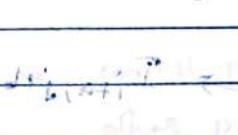


* Convolutional Neural Network (CNN)

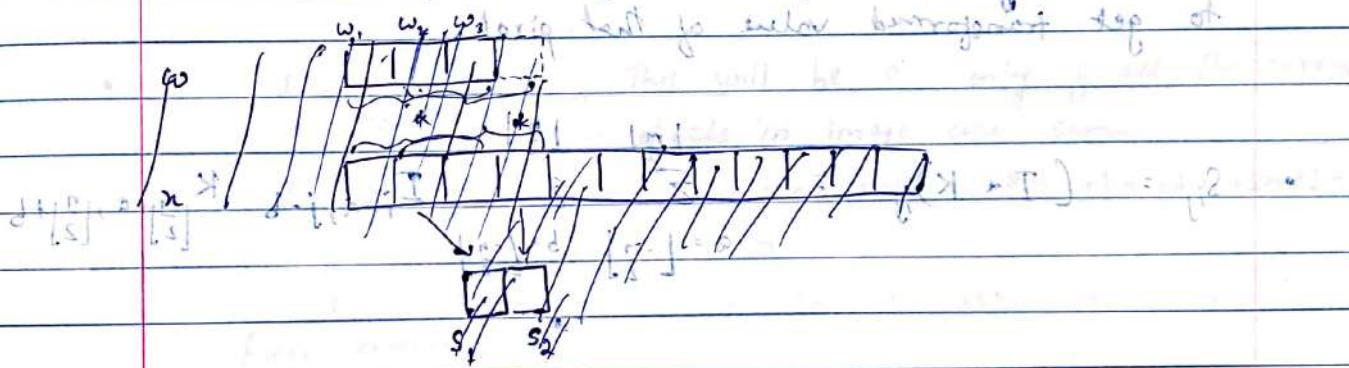
→ 1D CNN

$$s_t = \sum_{a=0}^{n-1} x_{t-a} w_a \quad \text{where } n = \text{size of filter}$$

input  filter 
 convolution

→ Let's say $n = 3$

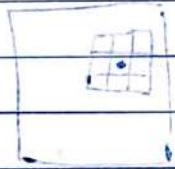
Let's say we are applying a filter of size 3 on a 5x5 input image.



Let's say we are applying a filter of size 3 on a 5x5 input image.

Let's say we are applying a filter of size 3 on a 5x5 input image.

$$\omega_1 \quad \omega_2 \quad \omega_3$$



Let's say we are applying a filter of size 3 on a 5x5 input image.

$$s_1 \quad s_2$$

Let's say we are applying a filter of size 3 on a 5x5 input image.

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

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

$$s_1 \quad s_2$$

$$s_1 \quad s_2$$

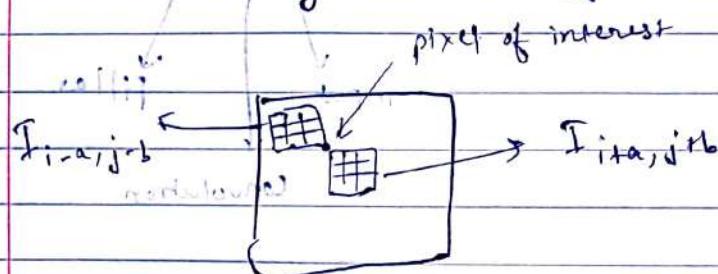
$$s_1 \quad s_2$$

Ishan Modi

→ 2D Convolutions (With Examples)

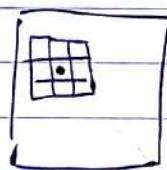
$$\cdot S_{ij} = (I * K)_{ij} = \sum_{a=0}^{m-1} \sum_{b=0}^{n-1} I_{i+a, j+b} K_{a,b}$$

(lets say we have image)



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

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



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

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

Ishan Modi

→ Types of Filters (sizing, average, median, etc.)

- $\begin{matrix} 1 & 1 & 1 \end{matrix}$ replacing this pixel in image with average of all pixels around it \rightarrow $\frac{1+1+1}{3} = 1$
- $\begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 1 \end{matrix}$ \rightarrow $\frac{1+1+1+1+1+1}{9} = 1$ \rightarrow weighted average

Blur: \rightarrow environments and the scene background will be blurred.

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

Sharpen: basic rule of blur is $r=9$. \rightarrow $r=9$ \rightarrow 0 increased with a factor of 9.

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

$$1 + 1 + 1 + (-8) + 1 + 1 + 1 = 1n + 1n = 0$$

between treatments, we have a smooth & good form of \rightarrow Edge detector

$$\text{channel mask} \quad \left[\begin{array}{ccc} 1 & 1 & 1 \\ 1 & -9 & 1 \\ 1 & 1 & 1 \end{array} \right] \quad \frac{1+1+1-9+1+1}{9} = 0$$

$$\rightarrow 3D \text{ Convolution} \quad \left[\begin{array}{ccc} 1 & 1 & 1 \\ 1 & -9 & 1 \\ 1 & 1 & 1 \end{array} \right] \quad \frac{1+1+1-9+1+1}{9} = 0$$

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

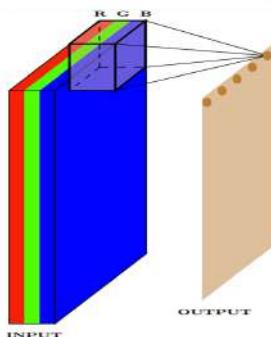
with assumption (1)

$$\left[\begin{array}{ccc} 1 & 1 & 1 \\ 1 & -9 & 1 \\ 1 & 1 & 1 \end{array} \right] \rightarrow \text{EXTRACTOR}$$

(2)

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

$3 \times 3 \times 3$



In simple terms, which is a 3D feature map

Ishan Modi

* 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 $W_2 \times H_2$

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

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

→ In ^{almost} _{all} cases we might have to add padding (^{Transform all pixels by taking them at ~~very less cost~~ ~~nearby cost~~ ~~nearby 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

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

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

Main formula

(Q)

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

② no. of neurons

∴ 2 images giving 96 filter will give 96 output.

$$W_2 \times H_2 \times D$$

Ans since filter depth & img depth are same

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

Ishan Modi

Since there are 96 filters (in working)

$$D = 96$$

→ given $s=4$, $p=0$

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

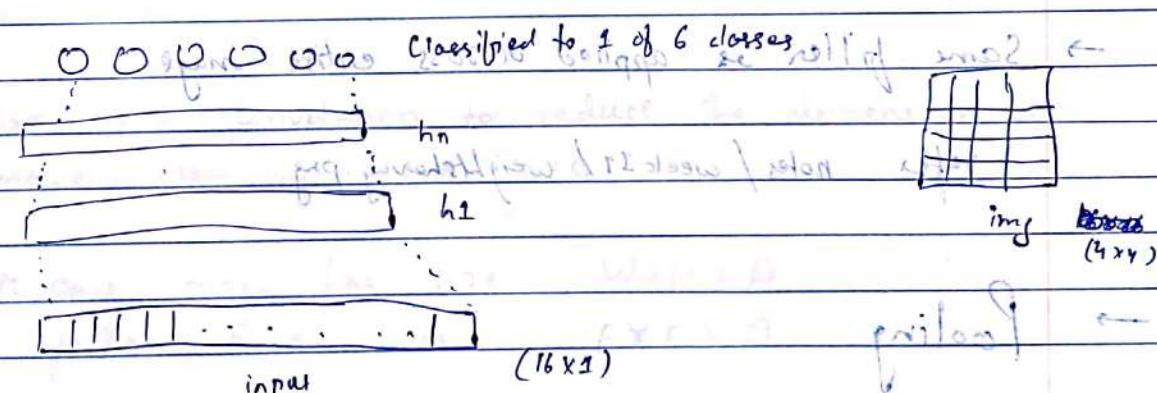
$$\therefore h_2 = \frac{227 - 11}{4} + 1 = 55$$

* 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.

Ishan Modi

Wipro Next Labs Graduate Researcher
Data Science

→ Convolution NN

2	4	7
2	5	8
3	6	9

⇒

2	3	4	5	6	7	8	9

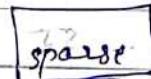
3x3 img.

3P = 0

9x12 vector

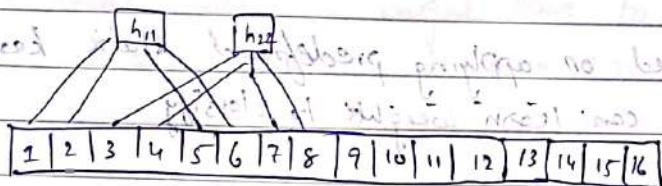
(1)

so this is comparatively sparse



E + N = 558

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



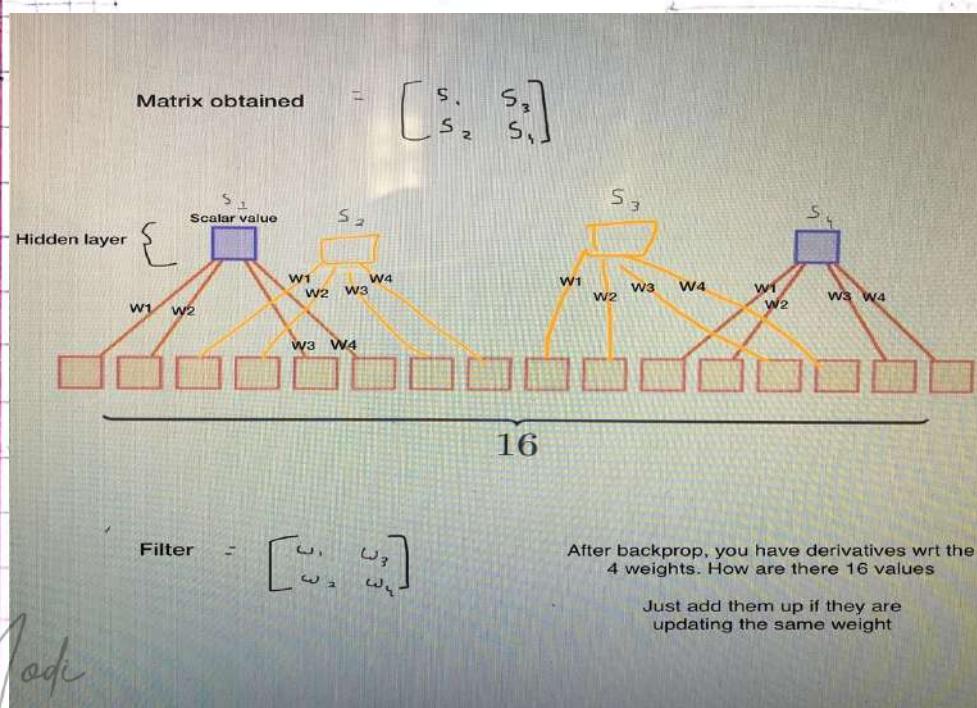
work and recall, writing back forward as backward process cost

forward pass of 24 h1

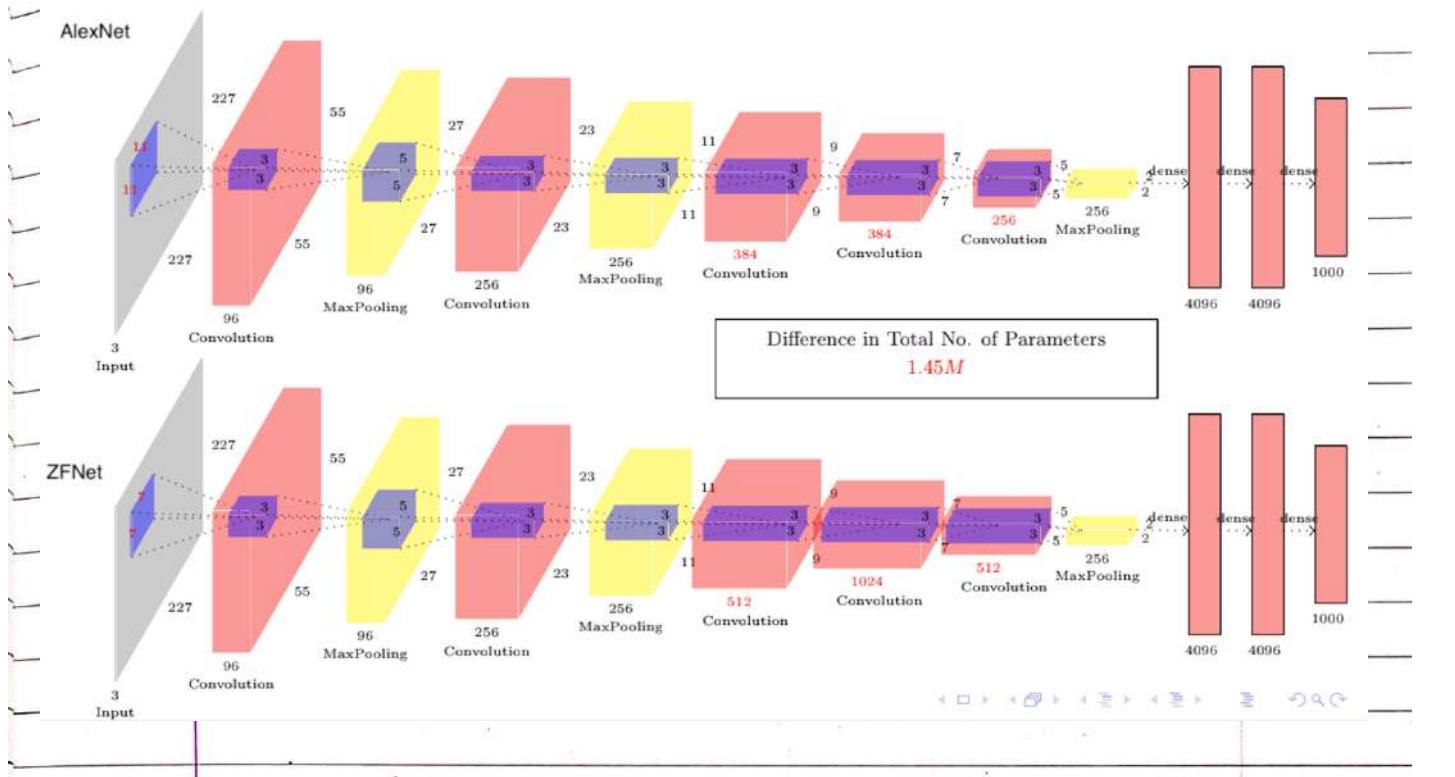
(2)

