

# CBCS SCHEME

USN 1 C R 2 L C I 0 0 3

BCA701

## Seventh Semester B.E./B.Tech. Degree Examination, Dec.2025/Jan.2026 Deep Learning

Time: 3 hrs.

Max. Marks: 100

*Note: 1. Answer any FIVE full questions, choosing ONE full question from each module.  
2. M : Marks , L: Bloom's level , C: Course outcomes.*

Module - 1			
Q.1	a.	Define Neural Network. Explain briefly the benefits of neural networks.	M 10 L L1 C CO1
	b.	Explain the role of weights, bias and activation functions in a neuron model.	10 L2 CO1
OR			
Q.2	a.	State and explain Perceptron Convergence Theorem.	10 L2 CO1
	b.	Discuss relation between the perceptron and bayes classifier for a gaussian environment.	10 L2 CO1
Module - 2			
3	a.	What is Multi Layer Perceptron (MLP) and explain how does it differ from a single layer perceptron.	10 L2 CO2
	b.	Explain Back propogation algorithm and its importance in training multilayer perceptrons.	10 L2 CO2
OR			
Q.4	a.	Explain how a multilayer perceptron trained by Back-propogation can solve the XOR problem, where as a single layer perceptron cannot.	10 L2 CO2
	b.	Discuss Heuristics for making the back-propagation algorithm perform better.	10 L2 CO2
Module - 3			
5	a.	What is regularization in the context of deep learning? Explain L2 parameter regularization in detail.	10 L2 CO3
	b.	Discuss the following : i) Dataset Augmentation ii) Semi Supervised Learning.	10 L2 CO3
OR			
	a.	Write a short note on : i) Ill conditioning ii) Plateaus, saddle points and other flat regions.	10 L2 CO3
	b.	Elaborate the challenges <u>local minima</u> , <u>long term dependencies</u> and <u>Cliffs</u> and exploding gradients involved in optimization for training deep model.	10 L2 CO3

Module - 4			
Q.7	a.	Explain convolution neural network in detail.	10 L2 CO4
	b.	Discuss pooling in CNN.	10 L2 CO4
OR			
Q.8	a.	Describe the variants of the basic convolution function in brief.	10 L2 CO4
	b.	With the context of convolutional networks, Explain structured output and data types.	10 L2 CO4
Module - 5			
Q.9	a.	Discuss Recurrent Neural Networks.	10 L6 CO5
	b.	Explain the idea of unfolding computational graph across a deep network structure.	10 L2 CO5
OR			
Q.10	a.	Explain the following: i) Bidirectional RNNs ii) Encoder-Decoder sequence – to – sequence architectures.	10 L2 CO5
	b.	Explain the modes long short-term memory and network based on the gated recurrent unit.	10 L2 CO5

# Deep Learning (BCA 701)

Answerkey - January 2026

S.No  
1a) Neural N/w → 5 Marks with diagram

→ is a massively parallel distributed processor made up of simple processing units that has a natural tendency for storing experiential knowledge and making it available for use.

Benefits - 5 Marks.

- ① Non linearity
- ② Input Output Mapping
- ③ Adeptivity
- ④ Exponential Response
- ⑤ Contented Information
- ⑥ Fault Tolerance
- ⑦ VLSI Implementation
- ⑧ Uniformity of Analysis & Design.
- ⑨ Neurobiology analogy.

1-b) Role of weights, bias and activation function.

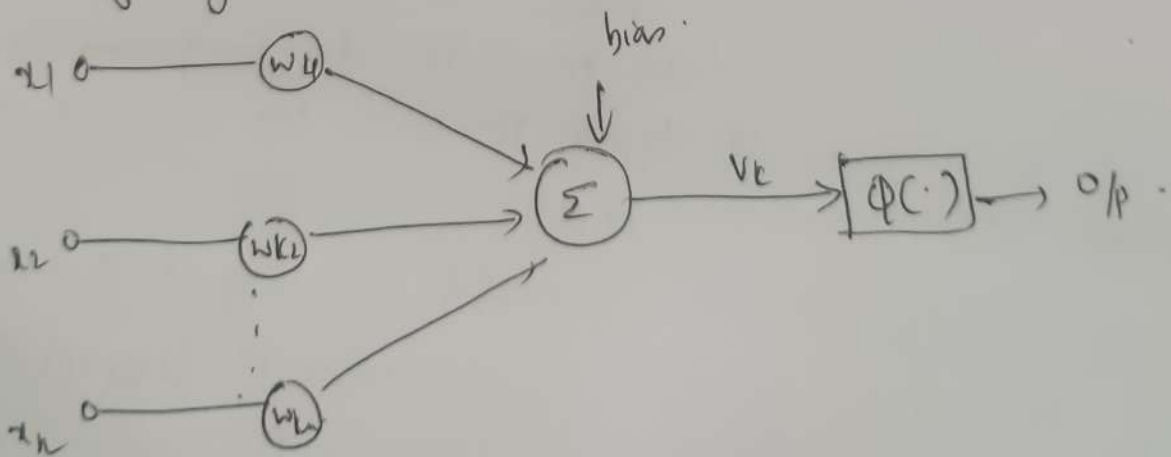


Fig: 4M

Weight + bias + Activation fn - 2 Marks each.

Shot on moto g54 5G

$$u_k = \sum_{j=1}^m w_{kj} z_j$$

$$y_k = \phi(\underbrace{u_k + b_k}_{v_k})$$

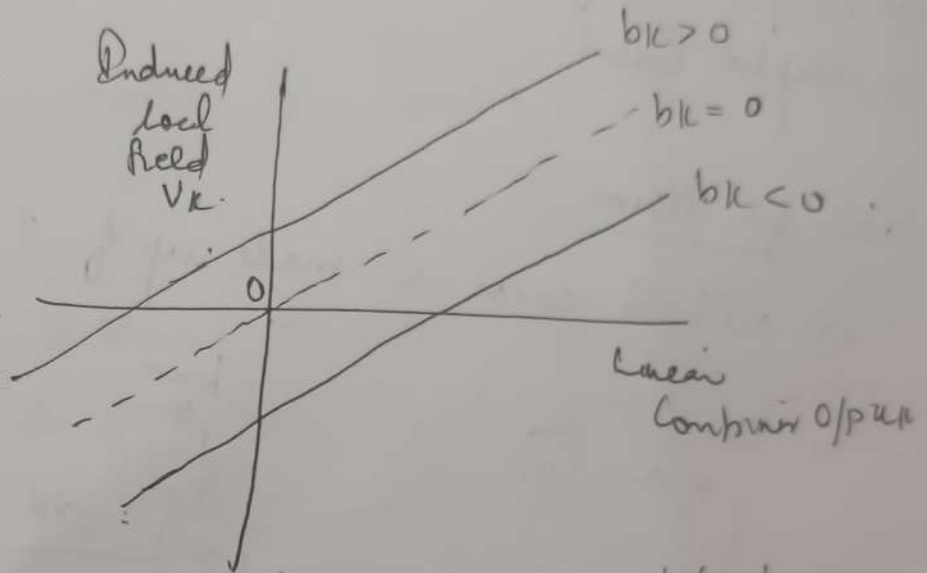
$\phi \rightarrow$  activation fun.

Affine Transformation  
- linear transformation  
+ a shift

$$y = Wx + b$$

$\uparrow$   
more  
shift if  
somehow  
else

Induced  
local  
field  
 $v_k$

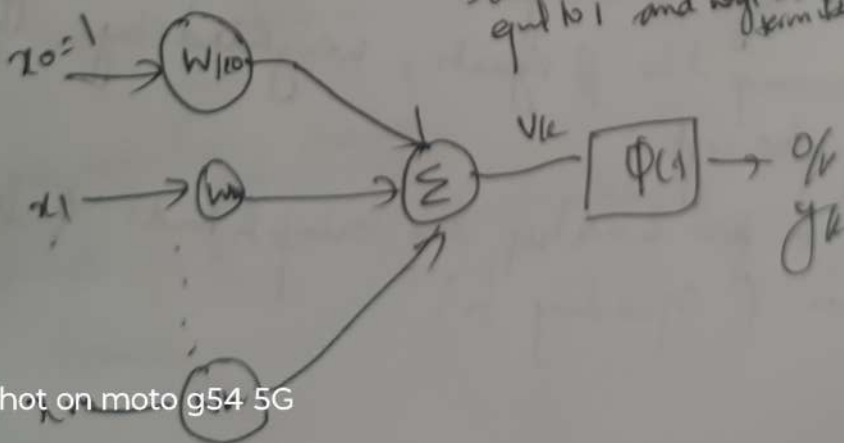


without bias, neurons always activate at the same pt ( $z=0$ )  
Bias lets a neuron shift its activation threshold so  
it can learn patterns.

Bias can be modeled as an extra fixed  $z_0$  with value +1  
( $z_0 = 1$ ) and wgt equal to bias ( $w_{k0} = b_k$ )

(bias can be thought  
as a special  $z_0$   
that is always  
equal to 1 and wgt = bias  
( $w_{k0} = b_k$ ))

Without bias, model learn  
that pass through (0,0)  
with bias, you can  
shift the  $z_0$



$$z = Wx + b$$

$$z = W^1 x^1$$

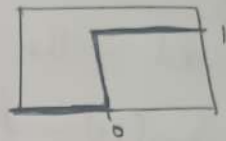
$$x^1 = [z_1, z_2, \dots, z_0, 1]$$

$$W^1 = [w_{11}, w_{12}, \dots, w_{10}, b_1]$$

## Types of Activation fn.

Threshold fn  $\rightarrow$  also called Heaviside fn.

$$\phi(v) = \begin{cases} 1 & \text{if } v \geq 0 \\ 0 & \text{if } v < 0 \end{cases}$$



$\rightarrow$  This model is the McCulloch and Pitts neural model.

$\rightarrow$  all or non property  $v = \sum_j w_j x_j + b_k$

## Sigmoid fn

$\rightarrow$  Commonly used fn. - Smooth curve b/w 0 and 1

$\rightarrow$  increasing fn that balance b/w linear and non-linear.

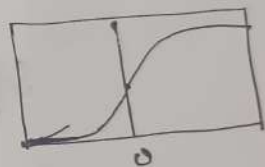
$$\phi(v) = \frac{1}{1 + e^{-av}}$$

$a \rightarrow$  slope parameter.

$a \rightarrow$  scalar.

$v \rightarrow$  4p.

$e \rightarrow 2.718$  o/p.  
 $\sigma(av)$  0 to 1



Signum fn.  $\rightarrow$  which is an odd fn  $\rightarrow$  extracts sign of a real fn.

$$\phi(v) = \begin{cases} 1 & \text{if } v > 0 \\ 0 & \text{if } v = 0 \\ -1 & \text{if } v < 0 \end{cases}$$

## Hyperbolic Tangent fn.

allows for an odd sigmoid type fn.

$$\phi(v) = \tanh(v)$$

range -1 to 1

$$\frac{e^x - e^{-x}}{e^x + e^{-x}}$$

ReLU fn

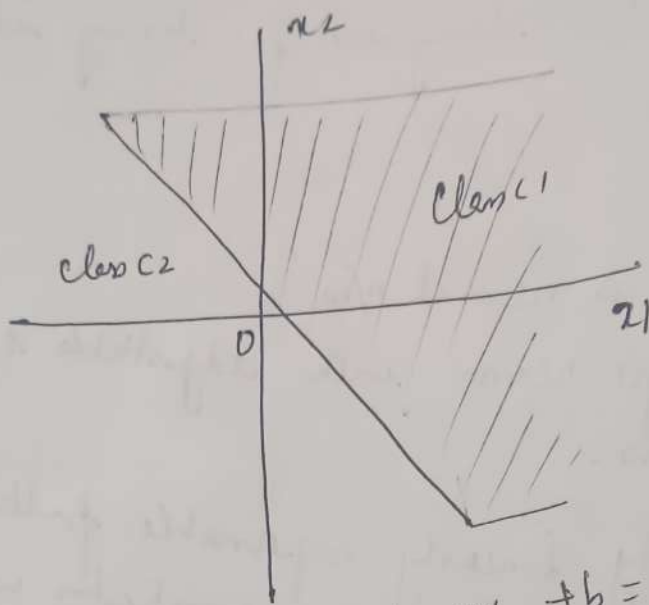
- used in hidden layers

$$\text{ReLU}(x) = \max(0, x)$$

$x > 0 \quad x$

$x \leq 0 \quad 0$

decision regions separation



$w_1 x_1 + w_2 x_2 + b = 0$   
 $(x_1, x_2)$  lies above the boundary line  $\rightarrow C1$   
 $(x_1, x_2)$  " below "  $\rightarrow C2$

### The Perceptron Convergence Theorem

$\rightarrow$  If the data is linearly separable, then the Perceptron algorithm is guaranteed to find a solution in a finite number of steps.

