Chapter 4. Neural Networks Based on Competition

- Competition is important for NN
 - Competition between neurons has been observed in biological nerve systems
 - Competition is important in solving many problems

To classify an input pattern into one of the m classes

- idea case: one class node has output 1, all other 0;
- often more than one class nodes have non-zero output



 If these class nodes compete with each other, maybe only one will win eventually (winner-takes-all). The winner represents the computed classification of the input

- Winner-takes-all (WTA):
 - Among all competing nodes, only one will win and all others will lose
 - We mainly deal with single winner WTA, but multiple winners WTA are possible (and useful in some applications)
 - Easiest way to realize WTA: have an external, central arbitrator (a program) to decide the winner by comparing the current outputs of the competitors (break the tie arbitrarily)
 - This is biologically unsound (no such external arbitrator exists in biological nerve system).

- Ways to realize competition in NN
 - Lateral inhibition (Maxnet, Mexican hat)
 - output of each node feeds to others through inhibitory connections (with negative weights)

– Resource competition

- output of x_k is distributed to y_i and y_j proportional to w_ki and w_kj, as well as y_i and y_j
- self decay
- biologically sound

• Learning methods in competitive networks

- Competitive learning
- Kohonen learning (self-organizing map, SOM)
- Counter-propagation net
- Adaptive resonance theory (ART) in Ch. 5



 $w_{ii} < 0$

Fixed-weight Competitive Nets

• Maxnet

- Lateral inhibition between
 - competitors



weights :
$$w_{ij} = \begin{cases} 1 & \text{if } i = j \\ -e & \text{otherwise} \end{cases}$$

activation function : $f(x) = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$

- Notes:
 - Competition:
 - iterative process until the net stabilizes (at most one node with positive activation)
 - 0 < e < 1/m, where *m* is the # of competitors
 - *e* too small: takes too long to converge
 - **e** too big: may suppress the entire network (no winner)

Mexical Hat

- Architecture: For a given node,
 - close neighbors: cooperative (mutually excitatory , w > 0)
 - farther away neighbors: competitive (mutually inhibitory, w < 0)
 - too far away neighbors: irrelevant (w = 0)



- Need a definition of distance (neighborhood):
 - one dimensional: ordering by index (1,2,...n)
 - two dimensional: lattice

weights

$$w_{ij} = \begin{cases} c_1 & \text{if } \text{distance}(i, j) < k \\ c_2 & \text{if } \text{distance}(i, j) = k \\ 0 & \text{if } \text{distance}(i, j) > k \end{cases}$$

activation function

$$f(\mathbf{x}) = \begin{cases} 0 & if \quad \mathbf{x} < 0 \\ \mathbf{x} & if \quad 0 \le \mathbf{x} \le \max \\ \max & if \quad \mathbf{x} > \max \end{cases}$$

ramp function :

example: x(0) = (0.0, 0.5, 0.8, 1.0, 0.8, 0.5, 0.0) x(1) = (0.0, 0.38, 1.06, 1.16, 1.06, 0.38, 0.9)x(2) = (0.0, 0.39, 1.14, 1.66, 1.14, 0.39, 0.0)

- Equilibrium:
 - negative input = positive input for all nodes
 - winner has the highest activation;
 - its cooperative neighbors also have positive activation;
 - its competitive neighbors have negative activations.

Hamming Network

- Hamming distance of two vectors, *x* and *y* of dimension n,
 - Number of bits in disagreement.
 - In bipolar:

 $x \cdot y = a - d$

where : a is number of bits in agreement in x and y

d is number of bits different in *x* and *y* d = n - a hamming distance $x \cdot y = 2a - n$ $a = 0.5(x \cdot y + n)$ larger $x \cdot y \Rightarrow$ larger $a \Rightarrow$ shorter Hamming distance

- Suppose a space of patterns is divided into k classes, each class has an exampler (representative) vector e_i .
- An input x belongs to class i, if and only if x is closer to e_i than to any other e_i , i.e., $x \cdot e_i \ge x \cdot e_i$, $\forall j \neq i$
- Hamming net is such a classifier:
 - Weights: let Y_i represent class j

 $\boldsymbol{w}_{\bullet j} = 0.5\boldsymbol{e}_j, \ \boldsymbol{b}_j = 0.5\boldsymbol{n}$

- The total input to Y_j $y_i n_j = b_j + \sum_i w_{ij} x_i$ $= \frac{1}{2} (x \cdot e_j + n) = a_j$



