

When single layer perceptron is sufficient and when do we need hidden layers?

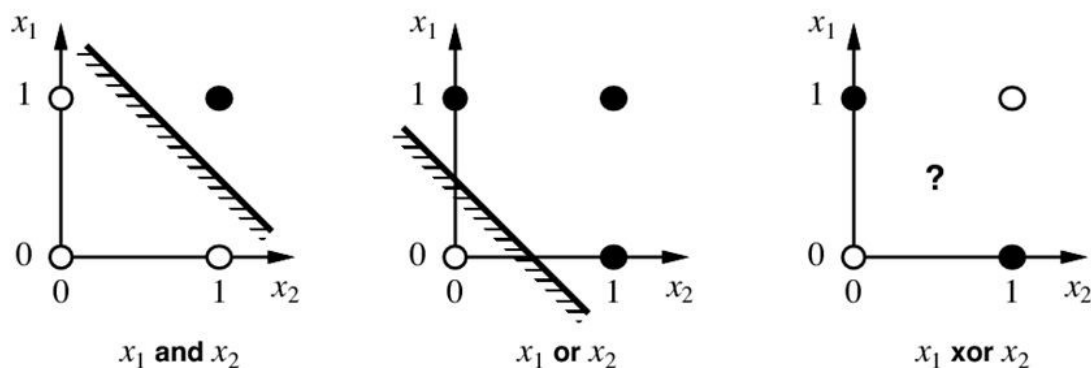
If the classes or patterns are linearly separable then single layer Perceptron is sufficient otherwise we need to incorporate hidden layers in the network to introduce non-linearity in the network.

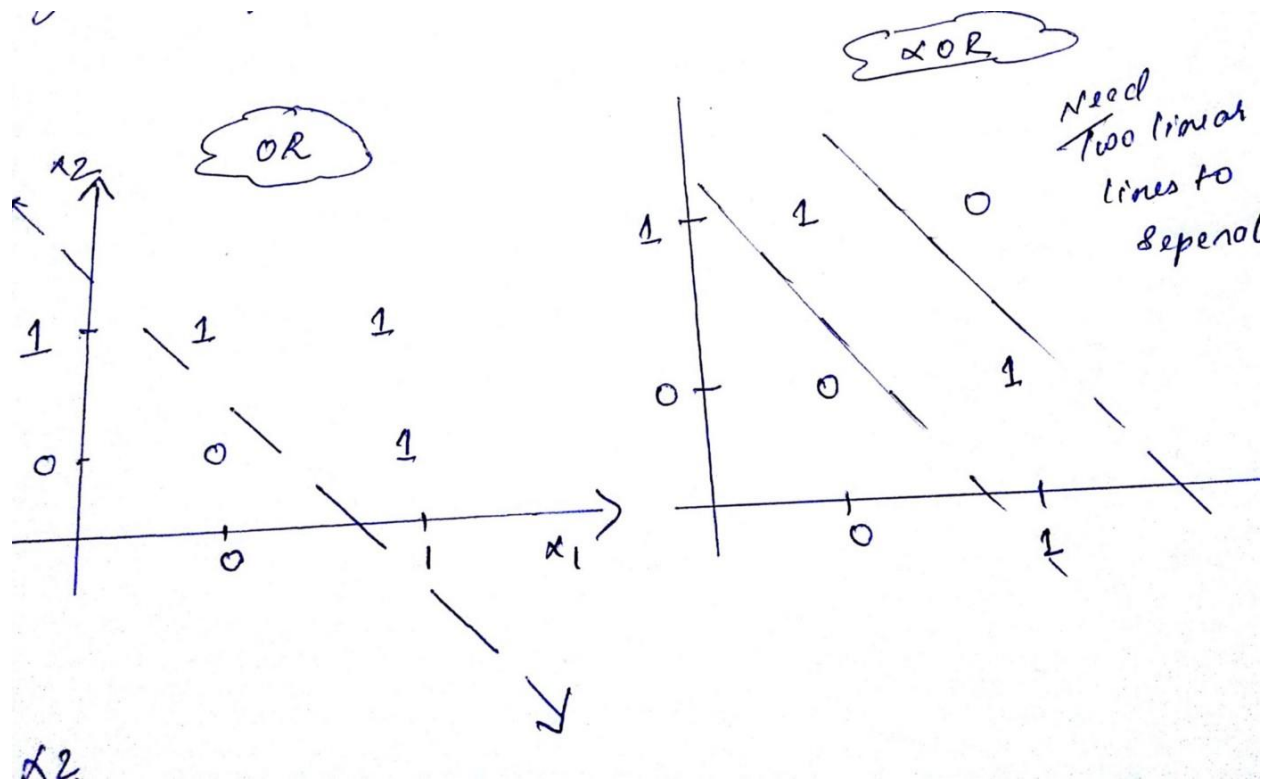
The hidden layer simply represents non-linear boundaries by a set of pairwise linear boundaries.

