

$$Au - \lambda u = 0$$

$$|A - \lambda I| = 0$$

PCA : Disadvantage :

↓ Only used when linear transformation
linear

1st Oct

DEEP LEARNING

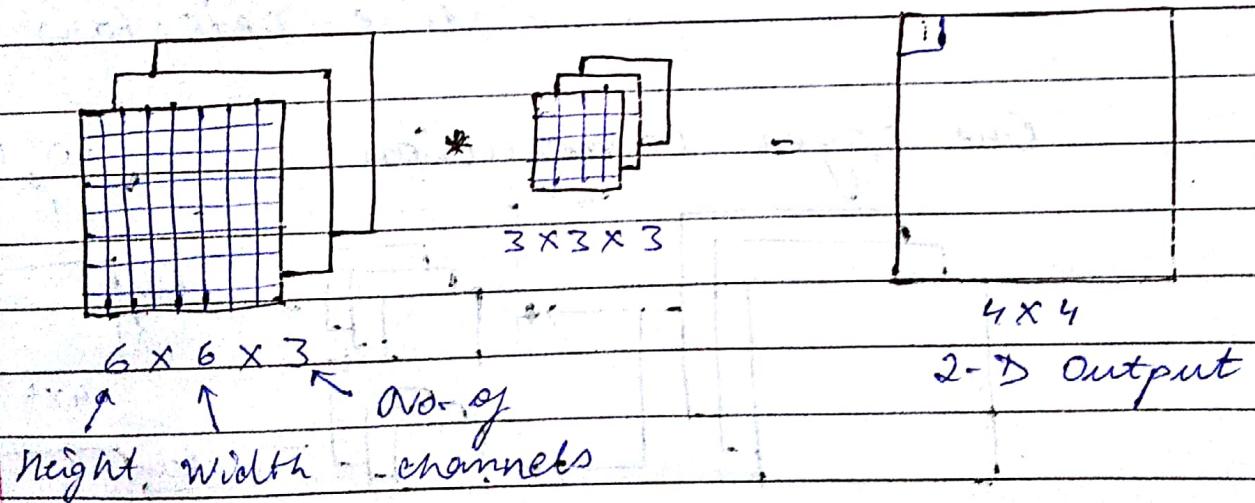
Boltzmann Machine

Auto-encoder

CNN

22nd Oct

Convolution Neural network

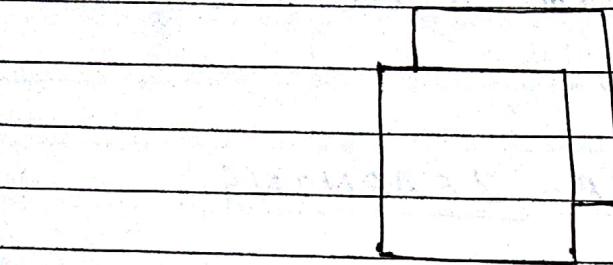


→ 27 elements are multiplied by 27 elements of filter and added to generate 1 element of output.

Vertical Edge:

	R	G	B
Get info from all channels	1 0 -1	1 0 -1	1 0 -1
	1 0 -1	1 0 -1	1 0 -1
	1 0 -1	1 0 -1	1 0 -1

If we ~~are~~ select more no. of filters (vertical & horizontal), those will be stacked of two 2D image $4 \times 4 \times 2$



For 2 filters

$$n \times n \times n_c * f \times f \times n_c$$

$$\Rightarrow (n-f+1) \times (n-f+1) \times n_c'$$

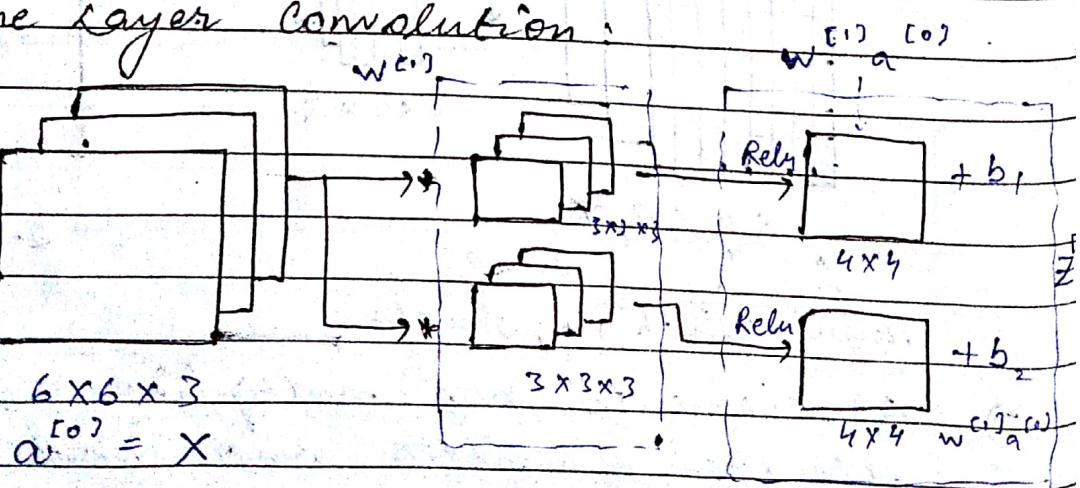
Padding = 0

No. of filters α

stride = 1

Feature Maps / Feature Detectors

One layer convolution:



$$z^{[1]} = w^{[1]} a^{[0]} + b,$$

$$a^{[1]} = g(z^{[1]})$$

$a^{[0]}$ to $a^{[1]}$

$6 \times 6 \times 3 \rightarrow 4 \times 4 \times 2$ after convolution

Down Sampling

How many parameters?

$$4 \times 4 \times 10 \Rightarrow (27+1) \times 10 \\ = 280$$

\Rightarrow Independent of Input image size.
less prone to overfitting

if l no. of layers are in Convolution layer,

$f^{[l]}$ = filter size of l^{th} layer

$P^{[l]}$ = padding " "

$s^{[l]}$ = stride " "

$$\text{Input} = n_h^{[L-1]} \times n_w^{[L-1]} \times n_c^{[L-1]}$$

$$\text{Output} : n_h^{[L]} \times n_w^{[L]} \times n_c^{[L]}$$

$$n_h^{[L]} = \left\lfloor \frac{n_h^{[L-1]} + 2P^{[L]} - f^{[L]} + 1}{s} \right\rfloor$$

$$n_w^{[L]} = \left\lfloor \frac{n_w^{[L-1]} + 2P^{[L]} - f^{[L]} + 1}{s} \right\rfloor$$

$n_c^{[L]}$ = no. of filters

Each filter is of size $f \times f \times n_c^{[L-1]}$

$$\text{Activations} : a^{[L]} \rightarrow n_h^{[L]} \times n_w^{[L]} \times n_c^{[L]}$$