- Upper layer: MAX net
 - it takes the y_in as its initial value, then iterates toward stable state
 - one output node with highest y_in will be the winner because its weight vector is closest to the input vector
- As associative memory:
 - each Y_i corresponds to a stored pattern;
 - pattern connection/completion;
 - storage capacity

total # of nodes: k

total # of patterns stored: k

capacity: k (or k/k = 1)

• Implicit lateral inhibition by competing limited resources: the activation of the input nodes

$$y_{1} y_{j} y_{m} y_{m$$

Competitive Learning

- Unsupervised learning
- Goal:
 - Learn to form classes/clusters of examplers/sample patterns according to similarities of these exampers.
 - Patterns in a cluster would have similar features
 - No prior knowledge as what features are important for classification, and how many classes are there.
- Architecture:
 - Output nodes:
 - Y_1,.....Y_m,
 - representing the m classes
 - They are competitors
 (WTA realized either by an external procedure or

by lateral inhibition as in Maxnet)



- Training:
 - Train the network such that the weight vector w.j associated with Y_j becomes the representative vector of the class of input patterns Y_j is to represent.
 - Two phase unsupervised learning
- competing phase:
 - apply an input vector \mathbf{x} randomly chosen from sample set.
 - compute output for all y: $y_j = x \cdot w_{\bullet j}$
 - determine the winner (winner is not given in training samples so this is unsupervised)
- rewarding phase:
 - the winner is reworded by updating its weights (weights associated with all other output nodes are not updated)
- **repeat** the two phases many times (and gradually reduce the learning rate) until all weights are stabilized.

• Weight update: $w_{\bullet j} = w_{\bullet j} + \Delta w_{\bullet j}$ – Method 1: Method 2

$$\Delta w_{\bullet j} = \mathbf{a} \ (\mathbf{x} - \mathbf{w}_{\bullet j}) \qquad \Delta w_{\bullet j} = \mathbf{a} \ \mathbf{x}$$

x+w_j

ax



In each method, $W_{\bullet i}$ is moved closer to x

- Normalizing the weight vector to unit length after it is updated

$$\boldsymbol{w}_{\bullet j} = \frac{\boldsymbol{w}_{\bullet j}}{\left\|\boldsymbol{w}_{\bullet j}\right\|}$$

- $W_{\bullet j}$ is moving to the center of a cluster of sample vectors after repeated weight updates
 - Three examplers:S(1), S(2) and S(3)
 - Initial weight vector $w_j(0)$
 - After successively trained by S(1), S(2), and S(3), the weight vector changes to w_j(1), w_j(2), and w_j(3)



Examples

- A simple example of competitive learning (pp. 172-175)
 - 4 vectors of dimension 4 in 2 classes (4 input nodes, 2 output nodes)

$$S(1) = (1, 1, 0, 0) \quad S(2) = (0, 0, 0, 1)$$

$$S(3) = (1, 0, 0, 0) \quad S(4) = (0, 0, 1, 1)$$

Initialization: **a** = 0.6, weight matrix: **W** =
$$\begin{bmatrix} .2 & .8 \\ .6 & .4 \\ .5 & .7 \\ .9 & .3 \end{bmatrix}$$

$$D(1) = \|S(1) - w_{\bullet 1}\|$$

$$= (.2 - 1)^{2} + (.6 - 1)^{2} + (.5 - 0)^{2} + (.9 - 0)^{2} = 1.86$$

$$D(1) = \|S(1) - w_{\bullet 2}\| = 0.98, \text{ class } 2 \text{ wins}$$

$$w_{\bullet 2} = \begin{bmatrix} .8 \\ .4 \\ .7 \\ .3 \end{bmatrix} + 0.6(\begin{bmatrix} 1 \\ .0 \\ .0 \end{bmatrix} - \begin{bmatrix} .8 \\ .4 \\ .7 \\ .3 \end{bmatrix}) = \begin{bmatrix} .92 \\ .76 \\ .28 \\ .12 \end{bmatrix} \text{ then } W = \begin{bmatrix} .2 & .92 \\ .6 & .76 \\ .5 & .28 \\ .9 & .12 \end{bmatrix}$$

_	Similarly, after training with	Γ.	.08
S	$\mathbf{S}(2) = (0, 0, 0, 1) ,$	11/	.24
	in which class 1 wins,	$\mathbf{v}\mathbf{v} = $.20
	weight matrix becomes	Į.	.96

- At the end of the first iteration (each of the 4 vectors are used), weight matrix becomes
- Reduce **a**

