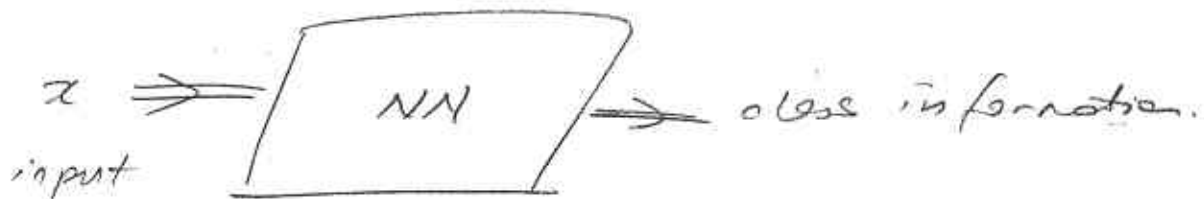


# UNSUPERVISED LEARNING

- CLUSTER DETECTION.

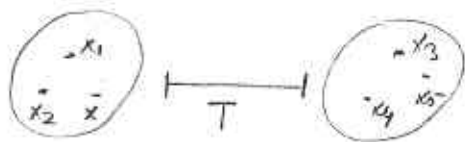


input space



PROBLEM: Given the training set  $\{x_1, \dots, x_N\}$ , but without the class information, find a NN architecture and learning rule, which detects the clusters.

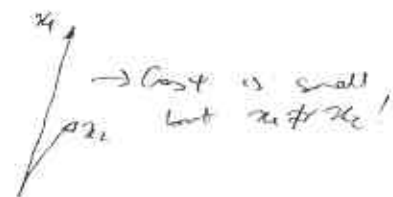
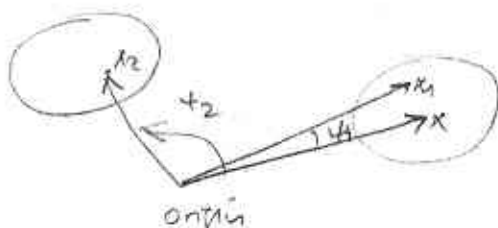
To define a cluster, we need to use a "similarity measure" to check the presence of patterns within a class. The most common measure is the Euclidean distance between patterns.



$$\|x - x_i\|^2 = (x - x_i)^t (x - x_i)$$

Here we also determine a smallest distance  $T$  to discriminate the clusters.

Another measure might be the cosine of the angle between the vectors.

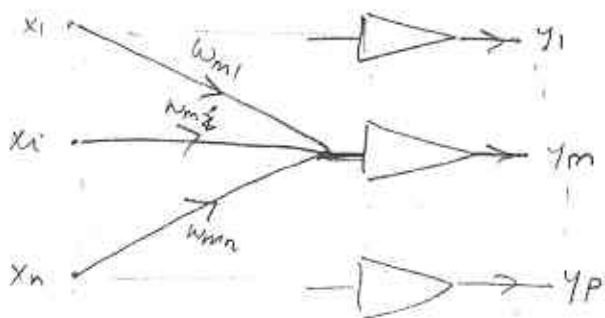


$$x^t x_i = \|x\| \cdot \|x_i\| \cdot \cos \phi_i \quad \cos \phi_i = \frac{x^t x_i}{\|x\| \cdot \|x_i\|}$$

The smaller the  $\cos \phi$ , the closer the vectors  $x$  and  $x_i$ . However, the vectors  $x$  and  $x_i$  should be comparable.

