

CSE 473
Pattern Recognition

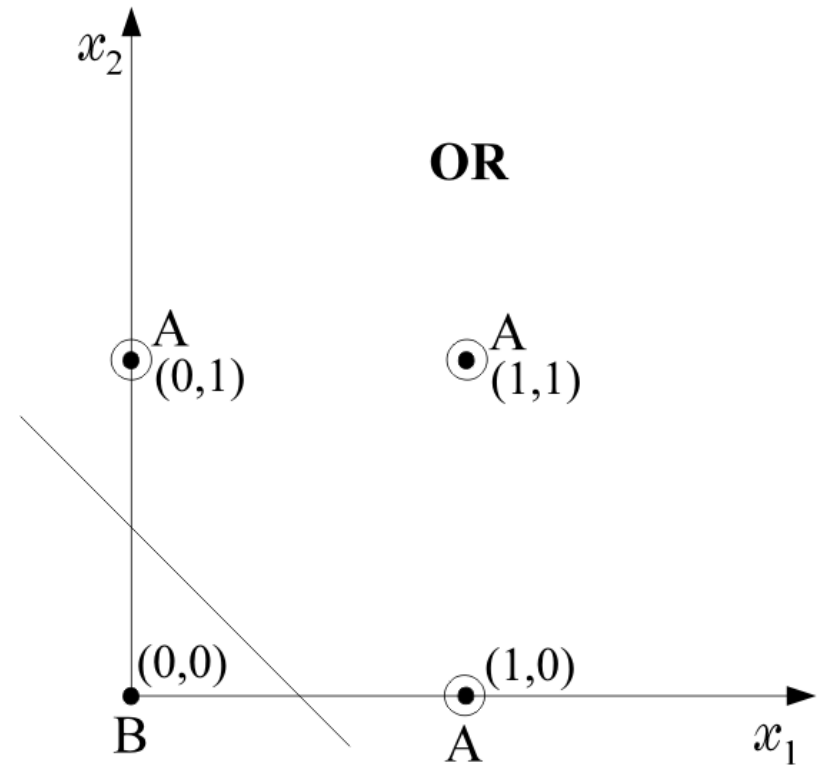
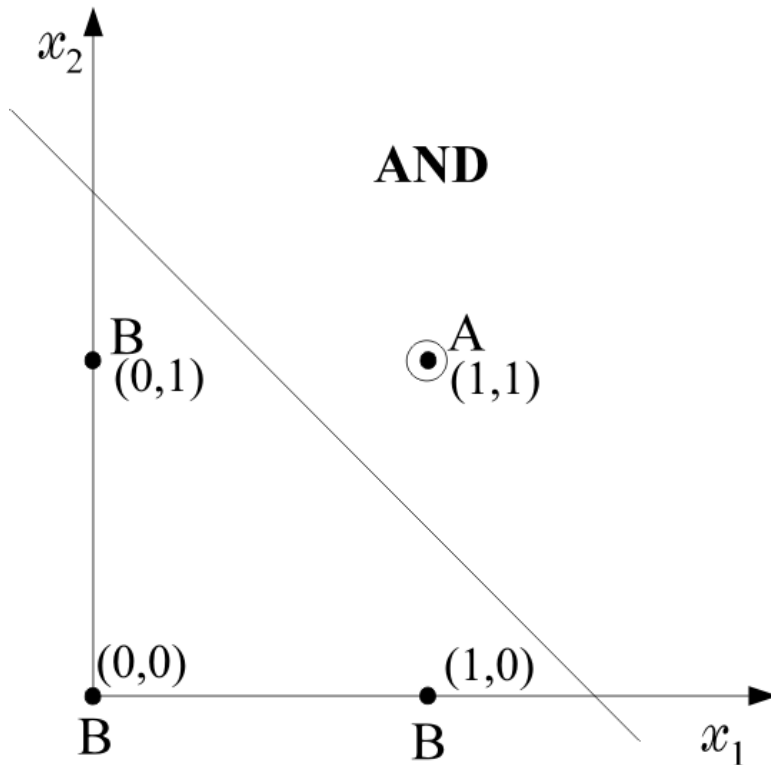
Non-Linear Classifier

