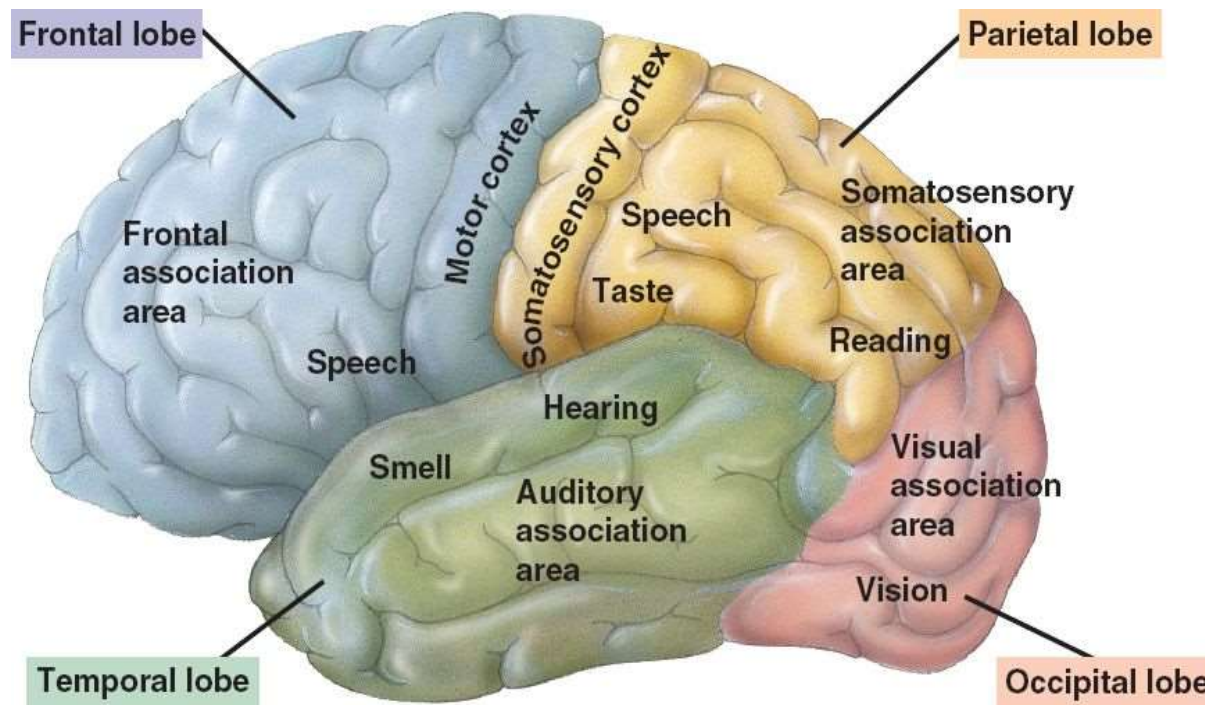# 4. Self-Organizing Map Neural Network

The development of the self-organizing map is inspired by the research in neurobiology. Consider the map of a human brain:



The map of a human brain

The diagram presents the map of the cerebral cortex. The different regions of the cortex are identified by the thickness of their layer and the types of neurons within them. Some of the example regions are:

(1) Motor cortex
(2) Visual cortex
(3) Auditory cortex
(4) Speech cortex

The brain map shows clearly that
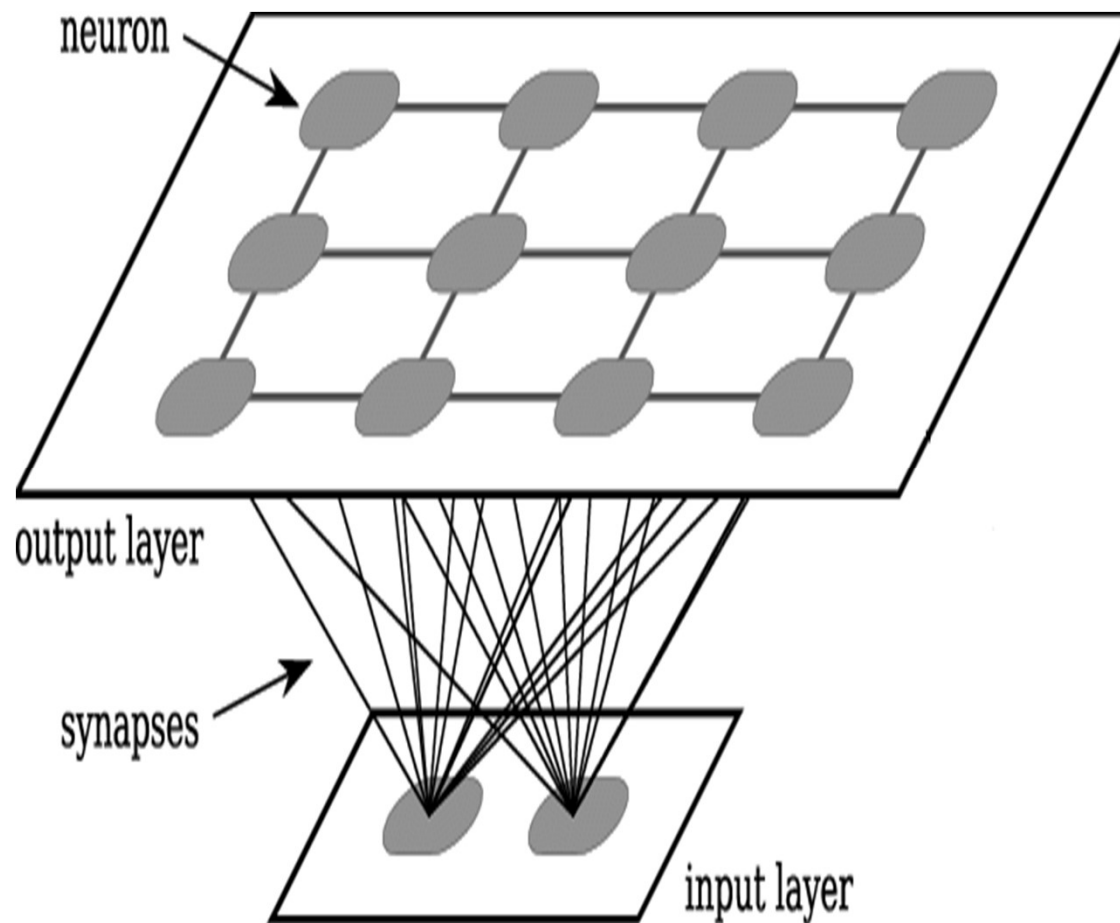
(1) Different sensory inputs (motor, visual, auditory etc.) are mapped to corresponding regions in an orderly fashion;
(2) These cortical maps are not entirely genetically pre-determined. Rather, they are sketched in during the early development of the nervous system.
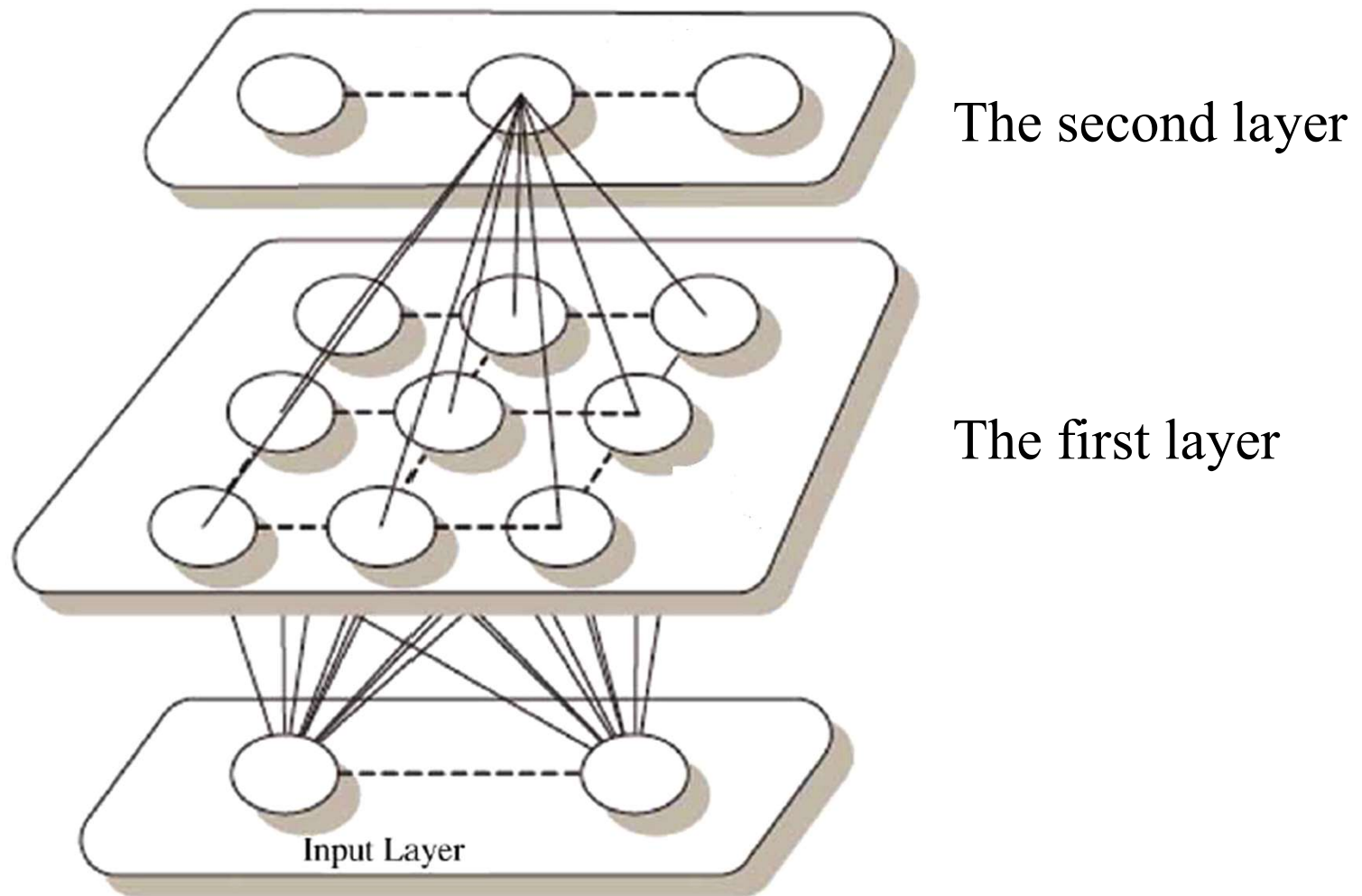
# Principle of topographic map

