# Chapter3 Pattern Association & Associative Memory

- Associating patterns which are
  - similar,
  - contrary,
  - in close proximity (spatial),
  - in close succession (temporal)
- Associative recall
  - evoke associated patterns
  - recall a pattern by part of it
  - evoke/recall with incomplete/ noisy patterns
- Two types of associations. For two patterns *s* and *t* 
  - hetero-association (s = t) : relating two different patterns
  - auto-association (s = t): relating parts of a pattern with other parts

- Architectures of NN associative memory
  - single layer (with/out input layer)
  - two layers (for bidirectional assoc.)
- Learning algorithms for AM
  - Hebbian learning rule and its variations
  - gradient descent
- Analysis
  - storage capacity (how many patterns can be remembered correctly in a memory)
  - convergence
- AM as a model for human memory

# **Training Algorithms for Simple AM**

- Network structure: single layer
  - one output layer of non-linear units and one input layer
  - similar to the simple network for classification in Ch. 2



- Goal of learning:
  - to obtain a set of weights w\_ij
  - from a set of training pattern pairs {s:t}
  - such that when s is applied to the input layer, t is computed at the output layer
  - for all training pairs  $s: t: t_j = f(s^T w_{\bullet j})$  for all j

## Hebbian rule

- Similar to hebbian learning for classification in Ch. 2
- Algorithm: (bipolar or binary patterns)
  - For each training samples s:t:  $\Delta w_{ij} = s_i \cdot t_j$
  - $-\Delta w_{ij} \text{ increases if both } s_i \text{ and } t_j$ are ON (binary) or have the same sign (bipolar)
- If  $\Delta w_{ij} = 0$  initiall. Then, after updates for all P training patterns  $w_{ij} = \sum_{P=1}^{P} s_i(p) t_j(p)$   $W = \{w_{ij}\}$
- Instead of obtaining *W* by iterative updates, it can be computed from the training set by calculating the outer product of *s* and *t*.

• Outer product. Let *s* and *t* be **row** vectors. Then for a particular training pair *s*:*t* 

$$\Delta W(p) = s^{T}(p) \cdot t(p) = \begin{bmatrix} s_{1} \\ s_{n} \end{bmatrix} [t_{1}, \dots, t_{m}] = \begin{bmatrix} s_{1}t_{1}, \dots, s_{1}t_{m} \\ s_{2}t_{1}, \dots, s_{2}t_{m} \\ s_{n}t_{1}, \dots, s_{n}t_{m} \end{bmatrix} = \begin{bmatrix} \Delta w_{11}, \dots, \Delta w_{1m} \\ \Delta w_{n1}, \dots, \Delta w_{nm} \end{bmatrix}$$

And 
$$W(P) = \sum_{p=1}^{P} s^{T}(p) \cdot t(p)$$

• It involves 3 nested loops p, i, j (order of p is irrelevant) p=1 to P /\* for every training pair \*/ i = 1 to n /\* for every row in W \*/ j = 1 to m /\* for every element j in row i \*/  $w_{ij} := w_{ij} + s_i(p) \cdot t_j(p)$ 

- Does this method provide a good association?
  - Recall with training samples (after the weights are learned or computed)
  - Apply s(k) to one layer, hope t(k) appear on the other, e.g. f(s(k)W) = t(k)
  - May not always succeed (each weight contains some information from all samples)

$$s(k)W = s(k)\sum_{p=1}^{P} s^{T}(p)t(p) = \sum_{p=1}^{P} s(k) \cdot s^{T}(p) \cdot t(p)$$
$$= s(k)s^{T}(k)t(k) + \sum_{p \neq k} s(k)s^{T}(p)t(p)$$
$$= \left\| s(k) \right\|^{2} t(k) + \sum_{p \neq k} s(k)s^{T}(p)t(p)$$
cross-talk term

