

//_

SECTION - B UNIT - I

ASSOCIATIVE MEMORY

An associative memory is a brain like distributed memory that learns by association. Association takes one form \rightarrow auto ^{or} ~~and~~ hetero.

- In auto-association, a NN is required to store a set of patterns by repeatedly presenting them to network. The network is subsequently presented with a partial description or distorted version of an original pattern stored in it. and the task is to retrieve or recall that particular pattern.
- Hetero-association differs in way \rightarrow an arbitrary set of input patterns is paired with another arbitrary set of output patterns.

* Auto-association uses unsupervised learning.

* Hetero- " " supervised " .

Let x_k denote a key pattern applied to an associative memory & y_k denote a memorised pattern.

The pattern association performed by the network is described as \rightarrow

$$x_k \rightarrow y_k \quad k=1, 2, 3, \dots, q$$

where q is no. of input patterns

In auto-associative memory, the input & output spaces have same dimensionality but in hetero-associative memory, it may or may not be the case.

$$y_k = x_k$$

$$y_k \neq x_k$$

OPERATION

2 Phases

1. The associative storage phase: where training happens in accordance with the mapping.

$$\underline{x_k \rightarrow y_k} \quad k=1, 2, 3, \dots, q$$

2. The recall / retrieval phase: involves retrieval of a memorised pattern.

In auto-associative

Let x represent a noisy / distorted version of x_j , consider response y is produced. For perfect recall, $y = y_j$ but if it is not then memory is said to have made an error in recall.

CAPACITY (STORAGE) OF NETWORK

STORAGE PHASE

$$x_k = [x_{k1}, x_{k2}, \dots, x_{km}]^T$$

$$y_k = [y_{k1}, y_{k2}, \dots, y_{km}]^T$$

$$y_k = W(k) x_k \quad k=1, 2, 3, \dots, q$$

$$y_{ki} = \sum_{j=1}^m w_{ij}(k) x_{kj} \quad i=1, 2, \dots, m$$

$$y_{ki} = \begin{bmatrix} w_{i1}(k) & w_{i2}(k) & \dots & w_{im}(k) \end{bmatrix} \begin{bmatrix} x_{k1} \\ x_{k2} \\ \vdots \\ x_{km} \end{bmatrix}$$

$$\begin{bmatrix} y_{k1} \\ y_{k2} \\ \vdots \\ y_{km} \end{bmatrix} = \begin{bmatrix} w_{11}(k) & w_{12}(k) & \dots & w_{1m}(k) \\ w_{21}(k) & w_{22}(k) & \dots & w_{2m}(k) \\ \vdots & \vdots & \ddots & \vdots \\ w_{m1}(k) & w_{m2}(k) & \dots & w_{mm}(k) \end{bmatrix} \begin{bmatrix} x_{k1} \\ x_{k2} \\ \vdots \\ x_{km} \end{bmatrix}$$

$\underbrace{\hspace{15em}}_{W(k)} \quad m \times m$

$$M = \sum_{k=1}^q W(k)$$

$$M_k = M_{k-1} + W(k)$$

$$\hat{M} = \sum_{k=1}^q y_k x_k^T$$

$$\hat{M} = \begin{bmatrix} y_1 & y_2 & \dots & y_q \end{bmatrix} \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_q^T \end{bmatrix}$$

$$\hat{M} = Y X^T$$

→ correlation matrix memory using outer product rule (Hebbian learning)

$$\hat{M}_k = \hat{M}_{k-1} + y_k x_k^T$$

$k=1, 2, 3, \dots, q$

$$X = \begin{bmatrix} x_1 & x_2 & x_3 & \dots & x_q \end{bmatrix}$$

$$Y = \begin{bmatrix} y_1 & y_2 & y_3 & \dots & y_q \end{bmatrix}$$

RECALL PHASE

$$y = \hat{M} x_j$$

$$y = \sum_{k=1}^q y_k x_k^T x_j = \sum_{k=1}^q (x_k^T x_j) y_k$$

$$y = \underbrace{(x_j^T x_j)}_{=1} y_j + \sum_{\substack{k=1 \\ k \neq j}}^q (x_k^T x_j) y_k$$

$$y = y_j + V_j$$

$$E_k = \sum_{l=1}^m x_{kl}^2$$

$$E_k = x_k^T x = 1 \quad \text{where } k=1, 2, 3, \dots, q$$

$$\cos(x_k, x_j) = \frac{x_k^T x_j}{\|x_k\| \|x_j\|}$$

$$\|x_k\| = (x_k^T x_k)^{1/2} = E_k^{1/2}$$

$$\cos(x_k, x_j) = x_k^T x_j$$

$$V_j = \sum_{\substack{k=1 \\ k \neq j}}^m \cos(x_k, x_j) y_k$$

CONDITION FOR PERFECT RECALL $\rightarrow \cos(x_k, x_j) = 0$
i.e. x_k and x_j are orthogonal