In general, there is always a possibility that the input feature vector which is not linearly separable on current dimensions can become linearly separable while projecting on higher dimensions.

Let's look at an example of OR and XOR gates output. We can see OR output can be linearly classified but the same is not true with XOR outputs as we need two linear lines to separate the boundaries of output (0 and 1).

Linear separability

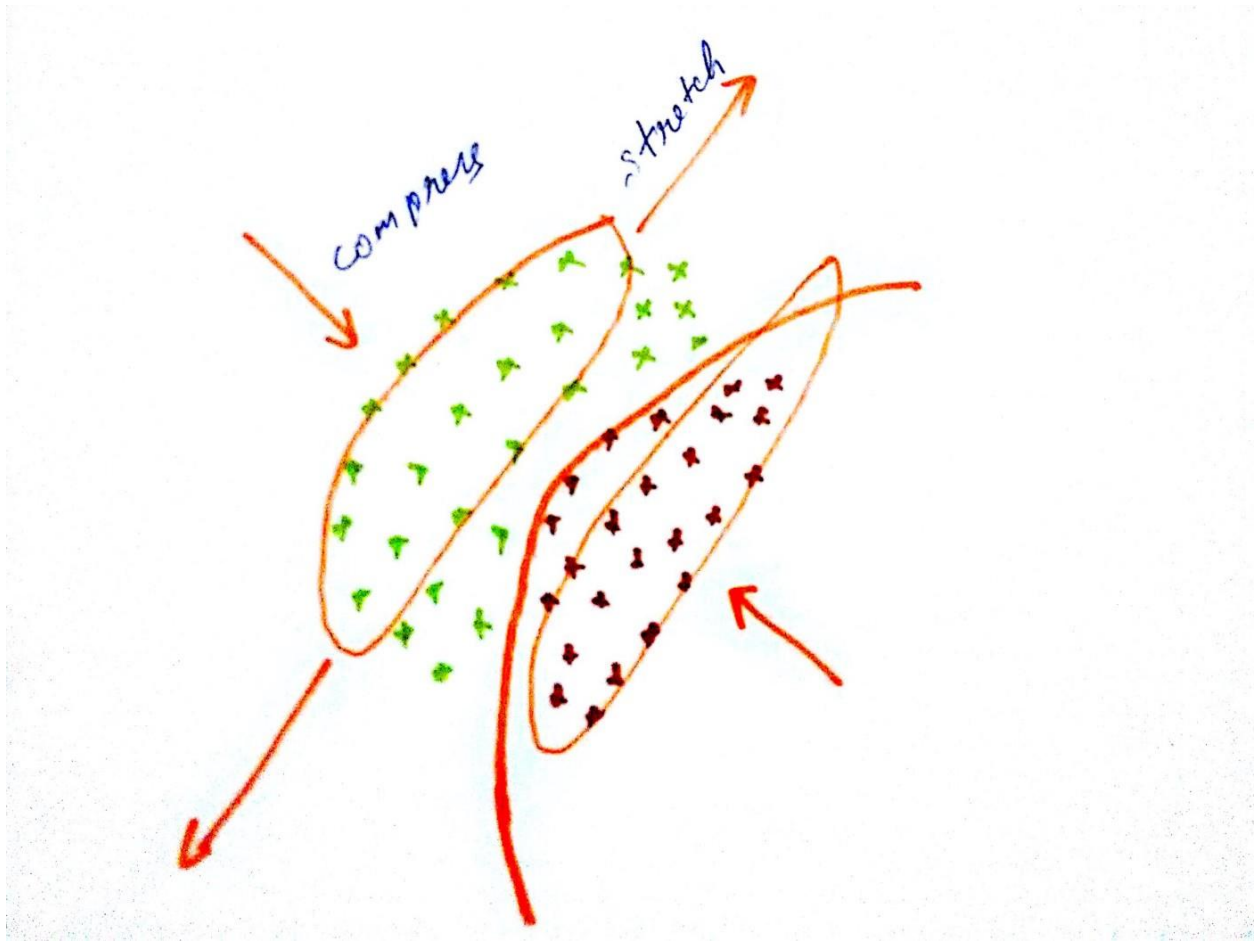




What is Radial Basis Function?

So coming to Radial Basis Function (RBF) what it does for our above problem of non-linear separable patterns. RBF performs nonlinear transformation over input vectors before they are fed for classification with help of the below transformations.