- Principal term gives the association between s(k) and t(k).
- Cross-talk represents correlation between *s*(*k*):*t*(*k*) and other training pairs. When cross-talk is large, *s*(*k*) will recall something other than *t*(*k*).
- If all s(p) are orthogonal to each other, then  $s(k) \cdot s^{T}(p) = 0$ , no sample other than s(k):t(k) contribute to the result.
- There are at most *n* orthogonal vectors in an *n*-dimensional space.
- Cross-talk increases when *P* increases.
- How many arbitrary training pairs can be stored in an AM?
  - Can it be more than *n* (allowing some non-orthogonal patterns while keeping cross-talk terms small)?
  - Storage capacity (more later)

### **Delta Rule**

- Similar to that used in Adaline
- The original delta rule for weight update:  $\Delta w_{ij} = \mathbf{a} (t_j y_j) x_i$
- Extended delta rule  $\Delta w_{ij} = \mathbf{a} (t_j y_j) x_i f'(y_i n_j)$ 
  - For output units with differentiable activation functions
  - Derived following gradient descent approach).  $E = \sum_{j=1}^{m} (t_j y_j)^2$

$$\frac{\partial E}{\partial w_{IJ}} = 2(t_J - y_J) \frac{\partial}{\partial w_{IJ}} (t_J - y_J) = -2(t_J - y_J) \frac{\partial y_J}{\partial w_{IJ}}$$
$$= -2(t_J - y_J) \frac{\partial f(y_J - in_J)}{\partial w_{IJ}}$$
$$= -2(t_J - y_J) f'(y_J - in_J) \cdot \frac{\partial}{\partial w_{IJ}} \sum w_{iJ} x_i$$
$$= -2(t_J - y_J) f'(y_J - in_J) \cdot x_I$$
$$\Delta w_{ij} \propto -2 \frac{\partial E}{\partial w_{ij}} = 2(t_J - y_J) f'(y_J - in_J) x_i$$

- same as the update rule for output nodes in BP learning.
- Works well if S are linearly independent (even if not orthogonal).

### **Example of hetero-associative memory**

- Binary pattern pairs s:t with |s| = 4 and |t| = 2.
- Total weighted input to output units:  $y_{ij} = \sum_{i} x_{i} w_{ij}$
- Activation function: threshold

$$\mathbf{y}_{j} = \begin{cases} 1 & if \quad \mathbf{y}_{j} = \mathbf{i}\mathbf{n}_{j} > 0 \\ 0 & if \quad \mathbf{y}_{j} = \mathbf{i}\mathbf{n}_{j} \le 0 \end{cases}$$

• Weights are computed by Hebbian rule (sum of outer products of all training pairs)

$$W = \sum_{p=1}^{P} s_i^{T}(p) t_j(p)$$

• Training samples:

s(p)t(p)
$$p=1$$
 $(1\ 0\ 0\ 0)$  $(1,\ 0)$  $p=2$  $(1\ 1\ 0\ 0)$  $(1,\ 0)$  $p=3$  $(0\ 0\ 0\ 1)$  $(0,\ 1)$  $p=4$  $(0\ 0\ 1\ 1)$  $(0,\ 1)$ 

$$s^{T}(1) \cdot t(1) = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} \begin{pmatrix} 1 & 0 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{pmatrix} \qquad s^{T}(2) \cdot t(2) = \begin{pmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} \begin{pmatrix} 1 & 0 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{pmatrix}$$
$$s^{T}(3) \cdot t(3) = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{pmatrix} \begin{pmatrix} 0 & 1 \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 1 \\ 0 & 1 \end{pmatrix} \qquad s^{T}(4) \cdot t(4) = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 1 \\ 0 \end{pmatrix} \begin{pmatrix} 0 & 1 \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 1 \\ 0 & 1 \end{pmatrix}$$
$$W = \begin{pmatrix} 2 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 2 \end{pmatrix} \qquad \text{Computing the weights}$$