Initial weights and bias

$$w = 0 \quad b = 0$$

Repeat for each training example

Compute activation

$$a = w^T x_i + b$$

Predict label:

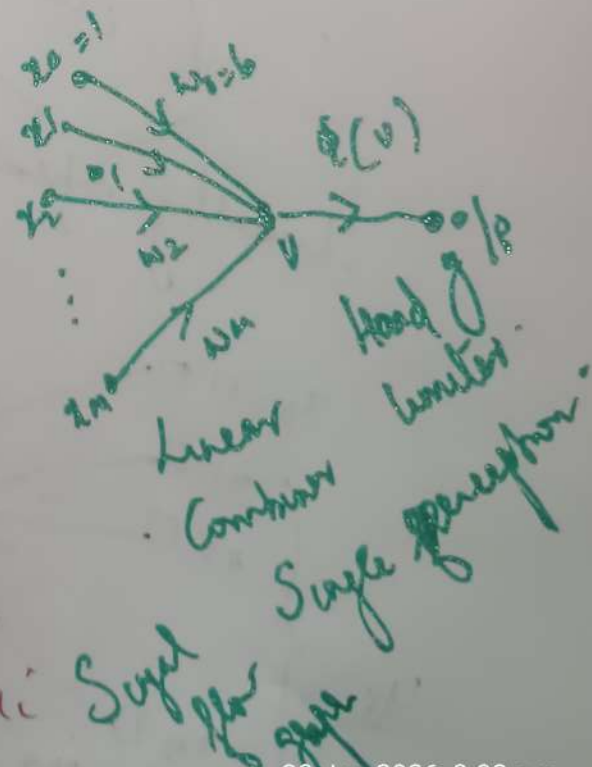
$$\hat{y}_i = \text{sign}(a)$$

if prediction is wrong update

$$w \leftarrow w + \eta \cdot y_i \cdot x_i$$

$$b \leftarrow b + \eta \cdot y_i$$

Repeat until all training examples are classified correctly



## Perceptron Convergence Algorithm

$$x(n) = [1, x_1(n), x_2(n), \dots, x_m(n)]^T$$

$$w(n) = [b, w_1(n), w_2(n), \dots, w_m(n)]^T$$

$b$  = bias

$y(n)$  = actual response

$d(n)$  = desired response

$\eta$  = learning rate

① Initialization: Set  $w(0) = 0$ . Then perform the full computations for time step  $n = 1, 2, \dots$

② Activation: At time step  $n$ , activate the perceptron by applying continuous valued input vector  $x(n)$  and desired response  $d(n)$

③ Computation of Actual Response: Compute the actual response of the perceptron as

$$y(n) = \text{sgn} [w^T(n) x(n)]$$

where  $\text{sgn}(\cdot)$  is the signum fn.

④ Adaptation of Weight Vector

$$w(n+1) = w(n) + \eta [d(n) - y(n)] x(n)$$

$$d(n) = \begin{cases} +1 & \text{if } x(n) \text{ belongs to class } C_1 \\ -1 & \text{if } x(n) \text{ belongs to class } C_2 \end{cases}$$



# Relation b/w the Perceptron and Bayes classifier for a gaussian Environment

## Bayes classification

→ is a probabilistic method used to decide which class an input belongs to. It is based on Bayes theorem which calculates the prob. of a class given some observed data.

email → spam/not spam words - higher probability.

## Gaussian Environment

- ~~each~~ data for each class follows a bell shaped curve.
- normal distribution
- height of students



## Perceptron

- model that tries to find a line to separate two classes.
- doesn't use probabilities.
- adjust weights based on whether it got the prediction right or wrong.

If the data follows Gaussian distribution, the Bayes classifier decision boundary is a linear function just like the Perceptron.

### ① Likelihood Ratio ( $\lambda(x)$ )

class 1 and class 2

$\lambda(x)$  is high  $x \rightarrow$  class 1  
 $\lambda(x)$  is low  $x \rightarrow$  class 2

looks more likely to

looks more likely to

Shot on moto g54 5G

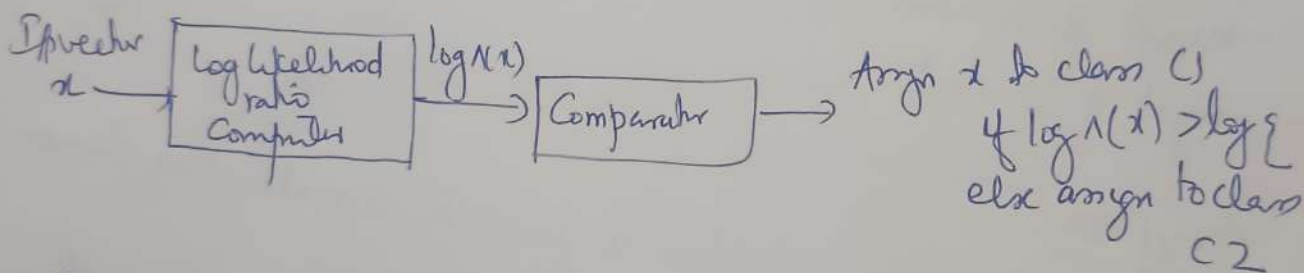
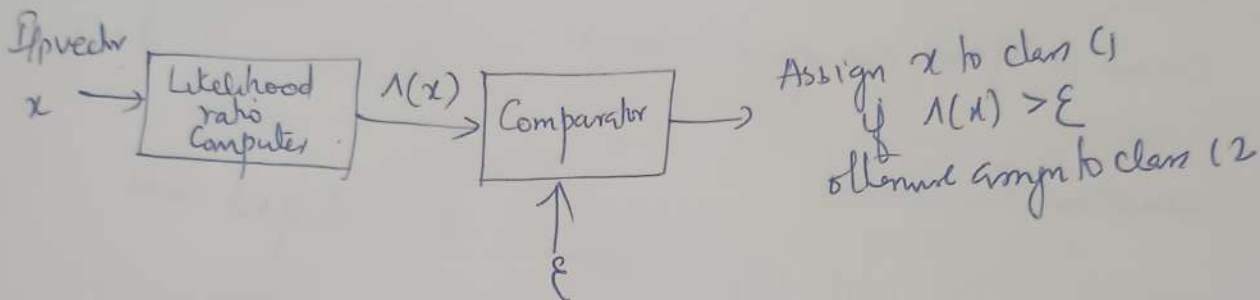
## Threshold $\epsilon$

- cutoff point for deciding between the class

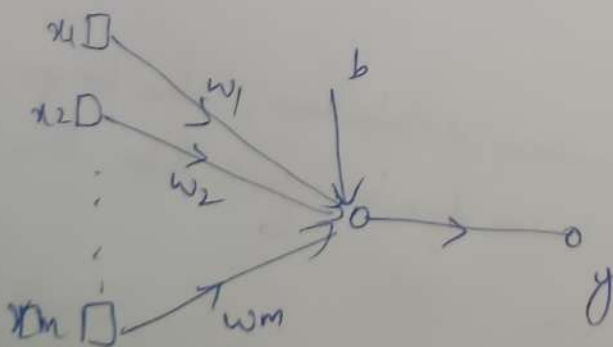
$\Lambda(x) > \epsilon$  assign  $x$  to class 1 ( $C_1$ )

$\Lambda(x) \leq \epsilon$  assign  $x$  to class 2 ( $C_2$ )

likelihood ratio Test



## Signal flow graph of Gaussian classifier



$$y = w^T x + b$$

$$y = \log \left( \frac{\lambda_1}{\lambda_2} \right)$$

log of likelihood ratio  
o/p can be based on the ratio of  
prob. of 2 classes.

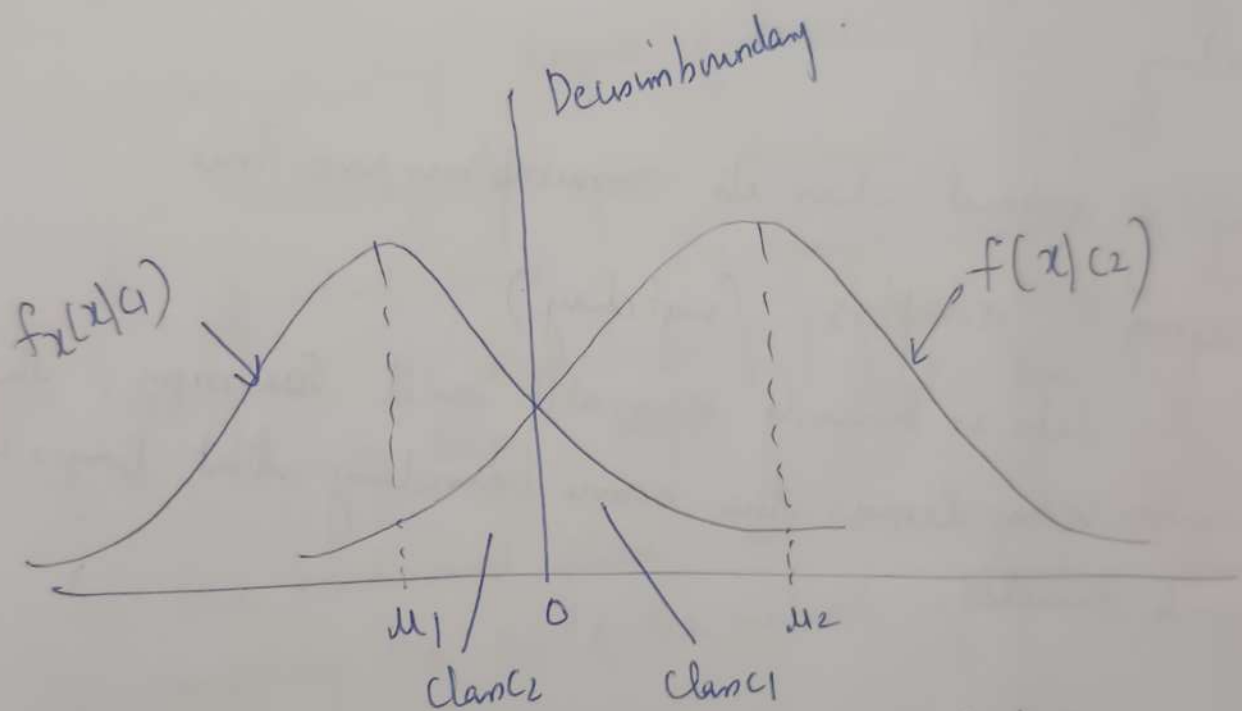


$$w = C^{-1} (\mu_1 - \mu_2)$$

diff b/w avg. means of 2 classes adjusted by C (Covariance)

$$b = \frac{1}{2} (\mu_2^T C^{-1} \mu_2 - \mu_1^T C^{-1} \mu_1)$$

bias values, which helps shift the decision boundary



Two overlapping gaussian distributions.

horizontal axis  $\rightarrow$  4p feature (height, weight)

vertical axis  $\rightarrow$  prob. density

2 curves overlap  $\rightarrow$  some values could belong to either class.

decision boundary  $\rightarrow$  at  $x=0$  where classifier switches from predicting class 1 to class 2.



## Backpropagation algorithm

→ core method for training MLP

→ 2 phases — Forward Pass (Signal Propagation)

— Backward Pass (Error propagation & wgt adjustment)

## → Initialisation

— all wghts and threshold (bias) are initialised to small random values.

## → 2 Presentation of training examples

Training data is given in epochs.

N/w performs — Forward Computation

— Backward "

## → Forward Computation

Input vector :  $x(n)$  is fed into the i/p layer.

Desired o/p :  $d(n)$  is provided at the o/p layer.

Last field of a neuron  $j$  in layer  $l$ .

$$v_j^{(l)}(n) = \sum_i w_{ji}^{(l)}(n) y_i^{(l-1)}(n)$$

O/p of neuron  $j$

$$y_j^{(l)}(n) = \phi_j(v_j^{(l)}(n))$$

## → Error Computation

$$e_j(n) = d_j(n) - o_j(n)$$

## → 5 Backward Computation

- Error is propagated backward to adjust weight  
Local gradient  $\delta$

$$\delta_j^{(l)}(n) = \begin{cases} e_j^{(l)}(n) \phi_j'(v_j^{(l)}(n)) & \rightarrow \text{for neuron } j \text{ in op layer } l \\ \phi_j'(v_j^{(l)}(n)) \sum_k \delta_k^{(l+1)}(n) w_{kj}^{(l+1)}(n) \end{cases}$$

## → 6. Weight Update Rule . generalized delta rule

$$w_{ji}^{(l)}(n+1) = w_{ji}^{(l)}(n) + \alpha [w_{ji}^{(l)}(n) - w_{ji}^{(l)}(n-1)] + \eta \delta_j^{(l)}(n) y_i^{(l-1)}(n)$$

## → 7 Iteration

Steps 3-6 are repeated until stopping criterion is met



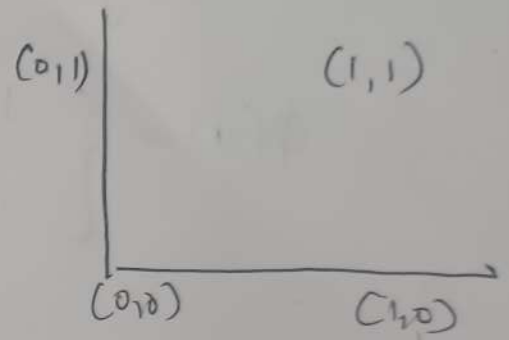
## XOR problem

→ In Rosenblatt's single layer perceptron, there are no hidden neurons.

→ it cannot classify 4p patterns that are not linearly separable.

XOR

$x_1$	$x_2$	$y$
0	0	0
0	1	1
1	0	1
1	1	0

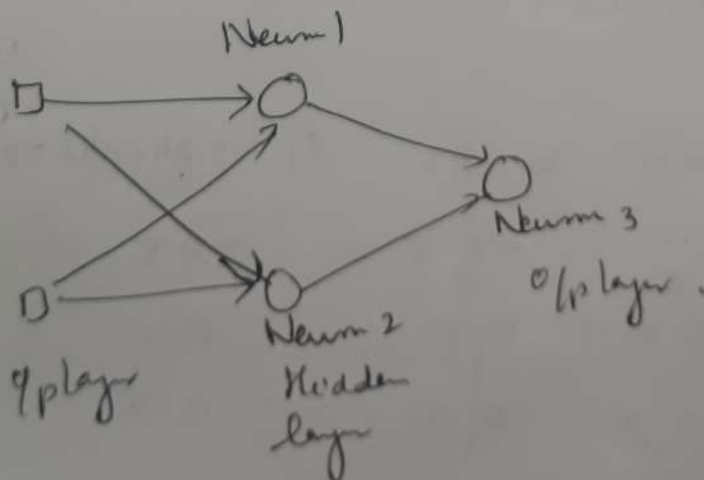


→ XOR not linearly separable.