 $a = 0.5 \cdot a = 0.5 \cdot 0.6 = 0.3$

Repeat training. After 10
 iterations, weight matrix becomes

$$W = \begin{bmatrix} 6.7e - 17 \ 1.000000 \\ 2.0e - 16 \ .4900000 \\ .5100000 \ 2.3e - 16 \\ 1.000000 \ 1.0e - 16 \end{bmatrix} \rightarrow \begin{bmatrix} 0.0 \ 1.0 \\ 0.0 \ 0.5 \\ 0.5 \ 0.0 \\ 1.0 \ 0.0 \end{bmatrix}$$

$$W = \begin{bmatrix} .00 & .92 \\ .24 & .76 \\ .20 & .28 \\ .96 & .12 \end{bmatrix}$$
$$W = \begin{bmatrix} .032 & .970 \\ .096 & .300 \\ .680 & .110 \\ .980 & .048 \end{bmatrix}$$

 0^{7}

- S(1) and S(3) belong to class 2
- S(2) and S(4) belong to class 1
- w_1 and w_2 are the centroids of the two classes

Comments

- 1. Ideally, when learning stops, each $w_{\bullet j}$ is close to the centroid of a group/cluster of sample input vectors.
- 2. To stabilize $w_{\cdot j}$, the learning rate **a** may be reduced slowly toward zero during learning.
- 3. # of output nodes:
 - too few: several clusters may be combined into one class
 - too many: over classification
 - ART model (later) allows dynamic add/remove output nodes
- 4. Initial $w_{\bullet j}$:
 - training samples known to be in distinct classes, provided such info is available
 - random (bad choices may cause anomaly)



*w*_{•1} will always win no matter
the sample is from which class *w*_{•2} is stuck and will not participate
in learning

unstuck:

let output nodes have some conscience temporarily shot off nodes which have had very high winning rate (hard to determine what rate should be considered as "very high")

Kohonen Self-Organizing Maps (SOM)

- Competitive learning (Kohonen 1982) is a special case of SOM (Kohonen 1989)
- In competitive learning,
 - the network is trained to organize input vector space into subspaces/classes/clusters
 - each output node corresponds to one class
 - the output nodes are not ordered: *random map*



- The topological order of the three clusters is 1, 2, 3
- The order of their maps at output nodes are 2, 3, 1
- The map does not preserve the topological order of the training vectors

• Topographic map

- a mapping that preserves neighborhood relations between input vectors, (topology preserving or feature preserving).
- if x1 and x2 are two neighboring input vectors (by some distance metrics),
- their corresponding winning output nodes (classes), i and j must also be close to each other in some fashion
- one dimensional: line or ring, node i has neighbors $i \pm 1$ or $i \pm 1 \mod n$
- two dimensional:grid.

rectangular: node(i, j) has neighbors: ($i, j \pm 1$), ($i \pm 1, j$), (or additional ($i \pm 1, j \pm 1$)) hexagonal: 6 neighbors

- Biological motivation
 - Mapping two dimensional continuous inputs from sensory organ (eyes, ears, skin, etc) to two dimensional discrete outputs in the nerve system.
 - Retinotopic map: from eye (retina) to the visual cortex.
 - Tonotopic map: from the ear to the auditory cortex
 - These maps preserve topographic orders of input.
 - Biological evidence shows that the connections in these maps are not entirely "pre-programmed" or "pre-wired" at birth. Learning must occur after the birth to create the necessary connections for appropriate topographic mapping.

SOM Architecture

- Two layer network:
 - Output layer:
 - Each node represents a class (of inputs)
 - Node activation : $y_j = x \cdot w_{\bullet j} = \sum_{i} w_{ij} x_i$
 - Neighborhood relation is defined over these nodes Each node cooperates with all its neighbors within distance R and competes with all other output nodes.
 - Cooperation and competition of these nodes can be realized by Mexican Hat model
 - R = 0: all nodes are competitors (no cooperative) → random map
 - R > 0: \rightarrow topology preserving map

SOM Learning

- 1. Initialize W_{i} for all output nodes, and **a** to a small value
- 2. For a randomly selected input sample/exampler \boldsymbol{x} determine the winning output node \boldsymbol{J} either $\boldsymbol{x} \cdot \boldsymbol{W}_{\boldsymbol{\cdot} \boldsymbol{I}}$ is maximum or

$$|\mathbf{x} - \mathbf{w}_{\bullet J}|| = \sum_{i} (\mathbf{w}_{ij} - \mathbf{x}_{i})^{2}$$
 is minimum

- 3. For all output node j with $|J j| \le R$, update the weight $w_{ij} = w_{ij} + \mathbf{a} (x w_{ij})$
- 4. Periodically reduce **a** and **R** slowly.
- 5. Repeat 2 4 until the network stabilized.

