

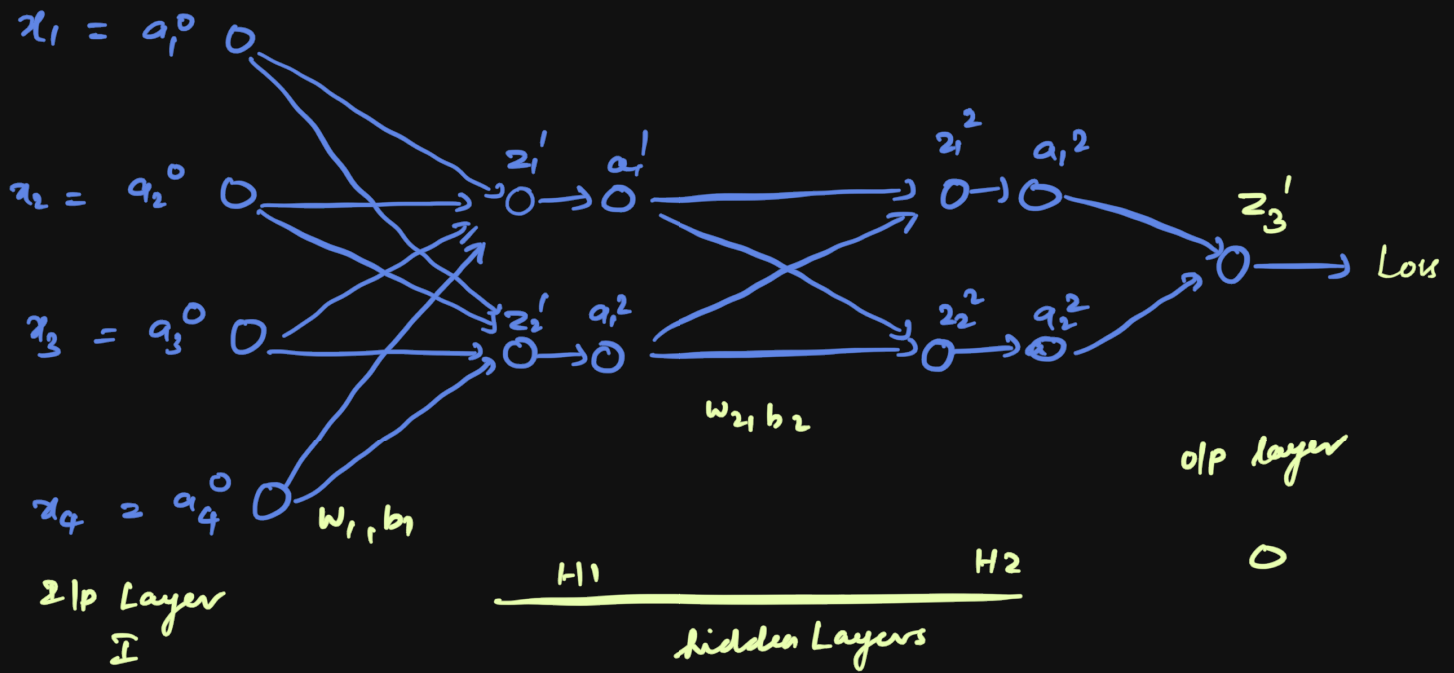
MLP & RBF

- LOKESH · N

(Thanks to Prof. Annie)

Because MLP and RBF are already covered in class. Let us go over Numericals to get a complete flavour of it

Disclaimer - There may be mistakes in derivations.
 pls cross check once.