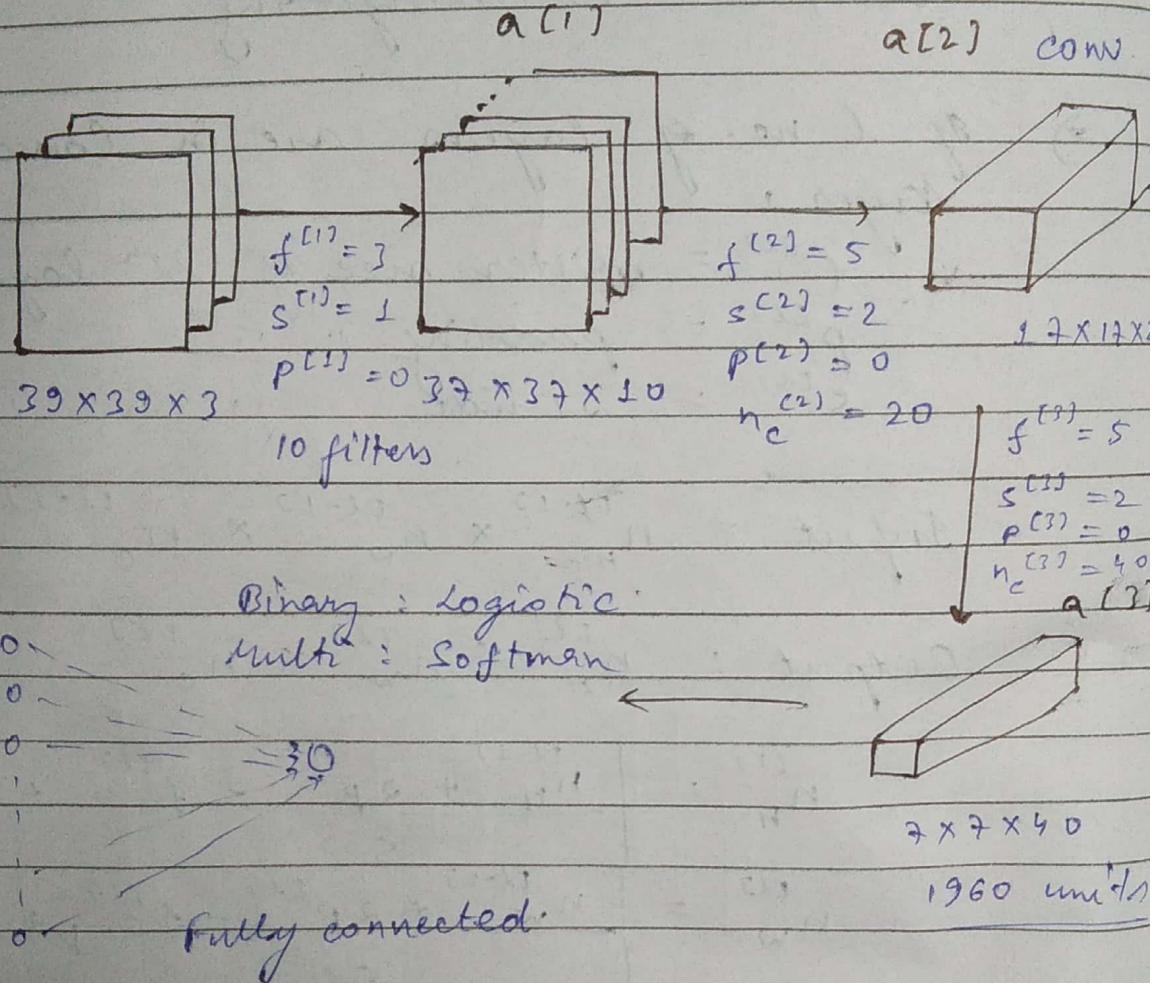
$$\underbrace{A^{[L]}}_I \rightarrow M \times n_h^{[L]} \times n_w^{[L]} \times n_c^{[L]}$$

Vector notation

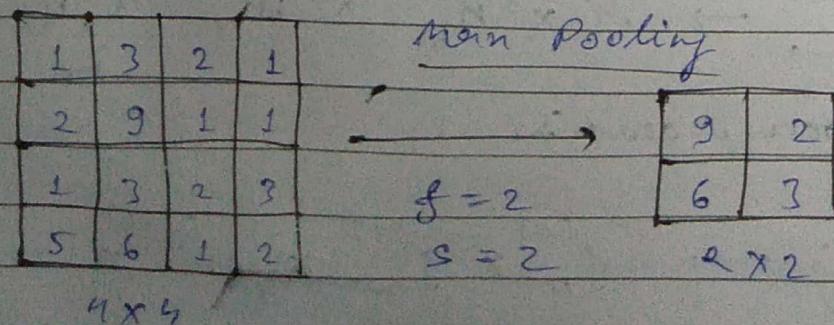
$$\text{Weight} = f_h^{[L]} \times f_w^{[L]} \times n_c^{[L-1]} \times n_c^{[L]}$$

Bias : a vector
 $[1, 1, \dots, n_c^{(1)}]$

Deeper Convolution NN :



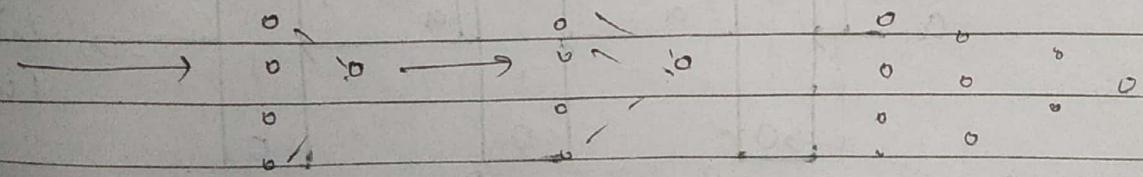
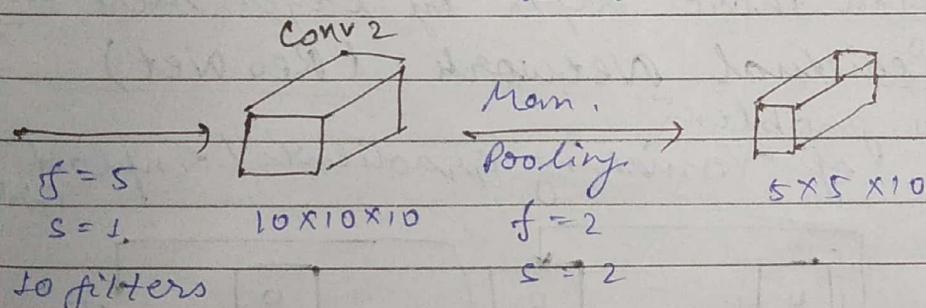
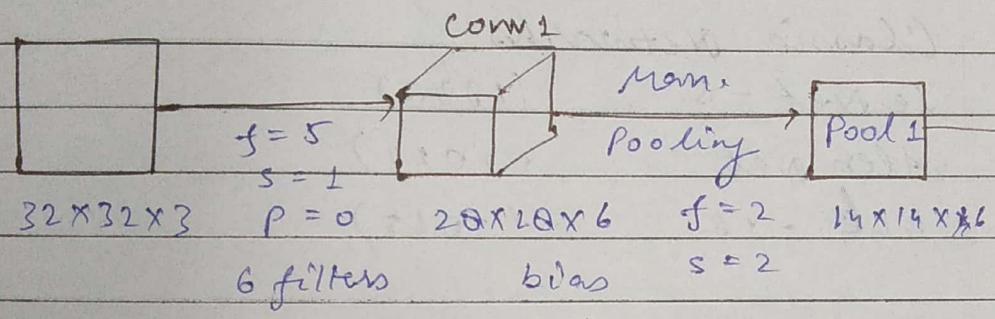
- ① Conv Layer
- ② Pooling layer → Reduce dimension
Robust feature detector
- ③ Fully connected layer