- Imposes non-linear transformation on input feature vector.
- Increases Dimensionality of the feature vector.



In above image green and red is the feature vector of two different classes which are quite evident that they are non linearly separable so once we impose RBF what it does that it performs a non-linear transformation on feature vectors which in effect compress them along the width and stretches them along length due to which it becomes linearly separable.

Along with the non-linear transformation, the dimensionality is increased up to the M dimension from the P original dimension with the below equations.

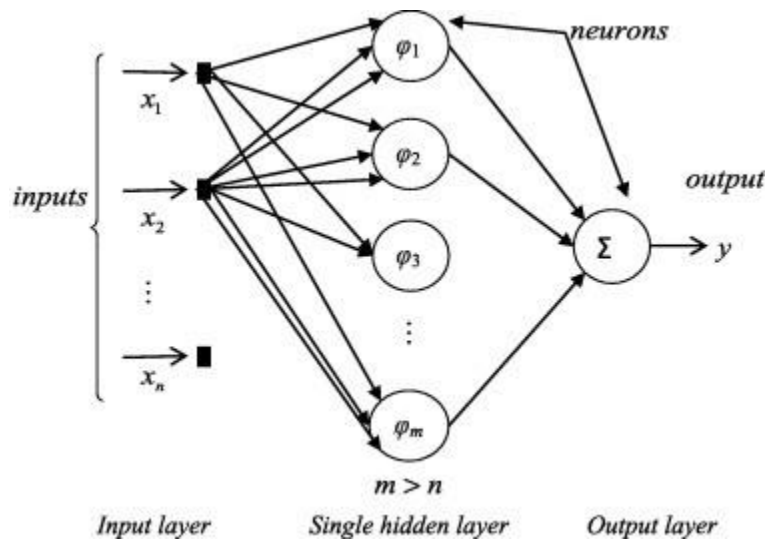
$$P \rightarrow M$$

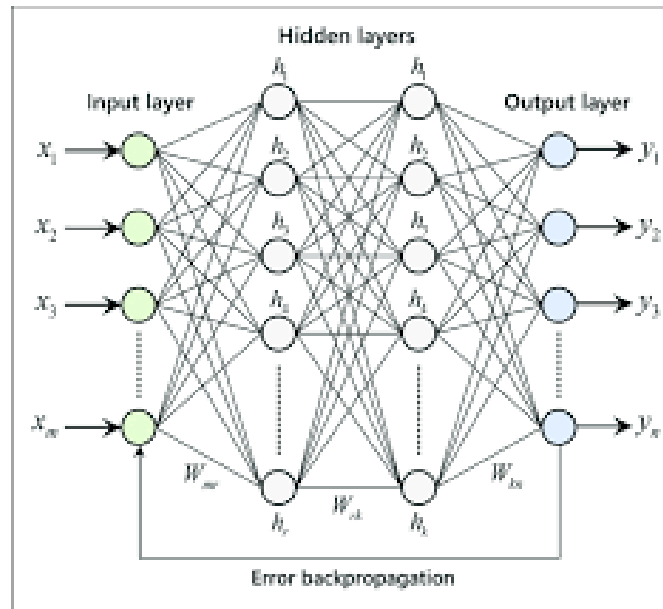
$$M > P$$

RBF transform P dim to M dim

$$\phi(x) = [\phi_1(x), \phi_2(x), \dots, \phi_m(x)]^T$$

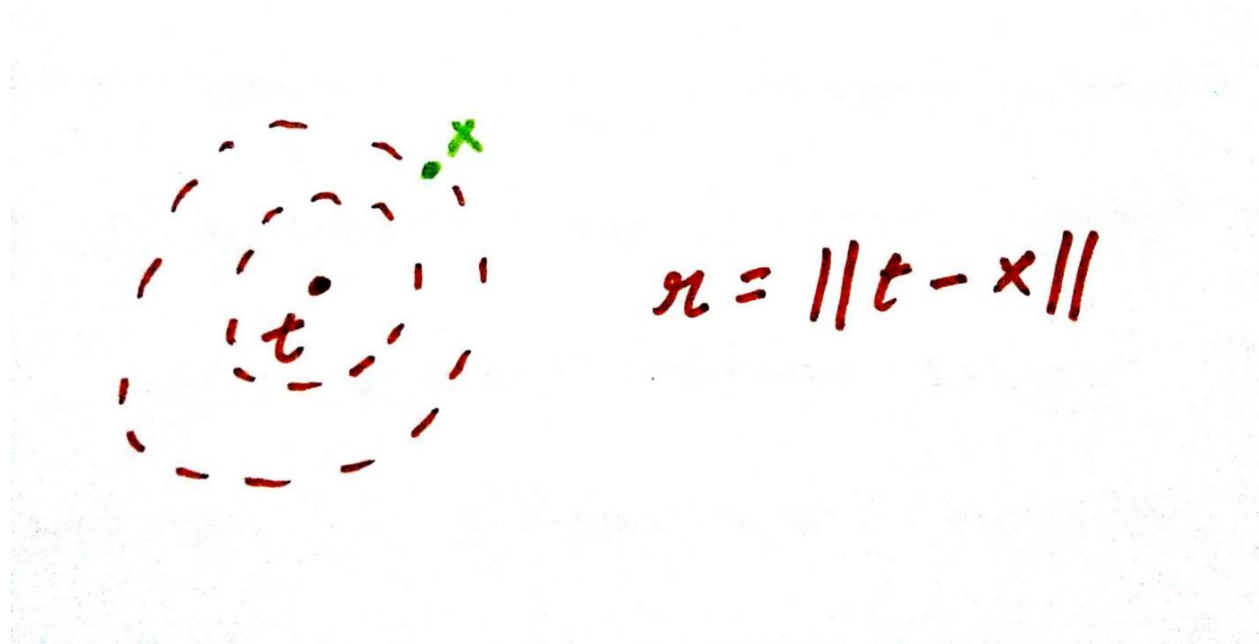
all values of ϕ are real



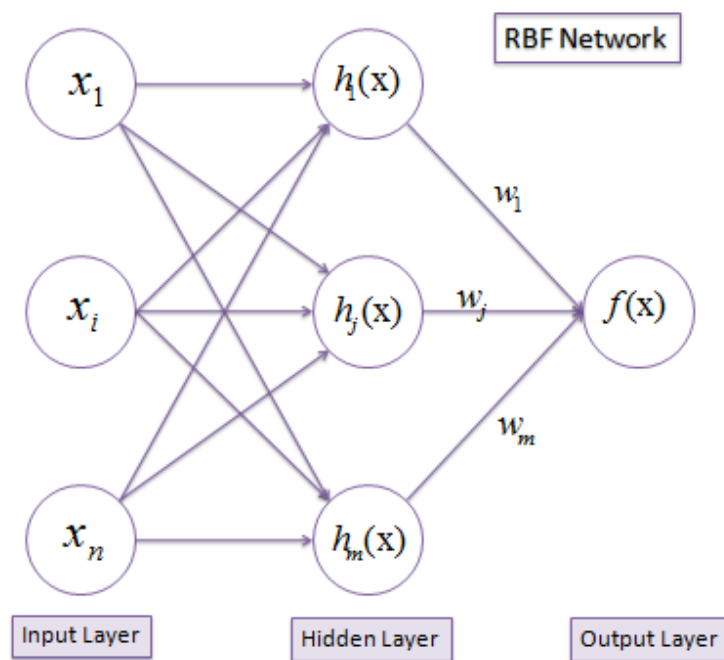