#### recall:

$$\mathbf{x} = (\mathbf{1} \ \mathbf{0} \ \mathbf{0} \ \mathbf{0}) \qquad \mathbf{x} = (\mathbf{0} \ \mathbf{1} \ \mathbf{0} \ \mathbf{0}) \text{ (similar to } \mathbf{S}(1) \text{ and } \mathbf{S}(2)$$

$$(1 \quad 0 \quad 0) \begin{pmatrix} 2 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 2 \end{pmatrix} = (2 \quad 0) \qquad (0 \quad 1 \quad 0 \quad 0) \begin{pmatrix} 2 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 2 \end{pmatrix} = (1 \quad 0)$$

$$\mathbf{y}_1 = 1, \quad \mathbf{y}_2 = 0$$

$$\mathbf{y}_1 = 1, \quad \mathbf{y}_2 = 0$$

$$\mathbf{x} = (\mathbf{0} \ \mathbf{1} \ \mathbf{1} \ \mathbf{0})$$
$$(0 \ 1 \ 1 \ 0) \begin{pmatrix} 2 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 2 \end{pmatrix} = (1 \ 1)$$
$$\mathbf{y}_1 = 1, \quad \mathbf{y}_2 = 1$$

(1 0 0 0), (1 1 0 0) class (1, 0) (0 0 0 1), (0 0 1 1) class (0, 1) (0 1 1 0) is not sufficiently similar to any class

delta-rule would give same or similar results.

## **Example of auto-associative memory**

- Same as hetero-associative nets, except t(p) = s(p).
- Used to recall a pattern by a its noisy or incomplete version. (pattern completion/pattern recovery)
- A single pattern *s* = (1, 1, 1, -1) is stored (weights computed by Hebbian rule outer product)

training pat.
 noisy pat
 missing info
 more noisy

$$(111-1) \cdot W = (444-4) \rightarrow (111-1)$$
  
(-111-1) \cdot W = (222-2) \rightarrow (111-1)  
(001-1) \cdot W = (222-2) \rightarrow (111-1)  
(-1-11-1) \cdot W = (0000) not recognized

- Diagonal elements will dominate the computation when multiple patterns are stored (= P).
- When P is large, *W* is close to an identity matrix. This causes output = input, which may not be any stoned pattern. The pattern correction power is lost.
- Replace diagonal elements by zero.

$$W_{0} = \begin{bmatrix} 0 & 1 & 1 & -1 \\ 1 & 0 & 1 & -1 \\ 1 & 1 & 0 & -1 \\ -1 - 1 - 1 & 0 \end{bmatrix}$$
  
(1 1 1 -1)W'= (3 3 3 -3)  $\rightarrow$  (1 1 1 -1)  
(-1 1 1 -1)W'= (3 1 1 -1)  $\rightarrow$  (1 1 1 -1)

$$(0 \quad 0 \quad 1 \quad -1)W' = (2 \quad 2 \quad 1 \quad -1) \rightarrow (1 \quad 1 \quad 1 \quad -1)$$
$$(-1 \quad -1 \quad 1 \quad -1)W' = (1 \quad 1 \quad -1 \quad 1) \rightarrow wrong$$

## Storage Capacity

- # of patterns that can be correctly stored & recalled by a network.
- More patterns can be stored if they are not similar to each other (e.g., orthogonal)

non-orthogonal  

$$(1 \ -1 \ -1 \ 1) \rightarrow W_0 = \begin{bmatrix} 0 & 0 & -2 & 2 \\ 0 & 0 & 0 & 0 \\ -2 & 0 & 0 & -2 \\ 2 & 0 & -2 & 0 \end{bmatrix}$$
 (1-1-11)· $W_0 = (1 \ 0 \ -1 \ 1)$   
It is not stored correctly  
orthogonal

• Adding one more orthogonal pattern (1 1 1 1) the weight matrix becomes:

- **Theorem**: an n by n network is able to store up to n-1 mutually orthogonal (M.O.) bipolar vectors of n-dimension, but not n such vectors.
- Informal argument: Suppose m orthogonal vectors *a*(1)..... *a*(*m*) are stored with the following weight matrix:

$$w_{ij} = \begin{cases} 0 & \text{if } i = j \text{ (zero diagonal)} \\ \sum_{p=1}^{m} a_i(p)a_j(p) & \text{otherwise (Hebbian rule)} \end{cases}$$

