

## MLP vs CNN

→ MLP is referred as Vanilla NN, when they have a single hidden layer.

→ In case of CNN, the neurons in a layer will ~~only be connected to~~ <sup>sparsely connected</sup> all of the neurons in a fully connected manner in MLP.

### (2) Dimension Information

CNN take into account the dimensional information of an image which gives them an advantage over MLP's for image classification tasks.

### (3) Better for Complex Images

CNNs perform better than MLP's on for complicated images because they can understand spatial relationships between pixels.

### (4) Parameter Sharing

In CNN, the same filter is used across the entire image during convolution operation. This significantly reduces the no. of parameters that need to be learned.

In MLP, each neuron on a layer is connected to all neurons on the previous layer.

### (5) Reduced memory use

CNNs reduce memory usage by sharing weights among spatial dimensions.

### (6) Faster Computation

CNNs compute faster than MLPs because of the shared weights.

~~Input (32x32x3)~~

~~Conv(5x5)~~. ~~@28, 28~~

$$S=1, G$$

poolmax | 64@ 14, 14

CONV2 ( $g=5, S=1$ ) 129 @ 10x10

Page 2

$$w_0 + \int_{\partial M} c_1(f) + \tau w_1(x) = 0$$

$$(-5)^2 = +25$$

$$pd + \int_p W \circ \gamma d) = \frac{1}{\pi} \delta$$

$$f((g)) \rightarrow g = g$$

( off 2-21pm 2 )

## Pooling :

- Pooling layers aim to gradually reduce the dimensionality of the representation, and thus further reduce the no. of parameters and the computational complexity of the model.
- A pooling function replaces the Op of the net at a certain location with a summary statistic of the nearby op.

① Max pooling layers

② Average pooling

③ Global pooling

## Pooling Layers :

- Max pooling condenses information in the I/p volume,
- we take only the maximum entry of each window.

12	20	30	0
8	12	2	0
31	70	37	4
112	100	25	12


$2 \times 2$  Max-pool.

20	30
112	37

$2 \times 2$  Avg-pool

13.25	8
79	19.5

- Average pooling condenses information in the volume,
- we take average of entries of each window

## Padding :

Problem with simple convolution layers

- For a gray scale ( $n \times n$ ) image and ( $f \times f$ ) filter/Kernel the dimensions of the image resulting from a convolution operation is  $(n-f+1) \times (n-f+1)$ .
- The pixel in the corner will only get covered once time but if you take the middle pixel it will get covered more than once basically what does that means is we have more info on that middle pixel so there are 2 main downsides
  - \* Shrinking outputs
  - \* Loosing information on corners of the image.

To overcome this, we can introduce padding to an image.

## Padding :

- Padding is simply a process of adding layers of zeros to our IP images so as to avoid the problems. This prevents shrinking.
- If  $P = \text{no. of layers of } 0's$  added to the border of the image, then our  $(n \times n)$  image becomes  $(n+2P) \times (n+2P)$  image after padding.
- So, applying convolution operation with  $f \times f$  filter outputs  $(n+2P-f+1) \times (n+2P-f+1)$  images.

0	0	0	0	0	0	0	0
0							
0							
0							

[zero padding]

## Type of Padding

### ⇒ Valid Padding

→ It implies no padding & all the IP image

in its valid/considered shape

$$\rightarrow [(n \times n) \text{ image}] * [(f \times f) \text{ filter}] \rightarrow [(n-f+1) \times (n-f+1) \text{ image}]$$

Here \* means convolution operation

### ⇒ Same Padding

→ In this, add ' $p$ ' padding layers such that the op image has the same dimensions as the IP image.

$$\therefore [(n+2p) \times (n+2p) \text{ image}] * [(f \times f) \text{ filter}] \rightarrow [(n \times n) \text{ image}]$$

$$\rightarrow \text{In general, } p = (f-1)/2 \quad (\text{because } n+2p-f+1 = n)$$

