\partial s_4} \frac{\partial s_4}{\partial w}$$

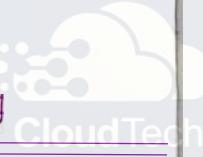
refer Backprop derivative is small for hidden layer w.r.t first hidden layer gradient w.r.t w

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

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

Here we have ignored σ for simplicity

DELL
PAGE NO.:
DATE:



$$\frac{\partial S_4}{\partial w} = \frac{\partial^+ S_4}{\partial w} + \underbrace{\frac{\partial S_4}{\partial S_3}}_{\text{Explicit}} \underbrace{\frac{\partial S_3}{\partial w}}_{\text{Implicit}}$$

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

$$= \frac{\partial^+ S_4}{\partial w} + \frac{\partial S_4}{\partial S_3} \left[\frac{\partial^+ S_3}{\partial w} + \underbrace{\left(\frac{\partial S_3}{\partial S_2} \frac{\partial S_2}{\partial w} \right)}_{\text{Explicit}} \underbrace{\left(\frac{\partial S_2}{\partial S_1} \frac{\partial S_1}{\partial w} \right)}_{\text{Implicit}} \right]$$

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

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

Leave out 3rd and 4th term

$$\frac{\partial S_4}{\partial w} = \sum_{k=1}^4 \frac{\partial S_4}{\partial S_k} \frac{\partial^+ S_k}{\partial w}$$

$$\therefore \frac{\partial L_4(\theta)}{\partial w} = \frac{\partial L_4(\theta)}{\partial S_4} \left[\sum_{k=1}^4 \frac{\partial S_4}{\partial S_k} \frac{\partial^+ S_k}{\partial w} \right]$$

$$(3) \quad [c_{10}^{(0)}, c_{11}^{(0)}, c_{12}^{(0)}, c_{13}^{(0)}, c_{14}^{(0)}] = 12$$

General

$$\frac{\partial L_f(\theta)}{\partial w} = \frac{\partial L_f(\theta)}{\partial S_2} \left[\sum_{k=1}^4 \frac{\partial S_k}{\partial S_2} \frac{\partial^+ S_k}{\partial w} \right]$$

(3) Gradient cost w
use chain rule $\nabla f(w) \rightarrow \nabla f(w)$ and find ∇

* Exploding & Vanishing Gradients

→ Now, we find

$$\frac{\partial S_i}{\partial S_{j-1}} = \frac{\partial S_i}{\partial S_k} \cdot \frac{\partial S_k}{\partial S_{k-1}} \cdot \dots \cdot \frac{\partial S_{j+1}}{\partial S_j}$$

$$\frac{\partial S_i}{\partial S_{j-1}} = \frac{\partial S_i}{\partial S_k} \cdot \frac{\partial S_k}{\partial S_{k-1}} \cdot \dots \cdot \frac{\partial S_{j+1}}{\partial S_j} = \prod_{j=k}^{i-1} \frac{\partial S_{j+1}}{\partial S_j}$$

$$\rightarrow \text{Thus } \frac{\partial S_i}{\partial S_{j-1}} = \frac{\partial S_i}{\partial S_1} \cdot \frac{\partial S_1}{\partial S_2} \cdot \dots \cdot \frac{\partial S_{j+1}}{\partial S_j}$$

$\frac{\partial S_i}{\partial S_{j-1}}$ is to be found

let,

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

$$S_j = \sigma(a_j)$$

$$a_j = [a_{j1}, a_{j2}, a_{j3}, \dots, a_{jd}]^T \quad \text{--- (1)}$$

$$S_j = [\sigma(a_{j1}), \sigma(a_{j2}), \dots, \sigma(a_{jd})]^T \quad \text{--- (2)}$$

→ Now

$$\frac{\partial S_j}{\partial S_{j-1}} = \frac{\partial S_j}{\partial a_j} \cdot \frac{\partial a_j}{\partial S_{j-1}} = \frac{\sigma'(a_j)}{w}$$

will be diagonal matrix $\times J$ less + will be w

from ① & ②

$$\frac{\partial \underline{s}_j}{\partial \underline{a}_j} = \begin{bmatrix} \frac{\partial s_{j1}}{\partial a_{j1}} & \frac{\partial s_{j2}}{\partial a_{j2}} & \dots \\ \vdots & \ddots & \vdots \\ 0 & 0 & 0 \end{bmatrix} \times w$$

multiplied

$$= \begin{bmatrix} 1 & \frac{\partial s_{j1}}{\partial a_{j1}} & \frac{\partial s_{j2}}{\partial a_{j2}} & \dots \\ \vdots & \sigma'(a_{j1}) & \sigma'(a_{j2}) & \vdots \\ 0 & 0 & 0 & \vdots \\ \vdots & \ddots & \ddots & \sigma'(a_{jd}) \end{bmatrix} \times w$$