The essential point of the discovery in the neurobiology lies in the principle of topographic map formation stated as follows:

*The spatial location of an output neuron in the topographic map corresponds to a particular domain or feature of the input data.*

Based on this principle, self-organizing map neural network is developed. The output neurons are usually arranged in a one-dimensional or two-dimensional lattice. This is to ensure that each neuron has a set of neighbors.

Model 1 of the self-organizing map neural network

The second layer

The first layer

Input Layer

Model 2 of the self-organizing map neural network

## Self-Organizing Map (SOM)

The objective of SOM is to transform an incoming signal pattern of arbitrary dimension into a one-dimensional (1D) or two-dimensional (2D) map, and to perform this transform adaptively in a topological ordered fashion. The essential ingredients of the SOM neural network include:

(1) A 1D or 2D lattice of neurons that computes simple discriminant functions of inputs received.
(2) A mechanism that compares these functions and selects the neurons with the largest discriminant function value.
(3) An adaptive process that enables the activated neurons to increase their discriminant values.
(4) An interactive network that activates the selected neuron and its neighbors.

The algorithm responsible for the formation of the SOM proceeds first by initializing the synaptic weights in the network through assigning them small random values. After initialization, three essential processes are involved in the formation (training) of the self-organizing map:

(1) **Competition**. For each input pattern, the neurons in the network compute their respective values of a discriminant function. This discriminant function provides the basis for competition among the neurons;

(2) **Cooperation.** The winning neuron determines the spatial location of a topological neighborhood of excited neurons in the 1D or 2D lattice;

(3) **Weights adaptation**. This mechanism enables the excited neurons to adjust their weights.

## Competition process

Assuming that an input pattern selected randomly from the input space is denoted by:

$$\mathbf{x} = [x_1, x_2, \cdots, x_n]^T$$

The weight vector of each neuron in the network has the same dimension as the input space. Let the weight vector of neuron j be denoted by:

$$\mathbf{w}_j = [w_{j1}, w_{j2}, \cdots, w_{jn}]$$

The similarity between $\mathbf{x}$ and $\boldsymbol{w_j}$ is used as the discriminant function for competition. To find the best match of the input vector $\mathbf{x}$ with the synaptic weight vector $\mathbf{w}_j$, we can simply use the inner products $\mathbf{w}_j\mathbf{x}$ for j=1,2,…,m, where m is the number of neurons. In practice, we find it is convenient to normalize the weight vector to constant Euclidean norm. In such a case, the

maximum inner product is equivalent to minimum Euclidean distance between vectors. If we use i(**x**) to denote the index of the neuron that best matches the input vector **x**, we may then determine i(**x**) by using:

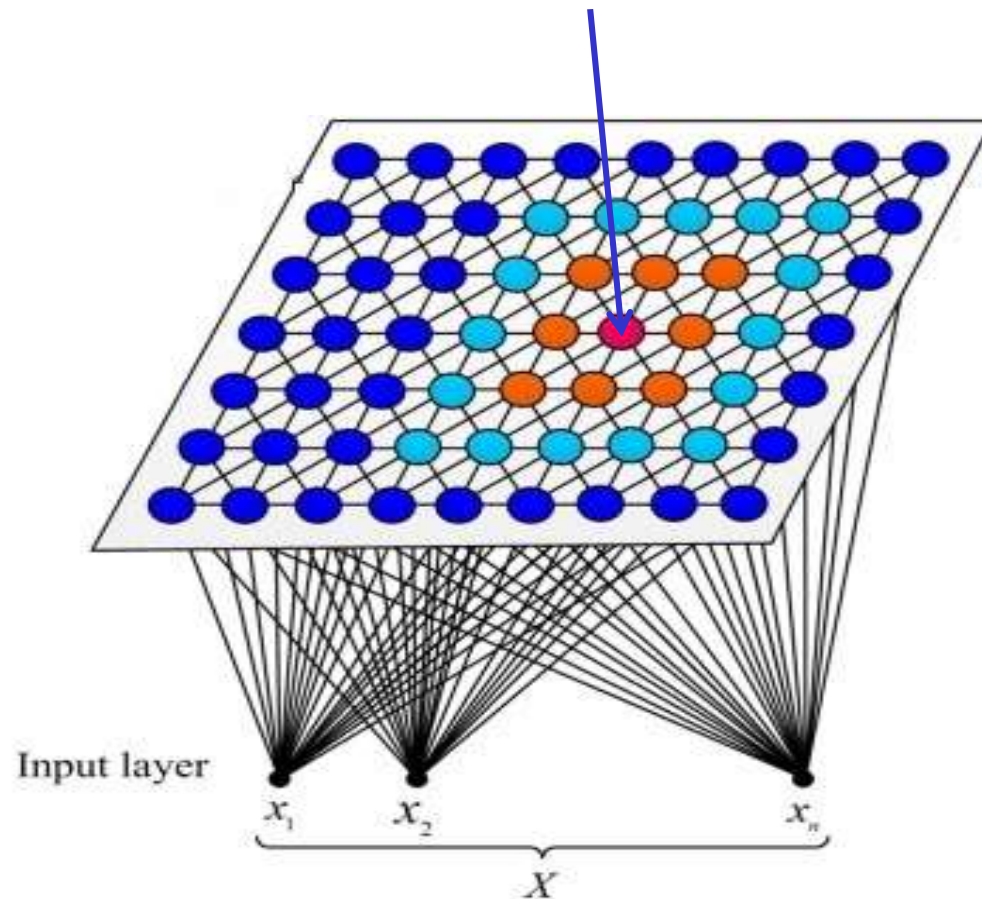$$i(\mathbf{x}) = \arg\min \left\| \mathbf{x} - \mathbf{w}_j^T \right\|$$

Where ||*|| denotes the Euclidean norm of the argument vector. The corresponding neuron is called the winning neuron or best match unit (BMU) for the input vector **x**.