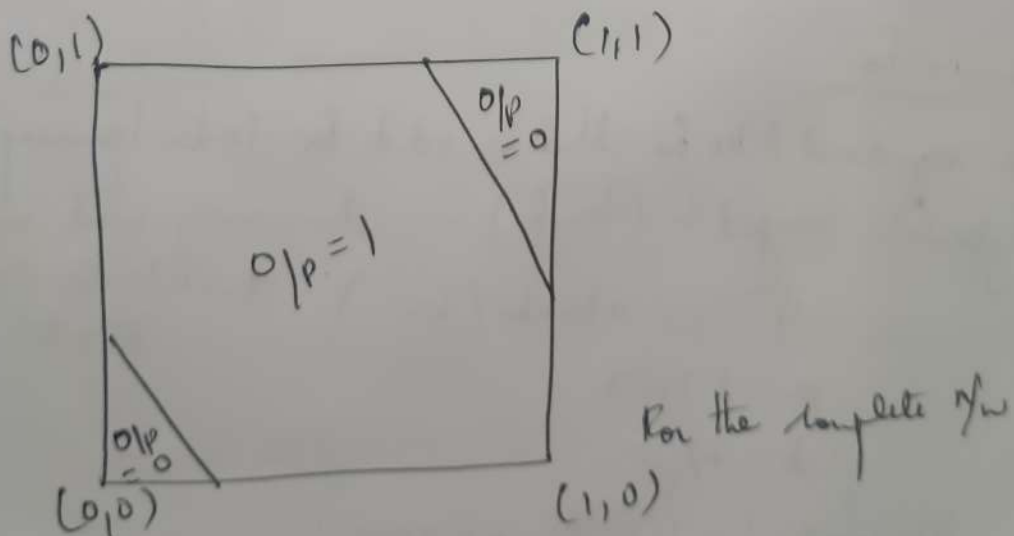
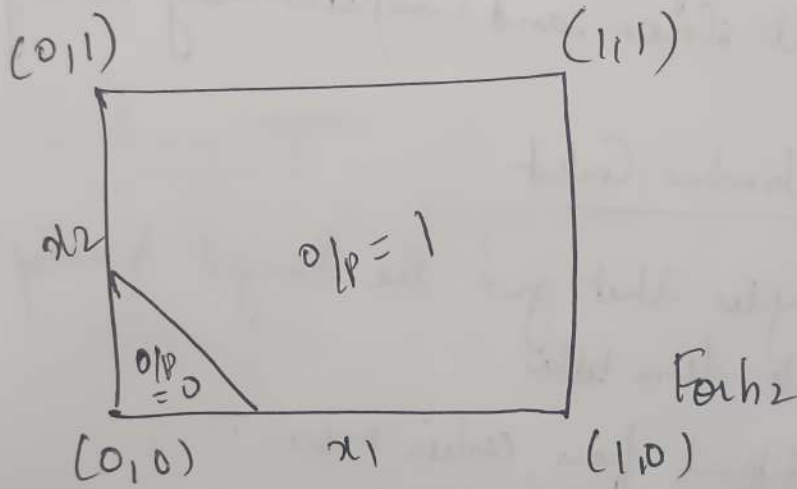
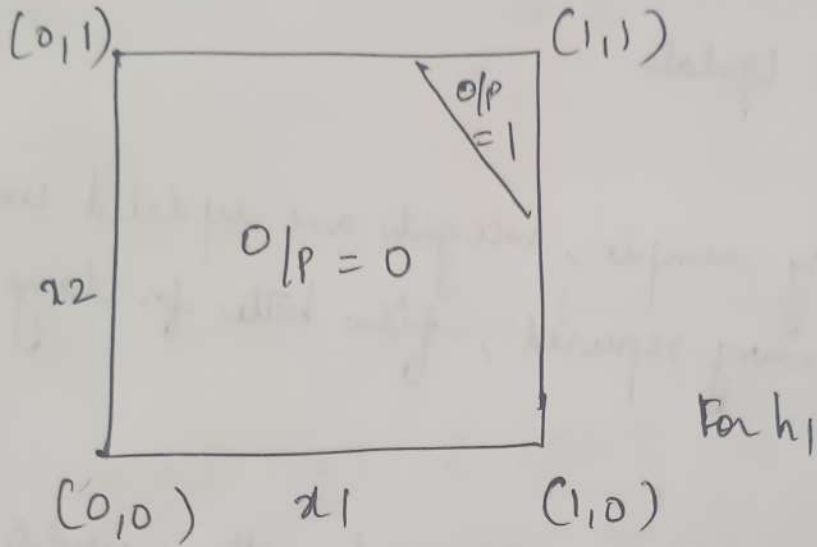
→ We cannot draw a single straight line to separate the 4p that yield different o/p.

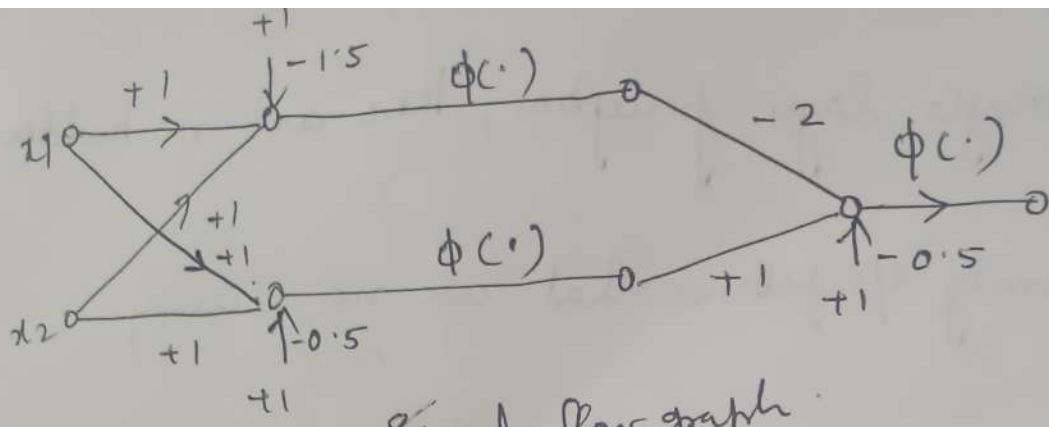
→ A single layer perceptron due to its linear nature, fails to model the XOR fn.

→ XOR problem can be solved using a single hidden layer with 2 neurons.



Decision boundary





Signal flow graph.

$$\phi(v) = \begin{cases} 1 & v \geq 0 \\ 0 & v < 0 \end{cases}$$

$$s_1 = x_1(+1) + x_2(+1) + 1(-1.5)$$

$$= x_1 + x_2 - 1.5$$

$$h_1 = \phi(s_1)$$

$$s_2 = x_1(+1) + x_2(+1) + 1(-0.5)$$

$$= x_1 + x_2 - 0.5$$

$$h_2 = \phi(s_2)$$

$$s_3 = h_1(-2) + h_2(+1) + 1(-0.5) = -2h_1 + h_2 - 0.5$$

$$y = \phi(s_3)$$

$x_1$	$x_2$	$s_1 \rightarrow h_1$	$s_2 \rightarrow h_2$	$s_3 = -2h_1 + h_2 - 0.5$	$y$
0	0	-1.5	-0.5	-0.5	0
0	1	-0.5	0.5	0.5	1
1	0	-0.5	0.5	0.5	1
1	1	0.5	1.5	-1.5	0



# Heuristics for making the Backpropagation Algorithm

## ① Stochastic vs Batch Update

### Stochastic

After each training sample, weights are updated immediately  
→ Faster, less memory required, often better for large datasets

### Batch Update

Wait until all samples are processed, then update weights  
More stable but slower and computationally heavy

## ② Maximizing Informative Content

Use training samples that give the largest training error  
forces the n/w to learn better

Use examples different from earlier ones

Datasets are shuffled to keep variety.

## ③ Activation fn

Use sigmoid like fn that is odd for faster learning

Hyperbolic tangent fn (tanh) - commonly used in neural n/w

$$\phi(v) = a \tanh(bv) \quad \text{due to its symmetry and bounded o/p}$$

$$a = 1.7159$$

$$b = 2/3$$

$$\text{o/p range } [-1.7159, +1.7159]$$

At  $v=0$  slope is 1

$$\phi'(v) = \phi(v)$$

Symmetric which helps faster

#### 4. Target Values

- Desired o/p should not be exactly at the limits of the activation fn ( $+a$  or  $-a$ )
    - enters a saturated state where gradients become small - slow learning
  - $a - \epsilon$  or  $-a + \epsilon$
  - This prevents neurons from saturating.
    - when the o/p of the activation fn becomes close to its max/min gradient becomes very small or zero
  - $a = 1.7159$   $\epsilon = .7159$
- Target values can then be conveniently chosen as  $+1$  and  $-1$ .

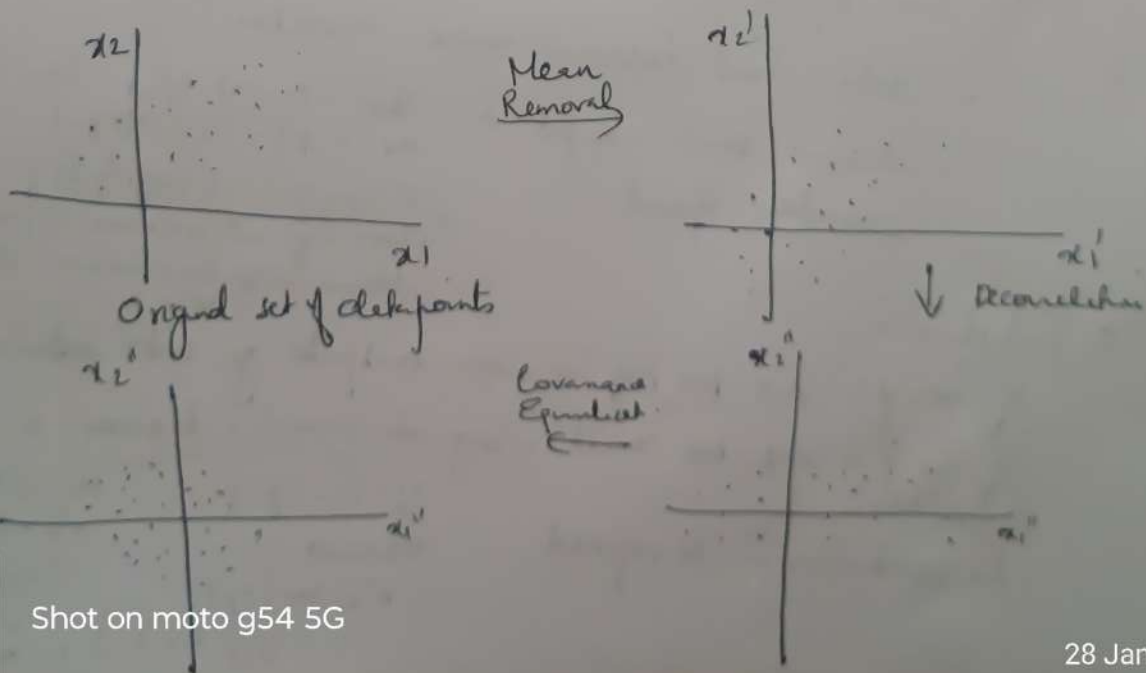
#### 5. Normalizing Inputs

Inputs should be preprocessed

$$\text{Mean} = 0$$

$$\text{S.D} = 1$$

If all inputs are large and always +ve, wghts may only grow in one direction - slow down learning  
Normalization assures balanced & faster emergence.



1. Original set of data points

Input features  $x_1, x_2$  are scattered

They may have different scales, means and correlations.

## 2. Mean Removal

Each feature's mean is shifted to zero.  $x'_i = x_i - \mu_i$

after this the data cloud is centered around the origin  $(0,0)$

But variables may still be correlated

## 3. Decorrelation

Using PCA, features are transformed

The cloud looks like an ellipse aligned with the coordinate axes.

## 4. Covariance Equalization

- scale the data so that variances along each axis are equal

- data cloud becomes more circular

- ensures each weight in the neural net learns at a similar speed

$$\rightarrow E[w_j y_i] = E[w_j] \cdot E[y_i]$$

Choose wgt's with zero mean  $E[w_j] = 0$

Normalized to have zero mean  $E[y_i] = 0$

$$\mu_V = 0$$

## 6. Initialization

If weights are too big neuron outputs go into saturation

If weights are too small - signals vanish, learning is very slow

So initialization is required.

$$V_j = \sum_{i=1}^m w_{ji} y_i$$

Variance,

$$\sigma^2 V = \text{Var}(V_j)$$

$\mu_V = 0$  - pre-activation mean is zero

$\sigma^2 V = E[V_j^2]$

$\sigma^2 V = m \cdot \text{variance}$

## 7. Learning from hints

- giving network extra guidance - certain properties of the target or  $f(\cdot)$
- Instead learning only from raw training data, the n/w can also learn from prior knowledge
- Symmetries, invariances, constraints
- Rules that o/p must satisfy.
- $f_n$  remain the same when  $o/p$  are rearranged.  $o/p$  of the  $f_n$  don't change under certain transformations.
- Translation Invariance - object in a image stays same if shifted.
- Rotation Invariance
- Scale Invariance
- If  $f_n$  is symmetric, we can design the n/w to respect that.
- If  $o/p$  are rotation invariant, we can add hints.
- Learning from hint  $\rightarrow$  <sup>not</sup> only learn from example but using clues about the target or  $f(\cdot)$

## 8. Learning rate

controls how fast weights are updated.