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

Taking magnitude

$$\left\| \frac{\partial \underline{s}_j}{\partial \underline{a}_{j-1}} \right\| = \left\| \text{diag}(\sigma'(a_j)) w \right\|$$

using property of norm

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

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

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

\downarrow
max value taken by derivative
is $\frac{1}{4}$

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

$$\leq \gamma \|w\|$$

x

length

$$\leq \gamma \lambda$$

Now,

$$\left\| \frac{\partial s_t}{\partial s_k} \right\| = \left\| \prod_{j=k+1}^t \frac{\partial s_j}{\partial s_{j-1}} \right\|$$

$$\leq \prod_{j=k+1}^t \gamma \lambda$$

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

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

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

* Important

$$(\text{chain rule}) \quad \gamma^{t-1} \geq \gamma^{t-1}$$

$$\frac{\partial L_t(\theta)}{\partial w} = \frac{\partial L_t(\theta)}{\partial s_2} \sum_{k=1}^{t-1} \frac{\partial s_k}{\partial s_{k-1}} \frac{\partial s_k}{\partial w}$$

$d \times d$

✓

$d \times d$

✓

$d \times d$

gradient

froced values are known to us

$$(\text{and } n \approx f) \quad \gamma = 1 \geq \gamma^{t-1}$$

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

$$\frac{\partial^+ s_{kp}}{\partial w_{qr}} \quad ; (p, q, r) \text{ element of tensor}$$

Note: If a_i is diff wrt w are 0

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

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

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

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

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

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

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

~~last number from 1 to n in previous list~~

$$\frac{\partial \text{stepout}_p}{\partial w_{qr}} = \sum_{i=1}^n \frac{w_{pi} s_{t-1,i}}{w_{qr}}$$

$$\begin{matrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{matrix}$$

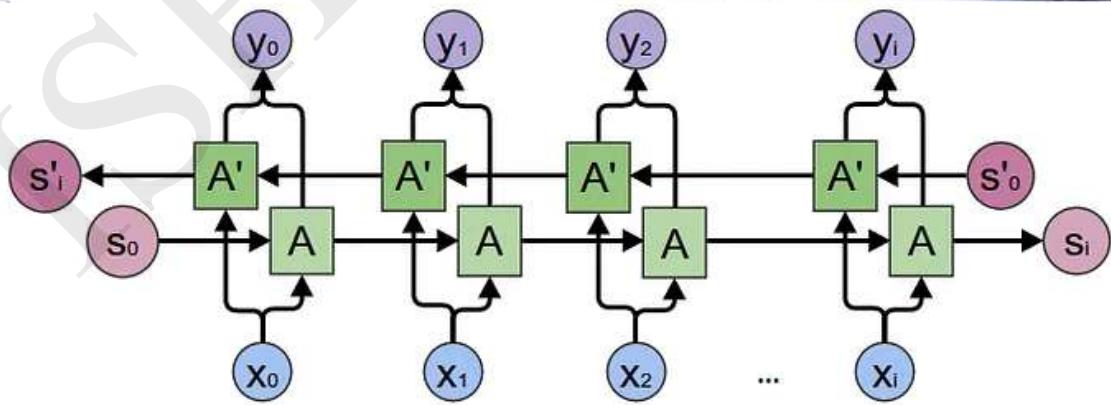
$$= s_{t-1,i} \quad \text{if } p = q \text{ and } i = r$$

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

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

$$\begin{cases} 1 & \text{if } t-1 = 0 \\ 0 & \text{otherwise} \end{cases}$$

* 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.

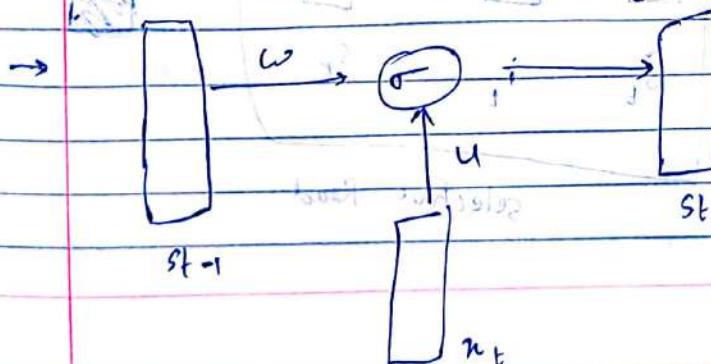
During Forward pass & Back pass info is tend to be forgotten.

