

Joseph John Gerlach - proposed nervous system is single continuous network - (biology) - Reticular theory

Camillo Golgi - staining technique.

Neuron - coined by - their Rich Wilhelm Ludwig von Waldeyer - Hartz.

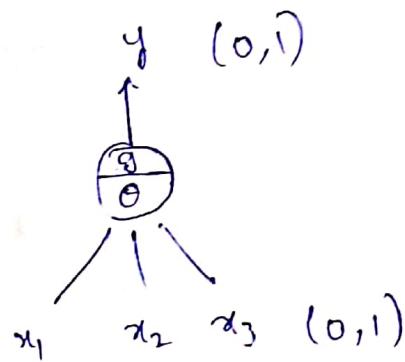
Perceptron - Frank Rosenblatt - idea initially pursued by Navy to translate languages

Minsky - Papert - Perceptual Paper - limitations highlighted

gradient descent - Cauchy.

Unsupervised pre training - 2006 - deep learning research boomed again.

McCulloch Pitts Neuron.

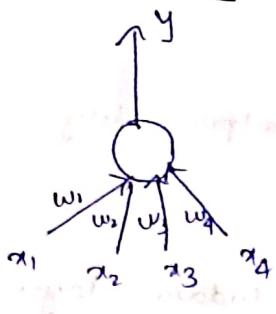


$$g(\theta) = x_1 + x_2 + x_3$$
$$y = 1 \text{ if } g(\theta) > 0$$
$$= 0 \text{ if } g(\theta) \leq 0$$

Single MPN can represent a boolean fn which is linearly separable.

Sifat

## Perceptron model



$$y = 1 \text{ if } g(x) > 0$$

$$g(x) = \sum w_i x_i$$

$$y = 1 \text{ if } \sum_0^n w_i x_i > 0$$

$$\text{where } w_0 = -\theta, x_0 = 1$$

$w_0$  is the bias.

## Perceptron learning Algorithm

till convergence (no errors)