$\eta$  is too large - oscillation, divergence.

$\eta$  is too slow - very slow learning.

Neurons closer to the o/p need smaller learning rates.  
" " input need larger " "



Regularization in DL is a set of techniques used to reduce overfitting by adding constraints or penalties to the model's parameters or by modifying the training process, so that the model generalizes better to unseen data.

## Types of Regularization in DL

### (1) L1 Regularization (Lasso)

- It adds absolute value of weights to the loss function.
- Encourages weights to become sparse (few features are kept, rest are dropped)
- Model becomes more interpretable.

### (2) L2 Regularization (Ridge)

- adds squared value of weights to the loss.
- encourages weights to be small but not exactly zero.
- Model becomes balanced & stable.

### (3) Dropout

- Randomly turns off some neurons during training.
- Prevents neurons from becoming too dependent on each other.

### (4) Data Augmentation

- Modify training images / text / audio.
- Give model more variety without collecting new data.
- Makes model more robust.

## Early Stopping

- stops training when validation error starts going up
- prevents overfitting
- model don't overheat

Suppose we have 3 points

x	y
1	3
2	5
3	7

$$y = wx + b$$

Without Regularization

$$w = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$b = \bar{y} - w\bar{x}$$

$$\bar{x} = 2 \quad \bar{y} = 5$$

$$w = 4/2 = 2$$

$$b = 1$$

$$y = 2x + 1$$

Fits perfectly

With L2 Regularization

$$L = \frac{1}{n} \sum (y_i - (wx_i + b))^2 + \lambda w^2$$

$$\lambda = 1$$
$$w = \frac{4}{2+1} = \frac{4}{3} \approx 1.333$$

$$b = \bar{y} - w\bar{x}$$
$$= 2.334$$

With regularization  $y = 1.333x + 2.334$

weight shrinks, avoids overfitting.  
Perform well with the unseen data.



Regularization  $\rightarrow$  - A central goal is to find an algo. that will perform well not just on the training data but also on new inputs.

- strategies used in ML.

- is defined as any modification we make to a learning algo. that is intended to reduce its generalization error but not its training error.

### Parameter Norm Penalty

$\downarrow$   
Weights

$\downarrow$   
size or magnitude

$\rightarrow$  extra cost

$$\text{Loss with Penalty} = \text{Original Loss} + \lambda \cdot \text{Norm of weights}$$

L1 or L2 measures how big weights are.

L1 norm  $|w_1| + |w_2| + \dots$

L2 norm  $w_1^2 + w_2^2 + \dots$

Regularization approaches are based on limiting the capacity of the models by adding norm penalty  $\Omega(\theta)$  to the objective fn.

It defines what model is being

$$\tilde{J}(\theta; x, y) = \underbrace{J(\theta; x, y)}_{\text{Original obj. fn.}} + \alpha \underbrace{\Omega(\theta)}_{\text{Penalty Term}}$$

$\alpha \in [0, \infty)$  is a hyperparameter that weights the relative contribution of the norm penalty term relative to std. obj. fn.

Penalty is also known as weight decay.

When our training algo minimizes the regularized obj fn  $\hat{J}$  will decrease both original objective  $J$  on the training data and some measure of the size of the parameters  $\theta$ .

## L<sup>2</sup> Parameter Regularization

L<sup>2</sup> parameter norm penalty  $\rightarrow$  weight decay.

Regularization term  $\rightarrow \Omega(\theta) = \frac{1}{2} \|W\|_2^2$  to the obj. fn.

Ridge Regression or Tikhonov regularization.

$\rightarrow \Omega(\theta) \rightarrow$  regularization penalty.

$\theta \rightarrow$  refers to the weights.

$$\|W\|_2^2$$

$$\|W\|_2 = \sqrt{w_1^2 + w_2^2 + \dots + w_n^2}$$

$$\|W\|_2^2 = \sum_{i=1}^n w_i^2$$

Sum of squares of all weights.

$\frac{1}{2} \rightarrow$  added for mathematical convenience.  
 $\frac{1}{2} w^2 \quad \frac{1}{2} \times 2w = w$

Weights are kept as small.

Simple model.

Better generalization.

Reduced overfitting.

# Dataset Augmentation

- improve a ML model - is by training it on more data
- To overcome ~~the~~ limited data, you can create fake (data) and add it to your training set.
- A classifier maps a complex high dimensional input  $x$  to a single category label  $y$ . The main challenge is that the classifier is invariant to many possible transformations of the input.
- By transforming the  $x$  data  $x \rightarrow x'$ , you can easily generate new  $(x', y)$  pairs for training.

## Limitations

- approach doesn't extend to all tasks.
- generating fake data for density estimation is hard.
- Data augmentation is especially useful for object recognition.
  - Translating images
  - Rotating
  - Scaling
- Don't apply transformations that change the correct label.
- Flipping or rotating letters could change 'b' to 'd' or '6' to '9' which alters the class label.
- Shear and 180° rotations are not appropriate.

cat → Rotate cat  
Image Flip cat  
Mirror cat  
Zoom cat

- Injecting noise into inputs of neural n/w can be seen as a form of data augmentation.   
 injecting randomness
- Neural n/w trained with noisy  $y_p$  improves robustness.   
 model able to see corrupted data during training
- Noise injection - denoising autoencoders.
- Injecting noise into hidden units can be viewed as data augmentation.
- augmentation can significantly reduce generalization error.
- application specific operations are treated as preprocessing step.

Original pixel = 120  
 Input gaussian noise mean 0  $\sigma = 10$   
 $120 + N(0, 10)$   
 $= 127$

## semisupervised Learning

Uses labeled and unlabeled data to improve model performance.

- We have unlabeled data from the distribution  $P(x)$  tells what  $y_p$  looks like.
- We have labelled data from the distribution  $P(x, y)$  how  $y_p$  relates to o/p.
- Estimate  $(y|x)$  directly predict  $y$  from  $x$ .   
 Conditional distribution.
- The goal is to learn a representation  $h = f(x)$  for the  $y_p$  data  $x$ .
- The idea is that samples belonging to the same class would have similar representations in the feature space.

# Challenges in Optimization

- Optimization is an extremely difficult task.
- Traditional ML: careful design of objective function and constraints to ensure convex optimization.

## ① Local Minima

In convex optimization, problem is one of finding a local minimum.

Some convex fn have a flat region rather than a global minimum pt.

Any pt within the flat region is acceptable.

With non convexity of neural nets many local minima are possible.

Many deep models are guaranteed to have an extremely large no. of local minima.

This is not necessarily a major problem.

## ② Ill conditioning of the Hessian.

The Hessian matrix  $H = \nabla^2 f(x)$  captures the curvature of the fn.

The Hessian of a scalar fn  $f(x)$  is the matrix of second partial derivatives.

$$H = \nabla^2 f(x) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \dots \\ \frac{\partial^2 f}{\partial x_1 \partial x_2} & \frac{\partial^2 f}{\partial x_2^2} & \dots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

Condition number of Hessian

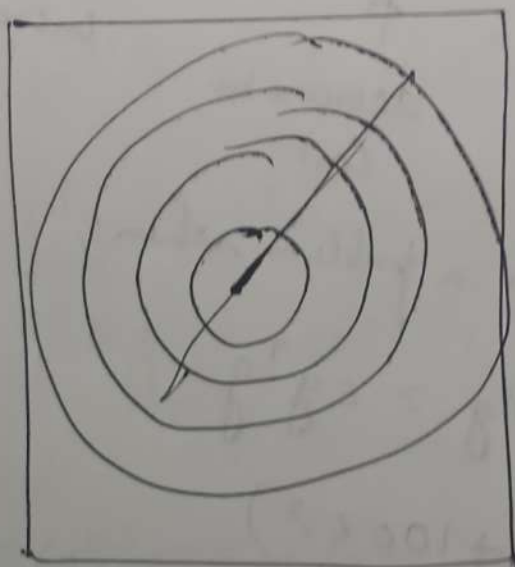
$$k = \frac{\lambda_{\max}}{\lambda_{\min}}$$

$\lambda_{\max}$  : largest eigen value of H.

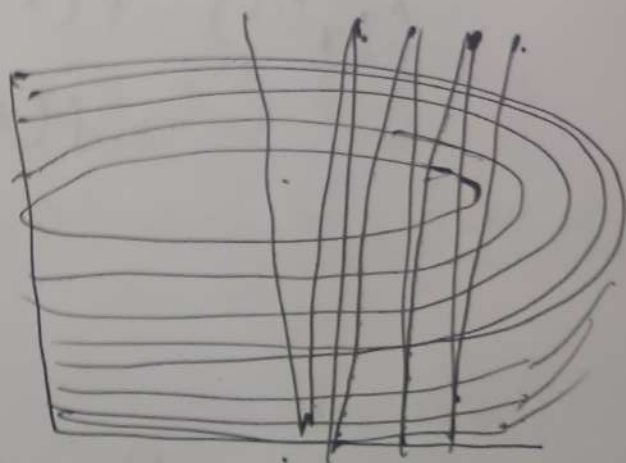
$\lambda_{\min}$  : smallest eigen value of H.

If  $k$  is large, the Hessian is ill conditioned.

This means the  $f(x)$ 's contours are very elongated (like thin ellipses) so gradient descent takes zig zag paths and converges very slowly.



Well conditioned case.  
Contours are round.  
Eigenvalues are equal.



Well conditioned - round contours - fast convergence  
ill conditioned - stretched " - slow "

Scaling cancels out

## Plateaus, Saddle Points and other flat Regions

A saddle point is a critical point where the gradient is zero but the Hessian has both positive and negative eigen values.

$$f(x, y) = x^2 - y^2$$

$$\text{Gradient } \nabla f = (2x, -2y)$$

At  $(0, 0)$  gradient is 0.

$$H = \begin{pmatrix} 2 & 0 \\ 0 & -2 \end{pmatrix}$$

Eigen values  $+2, -2 \rightarrow$  mixed signs saddle.

Along the  $x$  direction the function curves up along the  $y$  direction it curves down.

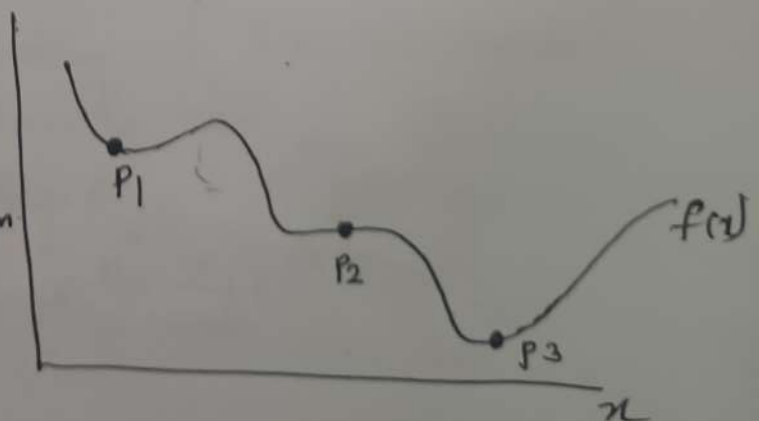
So the point is neither a pure min nor max.

In high dimensions saddle points become very common and can dominate training behaviour.

$P_1 \rightarrow$  local minimum

$P_2$ : Saddle point

$P_3$ : Global minimum



## Plateau

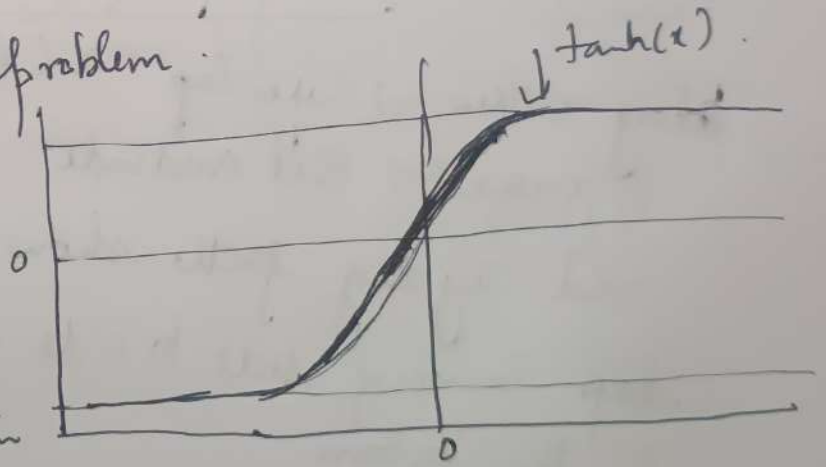
- is a region of the loss surface where the  $f_n$  is almost flat.
- Gradient = 0
- Loss doesn't change much no matter which direction you move.
- gradients are tiny, updates to weights are super small.
- optimizer gets stuck.
- for  $f_n$  like tanh or sigmoid, when  $x$  are very large or small derivatives  $\approx 0$ .
- Vanishing gradients problem.