→ Selective read, write, forget using white board analogy

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

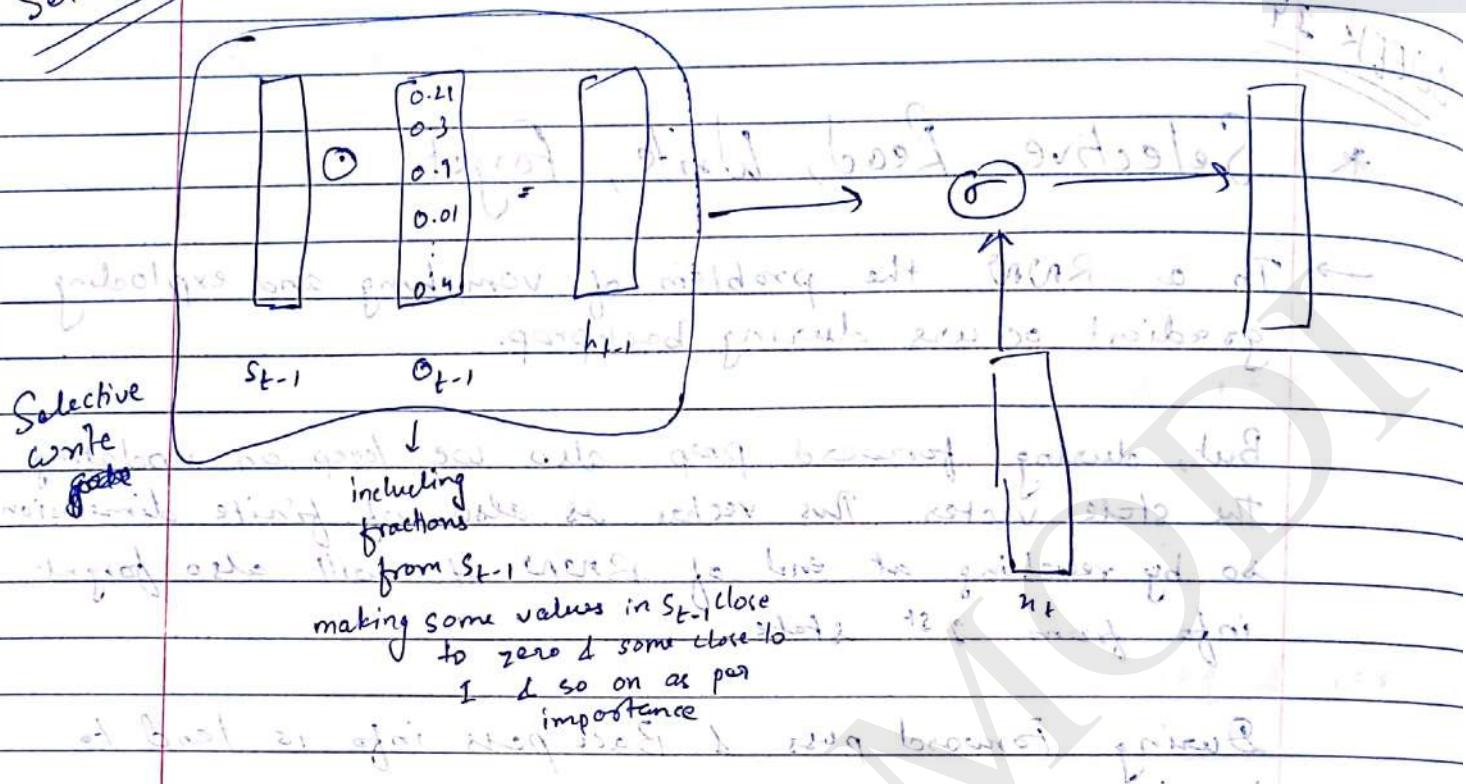
→ Let's consider task of sentiment analysis.