Review of Perceptron's Capability

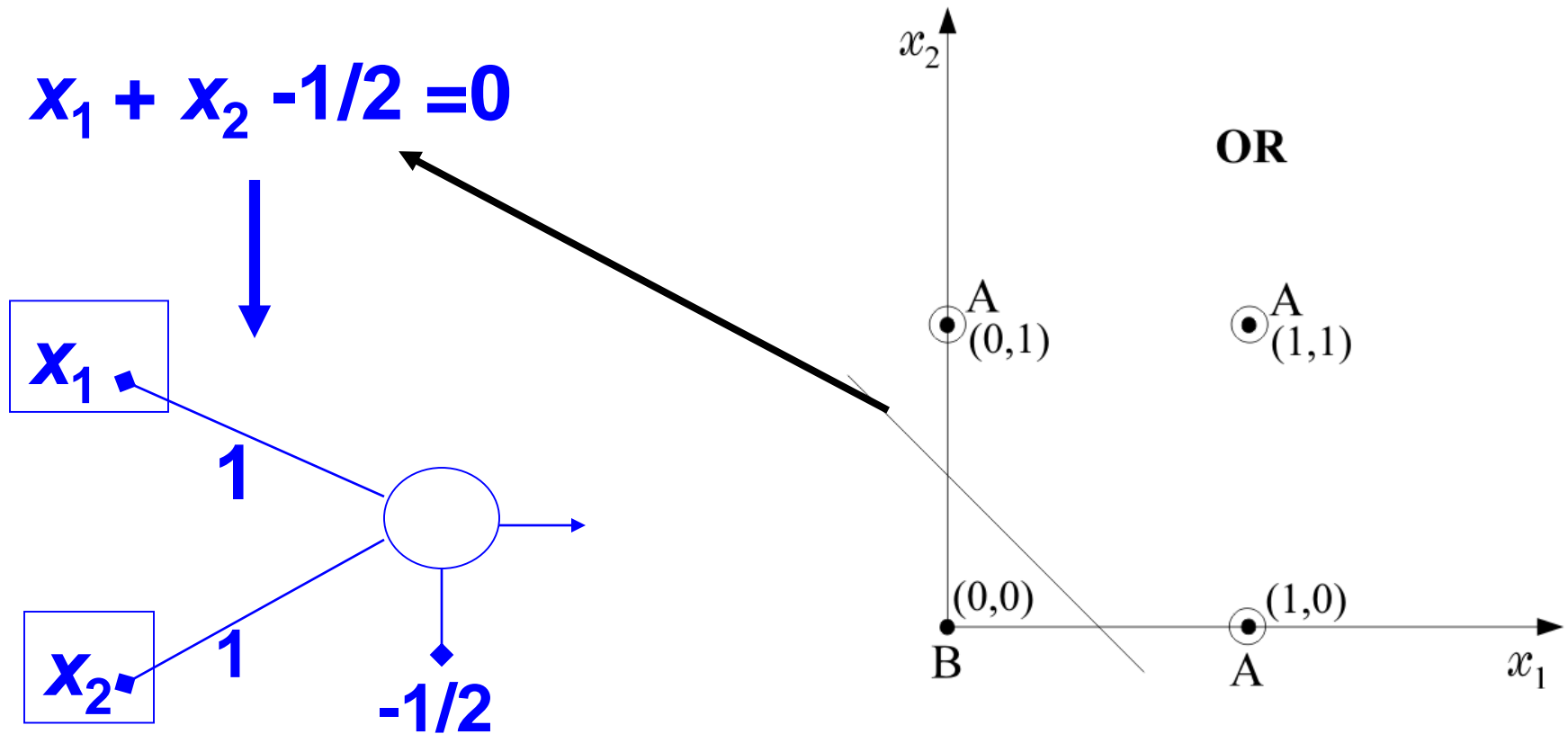
Recall the **AND** or **OR** functions

x_1	x_2	AND	Class	OR	Class
0	0	0	B	0	B
0	1	0	B	1	A
1	0	0	B	1	A
1	1	1	A	1	A

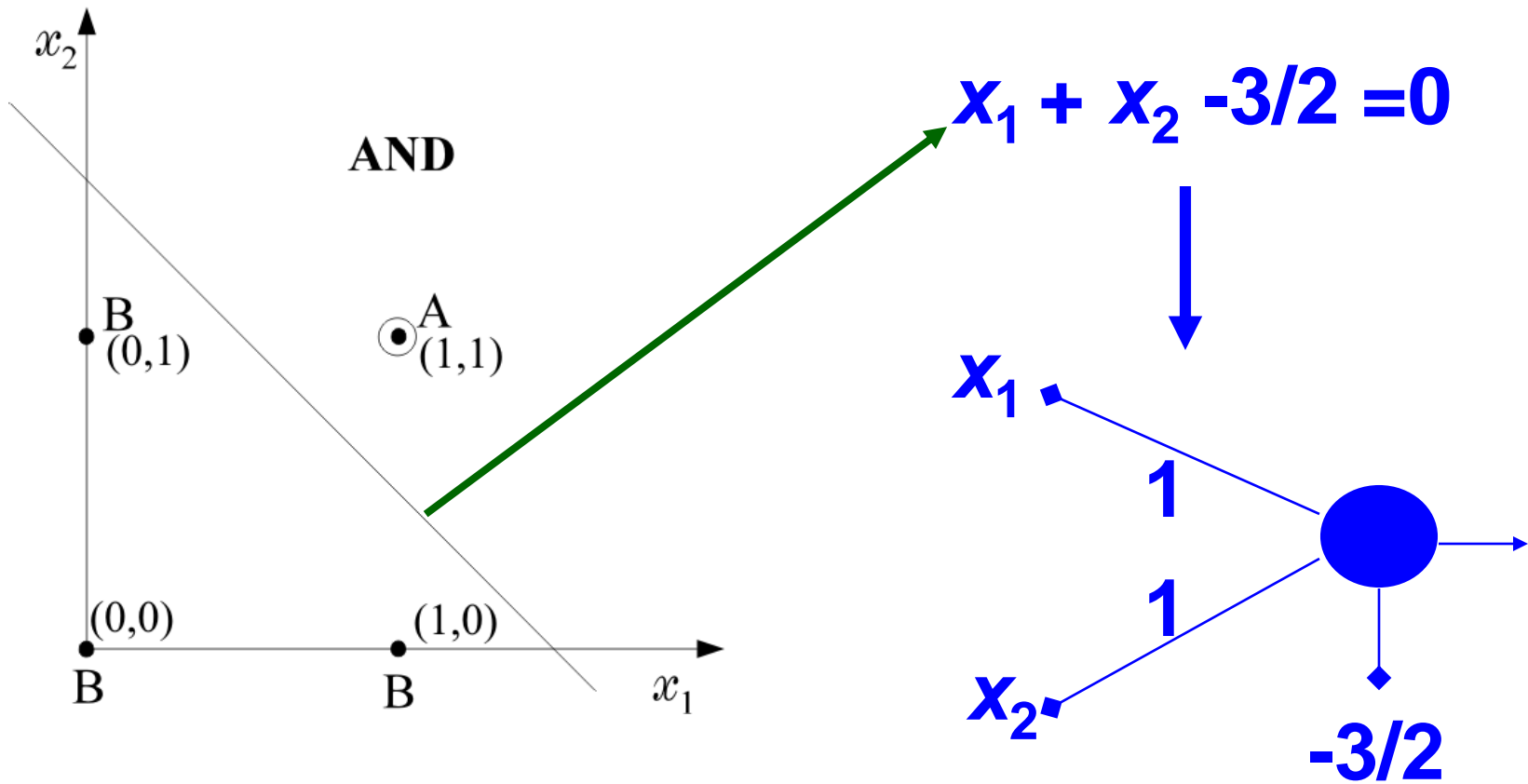
Review of Perceptron's Capability



Review of Perceptron's Capability



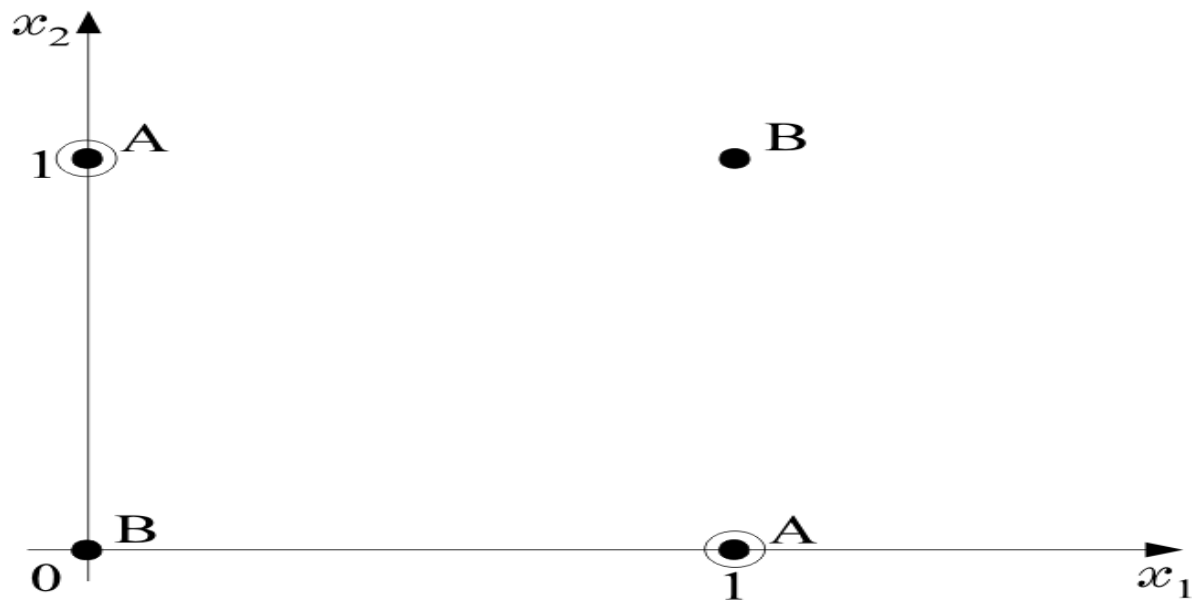
Review of Perceptron's Capability



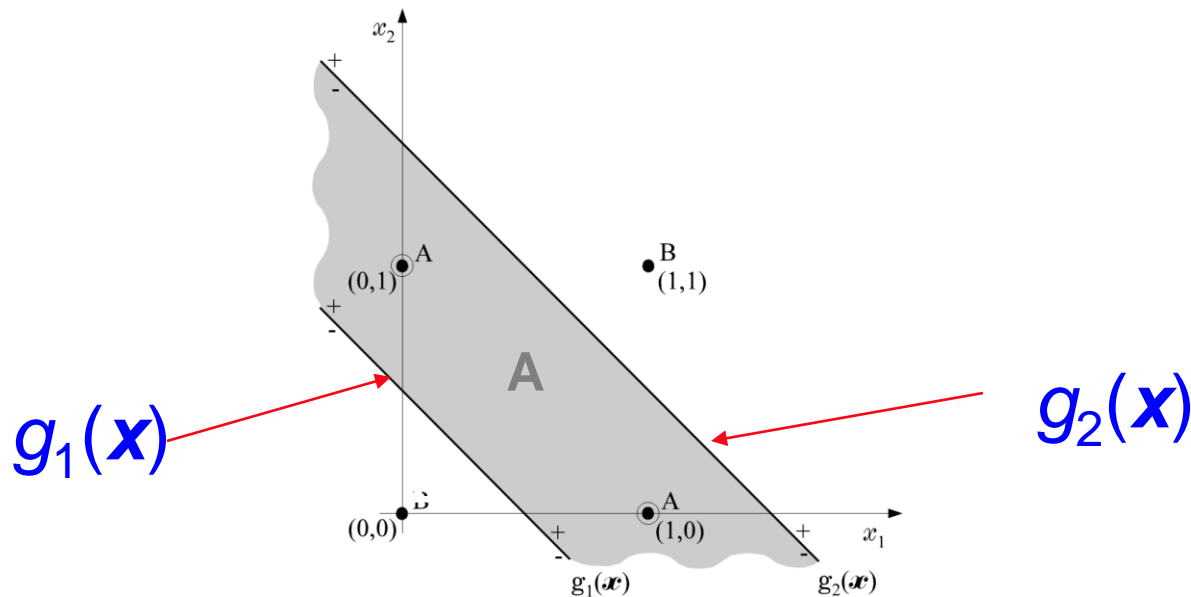
Review of Perceptron's Capability

Now recall
the **XOR**
function

x_1	x_2	XOR	Class
0	0	0	B
0	1	1	A
1	0	1	A
1	1	0	B



Review of Perceptron's Capability



- Each of them is realized by a perceptron.

$$y_i = f(g_i(\underline{x})) = \begin{cases} 0 \\ 1 \end{cases} \quad i = 1, 2$$