### WINNER-TAKES-ALL LEARNING.

$\{x_1, \dots, x_n\}$  are given as inputs to be classified. The number of distinct classes  $p$  is also known. The network is called Kohonen network.



$$Y = \begin{pmatrix} y_1 \\ y_m \\ y_p \end{pmatrix} = W X.$$

$$W = \begin{pmatrix} w_{m1}^t \\ w_{m2}^t \\ w_{pn}^t \end{pmatrix} = \begin{pmatrix} w_{m1} & w_{m2} & w_{mn} \\ w_{p1} & w_{p2} & w_{pn} \end{pmatrix}$$

$$w_{m1} = \begin{pmatrix} w_{m1} \\ w_{m2} \\ w_{mn} \end{pmatrix}$$

In the beginning, all weights are chosen to be random and having unit length:

$$\hat{w}_i = \frac{w_i}{\|w_i\|} \quad i=1, \dots, p$$

When an input pattern  $x$  is given, the procedure is to find the output, for which the following is minimum:

$$\|x - \hat{w}_m\| = \min_i \|x - \hat{w}_i\| \Rightarrow \text{x should be normalized as well}$$

Hence, the network finds the weight vectors which is closest to the input  $x$ .

$\|x - \hat{w}_m\|^2 = (x^t - \hat{w}_m^t)(x - \hat{w}_m) = \|x\|^2 - 2\hat{w}_m^t x + 1$   
 ( $\|\hat{w}_m\|=1$ ). Obviously  $\|x - \hat{w}_m\|$  is min. if  $w_m^t x$  is max. Hence equivalently, we may look for:

$$w_m^t x = \max_i w_i^t x$$

Hence, when  $x$  is present, compute  $Wx = \begin{pmatrix} w_1^t x \\ w_2^t x \\ \vdots \\ w_p^t x \end{pmatrix}$  and

find the output which is maximum. After selecting this output, corresponding weight vector should be adjusted so that it becomes more similar to the input pattern (i.e.  $\|x - \hat{w}_m\|$  becomes smaller). A standard technique to guarantee that is to update along the negative gradient of the following function.

$$E = \frac{1}{2} \|x - \hat{w}_m\|^2 = \frac{1}{2} (\|x\|^2 - 2w_m^t x + w_m^t w_m)$$

$$\frac{\partial E}{\partial w_m} = -(x - \hat{w}_m)$$

Hence the update rule is:

$$\Delta w_m = \alpha (x - \hat{w}_m) = -\alpha \frac{\partial E}{\partial w_m} \quad \begin{matrix} m: \text{winning neuron} \\ i: \text{non winning neuron} \end{matrix}$$

$$w_m(k+1) = w_m(k) + \Delta w_m, \quad w_m(k+1) \leftarrow \frac{w_m(k+1)}{\|w_m(k+1)\|} \text{ (normalize!)}$$

\* the rule should increase the chances of winning by  $m^{\text{th}}$  neuron when this update rule is used.