- 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



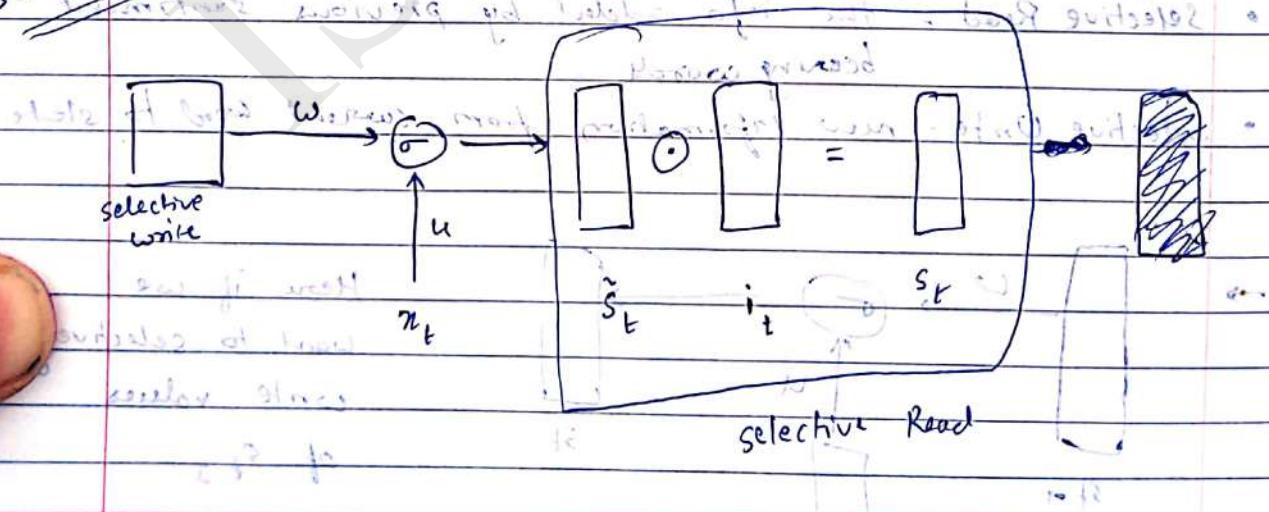
Now, we need to learn O_{t-1} (Output Gate)

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

$\rightarrow (W, V, b)$ are parameters of RNN

b_o, W_o are parameters of $RNN(O_{t-1})$

Selective Read



$$\tilde{s}_t = \sigma(w_{h_t} h_t + u_{x_t} x_t + b) \quad \text{Eq. 1}$$

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

$$i_t = \sigma(w_i h_{t-1} + u_i x_t + b_i) \quad \text{Eq. 2}$$

Note

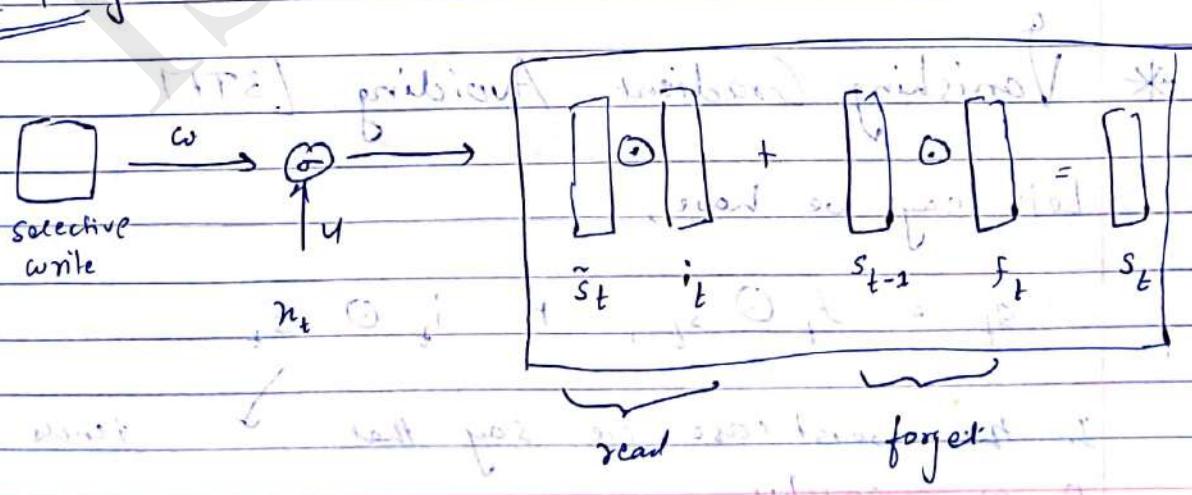
$\rightarrow w, u, b$ & w_i, u_i, b_i are different parameters

Till now

- Previous state $\rightarrow s_{t-1} = \sigma(w_{h_{t-1}} h_{t-1} + u_{x_{t-1}} x_{t-1} + b)$
- Output gate $\rightarrow o_{t-1} = \sigma(w_{o h_{t-1}} h_{t-1} + u_{o x_{t-1}} x_{t-1} + b_o)$
- Selective Write $\rightarrow h_{t-1} = o_{t-1} \odot \sigma(s_{t-1})$
- Current (temp) state $\rightarrow \tilde{s}_t = \sigma(w_{h_{t-1}} h_{t-1} + u_{x_t} x_t + b)$
- Input gate $\rightarrow i_t = \sigma(w_i h_{t-1} + u_i x_t + b_i)$

Selective Read $\rightarrow i_t \odot \tilde{s}_t$ will store forget value

Selective Forget



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

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

* GRU (Gated Recurrent Unit)

→ Gates

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

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

→ States

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

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

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 $(a) \rightarrow$ for vanishing all terms of s_t in s_{t+1} which prove backprop gradient vanishes with time.