- Find the position of \underline{x} w.r.t. both lines, based on the values of y_1, y_2 .

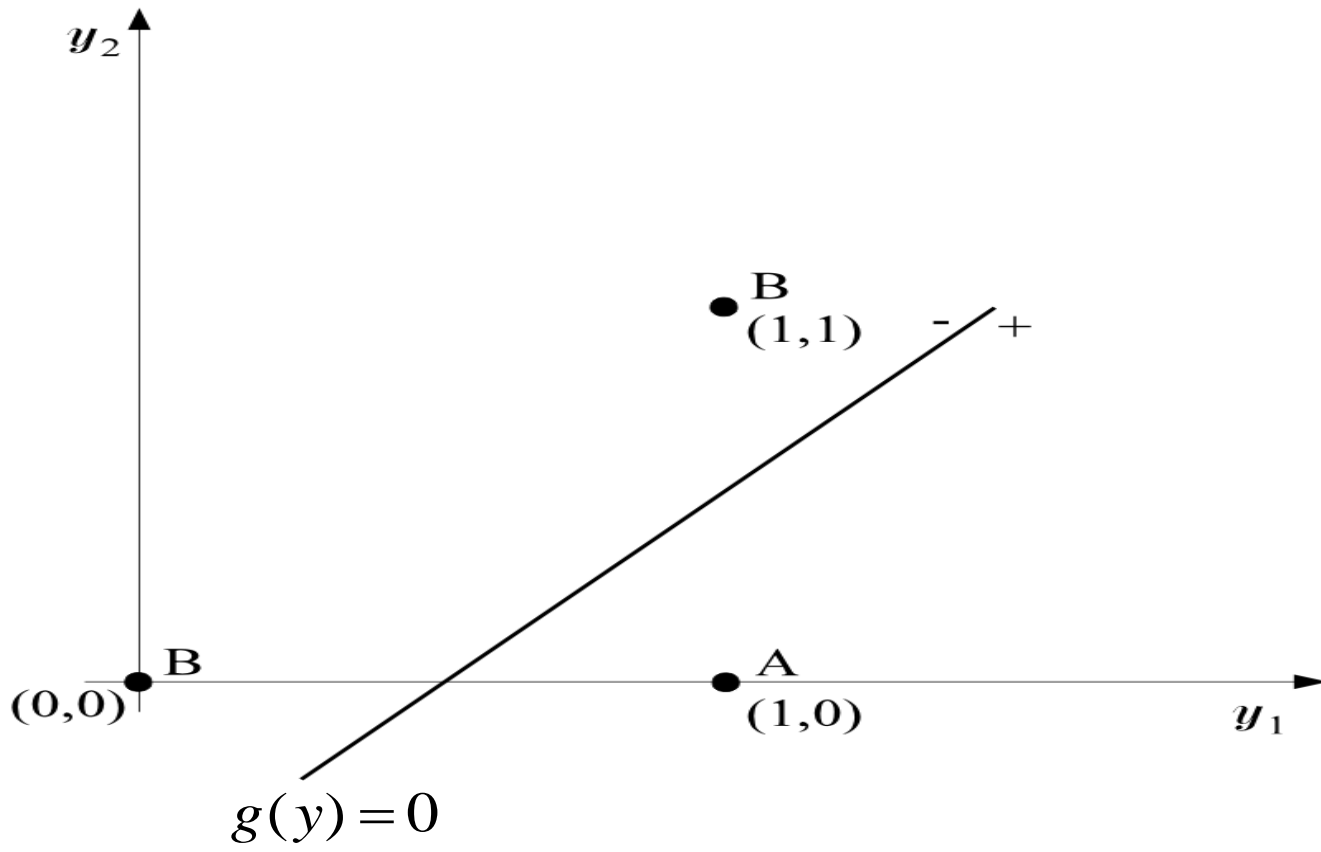
Review of Perceptron's Capability

1 st phase				2 nd phase
x_1	x_2	y_1	y_2	
0	0	0(-)	0(-)	B(0)
0	1	1(+)	0(-)	A(1)
1	0	1(+)	0(-)	A(1)
1	1	1(+)	1(+)	B(0)

- Equivalently: The computations of the first phase perform a mapping $\underline{x} \rightarrow \underline{y} = [y_1, y_2]^T$

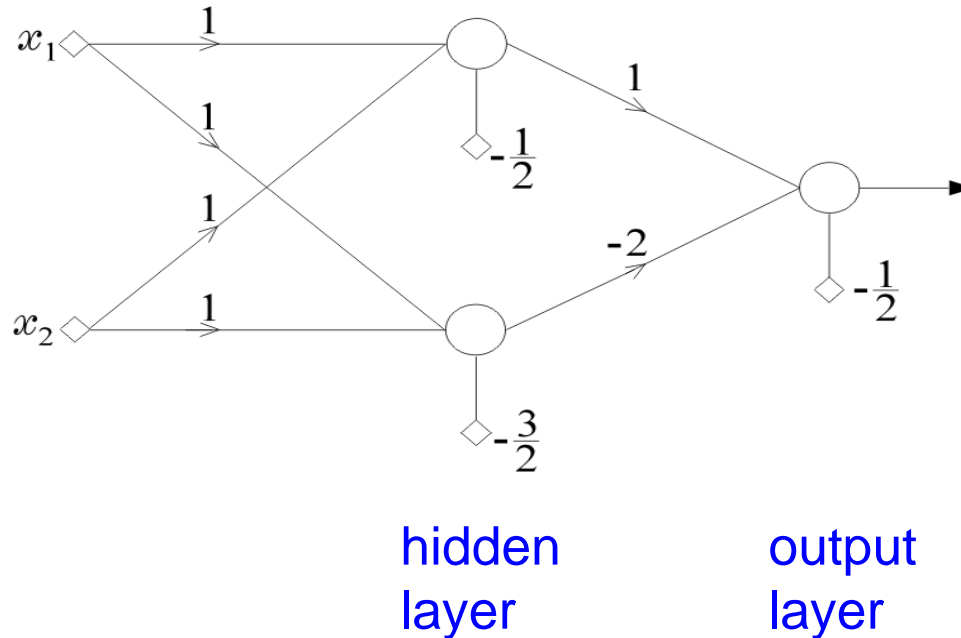
Review of Perceptron's Capability

The decision is now performed on the transformed \underline{y} data.



This can be performed via a second line, which can also be realized by a perceptron.