Notes

- 1. Initial weights: small random value from (-e, e)
- 2. Reduction of **a** :

Linear: $\boldsymbol{a}(t + \Delta t) = \boldsymbol{a}(t) - \Delta \boldsymbol{a}$

Geometric: $\boldsymbol{a}(t + \Delta t) = \boldsymbol{a}(t)\boldsymbol{b}$

 Δt may be 1 or greater than 1

3. Reduction of \mathbf{R} : $\mathbf{R}(t + \Delta t) = \mathbf{R}(t) - 1$ while $\mathbf{R}(t) > 0$

should be much slower than **a** reduction.

R can be a constant through out the learning.

4. Effect of learning

For each input x, not only the weight vector of winner J is pulled closer to x, but also the weights of J's close neighbors (within the radius of R).

- 5. Eventually, $w_{\bullet j}$ becomes close (similar) to $w_{\bullet j \pm 1}$. The classes they represent are also similar.
- 6. May need large initial **R**

Examples

- A simple example of competitive learning (pp. 172-175)
 - 4 vectors of dimension 4 in 2 classes (4 input nodes, 2 output nodes)

$$S(1) = (1, 1, 0, 0) \quad S(2) = (0, 0, 0, 1)$$

$$S(3) = (1, 0, 0, 0) \quad S(4) = (0, 0, 1, 1)$$

Initialization: $\mathbf{a} = 0.6$, weight matrix: $\mathbf{W} = \begin{bmatrix} .2 & .8 \\ .6 & .4 \\ .5 & .7 \\ .9 & .3 \end{bmatrix}$

$$D(1) = \|S(1) - \mathbf{w}_{\bullet 1}\|$$

$$= (.2 - 1)^{2} + (.6 - 1)^{2} + (.5 - 0)^{2} + (.9 - 0)^{2} = 1.86$$

$$D(1) = \|S(1) - \mathbf{w}_{\bullet 2}\| = 0.98, \text{ class } 2 \text{ wins}$$

$$w_{\bullet 2} = \begin{bmatrix} .8 \\ .4 \\ .7 \\ .3 \end{bmatrix} + 0.6(\begin{bmatrix} 1 \\ .0 \\ .0 \end{bmatrix} - \begin{bmatrix} .8 \\ .4 \\ .7 \\ .3 \end{bmatrix}) = \begin{bmatrix} .92 \\ .76 \\ .28 \\ .12 \end{bmatrix} \text{ then } \mathbf{W} = \begin{bmatrix} .2 & .92 \\ .6 & .76 \\ .5 & .28 \\ .9 & .12 \end{bmatrix}$$

How to illustrate Kohonen map

– Input vector: 2 dimensional

Output vector: 1 dimensional line/ring or 2 dimensional grid.

Weight vector is also 2 dimension

- Represent the topology of output nodes by points on a 2 dimensional plane. Plotting each output node on the plane with its weight vector as its coordinates.
- Connecting neighboring output nodes by a line output nodes: (1, 1) (2, 1) (1, 2) weight vectors: (0.5, 0.5) (0.7, 0.2) (0.9, 0.9)



Traveling Salesman Problem (TSP) by SOM

- Each city is represented as a 2 dimensional input vector (its coordinates (x, y)),
- Output nodes C_j form a one dimensional SOM, each node corresponds to a city.
- Initially, C_1, ... , C_n have random weight vectors
- During learning, a winner C_j on an input (x, y) of city I, not only moves its w_j toward (x, y), but also that of of its neighbors (w_(j+1), w_(j-1)).
- As the result, C_(j-1) and C_(j+1) will later be more likely to win with input vectors similar to (x, y), i.e, those cities closer to I
- At the end, if a node j represents city I, it would end up to have its neighbors j+1 or j-1 to represent cities similar to city I (i,e., cities close to city I).
- This can be viewed as a concurrent greedy algorithm





Two candidate solutions:

ADFGHIJBC ADFGHIJCB

Additional examples





Counter propagation network (CPN)

- Basic idea of CPN
 - Purpose: fast and coarse approximation of vector mapping y = f(x)
 - not to map any given x to its f(x) with given precision,
 - input vectors x are divided into clusters/classes.
 - each cluster of x has one output y, which is (hopefully) the average of f(x) for all x in that class.
 - Architecture: Simple case: FORWARD ONLY CPN,



from (input)from class tofeatures to class(output) features

- Learning in two phases:
 - training sample x:y where y = f(x) is the precise mapping
 - Phase1: $V_{\bullet j}$ is trained by competitive learning to become the representative vector of a cluster of input vectors x (use sample x only)