$y = +1$   
 $y = -1$   
tanh saturates

$x > 2 \rightarrow$  positive plateau

$x < -2 \rightarrow$  negative plateau

curve is almost flat — gradients vanish



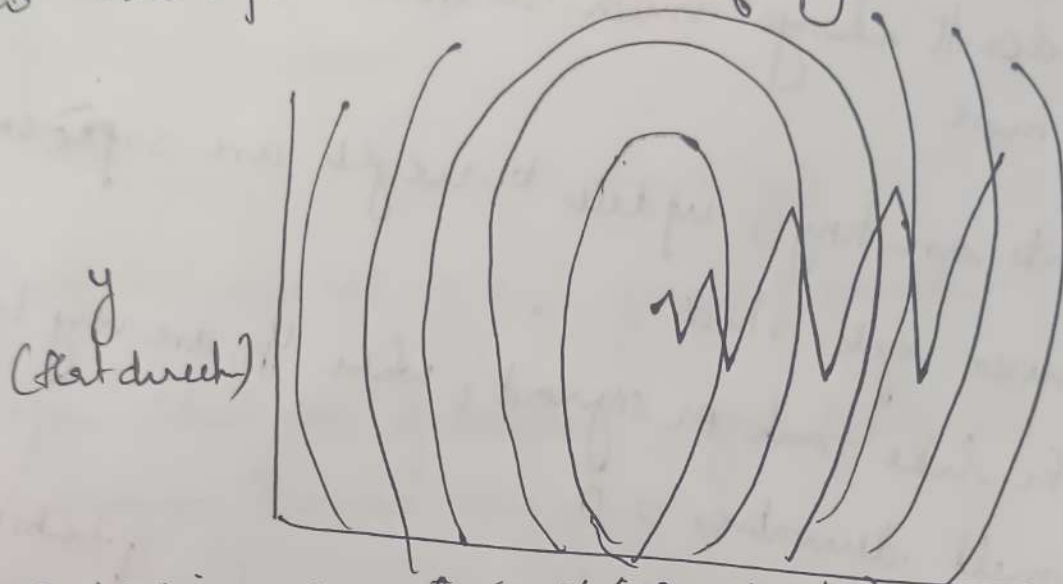
## Ravine

- is a very narrow, steep valley in the loss surface.
- It happens when the loss  $f_n$  has one direction with very steep slope.
- another direction with very gentle slope.

## Zig-zag path in gradient descent

Gradient descends steeply down the rainie walls instead of smoothly along the valley floor.

So the optimizer bounces left right ~~left~~ left.



Along x axis  $\rightarrow$  super steep  $x$  (steep direction)  
y axis  $\rightarrow$  flat and wide.

red zigzag path shows how gradient descent keeps bouncing side to side instead of going straight to the bottom.

## Cliffs and exploding gradients

- $\rightarrow$  In deep  $W$  multiplying many large weights creates steep regions in the loss surface like a cliff edge.
- $\rightarrow$  when gradient descent takes a step it can fall off the cliff.
- $\rightarrow$  Update is too large
- $\rightarrow$  optimizer skips to minima

## Long Term dependencies

- In deep computational graphs, gradients are multiplied repeatedly.
- Vanishing gradients → signal shrinks, early layers can't learn.
- Exploding gradients - signal grows uncontrollably.
- In RNN applying the same weights at each time step makes it worse - n/w struggles to remember information far in the past.
- ⇒ n/w becomes good at remembering recent information but struggles with long term dependencies → RNN forgets the thing happened many steps ago.  
LSTM, GRU



## Convolution N/w

- are also known as convoluted neural n/w or CNN
- CNN are a specialised kind of neural n/w for processing data that has a known, grid like topology
- 1D grid time series data
- 2D grid - image data (2D grid of pixels)
- network employs a mathematical operation convolution
- special kind of linear fashion
- Convolutional n/w are simply neural n/w that use convolve in place of general matrix multiplication in atleast one of their layers
- Convolve is an operation on 2 hrs of a real valued argument.

$$S(t) = (x * W)(t)$$

↑            ↑    ↑  
O/p → feature map    4p    kernel



In machine learning,

Input = a tensor (multidimensional array of nos.  
ie image pixels)

Kernel = a tensor of learnable parameters

O/p = another tensor, a smaller or same size  
depending on padding/stride.

For 2D Convolution

$$S(i, j) = (I * K)(i, j) = \sum_m \sum_n I(m, n) * K(i-m, j-n)$$

$$S(i, j) = \sum_m \sum_n I(i-m, j-n) k(m, n)$$

$S(i, j)$  = o/p pixel at row  $i$ , col  $j$   
 $I(m, n)$  pixel value at row  $m$  col  $n$

$K(i-m, j-n)$  - kernel value  
↓

reverse (flip) the kernel before sliding it  
across the image

CNN They  
slip the filter  
just do the  
same

Size of the o/p feature map is given by,

$$\left\lfloor \frac{H - kn + 2p}{s} \right\rfloor + 1$$

I/p size =  $32 \times 32$

Kernel =  $5 \times 5$

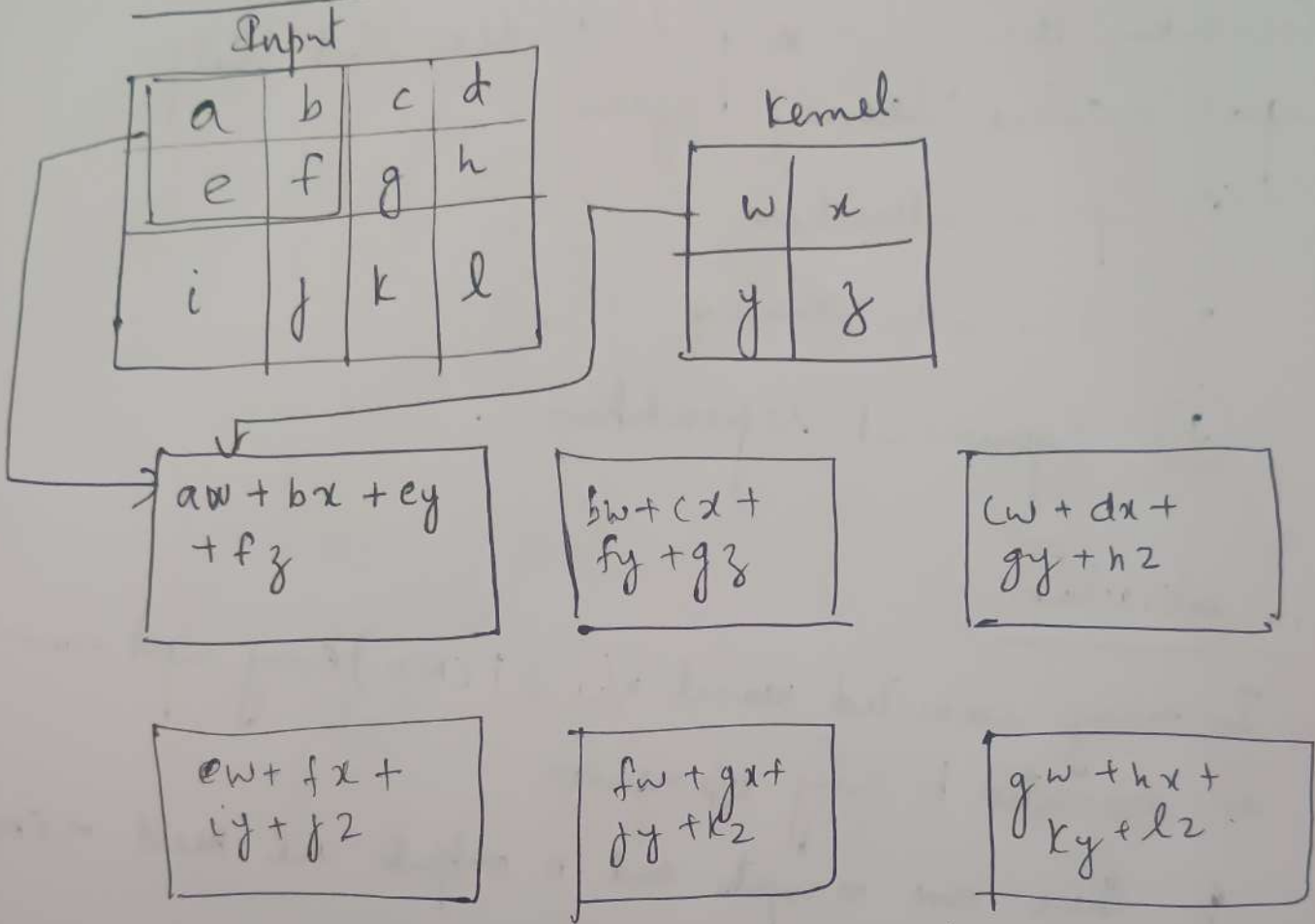
Stride = 1

$$\text{feature map} = \left\lfloor \frac{32 - 5 + 0}{1} \right\rfloor + 1$$

28 Jan 2026, 3:06 pm



# The Convolution Operation



Convolve image with a filter (kernel) matrix

3	0	1	2	7	4
1	5	8	9	3	1
2	7	2	5	1	3
0	1	3	1	7	8
4	2	1	6	2	8
2	4	5	2	3	9

6x6

x

1	0	-1
1	0	-1
1	0	-1

3x3

=

-5	-4	0	8
-10	-2	2	3
0	-2	-4	-7
-3	-2	-3	-16

4x4

$$6 - 3 + 1$$



## Pooling

→ summarizes nearby o/p's to make representation smaller and more stable to small translation.

→ Reduce Size (Dimensionality Reduction)

## Max Pooling

- keeps the max. in each window.

Feature maps from convolution are often big.

Pooling makes them small faster training, fewer parameters, less memory.

## Avg Pooling

- keeps the avg.

$$\begin{bmatrix} 1 & 1 & 2 & 4 \\ 5 & 6 & 7 & 8 \\ 3 & 2 & 10 & \\ 1 & 2 & 3 & 4 \end{bmatrix}$$

2x2 Max Pooling with stride 2

$$\{1, 1, 5, 6\} \quad \text{max} = 6$$

$$\{2, 4, 7, 8\} \quad \text{max} = 8$$

$$\{3, 2, 1, 2\} \quad \text{max} = 3$$

$$\{1, 0, 3, 4\} \quad \text{max} = 4$$

$$\begin{bmatrix} 6 & 8 \\ 3 & 4 \end{bmatrix}$$

Avg Pooling

## Pooling

→ A typical layer of Convolutional n/w consists of 3 stages

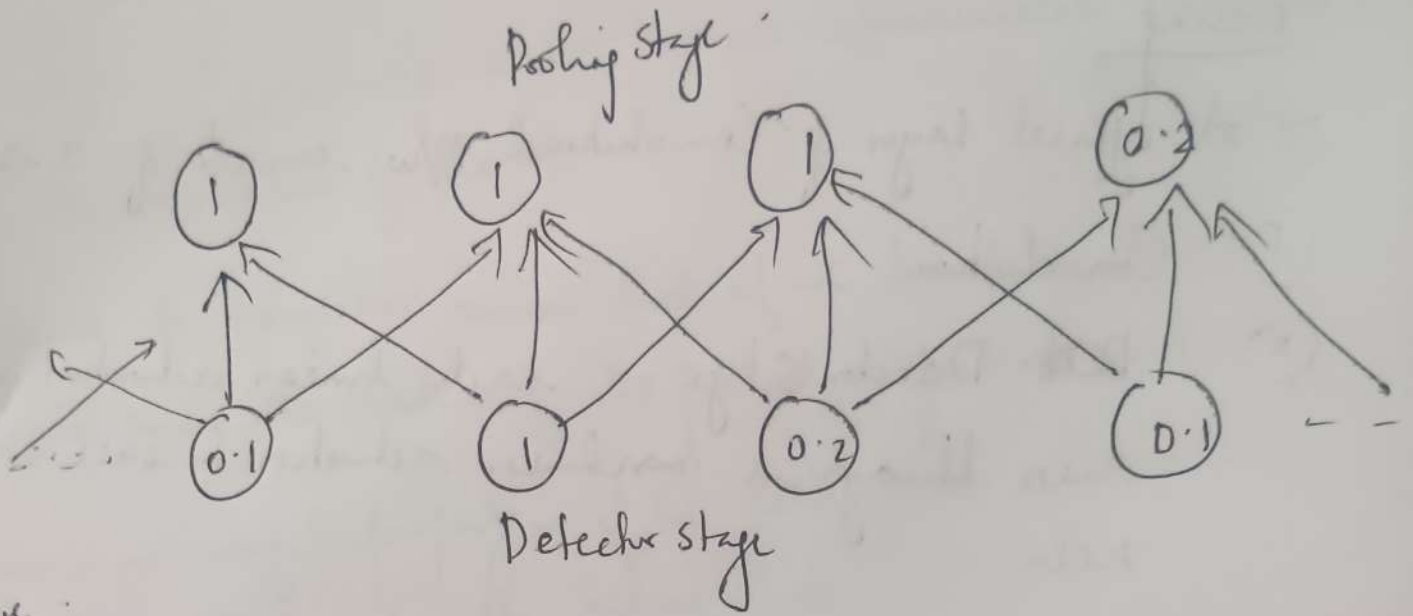
- ① Convolution
- ② ~~ReLU~~ Detector Stage → each linear activation is run through a non-linear activation fn such as relu.
- ③ Pooling function to modify the o/p of the layer.
- ④ A pooling fn replaces the o/p of the net at a certain location with a summary statistic of the nearby o/p (max pooling or avg pooling).

A pooling independently operates on each feature map to produce another Feature map.