By convention, $w_1 = 2 \times 4$ matrix

$x = 4 \times 1$ column vector

and b_1 bias vector - $[4 \times 1]$

Feed propagation

$$\begin{bmatrix} z_1^1 \\ z_2^1 \end{bmatrix} = \begin{bmatrix} w_{2 \times 4} \end{bmatrix} \begin{bmatrix} x \\ 4 \times 1 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$$

$$z_1 = w_1 x + b_1$$

$a_1 = \sigma(z_1) \rightarrow$ elementwise operation

$\sigma = \text{Relu} \mid \text{Sigmoid}$

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

$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$



$$\text{III } y \quad z_2 = w_2 a_1 + b_2 \quad \rightarrow [2 \times 2]$$

$$a_2 = \sigma(z_2)$$

$$\text{Finally } 0 = z_3' = w_3 a_2 + b_3 \quad \rightarrow [1 \times 2]$$

Loss = mean squared error

$$L = \frac{1}{n} \sum_{i=1}^n (y - o)^2 \quad \text{For Batched case}$$

$$L = (y - o)^2 \quad \text{For single example}$$

Cross Ent

$$L = y \ln o + (1-y) \ln (1-o)$$

Now let us derive Back prop rules.

Goal

$$\text{find } \frac{\partial L}{\partial w_3}, \frac{\partial L}{\partial b_3}, \frac{\partial L}{\partial w_2}, \frac{\partial L}{\partial b_2}, \frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial b_1}$$

and update rule is

$$\text{Parameters of the NN} = \Theta = \left\{ w_i, b_i \right\}_{i=1}^{\# \text{ Layers}}$$

$$\Theta_{\text{new}} = \Theta_{\text{old}} - \eta \nabla \Theta_{\text{old}}$$

what does chain rule say?

assume f and g are continuous & differentiable
(need not be, we may still use sub gradients.) — out of scope.

$$h(\theta) = f(g(\theta))$$

$$h'(\theta) = f'(g(\theta)) * g'(\theta) * 1$$

break back at fwd pass

$$x \rightarrow \underbrace{z_1 = w_1 x + b_1}_{f_1} \rightarrow \underbrace{a_1 = \sigma(z_1)}_{f_2} \rightarrow \underbrace{z_2 = w_2 a_1 + b_2}_{f_3} \downarrow$$

$$L = (y - 0)^2 \leftarrow \underbrace{0 = w_3 a_2 + b_3}_{f_5} \leftarrow \underbrace{a_2 = \sigma(z_2)}_{f_4}$$

If we interpret $L = (NN(x) - 0)^2$

$$NN(x) = f_5 (f_4 (f_3 (f_2 (f_1 (x)))))$$

\downarrow \downarrow \downarrow
 $w_3 b_3$ $w_2 b_2$ $w_1 b_1$

f_5 f_3 and f_1 are linear

$$\Rightarrow \frac{\partial f_5}{\partial w_3} = \frac{\partial 0}{\partial w_3} = \overrightarrow{a_2}$$

$$\frac{\partial f_5}{\partial b_3} = \frac{\partial 0}{\partial b_3} = \overrightarrow{1}$$

111^{ly}

$$\frac{\partial f_3}{\partial w_2} = a_i ? \quad \text{wrong}$$

$$f_3 = \begin{bmatrix} w_{11}^2 & w_{12}^2 \\ w_{21}^2 & w_{22}^2 \end{bmatrix} \begin{bmatrix} a_1^1 \\ a_1^2 \end{bmatrix}$$

$$= \begin{bmatrix} w_{11}^2 a_1^1 & w_{12}^2 a_1^2 \\ w_{21}^2 a_1^1 & w_{22}^2 a_1^2 \end{bmatrix}$$

$$\frac{\partial f}{\partial w_2} = \begin{bmatrix} a_1^1 & a_1^2 \\ a_1^1 & a_1^2 \end{bmatrix}$$

derivative of activation function

$$\text{If } \sigma = \text{relu} \quad \sigma' = \begin{cases} 0 & \text{If } \text{Ip} \leq 0 \\ 1 & \text{If } \text{Ip} > 0 \end{cases}$$

$$\text{If } \sigma = \text{sigmoid} \quad \sigma' = \sigma(1-\sigma)$$

derivative with loss

$$L = y \ln o + (1-y) \ln (1-o)$$

$$L' = y * \frac{1}{o} + (1-y) \frac{1}{1-o}$$

In this derivation let us use MSE Error

Recall

$$L = \frac{1}{2} (y-o)^2$$

$$L' = (y-o)(-1)$$

$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial o} * \frac{\partial o}{\partial w_3} = -(y-o) * a_2$$

$$\frac{\partial L}{\partial b_3} = \frac{\partial L}{\partial o} * \frac{\partial o}{\partial b_3} = -(y-o) * \vec{1}$$