Pooling \Rightarrow Only reduction, no parameter learning.

\rightarrow Average Pooling is also there.
 \hookrightarrow Used for very deep conv.

Fully connected layer:



n_H, n_W decreasing
 n_C increasing with layer

	Activation	Activation size	# Parameters
Input	$(32, 32, 3)$	3072	0
Conv 1 $(f=5, s=1)$	$(20, 20, 6)$		
Pool 1	$(14, 14, 6)$		
FC			

Adv. of conv. net over MLP :

- 1) Sharing of weights / parameters
- 2) Sparsity of connection \therefore Less no. of computations ; depends on small no. of parameters.

Classic networks

LeNet - 5 (1998)

AlexNet (2012)

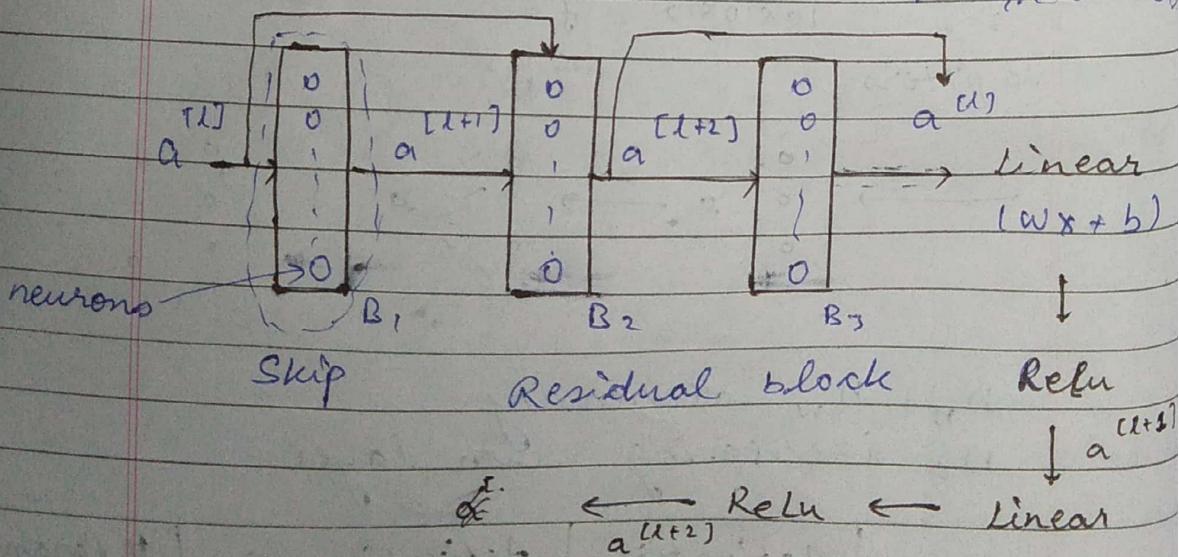
VGG - 16 (2015)

With more depth of layer, more 26th Oct

Residual Network (ResNet)

is problem

of Vanishing gradients / exploding gradients



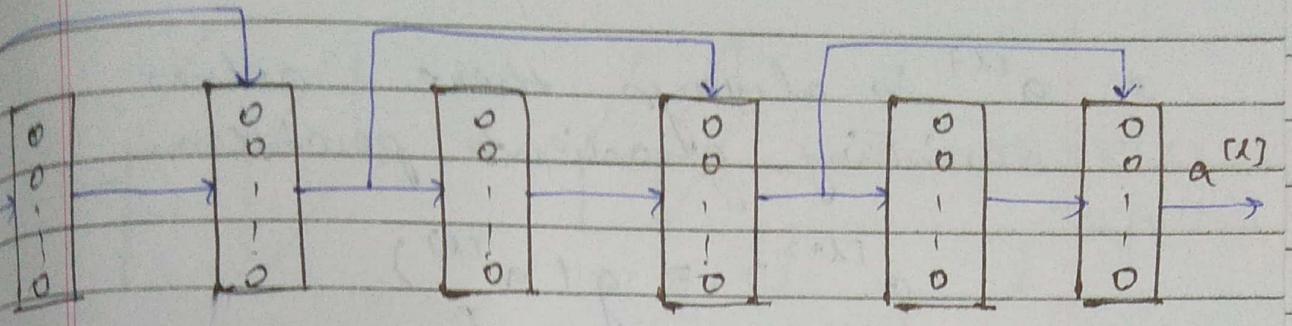
$$z^{(l+2)} = w^{(l+2)} a^{(l)} + b^{(l+2)}$$

$$a^{(l+2)} = g(z^{(l+2)})$$

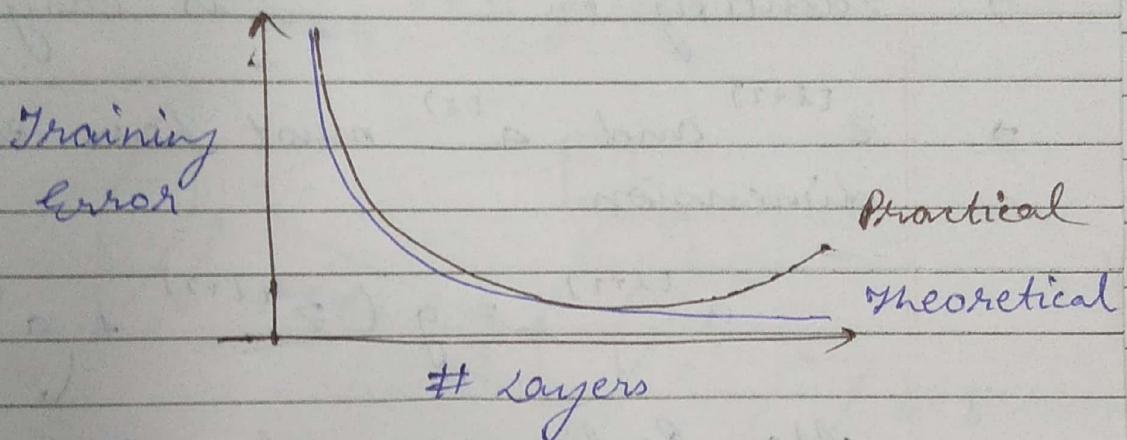
- In case of ResNet, $a^{(l)}$ is going through block B_1 as well as directly to B_2 .
- Here, B_2 is called Residual block $a^{(l)}$.
- So, for ResNet :