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

$$\frac{\partial s_t}{\partial s_{t-1}} = \text{diag}(f_t)$$

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

$$= (f_t)^T$$

Here also during back pass f_t vanishes

But during forward pass

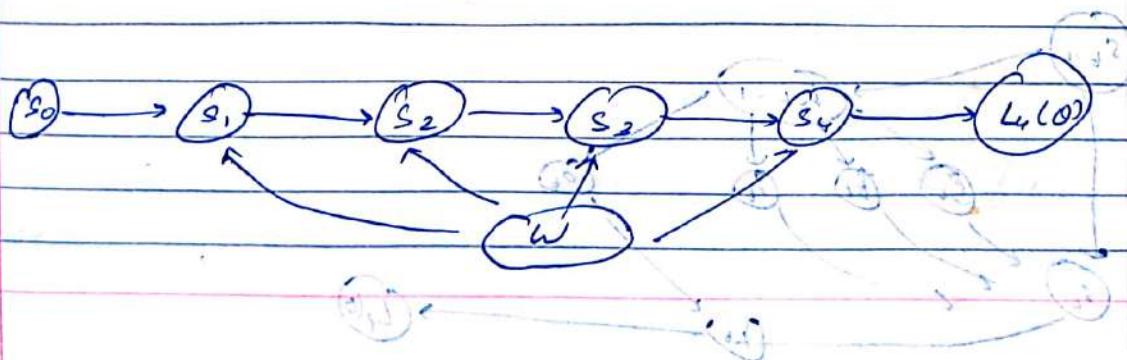
f_t information goes from 1 state to next