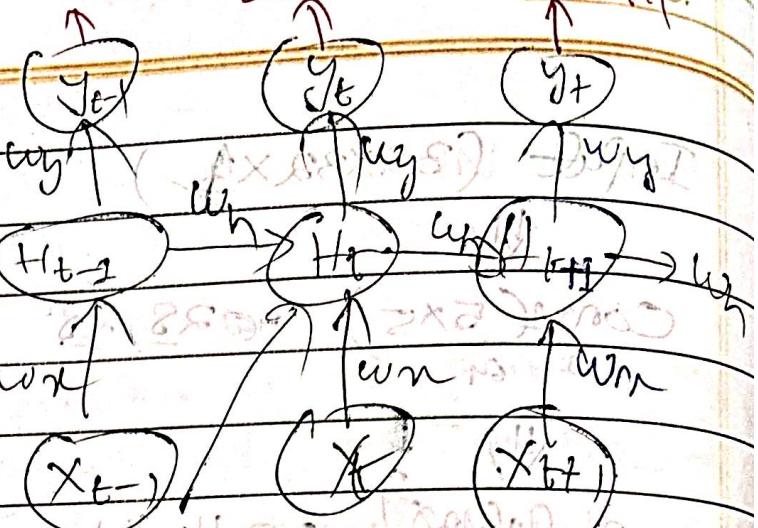
# RNN

Category

Date \_\_\_\_\_  
Page No. \_\_\_\_\_



~~Input~~ ~~Output~~



unfold Representations

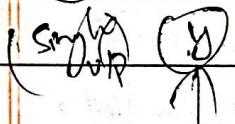
RNN is used to process sequential data and make predictions based on that data.

Types of RNN → uses a hidden layer that is specific in form about a sequence.

→ Has a memory that stores all information about a sequence.

① one to one

② one to many.



$$z_h = [x_t w_x + h_{t-1} w_h] + b_h$$

$$h_t = \sigma(z_h)$$



$$z_y = (h_t w_y) + b_y$$

$$y_t = \sigma(z_y)$$

ve e.g.   
con lang Translat  
 (English French)  
 smile w/p  
 scimile w/p

(single step)



(multiple op)

(single step)

(image to text)

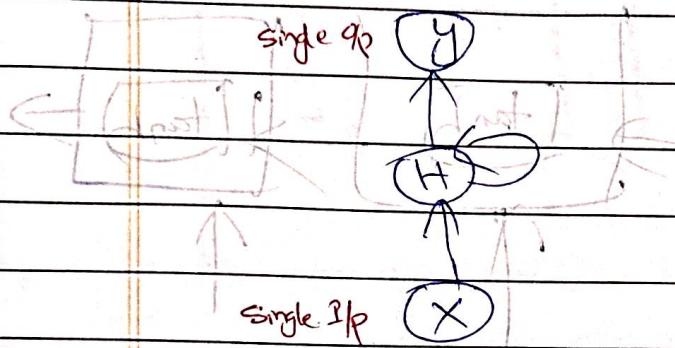
Use Case

Image captioning  
Text description of image  
Text description of image

## Types of RNN

- ① One to one
- ② One to many
- ③ Many to one
- ④ Many to many

### ⇒ One to one RNN :

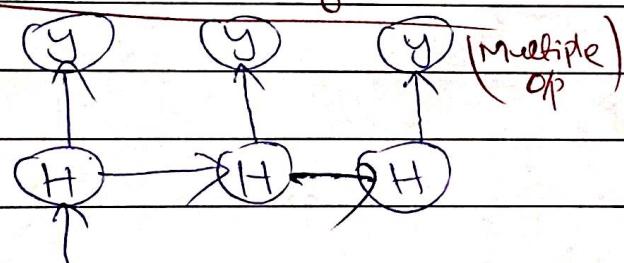


Use case

Vanilla RNN

Used for regular timestamp data

### ⇒ One to many



Use case

→ Image Capturing  
→ Music generation

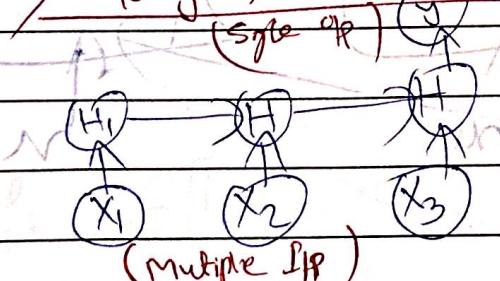
Op is sequence

(X) (single I/p)

I/p :- Image

O/p :- Text descriptn of I/p (img)