$$\checkmark \quad a^{(l+2)} = g(z^{(l+2)})$$

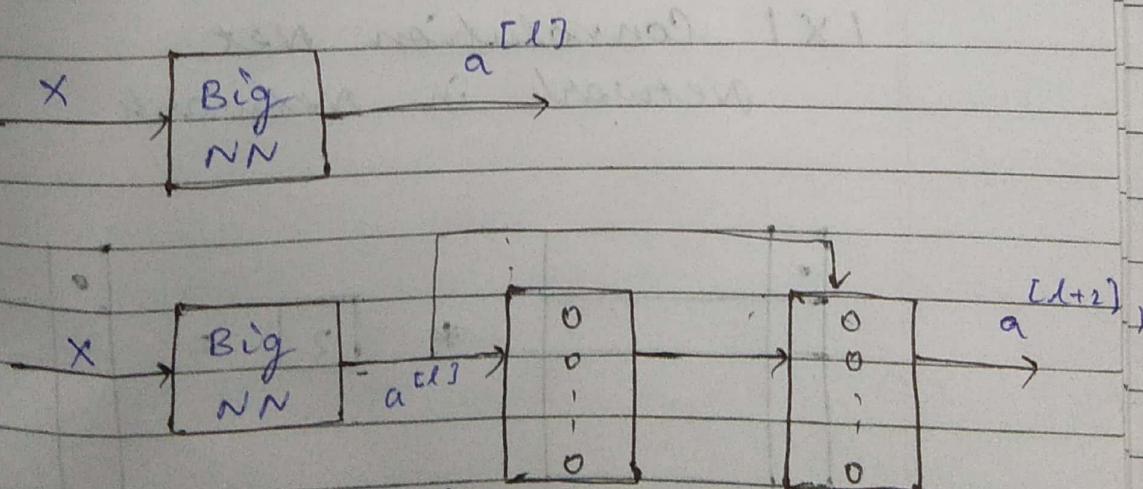
$$\checkmark \quad a^{(l+2)} = g(z^{(l+2)} + a^{(l)})$$



Stacked of Residual block



Deep Conv. Neural Net



$$\begin{aligned}
 a^{[l+2]} &= g(z^{[l+2]}) \\
 &= g(w^{[l+2]} a^{[l+1]} + b^{[l+2]} + a^{[l]}) \\
 &\quad \text{Assume bias = 0 after certain layers}
 \end{aligned}$$

→ $a^{(l)}$ is always there, solves vanishing gradient problem.

$$a^{(l+2)} = g(a^{(l)})$$

Relu : $a \geq 0$

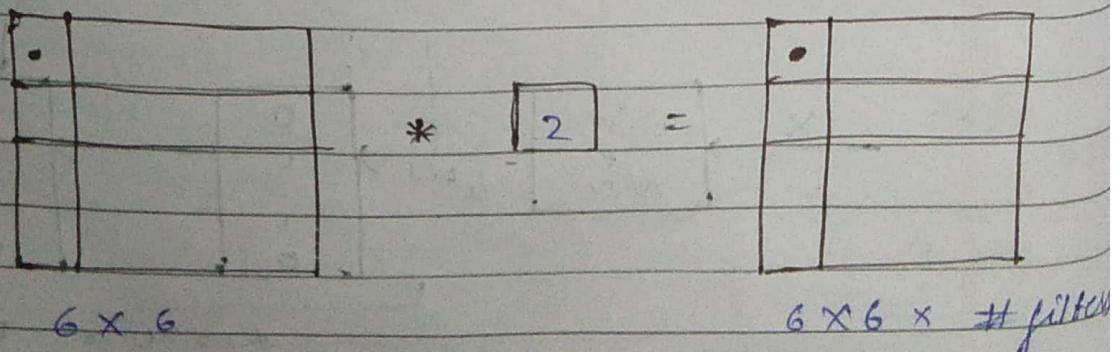
→ Identity fn: is easy to learn

→ $z^{(l+2)}$ and $a^{(l)}$ must have equal dimension.