Weight Sharing

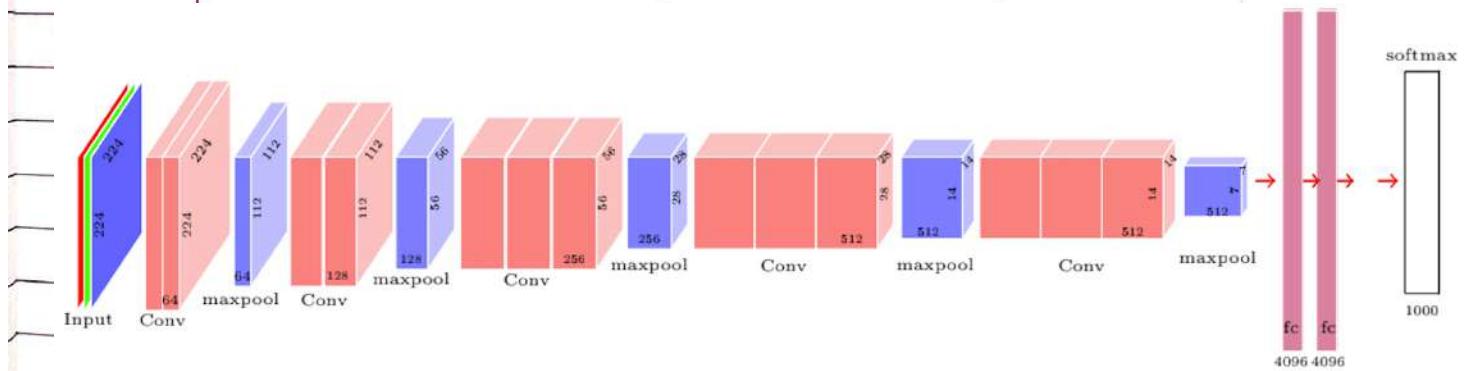
→ Same filter is applied across entire image



---> ZFNet



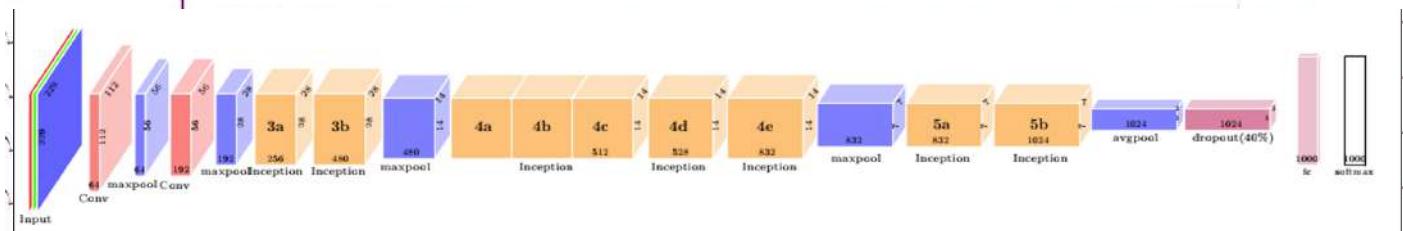
---> VGGNet



- Kernel size is 3×3 throughout
- Total parameters in non FC layers = $\sim 16M$
- Total Parameters in FC layers = $(512 \times 7 \times 7 \times 4096) + (4096 \times 4096) + (4096 \times 1024) = \sim 122M$
- Most parameters are in the first FC layer ($\sim 102M$)

Ishan Modi

---> 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 $F \times F \times D$

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

multiplied by each element of o/p

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

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

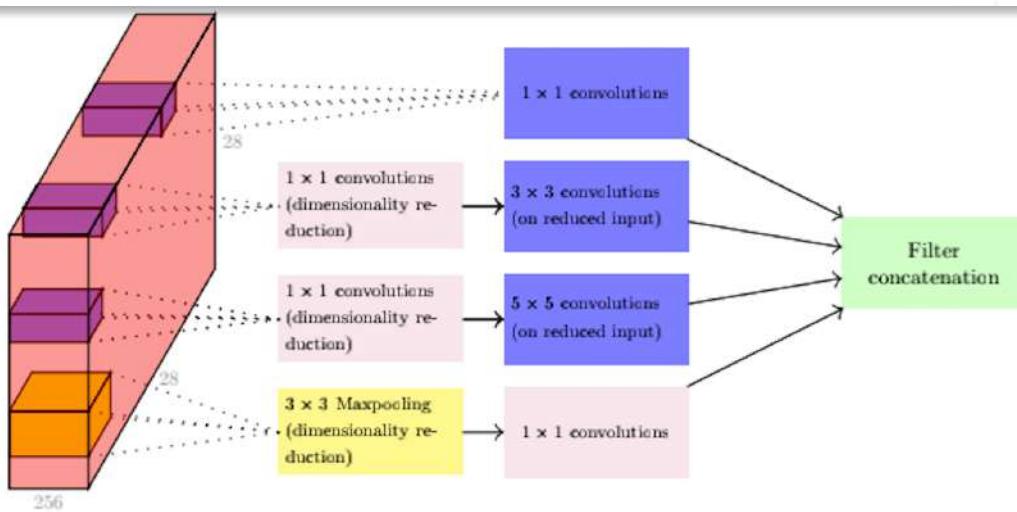
If $D' \ll D$, computations are reduced

now applying filter size $F \times F \times D'$

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

Ishan Modi

---> GoogleNet has various inception modules with basic structure as follows



Filter concatenation adds up all feature maps after applying appropriate padding

- In the last convolution layer

$$\begin{matrix} 7 \times 7 \times 1024 \\ \xrightarrow{\text{avg pool}} \end{matrix} \Rightarrow \underbrace{1 \times 1 \times 1024}_{\downarrow}$$

equivalent to a vector can be connected with dense layer with less parameters than earlier

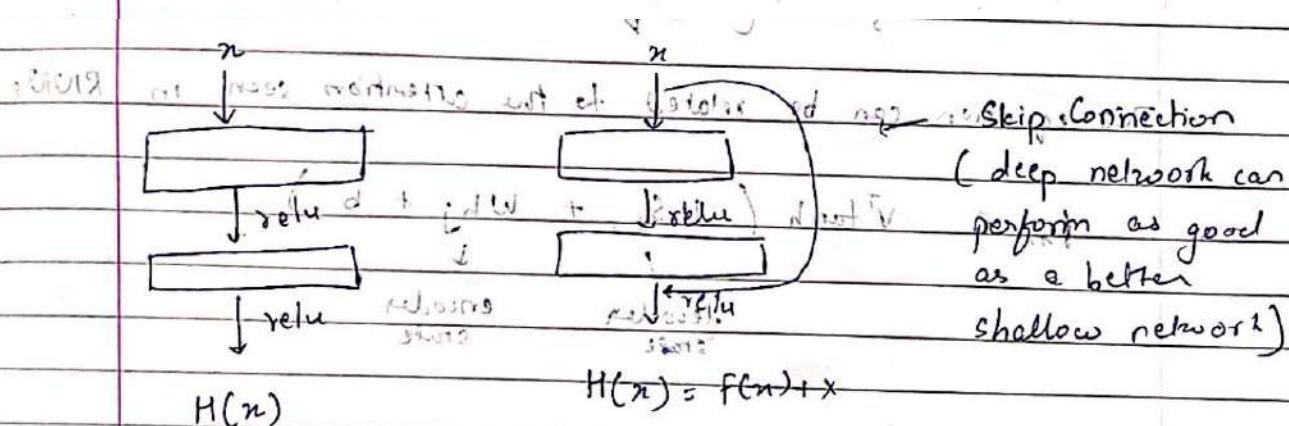
---> ResNet



Bag of tricks

- Batch Normalization after every CONV layer
- Xavier/2 initialization from [He et al]
- SGD + Momentum(0.9)
- Learning rate: 0.1, divided by 10 when validation error plateaus
- Mini-batch size 256
- Weight decay of 1e-5
- No dropout used

ResNet, 152 layers



Arshita Subedi (S)

In residual network, in forward pass current layer that has finer features is aggregated with previous layers [of 2 or 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.

Ishan Modi

~~WEEK 12~~~~Dimensionality reduction~~

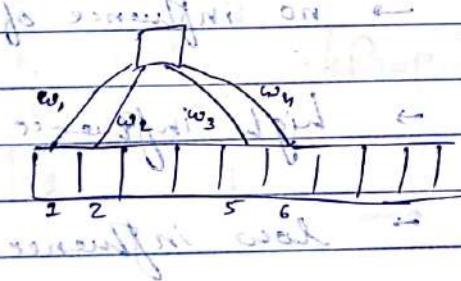
* Visualizing Patches

→ 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 bias.

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



$$\begin{bmatrix} w_1 & w_2 \\ w_3 & w_4 \\ w_5 & w_6 \end{bmatrix} \Rightarrow [w_1, w_2, w_3, w_4, w_5, w_6]$$

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

Choosing all weights to be zero except one.

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

Ishan Modi

15
Date: 6/6/2023

* 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
dpf has a very less impact on op. i.e. it is doing a
good job

Similarly do it for all patches to understand important
patches for classification. To profit 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 $\frac{\partial h_j}{\partial x_i} = 0$ → no influence of x_i on h_j

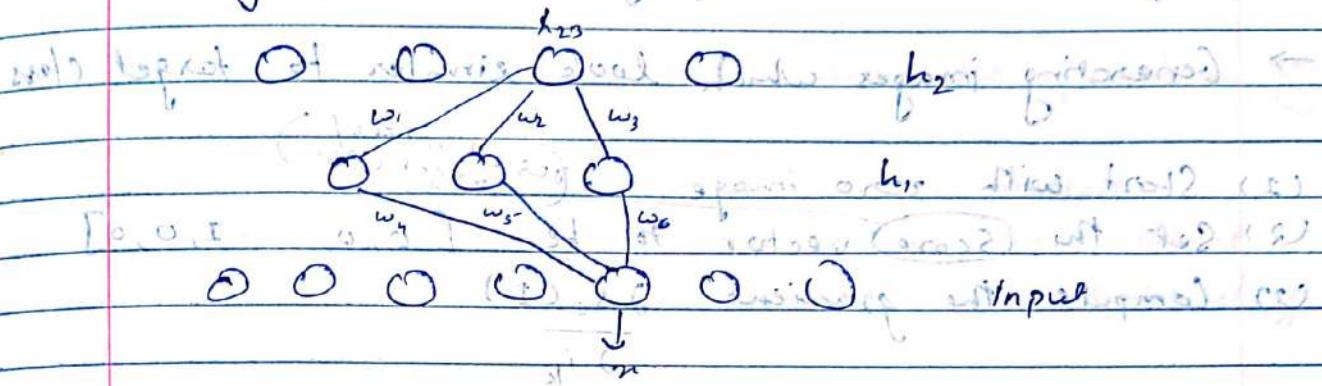
$\frac{\partial h_j}{\partial x_i}$ = large → high influence

$\frac{\partial h_j}{\partial x_i}$ = small → 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