1. For a chosen x, feedforward to determined the winning z_i

2.
$$v_{ij}(new) = v_{ij}(old) + a(x_i - v_{ij}(old))$$

3. Reduce \boldsymbol{a} , then repeat steps 1 and 2 until stop condition is met

Phase 2: w_j, is trained by delta rule to be an average output of f(x) where x is an input vector that causes z_j to win (use both x and y).
1. For a chosen x, feedforward to determined the winning z_i

2.
$$v_{ij}(new) = v_{ij}(old) + a(x_i - v_{ij}(old))$$
 (optional)

3.
$$w_{jk}(new) = w_{jk}(old) + a(y_k - w_{jk}(old))$$

4. Repeat steps 1 - 3 until stop condition is met

Notes

- A combination of both unsupervised learning (for $v_{\bullet j}$ in phase 1) and supervised learning (for $w_{j\bullet}$ in phase 2).
- After phase 1, clusters are formed among sample input x, each $v_{\bullet j}$ is a representative of a cluster (average).
- After phase 2, each cluster j maps to an output vector y, which is the average of {f(x): x ∈ cluster j}
- View phase 2 learning as following delta rule

 $w_{jk} = w_{jk} + \mathbf{a} (y_k - w_{jk}) \text{ where } y_k - w_{jk} \propto -\frac{\partial E}{\partial w_{jk}}, \text{ because}$ $\frac{\partial E}{\partial w_{ik}} = \frac{\partial}{\partial w_{ik}} \sum_{ik} (y_k - w_{jk} z_j)^2 = -2(y_k - w_{jk} z_j) \text{ when } z_j \text{ wins}$

• It can be shown that, when $t \to \infty$, $v_{\bullet j}(t) \to \langle x \rangle$ and $w_{j\bullet}(t) \to \langle f(x) \rangle$ where $\langle x \rangle$ is the mean of all training samples that make x win Show only on v_{ij} (proof of w_{jk} is similar.) Weightupdate rule can be rewriteenas $v_{ij}(t+1) = (1-a)v_{ij}(t) + a x_i(t+1)$

$$\begin{aligned} \mathbf{v}_{ij}(t+1) &= (1-\mathbf{a})\mathbf{v}_{ij}(t) + \mathbf{a} \, x_i(t+1) \\ &= (1-\mathbf{a})((1-\mathbf{a})\mathbf{v}_{ij}(t-1) + \mathbf{a} \, x_i(t)) + \mathbf{a} \, x_i(t+1) \\ &= (1-\mathbf{a})^2 \mathbf{v}_{ij}(t-1) + (1-\mathbf{a})\mathbf{a} \, x_i(t) + \mathbf{a} \, x_i(t+1) \\ &\approx \mathbf{a} \left[x_i(t+1) + x_i(t)(1-\mathbf{a}) + x_i(t-1)(1-\mathbf{a})^2 + \dots x_i(1)(1-\mathbf{a})^t \right] \end{aligned}$$

If
$$\boldsymbol{x}$$
 are drawn randomly from the training set, then

$$E[\boldsymbol{v}_{ij}(t+1)] = E[\boldsymbol{a} (\boldsymbol{x}_i(t+1) + \boldsymbol{x}_i(t)(1-\boldsymbol{a}) + \dots \boldsymbol{x}_1(1-\boldsymbol{a})^t]$$

$$= \boldsymbol{a} [E(\boldsymbol{x}(t+1)) + (1-\boldsymbol{a})E(\boldsymbol{x}(t)) + \dots (1-\boldsymbol{a})^t E(\boldsymbol{x}_i(1))]$$

$$= \boldsymbol{a} \langle \boldsymbol{x} \rangle [1 + (1-\boldsymbol{a}) + \dots (1-\boldsymbol{a})^t]$$

$$\Rightarrow \boldsymbol{a} \langle \boldsymbol{x} \rangle \frac{1}{1-(1-\boldsymbol{a})} = \langle \boldsymbol{x} \rangle$$

- After training, the network works like a look-up of math table.
 - For any input *x*, find a region where *x* falls (represented by the wining *z* node);
 - use the region as the index to look-up the table for the function value.
 - CPN works in multi-dimensional input space
 - More cluster nodes (z), more accurate mapping.

Full CPN

- If both y = f(x) and its inverse function $x = f^{-1}(y)$ exist we can establish bi-directional approximation
- Two pairs of weights matrices: V(x to z) and U(z to y) for approx. map x to y = f(x)W(y to z) and T(z to x) for approx. map y to $x^* = f(x)$
- When *x*:*y* is applied (*x* on *X* and *y* on *Y*), they can jointly determine the winner J or separately for z_{Jx}, z_{Jy}
- pp. 196–206 for more details