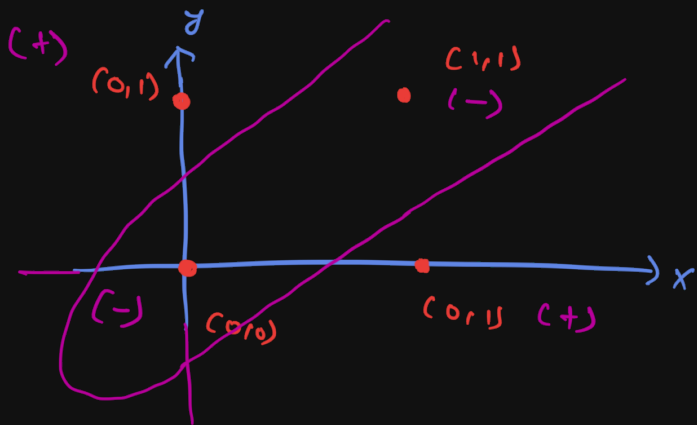
$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial o} * \frac{\partial o}{\partial a_2} * \frac{\partial a_2}{\partial z_2} * \frac{\partial z_2}{\partial w_2}$$

$$= -(y-o) * w_3^T * \sigma(z_2)(1-\sigma(z_2)) * \begin{bmatrix} a_1 \\ -a_1 \end{bmatrix}$$

you can keep unfolding this and immediately see an opportunity for "Dynamic programming"

H.W - Complete the derivation using 'DP' and implement the DP formula you get

classical XOR problem

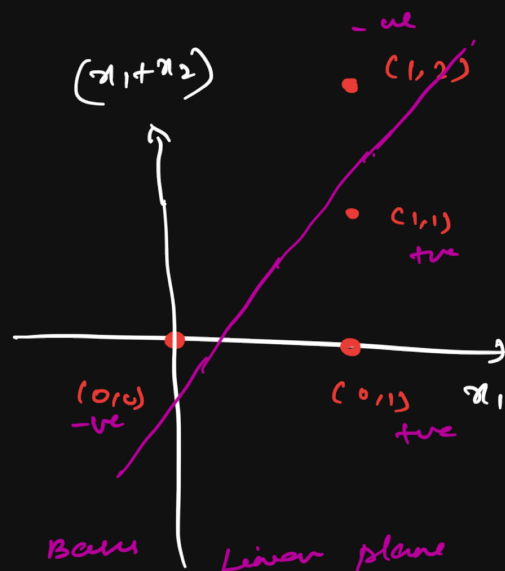


Any linear plane cannot separate the above simple 4 points dataset!

Logistic Regression / Perceptron will Fail

Solution - Basis Expansion

x_1	x_2	x_1+x_2	y
0	0	0	-ve
0	1	1	ve
1	0	1	ve
1	1	2	-ve



In new Basis Linear plane can separate

now Log Reg / Perceptron can win in Basis (x_1, x_1+x_2)

can a Human Hand craft such Basis / feature function $\begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}$

Beauty of Neural Networks

We utilize $NN(x)$ as such very complicated feature functions

NN learn them automatically through back propagation

Universal Approximation Theorem

(Vaguely) under certain conditions, a 3-layer NN

with p hidden nodes can approximate any function

→ But no one knows how to find 'p'

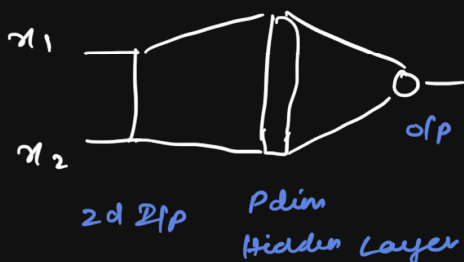
can $p \rightarrow \infty$?