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



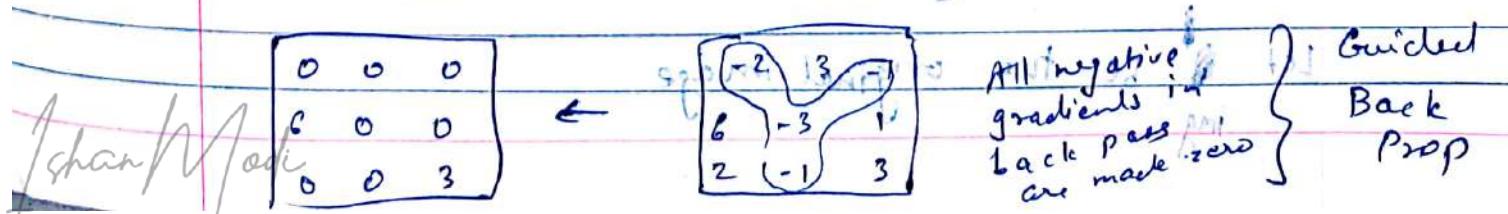
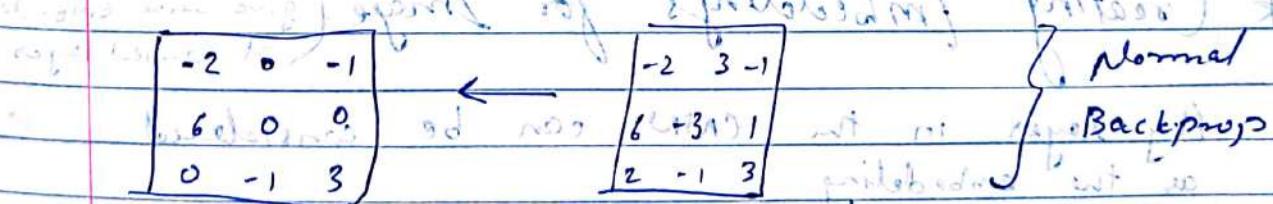
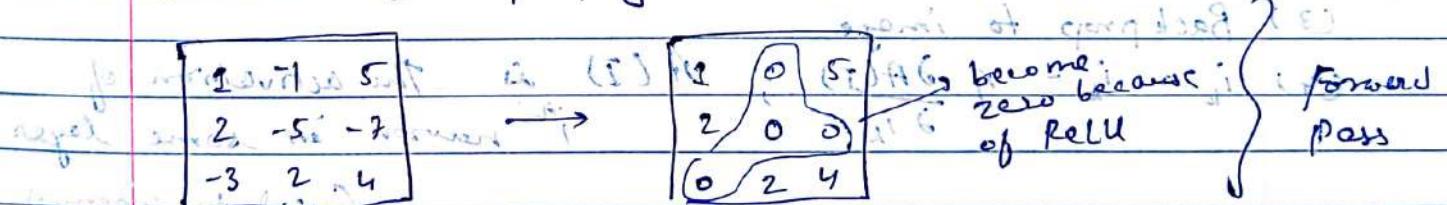
so we want to define $\frac{\partial h_{23}}{\partial x}$ w.r.t. Input (i)

$$\frac{\partial h_{23}}{\partial x} = \frac{\partial h_{23}}{\partial h_{11}} \frac{\partial h_{11}}{\partial x} + \frac{\partial h_{23}}{\partial h_{12}} \frac{\partial h_{12}}{\partial x} + \frac{\partial h_{23}}{\partial h_{13}} \frac{\partial h_{13}}{\partial x}$$

$$= w_1 w_3 + w_2 w_5 + w_3 w_6$$

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

* Guided BackPropagation (impact of image pixels on result)



Ishan Modi

* 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

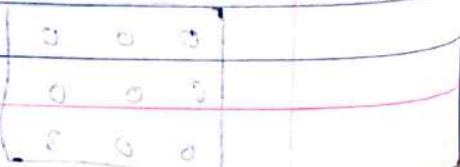
(4) $i_k = i_k + \eta \frac{\partial A(j)}{\partial i_k}$, if $A(j)$ is the activation of j^{th} neuron in some layer

* Creating embeddings for Image

Any layer in the CNN can be considered as the embedding

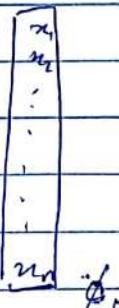
Let img be the original image

Ishan Modi



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

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

obtained embedding required embedding

update rule

$$\frac{\partial L(c_i)}{\partial u}$$

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

Similar is possible for convolution layers as well.

* Deep Dream

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

Ishan Modi

So if h_{ij} fires a neuron in convolution layer,

the neuron to top would learn what is the context we want to maximize $L(I)$

better and better context

thus, loss maxed $L(I)$ when

gradient where h_{ij} is also for

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

During backprop. update pixel of image

function of net loss again relation to what we want

gradient $\frac{\partial L(I)}{\partial I}$ is $\frac{\partial L(I)}{\partial h_{ij}}$ ~~if h_{ij} has best down in~~

~~any dimm (\rightarrow) with h_{ij} as g_i in $g_i \otimes g_j$ $\otimes g_k$ $\otimes g_l$ $\otimes g_m$ $\otimes g_n$ $\otimes g_o$ $\otimes g_p$ $\otimes g_q$ $\otimes g_r$ $\otimes g_s$ $\otimes g_t$ $\otimes g_u$ $\otimes g_v$ $\otimes g_w$ $\otimes g_x$ $\otimes g_y$ $\otimes g_z$~~

update rule

$$\| \text{within} \| = 1 \text{ in } \frac{\partial L(I)}{\partial h_{ij}} \text{ w.r.t. } \| = (.)$$

~~within~~ $\frac{\partial L(I)}{\partial h_{ij}}$ $\frac{\partial h_{ij}}{\partial \text{pixel}}$

the slope

→ 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.

also makes a relation with the other neurons

Ishan Modi

* Deep Art

referred book: DCGAN

→ 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 assigned ~~are~~ ~~new~~ embeddings generated using new image are similar to embeddings generated by actual image

(2) Style (cartoon)

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

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

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

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

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

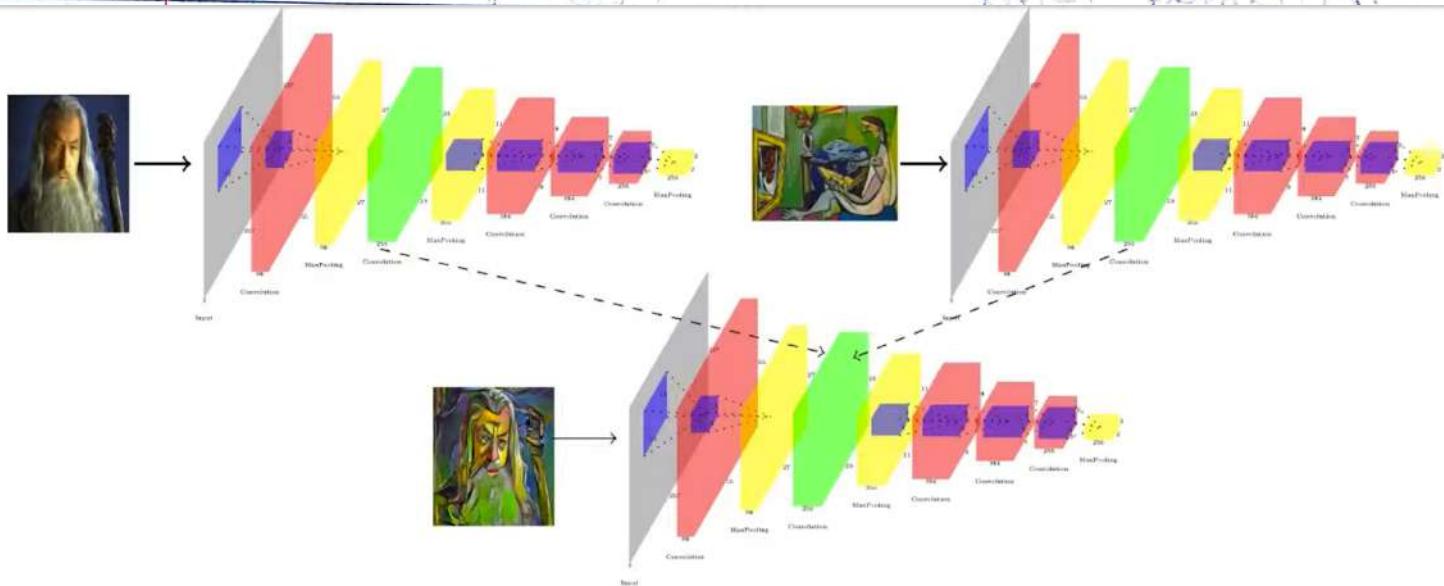
[Style matrix generated from embedding obtained from original image]

[Style matrix generated from embedding obtained from new image]

Ishan Mod

→ Total loss function

$$L_{total}(\vec{p}, \vec{a}, \vec{x}) = \alpha L_{content}(\vec{p}, \vec{x}) + \beta L_{style}(\vec{a}, \vec{x})$$



- The total loss is given by :-

$$\mathcal{L}_{total}(\vec{p}, \vec{a}, \vec{x}) = \alpha \mathcal{L}_{content}(\vec{p}, \vec{x}) + \beta \mathcal{L}_{style}(\vec{a}, \vec{x})$$

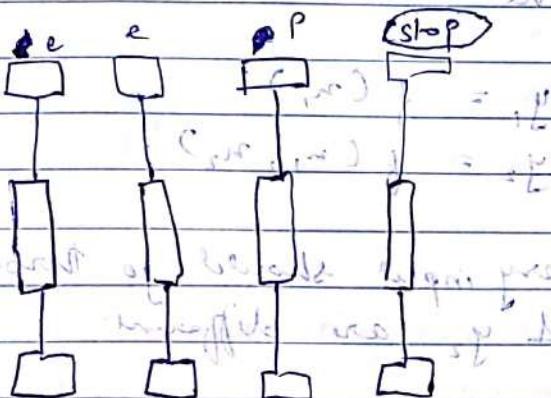
* Fooling CNNs

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

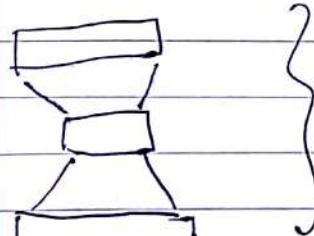
Ishan Modi

WEEK 13

* Sequence Learning Problems

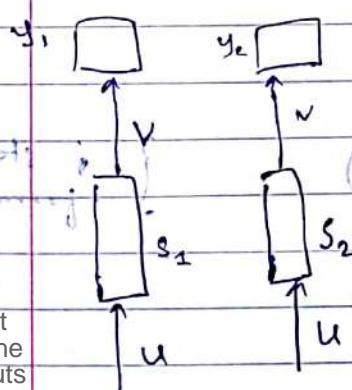


Predicting next character in a word depends on previous characters



Each one of them is a feed forward NN.

* Recurrent Neural Networks (RNN)



$$s_i = \sigma(u s_i + b)$$

$$(d + i \cdot 2w + b) \rightarrow i = 1, 2$$

$$y_i = o(v s_i + c)$$

$c = 2w$ activation of v

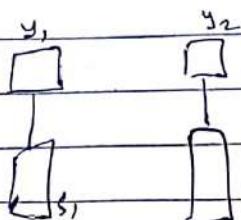
$i = \text{timestamp}$

$$n_1 (n_2 (d, v, w, b, w)) \rightarrow i = 1$$

This does not account for the previous inputs and thus, does not satisfy the wishlist

Wishlist

- Account for dependence between inputs
- Account for variable number of inputs
- Make sure that the function executed at each time step is the same
- We will focus on each of these to arrive at a model for dealing with sequences



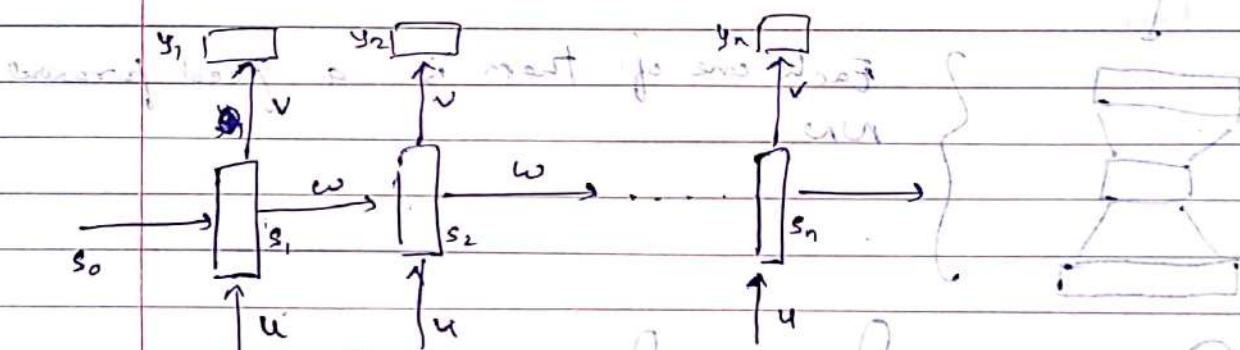
and if we connect it like this?

then what is?