### ⇒ Many to one



Use case

Sequence analysis

Mice eats cheese

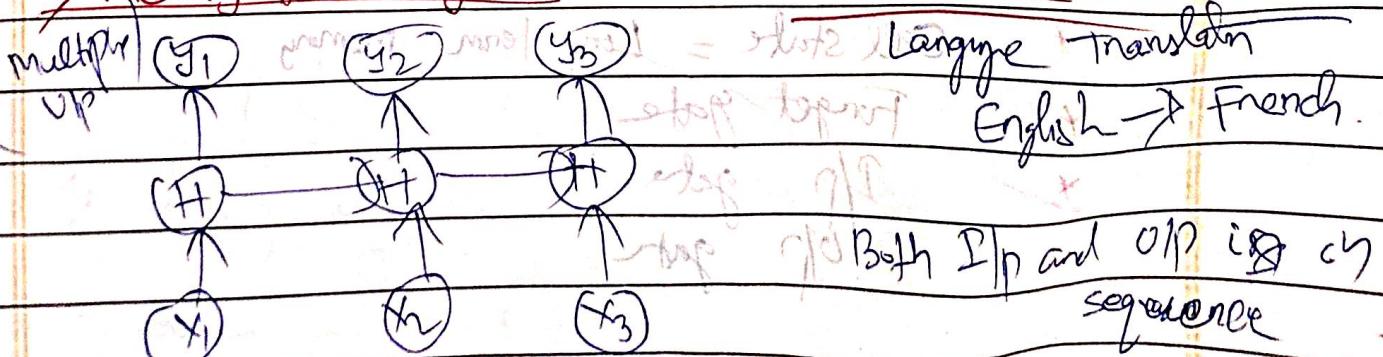
I/p is sequence

English

O/p is fixed

numbers 1 2 3 M T S A size restr.

### ⇒ Many-to many



Use case

Language Translation

English → French

stop stop

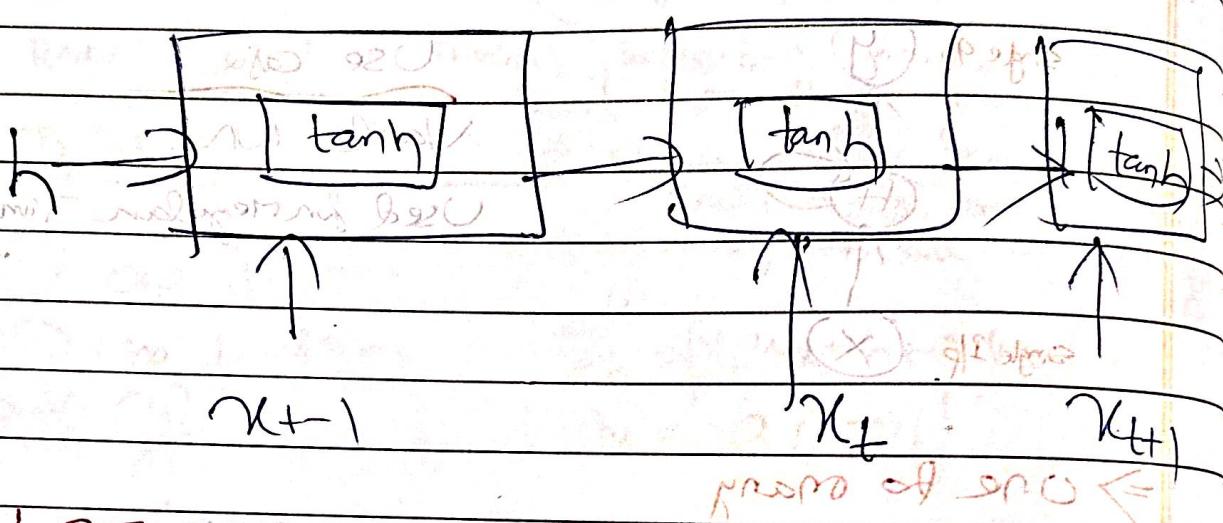
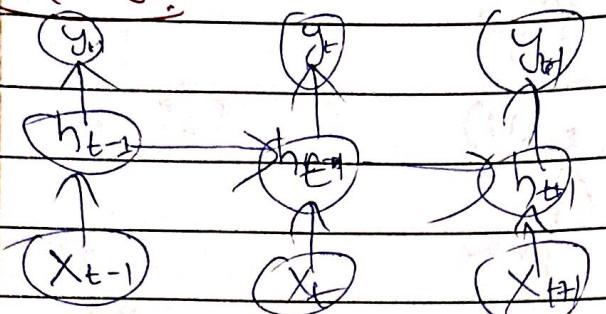
stop stop