Explain the RBF behaviour.

Each RBF function has a receptor and values of the function either increases or decreases while moving away from receptor t .



r is the distance between receptor t and any input feature vector X



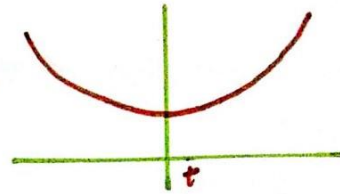
$$f(x) = \sum_{j=1}^m w_j h_j(x)$$

$$h(x) = \exp\left(-\frac{(x-c)^2}{r^2}\right)$$

Choices of RBF 1-

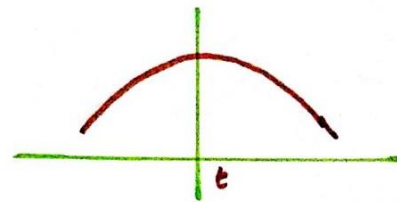
① MULTIQUADRICS:

$$\phi(r) = (r^2 + c^2)^{1/2} \quad c > 0$$



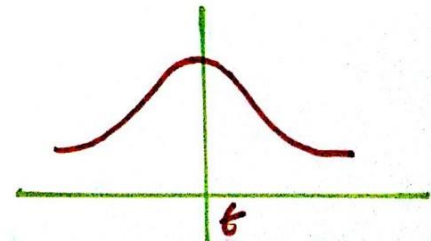
② INVERSE MULTIQUADRICS:

$$\phi(r) = \frac{1}{(r^2 + c^2)^{1/2}} \quad c > 0$$



③ GAUSSIAN FUNCTION:

$$\phi(r) = \exp\left[-\frac{r^2}{2\sigma^2}\right] \quad \sigma > 0$$



Commonly used types of radial basis functions include (writing $r = \|\mathbf{x} - \mathbf{x}_i\|$ and using ε to indicate a *shape parameter* that can be used to scale the input of the radial kernel^[11]):