Let's try to recall one of them, say  $a(k) = (a_1(k)....a_n(k))$   $a(k)W = a(k)(w_{\bullet 1}, w_{\bullet 2}, ...., w_{\bullet n})$   $= (a(k) \cdot w_{\bullet 1}, a(k)w_{\bullet 2}, ....a(k)w_{\bullet n})$  $= (\sum_{i=1}^n a_i(k)w_{i1}, \sum_{i=1}^n a_i(k)w_{i2}, ...., \sum_{i=1}^n a_i(k)w_{in})$ 

the jth component :

$$\sum_{i=1}^{n} a_i(k) w_{ij} = \sum_{i \neq j} a_i(k) \cdot \sum_{p=1}^{m} a_i(p) a_j(p)$$
  

$$= \sum_{p=1}^{m} a_j(p) \cdot \sum_{i \neq j} a_i(k) a_i(p)$$
  

$$\sum_{i \neq j} a_i(k) a_i(p) = \sum_{i=1}^{n} a_i(k) a_i(p) - a_j(k) a_j(p)$$
  

$$= \begin{cases} -a_j(k) a_j(p) & k \neq p \text{ (since } a(k) \text{ and } a(p) \text{ are M.O.)} \\ n-1 & k = p \text{ (since } a^T(p) \cdot a(p) = n) \end{cases}$$

$$\begin{split} &\sum_{p=1}^{m} a_{j}(p) \sum_{i \neq j} a_{i}(k) a_{i}(p) = \sum_{p \neq k} a_{j}(p) \Big[ -a_{j}(k) a_{j}(p) \Big] + a_{j}(k) (n-1) \\ &= \sum_{p \neq k} -a_{j}(k) + a_{j}(k) (n-1) \\ &= -(m-1)a_{j}(k) + a_{j}(k) (n-1) \\ &= (n-m)a_{j}(k) \end{split}$$

Therefore, a(k)W = (n-m)a(k)

- When m < n, *a*(*k*) can correctly recall itself when m = n, output is a *0* vector, recall fails
- In linear algebraic term, *a(k)* is a eigenvector of *W*, whose corresponding eigenvalue is (n-m).
  when m = n, *W* has eigenvalue zero, the only eigenvector is *0*, which is a trivial eigenvector.

How many mutually orthogonal bipolar vectors with given dimension n?
 n can be written as n = 2<sup>k</sup> m, where m is an odd integer.

Then maximally:  $2^k$  M.O. vectors

- Follow up questions:
  - What would be the capacity of AM if stored patterns are not mutually orthogonal (say random)
  - Ability of pattern recovery and completion.
     How far off a pattern can be from a stored pattern that is still able to recall a correct/stored pattern
  - Suppose x is a stored pattern, x' is close to x, and x"=
    f(xW) is even closer to x than x'. What should we do?
    Feed back x", and hope iterations of feedback will lead to x

### **Iterative Autoassociative Networks**

• Example:

Example:  $\mathbf{x} = (1, 1, 1, -1)$   $\mathbf{W} = \begin{bmatrix} 0 & 1 & 1 & -1 \\ 1 & 0 & 1 & -1 \\ 1 & 1 & 0 & -1 \\ -1 & -1 & -1 & 0 \end{bmatrix}$  Output units are threshold units

An incomplete recall input : 
$$x' = (1, 0, 0, 0)$$
  
 $x'W = (0, 1, 1, -1) = x''$   
 $x''W = (3, 2, 2, -3) \rightarrow (1, 1, 1, -1) = x$ 

• In general: using current output as input of the next iteration x(0) = initial recall input  $x(I) = f(x(I-1)W), \quad I = 1, 2, \dots$ 

until x(N) = x(K) where K < N

- Dynamic System: state vector x(I)
  - If k = N-1, x(N) is a stable state (fixed point) f(x(N)W) = f(x(N-1)W) = x(N)
    - If x(K) is one of the stored pattern, then x(K) is called a *genuine memory*
    - Otherwise, x(K) is a *spurious memory* (caused by cross-talk/interference between genuine memories)
    - Each fixed point (genuine or spurious memory) is an **attractor** (with different attraction basin)
  - If k != N-1, limit-circle,
    - The network will repeat

 $x(K), x(K+1), \dots, x(N)=x(K)$  when iteration continues.