pick random  $x \in P \cup N$

$$\text{if } x \in P \text{ if } \sum x_i w_i \leq 0$$

$$x = x + w$$

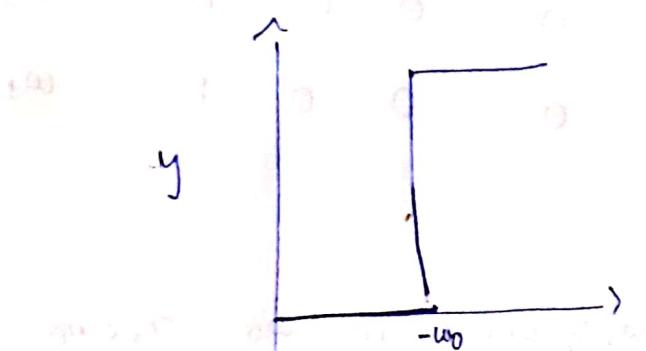
else

$$x \in N \text{ if } \sum x_i w_i > 0$$

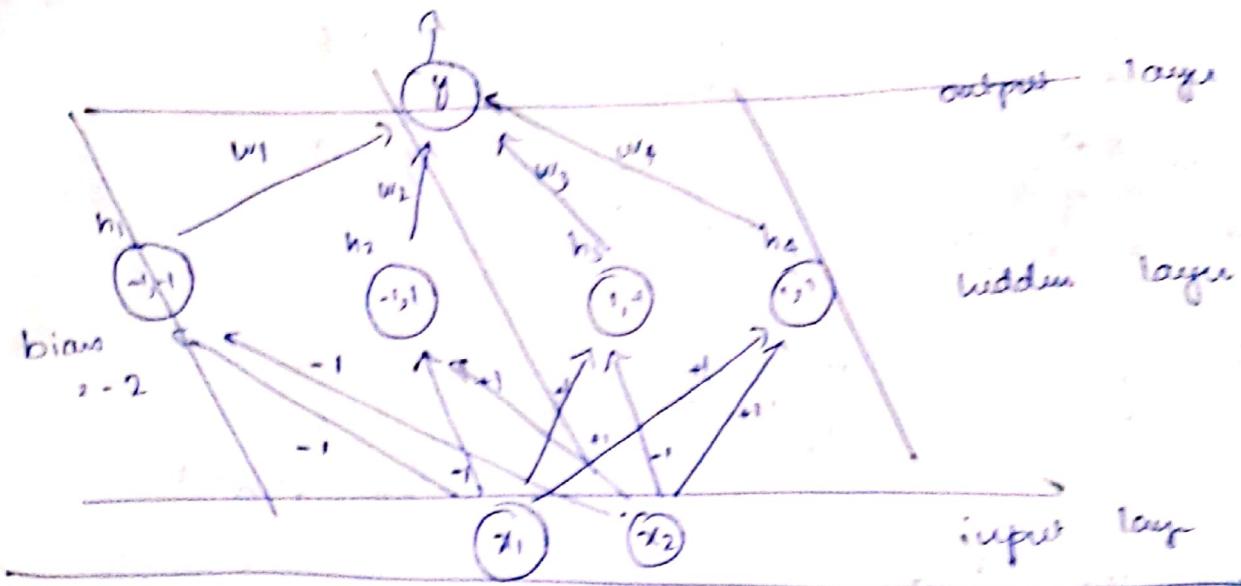
$$x = x - w$$

Network of perceptron can deal with non-linearly separable data.

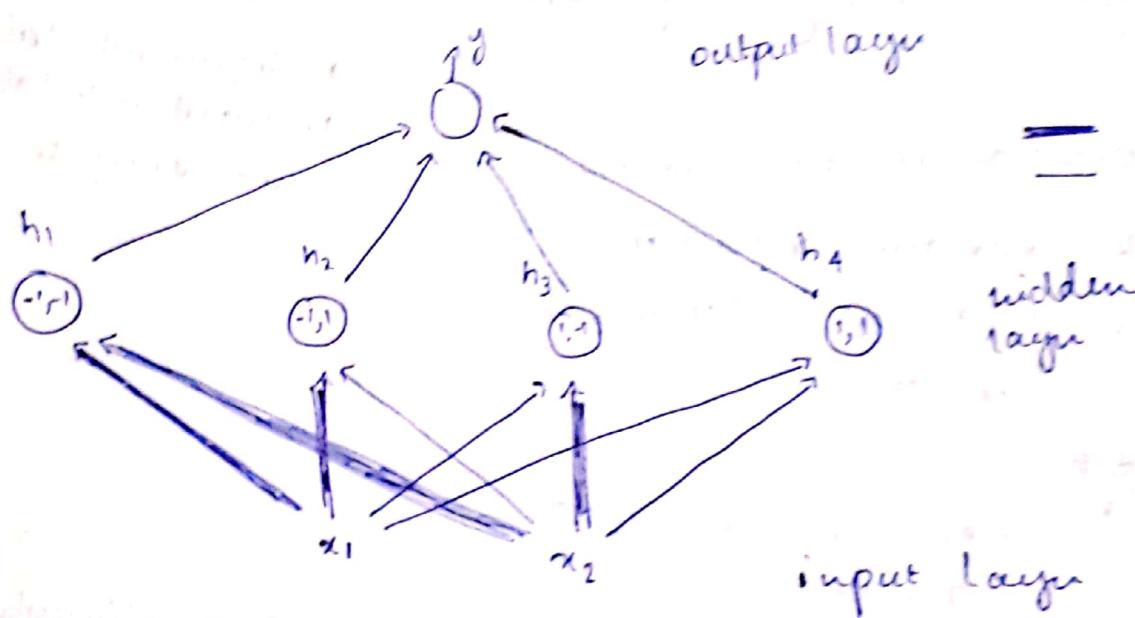
for two inputs, two input combinations are not linearly separable (XOR & !XOR)



Network of perceptrons



Network of perceptrons - Multilayer Perceptrons



XOR for this

$x_1$	$x_2$	XOR	$h_1$	$h_2$	$h_3$	$h_4$	$\Sigma w_i h_i$
0	0	0	1	0	0	0	$w_1$
0	1	1	0	1	0	0	$w_2$
1	0	1	0	0	1	0	$w_3$
1	1	0	0	0	0	1	$w_4$

for XOR,

$$w_1 < w_0$$

$$w_2 \geq w_0$$

very neuron in hidden layer takes for each combination of the input and weights  $w_i$  can be adjusted to get the required output for the case.