H.W

Solve XOR problem using 2-layer Neural N/w, Logistic Regression and 3 layer Neural N/w and plot the Loss curves

Data = XOR table Label = XOR output

FN 3 Layered NN



We can interpret the activations of Hidden Layer as Basis function

For Each of the 4 data points take activations and train Logistic Regression model on activations.

Radial Basis Networks

① These are just 3 Layered N/Ws

known M I/P nodes

p Hidden nodes

1 o/p node.

② But we need the activations to be Radially Symmetric.

↳ Each neuron should account 'only' locally.

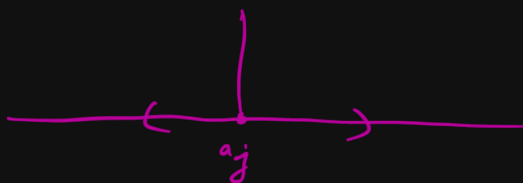
Earlier we had $\sigma(w^T a)$

and activation fixed for a long range on 'a'

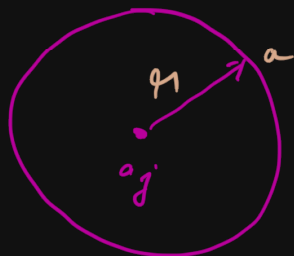
Now Each Neuron fixes a ' a_j ' and accounts

for only certain distance around ' a_j '

In 1-d



In 2-d



accounts for a radius r_1 around a_j

Hence the name Radially

Symmetric

The accountability of a neuron with center ' a_j ' decreases as distance of a from it increases

Let us call such an activation as ϕ

Then possibilities of ϕ are

$$\textcircled{1} \phi(a) = \exp\left(-\frac{1}{2} \frac{\|a - a_j\|^2}{\sigma^2}\right)$$

$$\text{if } \sigma = 1$$

$$a = a_j \Rightarrow \phi(a) = 1$$

$$a = a_j + \xi \Rightarrow \phi(a + \xi) = \exp\left(-\frac{1}{2} \|\xi\|^2\right) < 1$$

$$a = a_j - \xi \Rightarrow \phi(a - \xi) = \exp\left(-\frac{1}{2} \|\xi\|^2\right) < 1$$