## Cooperation process

The winning neuron locates the center of a topological neighborhood of cooperating neurons. The key question is: how we define a topological neighborhood that is neurobiological correct? Research in neurobiology reveals that when a neuron is

excited, its immediate neighborhood tends to be excited more than those that are far away from it.
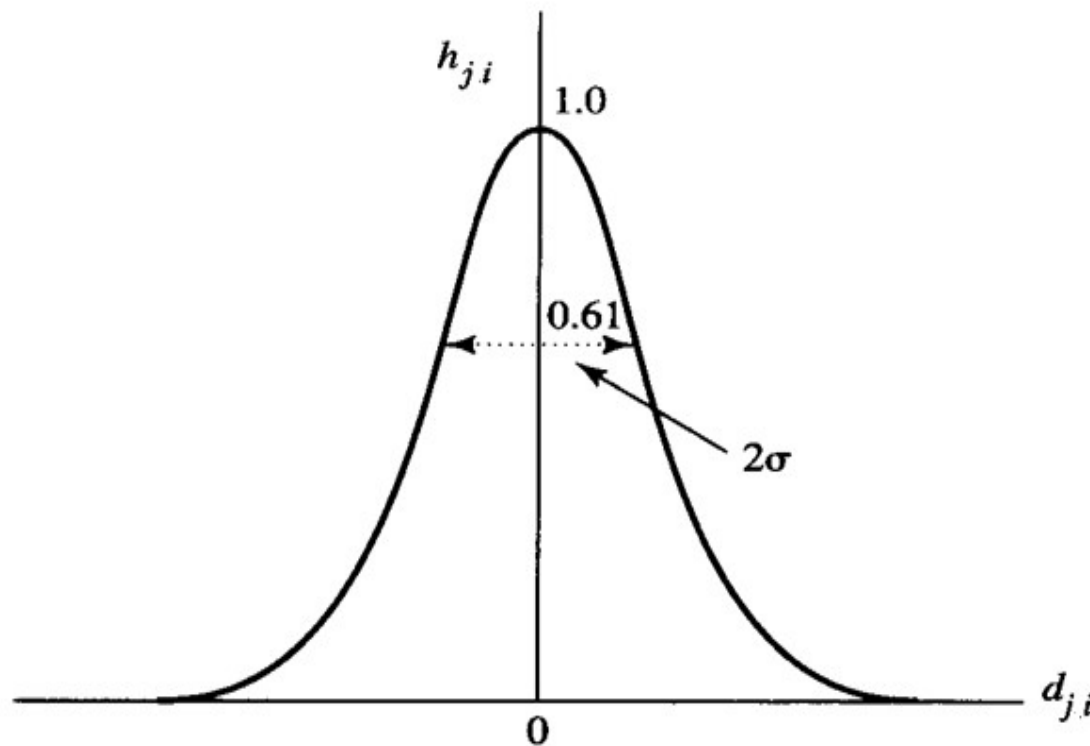
winning neuron
and its neighbourhood



Input layer

$x_1$    $x_2$    $x_n$

$X$

Let $h_{ji}$ denote the topological neighborhood centered on winning neuron i, neuron j is one of the set of excited (cooperating) neurons. The lateral distance between winning neuron i and excited neuron j in the 1D or 2D lattice is denoted by $d_{ji}$. Then we may assume that the neighborhood $h_{ji}$ is a function of the lateral distance $d_{ji}$ that satisfies two distinct requirements:

(1) The topological neighborhood $h_{ji}$ is symmetric about the maximum point defined by $d_{ji}$. In other words, it attains its maximum value at the winning neuron.

(2) The amplitude of the topological neighborhood $h_{ji}$ decreases monotonically with increasing lateral distance $d_{ji}$, and decays to zero when $d_{ji} \to \infty$.

A typical choice of $h_{ji}$ that satisfies these requirements is the Gaussian function:

$$h_{ji} = \exp\left(-\frac{d_{ji}^2}{2\sigma^2}\right)$$

Parameter σ is the width of the Gaussian function, it determines the degree to which excited neurons in the vicinity of the winning neuron participate in the learning process.
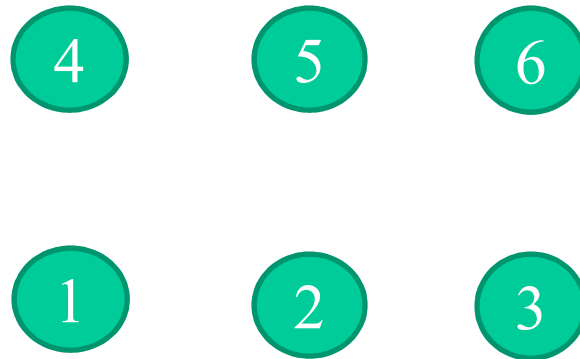
For cooperation among neighboring neurons to hold, it is necessary that the topological neighborhood $h_{ji}$ be dependent on distance between winning neuron $i$ and excited neuron $j$ in the output space, i.e. 2D lattice, rather than the distance measure in the original input (measurement) space.

In the case of a 2D lattice, $d_{ji}$ is defined as:

$$d_{ji} = \left\| \mathbf{r}_j - \mathbf{r}_i \right\|$$

Where $\mathbf{r}_i$ and $\mathbf{r}_j$ denotes the discrete position of winning neuron i and excited neuron j measured in the output space respectively.

Example: 2x3 lattice



Neuron 1 position is $\mathbf{r}_1 = (0,0)$
Neuron 2 position is $\mathbf{r}_2 = (1,0)$
Neuron 6 position is $\mathbf{r}_6 = (2,1)$

The lateral distance between neurons 1 and 6 is:

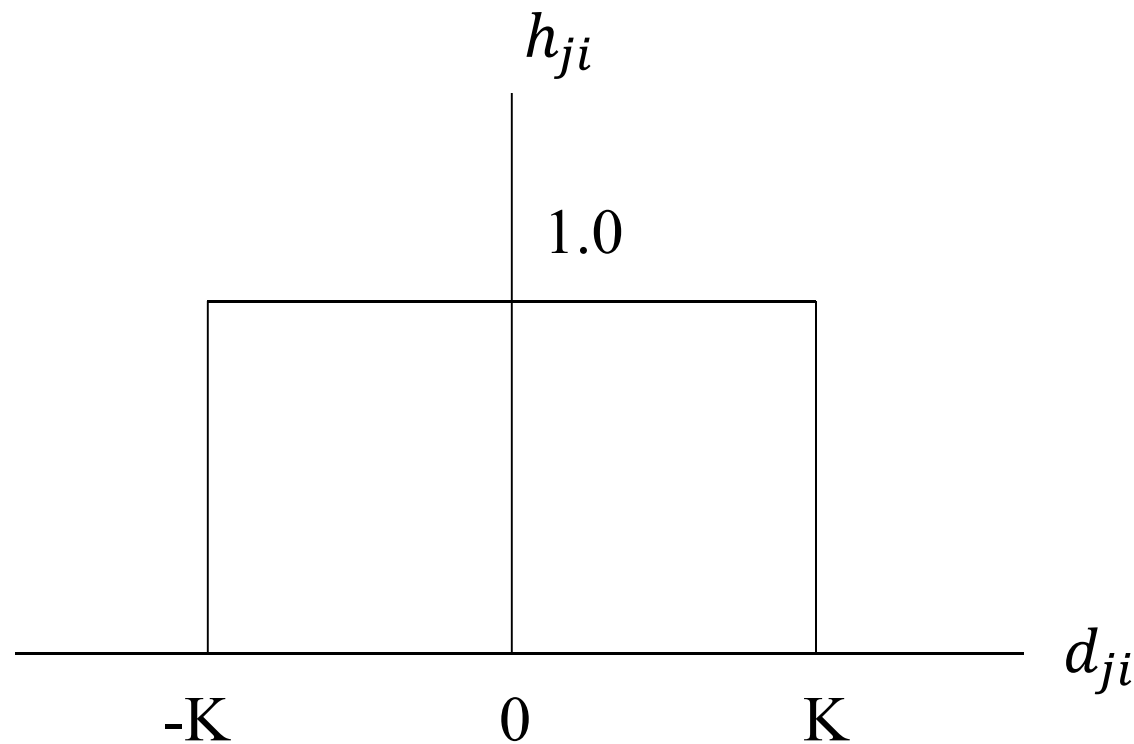$$d_{15} = \|\mathbf{r}_1 - \mathbf{r}_6\| = \sqrt{(2-0)^2 + (1-0)^2} = 2.2361$$

Another unique feature of SOM is that the size of the topological neighborhood shrinks with the learning process. This can be implemented by making the width parameter $\sigma$ of the topological neighborhood function $h_{ji}$ decrease with the learning process. A typical choice for the dependence of $\sigma$ on the learning process is the exponential decay described by:

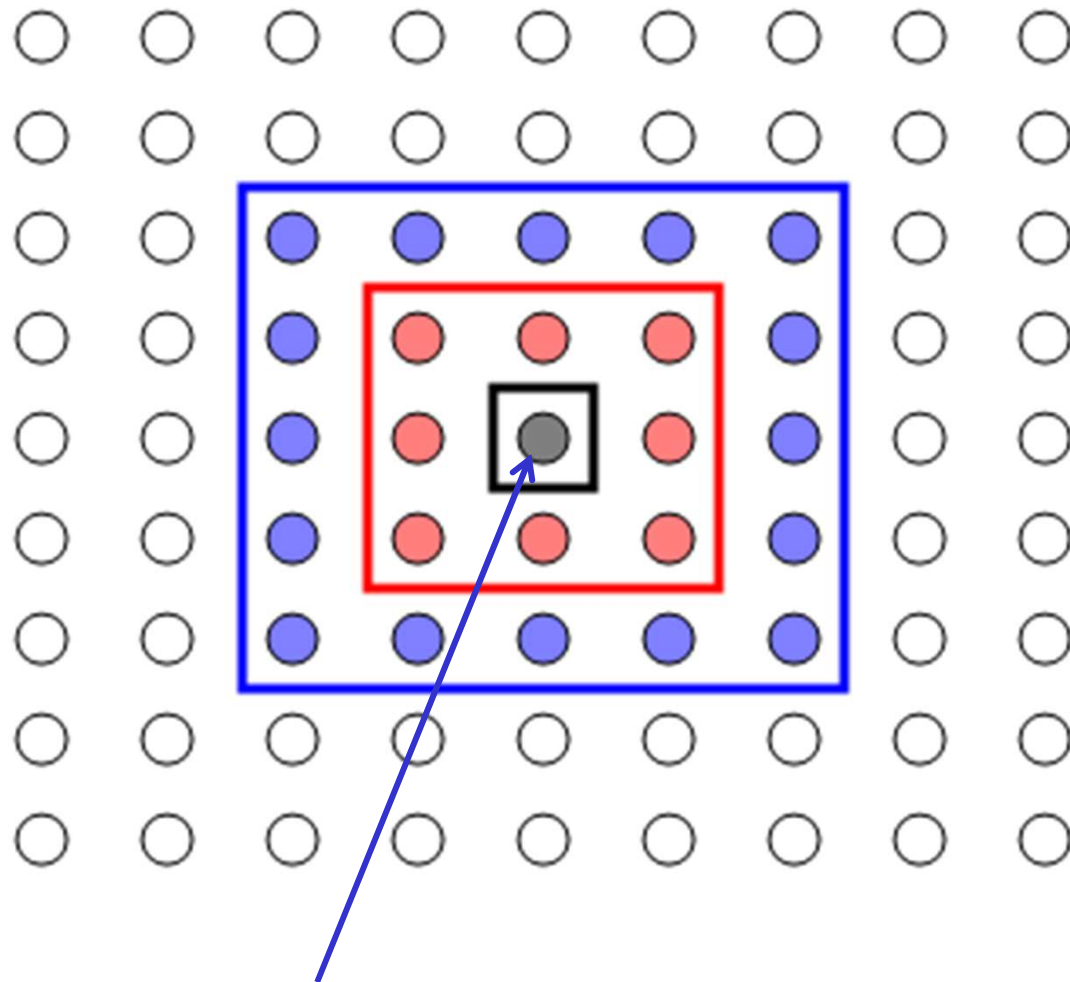$$\sigma(n) = \sigma_0 \exp\left(-\frac{n}{\tau_1}\right)$$

Where $\sigma_0$ is the value at the initiation of the SOM algorithm, and $\tau_1$ is a time constant, n is the number of iterations. Thus, the width and neighborhood shrink with the progress of the learning process:

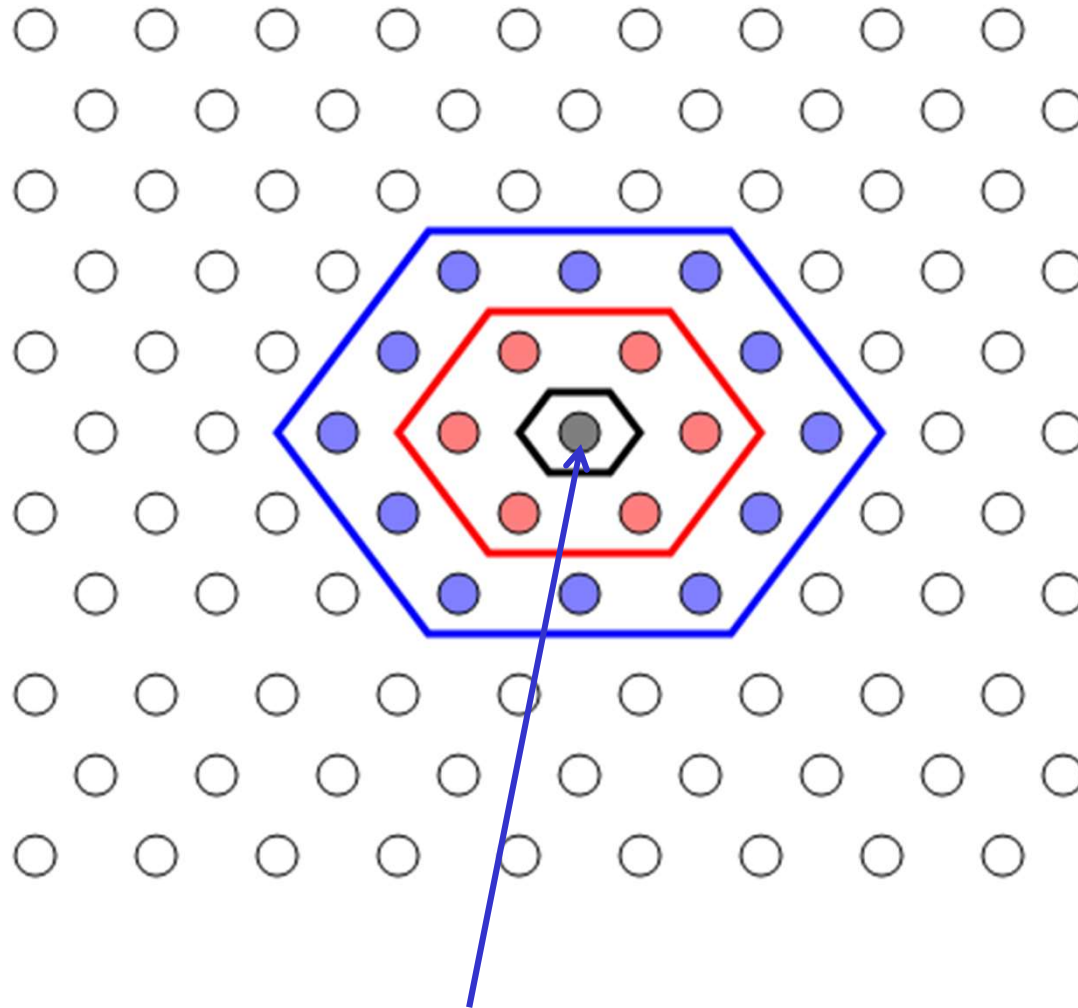$$h_{ji}(n) = \exp\left(-\frac{d_{ji}^2}{2\sigma^2(n)}\right)$$

Besides the Gaussian neighborhood function, hexagon and square neighborhood can also be used in SOM. Neurons inside the defined neighborhood K are excited. In such cases, the $h_{ji}$ is in the following form:

winning neuron and its square neighborhood

winning neuron and its hexagon neighborhood

## Adaptive process

In the weight adaption process, the weight vector $\mathbf{w}_j$ of neuron j in the network is required to change in relation to the input vector $\mathbf{x}$. In the SOM algorithm, the change to the weight vector of neuron j, which is inside the topological neighborhood of winning neuron i, is as follow:
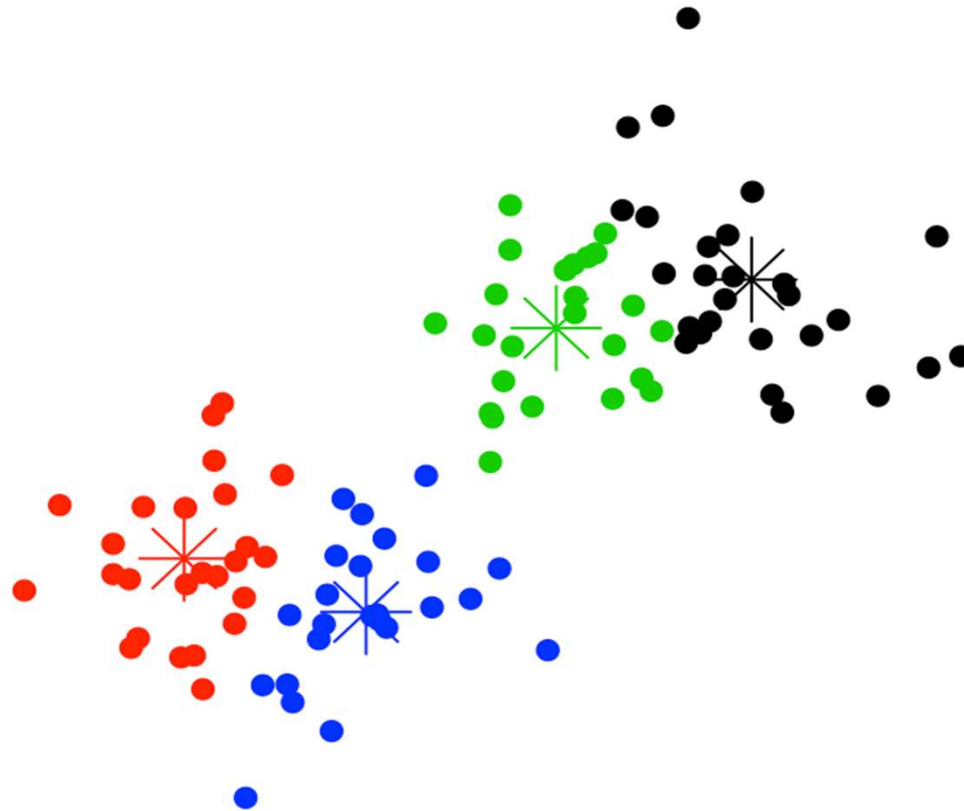
$$\Delta \mathbf{w}_j(n) = \eta(n) h_{ji(x)}(n)[\mathbf{x} - \mathbf{w}_j(n)]$$

Where $\eta(n)$ is the learning rate at iteration n. After change, the weight vector becomes:

$$\mathbf{w}_j(n+1) = \mathbf{w}_j(n) + \eta(n) h_{ji(x)}(n)[\mathbf{x} - \mathbf{w}_j(n)]$$

The above equation has the effect of moving the weight vector of winning neuron and its neighbors toward the input vector $\mathbf{x}$.

Upon repeated presentation of the training data, the weight vectors (denoted by $*$) tend to follow the distribution of the input vectors.

The learning-rate parameter η(n) should be time-varying. It should start at an initial value and then decrease gradually with the increase of iteration number $n$ as shown below:

$$\eta(n) = \eta_0 \exp\left(-\frac{n}{\tau_2}\right)$$

Where $\eta_0$ is the initial value, and $\tau_2$ is another time constant.

## Two phases of the Adaptive process

Starting from an initial state of complete disorder, the SOM algorithm gradually leads to an organized representation of activation patterns drawn from the input space.

The adaptation process of weight vectors in the network can be divided into two phases: an ordering or self-organizing phase followed by a convergence phase as described below:

## (1) Self-organizing

It is in this first phase of the adaptive process that topological ordering of the weight vectors takes place. The ordering phase may take as many as 1000 iterations, or even more, of the SOM algorithm, depending on the number of training samples. In this phase, learning rate parameter and neighborhood functions must be carefully selected.

## Learning Rate

It is found experimentally that the learning rate parameter $\eta(n)$ should begin with a value close to 0.1; and the value decreases gradually, but remains above 0.01. These desirable values are satisfied by the setting $\eta_0=0.1$ and $\tau_1=1000$. Thus, the varying learning rate formula becomes:

$$\eta(n) = 0.1\exp\left(-\frac{n}{1000}\right)$$

## Neighborhood Function

The neighborhood function $h_{ji}(n)$ should initially include almost all neurons in the network centered on the winning neuron i, and then shrink slowly with iterations.

Specifically, during the ordering phase that may occupy 1000 iterations or more, $h_{ji}(n)$ is permitted to reduce to a small value of only a couple of neighboring neurons around the winning neurons, or even to the winning neuron itself. We may set the size of $\sigma_0$ to the radius of the two-dimensional lattice, and set the time constant $\tau_1$ as follow:

$$\sigma(n) = \sigma_0 \exp\left( -\frac{n}{\tau_1} \right)$$

$$\tau_1 = \frac{1000}{\ln \sigma_0}$$

Example: 4x4 two-dimensional lattice

neuron 16

neuron 1

Assuming the position of neuron 1 is (0,0), and the position of neuron 16 is (3,3), then the lateral distance between the two neurons is:

$$d = \sqrt{(3-0)^2 + (3-0)^2} = 4.242$$

The initial value of the width can be set as:

$$\sigma_0 = 2.121$$

## (2) Convergence phase

This second phase of the adaptive process is needed to fine tune the map and therefore provide an accurate statistical quantification of the input space. As a general rule, the number of iterations of the convergence phase must be at least 500 times the number of neurons in the network. Thus, the convergence phase may have to go on for thousands and possibly tens of thousands of iterations.

For good performance, the learning rate parameter $\eta(n)$ should be maintained during the convergence phase at a small value on the order of 0.01, but should not be zero.

We may, for example, simply set the learning rate to a small constant:

$$\sigma(n) = 0.01$$

The neighborhood function $h_{ji}(n)$ should contain only the nearest neighbors of a winning neuron, which may eventually reduce to one or zero neighboring neurons:

$$h_{ji}(n) = \begin{cases} 1 & \text{if neuron j is the winning neuron i} \\ 0 & \text{otherwise} \end{cases}$$

## Summary of the SOM algorithm

There are four basic steps involved in the algorithm, namely, initialization, sampling, competition and updating (adaption and cooperation). The three steps after initialization are repeated until the map function is completed. The algorithm is summarized as follows.

(1) Initialization. Choose random values for the initial weight vectors $\mathbf{w}_j(0)$. The only restriction on the initialization is that the initial weight vectors must be different from each other, and it is desirable to keep the magnitude of the weights small.

(2) Sampling. Draw a sample $\mathbf{x}$ from the input distribution with a certain probability. Usually, the $\mathbf{x}$ is drawn from the given training samples set.

(3) Competition. Find the winning neuron $i(\mathbf{x})$ at the step $n$ using the minimum-Euclidean distance criterion:

$$i(\mathbf{x}) = \arg \min \left\| \mathbf{x} - \mathbf{w}_j(n) \right\|$$

where j=1,2,…,N.

(4) Updating (cooperation and adaption). Adjust the weight vectors of all neurons using the updating formula:

$$\mathbf{w}_j(n+1) = \mathbf{w}_j(n) + \eta(n)h_{ji(x)}(n)[\mathbf{x} - \mathbf{w}_j(n)]$$

(5) Continuation. Repeat steps (2)-(4) until the stopping criterion is satisfied. The stopping criterion can be:

    (i)  Pre-defined iteration number;
    (ii) No noticeable changes in the map

The SOM neural network training process goes through two phases. In the self-organizing phase, random values are assigned to the initial weights. In the convergence phase, the weights learned in the self-organizing phase will be used as initial weights.

## Example 1 (for demonstration purpose only)

Given 4 training samples:

$$\mathbf{x}_1 = [1,1,0,0]^{\mathrm{T}} \qquad\qquad \mathbf{x}_2 = [0,0,0,1]^{\mathrm{T}}$$

$$\mathbf{x}_3 = [1,0,0,0]^{\mathrm{T}} \qquad\qquad \mathbf{x}_4 = [0,0,1,1]^{\mathrm{T}}$$

We wish to find 2 clusters of the training samples. Suppose the learning rate is

$$\eta(0) = 0.1$$
$$\eta(n+1) = \eta(0)\exp(-n/1000)$$

The neighborhood function is so set that only the winning neuron is updated with its weights at each step.