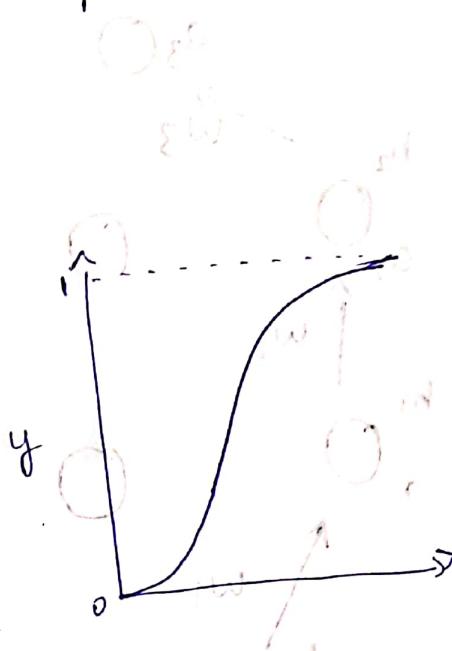
(minimum of combination of inputs)

thus all  $n$  inputs,  $2^n$  perceptrons in hidden layer. (for boolean functions)

Sigmoid neurons

logistic function

$$y = \frac{1}{1 + e^{-w^T x}}$$



Output  $y$  is no longer binary.

Error function

$$\text{loss} = \frac{1}{2} (\sum y_i - f(x_i))^2$$

$\text{loss} = \frac{1}{n} \sum (y_i - f(x_i))^2$

gradient descent

$$\theta_{t+1} = \theta_t - \eta \text{loss}$$

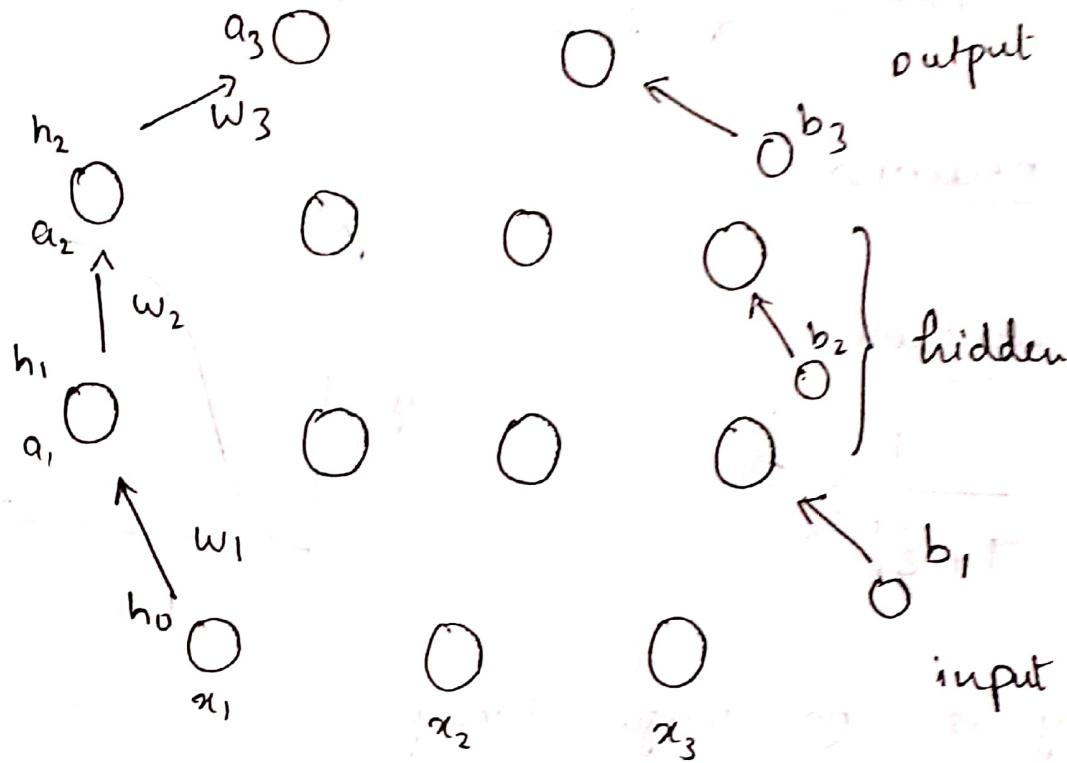
gradient descent for sigmoid

$f'(x) = y(1-y)$

$\text{loss} = \frac{1}{n} \sum (y_i - f(x_i))^2$

## Feed forward Neural Network

Each neuron in hidden layer has 2 layers (pre activation & activation)



$$a_1 = b_1 + w_1 h_0$$

$$a_i = b_i + w_i h_{i-1}$$

$$h_i = g(a_i(x))$$

$g(x)$  is the activation function.

$$f(x) = D(a_3(x)) \quad \text{Output function}$$

$O$  is Output Activation function.

$$y = f(x)$$

$$= O(w^3 g(w^2 g(w^1 x + b_1) + b_2) + b_3)$$

## Selecting o/p function:

when output needs to be between zero and one, we can use sigmoid function.

for large output values, we can use linear fn.

$$f(x) = O(a_0) \quad (\text{used for regression in this case})$$

$$= w_0 a_0 + b_0$$

for classification problems, softmax function can be used as output function

In cases where  $y$  is a probability distribution, we use softmax function to keep  $y$  in a probability distribution.