$$\Rightarrow w_m^t x \leq (w_m + \Delta w_m)^t x$$

$$\Rightarrow \Delta w_m^+ x \geq 0 \Leftrightarrow x(x - w_m)^+ \geq 0$$

$$\Leftrightarrow x^T x - w_m^+ x = \|x\|^2 - w_m^+ x \geq 0$$

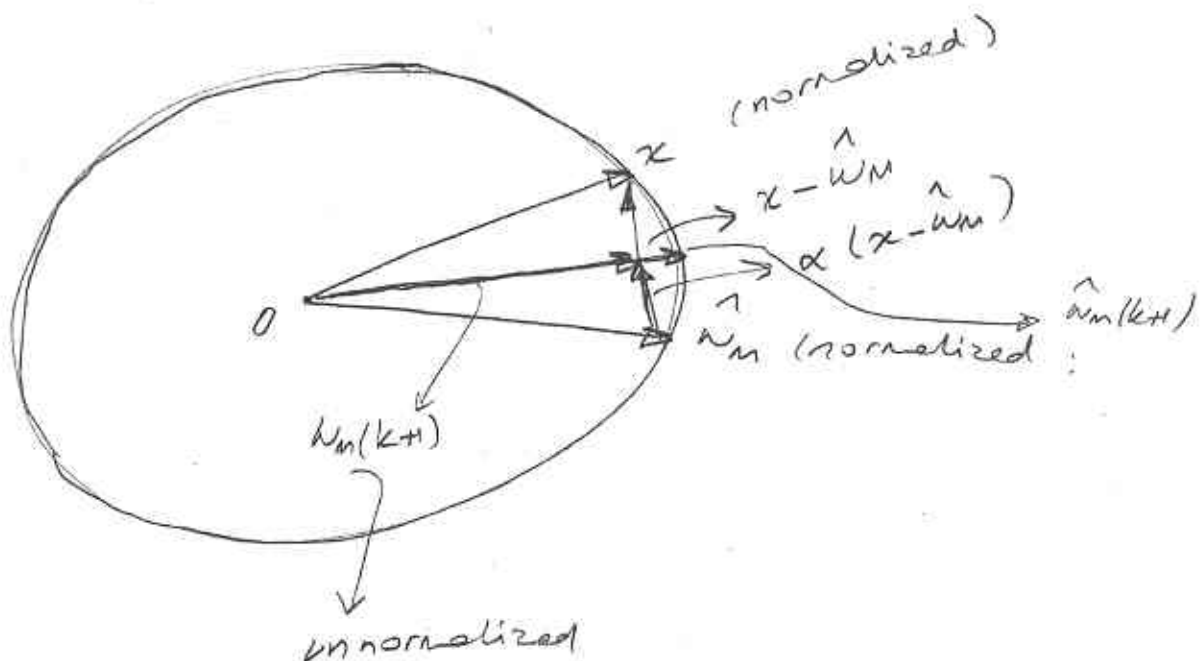
$$w_m^+ x = \|w_m\| \cdot \|x\| \cdot \cos(\theta)$$



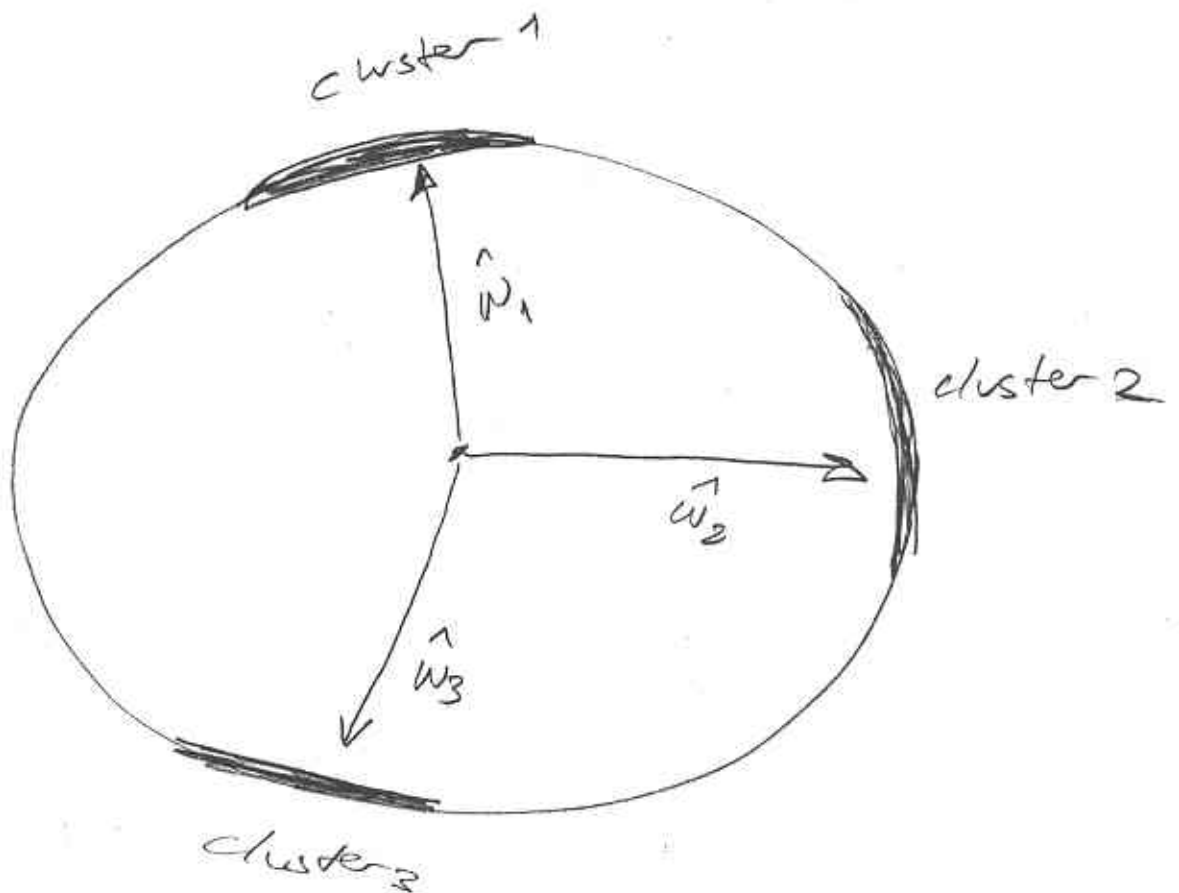
$$\|x\|^2 - w_m^+ x = \|x\| (\|x\| - \|w_m\| \cos(\theta)) \geq 0$$