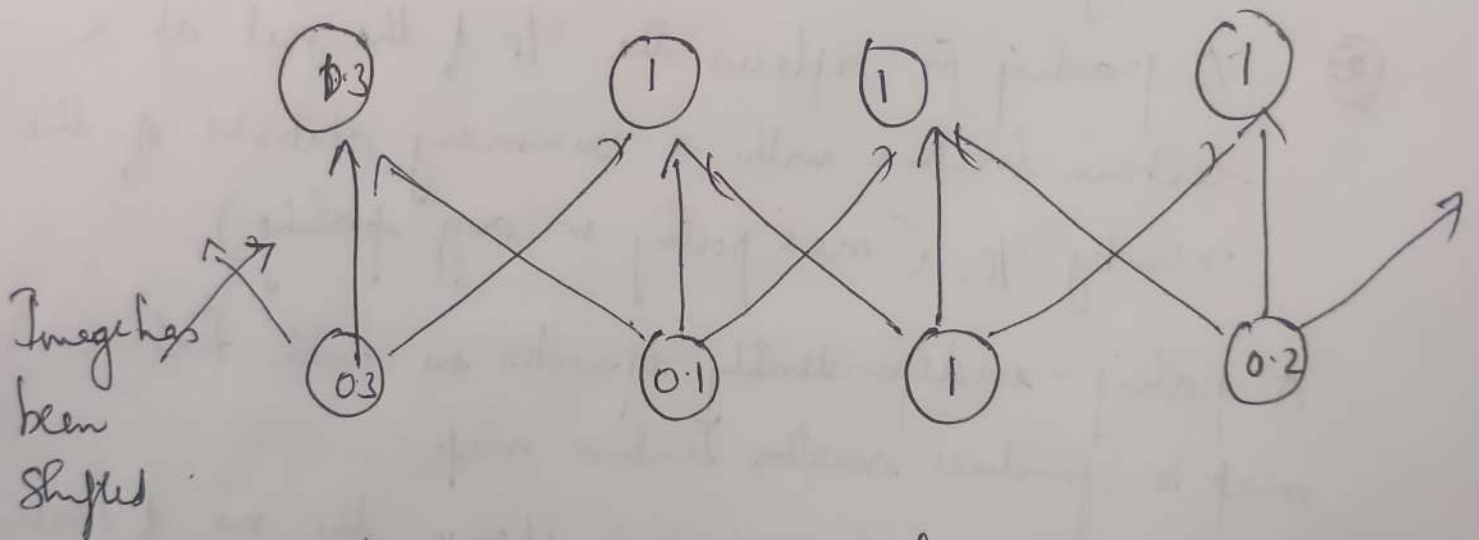
Operation of pooling doesn't change the no. of feature maps.

An earliest convolutional n/w LeNet-5 → avg pooling  
Max Pooling is more popular than avg. pooling

Pooling reduces the spatial size of the feature and only few operations are required.



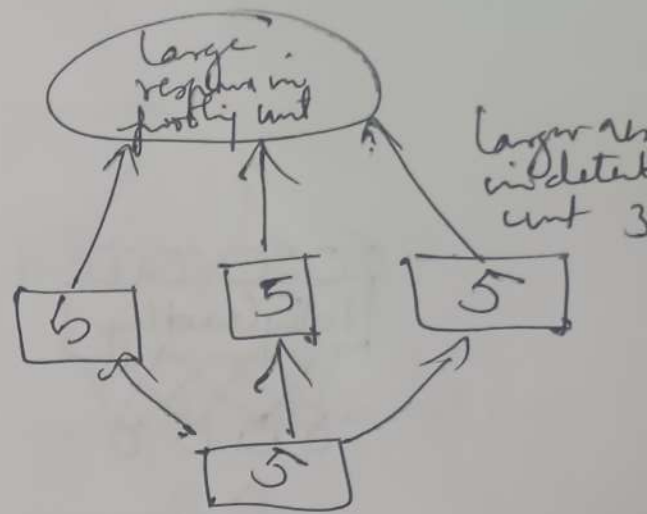
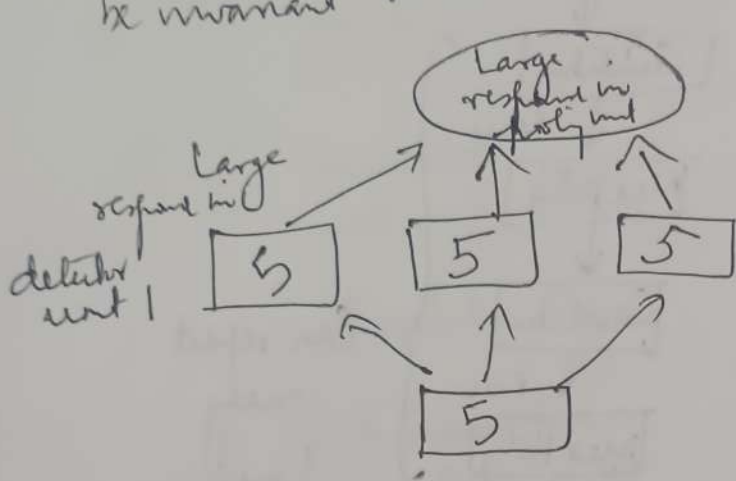
Image



Every value in the bottom row has changed, but only half of the values in the top row have changed. More pooling introduces invariance.



The use of pooling can be viewed as adding an invariance strategy prior that the function the layer learns must be invariant to small translations.



A pooling unit that pools over multiple features that are learned with separate parameters can learn to be invariant to translations of the  $\mathbb{Y}$ .

3 filters are intended to detect a handwritten image 5

Pooling is essential for handling  $\mathbb{Y}$  of varying size. For example if we want to classify images of variable size the  $\mathbb{Y}$  to the classification layer must have a fixed size.

Varying the size of an offset b/w pooling regions so that classification layer receives the same no. of summary statistics.



## Variants of the Basic Convolution

- when we refer to convolution in the context of neural n/w we actually mean an operation that consists of many applications of convolution in parallel.
  - Convolution with single kernel can only extract one kind of feature.
  - We use many filters at once.
  - Each filter extracts a diff. kind of feature (edges, colors, textures).
  - CNN perform many convolutions together — that's gives multiple feature maps.
  - Each layer of network extract many kinds of features at many locations — many kernels per layer.
  - maps are not single value grids. Each pixel has multiple values (R, G, B).
  - Each pixel is like a vector.
  - Each pixel = R, G, B.
  - input as a 3D tensor  $w \times h \times 3$  channels.
  - In a multilayer CN, up to the second layer is the output of first layer. Each layer receives multiple up and applies new filters.
- CNN works on 3D data. one dimension: R, G, B

## Multi-channel Convolution

CNN → deal with multichannel data

$Y_p$  → several channels  
 Kernels → handle all these channels

Input RGB

Operations are not guaranteed to be commutative even if kernel flipping is used

$$a + b = b + a \quad \times$$

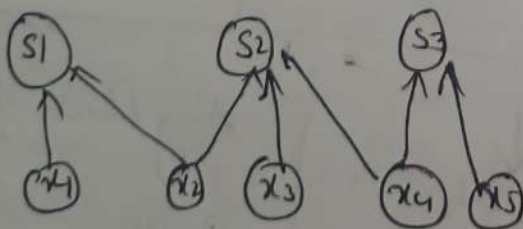
Multi-channel operations are commutative if

no. of  $Y_p$  channels = no. of  $Y_p$  channels.

$$Z_{ijk} = \sum_{l,m,n} V_{l+j-1, k+n-1} \cdot K_{l,i, m,n}$$

We may want to skip over some portions of the kernel in order to reduce the computation cost.

Stride



## Zero Padding

One essential feature when we do convolution the image becomes smaller after every layer. To prevent it shrinking too fast we add extra zeros around the borders of the image.

$Y_p$  5x5    kernel 3x3    o/p 3x3

Without zero padding o/p will shrink the image size.

### 3 types

① No zero padding - Valid Convolution

We don't add any zeros. O/p becomes smaller than  $Y_p$ . Valid convolution - because it only uses valid pixels inside the image.

② Enough zero padding - Same convolution

We add zeros to keep

$$O/P = Y_p \text{ size}$$

$$5 \times 5 - 3 \times 3 \text{ kernel} + \text{padding} = 1$$

$$O/P \ 5 \times 5$$

③ Full convolution

We add enough zeros. multiple times

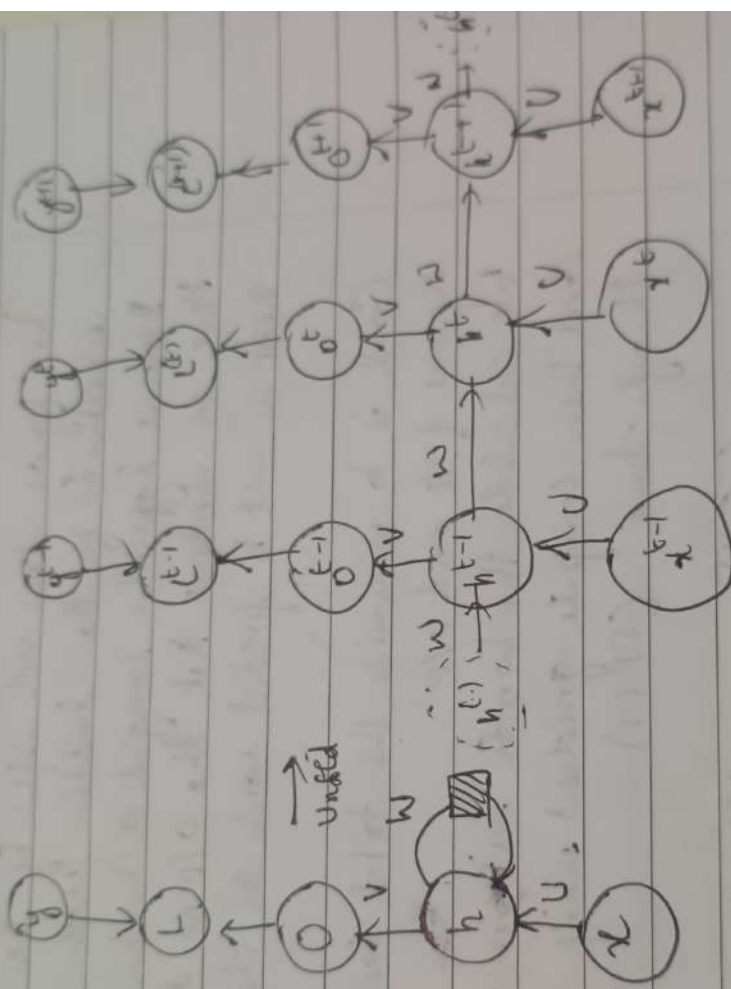
larger o/p than  $Y_p$

$$O/p \text{ width} = m + k$$



Recurrent Neural Network

Recurrent n/w that produce an output at each time step and have recurrent connections b/w hidden units.



4p sequence of  $x$  values are mapped to the corresponding sequence of o/p values.  
 Loss is calculated from  $o$  and  $y$ .

$$\hat{y} = \text{softmax}(o)$$

$$y = \text{target}$$

$$o, v, w \rightarrow \text{weight matrix}$$



Forward Propagation begins with a specified (to avoid state  $h^{(0)}$ ) time step  $t=1$  to  $t=T$  / weights

$$a^{(t)} = b + w \cdot h^{(t-1)} + Ux^{(t)}$$

$$h^{(t)} = \tanh(a^{(t)})$$

$$o^{(t)} = c + Vh^{(t)}$$

$$y^{(t)} = \text{softmax}(o^{(t)})$$

$b, c \rightarrow b, c$   
 $U, V, w \rightarrow \text{weights}$   
 $x^{(t)}$  vectors

The figure does not specify the choice of activation fn for the hidden units

Assume differentiable tangent activation fn

Here we assume that the o/p is discrete so if the RNN is used to predict words or characters

Goal is to minimize the total loss over a square of  $4p$  and  $o/p$

We calculate the loss fn across time steps

And then compute gradient using Backpropagation through time (BPTT)

$$L(\{x^{(1)}, \dots, x^{(T)}\} \{y^{(1)}, \dots, y^{(T)}\})$$

$$= \sum_t L^{(t)}$$

$L^{(t)}$  loss at time step  $t$   
 $T$ : total no. of time steps



## Model Prediction

$$\hat{y}^{(t)} = \text{model}(y^{(t)} | \{x^{(t)}, \dots, x^{(1)}\})$$

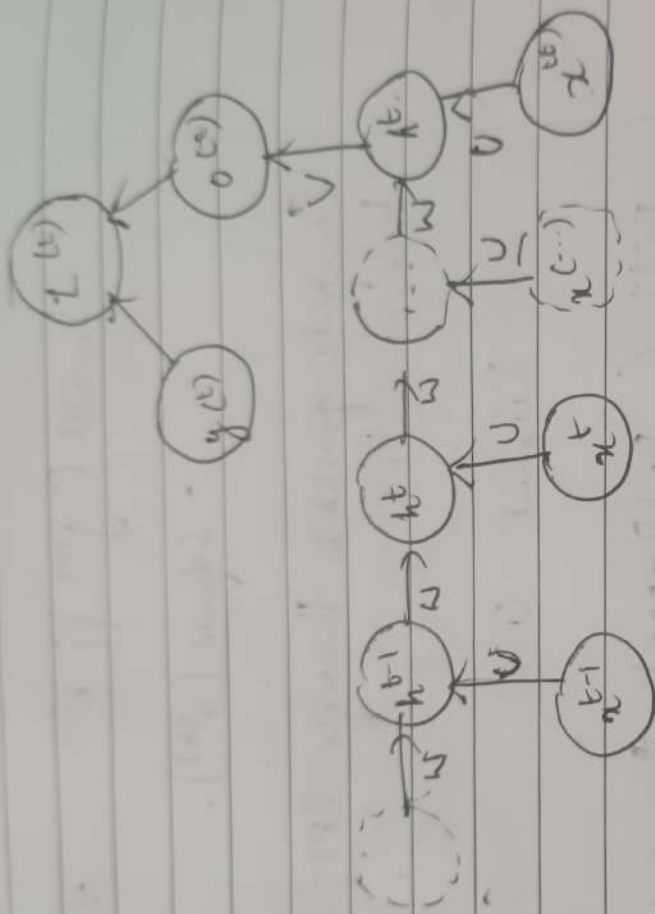
$$= \text{sigmoid}(\delta^{(t)})$$

Compute gradients through BPTT and update  $U, V, W, b, c$ .