- Infinitely Smooth RBFs

These radial basis functions are from $C^\infty(\mathbb{R})$ and are strictly **positive definite functions**^[12] that require tuning a shape parameter ε

- Gaussian:

$$\varphi(r) = e^{-(\varepsilon r)^2}, \quad (2)$$

- Multiquadric:

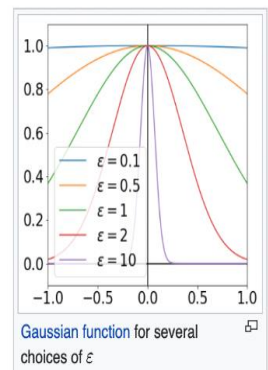
$$\varphi(r) = \sqrt{1 + (\varepsilon r)^2}, \quad (3)$$

- Inverse quadratic:

$$\varphi(r) = \frac{1}{1 + (\varepsilon r)^2}, \quad (4)$$

- Inverse multiquadric:

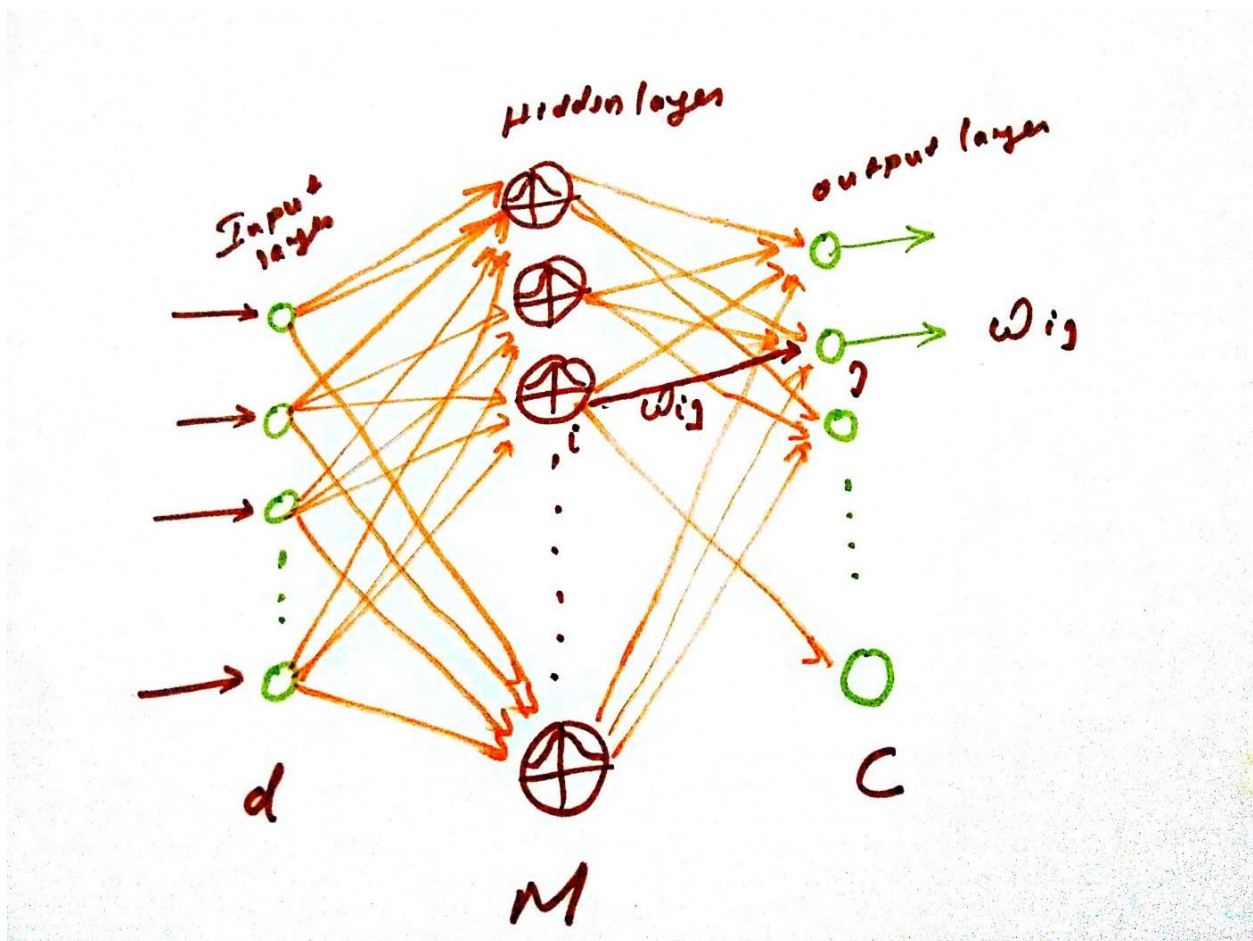
$$\varphi(r) = \frac{1}{\sqrt{1 + (\varepsilon r)^2}}, \quad (5)$$



Above are the general choice for RBF but gaussian is the most popular one, we can see the function decays while going away from the receptor.