Two Layer Perceptron



nodes realizes
hyper planes:

$$g_1(\underline{x}) = x_1 + x_2 - \frac{1}{2} = 0$$

$$g_2(\underline{x}) = x_1 + x_2 - \frac{3}{2} = 0$$

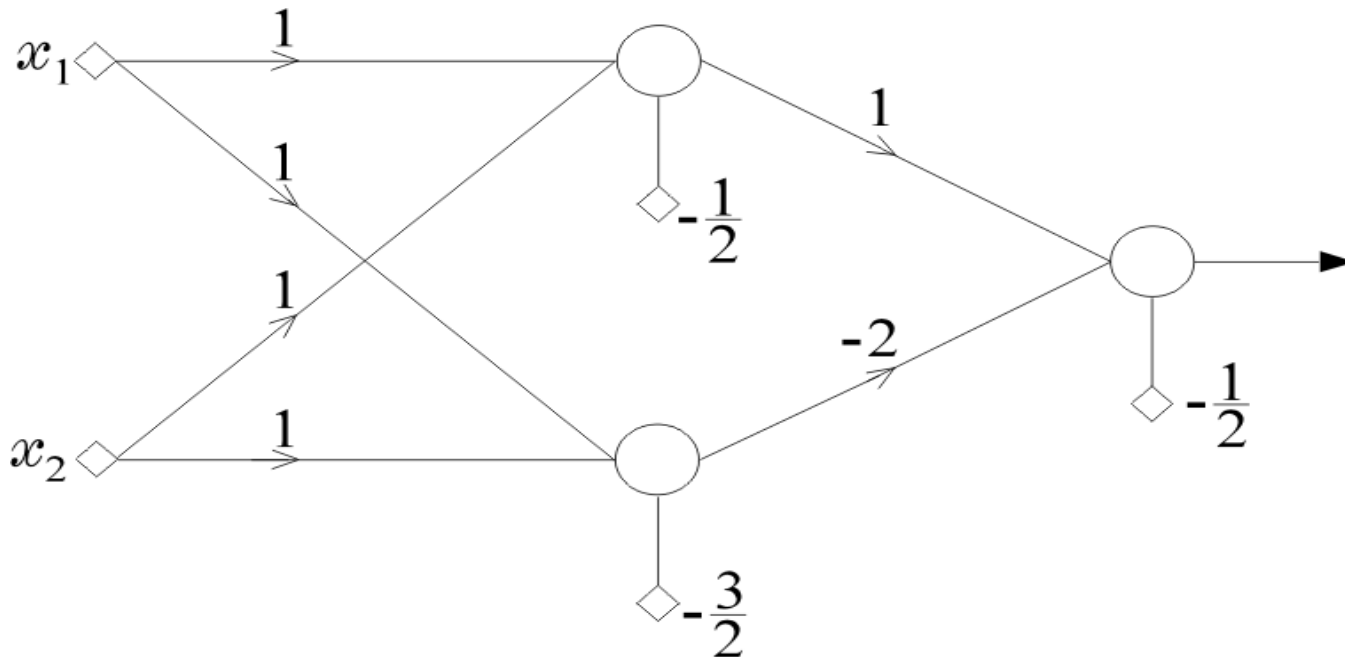
$$g(\underline{y}) = y_1 - 2y_2 - \frac{1}{2} = 0$$

Activation
function:

$$f(.) = \begin{cases} 0 \\ 1 \end{cases}$$

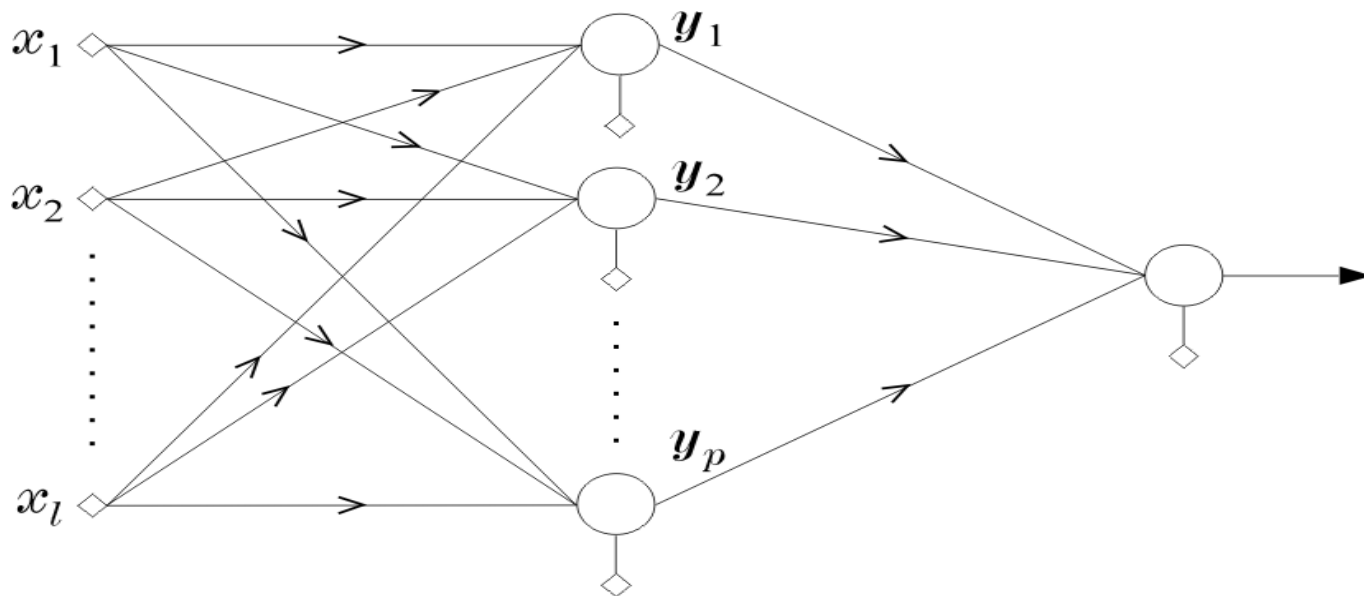
Classification Capabilities of Two Layer Perceptron

- The mapping performed by the first layer neurons is onto the vertices of the unit side square, e.g., $(0, 0)$, $(0, 1)$, $(1, 0)$, $(1, 1)$.