Because $\phi(a_j + \xi) = \phi(a_j - \xi) \quad \forall \xi$

ϕ = radially symmetric

H.W → check relation b/w

$$\phi(a_j + \xi_1) \text{ and } \phi(a_j + \xi_2)$$

$$\text{if } \xi_1 > \xi_2$$

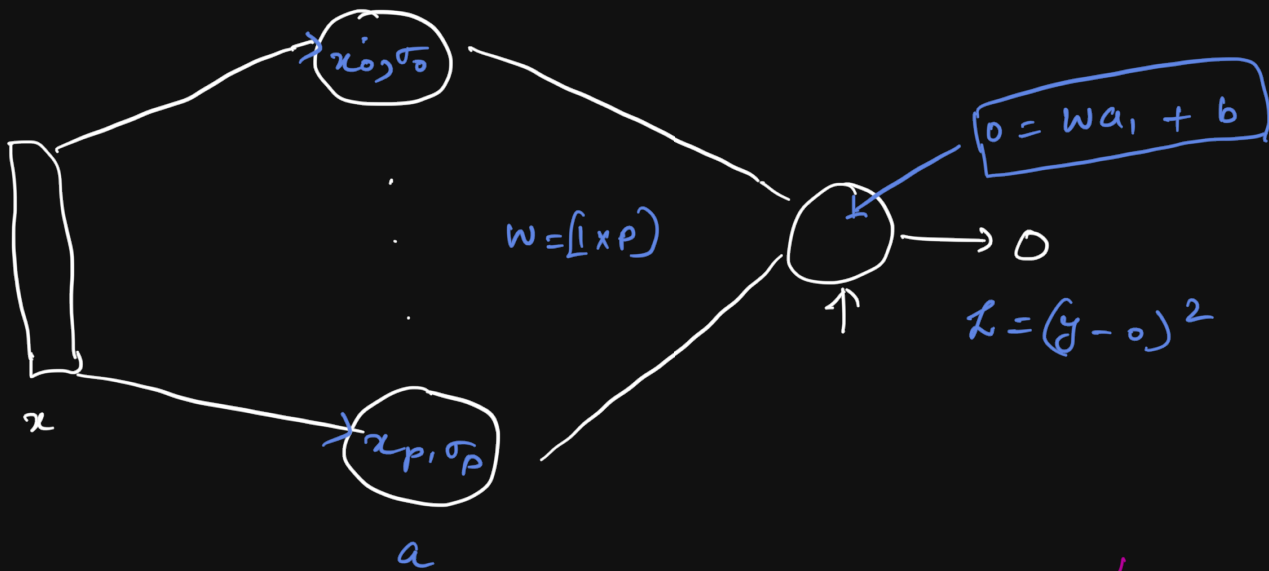
H.W 2

Repeat the above process for the following two choices of ϕ

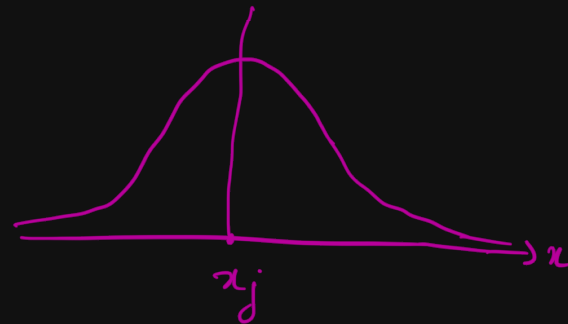
$$\textcircled{1} \phi(a) = \|a_j - a\|^2$$

$$\textcircled{2} \phi(a) = \frac{1}{\|a_j - a\|^2} \quad \forall a_j \neq a$$

RBF Network



$$a_j = \exp\left(-\frac{1}{2} \frac{\|x - x_j\|^2}{\sigma_j^2}\right)$$



Perfect Interpolation

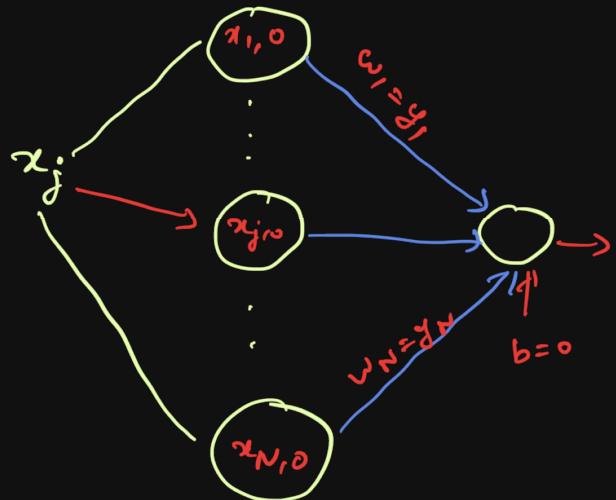
given $D = \left\{ (x_i, y_i) \right\}_{i=1}^N \exists$ RBF with Loss = 0

take $p = N$

$$x_j = x_i$$

$$\sigma_j = 0$$

and assume that $\frac{0}{0} = 0$



H/W \rightarrow check the above N/w products
Loss = 0.

You have already seen a solution for

How RBF N/ws train on XOR data

Caveat

↳ ML is not about getting '0' error on D (f)

we aspire to get '0' error on unseen test data

Final comments

① If x_j and σ_j are fixed then

RBF's are Linear classifiers

② How do we fix x_j and σ_j

one naive approach

* fix $\sigma_j = 1$ (blindly)

* fix $x_j =$ means of p means algorithm run on D

How to fix p ?

This is a question where many practitioners suffer to answer

Heuristic - Kmeans Elbow technique (out of scope)

Training RBF

Because it is only a Linear Network

$\frac{\partial L}{\partial w}$ and $\frac{\partial L}{\partial b}$ can be obtained at each (f)

Code walk through!

Side Note:

For RBF (or any Linear classifiers)

⇒ many cute optimization algorithms that work elegantly in practice other than gradient descent (SUM, LMS algo)

↳ out of scope