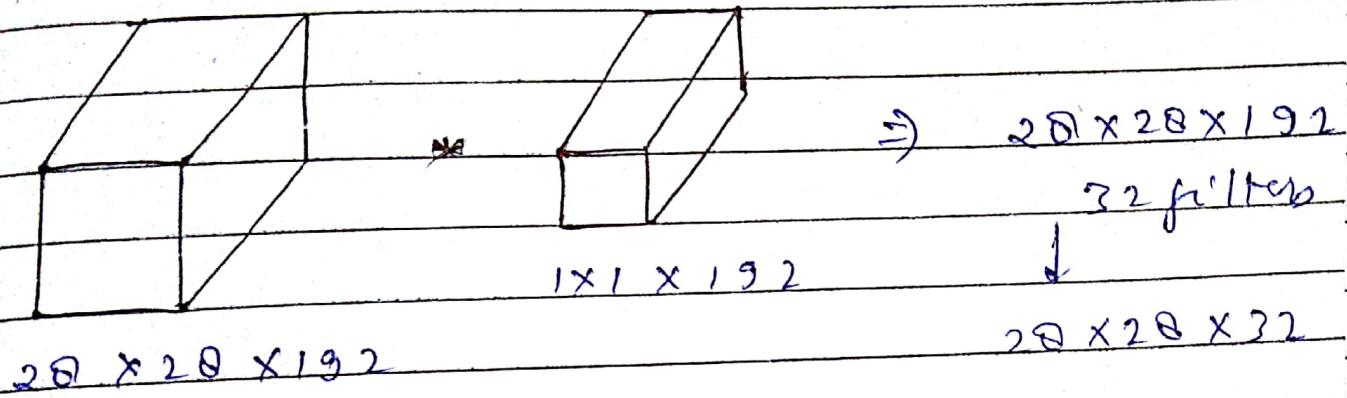
$$a^{(l+2)} = g(z^{(l+2)} + a^{(l)})$$

Max Pooling can also be applied

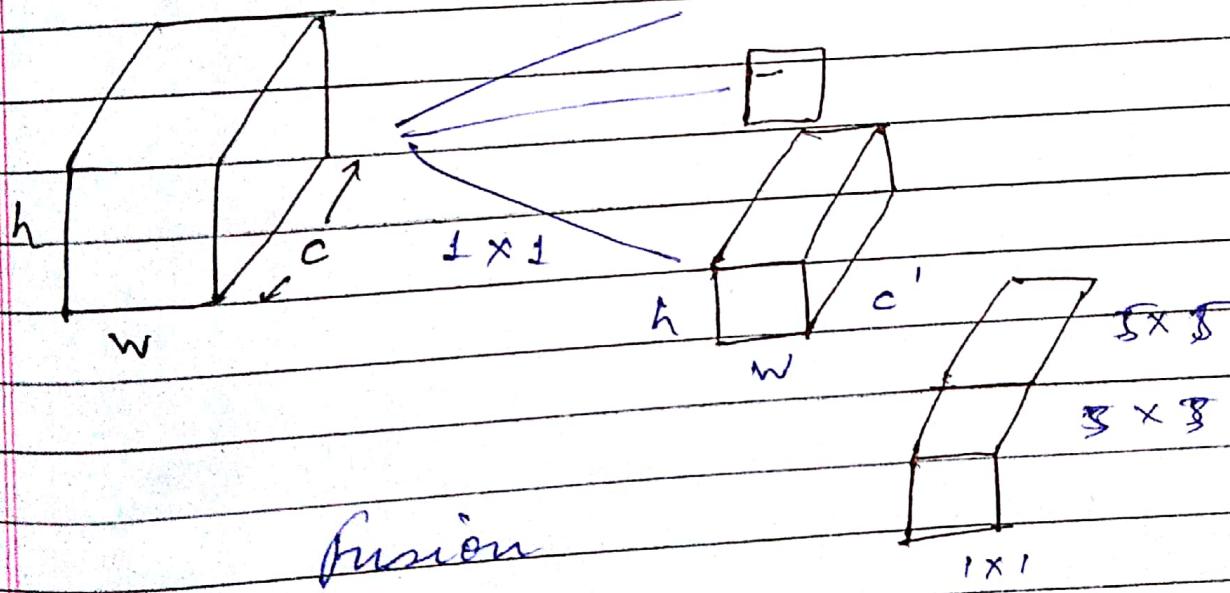
1×1 Convolution Net
Network in Network



Full connected layers act as a conv. net.



Shrink \rightarrow Another application of
 1×1 filter



Sequence Modelling

Feedforward AN

Inputs are given randomly.

(Not appropriate)

Bag of words : vector of unique words in sentence

$$\begin{bmatrix} 1 & 0 & 0 & 0 & \dots & 1 & 0 & 0 \end{bmatrix}^T$$

ex Food is good, not so bad

Food is bad, not so good

Markov Model : Predict t using state

$$y^{(t-1)} \rightarrow y^{(t)}$$

long term dependence

I was in France. time was beautiful

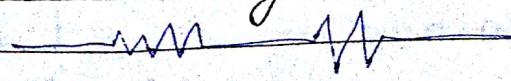
I try to learn

predict (should be French)

→ This model fails here.

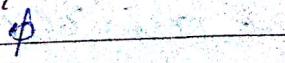
Recurrent AN : It is transformed to sequence modelling

① Speech Recognition :

 It is raining

X Y

② Music Generation :

 - ||| //

3) Sentiment classification

4) MIC Translation

5) DNA Sequence Analysis

ACCTGG --> Protein

6) Name Entity Recognition (NER)
Proper nouns.

→ RNN is supervised model.

Notation:

$x := \text{Bill Clinton was president of}$
Index :- $x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad x^{<4>} \quad x^{<5>} \quad x^{<6>}$
 $\text{America} \quad \quad \quad \quad \quad \quad \quad$

$$T_x = 6$$

$y := [1 \ 1 \ 0 \ 0 \ 0 \ 1] \quad T_y = 6$

NLP → how to represent the words

Vocabulary / Dictionary

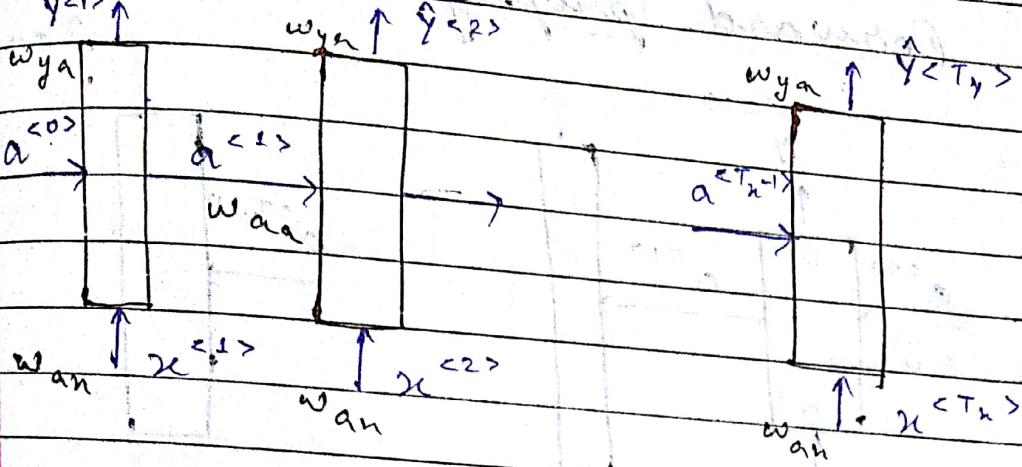
a	1	0
an	0	1
and	0	0
:	0	1
Bill	1	0
chose	0	0
Clinton	0	1
	10,000	

One-hot vector

$(x^{<1>}, y^{<1>}), (x^{<2>}, y^{<2>})$

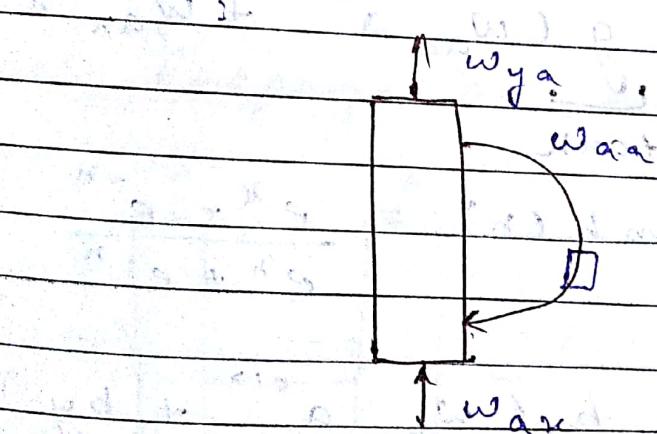
LDA : Latent Dirichlet Allocation

RNN Structure:



Weights are shared.