Classification Capabilities of Two Layer Perceptron

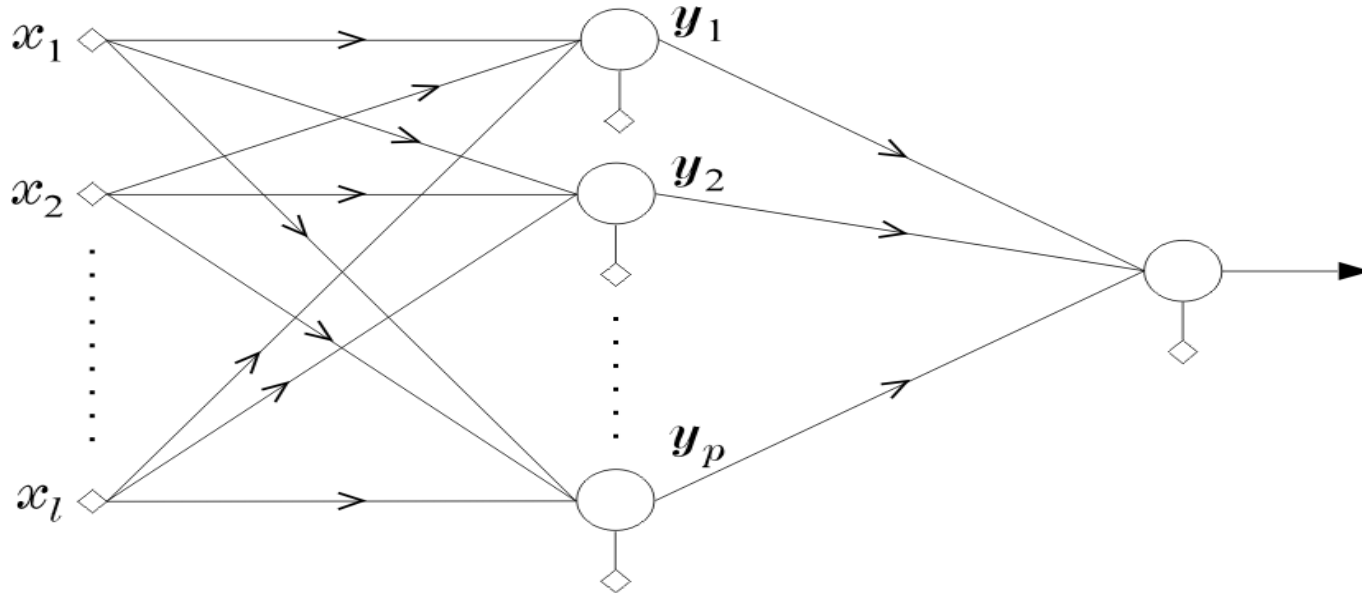
- Consider a more general case,



$$\underline{x} \in R^l$$

$$\underline{x} \rightarrow \underline{y} = [y_1, \dots, y_p]^T, y_i \in \{0, 1\} \quad i = 1, 2, \dots, p$$

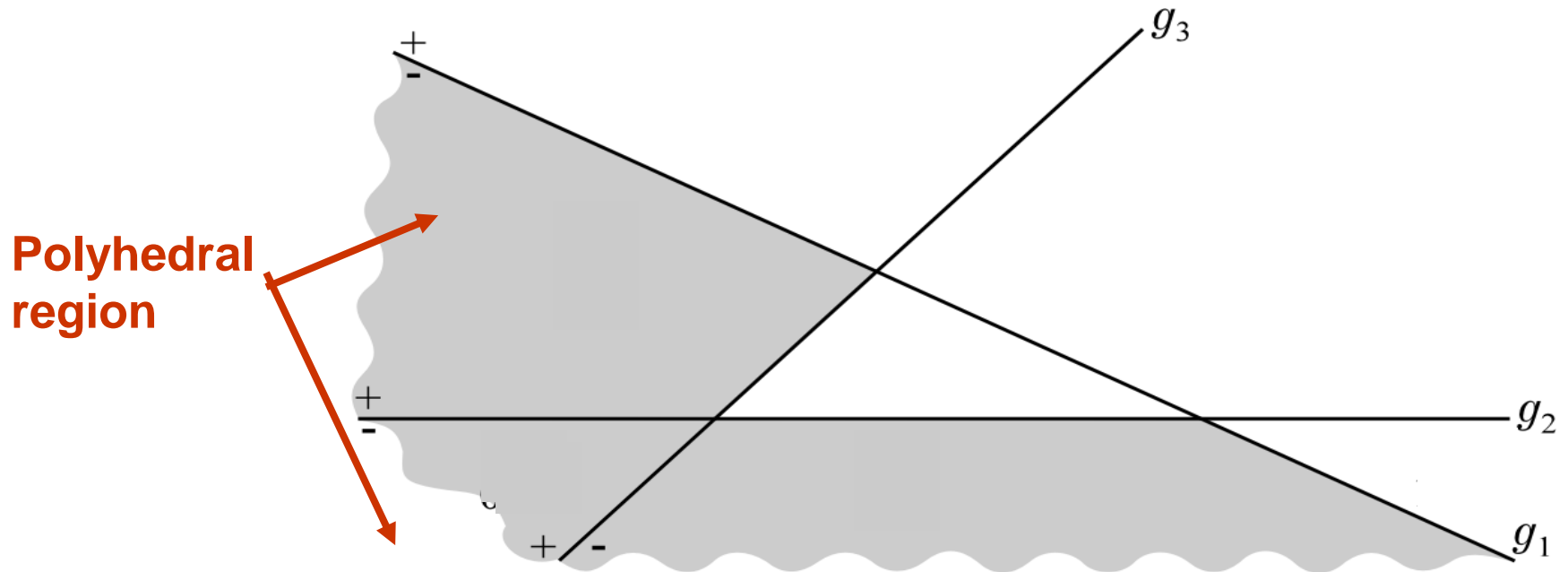
Classification Capabilities of Two Layer Perceptron



- maps a vector onto the vertices of the unit side hypercube, H_p
- mapping is through p neurons each realizing a hyper plane.
- The output of each of these neurons is 0 or 1

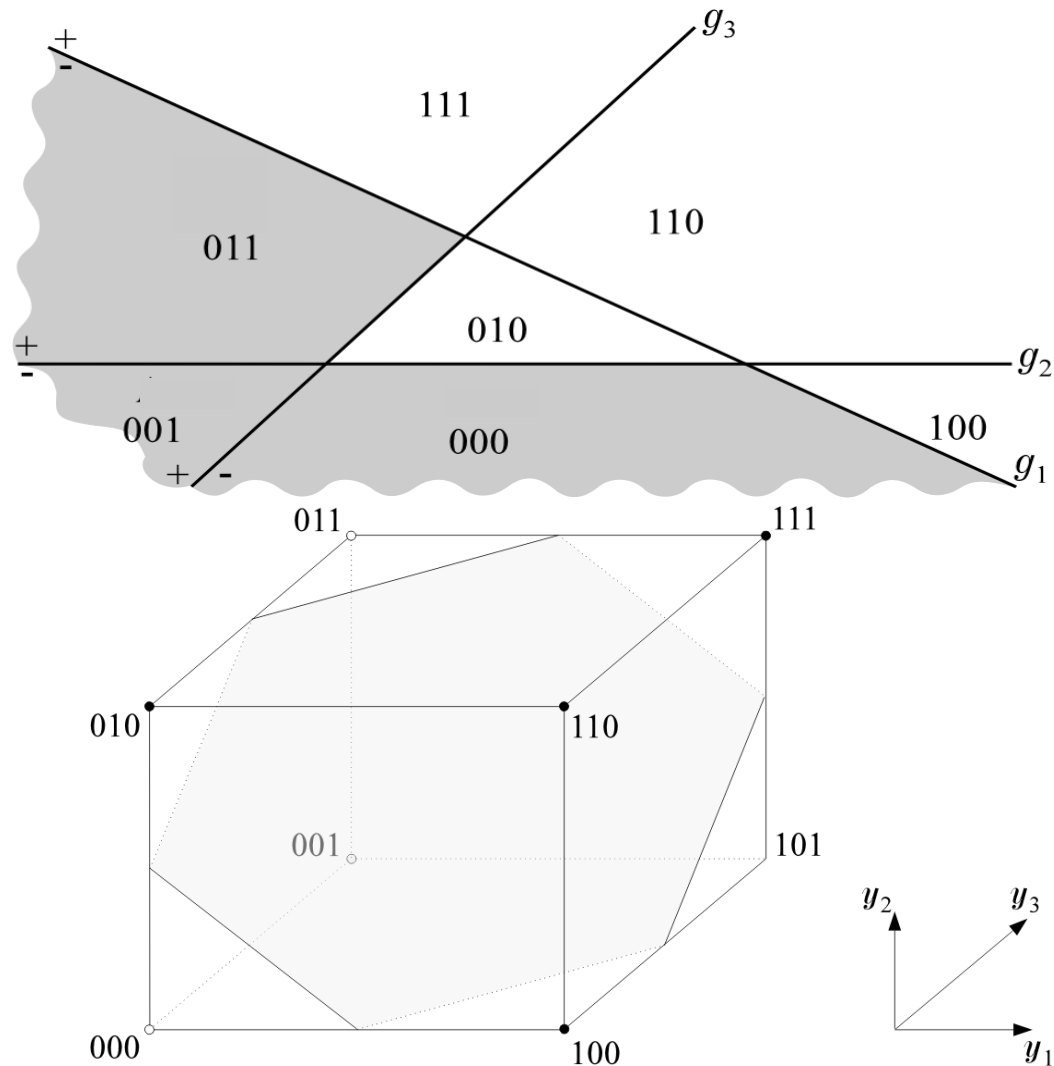
Classification Capabilities of Two Layer Perceptron

- Intersections of hyperplanes form regions.