RADIAL - BASIS FUNCTION NETWORKS (RBF networks)

In RBF networks, we take a completely different approach by viewing the design of a neural network as a curve fitting problem in a high-dimensional space. Acc. to this new point, learning is equivalent to finding a surface in a multi-dimensional space that provides a best fit to the training data with the criteria for best fit being measured in some statistical sense.

Correspondingly, generalization is equivalent to the use of this multi-dimensional surface to interpolate the test data. Such a view point is the motivation behind RBF networks.

In the context of a NN, the hidden units provide a set of functions that constitute an arbitrary basis for the input patterns when they are expanded into the hidden space. These are called radial-basis functions.

The construction of a RBF network in its most basic form involves 3 layers.

- INPUT LAYER \rightarrow Made of source ~~as~~ nodes that connect network to its environment.
- HIDDEN LAYER \rightarrow Applies non-linear transformation from input space to hidden space. In most applications, the hidden space is of high dimensionality.
- OUTPUT LAYER \rightarrow It is linear that supplies response of network to activation pattern applied to input layer.

* The mathematical justification for the rational of a non-linear transformation followed by a linear transformation is due to COVER'S THEOREM.

It states a pattern classification problem cast in a high dimensional space is more likely to be linear separable than in low dimensional space.

MATHEMATICAL MODEL

Consider a family of surfaces where each naturally divides an input space into 2 regions.

Let X denote a set of N patterns.

$$x_1, x_2, \dots, x_N \quad \checkmark$$

each of which is assigned to one to the two classes X_1 or X_2 . \checkmark

This binary partition of the points is said to be separable wrt family of surfaces if a surface exists that separates the points in class X_1 from those in class X_2 .

For each pattern, $x \in X$, define a vector made up of a set of real valued functions \rightarrow

$$\{\phi_i(x) \mid i=1, 2, \dots, m_1\}$$

$$\Phi(x) = [\phi_1(x), \phi_2(x), \dots, \phi_{m_1}(x)]$$

Suppose that pattern x is a vector in m_0 dimensional input space. The vector $\Phi(x)$ then maps points in m_0 dim. input space into corresponding points in a new space of dimension m_1 .

//_

We refer to $\Phi_i(x)$ as hidden function.

Correspondingly, the space spanned by the set of hidden functions ^{is} referred to as hidden space / feature space.

The binary partition X_1, X_2 of X is said to be Φ -separable if there exists a m_2 dimensional vector w such that \rightarrow

$$w^T \Phi(x) > 0 \quad x \in X_1$$

$$w^T \Phi(x) < 0 \quad x \in X_2$$

Hyperplane / separating surface $\rightarrow w^T \Phi(x) = 0$
in hidden space.

Consider a natural class of mappings obtained by using a linear combination of x -wise products of the pattern vector coordinates. The separating surfaces corresponding to such mappings are referred to as x^{th} order rational varieties.

A rational variety of order x in a space of dim. m_0 is described by an x^{th} degree homogenous eqⁿ in the coordinates of input vector x given by \rightarrow

$$\sum_{0 \leq i_1 \leq i_2 \leq \dots \leq i_x \leq m_0} a_{i_1 i_2 \dots i_x} x_{i_1} x_{i_2} \dots x_{i_x} = 0$$

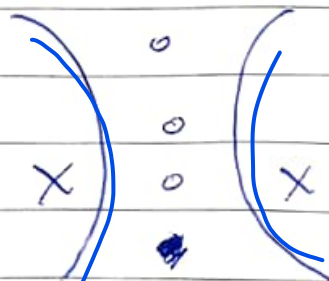
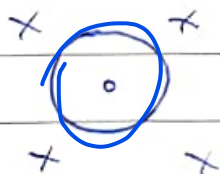
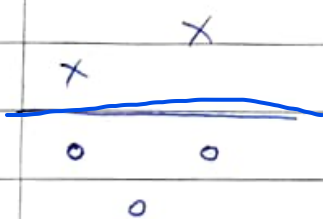
where x_i is the i^{th} component of input vector x .