To find diff between two probability distributions, we can use cross entropy.

$$L(\theta) = \text{Loss} = -\sum y \log \hat{y}$$

to minimize this loss function, minimize  $L(\theta)$

$$\text{we should minimize } L(\theta) = -\log \hat{y}$$

$$\text{with } \theta \text{ minimize } -L(\theta) = \log \hat{y}$$

where  $\hat{y}$  is the linear scaling softmax of linear combination of all weights in the previous layers.

$$\hat{y} = O(w_3 g(w_2 g(w_1 x + b_1) + b_2) + b_3)$$

$\hat{y}$  is the probability

$\log \hat{y}_x$  is the log likelihood of the data

		Outputs	Probabilities
Output	Activation	Linear	Softmax
Loss function		Squared error	Cross entropy

$$g(z) = \frac{1}{1+e^{-z}}$$

$$g'(z) = g(z)(1-g(z))$$

$$g(z) = \tanh(z)$$

$$g'(z) = 1 - (g(z))^2$$

Nesterov Accelerated gradient descent (NAG).

the next step gradients are found beforehand

so convergence is quicker

Oscillations are smaller

missing minima is , chance of

normal

gradient

descent

all data calculation is done and then

weights updated,

actual answer no approx.

updated after each data point

This is an approx gradient

Stochastic L.D. weights

mini-batch SGD updated for set no of data from

To tune learning rate:

Try different values on log scale ( $0.0001, 0.001, \dots$ )

Run few epochs of each & figure out what works best. Then do fine search around this value (if  $\delta=0.1$  then see for 0.05, 0.2 etc)

adaptive learning rate:  
sup. decay:

halve  $\alpha$  every 5 epochs

- halve  $\alpha$  after an epoch if CV error is more than what is at end of previous

exponential decay:

$$m^2 = m_0 e^{-kt}$$

$t \rightarrow$  step no

$1/t$  decay:

$$m^2 = \frac{m_0}{1+kt}$$

$$\alpha = \min\left(1 - 2^{-1 - \log_2(1 + \frac{1}{2\delta}) + 1}, \alpha_{\max}\right)$$

$$\alpha_{\max} = \{0.999, 0.995, 0.99, 0.9, 0\}$$

## adaptive gradient descent (adagrad)

$$v_t = v_{t-1} + (\nabla w_t)^2$$

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{v_t + \epsilon}} * \nabla w_t$$

adagrad decreases learning rate aggressively as the denominator grows.

to avoid this, denominator should be decayed.

$$v_t = \beta * v_{t-1} + (1-\beta) (\nabla w_t)^2$$

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{v_t + \epsilon}} * \nabla w_t$$

## Adam

$$m_t = \beta_1 * m_{t-1} + (1-\beta_1) * \nabla w_t$$

$$v_t = \beta_2 * v_{t-1} + (1-\beta_2) * (\nabla w_t)^2$$

$$\hat{m}_t = \frac{m_t}{1-\beta_1^t}$$

$$\hat{v}_t = \frac{v_t}{1-\beta_2^t}$$

$$\beta_1 = 0.9 \quad \beta_2 = 0.999$$

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{v_t + \epsilon}} * \hat{m}_t$$

Adam is best choice

## Eigen vectors

$$Ax = \lambda x$$

Eigen vectors corresponding to different eigen values are linearly independent.

Eigen values of a square symmetric matrix are orthogonal.

## Principal Component Analysis

dimension reduction

In new dimension are non-redundant.

(low covariance) & not noisy (high variance)

## Singular Value Decomposition

used for data compression

SVD factorises a matrix into singular vectors and singular values.

## Autoencoder

special type of feed forward neural network  
encoder takes input into a hidden rep

$$h = g(wx + b)$$

decodes again from its hidden rep

$$x_i \rightarrow f(w^T h + c)$$

model is trained to minimise certain loss fn  
which will ensure  $\hat{x}_i$  is close to  $x_i$

if  $\dim(h) < \dim(x)$

$h$  is loss the encoding of  $x_i$

encoder with  $\dim(h) < \dim(x)$  is called

under complete autoencoder.