Runtime  $O(t)$  Sequential - cannot parallelize

## BPTT (Backpropagation through time)

- applied to the unrolled RNN w/w.
- Compute total loss  $L = \sum L^{(t)}$ .
- Perform forward pass through time.
- Compute gradients  $\frac{\partial L}{\partial w}, \frac{\partial L}{\partial v}, \frac{\partial L}{\partial U}$ .
- update parameters using gradient descent.



Time Unfolded recurrent neural netw with a  
 single output the end of the sequence.  
 Used to summarize a sequence and produce a  
 fixed size representation  
 which can be used for further processing.

## Feedforward vs CNN vs RNN

FF

Learns separately for each up.

CNN

Shares parameters locally.

RNN

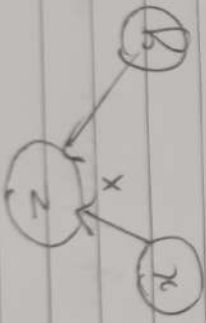
Shares parameters through time.  
Each o/p depends on the prev. o/p.

## Unfolding Computational graphs.

- A computational graph is a visual representation of how computations are performed step by step.
- It helps us formalize and visualize how up are transformed into o/p.
- Each node in the graph represents a variable. (scalar, vector, matrix or tensor)
- Each edge represents an operation applied to variables.

$x \rightarrow y$   
y is computed from x.

$$z = x \times y$$



Recurrent Computation is defined as

$$s^{(t)} = f(s^{(t-1)}; \theta)$$

↑  
state at time step

↑  
previous state

↑  
parameters weight shared across all the steps

Unfolding means expanding the recurrence to show how each state depends on earlier ones

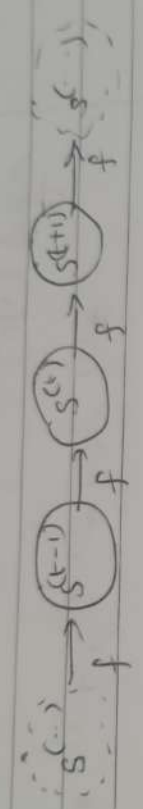
$$\begin{aligned} s^{(t)} &= f(s^{(t-1)}; \theta) \\ &= f(f(s^{(t-2)}; \theta); \theta) \\ &= f(f(f(s^{(t-3)}; \theta); \theta); \theta) \\ &= f(f(f(\dots f(s^{(1)}; \theta) \dots)) \end{aligned}$$

info flows through many layers over time  
creating a deep chain of computation  
 $\theta \rightarrow$  parameters

Time delay  
by an operator  
f(s) → f(s) + θ



Unfold



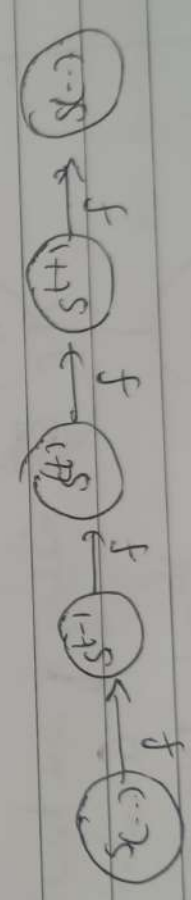
A chemical form of a dynamical system

$$s^{(t)} = f(s^{(t-1)}; \theta)$$

Unfold 3 times

$$s^{(t)} = f(f(f(s^{(t-3)}; \theta); \theta); \theta)$$

Recurrent models reuse parameters (same  $\theta$ ) over many time steps.



## Bidirectional RNNs

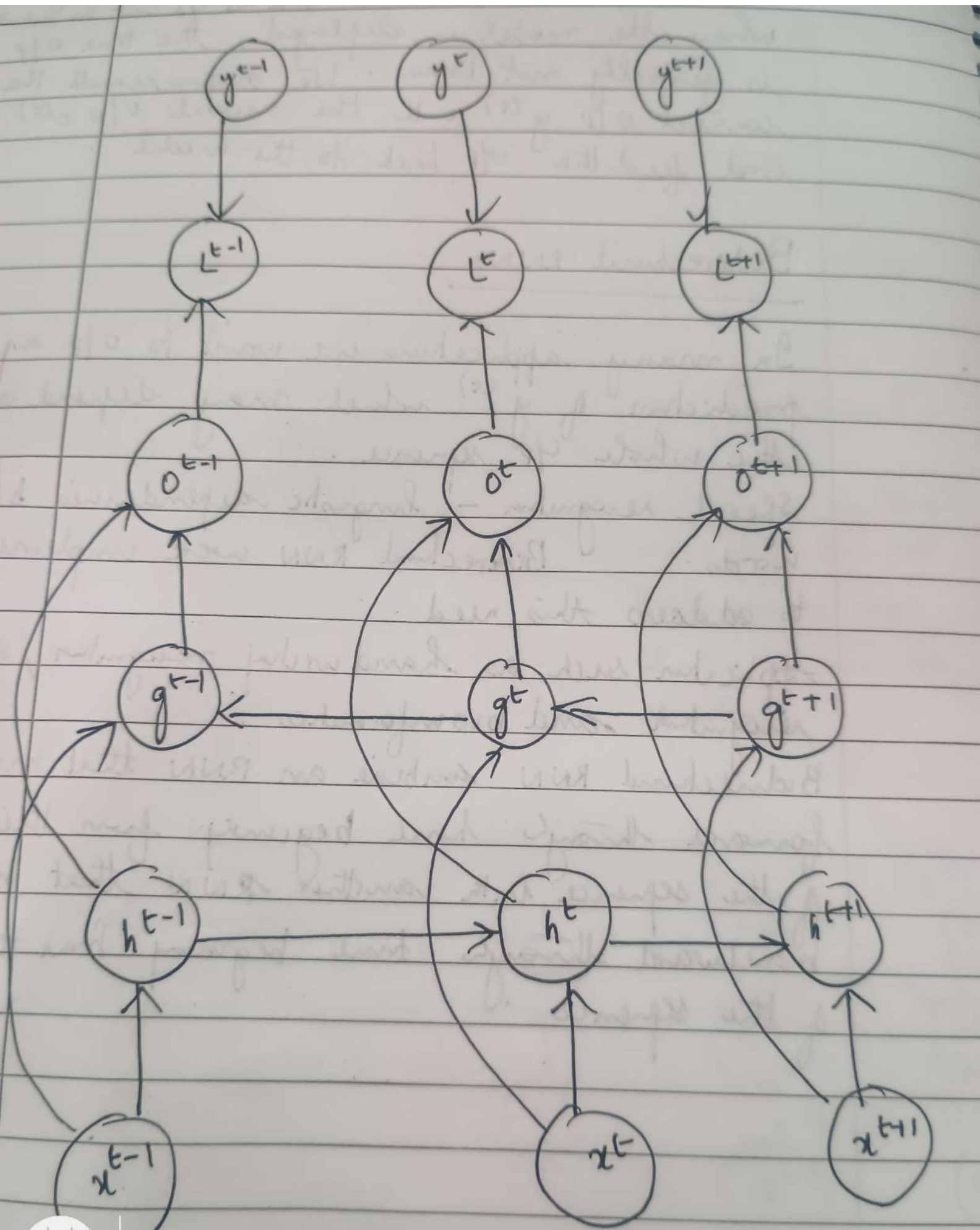
In many applications we want to o/p a prediction of  $y^{(t)}$  which may depend on the whole  $x^p$  sequence.

Speech recognition - linguistic dependencies b/w nearby words. Bidirectional RNN were implemented to address this need.

Applications such as handwriting recognition, speech recognition and bioinformatics.

Bidirectional RNN combine an RNN that moves forward through time beginning from the start of the sequence with another RNN that moves backward through time beginning from the end of the sequence.





is 0  
info from hidden state at prev. time step  
 $t-1$

## Encoder - Decoder Sequence to sequence architecture

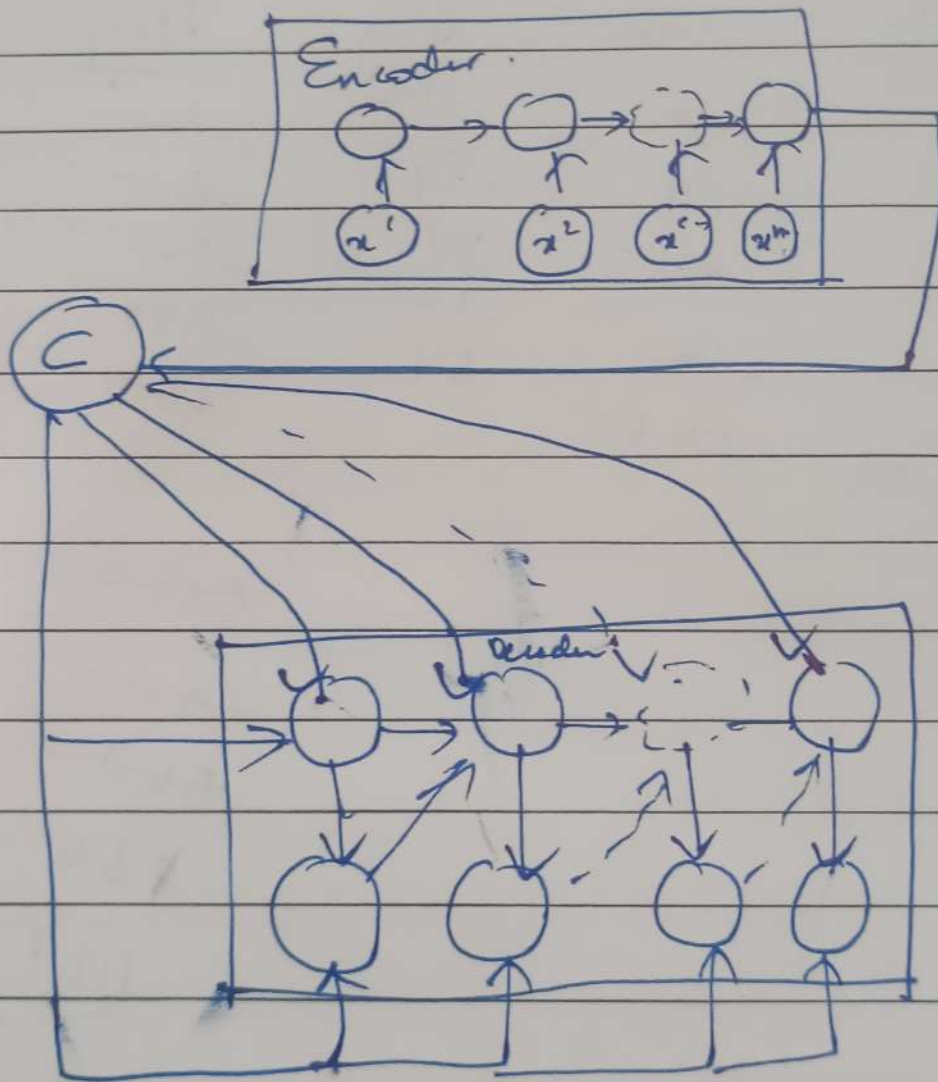
Discuss how an RNN can be trained to map an  $Y_P$  sequence to an  $O/P$  sequence which is not necessarily of the same length.

- In speech recog, machine translation,  $Y_P$   $O/P$  sequence in the training set are not of same length.
- Representation of context  $c$ .
- The context  $c$  might be a vector or sequence of vectors that summarise the input sequence  $x = x^{(1)} \dots x^{(n)}$ .
- An encoder or reader or  $Y_P$  RNN processes the  $Y_P$  sequences. The encoder emits the context  $c$ , usually as a sample of its final hidden state.