If  $\|x\|$  and  $\|w_m\|$  are normalized this is always true!

GEOMETRIC IDEA:



IDEA:



- success may depend on how the clusters are separated
- initial weights

Trick: If we know the cluster patterns, initial weights may be chosen as the cluster centers

ie.  $C_i = \{x_{i1}, \dots, x_{im}\}$

choose  $w_i = \frac{1}{m} \{x_{i1} + \dots + x_{im}\}$   $i = 1, 2, \dots, p$

# ALGORITHM:

Given the training set:

$$\{x_1, x_2, \dots, x_N\} \in \mathbb{R}^n$$

(in our case  $N=100$  (patterns)  
 $n=100$   $100 \times 100$  grid)

And  $p$  distinct classes ( $p=10$  in our case)

- Choose the initial weights randomly and uniformly distributed on the unit hypercube in  $\mathbb{R}^n$ .

(suggestion:

$$W = \begin{pmatrix} w_{11} & w_{12} & \dots & w_{1n} \\ w_{21} & w_{22} & \dots & w_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{p1} & w_{p2} & \dots & w_{pn} \end{pmatrix}$$

choose  $w_{ij} \in \{-1, 1\}$  uniformly distributed and randomly.

$$\Rightarrow \hat{w}_i = \frac{w_i}{\|w_i\|} \quad \|w_i\| = \sqrt{w_{i1}^2 + \dots + w_{in}^2}$$

- Normalize the pattern set.

$$y_i = \frac{x_i}{\|x_i\|} \quad (\text{since } x_{ij} = \pm 1 \rightarrow \|x_i\| = \sqrt{1 + \dots + 1} = \sqrt{100} = 10)$$

For each pattern  $y_i$  in the pattern set

$$(i) \rightarrow (i=1, 2, \dots, N)$$

- find arg max  $\sum_{k=1}^p w_k y_k$

(i.e.  $w_j$  → find max entry)

let  $m$  be the max index. (winning neuron)

- update only  $w_m$

$$w_m(i+1) = w_m(i) + \alpha (y_i - w_m(i))$$

- Do not update the remaining neuron weights.

- Normalize  $w_m(i+1)$

$$\hat{w}_m(i+1) \leftarrow \frac{w_m(i+1)}{\|w_m(i+1)\|}$$

- Go to (1)

if all patterns are exhausted

— How to reduce  $\alpha$ ?

$0.1 < \alpha < 0.7$  initially

as epochs increase  $\alpha \rightarrow 0.01$  (or even smaller)

This could be done in various ways.

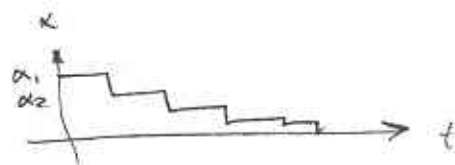
• In discrete epoch steps. i.e.

let  $t_1, t_2, \dots, t_n$  be discrete epoch numbers

$t$ : current epoch number.