The x^{th} order product of entries x_i is called a monomial.

For an input space of dim. m_0 , there are $\frac{(m_0 + x)!}{m_0! x!}$ monomials.

The separating surfaces described by above eqⁿ, are \rightarrow

- Hyperplanes (1st order)
- Quadrics (2nd order)
- Hyperspheres (quadrics with certain linear constraints on coefficients)



Linearly separable | Spherically separable | Quadratically separable

Let m_0 denote the dim. of input space, then in the overall fashion, the network represents a map from m_0 dimensional input space to a single dim. output space given as: $S: R^{m_0} \rightarrow R^1$

We can think of this map as a hypersurface graph which is a subset of (m_0+1) dim. space. \rightarrow

$$\Gamma \subset R^{m_0+1}$$

The training and generalisation phase of the learning process can be viewed as \rightarrow

1. The training phase constitutes the optimization of a fitting ^{procedure for} surface Γ based on ^{known} ~~real~~ data points presented to network in form of input/output patterns.
2. The generalisation phase is synonymous with interpolation b/w the data points, with the interpolation being performed along the constrained surface generated by the fitting procedure as the optimum approximation to the true surface Γ .

Now we are having a multi-variable interpolation problem in high dim. space. It can be stated as \rightarrow

INTERPOLATION PROBLEM.

Given a set of N different points $\{x_i \in \mathbb{R}^{m_0} \mid i=1,2,\dots,N\}$ and a corresponding set of N real numbers $\{d_i \in \mathbb{R}^1 \mid i=1,2,\dots,N\}$. Find a function $F: \mathbb{R}^N \rightarrow \mathbb{R}^1$ which satisfies interpolation condition $F(x_i) = d_i$ where $i=1,2,\dots,N$. For strict interpolation, the interpolating surface is constrained to pass through all the training data points.

The RBF ~~technique~~ techniques consist of choosing a function that is of form \rightarrow

$$F(x) = \sum_{i=1}^N w_i \phi(\|x - x_i\|)$$

$\phi(\|x - x_i\|)$ is set of N arbitrary RBFs.
 $x_i \in \mathbb{R}^{m_0}$ are centers to the RBFs

ϕ_{11}	ϕ_{12}	\dots	ϕ_{1N}	w_1	d_1
ϕ_{21}	ϕ_{22}	\dots	ϕ_{2N}	w_2	d_2
\vdots	\vdots		\vdots	\vdots	\vdots
\vdots	\vdots		\vdots	\vdots	\vdots
ϕ_{N1}	ϕ_{N2}	\dots	ϕ_{NN}	w_N	d_N

$=$

$\Phi W = X$
 $W = \Phi^{-1}(X)$

There is large class of RBFs that is covered by
 Micchelli's theorem and considers given forms \rightarrow

1. Multiquadric

$$\Phi(x) = \sqrt{x^2 + c^2}$$

2. Inverse multiquadric

$$\Phi(x) = \frac{1}{\sqrt{x^2 + c^2}}$$

3. Gaussian

$$\Phi(x) = e^{\left(\frac{-x^2}{2\sigma^2}\right)}$$

• Solution to XOR using RBF

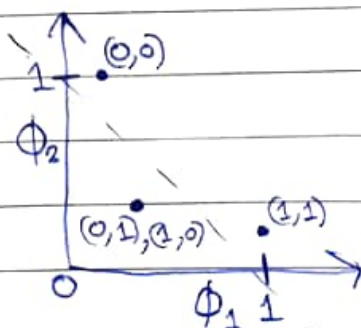
$$\Phi_1(x) = e^{\|x - t_1\|^2}$$

$$t_1 = \begin{bmatrix} 1 & 1 \end{bmatrix}^T$$

$$\Phi_2(x) = e^{\|x - t_2\|^2}$$

$$t_2 = \begin{bmatrix} 0 & 0 \end{bmatrix}^T$$

x	$\Phi_1(x)$	$\Phi_2(x)$
(0, 0)	0.1353	1
(0, 1)	0.3678	0.3678
(1, 0)	0.3678	0.3678
(1, 1)	1	0.1353



Supervised learning is illposed hyperspace reconstruction.