stop stop

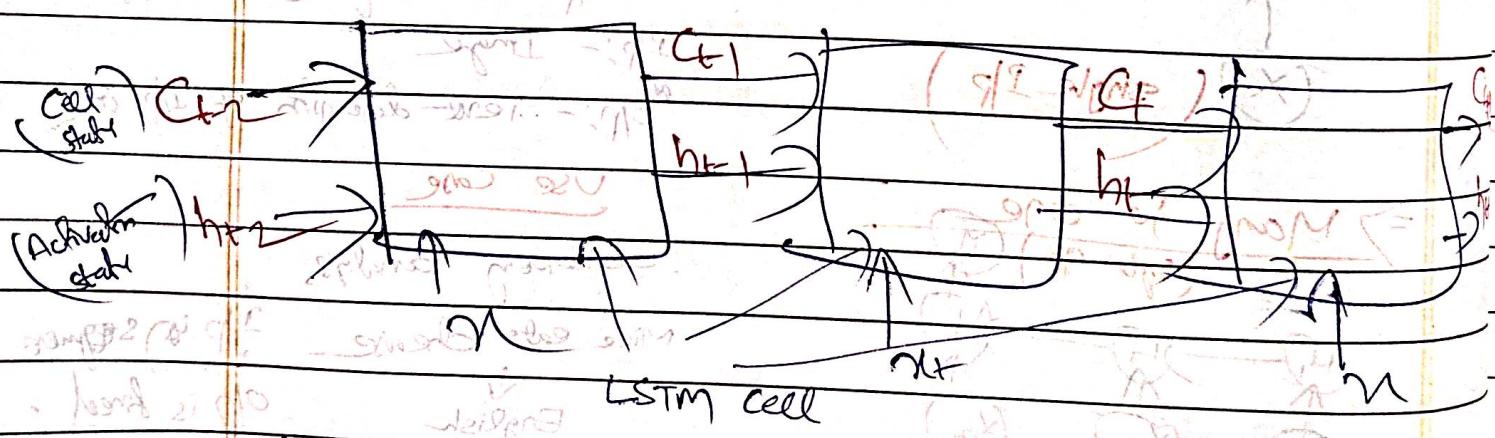
Both I/p and O/p is sequence

sequence

RNN:



LSTM:



A LSTM cell contains

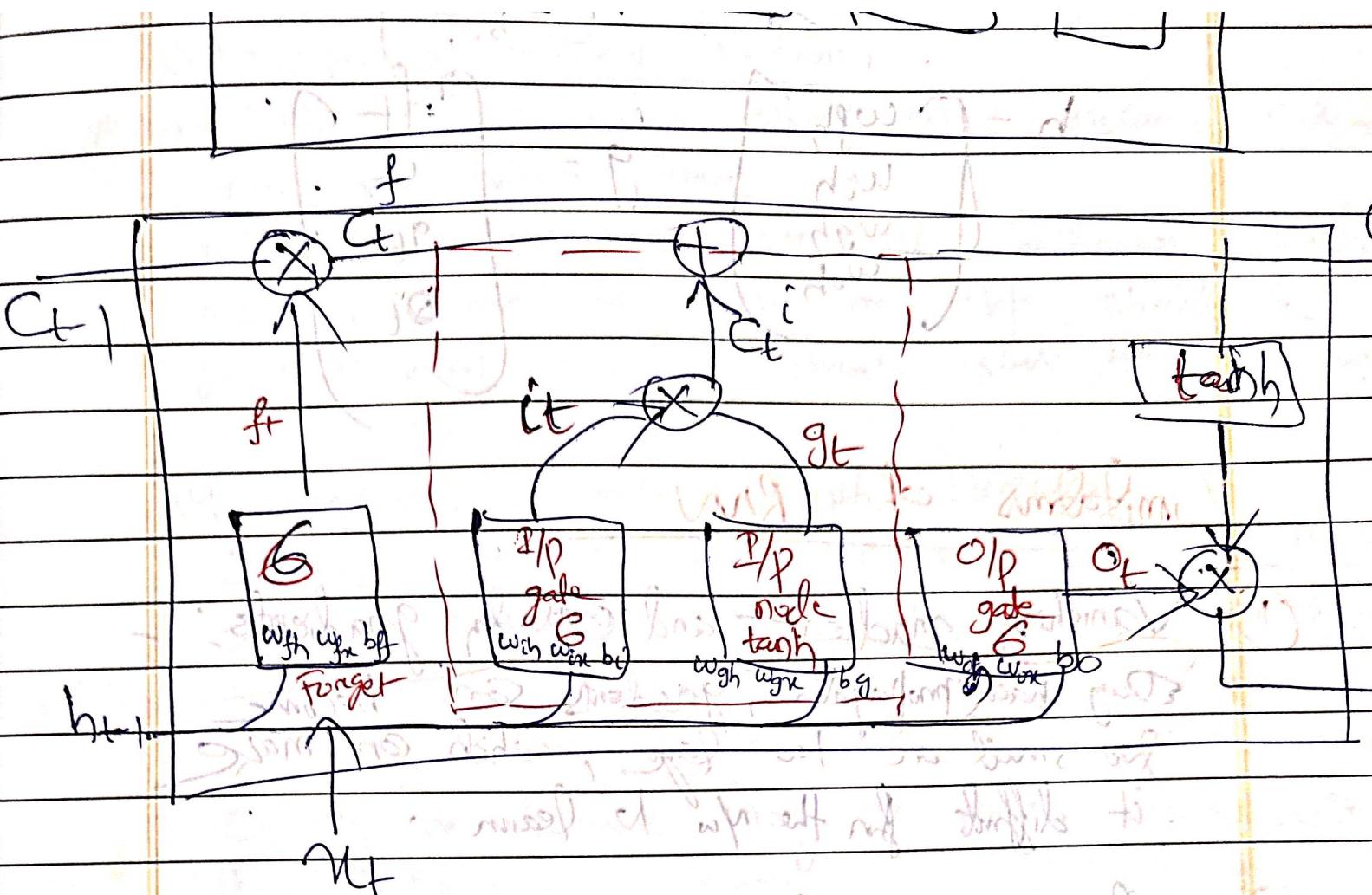
\* A simple RNN cell

\* Cell state = Long term memory

Forget gate

\* I/P gate

O/P gate



$$f_t = G \left[ (w_{fh} * h_{t-1}) + (w_{fx} * x_t) + b_f \right]$$

$$i_t = G \left[ (w_{ih} * h_{t-1}) + (w_{ix} * x_t) + b_i \right]$$

$$g_t = G \left[ (w_{gh} * h_{t-1}) + (w_{gx} * x_t) + b_g \right]$$

$$C_t = i_t * g_t$$

$$C_t = C_{t-1} + C_t^f$$

$$O_t = G \left( w_{oh} * h_{t-1} \right) + \left( w_{ox} * x_t \right) + b_o$$

$$h_t = \tanh(c_t) * o_t$$

$$w = \begin{bmatrix} w_f \\ w_i \\ w_g \\ w_o \end{bmatrix}, b = \begin{bmatrix} b_f \\ b_i \\ b_g \\ b_o \end{bmatrix}$$

$$h = \begin{bmatrix} w_{fh} \\ w_{ih} \\ w_{gh} \\ w_{oh} \end{bmatrix}, g = \begin{bmatrix} f_t \\ i_t \\ g_t \\ o_t \end{bmatrix}$$

## Limitations of RNN

(1) Vanishing gradient and exploding gradients: -  
During back propagation, gradients can become too small or too large, which can make it difficult for the network to learn.

(2) Difficulty with long sequences

RNNs can have trouble capturing long-term dependencies and often forget earlier information in long sequences.

(3) RNNs can be slow to train compared to feed forward

(4) Difficulty often pretty O/P

The O/P of an RNN can be difficult to understand especially while dealing with complex inputs.