The network architecture of RBF

Radial basis function (RBF) networks typically have three layers: an input layer, a hidden layer with a non-linear RBF activation function and a linear output layer. The input can be modelled as a vector of real numbers. The output of the network is then a scalar function of the input vector.



d is the dimensionality of input feature space, M is the dimensionality of transformed feature space where we have imposed our RBF, C is the number of classes to be identified.

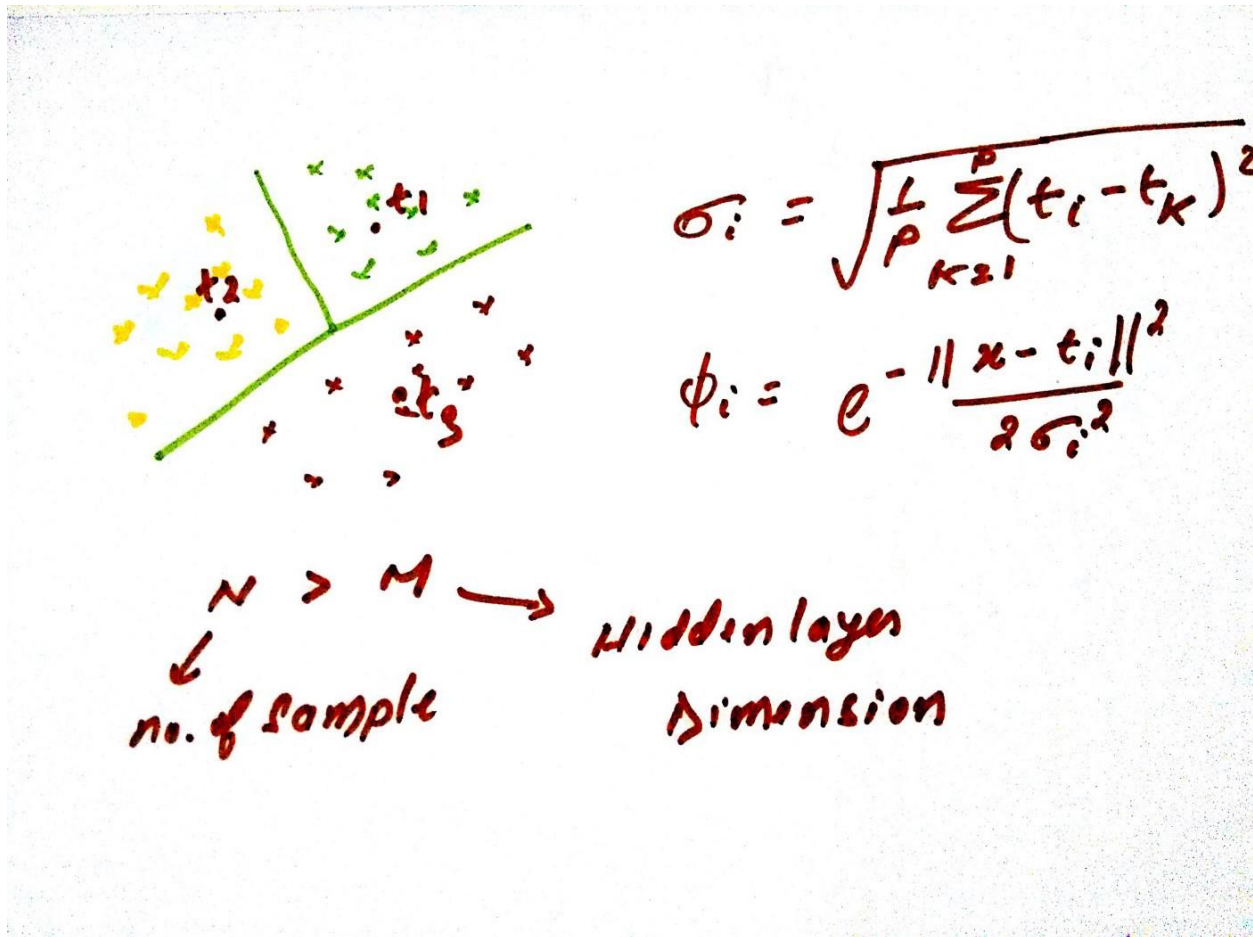
Training Comprises for this kind of network in two phases:

- a) Training Hidden layer which comprises of M RBF functions, the parameters to be determined for RBF function are receptor position t and the Sigma in case of Gaussian RBF.
- b) Training weight vectors W_{ij} for the output layer.