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

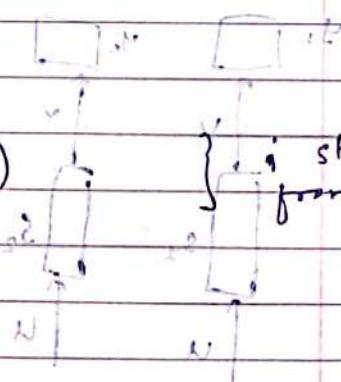


Now here

$$(s_i + v) \rightarrow = \dots$$

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

$$y_i = f(v s_i + c)$$



OR

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

Note

s_i is state of network at timestamp i

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

Ishan Modi

* Back Propagation Through Time (BPTT)

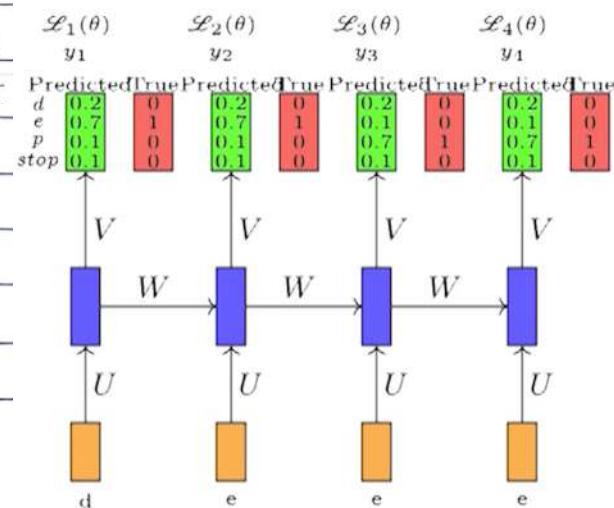
lets say,

input words are

$$x_i \in R^n \quad U \in R^{n \times d}$$

$$s_i \in R^d \quad W \in R^{d \times d}$$

$$y_i \in R^k \quad V \in R^{d \times k}$$



Here in this example we are predicting the next character in the word using RNN

Activation at output layer = softmax

Loss

$$L_t(\theta) = -\log(y_{t+1}) \quad \left. \begin{array}{l} \text{cross entropy} \\ \text{at each time stamp} \end{array} \right.$$

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

y_{t+1} = predicted probability of true character at time-step;

T = number of timesteps

→ Here during backprop we need to find gradients for

V, W, U ($t+1, t, t-1$)

... $t-2$... $t-1$ t $t+1$ $t+2$... $t+K$

... $t-2$... $t-1$ t $t+1$ $t+2$... $t+K$

Ishan Modi

(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

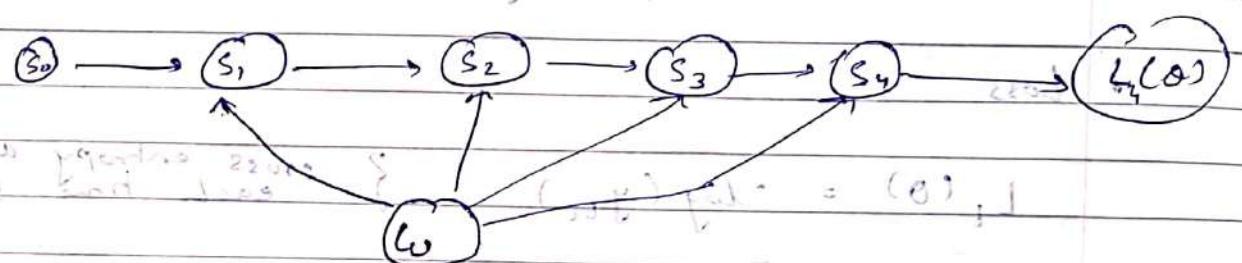
size 4×4

(2) Gradient w.r.t W

$$\frac{\partial L(\theta)}{\partial w} = \sum_{t=1}^T \frac{\partial L_t(\theta)}{\partial w} \quad \text{for gradient w.r.t weight} \rightarrow$$

gradient of loss function w.r.t weight

For the network = output layers of the architecture.



Now,

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

we will do $\frac{\partial L_n(\theta)}{\partial s_4}$ and $\frac{\partial s_4}{\partial w}$ both using backpropagation

refer Backprop $\left. \begin{array}{l} \text{derivative is small} \\ \text{w.r.t hidden layer} \\ \text{first hidden layer} \end{array} \right\}$ for minimum $\Rightarrow \nabla$
using of backprop

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

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

Ishan Modi

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)

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

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

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

Converg set of si.

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

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

$$(3) \quad [c_{10}^{(0)} \partial \dots \partial c_{10}^{(0)}, c_{10}^{(0)} \partial \dots \partial c_{10}^{(0)}] = 0$$

General

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

(3) Gradient cost w respect to w

Ishan Modi use chain rule know L v i w and find α



Scanned with OKEN Scanner

* Exploding & Vanishing Gradients

→ Now, we find

$$\frac{\partial S_i}{\partial S_j} = \frac{\partial S_i}{\partial S_k} \cdot \frac{\partial S_k}{\partial S_l} \cdot \frac{\partial S_l}{\partial S_m} \dots \frac{\partial S_{i-1}}{\partial S_j}$$

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

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

$\frac{\partial S_i}{\partial S_j}$ is to be found

let,

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

$$S_j = \sigma(a_j)$$

$$a_j = [w a_{j-1}, w a_{j-1}, w a_{j-1}, \dots, w a_{j-1}] \cdot a_{j-1} - ①$$

$$S_j = [\sigma(a_{j-1}), \sigma(a_{j-1}), \dots, \sigma(a_{j-1})] - ②$$

→ Now

$$\frac{\partial S_j}{\partial S_{j-1}} = \frac{\partial S_j}{\partial a_j} \cdot \frac{\partial a_j}{\partial S_{j-1}}$$

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

Ishan Modi

from ① & ②

$$\frac{\partial \delta_j}{\partial a_j} = \begin{bmatrix} \frac{\partial \delta_{j1}}{\partial a_{j1}} & \frac{\partial \delta_{j2}}{\partial a_{j2}} & \dots \\ \vdots & \ddots & \vdots \\ 0 & 0 & 0 \end{bmatrix} \times w$$

multiplied

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

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

Taking magnitude

$$\left\| \frac{\partial \delta_j}{\partial a_{j-1}} \right\| = \left\| \text{diag}(\sigma'(a_j)) w \right\|$$

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 $\frac{1}{4}$

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

Ishan Modi

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

$\omega \times$

$\gamma \lambda$

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

$$(\text{using } \frac{\partial L_t(\theta)}{\partial w}) \gamma^{t-1} \geq c_{t-1}^{1-\alpha}$$

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

$d \times d$

✓

$d \times d$

✓

$d \times d$

gradient of weight

froced values are known to us

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

Ishan Modi

→ 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 $a_i = \sum_j w_{ij} s_{jkl}$, then if w_{ij} are 0, then a_i is 0.

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

$$s_k = \sigma(a_k)$$

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

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

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

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

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

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

Ishan Modi

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

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

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

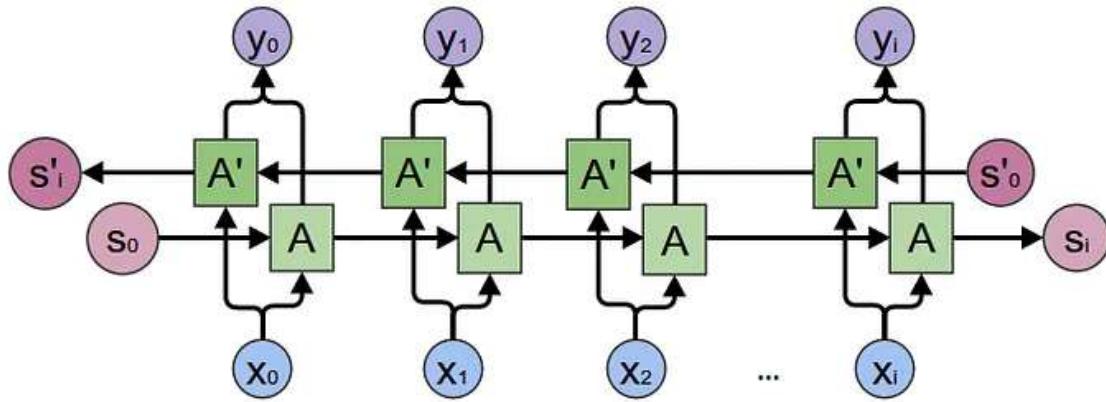
$$= s_{t-1,i} \quad \text{if } p = q \text{ and } i = r$$

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

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

$$\begin{cases} 1 & \text{if } t-1 = 0 \\ 0 & \text{otherwise} \end{cases} \quad \begin{matrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{matrix} \rightarrow \begin{matrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{matrix}$$

* 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.

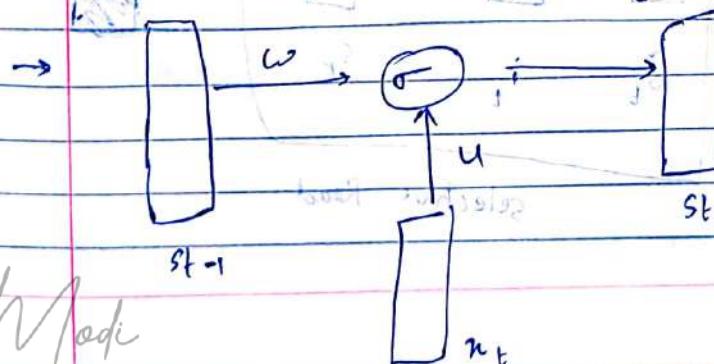
During Forward pass & Back pass info is tend to be forgotten.

→ Selective read, write, forget using white board analogy

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

→ Let's consider task of sentiment analysis.

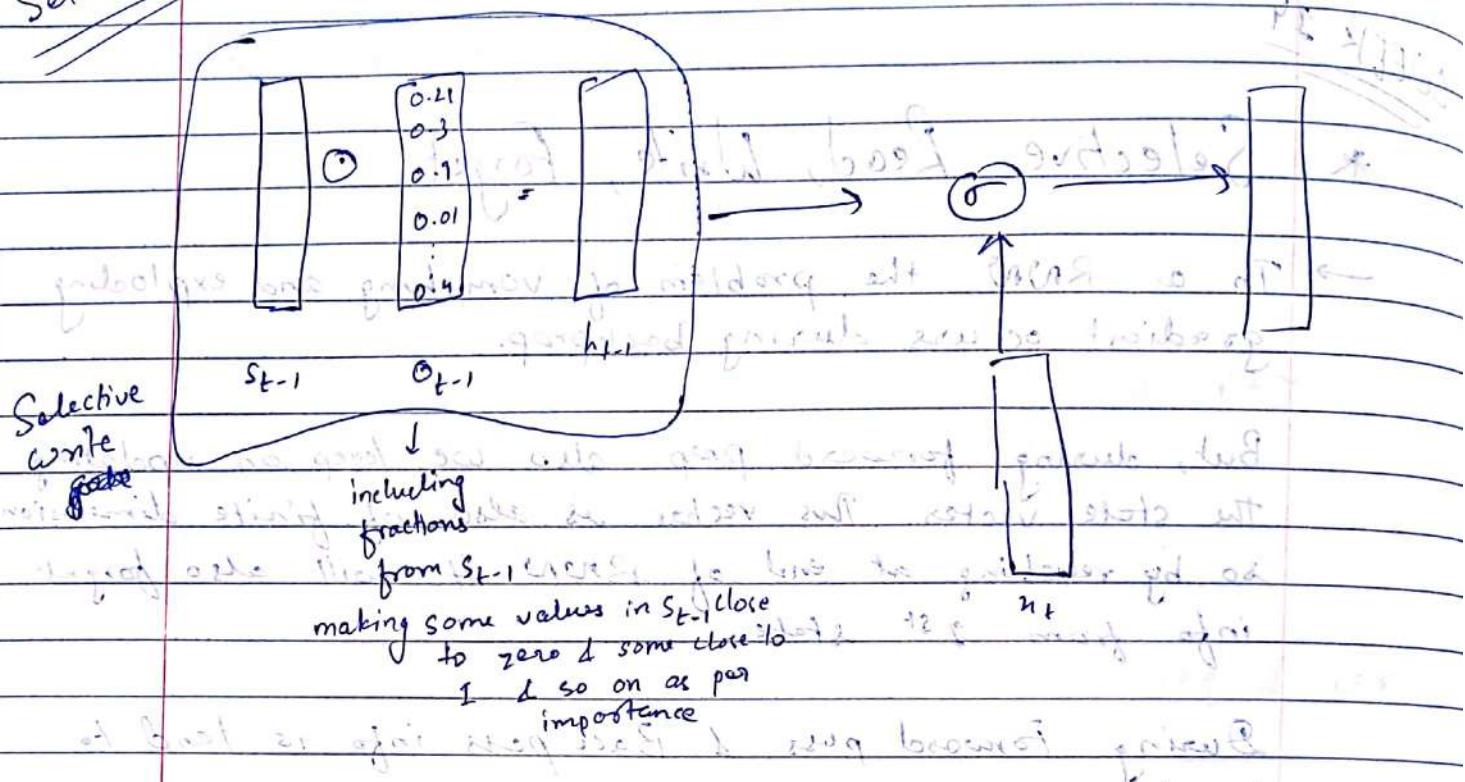
- Forget - the information like stop words
- Selective Read - the info added by previous sentiment bearing words
- Selective Write - new information from current word to state.