*

RBF

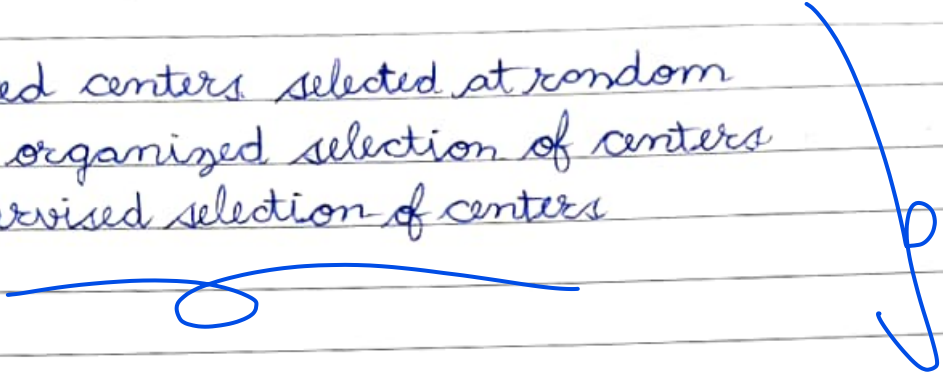
MLP

- Single hidden layer
- Computational nodes

- 1 or more hidden layer

- Computational nodes of MLP share a common model. For RBF, these ^{computational} nodes are quite different from others and serve a different purpose.
- The argument of activation function of each hidden unit in RBF network computes Euclidean norm whereas activation fn of hidden unit in MLP computes the inner product.
- MLPs construct global approximations and RBFs using exponential decaying localised non linear functions construct local approximations.

• Learning strategies in RBF

1. Fixed centers selected at random
 2. Self organized selection of centers
 3. Supervised selection of centers
- 

- * Global order can arise from local interactions.
- * Modifications in synaptic weights tend to self-amplify.
- * Limitation of resources leads to competition among synapses and thus the selection of most fittest at the expense of others.
- * Modifications in synaptic weights tend to co-operate.
- * Order and structure in the activation patterns represent redundant information that is acquired by network in form of knowledge.

Basic philosophy of PCA

Let X denote a m -dimensional vector representing

We assume the vector X to have 0 mean.

$$E[X] = 0$$

Let q denote unit vector also of dimension m onto which vector X is to be projected. This projection is defined by eqⁿ \rightarrow