$t_1 \leq t < t_2 \rightarrow \alpha = \alpha_1$

$t_2 \leq t < t_3 \rightarrow \alpha = \alpha_2$



Or:

$$\alpha = \alpha_0 + \alpha_1 e^{-t/\tau}$$

$\tau$ : # of epochs

$\alpha_0 = 0.01$

$0.1 \leq \alpha_1 \leq 0.7$

— IF CONVERGED

— Calibration:

Now it is not known which output corresponds to which class. To determine this, we could do calibration: as follows:

— Select candidates from the training set representing classes. i.e. let  $y_i \in C_i \Rightarrow$  apply  $A \Rightarrow$  find the max. output  $\Rightarrow$  This output represents class  $C_i$ .

Do this for any other class as well.

TRICK For initial weights, instead of choosing  $w_i$  randomly, we could choose  $w_i$  as the arithmetic average of patterns in class  $C_i$ .

KOHONEN'S SELF ORGANIZING FEATURE MAPS

In its basic form, we have  $p$  output neurons representing  $p$  classes. We could also increase the number of neurons, and allow the output neurons to compete among themselves.

procedure is the same as before:

Given the training patterns and  ~~$p$~~  output neurons

$$\{y_1, y_2, \dots, y_q\}$$

$p$ : # of classes  
 $I$ : # of output neurons

in epoch  $t$   
 - Find the max. output:

$$\underline{I > P}$$

$$\arg \max_{1 \leq k \leq I} w_k^T y$$

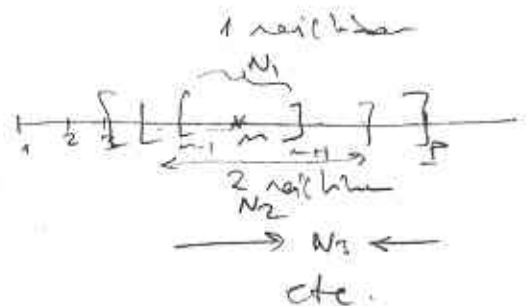
Let  $m$  be the winning neuron

Define a neighborhood of winning neuron  $N_m(t)$   
 ( $t$ : epoch number)

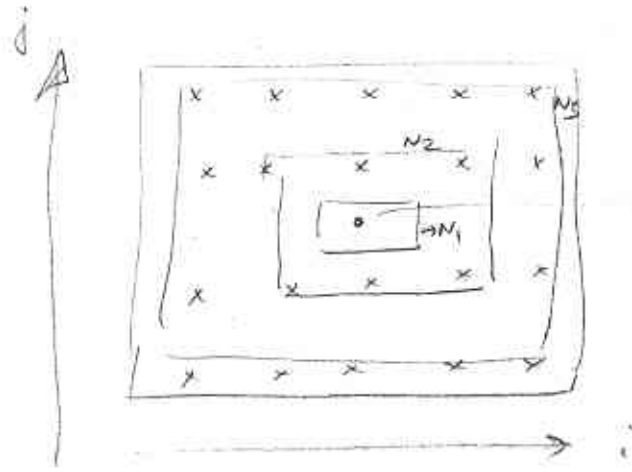
$$w_j(n+1) = w_j(n) + \alpha(t) (y_j - w_j(n)) \quad j \in N_m(t)$$

i.e. all neurons in the neighborhood of the winning neuron are updated.

How to choose the neighborhoods of neurons?  
 That depends on the positions of output neurons.



If the outputs are selected as 2-dimensional



→ winning neuron.