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

Ishan Modi

Selective Write



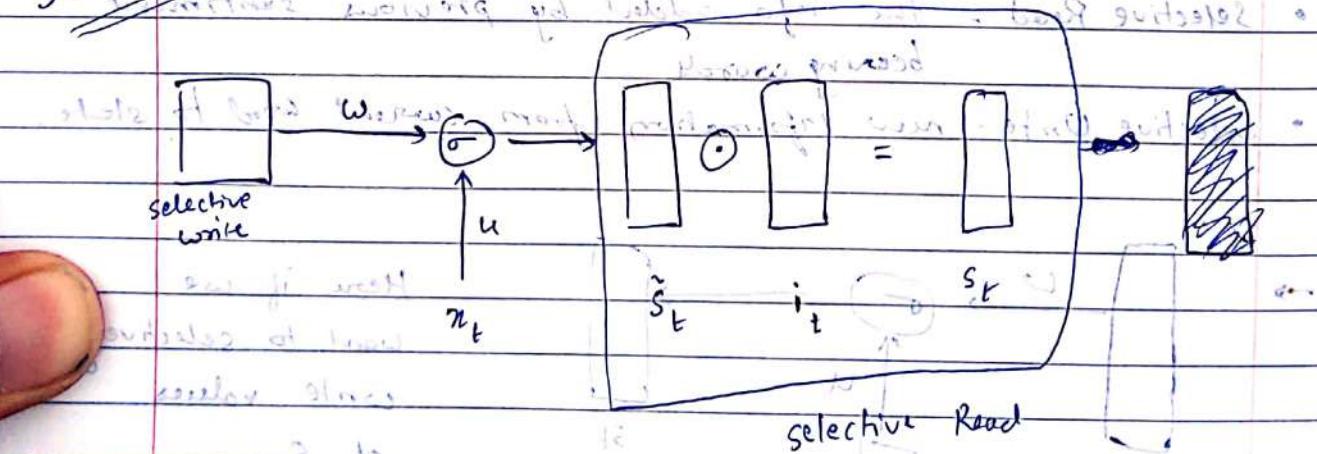
Now, we need to learn O_{t-1} (Output Gate)

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

→ (W, V, b) are parameters of RNN

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

Selective Read



Ishan Modi

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

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

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

Note

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

Till now

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

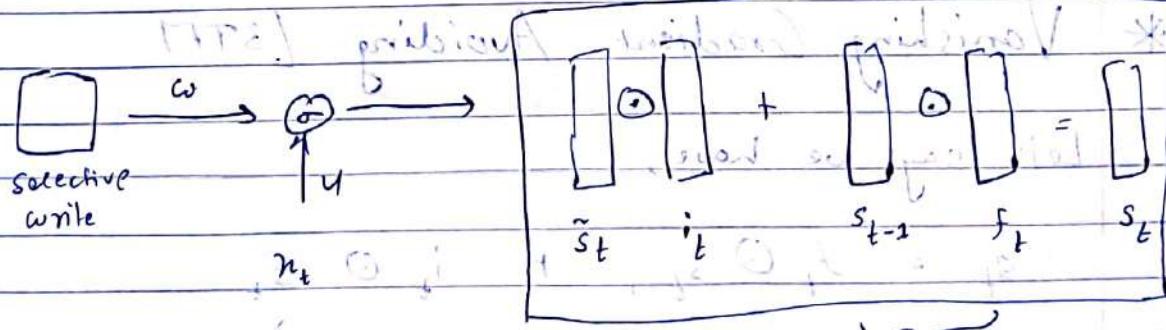
- Output gate $\rightarrow o_{t-1} = \sigma(w_o h_{t-1} + u_{ox_{t-1}} + b_o)$

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

- Current (temp) state $\rightarrow \tilde{s}_t = \sigma(w_{h_t} h_t + u_{x_t} + b)$

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

Input \rightarrow Selective Read $\rightarrow i_t \odot \tilde{s}_t$ (if true) else s_{t-1}

Selective Forget

Ishan Modi

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

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

* GRU (Gated Recurrent Unit)

→ Gates

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

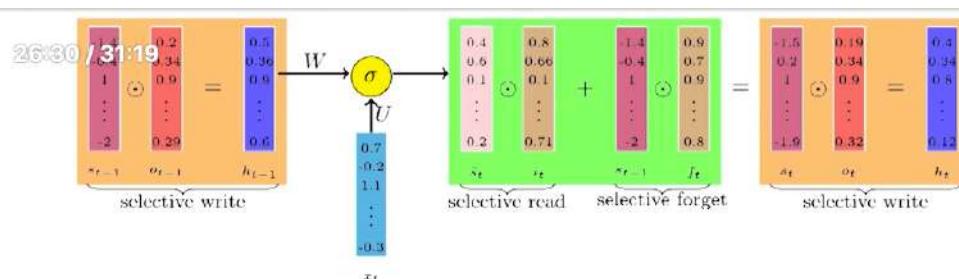
$$\bullet i_t = \sigma(W_i s_{t-1} + U_i v_t + b_i)$$

→ States

$$\bullet \tilde{s}_t = \sigma(W(o_t \odot s_{t-1}) + U v_t + b)$$

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

LSTM (Diagram - 1)



- We now have the full set of equations for LSTMs
- The green box together with the selective write operations following it, show all the computations which happen at timestep t

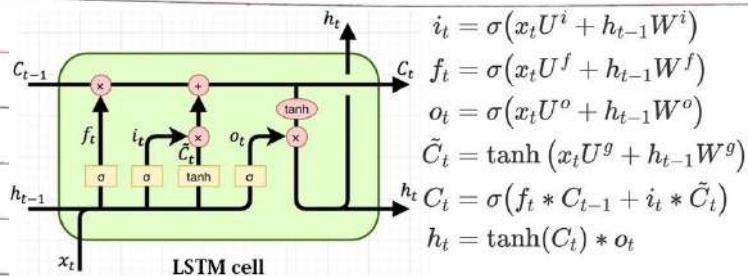
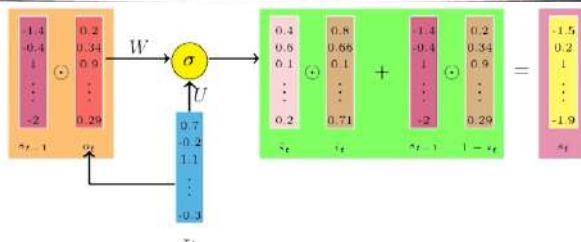
Gates:

$$\begin{aligned} o_t &= \sigma(W_o h_{t-1} + U_o x_t + b_o) \\ i_t &= \sigma(W_i h_{t-1} + U_i x_t + b_i) \\ f_t &= \sigma(W_f h_{t-1} + U_f x_t + b_f) \end{aligned}$$

States:

$$\begin{aligned} \tilde{s}_t &= \sigma(W h_{t-1} + U x_t + b) \\ s_t &= f_t \odot s_{t-1} + i_t \odot \tilde{s}_t \\ h_t &= o_t \odot \sigma(s_t) \text{ and } rnn_{out} = h_t \end{aligned}$$

Ishan Modi

LSTM (Diagram - 2)**GRU (Diagram - 1)**

The full set of equations for GRUs
Gates:

$$o_t = \sigma(W_o s_{t-1} + U_o x_t + b_o)$$

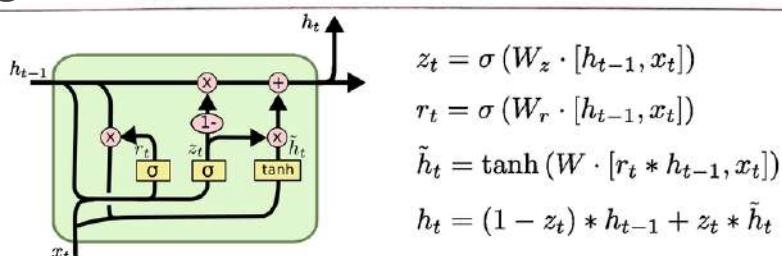
$$i_t = \sigma(W_i s_{t-1} + U_i x_t + b_i)$$

States:

$$\tilde{s}_t = \sigma(W(o_t \odot s_{t-1}) + U x_t + b)$$

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

- No explicit forget gate (the forget gate and input gates are tied)
- The gates depend directly on s_{t-1} and not the intermediate h_{t-1} as in the case of LSTMs

GRU (Diagram - 2)

* Vanishing Gradient Avoiding LSTM

Let's say we have,

$$S_t = f_t \odot S_{t-1} + i_t \odot \tilde{S}_t$$

In worst case we say that \downarrow tends to 0 or vanishes

Ishan Modi

Thus we have $(a)_{t-1}$ for backprop all along at each step since the loss does depend primarily on s_t and a_t .

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

$$\frac{\partial L}{\partial s_t} = \text{grad } f_t$$

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

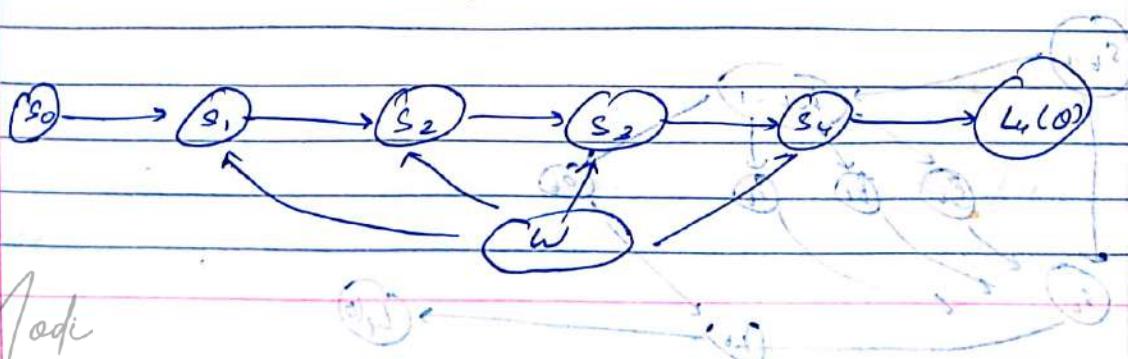
f_t information goes from 1 state to next

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

Conclusion

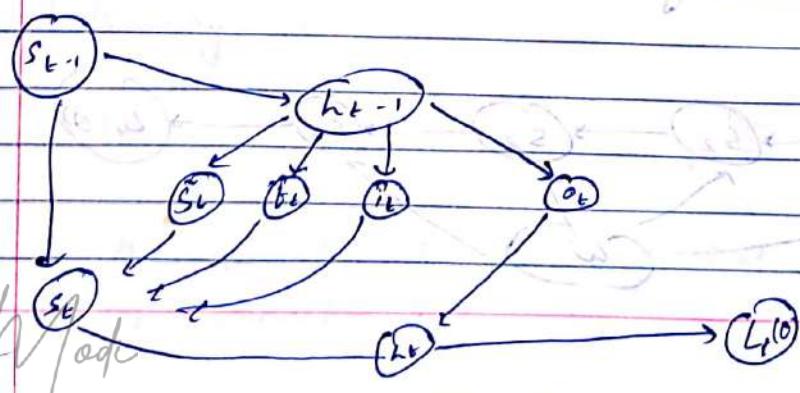
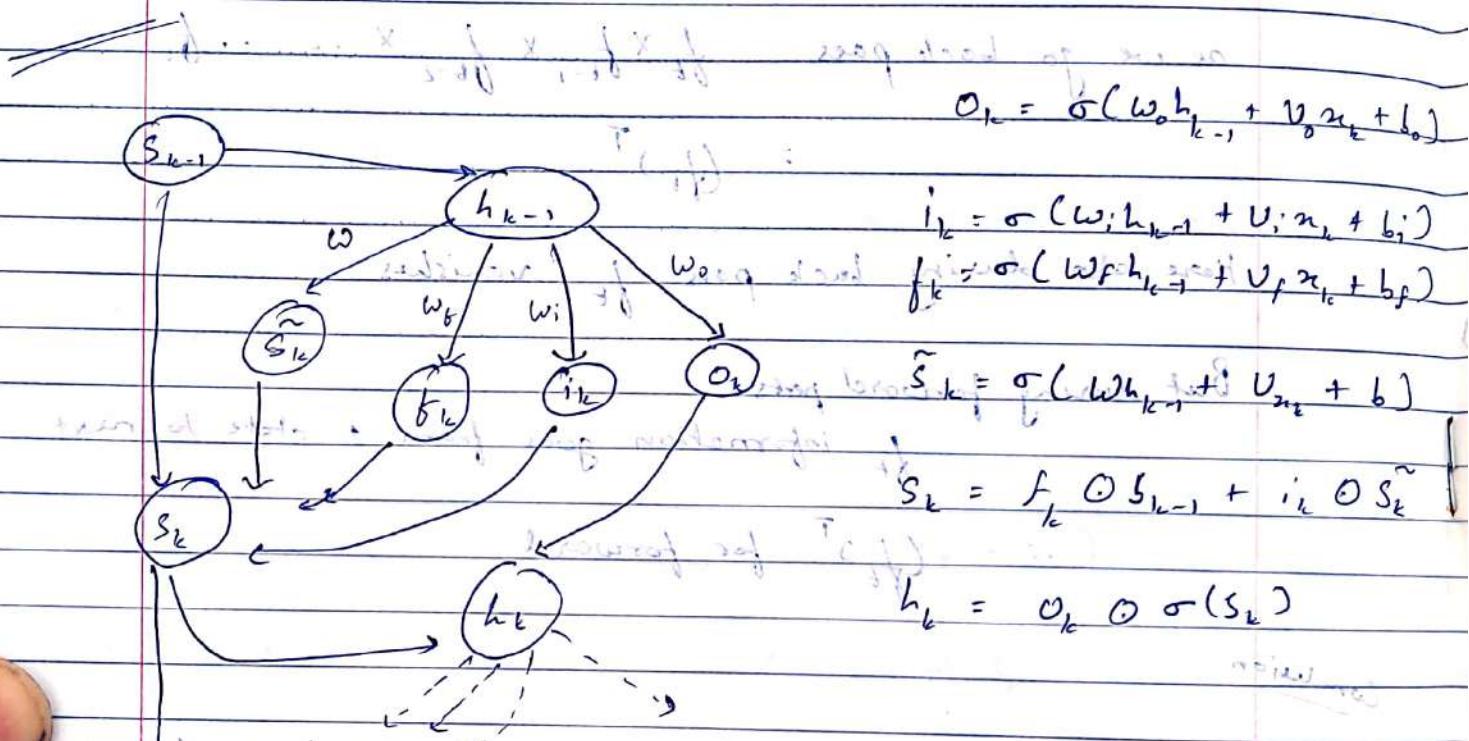
- If the state at time $t-1$ did not contribute to state at time t (ie, if $\|f_t\| \approx 0$ and $\|o_{t-1}\| \approx 0$) then during backprop the gradients flowing into s_{t-1} will vanish.

* Proof (Vanishing doesn't occur)



Ishan Modi

- 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.



→ Now let's say w_f causes error in output & we want to backprop gradient to (w_f) .