$y^{(i) < t>} \text{, } i^{\text{th}}$ instance at time t

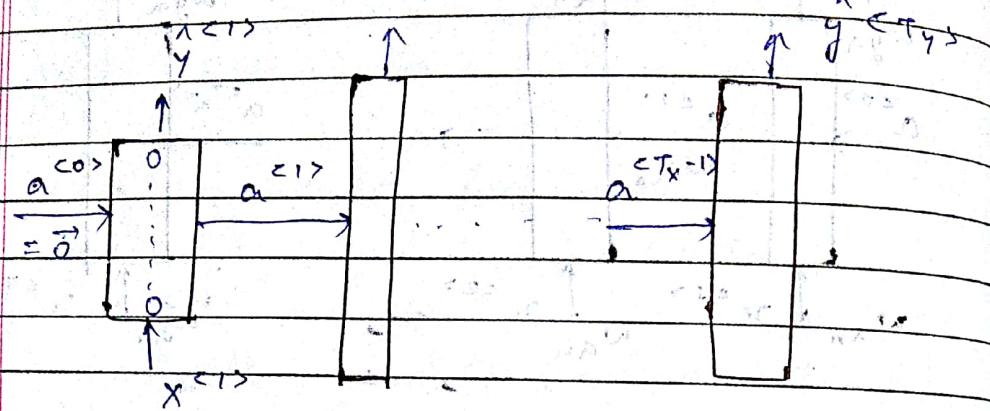


$y^{<3>}$ requires $x^{<1>} x^{<2>} x^{<3>}$

Apple is a computer (in discount)
Apple is a fruit which is delicious

For predicting computer / fruit, we need to know the later part of the sentence. (Drawback of this basic RNN)

Forward propagation



$$a^{<1>} = g(w_{aa} a^{<0>} + w_{an} x^{<1>} + b_a)$$

tanh

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$y^{<1>} = h(w_{ya} a^{<1>} + b_y)$$

sigmoidal / logistic or softmax

$$a^{<t>} = g(w_{aa} a^{<t-1>} + w_{an} x^{<t>} + b_a)$$

$$y^{<t>} = g(w_a [a^{<t-1>}], x^{<t>}) + b_a$$

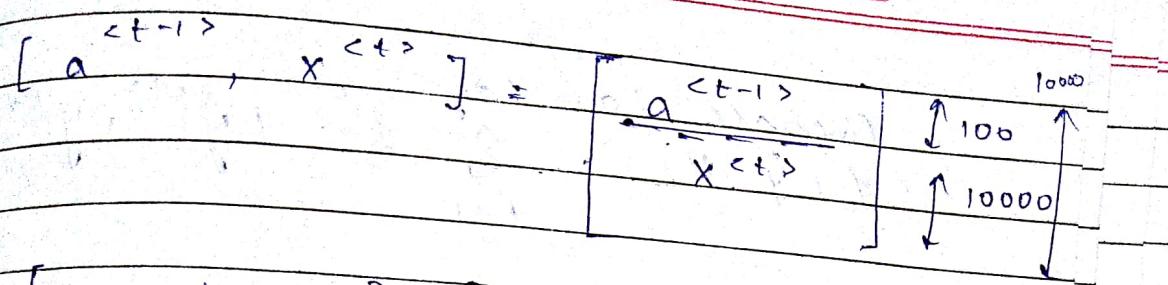
$$w_a = \begin{bmatrix} 100 \\ [w_{aa} \mid w_{an}] \end{bmatrix}$$

100 10,000

$$w_{aa} = (100, 100)$$

$$w_{an} = (100, 10000)$$

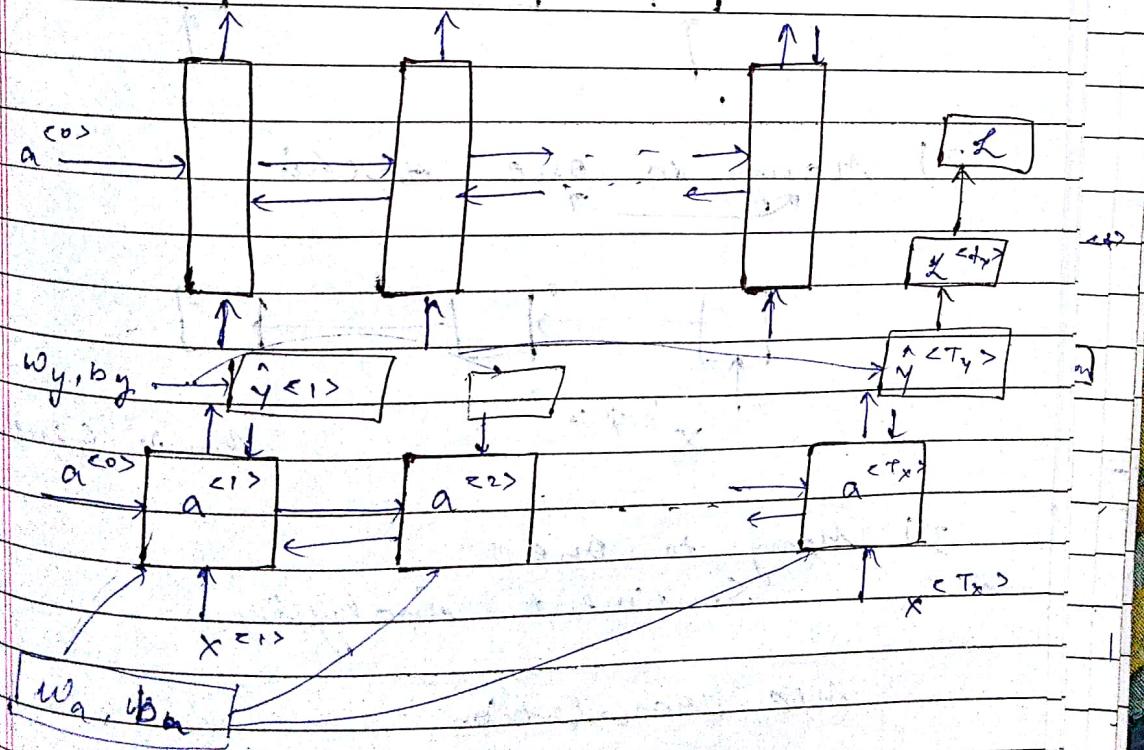
[Assume, dim (a<0>) = 100]



$$[W_{aa} \cdot [W_{an}]] \cdot [a^{t-1}] + W_{ax} \cdot x^{t+} = W_{aa} \cdot a^{t-1} + W_{an} \cdot x^{t+}$$

$$\hat{y}^{t+} = h[W_a(a^{t-1}) + b_y]$$

Back Propagation through time :



$$L^{t+} (\hat{y}^{t+}, y^{t+})$$

Single Prediction

$$= -y^{t+} \log \hat{y}^{t+} - (1-y^{t+}) \log (1-\hat{y}^{t+})$$

Cross Entropy Loss

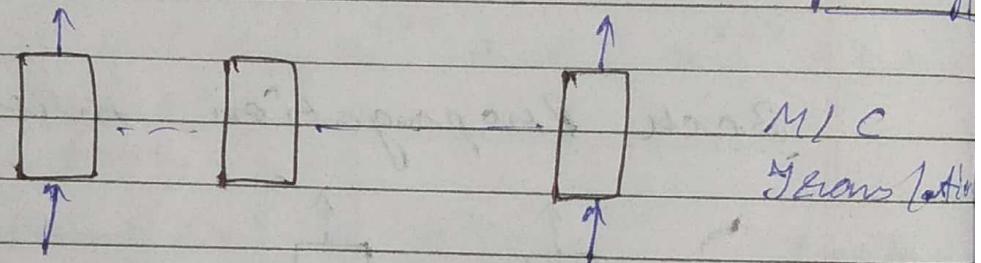
$$\text{Overall loss : } L(\hat{y}, y) = \sum_{t=1}^{T_y} L(\hat{y}^{(t)}, y^{(t)})$$