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

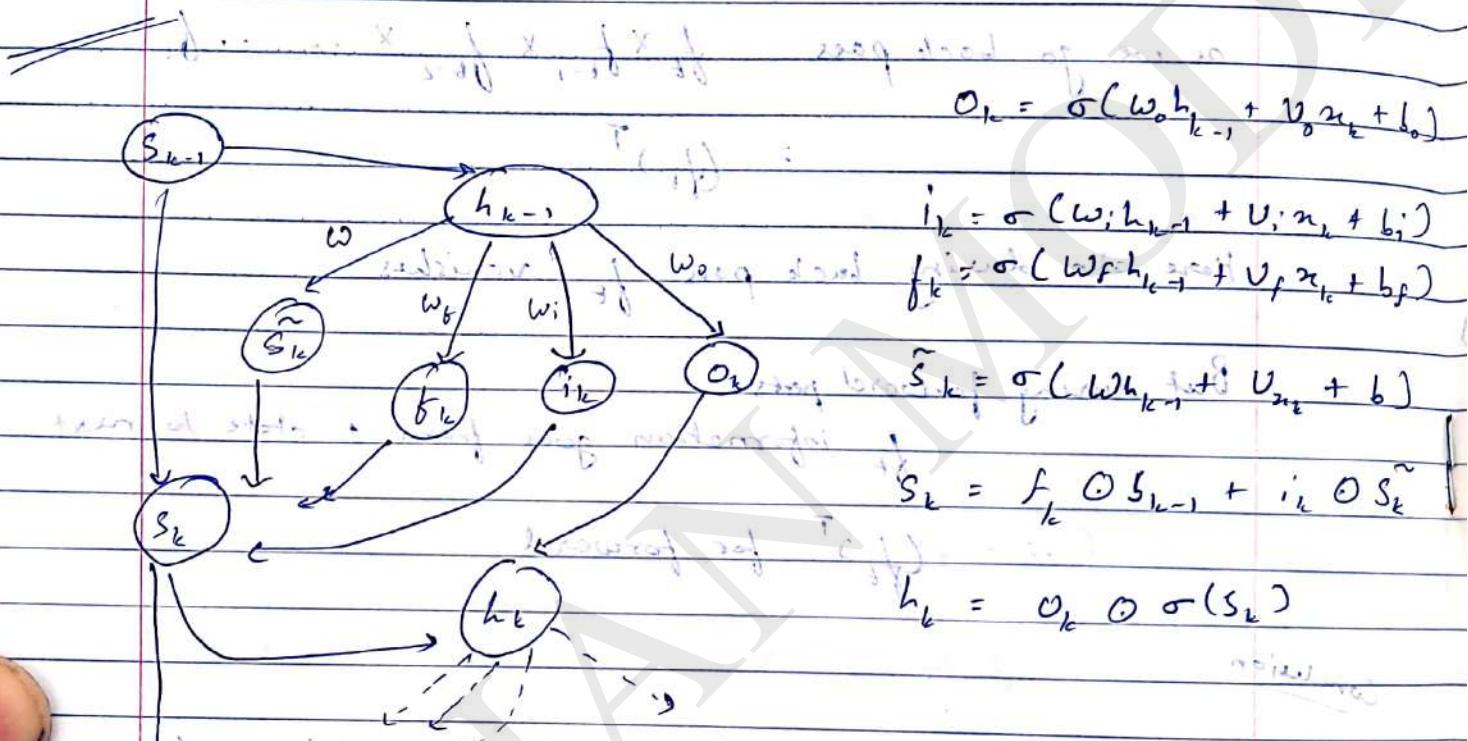
Conclusion

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

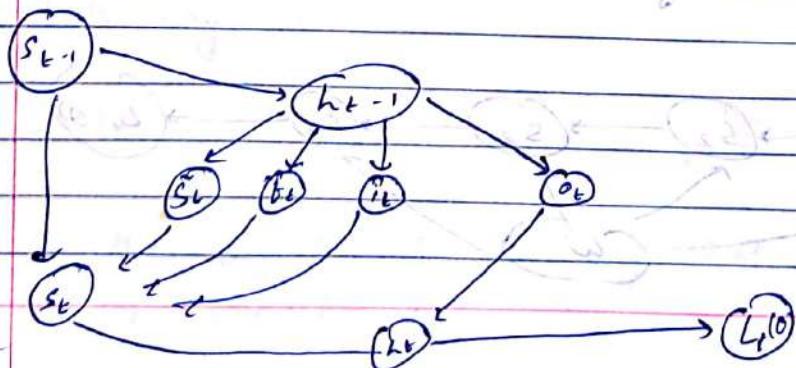
* Proof (Vanishing doesn't occur)



- 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.



At standard loss like L1-L2 and to some extent for cross entropy loss, it is well known that backpropagation of error is linear. This is because the derivative of the error function with respect to the weights is linear.



→ 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 doesn't vanish.

$$\frac{\partial L_t(\theta)}{\partial S_{1:t}} = \frac{\partial L_t(\theta)}{\partial h_t} \frac{\partial h_t}{\partial S_t} \left(\frac{\partial S_{1:t}}{\partial S_{t-1}} \dots \frac{\partial S_{1:t}}{\partial S_1} \right)$$

(~~longest~~
could
found) \rightarrow ~~grayscale~~ \downarrow (2)

$$\left(\begin{array}{cc} 06 & 16 \\ 16 & 06 \end{array} \right) = \cancel{16} \left(\begin{array}{cc} 06 & 16 \\ 16 & 06 \end{array} \right) \frac{(8)}{\cancel{16}} =$$

$$\textcircled{1} \quad \frac{\partial h_t}{\partial s_t} = \delta \frac{\partial_t \sigma(s_t)}{\partial s_t} = \text{Diagonal}(\sigma_t, \sigma_t' \underbrace{\sigma_t}_{\text{matrix}})$$

$$(3) \quad \frac{\partial s_t}{\partial s_{t-1}} = \frac{\partial f_t \circ s_{t-1}}{\partial s_{t-1}} + (i_t \circ s_t) \xrightarrow{\text{ignoring worst case}}$$

$$= D(f_t) \quad \text{!}$$

classical form diagonal matrix will be called as pentagonal

$\left(\frac{\partial S_i}{\partial S_{k+1}} \text{ and } -\frac{\partial S_{k+1}}{\partial S_i} \right) \text{ if } i \neq k$

1. $z_1 \leftarrow \sigma(s_1) \rightarrow z_2 \leftarrow \sigma(s_2) \rightarrow \dots \rightarrow z_t \leftarrow \sigma(s_t)$

$$\therefore \text{Let } \alpha_t \text{ be the error at the } t\text{-th layer} \rightarrow \\ = L_t'(h_t) \cdot D(\sigma_t \circ \sigma'(s_t))$$

Need to take next of $D(\sigma_{t+1}^t \circ f_t)$ will
mislead towards last of (σ_t) may stop here

Thus gradient vanish during backprop only if it
vanished in forward prop. If vanish only
when it needs to this pointer

* Proof (Exploding can be stopped)

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

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

$$= L_t'(h_t) \cdot (D(\sigma(s_t) \circ \sigma'(s_t)) \cdot w_o)$$

$$(D(\sigma(s_k) \circ \sigma'(s_k)) \cdot w_o)$$

$$\|z_t\| \leq \|L_t'(h_t)\| (\|w_o\| \|w_o\|) \dots$$

Depending on norm of $\|w_o\|$ gradient may explode

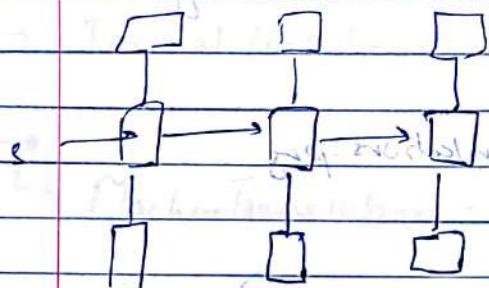
Thus use Gradient Clipping: If norm exceeds
certain value bring it down to a threshold.