Then we need to prove that there exist at least one path from $L(\theta)$ to s_k that allows gradient to flow (so that it doesn't vanish) without any loss.

$$\frac{\partial L(\theta)}{\partial s_k} = \underbrace{\frac{\partial L(\theta)}{\partial h_t}}_{\text{constant}} \frac{\partial h_t}{\partial s_t} \left(\frac{\partial s_{t+1}}{\partial s_t} + \dots + \frac{\partial s_k}{\partial s_t} \right)$$

(constant of proportionality) \downarrow
bound (2)

\rightarrow ignoring worst case

$$\left(\begin{array}{cc} 0.6 & 0.1 \\ 0.1 & 0.6 \end{array} \right) \dots \left(\begin{array}{cc} 0.6 & 0.1 \\ 0.1 & 0.6 \end{array} \right) \left(\begin{array}{c} 0.16 \\ 0.16 \end{array} \right) = \left(\begin{array}{c} 0.16 \\ 0.16 \end{array} \right)$$

$$\textcircled{1} \quad \frac{\partial h_t}{\partial s_t} = \frac{\partial}{\partial s_t} \left(w_f \cdot \left(\sum_i s_i \right) + b \right) = \text{Diagonal} \left(w_f, 0 \dots \overset{(s_t)}{w_f}, 0 \dots \right) \text{ matrix}$$

($w_f \cdot (0 \dots 0 \dots w_f) \approx 0$) ignoring worst case

$$\textcircled{2} \quad \frac{\partial s_t}{\partial s_{t-1}} = \frac{\partial}{\partial s_{t-1}} \left(f_t \odot s_{t-1} + (i_t \odot s_t) \right) \rightarrow s_{t-1}$$

$$\text{iff } \left(\| (f_t \odot s_{t-1}) \| \| + \| (i_t \odot s_t) \| \right) \leq \| s_t \|$$

$= D(f_t)$

check your notes well for more info

$$\left(\frac{\partial s_t}{\partial s_{t-1}} \text{ mod. } \frac{\partial s_{t+1}}{\partial s_t} \text{ mod. } \dots \right) D(f_t \odot s_{t-1}) + (i_t \odot s_t)$$

Ishan Modi

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

Need to take limit of $D(\sigma_{b=k+1}^t f_i)$ and

if gradient vanishes in backprop only if it vanished in forward prop. If vanish only when it needs to then prove

* Proof (Exploding can be stopped)

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

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

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

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

$$\|z\| \leq \|L_t'(h_t)\| (\|w_o\|)^{t-k+1}$$

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

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

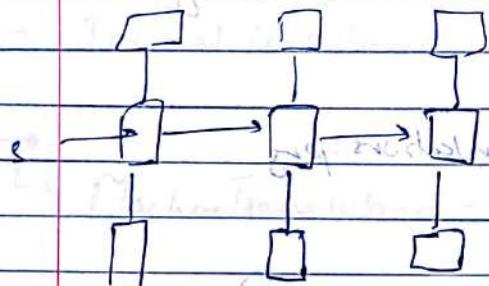
Ishan Modi

~~WEEK 15~~ What will the next token be predicted given the sequence of tokens seen so far?

* Introduction to Encoder-Decoder networks

→ 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_1, y_2, \dots, y_{t-1} are inputs to softmax function \rightarrow jth element of vector

$j \in V$, V is vocabulary

(vector of size equal to vocab).

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

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

∴ removed y_1^{t-1}

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

nothing changes

Ishan Mod

We will now use the following representations of RNNs, GRUs and LSTMs

$$\begin{array}{lll} s_t = \sigma(Ux_t + Ws_{t-1} + b) & \tilde{s}_t = \sigma(W(o_t \odot s_{t-1}) + Ux_t + b) & \tilde{s}_t = \sigma(Wh_{t-1} + Ux_t + b) \\ & s_t = i_t \odot s_{t-1} + (1 - i_t) \odot \tilde{s}_t & s_t = f_t \odot s_{t-1} + i_t \odot \tilde{s}_t \\ & & h_t = o_t \odot \sigma(s_t) \end{array}$$

$$s_t = \text{RNN}(s_{t-1}, x_t) \quad s_t = \text{GRU}(s_{t-1}, x_t) \quad h_t, s_t = \text{LSTM}(h_{t-1}, s_{t-1}, x_t)$$

Now {

rest of the slides are taken
and this day was also taken }

→ Image Captioning (Encoder - Decoder)

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

$i(s + \epsilon)$ variable : (i_p, i_f, i_g)

Here we have image as well

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

cannot take entire image
we need representation of image

$$P(y_t | s_t, f_{c_7}(I)) \rightarrow \text{Fully connected layer in VGG.}$$

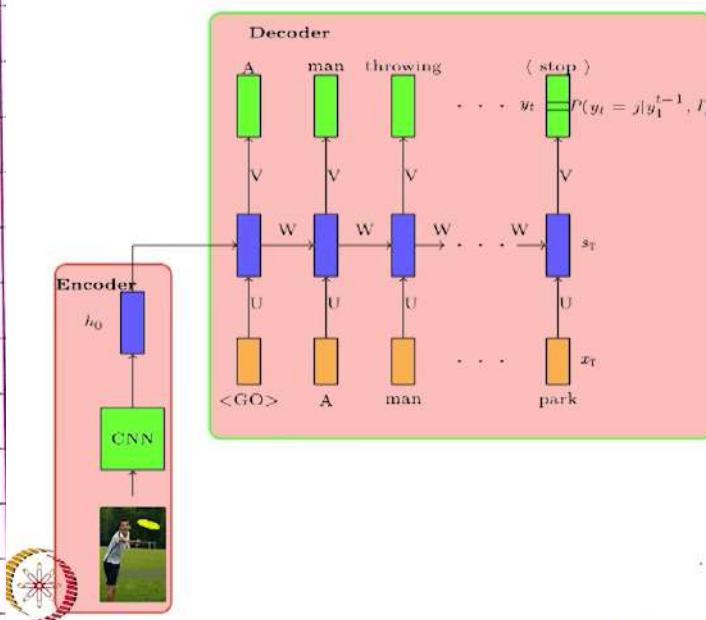
This becomes s_0 ,

$$\therefore s_0 = f_{c_7}(I) \text{ now we feed it to RNN}$$

for caption generation

Ishan Modi

We can Also feed FC7(I) to all states,

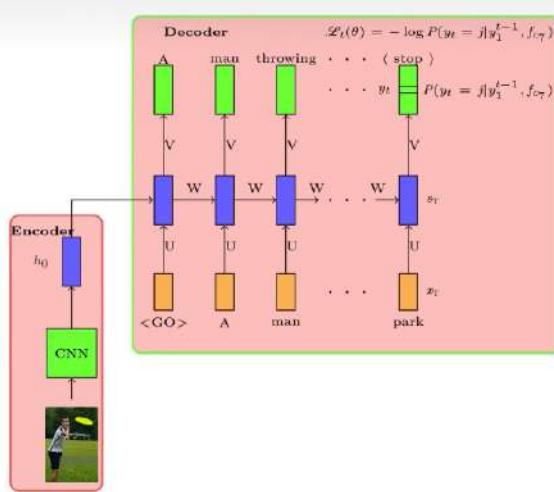


- Let us look at the full architecture
- A CNN is first used to **encode** the image
- A RNN is then used to decode (generate) a sentence from the encoding
- This is a typical **encoder-decoder architecture**



* Application of Encoder-Decoder

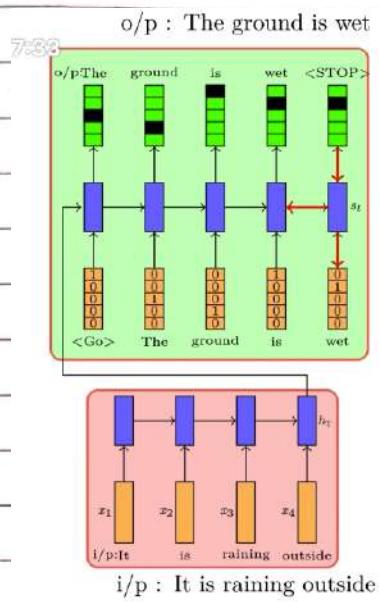
---> Image Captioning



- Task:** Image captioning
- Data:** $\{x_i = \text{image}_i, y_i = \text{caption}_i\}_{i=1}^N$
- Model:**
 - Encoder:** $s_0 = \text{CNN}(x_i)$
 - Decoder:** $s_t = \text{RNN}(s_{t-1}, e(\hat{y}_{t-1}))$
 $P(y_t | y_1^{t-1}, I) = \text{softmax}(Vs_t + b)$
- Parameters:** $U_{dec}, V, W_{dec}, W_{conv}, b$
- Loss:** $\mathcal{L}(\theta) = \sum_{i=1}^T \mathcal{L}_i(\theta) = - \sum_{t=1}^T \log P(y_t = \ell_t | y_1^{t-1}, I)$
- Algorithm:** Gradient descent with back-propagation

Ishan Modi

---> Text Entailment



- **Task:** Textual entailment

- **Data:** $\{x_i = \text{premise}_i, y_i = \text{hypothesis}_i\}_{i=1}^N$

- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

$$s_t = RNN(s_{t-1}, e(\hat{y}_{t-1}))$$

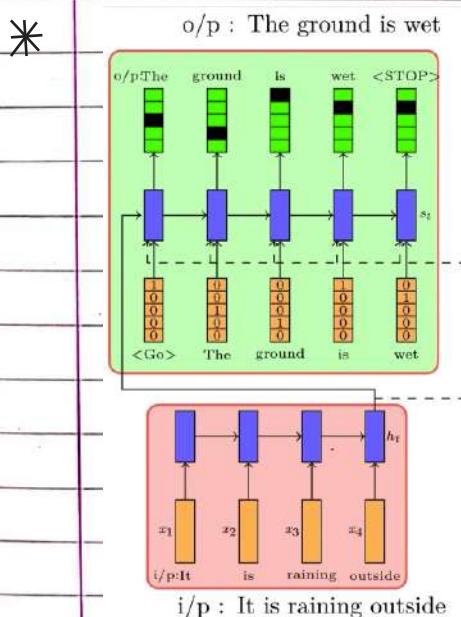
$$P(y_t|y_1^{t-1}, x) = \text{softmax}(Vs_t + b)$$

- **Parameters:** $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$

- **Loss:**

$$\mathcal{L}(\theta) = \sum_{i=1}^T \mathcal{L}_i(\theta) = - \sum_{t=1}^T \log P(y_t = \ell_t | y_1^{t-1}, x)$$

- **Algorithm:** Gradient descent with back-propagation



- **Task:** Textual entailment

- **Data:** $\{x_i = \text{premise}_i, y_i = \text{hypothesis}_i\}_{i=1}^N$

- **Model (Option 2):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

$$s_t = RNN(s_{t-1}, [h_T, e(\hat{y}_{t-1})])$$

$$P(y_t|y_1^{t-1}, x) = \text{softmax}(Vs_t + b)$$

- **Parameters:** $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$

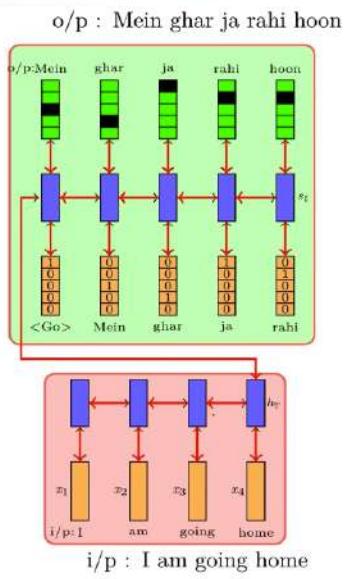
- **Loss:**

$$\mathcal{L}(\theta) = \sum_{i=1}^T \mathcal{L}_i(\theta) = - \sum_{t=1}^T \log P(y_t = \ell_t | y_1^{t-1}, x)$$

- **Algorithm:** Gradient descent with back-propagation

Ishan Modi

---> Machine Translation



- **Task:** Machine translation

- **Data:** $\{x_i = \text{source}_i, y_i = \text{target}_i\}_{i=1}^N$

- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

$$s_t = RNN(s_{t-1}, e(\hat{y}_{t-1}))$$

$$P(y_t | y_1^{t-1}, x) = \text{softmax}(Vs_t + b)$$

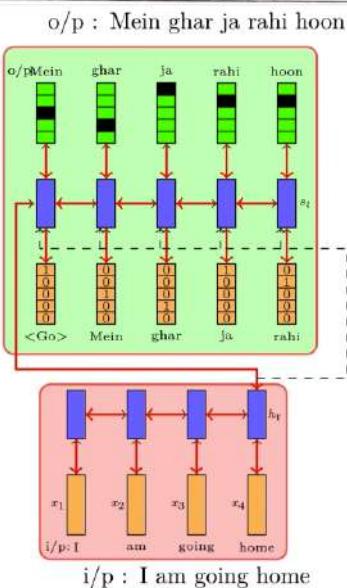
- **Parameters:** $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$

- **Loss:**

$$\mathcal{L}(\theta) = \sum_{i=1}^T \mathcal{L}_i(\theta) = - \sum_{i=1}^T \log P(y_t = \ell_t | y_1^{t-1}, x)$$