- Iteration will eventually stop because the total number of distinct state is finite (3<sup>n</sup>) if threshold units are used.
- If sigmoid units are used, the system may continue evolve forever (chaos).

## **Discrete Hopfield Model**

- A single layer network
  - each node as both input and output units
- More than an AM
  - Other applications e.g., combinatorial optimization
- Different forms: discrete & continuous
- Major contribution of John Hopfield to NN
  - Treating a network as a dynamic system
  - Introduce the notion of energy function & attractors into NN research

## **Discrete Hopfield Model (DHM) as AM**

### • Architecture:

- single layer (units serve as both input and output)
- nodes are threshold units (binary or bipolar)
- weights: fully connected, symmetric, and zero diagonal

 $w_{ij} = w_{ji}$  $w_{ii} = 0$ 

 $- x_i$  are external inputs, which may be transient or permanent



Figure 3.7 Discrete Hopfield net.

#### • Weights:

To store patterns s(p), p=1,2,...P  
**bipolar:** 
$$w_{ij} = \sum_{p} s_i(p) s_j(p)$$
  $i \neq j$   
 $w_{ii} = 0$   
same as Hebbian rule (with zero diagonal)  
**binary:**  $w_{ij} = \sum_{p} (2s_i(p) - 1)(2s_j(p) - 1)$   $i \neq j$   
 $w_{ii} = 0$ 

converting s(p) to bipolar when constructing W.

#### • Recall

- Use an input vector to recall a stored vector (book calls the application of DHM)
- Each time, randomly select a unit for update

### **Recall Procedure**

1. Apply recall input vector  $\boldsymbol{x}$  to the network:  $\boldsymbol{y}_i \coloneqq \boldsymbol{x}_i$   $i = 1, 2, ..., \boldsymbol{n}$ 

- 2. While convergence = fails do
  - 2.1.Randomly select a unit

2.2. Compute 
$$y_i n_i = x_i + \sum_{j \neq i} y_j w_{ji}$$
  
2.3. Determine activation of Yi

$$\mathbf{y}_{i} = \begin{cases} 1 & \text{if } \mathbf{y}_{i} \mathbf{n}_{i} > \mathbf{q}_{i} \\ \mathbf{y}_{i} & \text{if } \mathbf{y}_{i} \mathbf{n}_{i} = \mathbf{q}_{i} \\ -1 & \text{if } \mathbf{y}_{i} \mathbf{n}_{i} < \mathbf{q}_{i} \end{cases}$$

2.4. Periodically test for convergence.

#### • Notes:

- 1. Each unit should have equal probability to be selected at step 2.1
- 2. Theoretically, to guarantee convergence of the recall process, only one unit is allowed to update its activation at a time during the computation. However, the system may converge faster if all units are allowed to update their activations at the same time.
- 3. Convergence test:  $y_i(current) = y_i(next) \quad \forall i$
- 4.  $\boldsymbol{q}_i$  usually set to zero.
- 5.  $x_i$  in step 2.2  $(y_i n_j = x_i + \sum_j y_j w_{ji})$  is optional.