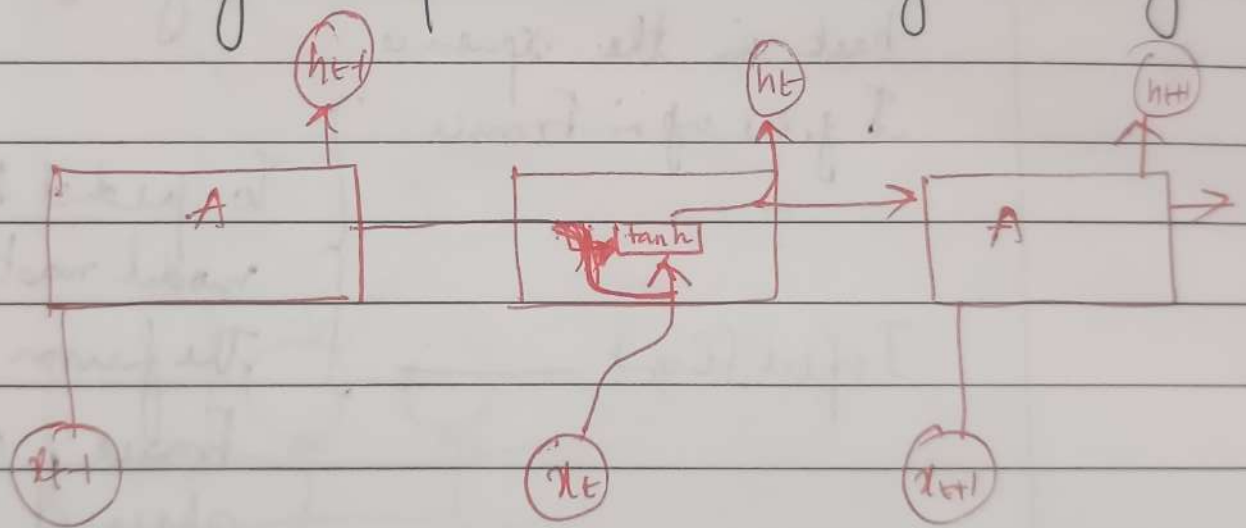
- A decoder or writer w/ o/p of RNN is fixed length to generate the o/p sequence  
 $Y = y^{(1)} \dots y^{(K)}$

Limitation: - Dimension of context  $C$  can be too small to summarize a long sequence.



# LSTM (Long Short Term Memory)

- special kind of RNN
- made to remember info for long time periods and avoid the long term dependency problem.
- it can keep both short-term and long-term memories easily.
- All recurrent neural n/w have the form of a chain of repeating modules of neural n/w.
- In std RNN this repeating module will have very simple structure i.e. single tanh layer.



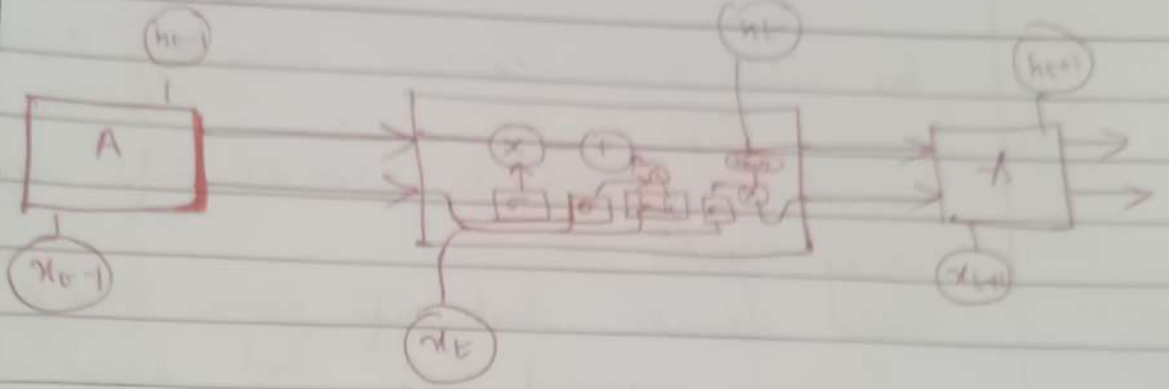
Each green box A is one cell or time step in the n/w

$x_t$  → i/p at each time step

$h_t$  → o/p or hidden states passed to the next step.

Simple RNN → box - one operation - usually tanh layer that updates the hidden state.

# LSTM Cell



$+ \times \rightarrow$  math operations  
 $\circ$  tank  $\rightarrow$  neural net layers  
 black lines info flow

Takes input  $x_t$  and previous hidden state  $h_{t-1}$

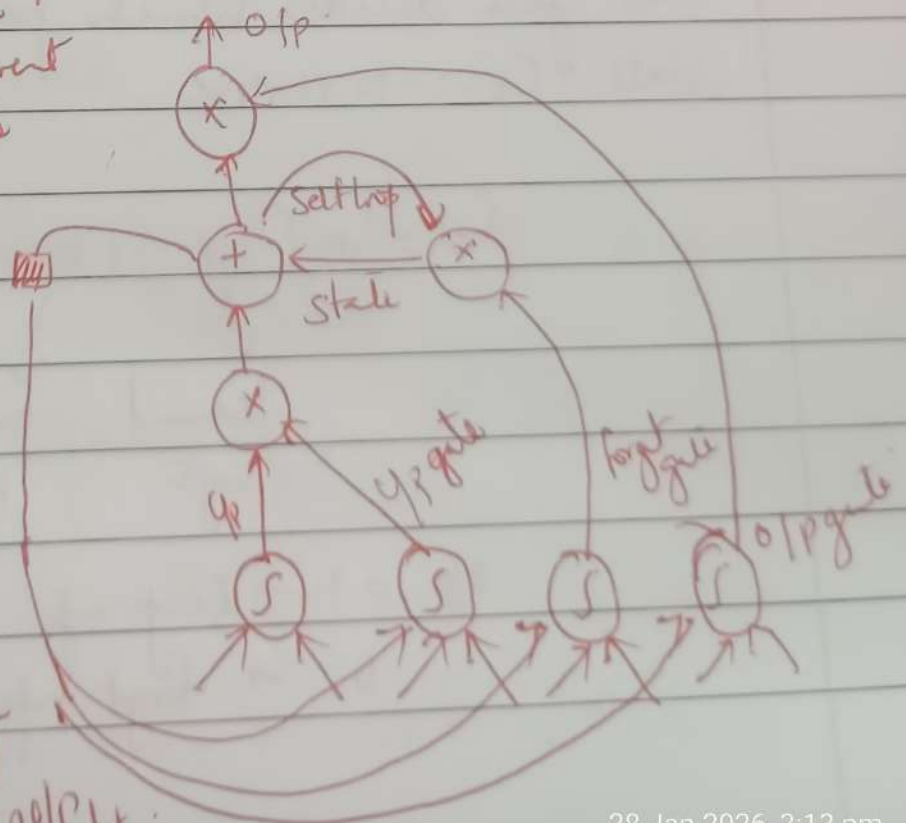
Uses gate to decide,  
 what to forget (forget gate)  
 what new info to add (input gate)  
 what to output (output gate)

$C_t \rightarrow$  updated, new  $h_t$  is sent to next state

**Input gate**  $\rightarrow$  controls how much new info from current input should enter the cell state

**Forget gate**  $\rightarrow$  decides what part of old info to erase from the cell state

**Output gate**  $\rightarrow$  decides what part of the cell state should be sent as output



The cell state carries memory across time steps.  
LSTM can add or remove info from it using gates.

Gates decide what info passes through. They use sigmoid ( $\sigma$ ) and pointwise multiplication ( $\times$ ) operators.

$\sigma$  - o/p value b/w 0 and 1,  
0  $\rightarrow$  block everything  
1  $\rightarrow$  let everything pass.

## Step by step LSTM Walkthrough.

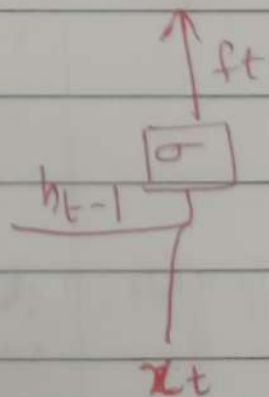
### Forget gate.

This gate decides which part of the old cell state  $C_{t-1}$  to remove.

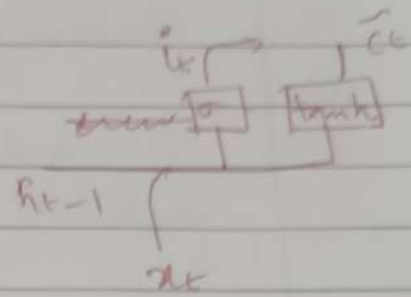
$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

$\downarrow$  learned wghts       $\downarrow$  prev. hidden state       $\downarrow$  bias.

o/p = 1  $\rightarrow$  keep info.  
0  $\rightarrow$  forget info.



## Input gate



## Input Gate Eqs

$$i_t = \sigma(w_i \cdot [h_{t-1}, x_t] + b_i)$$

## Candidate value Eqn.

$$\tilde{c}_t = \tanh(w_c \cdot [h_{t-1}, x_t] + b_c)$$

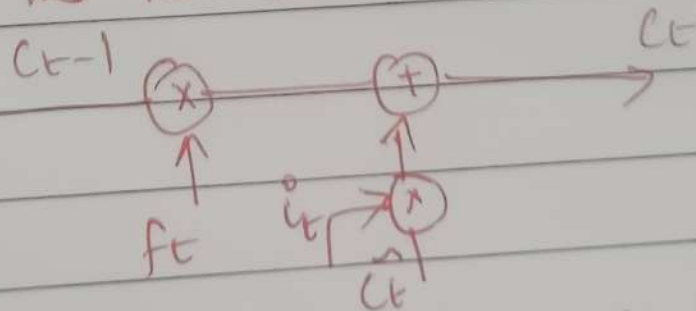
Creates new candidate memory values.

Combine both

$$i_t * \tilde{c}_t \rightarrow \text{new update to the cell.}$$

## Update Cell state

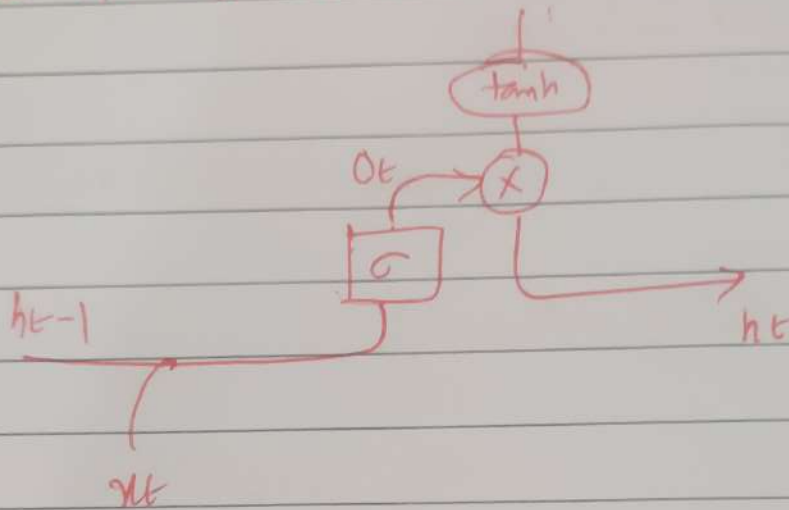
We use the forget gate and  $\frac{1}{2}$  gate together to form the new cell state.



$$c_t = f_t * c_{t-1} + i_t * \tilde{c}_t$$

## Output gate

This gate decides what part of the cell state to o/p as hidden state  $h_t$ .



$$o_t = \sigma (w_o \cdot [h_{t-1}, x_t] + b_o)$$

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

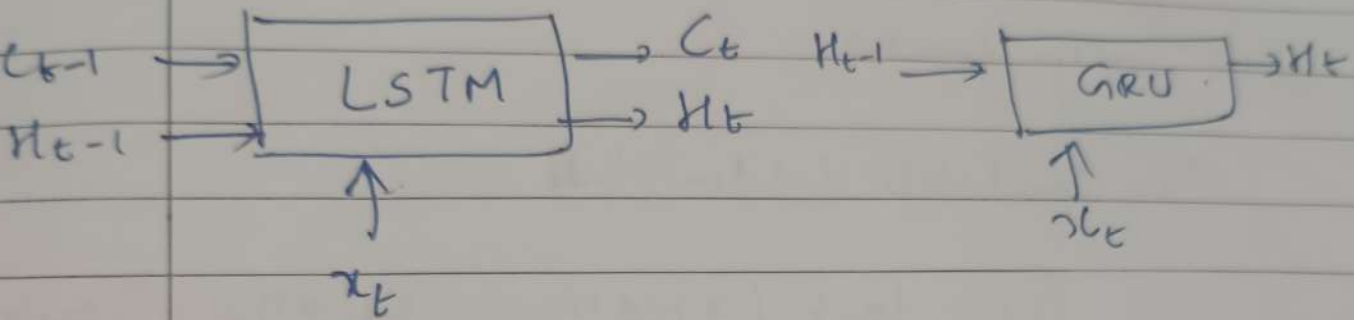
$o_t \rightarrow$  how much of the memory should be visible outside.

$h_t$  - becomes the o/p of this time step and also goes to the next one.



# GRU (Gated Recurrent Unit)

- Similar to LSTM
- uses gates to control the flow of info.
- simple architecture.



GRU doesn't have a separate cell state.  
Easier to train.

At timestamp  $t$ , it takes an input  $x_t$  and hidden state  $H_{t-1}$  from the previous timestamp. Later it outputs a new hidden state  $H_t$  which is passed to the next timestamp.

- GRU has 2 gates: Reset gate, Update gate

## Reset gate (Short term memory)

- is responsible for the short term memory of the n/w.

$$r_t = \sigma(x_t * U_r + H_{t-1} * W_r)$$

## Update gate (Long term memory)

$$U_t = \sigma(x_t * U_u + H_{t-1} * W_u)$$

To find the hidden state  $H_t$  is a 2 step process.

### ① Candidate hidden state

$$\hat{H}_t = \tanh(x_t * U_g + (r_t * H_{t-1}) * W_g)$$

if  $r_t = 1 \Rightarrow$  entire info from the previous hidden state  $H_{t-1}$  is considered.

$r_t = 0 \rightarrow$  info from the prev. hidden state is completely ignored

Hidden state  $H_t$

$u_t \rightarrow$  update gate

$$H_t = u_t * H_{t-1} + (1 - u_t) * \hat{H}_t$$

$u_t = 0$  first term vanishes.

new hidden state will not have much