$\dim(h) \geq \dim(x) \rightarrow$  over complete

$g$  is typically chosen as sigmoid

encoder part = PCA wif. sigmoid

- linear encoder

- linear decoder

- squared error loss fn.

- normalizing ip to 0

$$\hat{x}_{ij} = \frac{1}{\sqrt{m}} \left( x_{ij} - \frac{1}{m} \sum x_{kj} \right)$$

---

Sparse Autoencoder

tries to ensure most of the time neuron is inactive

Contractive Autoencoder

tries to prevent an overcomplete autoencoder from learning identity fn.

generally;

Simple model : high bias, low variance  
Complex model : low bias, high variance

Both b and v contribute to MSE

### Different types of regularisation:

- Data augmentation
- Parameter sharing & tying
- Adding noise to i/p
- Adding noise to o/p.
- Early stopping
- Ensemble methods
- dropouts

12

$$L(w) = L(w) + \frac{\alpha}{2} \|w\|^2$$

$$w_{t+1}, w_t - \eta \nabla L(w_t) - \eta d w_t$$

### Data augmentation

eg for image, rotate image in different angles

### Parameter sharing & tying

Used in CNNs

Same filter applied at diff. positions of image

Same weight matrix acts on diff input neurons

## Adding noise to I/P

Eg: Autoencoder

Adding noise to I/P

The true labels might be wrong  
Use soft targets. taking  $\epsilon \rightarrow 0$  small the count

a class:

apple	orange	banana	lemon
1	0	0	0

hard target

apple	$1-\epsilon$	$\epsilon/3$	$\epsilon/3$	$\epsilon/3$
-------	--------------	--------------	--------------	--------------

soft target

## Early stopping

Track validation error

If there are no improvements in cv error  
in last  $p$  steps, stop training.

## Ensemble methods

Combine output of different models

### dropout

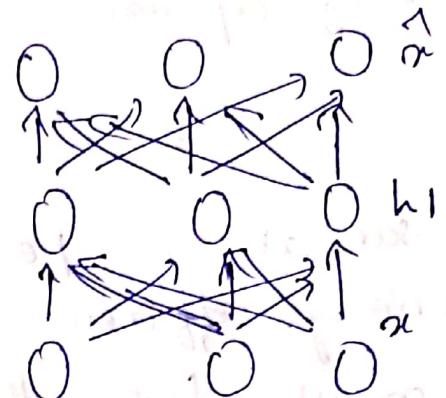
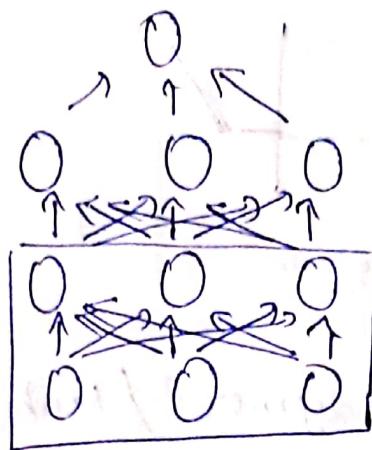
temporarily remove a node and  
all its connections resulting in  
thinned h/w

share weights across all h/w's.

Sample different h/w for each t instance

Unsupervised Pre training

deep learning became popular again due to this.  
for example



} unsupervised.

the weights in Layer 1 can be trained such that  
 $h_1$  captures an abstract rep of  $x$ .  
Fix the weights in Layer 1 and repeat it for

Layer 2.

repeat this process till last hidden layer.  
y is not involved at all.

After this layer will be pre training, add  
the output layer.  
pre training leads to better weight initialisations

Sigmoids

Saturated neurons can gradient to vanish

Sigmoid and tanh unit

Sigmoids are computationally expensive

## tanh

- zero centered
- gradient vanishes at saturation
- computationally expensive

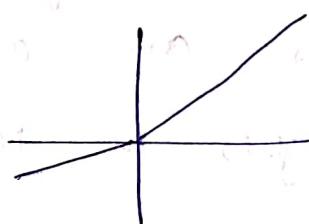


## ReLU

- does not saturate in positive region
- computationally efficient
- converges much faster than sigmoid/tanh
- units will die if  $\alpha$  is high

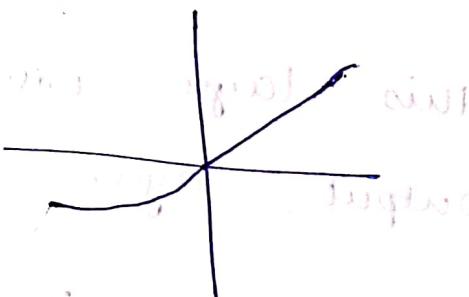
## Leaky ReLU

- no saturation
- no units will die
- computationally efficient
- close to zero centered o/p