## Advantages of LSTM

### ① Long term dependency learning

LSTM can capture long-term dependencies and can capture complex patterns in sequential data.

### ② Vanishing gradient problem

→ LSTM can avoid vanishing gradient problem.

→ LSTM uses gating mechanisms to control the flow of information and gradients, allowing them to learn and retain information over many time steps.

### ③ Selective memory retention

→ LSTM can selectively retain or discard information from the hidden state.

\* The forget gate in LSTM cells determines which information from the previous time step should be forgotten and which information should be retained.

### ④ Handling Variable Length Sequence

LSTM can handle variable length I/P

sequences by dynamically adjusting their internal state.

### ⑤ Handling noisy or missing data better

LSTM can handle noisy or missing data better than RNNs.

## Attribute

## CNN

(NT2) RNN sequential b/t

① Datatype :- spatial data

like image, signal

sequential data

series data

another - filtering proteins (c)

② I/p & O/p :- In CNN, the I/p size & O/p size are fixed. Here, I/p and O/p size can vary.

③ Architecture :- Feed forward NN and Recurrent feed forward and Recurrent

filtering proteins symbols (d)

④ Parallel processing possible Sequential processing states inherit soft max.

⑤ Memory :- Limited memory of part of buffers and memory cell may be lost.

⑥ Use Cases :- Facial Recognition, NLP, sentiment analysis, speech classification.

→ It is Batch Normalization (BN)

→ It is a technique in NN that improves speed and stability by centering data values around a common scale.

→ BN helps to reduce overfitting and improve the generalization ability of the model.

\* It is particularly useful for very deep networks, as it can help to reduce the internal covariate shift that can occur during training.

→ (Change in the distribution of your actions due to change in parameters during training)

## Benefits of BN

- improves the efficiency and reliability of the model.
- Reduce the internal covariate shift that is
- Decreases the drop out rate (amount of info lost before going to next layer)
- Allow for use of higher learning rates (to combat the risk of divergence)

Q- How non-linearity is introduced in CNN network?

- Non-linearity is introduced into a CNN through activation functions that are applied to the output of a neuron.
- Activation functions are typically applied after each convolution layer and FC layer, but not after fully connected layer.

## Advantages of MLP

- Complex modelling :- Handling non-linear relationships & Complex Patterns (X, Y, H1, H2)
- Flexibility :- Adaptable to diff. problem domains and achieves higher accuracy.

## Batch Normalization

→ BN normalizes the inputs to each layer so that they have zero mean and unit variance during training.

→ It uses the statistics (mean and variance) of each mini batch to scale and shift activations through learnable parameters  $\gamma$  and  $\beta$ .

$$\text{For a batch of activations } \mathbf{x}: \quad \mathbf{x} = \mathbf{x} - \bar{\mathbf{x}}_{\text{batch}}$$

$$y = \gamma x + \beta \quad \gamma = \sqrt{\sigma^2_{\text{batch}} + \epsilon}$$

where,  $\bar{\mathbf{x}}_{\text{batch}}, \sigma^2_{\text{batch}}$  are the mean and variance computed over mini batch.

∴  $\gamma$  and  $\beta$  are learnable scale/shift parameters.

- The IP Layer has no parameter (ie,  $m=0$ )
- The parameters in the CONV Layer = number of filters \* shape of height of filter  
 $= (\text{shape of width of filter} * \text{shape of height of filter}) * (\text{number of filters in the prev. layer})$
- Pooling Layer has no parameter (ie,  $m=0$ )

→ Parameters in FC Layer = current layer  
 $= (\text{current layer } C * \text{prev layer } P_2) + 1 * C$

Example:  $\# \text{filters} * \text{shape of filter} * \text{activation function} = \# \text{parameters}$

1. Input Layer  $32 \times 32 \times 3$   $3072$

2. CONV1 ( $f=5, s=1$ )  $(28 \times 28, 8)$   $6272$   $608$

3. Pool L1  $(14, 14, 8)$   $1568$   $10$

4. CONV2 ( $f=5, s=1$ )  $(10, 10, 16)$   $1600$   $3216$

5. Pool L2  $(5, 5, 16)$   $400$

6. FC3  $(120, 1)$   $120$

7. FC4  $(84, 1)$   $84$

8. Softmax  $(10, 1)$   $10$

$$10 + 84 + 1 = 100$$

## Advantages of MLP

- MLP can learn complex pattern. — Non-linear relationships in data (single layer perception can't)
- Real time training all MLP can train models on real time
- Can handle large amount of data.
- MLP can make quick prediction
- MLP can achieve accuracy with small samples.
- MLP are versatile tools and can be used for various of tasks (eye & speech recognition, classification, regression, etc.)