$$A = X^T q = q^T X$$

$$\|q\| = (q^T q)^{1/2} = 1$$

$$E[A] = q^T E[X] = 0 \quad \checkmark$$

$$\sigma^2 = E[A^2]$$

$$\sigma^2 = E[(q^T X)(X^T q)]$$

$$\sigma^2 = q^T E[XX^T] q$$

$$\boxed{\sigma^2 = q^T R q}$$

where $R = E[XX^T]$

\rightarrow correlation matrix

$$\boxed{\sigma^2 = \phi(q)}$$

$$\boxed{Rq = \lambda q}$$

λ is the scaling factor representing the eigen value of R

Matrix form \rightarrow

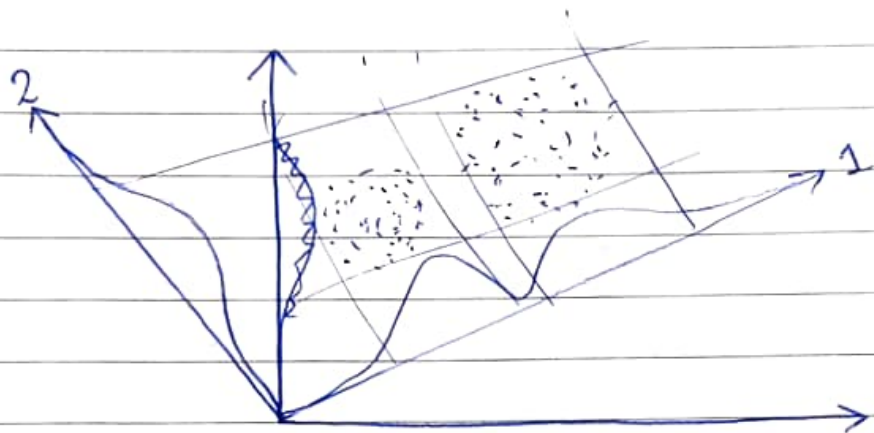
$$Q^T R Q = \Lambda$$

$$R = \sum_{i=1}^m \lambda_i q_i q_i^T$$

The important findings from mathematical analysis reveal \rightarrow

1. Eigenvectors of correlation matrix ~~R~~ pertaining to zero mean random vector X , define the unit vectors q_i representing principal directions along which $\phi(q)$ have the extremal values. ~~The above~~

2. The associated eigen values define the extremal values of $\Phi(q)$.



With m possible solutions for unit vector q , we find that there are m possible projections of data. \rightarrow

$$a_j = q_j^T X = X^T q_j \quad : j=1, 2, \dots, m$$

$$a = [a_1 \ a_2 \ \dots \ a_m]^T$$

$$a = [X^T q_1 \ X^T q_2 \ \dots \ X^T q_m]^T$$

$$a = Q^T X$$

$$X = Qa$$

$$X = \sum_{j=1}^m a_j q_j$$

$$\hat{X} = \sum_{j=1}^k a_j q_j$$

$$\hat{X} = [q_1 \ q_2 \ \dots \ q_l] \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_l \end{bmatrix}$$

$$l \leq m$$

$$\text{so, } \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_l \end{bmatrix} = \begin{bmatrix} q_1^T \\ q_2^T \\ \vdots \\ q_l^T \end{bmatrix} X$$

~~ENCODER~~

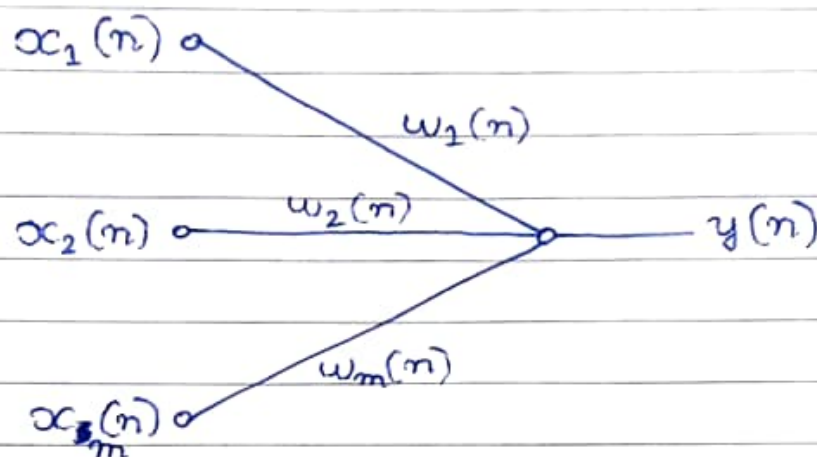
ENCODER

$$\begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_m \end{bmatrix} \Rightarrow \begin{bmatrix} q_1^T \\ q_2^T \\ \vdots \\ q_l^T \end{bmatrix} \Rightarrow \begin{bmatrix} \cancel{q_1^T} \\ \cancel{q_2^T} \\ \vdots \\ \cancel{q_l^T} \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_l \end{bmatrix}$$

$$\begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_l \end{bmatrix} \begin{bmatrix} \cancel{q_1^T} \\ \cancel{q_2^T} \\ \vdots \\ \cancel{q_l^T} \end{bmatrix} \Rightarrow [q_1 \ q_2 \ \dots \ q_l] \Rightarrow \begin{bmatrix} \hat{X}_1 \\ \hat{X}_2 \\ \vdots \\ \hat{X}_l \end{bmatrix}$$

DECODER

$$y = \sum_{i=1}^m w_i x_i$$



$$w_i(n+1) = w_i(n) + \eta y(n) x_i(n) \quad : i=1, 2, 3, \dots, m$$

$$w_i(n+1) = w_i(n) + \eta y(n) [x_i(n) - y(n) w_i(n)]$$

$$y(n) = X^T(n) W(n) = W^T(n) X(n)$$

$$W(n+1) = W(n) + \eta [X(n) X^T(n) W(n) - W^T(n) X^T(n) W(n)]$$