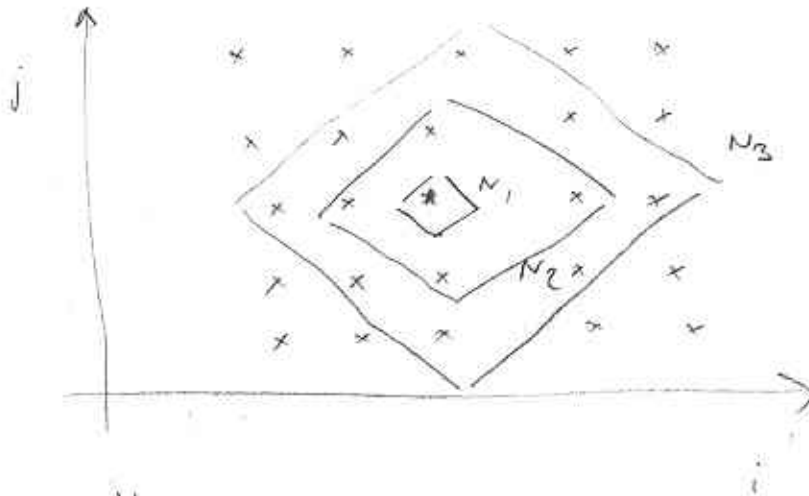
$(m_1, m_2)$  location

$N_1 \rightarrow$  only  $(m_1, m_2)$

$N_2 \rightarrow \{(i, j) \mid \begin{matrix} i = m_1 - 1, m_1, m_1 + 1 \\ j = m_2 - 1, m_2, m_2 + 1 \end{matrix}\}$

$N_3 \rightarrow \{(i, j) \mid \begin{matrix} i = m_1 - 2, \dots, m_1 + 2 \\ j = m_2 - 2, \dots, m_2 + 2 \end{matrix}\}$

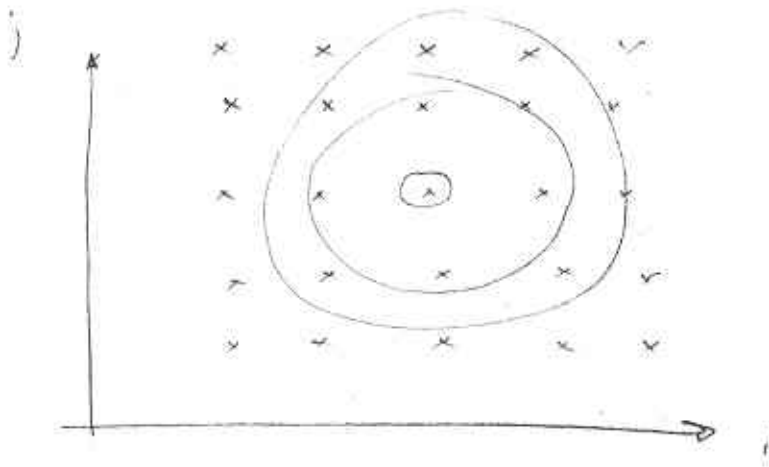
rectangle



winning neuron  $(m_1, m_2)$

$N_2 = \{(i, j) \mid |i + j| \leq L\}$

diamond



circular neighborhood.

- As a result, as epoch number  $t$  increases, we decrease both the neighborhood size and  $\alpha$ .

- Use predetermined epoch numbers  $t_1, t_2, t_3, \dots$  and decrease the neighborhood size as well as  $\alpha$

as  $t \rightarrow \infty$  on  $t_1, t_2, \dots$

$i = 0$

$1 \leq t \leq t_1$

$N_n = N_3$

$\alpha = \alpha_3$

$t_1 < t \leq t_2$

$N_n = N_2$

$\alpha = \alpha_2$

$\vdots$

OR  $\alpha$   $t$  automatically:

$$w_j(n+1) = w_j(n) + \alpha(t) \cdot e^{-\frac{\|r_j - r_m\|}{r(t)}} \quad (f - w_j(n)) \quad j \in N_m(t)$$

$r_j$ : position of  $j$ th neuron w.r.t. winning neuron

$$\alpha(t) = \alpha_0 e^{-t/\tau_2}$$

$$r(t) = r_0 e^{-t/\tau_1}$$

- After convergence  
Do calibration