Main Feature of RNN:
Hidden state

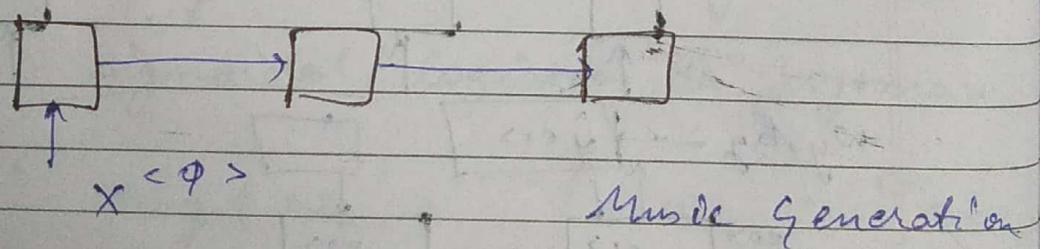
1) Different Architecture:

many to many architecture

$T_n \neq T_y$



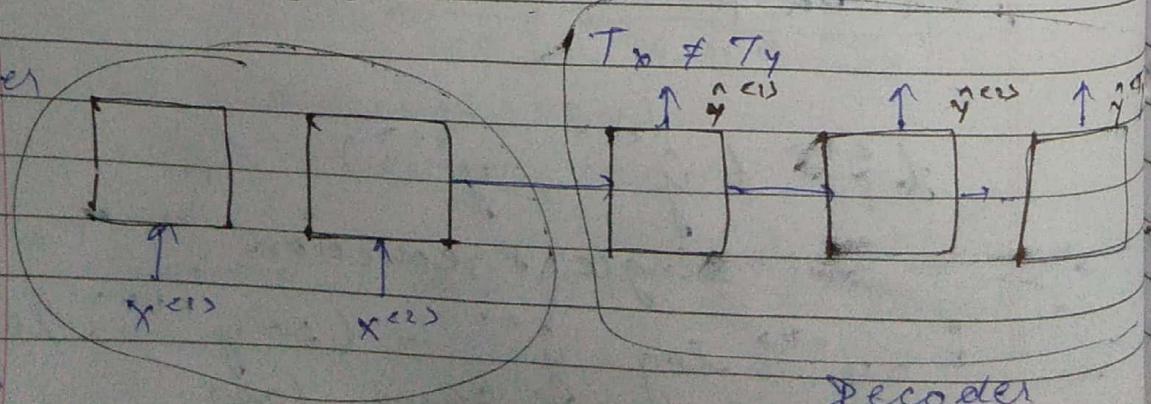
2) Many to one Archi



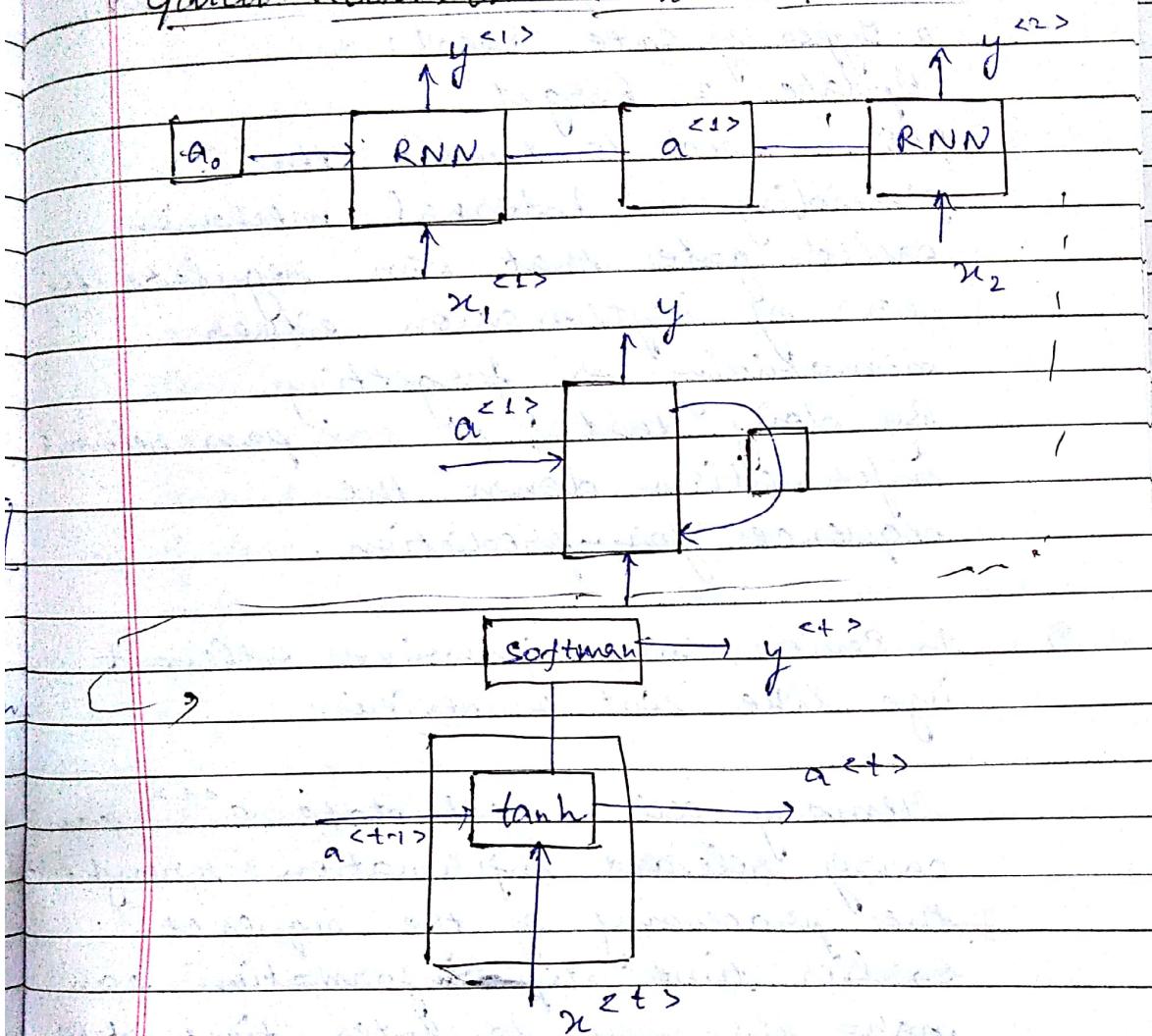
3) Many to one:
Sentiment Analysis.