## Exponentially linear unit

- all benefits of ReLU
- computationally expensive



## Maxout Neuron

$$\max(w_1 T_x + b_1, w_2 T_x + b_2)$$

- generates ReLU & Leaky ReLU
- No Saturation
- doubles no of parameters

Professor's notes:  
Sigmoids are bad  
ReLU ~ Standard for CNNs  
explore LReLU/Masout/ELU  
tanh ~ LSTMs/RNNs

## initialisation strategies

if all weights are initialised zero, all neurons in layer 1 fire up.  
this symmetry stays same for all layers

### Batch Normalisation

Add a normalisation layer after a neuron  
(this is differentiable wrt backpropagation)

One-hot encoding  
representing a single representation of vector (mathematical module) in terms of columns

Co-occurrence matrix  
It is a terms x terms matrix which captures the number of times a term appears in the context of another term

SVD for learning word representations.

SVD gives best rank-k approx of original data

To obtain filter map, sum filter over the input

Parameters involved:

- width of input ( $w_1$ )
- height of input ( $h_1$ )
- depth of input ( $d_1$ )

Stride(s)

- no. of filters ( $K$ )
- filter dimension ( $F$ )
- Output  $(w_2 \times h_2 \times d_2)$

In general,

$$w_2 = w_1 - F + 1$$

$$h_2 = h_1 - F + 1$$

To have output same dim of input, use padding

Stride - area intervals at which filter is applied

$$d_2 = K$$

$$w_2 = \frac{w_1 - F + 2P}{S} + 1$$

$$h_2 = \frac{h_1 - F + 2P}{S} + 1$$

$P$  = padding

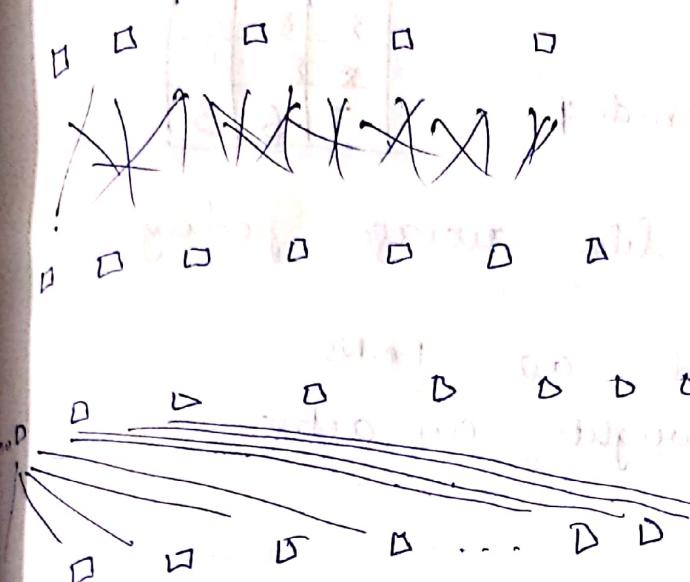
Job before applying bias

## Convolutional Neural Networks:

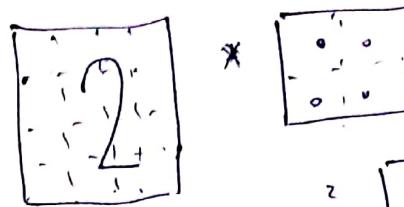
Along with learning weights of classifier, meaningful kernels, filters are learnt.