**Step 1**: initialization. Initialize two weight vectors:

$$\mathbf{w}_1(0) = [0.2, 0.6, 0.5, 0.9]^{\mathrm{T}}$$

$$\mathbf{w}_2(0) = [0.8, 0.4, 0.7, 0.3]^{\mathrm{T}}$$

**Step 2**: randomly select a sample, say $\mathbf{x}_1$, we have:

$$d_1 = \|\mathbf{w}_1(0) - \mathbf{x}_1\| = \sqrt{(0.2-1)^2 + (0.6-1)^2 + (0.5-0)^2 + (0.9-0)^2} = 1.3638$$

$$d_2 = \|\mathbf{w}_2(0) - \mathbf{x}_1\| = 0.99$$

Neuron 2 is the winning neuron, and its weights are updated:

$$\mathbf{w}_2(1) = \mathbf{w}_2(0) + \eta(0)[\mathbf{x}_1 - \mathbf{w}_2(0)] = [0.82, 0.46, 0.63, 0.27]^T$$

The weights of neuron 1 remain unchanged:

$$\mathbf{w}_1(1) = \mathbf{w}_1(0) = [0.2, 0.6, 0.5, 0.9]^T$$

**Step 3**: randomly select another sample, say $x_2$, we have:

$$d_1 = \|\mathbf{w}_1(1) - \mathbf{x}_2\| = 0.8124$$

$$d_2 = \|\mathbf{w}_2(1) - \mathbf{x}_2\| = 1.3468$$

Neuron 1 wins, and its weights are updated:

$$\eta(1) = \eta(0) \times \exp(-1/1000) = 0.0999$$

$$\mathbf{w}_1(2) = \mathbf{w}_1(1) + \eta(1)[\mathbf{x}_2 - \mathbf{w}_1(1)]$$

$$= [0.18, 0.5401, 0.4501, 0.91]^T$$

The weights of Neuron 2 remain unchanged:

$$\mathbf{w}_2(2) = \mathbf{w}_2(1) = [0.92, 0.76, 0.28, 0.12]^T$$

The above process is continued until the stopping criterion is satisfied.

# Example 2

Initial weight vectors:

# Weight vectors at the end of the self-organizing phase

# Weight vectors at the end of the convergence phase

## Example 3

Consider a SOM neural network with 100 neurons arranged in a 2D lattice with 10 rows and 10 columns. The network is trained with a two-dimensional input vector **x**, whose elements $x_1$ and $x_2$ are uniformly distributed in a region of square:

$$0 < x_1 < 1$$

$$0 < x_2 < 1$$

# 10000 input vectors uniformly distributed in the region

# After training, the 100 weight vectors divide the region uniformly

## Properties of the Feature Map

Once the SOM algorithm is converged, the feature map produced by the algorithm displays important statistical characteristics of the input space.

(1) **Approximation of the input space**. The feature map represented by a set of weight vectors in the output space provide a good approximation to the input space. Thus, a large set of input vectors can be represented by a smaller set of prototypes.

In the following example, we have 1000 data points. If we want to use a small set of prototypes, say 16, 25, or 100, to represent the distribution of the 1000 data points, we may use the SOM algorithm to find the prototypes as shown in the following figures.

# Example 4



1000 data points

16 prototypes
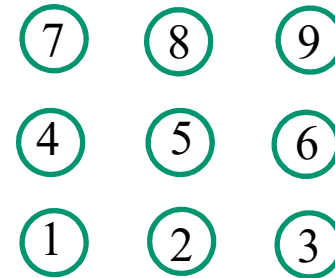
25 prototypes

100 prototypes

43

(2) **Topological ordering.** The feature map computed by the SOM algorithm is topologically ordered in the sense that the spatial location of a neuron in the lattice corresponds to a particular domain of the input patterns.
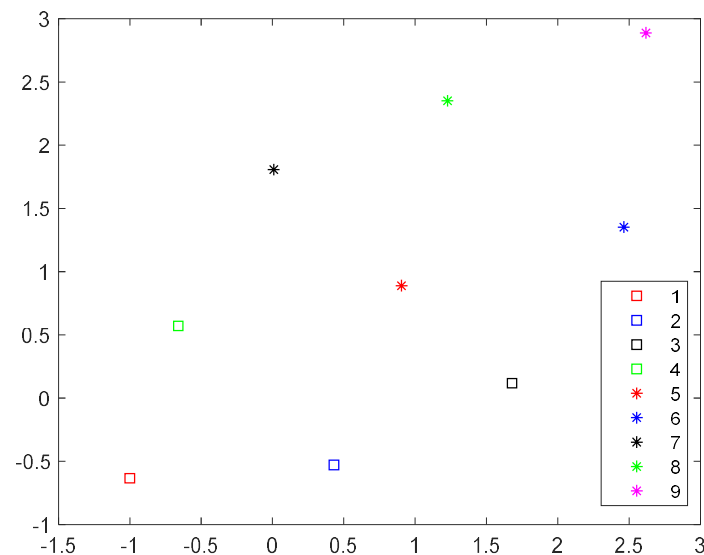
**Example 5**
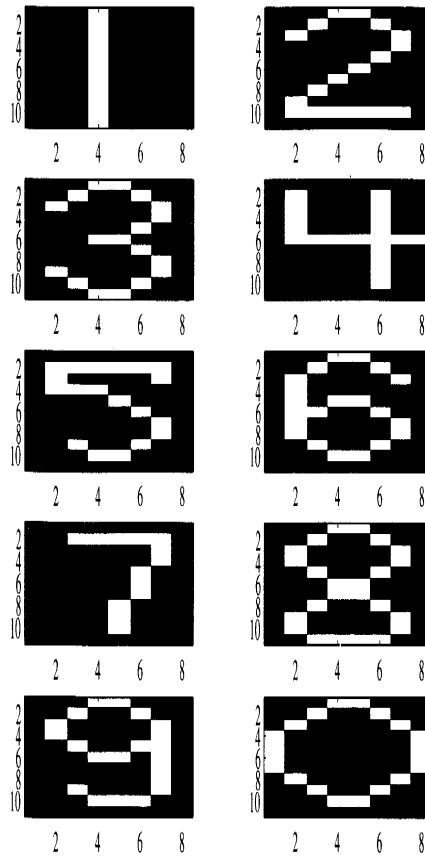
The location of 9 neurons in 2D lattice:

⑦ ⑧ ⑨

④ ⑤ ⑥

① ② ③

The location of the weight vectors of the 9 neurons in input space (after training with the 1000 data points)
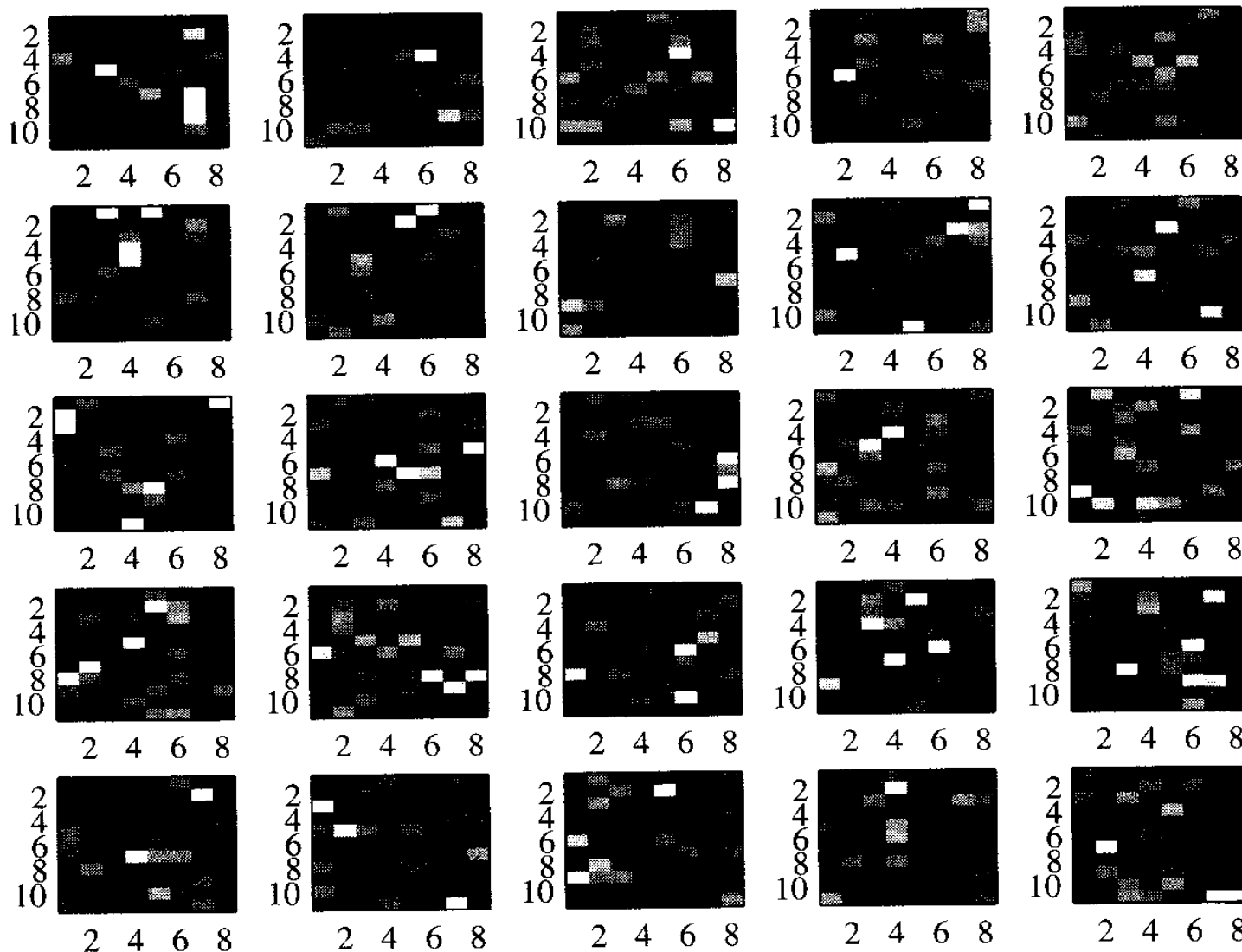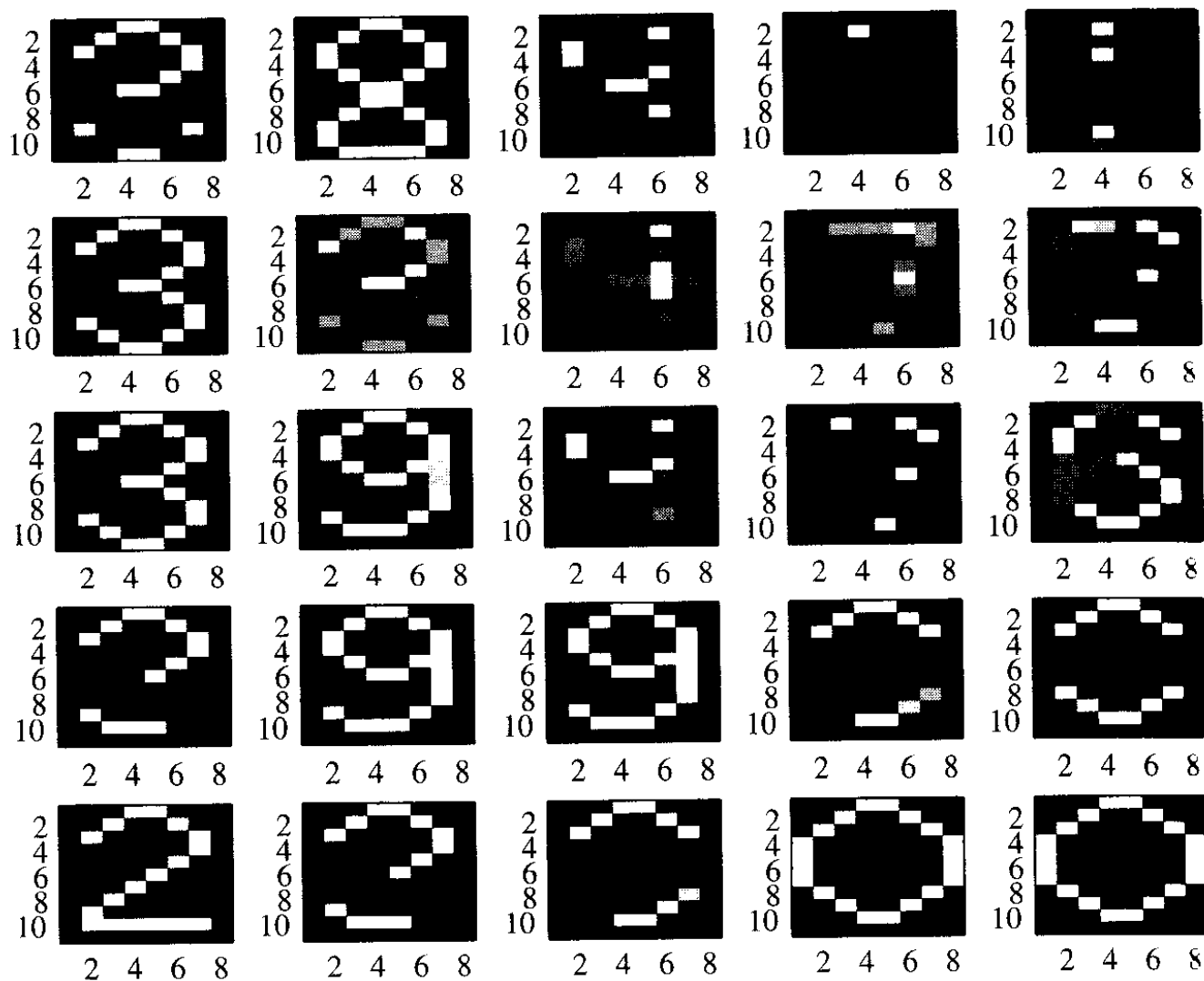
## Example 6

Consider 10 digits, represented by 11×8 dot matrices, we build a 5 ×5 2-D SOM map.