- **Algorithm:** Gradient descent with back-propagation

★



- **Task:** Machine translation

- **Data:** $\{x_i = \text{source}_i, y_i = \text{target}_i\}_{i=1}^N$

- **Model (Option 2):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

$$s_t = RNN(s_{t-1}, [h_T, e(\hat{y}_{t-1})])$$

$$P(y_t | y_1^{t-1}, x) = \text{softmax}(Vs_t + b)$$

- **Parameters:** $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$

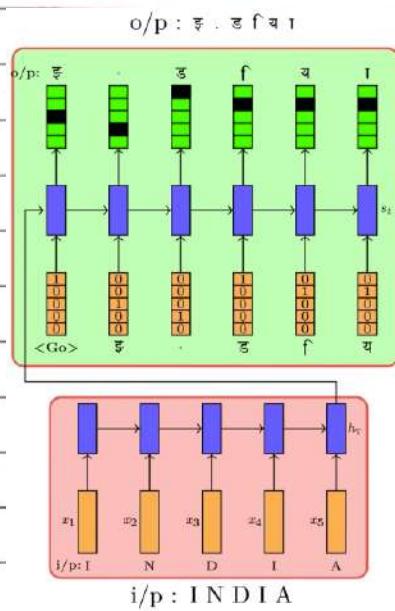
- **Loss:**

$$\mathcal{L}(\theta) = \sum_{i=1}^T \mathcal{L}_i(\theta) = - \sum_{i=1}^T \log P(y_t = \ell_t | y_1^{t-1}, x)$$

- **Algorithm:** Gradient descent with back-propagation

Ishan Modi

---> Transliteration



- Task: Transliteration

- Data: $\{x_i = \text{srcword}_i, y_i = \text{tgtword}_i\}_{i=1}^N$

- Model (Option 1):

 - Encoder:

$$h_t = RNN(h_{t-1}, x_{it})$$

 - Decoder:

$$s_0 = h_T \quad (T \text{ is length of input})$$

$$s_t = RNN(s_{t-1}, e(\hat{y}_{t-1}))$$

$$P(y_t|y_1^{t-1}, x) = \text{softmax}(Vs_t + b)$$

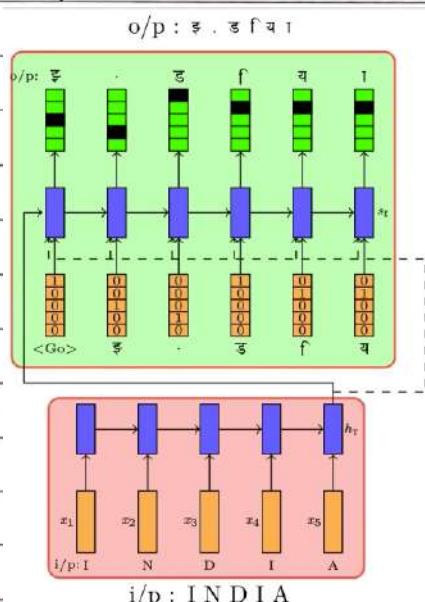
- Parameters: $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$

- Loss:

$$\mathcal{L}(\theta) = \sum_{i=1}^T \mathcal{L}_t(\theta) = - \sum_{t=1}^T \log P(y_t = \ell_t | y_1^{t-1}, x)$$

- Algorithm: Gradient descent with back-propagation

*



- Task: Transliteration

- Data: $\{x_i = \text{srcword}_i, y_i = \text{tgtword}_i\}_{i=1}^N$

- Model (Option 2):

 - Encoder:

$$h_t = RNN(h_{t-1}, x_{it})$$

 - Decoder:

$$s_0 = h_T \quad (T \text{ is length of input})$$

$$s_t = RNN(s_{t-1}, [e(\hat{y}_{t-1}), h_T])$$

$$P(y_t|y_1^{t-1}, x) = \text{softmax}(Vs_t + b)$$

- Parameters: $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$

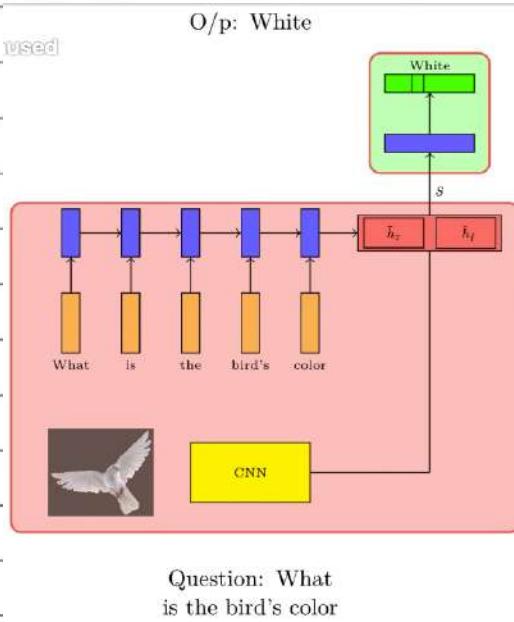
- Loss:

$$\mathcal{L}(\theta) = \sum_{i=1}^T \mathcal{L}_t(\theta) = - \sum_{t=1}^T \log P(y_t = \ell_t | y_1^{t-1}, x)$$

- Algorithm: Gradient descent with back-propagation

Ishan Modi

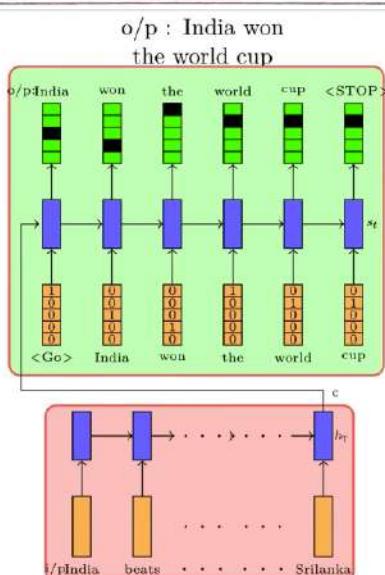
---> Image Q&A



- Task: Image Question Answering
- Data: $\{x_i = \{I, q\}_i, y_i = \text{Answer}_i\}_{i=1}^N$
- Model:
 - Encoder: $\hat{h}_I = \text{CNN}(I), \hat{h}_t = \text{RNN}(\hat{h}_{t-1}, q_i)$
 $s = [\hat{h}_T; \hat{h}_I]$
 - Decoder: $P(y|q, I) = \text{softmax}(Vs + b)$
- Parameters: $V, b, U_q, W_q, W_{conv}, b$

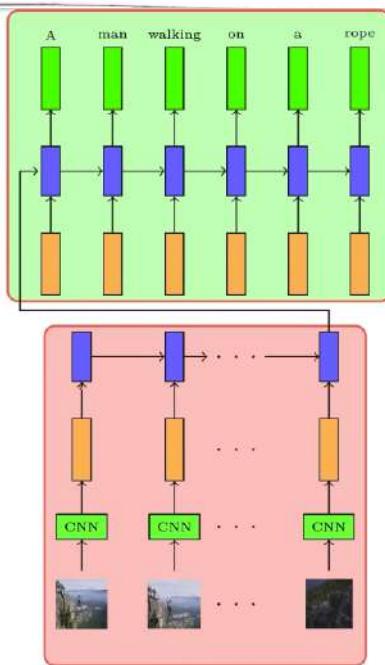


* Document Summarization



- Task: Document Summarization
- Data: $\{x_i = \text{Document}_i, y_i = \text{Summary}_i\}_{i=1}^N$
- Model:
 - Encoder: $h_t = \text{RNN}(h_{t-1}, x_{it})$
 - Decoder:
 - $s_0 = h_T$
 - $s_t = \text{RNN}(s_{t-1}, e(\hat{y}_{t-1}))$
 - $P(y_t|y_1^{t-1}, x) = \text{softmax}(Vs_t + b)$
- Parameters: $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$
- Loss: $\mathcal{L}(\theta) = - \sum_{t=1}^T \log P(y_t = \ell_t | y_1^{t-1}, x)$
- Algorithm: Gradient descent with backpropagation

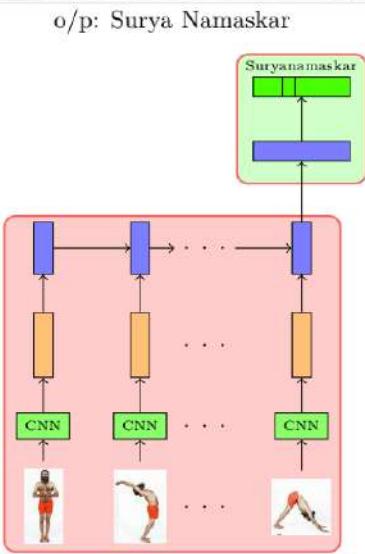
---> Video Captioning



- **Task:** Video Captioning
- **Data:** $\{x_i = \text{video}_i, y_i = \text{desc}_i\}_{i=1}^N$
- **Model:**
 - **Encoder:** $h_t = \text{RNN}(h_{t-1}, \text{CNN}(x_{it}))$
 - **Decoder:** $s_0 = h_T$
 $s_t = \text{RNN}(s_{t-1}, e(\hat{y}_{t-1}))$
 $P(y_t|y_1^{t-1}, x) = \text{softmax}(Vs + b)$
- **Parameters:** $U_{dec}, W_{dec}, V, b, W_{conv}, U_{enc}, W_{enc}, b$
- **Loss:** $\mathcal{L}(\theta) = -\sum_{i=1}^T \log P(y_t = \ell_t | y_1^{t-1}, x)$
- **Algorithm:** Gradient descent with backpropagation

*

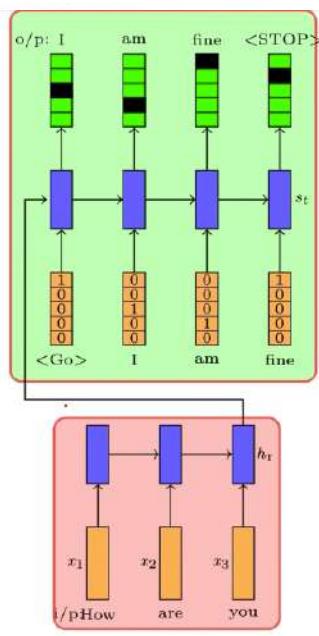
---> Video Classification



- **Task:** Video Classification
- **Data:** $\{x_i = \text{Video}_i, y_i = \text{Activity}_i\}_{i=1}^N$
- **Model:**
 - **Encoder:** $h_t = \text{RNN}(h_{t-1}, \text{CNN}(x_{it}))$
 - **Decoder:** $s = h_T$
 $P(y|I) = \text{softmax}(Vs + b)$
- **Parameters:** $V, b, W_{conv}, U_{enc}, W_{enc}, b$
- **Loss:** $\mathcal{L}(\theta) = -\log P(y = \ell | \text{Video})$
- **Algorithm:** Gradient descent with backpropagation

Ishan Modi

---> Dialog



i/p: How are you

- Task: Dialog
- Data: $\{x_i = \text{Utterance}_i, y_i = \text{Response}_i\}_{i=1}^N$

- Model:

- Encoder:

$$h_t = RNN(h_{t-1}, x_{it})$$

- Decoder:

$$s_0 = h_T \quad (T \text{ is length of input})$$

$$s_t = RNN(s_{t-1}, e(\hat{y}_{t-1}))$$

$$P(y_t|y_1^{t-1}, x) = \text{softmax}(Vs_t + b)$$

- Parameters: $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$

- Loss:

$$\mathcal{L}(\theta) = \sum_{i=1}^T \mathcal{L}_i(\theta) = - \sum_{t=1}^T \log P(y_t|x)$$