regular feed forward

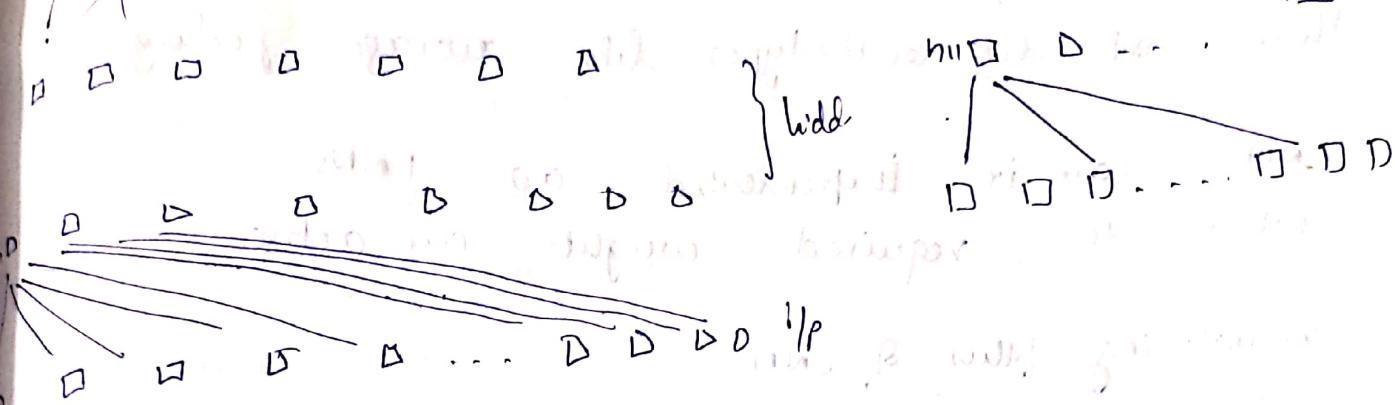
5 classes



CNNs.



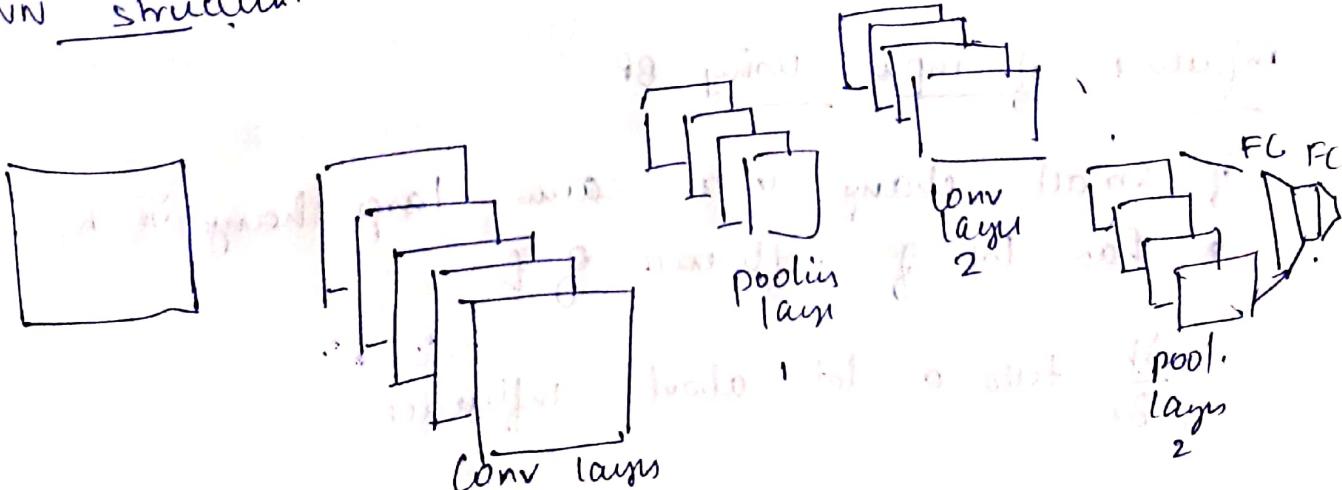
hidden



only few neurons contribute to '2'

There are many kernels but they are shared by all locations in the image.  
This is called weight sharing.

CNN structure:



Input \* filter

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

maxpool  
 $(2 \times 2) \text{ (2 stride)}$

8	4
7	5

stride 4

8	8	4
8	8	5
7	6	5

thus add different types like average pooling

CNN can be implemented as FNNs  
when the required weights are active

### Visualising filters of CNN

finding input which maximally excites a neuron,

$$\approx \frac{w}{\|w\|} \text{ where } w \text{ is the filter}$$

so thus filters are regarded as pattern detectors.

### Occlusion Experiments

graying out certain regions to see how it affects probability of output

### Influence of input using BP

if small change in  $x$  causes large change in  $h$ ,  
it has lot of influence of  $h$ .

$\frac{\partial h}{\partial x}$  tells a lot about influence

## guided back propagation

- feed i/p. & do forward pass
- consider some neurons in some feature map in some layer
- find influence of i/p at this neuron
- set all other neurons in region to 0
- back propagate to i/p
- During forward pass, ReLU pass only the units
- during back pass, no gradient passes through the dead ReLU units
- all -ve gradients are set to 0
- plotting this gives a better picture of true influence of each node of map on output.