# Initial weights are generated randomly and are shown below:

The resulted SOM, with weights viewed as digits, is shown below. Obviously, weights of adjacent neurons are similar.
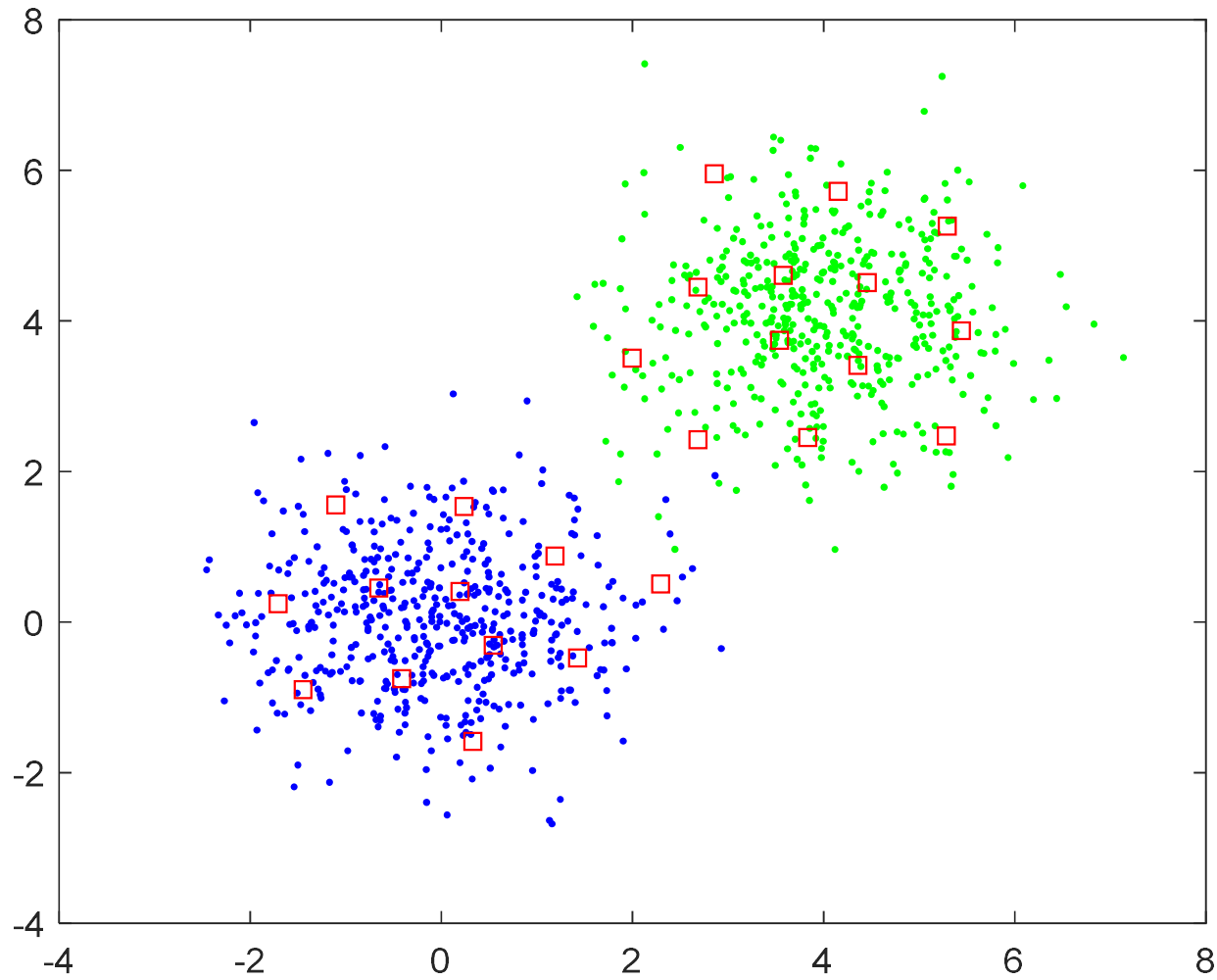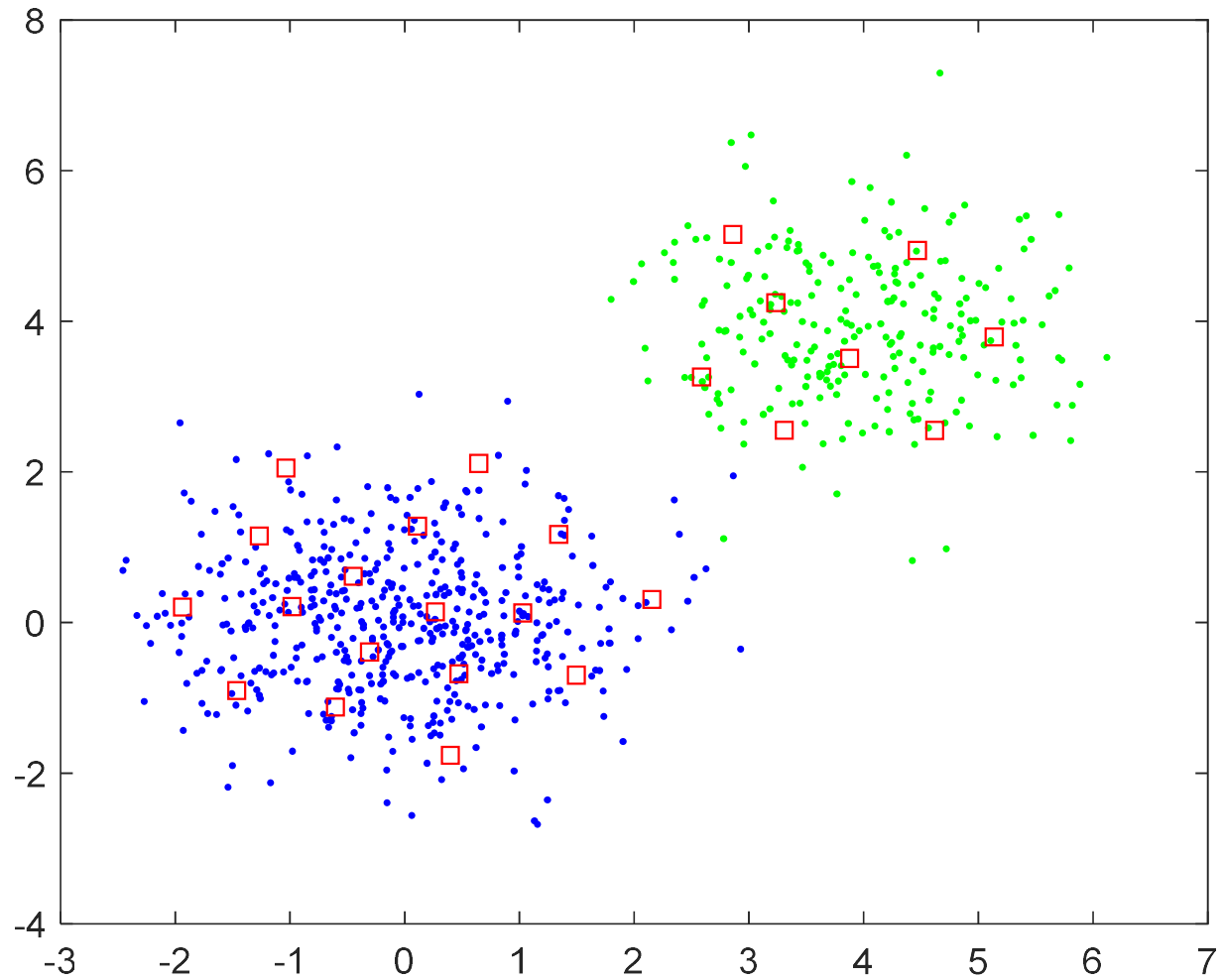
(3) **Density Matching**. The map reflects variations in the statistics of the input distribution: regions in the input space from which sample are drown with a high probability of occurrence are mapped onto larger domains of the output space, and therefore with better resolution than regions in the input space from which samples are drawn with a low probability of occurrence.

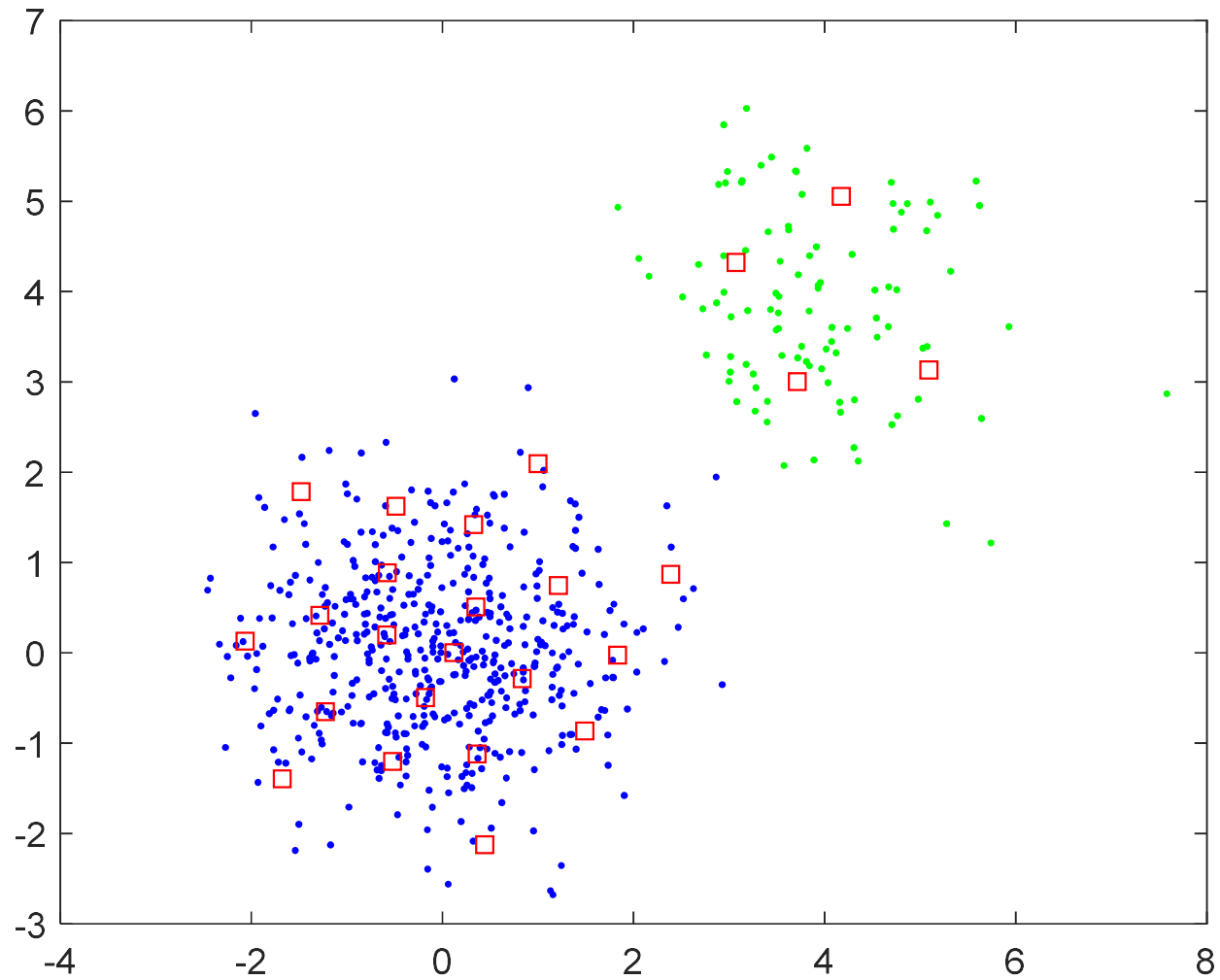Consider the example shown below, we have samples (data points) from two clusters.
1) Case 1: 500 samples from each of the two clusters (cluster 1 to cluster 2 ratio is 1:1)
2) Case 2: 500 samples from cluster 1, 250 samples from cluster 2 (cluster 1 to cluster 2 ratio is 2:1)
3) Case 3: 500 samples from cluster 1, 100 samples from cluster 2 (cluster 1 to cluster 2 ratio is 5:1)

Sample ratio is 1:1, prototype ratio is 11:13

Sample ratio is 2:1, prototype ratio is 17:8

Sample ratio is 5:1, prototype ratio is 21:4