Training Hidden layer:

So for training hidden layers, there are different approaches, let us assume for now we are dealing with Gaussian RBF so we need to determine receptor t and Spread ie Sigma. One of the approaches is to randomly select the M number of the receptor from the N number of sample feature vectors but this does not seem logical so we can go ahead with a clustering mechanism to determine receptors t_i .

As we have M nodes in the hidden layer and N samples so for clustering to work here $N > M$.



Calculation of receptors:

Let's look at the above example where we have $M=3$ so we need to determine three t 's. so initially we divide out feature vector space into three arbitrary clusters and took their means as the initial receptors, then we need to iterate for every sample feature vector and perform the below steps:

- From the selected input feature vector x determine distances of means (t_1, t_2, t_3) of three different clusters whichever distance mean is minimum sample x will get assigned to that cluster.
- After x got assigned to different cluster all the means (t_1, t_2, t_3) gets recomputed
- Perform step 1 and step 2 for all sample points

Once the iteration finishes we will get the optimal t_1, t_2 and t_3 .

Calculation of Sigma:

Training Weight Vectors

$x_k: k = 1, 2, 3, \dots, N$

$\phi_i(x_k) \rightarrow \phi_{ik} \quad \sum_{i=1}^M \omega_{ij} \phi_{ik} = +1 \text{ if } x_k \in \omega_j$
 $0 \text{ if } x_k \notin \omega_j$

Matrix Eqn:

$$\begin{bmatrix} \phi_{11} & \phi_{21} & \dots & \phi_{M1} \\ \phi_{12} & \phi_{22} & \dots & \phi_{M2} \\ \vdots & \vdots & \ddots & \vdots \\ \phi_{1N} & \phi_{2N} & \dots & \phi_{MN} \end{bmatrix} \begin{bmatrix} \omega_{1j} \\ \omega_{2j} \\ \vdots \\ \omega_{Mj} \end{bmatrix} = \begin{bmatrix} b_{1j} \\ b_{2j} \\ \vdots \\ b_{Nj} \end{bmatrix}$$

$b_{ij} = 1 \text{ if } x_i \in \omega_j$
 $0 \text{ if } x_i \notin \omega_j$

$\phi \omega_j = b_j$

Let us assume we don't have a perfect solution so we will take the error e and try to optimize our criteria function $J(w_j)$ using a closed-form solution to get our optimum W_j .

$$E = \frac{1}{2} * \sum (t - o)^2 \quad (\text{error})$$

$$\Delta w_{jk} = -1 * \underset{\substack{\uparrow \\ \text{learn rate}}}{\eta} * \left[h_j * \underbrace{(o_k - t_k)}_{\substack{\text{derivative} \\ \text{of error}}} * \underbrace{o_k * (1 - o_k)}_{\substack{\text{derivative of} \\ \text{softmax}}} \right] \quad \frac{\partial E}{\partial w_{jk}} \quad (\text{gradient})$$

The above W_j calculation can be done for every output node $j = 1, 2, \dots, C$ to get the respective weight vectors for them.

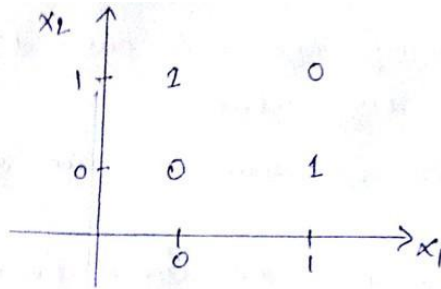
Advantages of RBF

- 1) In Comparison, to MultiLayer Perceptron the training phase is faster as there is no backpropagation learning involved.
- 2) The interpretation roles of Hidden layer nodes are easy in comparison to multi-layer perceptron.
- 3) The number of hidden layers and the number of nodes in the hidden layer are decided in the case of the RBF network but in multi-layer perceptron, there is no analytical approach to decide the number of nodes in the hidden layer or number of the hidden layers.

Disadvantages of RBF

Although the training is faster in the RBF network the classification is slow in comparison to Multilayer Perceptron due to fact that every node in the hidden layer has to compute the RBF function for the input sample vector during classification.