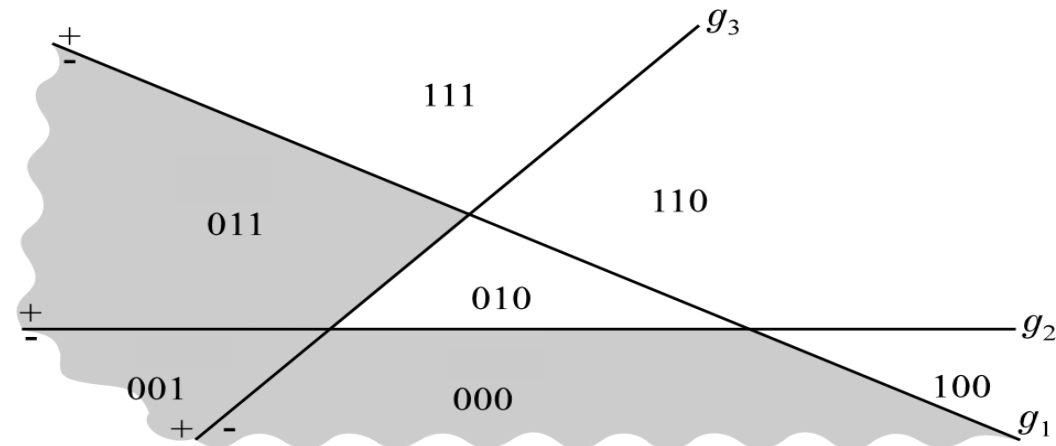
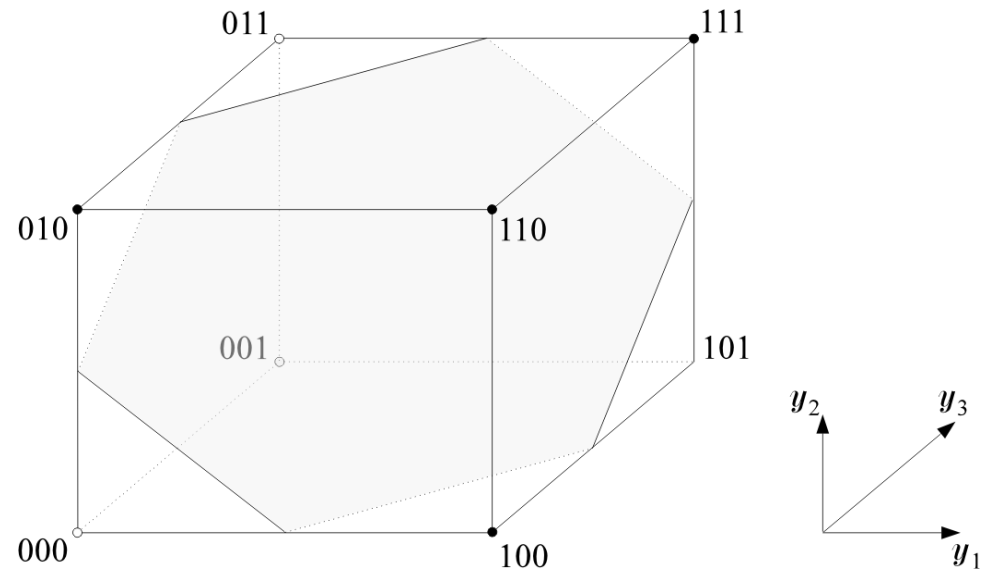
Classification Capabilities of Two Layer Perceptron

- Intersections of hyperplanes form regions.
- Each region corresponds to a vertex of the H_p unit hypercube.

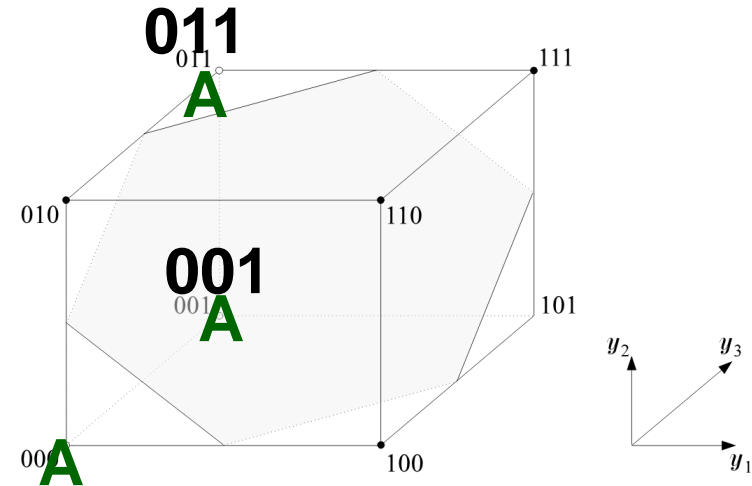


Classification Capabilities of Two Layer Perceptron

- For example, the 001 vertex corresponds to the region which is located
to the (-) side of $g_1(\underline{x})=0$
to the (-) side of $g_2(\underline{x})=0$
to the (+) side of $g_3(\underline{x})=0$