MIC Translation.

Encoder



Decoder

Gated Recurrent Unit (GRU)

The cat - - - - - (was) running
were

Back propagation training.

weight adjustment -

Failed to memorize the sequence,
resulting in bad prediction.

Vanishing Gradient : No or less
change.

Exploding Gradient : Large change

⇒ Normal RNN → short term memory.

Solution is Gated Recurrent Unit,
 & Types of gate used:

Update & forget

GRU is used to handle the
 providing an internal mechanism
 called gate that can regulate the
 flow of information either
 memorizing or forgetting.

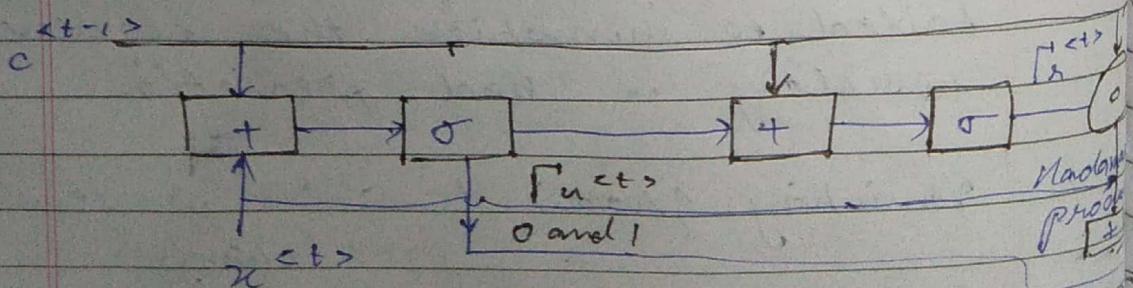
By doing that, it can pass relevant
 information down the
 sequence for prediction.

- ⇒ In Review, only memorize relevant info like cost & review.

Memory cell, cell state $c^{<t>}$ can
 carry relevant information throughout
 the processing of the sequence.

Earlier time step information can
 make its way to later time steps,
 reducing the effect of short term
 memory.

Gates are nothing but different neural
 network.



$$\#1 \text{ Update Gate: } f^{<t>} = \sigma(w_u[n^{<t>}])$$

Start calculating update
for time state t using

- 1) $n^{t+>}$ is plugged into network, multiplied by its own weight
- 2) Then c^{t+1} is mult. with the activation weight.
- 3) Results is & a sigmoid function is applied to obtain results b/w 0 & 1.
4. Update gate determines how much info. from the previous time state is learnt & passed along the future sequence cycle.

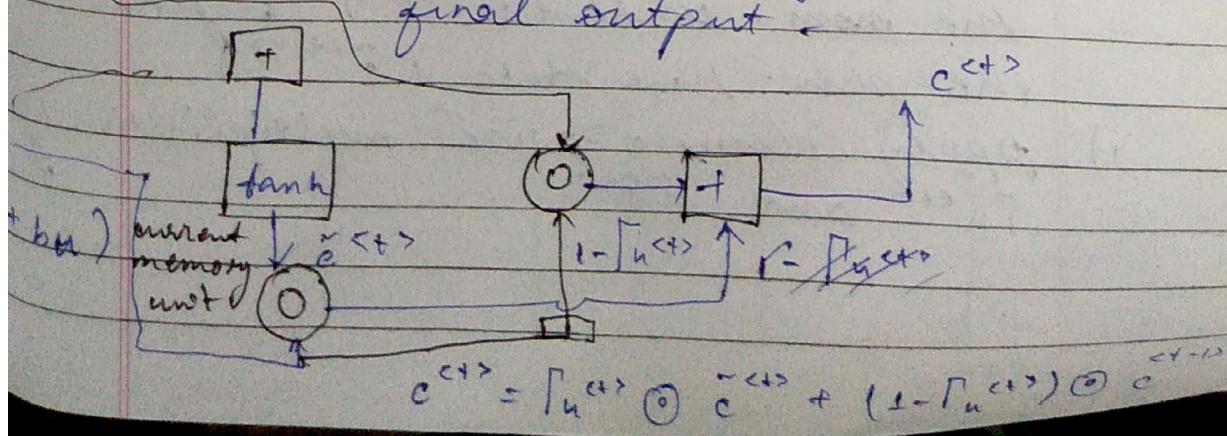
#2 Reset Gate:

$$r^{t+>} = \sigma(w_r [c^{t+1}, n^{t+>}] + b_r)$$

This gate is used to decide how much of the past information to forget.

$$\#3. c^{t+>} = \tanh(w_c [r^{t+>} \odot c^{t+1}, n^{t+>}] + b_c)$$

Current implies how exactly the gates (update & reset) now influence the final output.



Reset Gate - How to use:

Net's value of memory ($\tilde{c}^{<+>}$) which used the reset gate to store the relevant info. from past.

- 1) multiply the input $x^{<+>}$ with weight w_{xR} $c^{<+>}$ with corresponding weight.
- 2) calculate the Hadamard product b/w the reset gate & $w_{cR} c^{<+>}$. It determines what to remove from the previous time state.
- 3) 1 & 2 are added & finally tanh is applied to produce $\tilde{c}^{<+>}$.
- 4) Final memory at current time step: $c^{<+>}$ which requires $\tilde{c}^{<+>}$.

So, the

The network needs to collect from the current input & passes it down to the network. To do that the update gate is needed. It determines what to collect from current memory content ($\tilde{c}^{<+>}$) & previous mem-content $c^{<-1>}$.

Now to use Update Gate:

It determines what to collect from the mem-content ($\tilde{c}^{<+>}$) & from the prev-time state ($c^{<-1>}$)

- 1) Apply element-wise multiplication of Γ_{uR} and $\tilde{c}^{<+>}$.

LSTM : It also use Output gate.

LSTM is more powerful than GRU.

as it does not consider assumption

$$e^{<t>} = a^{<t>}$$

→ Long short term memory.

$$\tilde{c}^{<t>} = \tanh(w_c [a^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u^{<t>} = \sigma(w_u [a^{<t-1>}, x^{<t>}] + b_u)$$

$$\Gamma_f^{<t>} = \sigma(w_f [a^{<t-1>}, x^{<t>}] + b_f)$$

$$\Gamma_o^{<t>} = \sigma(w_o [a^{<t-1>}, x^{<t>}] + b_o)$$

$$c^{<t>} = \Gamma_u^{<t>} \odot \tilde{c}^{<t>} + \Gamma_f^{<t>} \odot c^{<t-1>}$$

$$a^{<t>} = \Gamma_o^{<t>} \odot \tanh c^{<t>}$$