#### • Example:

Store one pattern:<br/>binary pattern (1, 1, 1, 0)<br/>(bipolar counterpart (111-1) $W = \begin{bmatrix} 0 & 1 & 1 & -1 \\ 1 & 0 & 1 & -1 \\ 1 & 1 & 0 & -1 \\ -1 & -1 & -1 \\ -1 & -1 & 0 \end{bmatrix}$ 

Recall input x = (0, 0, 1, 0), first two bits are wrong  $Y_1$  is selected  $y_-in_1 = x_1 + \sum y_1 \cdot w_{j_1} = 0 + 1 = 1$   $y_-in_4 = x_4 + \sum y_4 \cdot w_{j_4} = 0 + (-2) = -2$   $y_1 = 1$   $y_4 = -2$ Y = (1, 0, 1, 0) Y = (1, 0, 1, 0)

$$\frac{Y_{3} \text{ is selected}}{y_{1} i n_{3} = x_{3} + \sum y_{3} \cdot w_{j_{3}} = 1 + 1 = 2} \qquad \frac{Y_{2} \text{ is selected}}{y_{1} i n_{2} = x_{2} + \sum y_{2} \cdot w_{j_{2}} = 0 + 2 = 2}$$

$$y_{3} = 1 \qquad \qquad y_{2} = 1$$

$$Y = (1, 0, 1, 0) \qquad \qquad Y = (1, 1, 1, 0)$$

The stored pattern is correctly recalled

## Convergence Analysis of DHM

### • Two questions:

- 1.Will Hopfield AM converge (stop) with any given recall input?
- 2.Will Hopfield AM converge to the stored pattern that is **closest** to the recall input ?
- Hopfield provides answer to the first question
  - By introducing an energy function to this model,
  - No satisfactory answer to the second question so far.

### • Energy function:

- Notion in thermo-dynamic physical systems. The system has a tendency to move toward lower energy state.
- Also known as Lyapunov function. After Lyapunov theorem for the stability of a system of differential equations.

• In general, the energy function E(y(t)), where y(t) is the state of the system at step (time) t, must satisfy two conditions

1. 
$$E(t)$$
 is bounded from below  $E(t) \ge c \quad \forall t$ 

2. E(t) is monotonically nonincreasing.

 $\Delta E(t+1) = E(t+1) - E(t) \le 0 \text{ (in continuous version : } \dot{E}(t) \le 0)$ 

• The energy function defined for DHM

$$\boldsymbol{E} = -0.5 \sum_{i \neq j} \sum_{j} y_{i} y_{j} w_{ij} - \sum_{i} x_{i} y_{i} + \sum_{i} \boldsymbol{q}_{i} y_{i}$$

• Show 
$$\Delta E(t+1) \leq 0$$
  
At t+1, $Y_k$  is selected for update  
 $\Delta y_k(t+1) = y_k(t+1) - y_k(t)$   
Note:  $\Delta y_j(t+1) = 0$   $j \neq k$  (only one unit can update at a time)  
 $E(t+1) - E(t)$   
 $= (-0.5\sum_{i\neq j}\sum_j y_i(t+1)y_j(t+1)w_{ij} - \sum_i x_i y_i(t+1) + \sum_i q_i y_i(t+1))$   
 $-(-0.5\sum_{i\neq j}\sum_j y_i(t)y_j(t)w_{ij} - \sum_i x_i y_i(t) + \sum_i q_i y_i(t))$ 

terms which are different in the two parts are those involving  $y_k$ 

$$\sum_{j} y_{k} y_{j} w_{jk}, \qquad \sum_{i} y_{i} y_{k} w_{ki}, \qquad x_{k} y_{k}, \qquad \mathbf{q}_{k} y_{k}$$
$$\Delta E(t+1) = -\left[\sum_{j \neq k} y_{j}(t) w_{kj} + x_{k} - \mathbf{q}_{k}\right] \Delta y_{k}(t+1)$$

cases:

if 
$$y_k(t) = 1 \& y_k(t+1) = -1 \quad \Delta y_k(t+1) = -2$$
  
 $\Rightarrow y_i n_k < \mathbf{q}_k \Rightarrow \Delta E(t+1) < 0$   
if  $y_k(t) = -1 \& y_k(t+1) = 1 \quad \Delta y_k(t+1) = 1$   
 $\Rightarrow y_i n_k > \mathbf{q}_k \Rightarrow \Delta E(t+1) < 0$   
otherwise,  $y_k(t+1) = y_k(t) \Rightarrow \Delta y_k(t+1) = 0 \Rightarrow \Delta E(t+1) = 0$ 

• Show E(t) is bounded from below, since  $y_i, x_i, q_i, w_{ij}$  are all bounded, E is bounded.