~~WEEK 15~~ What will the next token be given previous tokens?

* Introduction to Encoder-Decoder models

→ 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

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

$y \in V$, V is vocabulary

softmax gives scores in size equal to vocab.

$$P(y_t=j | y_1^{t-1}) = P(y_t=j | s_t) \quad (s_t)$$

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

Thus model is

refer notes/week15/lec1/model.png

Links at <http://tiny.cc/meyarw>

Machine learning models

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/week15/dec1/testcase.py

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

refer notes/week15/dec1/representations.py

→ Image Captioning (Encoder - Decoder)

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

Here we have image as well

$P(y_t | s_t, I)$ → cannot take entire image
we need representation of image

$$P(y_t | s_t, f_{C_I}(I)) \quad \# \text{Fully connected layer in VGG.}$$

This becomes s_0

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

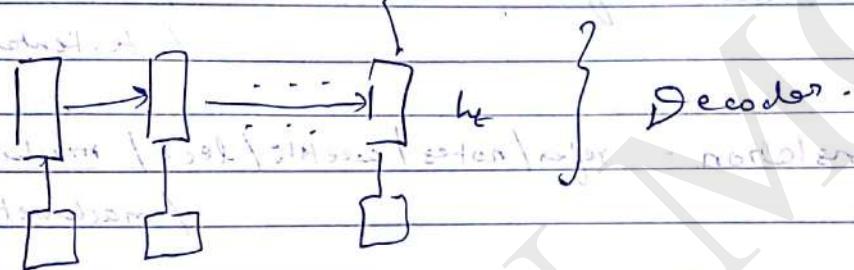
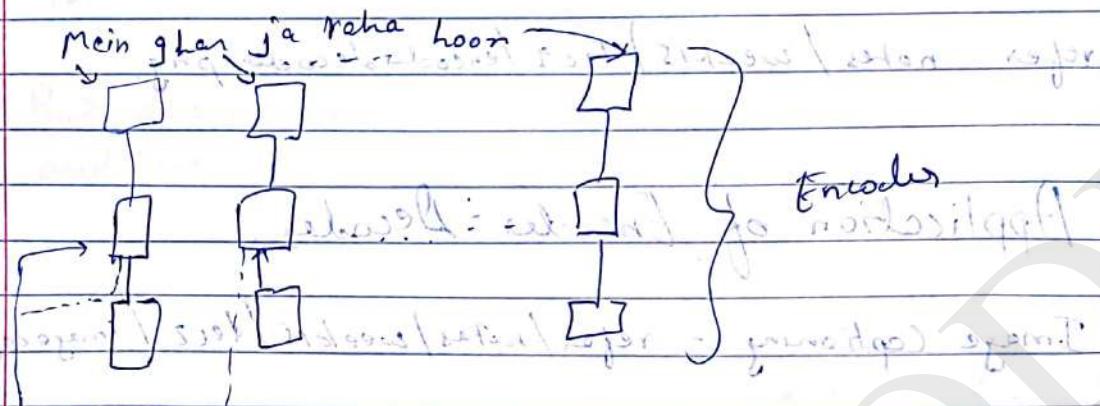
We can also feed in fig (ii) to all states.