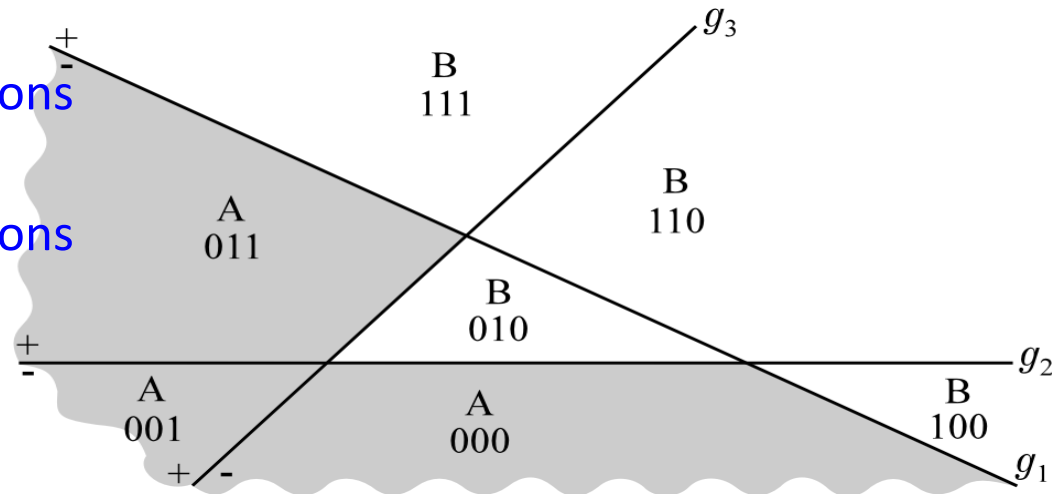


Classification Capabilities of Two Layer Perceptron

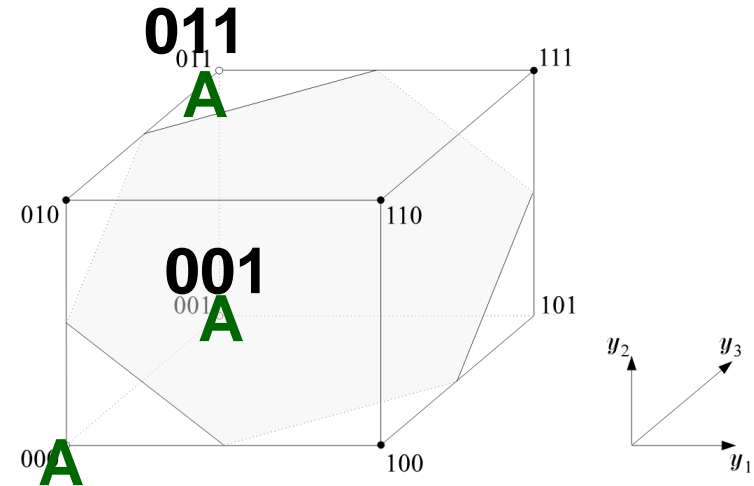
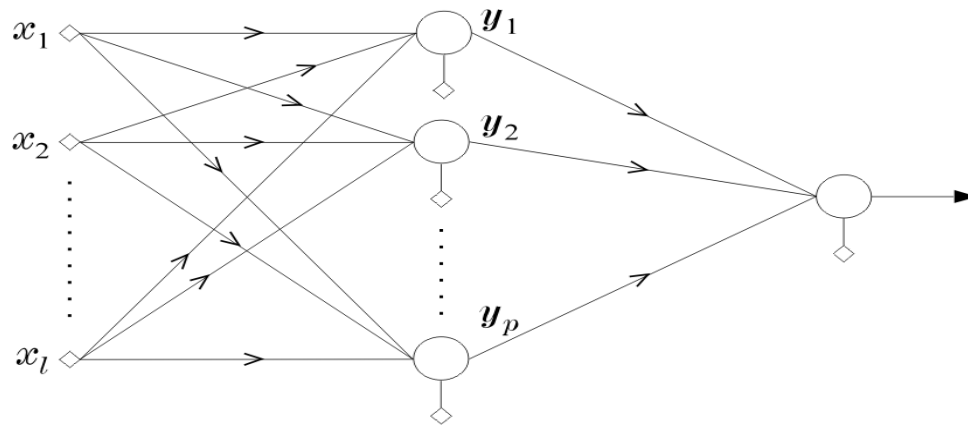


- A two-class problem

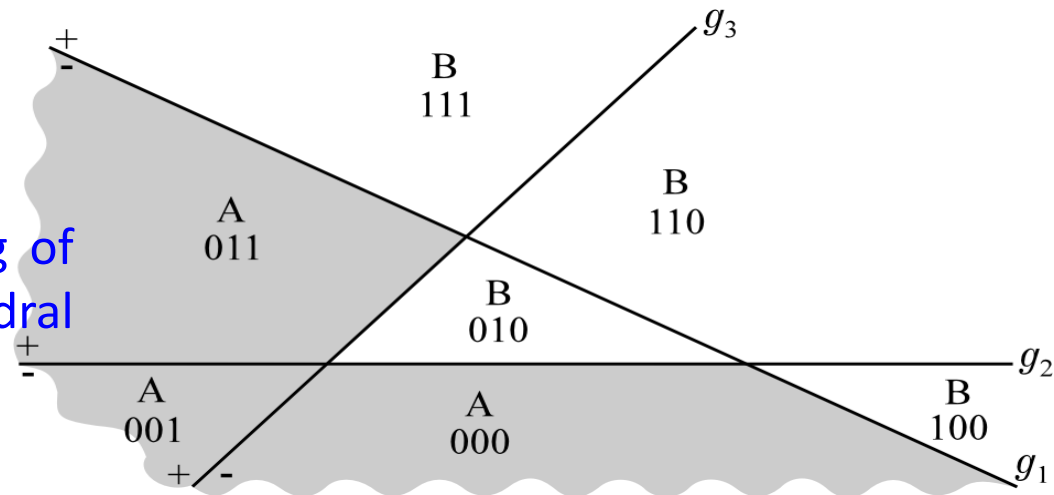
- **Class A patterns** from regions marked as A
- **Class B patterns** from regions marked as B



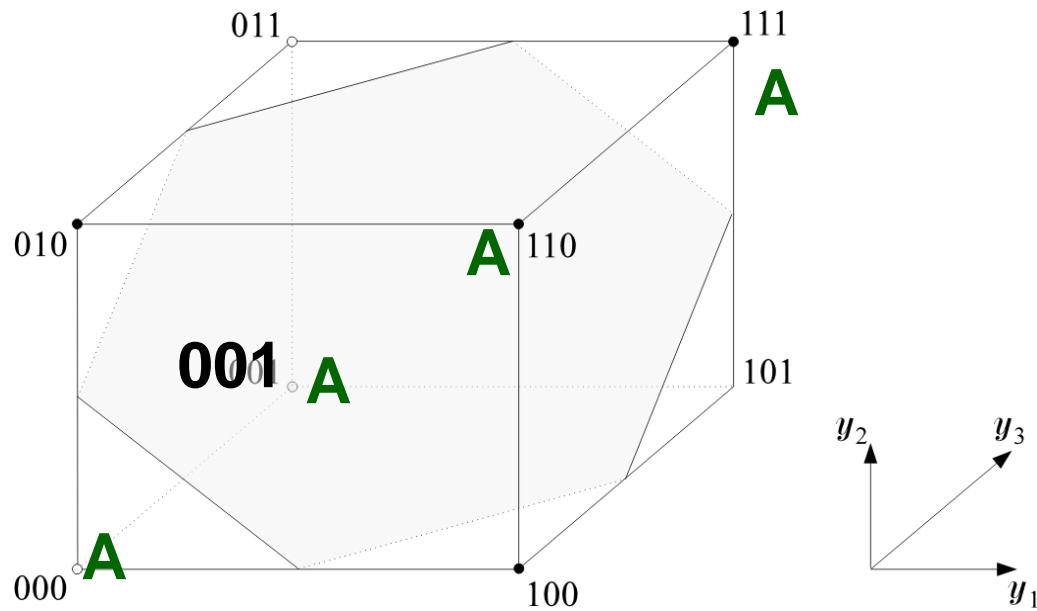
Classification Capabilities of Two Layer Perceptron



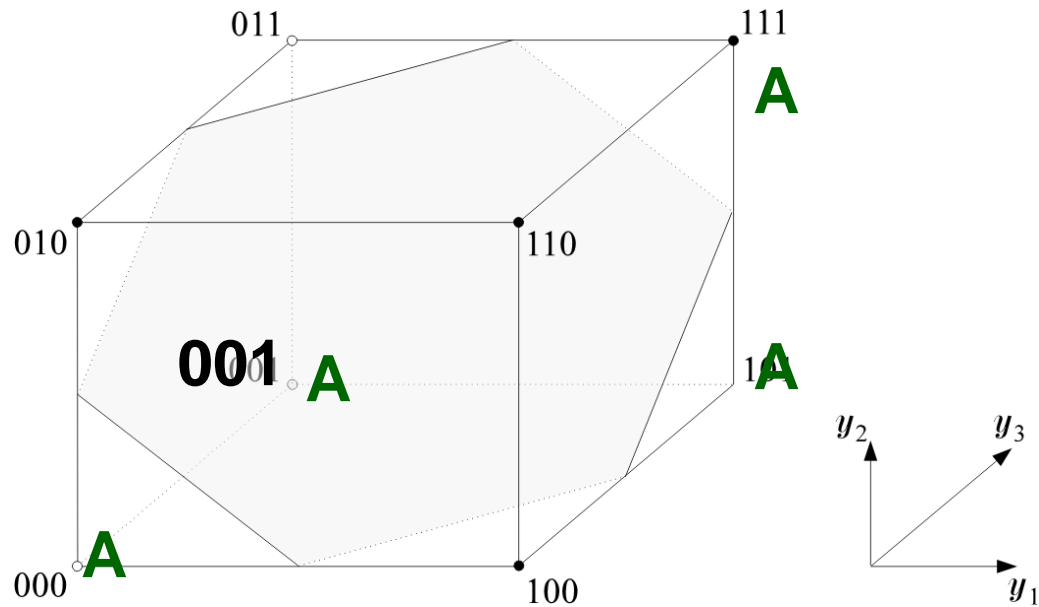
- The output neuron
 - realizes another hyperplane
 - separates the hypercube.
 - can classify vectors consisting of some unions of polyhedral regions.