Example: Solution for XOR gate using RBF Network



$P=2$

Radial basis function

$$\phi_1 \rightarrow t_1 = (0,0) ; \sigma_1 = 1$$

$$\phi_1(x) = e^{-\frac{\|x-t_1\|^2}{2}}$$

$$\phi_2 \rightarrow t_2 = (0,1) ; \sigma_2 = 1$$

$$\phi_2(x) = e^{-\frac{\|x-t_2\|^2}{2}}$$

$$\phi_3 \rightarrow t_3 = (1,0) ; \sigma_3 = 1$$

$$\phi_3(x) = e^{-\frac{\|x-t_3\|^2}{2}}$$

$$\phi_4 \rightarrow t_4 = (1,1) ; \sigma_4 = 1$$

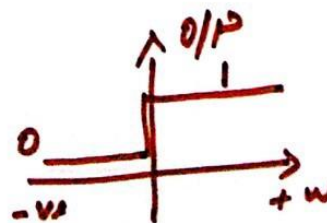
$$\phi_4(x) = e^{-\frac{\|x-t_4\|^2}{2}}$$

* converting 2 dim to 4 dim

Input	ϕ_1	ϕ_2	ϕ_3	ϕ_4	$\sum w_i \phi_i$	Output
0 0	1.0	0.6	0.6	0.4	-0.2	0
0 1	0.6	1.0	0.4	0.6	0.2	1
1 0	0.6	0.4	1.0	0.6	0.2	1
1 1	0.4	0.6	0.6	1.0	-0.2	0

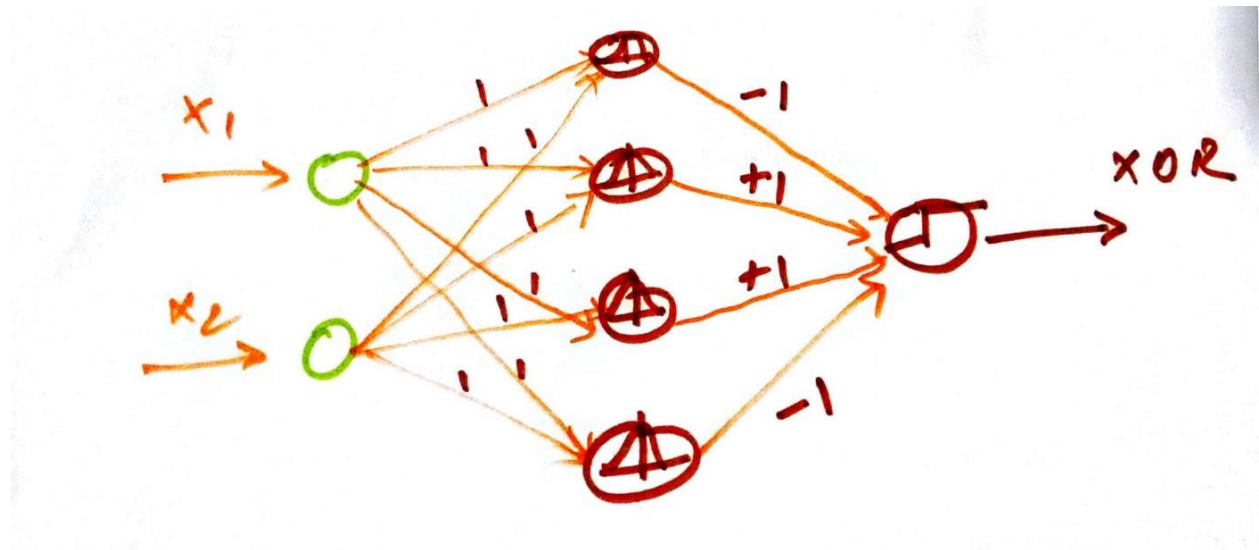
Weights $\rightarrow -1 \quad +1 \quad +1 \quad -1$

Hard threshold



Here optimal weights are -1,+1,+1 and -1 and at the output, we have enforced hard thresholding function as if the value is -ve o/p is 0 if it's +ve o/p is 1

Architecture diagram for XOR



References

https://www.youtube.com/watch?v=KiVJkqac82Q&t=1891s&ab_channel=nptelhrd
https://www.youtube.com/watch?v=nOt_V7ndmLE&t=1669s&ab_channel=nptelhrd
<https://medium.com/geekculture/radial-basis-function-network-952785127d61>

Implementation

<https://www.kaggle.com/residentmario/radial-basis-networks-and-custom-keras-layers/notebook>
<https://towardsdatascience.com/most-effective-way-to-implement-radial-basis-function-neural-network-for-classification-problem-33c467803319>