## Optimisation over Images

keep weights fixed  
vary input parameters till image is obtained

## Creating image from embeddings.

reconstruction of image after different layers.

## Deep Dream

start with some image

change the image so that some neurons fire even more

## Deep art

add an image to another image to add colors and textures.

fooling CNNs, and to get away from it

by adding some noise, CNN can be misguided to classify

## Sequence Learning Problem

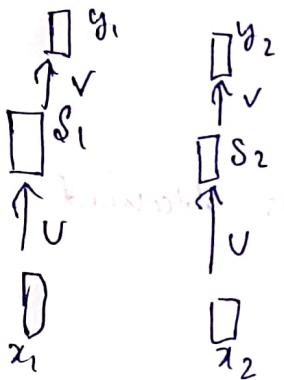
In FFNN, CNN, size of input was fixed.

each input was independent of previous

& Examples for sequence learning.

Predicting next letter.

annotation: put of speech tag to each word



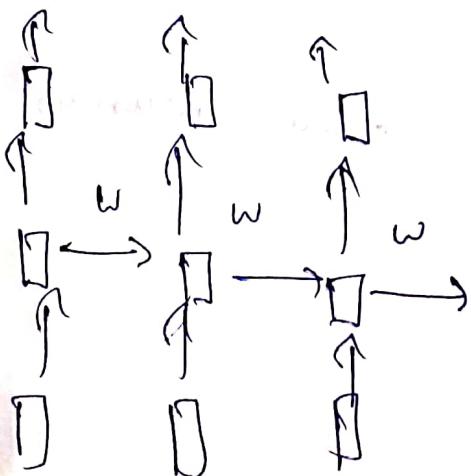
$$s_i = \sigma(Ux_i + b)$$

$$y_i = O(Vs_i + c)$$

is timestep

at each time step, feed all previous inputs  
(this won't work)

Add recurrent connections to  $n/w$ .



backpropagation over all previous time steps) through time (BPTT)

to avoid vanishing gradients, we truncated

RNN has a full state size, it needs to be selectively read, written & forgotten.

## Long - Short Term Memory (LSTM)

Selective write:  
decide what fraction of each element of previous state to be passed to next state decided by vector o vanishing bw. o E.g. RNN learns o along with other parameters.

Selective read:  
another function to eliminate some info passed over

Selectively forget:  
another gate added to forget when needed

$$o_t = \sigma(w_{oh} h_{t-1} + U_o x_t + b_o)$$

$$i_t = \sigma(w_{ih} h_{t-1} + U_i x_t + b_i)$$

$$f_t = \sigma(w_{fh} h_{t-1} + U_f x_t + b_f)$$

of before time step whether to keep or not

$$\hat{s}_t = \sigma(w_{ht-1} + U_{xt} + b)$$

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

$$h_t, o_t \odot \sigma(s_t) \text{ & } \text{rnnout} \rightarrow h_t$$

hSTM has many variables which include diff no of gates & also diff arrangement of gates.

### a popular variant - GRU (gated Recurrent Unit)

No explicit forget gate (the forget gate & input gates are tied)

LSTMs avoid vanishing gradients

During forward pass, gates control flow of info prevent irrelevant info from being written

During backward pass the gradients will get multiplied by gate

gradient vanish only when they should (when info is not needed)

But it does not deal with exploding gradients.

To overcome that, we gradient clipping

while back propagating, if norm of gradient exceeds a certain value, it is scaled to keep its norm within an acceptable threshold

generating sentence given an image (encoder-decoder)  
we have find  $P(y_t | y_{t-1}, I)$   
where previous outputs is the present output  
in previous outputs as given tag result taking  
as input (sequence)

CNN is used to encode the image,  
an RNN is used to decode into a  
sentence with

Applications:

text entailment

Machine translation

Translation

Image QA answering

Document Summarisation

Video captioning

Video Classification

Dialog

Attention Mechanism

defining a function  $\text{att}$  to specify which words need more weightage

Represent every location (time step) with a weight

for image, are output of one of the conv layers  
which has spatial information

~~(e.g. neural)~~ hierarchical attention: with

hierarchical attention:

for sequence of sequences (e.g. dialog between user and bot),

use 2 level RNN.

first level operates on sequences of words, each utterance is processed by another RNN which encodes this and gives single representation for sequences.

document summarisation can be done using this

methodology:

clustering of document

articles available in document

clustering of words

documents with similar

topics

clustering of subjects

clustering of words and read condition in document

clustering of users based on their

topic or other document related properties

clustering of documents based on their properties

clustering of users based on their properties