#### • Comments:

1.Why converge.

- Each time, E is either unchanged or decreases an amount.
- E is bounded from below.
- There is a limit E may decrease. After finite number of steps, E will stop decrease no matter what unit is selected for update.

 $\forall k \text{ either } y_k(t+1) = y_k(t) \Longrightarrow \Delta y_k = 0$ or  $y_i in_k = \mathbf{q} \Longrightarrow \Delta y_k = 0$ 

2. The state the system converges is a stable state.

Will return to this state after some small perturbation. It is called an **attractor** (with different attraction basin)

- 3.Error function of BP learning is another example of energy/Lyapunov function. Because
  - It is bounded from below (E>0)
  - It is monotonically non-increasing (W updates along gradient descent of E)

## Capacity Analysis of DHM

- *P*: maximum number of random patterns of dimension *n* can be stored in a DHM of *n* nodes
- Hopfield's observation:  $P \approx 0.15n$ ,  $\frac{P}{a} \approx 0.15$
- Theoretical analysis:  $P \approx \frac{n}{2\log_2 n}, \quad \frac{n}{n} \approx \frac{1}{2\log_2 n}$

*P/n* decreases because larger n leads to more interference between stored patterns.

• Some work to modify HM to increase its capacity to close to *n*, W is trained (not computed by Hebbian rule).

### My Own Work:

- One possible reason for the small capacity of HM is that it does not have hidden nodes.
- Train feed forward network (with hidden layers) by BP to establish pattern auto-associative.
- Recall: feedback the output to input layer, making it a dynamic system.
- Shown 1) it will converge, and 2) stored patterns become genuine memories.
- It can store many more patterns (seems O(2^n))
- Its pattern complete/recovery capability decreases when n increases (# of spurious attractors seems to increase exponentially)



## **Bidirectional AM(BAM)**

#### • Architecture:

- Two layers of non-linear units: X-layer, Y-layer
- Units: discrete threshold, continuing sigmoid (can be either binary or bipolar).



Figure 3.8 Bidirectional associative memory.

- Weights:  $- W_{n \times m} = \sum_{p=1}^{P} s^{T}(p) \cdot t(p) \quad \text{(Hebbian/outer product)}$   $- \text{Symmetric: } w_{ij} = w_{ji}$ 
  - Convert binary patterns to bipolar when constructing W
- Recall:
  - Bidirectional, either by X (to recall a Y) or by Y (to recall a X)
  - Recurrent:  $y(t) = (f(y_in_1(t),...,f(y_in_m(t)))$

where 
$$y_{i_{j}}(t) = \sum_{i=1}^{n} w_{i_{j}} \cdot x_{i}(t-1)$$
  
 $x(t+1) = (f(x_{i_{1}}(t+1),...,f(x_{i_{n}}(t+1))))$   
where  $x_{i_{i}}(t+1) = \sum_{i=1}^{m} w_{i_{j}} \cdot y_{j}(t)$ 

 Update can be either asynchronous (as in HM) or synchronous (change all Y units at one time, then all X units the next time)

- Analysis (discrete case)
  - Energy function: (also a Lyapunov function)

$$L = -0.5(XWY^{T} + YW^{T}X^{T}) = -XWY^{T}$$
$$= -\sum_{i=1}^{m} \sum_{i=1}^{n} x_{i}w_{ij}y_{j}$$

- The proof is similar to DHM
- Holds for both synchronous and asynchronous update (holds for DHM only with asynchronous update, due to lateral connections.)
- Storage capacity:  $\boldsymbol{o}(\max(\boldsymbol{n}, \boldsymbol{m}))$