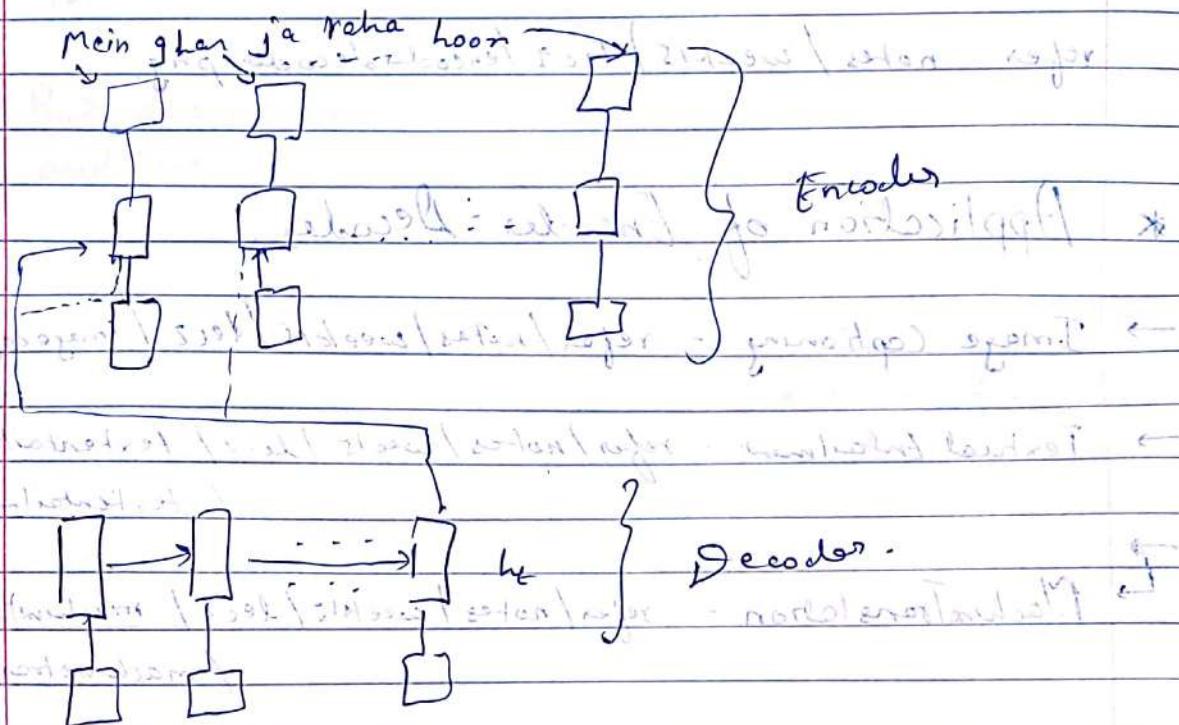
- Algorithm: Gradient descent w/gation



*

Ishan Modi

* Attention Mechanism

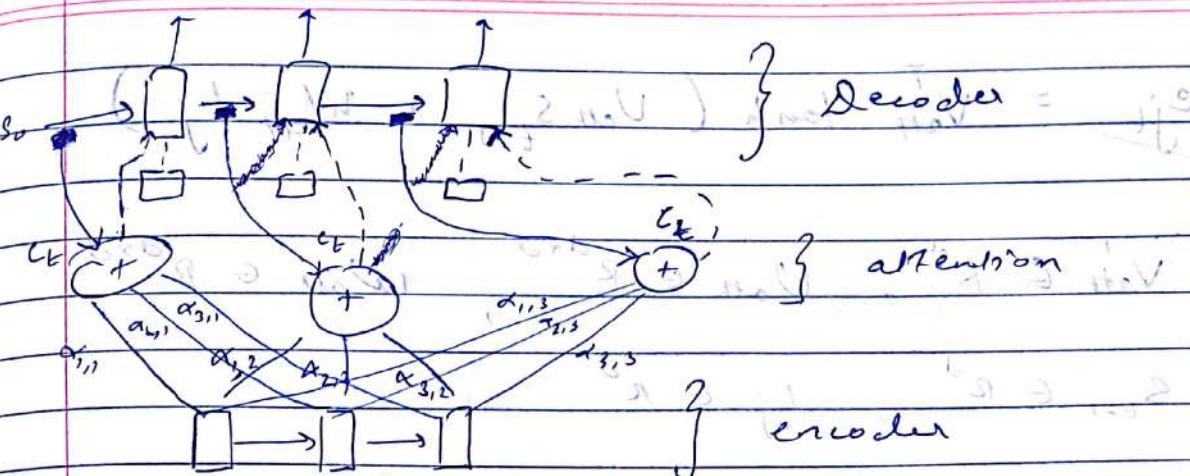


So here we are sort of bypassing all bypasses - So if you pass the network, this h_2 contains all the information ie (I am going home) but we can observe that you have used I as well as home

Travelling → Mein ghar ja raha hoon - intended after - priority
 home → ghar These are words that are highly similar & closely related
 am → raha intended after
 going → ja raha intended after

So the intuition is why not pass this relationship instead of h_2

Ishan Modi



small weight matrix α is to generate extra point $\alpha_{j,t}$.

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~~ ^{decoder state} | ~~encoder state~~ ^{encoder state}

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

argmax over j in M

j^{th} word in encoder t^{th} time stamp in decoder of argmax of $e_{j,t}$ $\alpha_{j,t}$ $\exp(e_{j,t})$

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

marking for $j=1$ to M

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

Thus $\alpha_{j,t}$ is measure of how important is j^{th} word with respect to t^{th} time stamp.

IshanModi

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

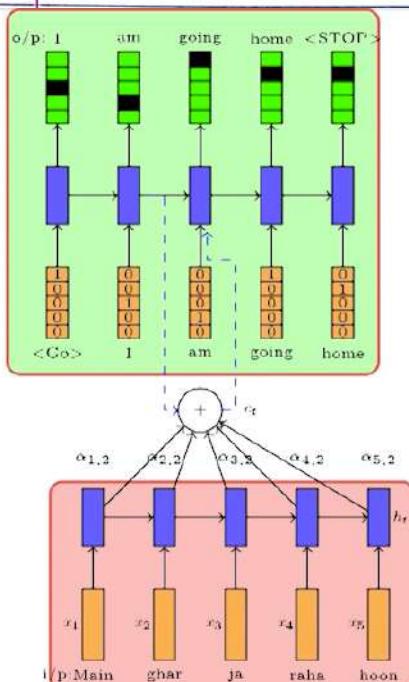
scales

$$V_{attn} \in \mathbb{R}^{d \times d}, \quad U_{attn} \in \mathbb{R}^{d \times d}, \quad W_{attn} \in \mathbb{R}^{d \times d}$$

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

This thing works because it is better modeling choice

Also, c_t is $\sum_{j=1}^T \alpha_{jt} h_j$ for encoder.



- Task: Machine Translation

- Data: $\{x_i = \text{source}_i, y_i = \text{target}_i\}_{i=1}^N$

- Encoder:

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

$$s_0 = h_T$$

- Decoder:

$$e_{jt} = V_{attn}^T \tanh(U_{attn}h_j + W_{attn}s_t)$$

$$\alpha_{jt} = \text{softmax}(e_{jt})$$

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

$$s_t = RNN(s_{t-1}, [e(\hat{y}_{t-1}), c_t])$$

$$\ell_t = \text{softmax}(Vs_t + b)$$

- Parameters: $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b, U_{attn}, V_{attn}$

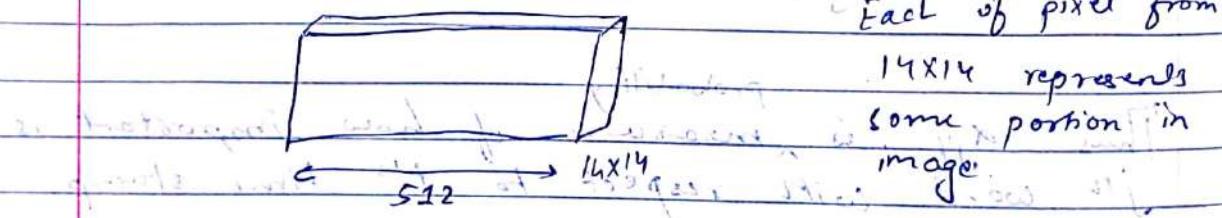
- Loss and Algorithm remains same

Ishan Modi

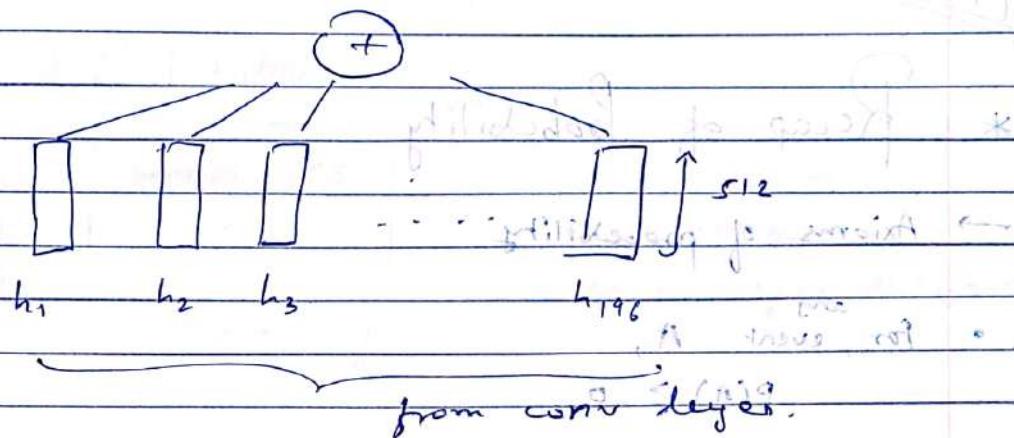
* Attention Over Images

Unlike RNN, where each state gives word CNN for image gives h_{t-1} which doesn't give us which entry is for which pixel in image

Thus we use the convolution layers before FC,



$$14 \times 14 = 196$$



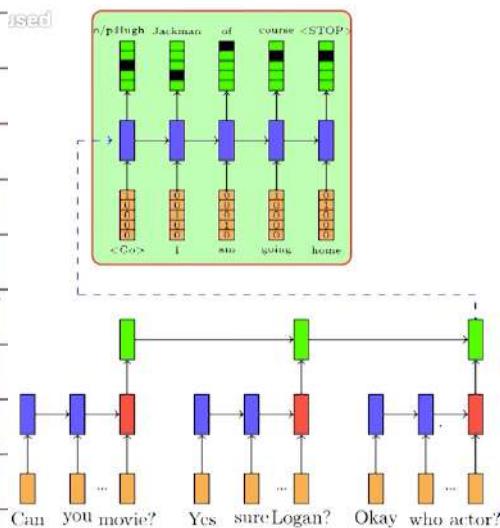
$$\alpha_{tj} = f_{ATT}(s_{t-1}, h_j)$$

$$f_{ATT} = v^T \tanh(u^T s_{t-1} + w h_j + b)$$

* Hierarchical Attention

---> Understanding hierarchical RNN structure

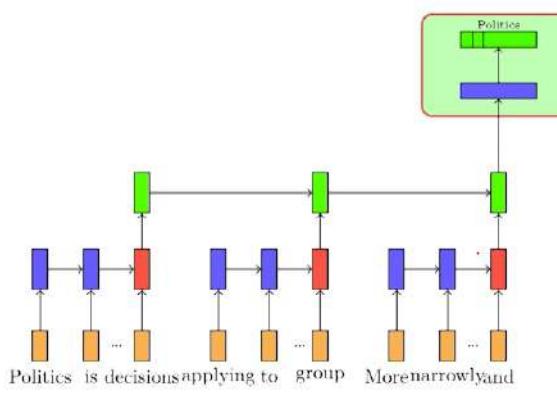
Dialog hierarchical RNN



- We could think of a two level hierarchical RNN encoder
- The first level RNN operates on the sequence of words in each utterance and gives us a representation
- We now have a sequence of utterance representations (red vectors in the image)
- We can now have another RNN which encodes this sequence and gives a single representation for the sequences of utterances
- The decoder can then produce an output sequence conditioned on this utterance

Document hierarchical RNN

Politics is the process of making decisions applying to all members of each group.
More narrowly, it refers to achieving and ...



- Data:** $\{Document_i, class_i\}_{i=1}^N$
- Word level (1) encoder:**

$$h_{ij}^1 = RNN(h_{ij-1}^1, w_{ij})$$

$$s_i = h_{iT_i}^1 \quad [T \text{ is length of sentence } i]$$
- Sentence level (2) encoder:**

$$h_i^2 = RNN(h_{i-1}^2, s_i)$$

$$s = h_K^2 \quad [K \text{ is number of sentences}]$$
- Decoder:**

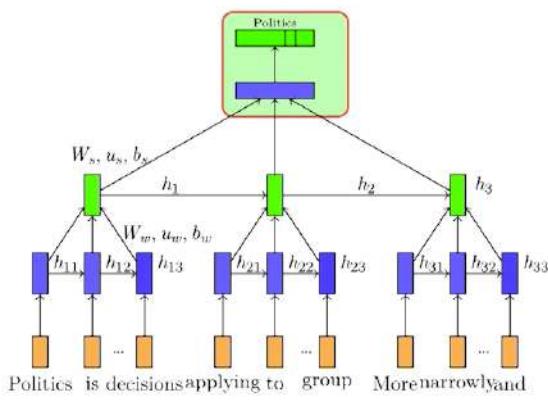
$$P(y|document) = softmax(Vs + b)$$
- Params:** $W_{enc}^1, U_{enc}^1, W_{enc}^2, U_{enc}^2, V, b$
- Loss:** Cross Entropy
- Algorithm:** Gradient Descent with backpropagation

Ishan Modi

Now we use attention mechanism in these hierarchical module

Encoder (heirarchical attention)

Politics is the process of making decisions applying to all members of each group.
More narrowly, it refers to achieving and ...



- Data: $\{Document_i, class_i\}_{i=1}^N$

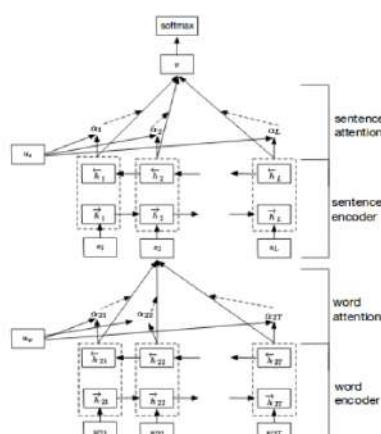
- Word level (1) encoder:

$$\begin{aligned} h_{ij} &= RNN(h_{ij-1}, w_{ij}) \\ u_{ij} &= \tanh(W_u h_{ij} + b_u) \\ \alpha_{ij} &= \frac{\exp(u_{ij}^T u_w)}{\sum_t \exp(u_{it}^T u_w)} \\ s_t &= \sum_j \alpha_{ij} h_{ij} \end{aligned}$$

- Sentence level (2) encoder:

$$\begin{aligned} h_i &= RNN(h_{i-1}, s_i) \\ u_i &= \tanh(W_s h_i + b_s) \\ \alpha_i &= \frac{\exp(u_i^T u_s)}{\sum_i \exp(u_i^T u_s)} \\ s &= \sum_i \alpha_i h_i \end{aligned}$$

Decoder (heirarchical attention)



- Decoder:

$$P(y|document) = \text{softmax}(Vs + b)$$

- Parameters:

$$W_w, W_s, V, b_w, b_s, b, u_w, u_s$$

- Loss: cross entropy

- Algorithm: Gradient Descent and backpropagation

Figure: Hierarchical Attention Network
[Yang et al.]

Ishan Modi