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

* Attention Mechanism



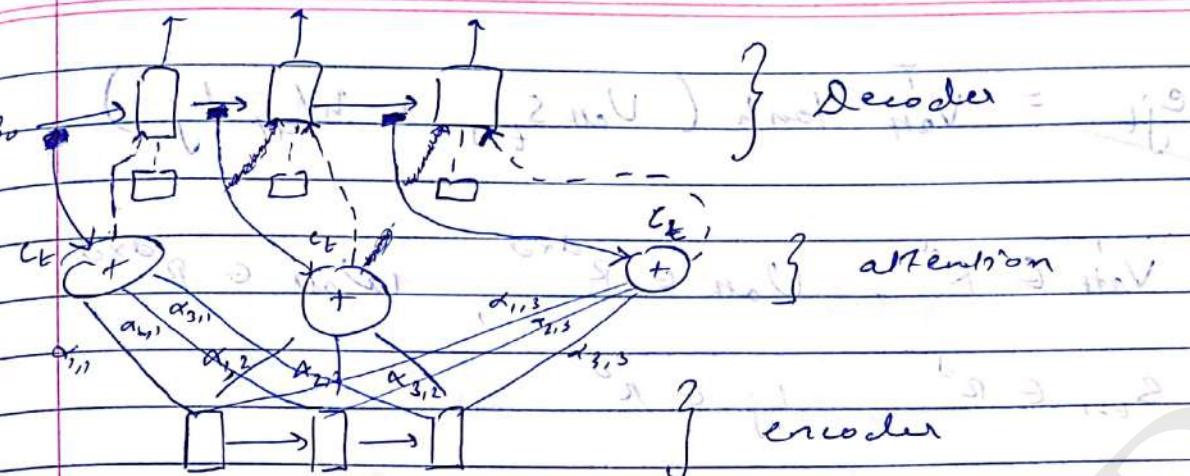
e.g. I am going home | (I am going home) - modafinil
e.g. I am going home |

So here we are sort of passing all the information to the next layer - So if you pass the entire network, this layer contains all the information ie (I am going home) but we can observe that h_t contains the words that are highly related to the input words.

home → Mein | ghar → ja raha → These are words that are highly related to the input words.

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 each input.

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~~

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

j^{th} word in encoder at time t^{th} stamp in decoder of j^{th} word in t^{th} time stamp.

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

thus α_{jt} is probability

Thus α_{jt} is measure of how important is j^{th} word with respect to t^{th} time stamp.

$$e_{jt} = V_{att}^T \tanh(V_{att} s_{t-1} + W_{att} h_j)$$

$V_{att} \in \mathbb{R}^{d \times d}$, $V_{att} \in \mathbb{R}^{d \times d}$, $W_{att} \in \mathbb{R}^{d \times d}$

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

This thing works because it is better modeling choice

Also, c_t is computed as

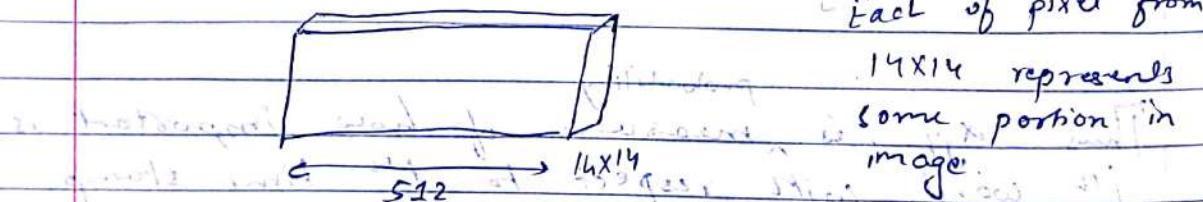
$$c_t = \sum_{j=1}^T \alpha_{jt} h_j \quad \text{for encoder.}$$

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

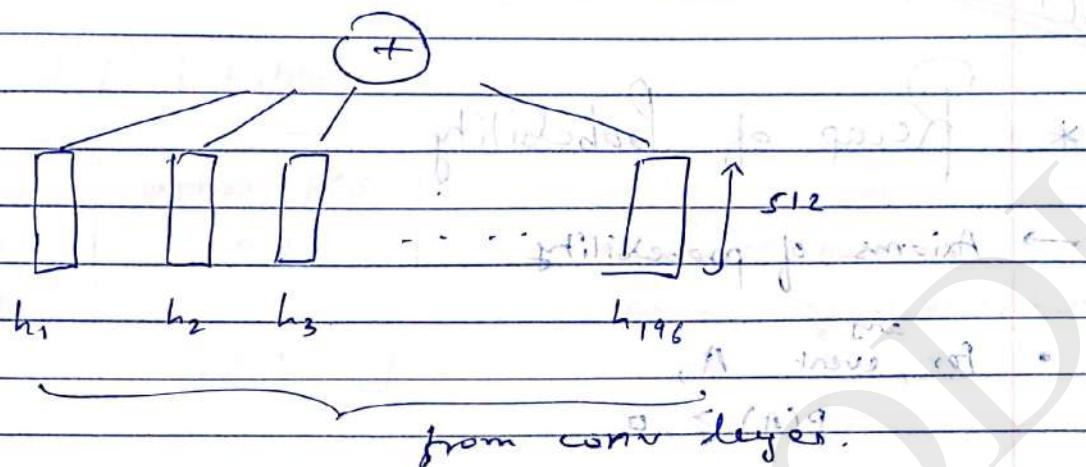
* Attention Over Images

Unlike RNN, where each state gives word context for image gives 14×14 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 goes to lower level prior state