## Batch Normalization in CNN

→ BN in CNN is normalizing the activations within each layer, thereby stabilizing the training process & significantly accelerating convergence.

### Benefits of BN in DL

- i) Improved training stability : — BN stabilizes the training process (by reducing internal covariate shift) making it less sensitive to initial weights and learning rate.
- ii) Faster convergence : Normalized activations allow the network to converge more quickly, reducing the no. of epochs.
- iii) Regularization : BN introduces a slight regularization effect by adding noise to the activations, helping to prevent overfitting.
- iv) Enhanced generalization : Model trained with BN often generalize better to unseen data, leading to improved performance on test sets.

## Data Augmentation

→ It is a technique used to artificially increase the training dataset by applying various transformations to the existing data.

e.g. ~~Foggy~~

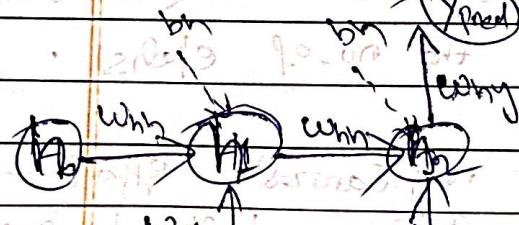
- \* Flipping (horizontal / vertical)
- \* Rotation, zooming, cropping, color jittering (brightness, contrast), adding noise.
- \* GAN (Generative Adversarial Network)

→ Data augmentation prevents overfitting and improves model generalization.

Q - Draw a RNN with  $I/P \rightarrow C2 \rightarrow \text{hidden} \rightarrow o/p$ . Calculate  $X_{\text{pred}}$  for the following given parameter values in the Table with

$$I/P \quad X = [1, 0], \quad X_0 = [0, 1], \quad h_0 = [0, 0]$$

Input	Layer	Weights	Bias	Output
$I/P$ to hidden		$\begin{bmatrix} 0.5 & 0.8 \\ -0.3 & 0.2 \end{bmatrix}$	$0.5, 0.2$	
Hidden to hidden		$\begin{bmatrix} 0.1 & -0.4 \\ 0.9 & 0.3 \end{bmatrix}$	$0.0, 0.1$	
Hidden to o/p		$\begin{bmatrix} 1.0 & -1.0 \end{bmatrix}$	$0.0$	



Step 1  
Compute 1st hidden state  $h_1$ .

$$o_1 = W_{h1}x_1 + W_{h1}h_0 + b_1$$

$$= \begin{bmatrix} 0.5 & 0.8 \\ -0.3 & 0.2 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} + \begin{bmatrix} 0.1 & -0.4 \\ 0.9 & 0.3 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} + \begin{bmatrix} 0.0 \\ 0.1 \end{bmatrix}$$

$$= \begin{bmatrix} 0.5 \\ -0.3 \end{bmatrix} + \begin{bmatrix} 0.0 \\ 0.1 \end{bmatrix} = \begin{bmatrix} 0.5 \\ -0.2 \end{bmatrix}$$

$$h_1 = \tanh(a_1) = [\tanh(0.5), \tanh(-0.2)] \\ \approx [0.462, -0.197]$$

Step-2

Compute 2nd hidden state  $h_2$

$$a_2 = W_{xh} \cdot x_2 + W_{hh} \cdot h_1 + b_h$$

$$= [0.5 \ 0.8] [0.1] + [0.1 \ -0.4] [0.462] + [-0.3 \ 0.2] [-0.197]$$

$$= [0.2 \ 0.8] + [0.175 \ 0.175] + [0.1 \ 0.1] \\ = [0.925 \ 0.333]$$

$$h_2 = [\tanh(a_2)] \approx [\tanh(0.925), \tanh(0.333)] \\ = [0.729, 0.321]$$

Step-3

Compute Ypred.  $W_{hy} = [0.0, -1.0]$

$$y_{pred} = W_{hy} \cdot h_2 + b_y \\ = [0.0 \ -1.0] \cdot [0.729 \ 0.321]$$

$$= 0.729 - 0.321 = 0.408$$

## Limitations of BN

### ① Sensitive to batch size

BN estimates mean/variance across mini batch

\* with large batch sizes, estimates are stable

\* with small batch sizes (e.g. 1-4) statistics become

noisy and unreliable, harming training stability and performance

e.g.  $\text{mean} = \frac{1}{n} \sum_{i=1}^n x_i$ ,  $\text{var} = \frac{1}{n} \sum_{i=1}^n (x_i - \text{mean})^2$

In medical imaging or high-resolution image models where GPU memory restricts batch size to small values, BN's estimation error increases, slowing learning or degrading accuracy.

### ② Training vs Inference mismatch

→ During training, BN uses mini-batch statistics, but during inference, it uses running (moving average) statistics

→ If the data distribution shifts b/w training and testing time, this mismatch can hurt performance.

### ③ Not ideal for sequential/variable length data

In recurrent architectures (like RNNs or LSTM), sequences may have different lengths and temporal dependencies, making it harder to compute meaningful batch statistics at each time step.

## Layer Normalization

→ It normalizes all features of each individual sample, rather than across the batch.

→ It computes mean and variance per data point across its features, meaning it is independent of batch size.

How it works?

→ suppose a neural layer has a features vector for one sample

$$\mathbf{x} = [3, 5, 2, 8]$$

LayerNorm Computation

\* Mean and Variance across these 4 features for this single sample

\* Normalizes each activation accordingly

\* Applies learned scale / shift

Benefits over Batch Normalization

→ works well even with batch size = 1

→ Ideal for RNNs and Transformers

Sequence length varies

→ No mismatch between training / inference — same computing always.

### Group Normalization

→ Group Norm breaks channels into groups and normalized within each group for each sample.

instead of

\* Normalizing across batch (BN)

\* Normalizing across all features (LayerNorm)

→ Group Norm splits channels into  $G$  groups and computes statistics within each group for each sample. This balances feature locality and batch independence.

e.g.

for a layer with  $G=8$  channels, and  $G=8$  groups.

\* Each group has 8 channels.

\* Compute mean / variance within each group across features in that group.

\* Normalize and apply group specific learnable parameters.

## Benefits of GroupNorm

### ① Batch size Independence

GroupNorm compute statistics within a single sample rather than across a batch, so, they work reliably even with small or varying batch sizes.

### ② Better for Sequences

Ideas for RNNs and transformers where sequence length varies over time.

### ③ Stable during Inference

The same normalization is applied at training and inference time - no distribution mismatch.

## Loss functions in Autoencoder

### ① Reconstruction Loss

This is the primary loss function used in autoencoders. It measures how well the network reconstructs the input to output.

#### (i) Mean Squared Error (MSE)

Used for continuous I/Os.

$$L_{MSE} = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{x}_i)^2$$

- \* Penalizes large reconstruction errors.
- \* Encourages the o/p to match the I/P closely.

#### (ii) Binary Cross Entropy (BCE)

gt. is used for binary or normalized data (e.g. pixel values in [0, 1]).

$$L_{BCE} = \sum_{i=1}^n x_i \log(\hat{x}_i) + (1-x_i) \log(1-\hat{x}_i)$$

- \* Suitable when o/p's represent probability.

## Variants of Loss Functions

### (2) Regularized Auto encoder

Some auto encoder Variants add additional terms to the basic reconstruction loss.

#### (i) Sparse Auto encoder loss

Adds a sparsity penalty (e.g.  $L_1$ ) to encourage most latent activations to be zero.

$$L = L_{\text{reconstruction}} + \lambda \cdot \text{Sparsity Penalty}$$

This forces the model to pick only the most relevant features.

#### (ii) Denoising Auto encoder loss

The  $x'$  is 1st corrupted with noise  $\eta'$  and the model is trained to reconstruct the original  $x$ .

$$L = L(x, \hat{x})$$

$$\text{where } \hat{x} = g(f(x'))$$

This improves robustness to noise.