Classification Capabilities of Two Layer Perceptron



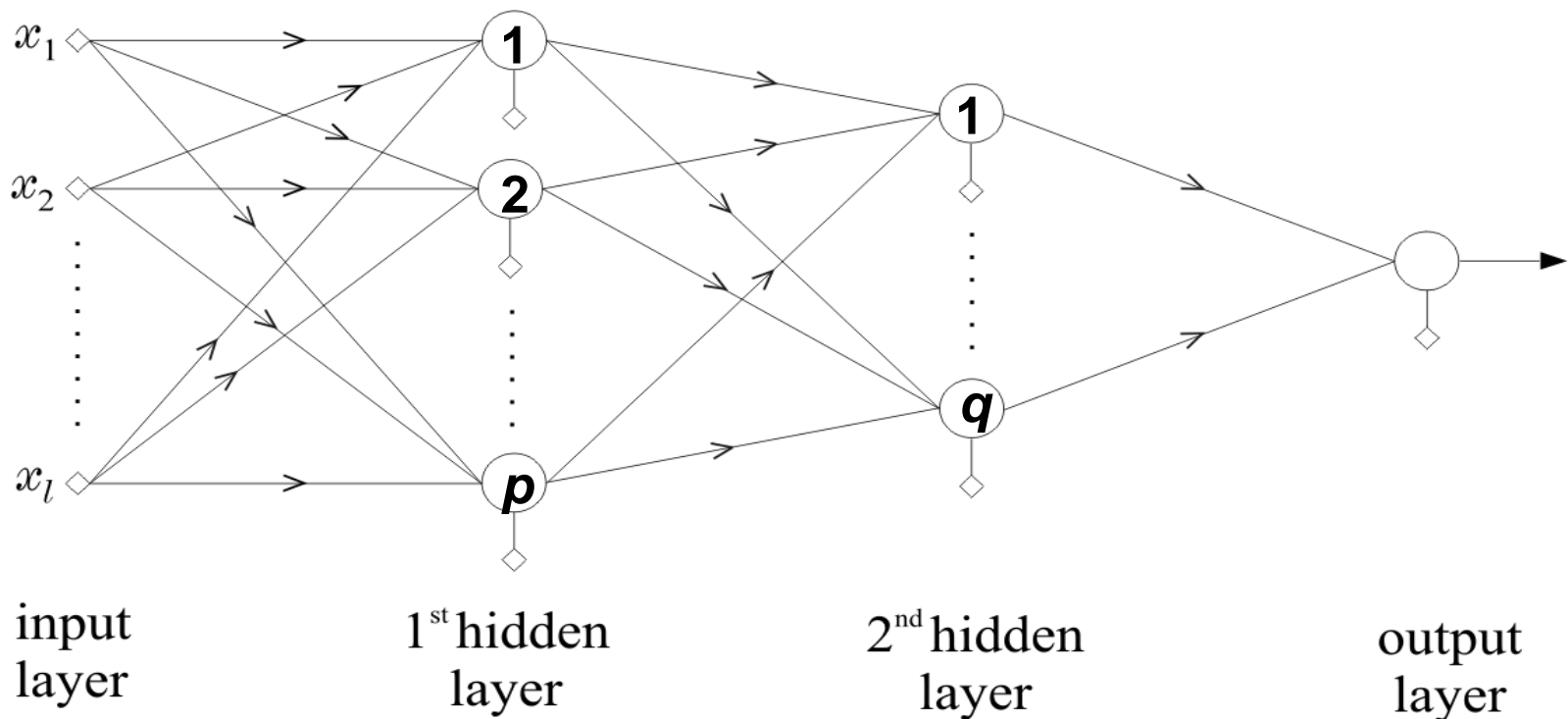
Classification Capabilities of Two Layer Perceptron



- The output neuron, *i.e.*, a 2 layer perceptron
 - cannot classify vectors consisting of arbitrary unions of polyhedral regions.

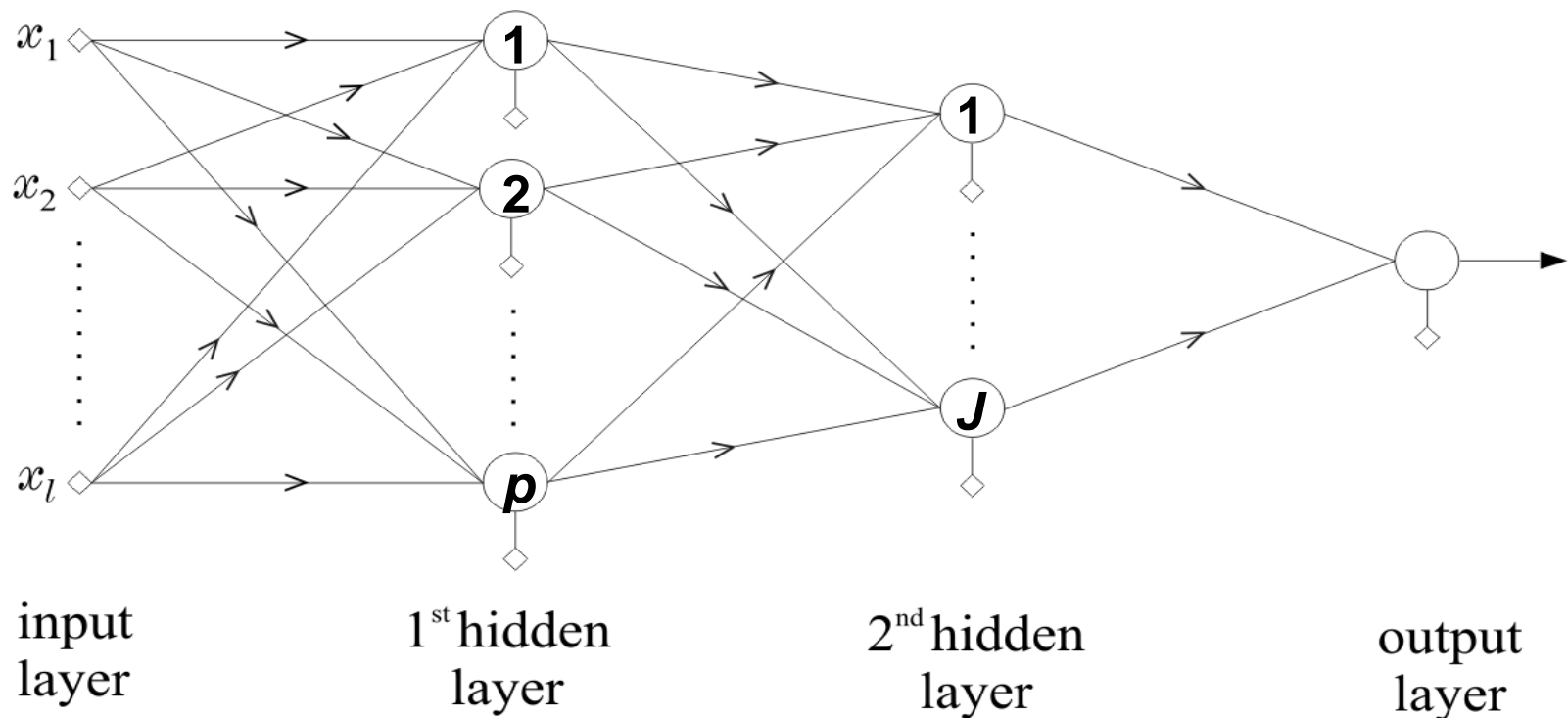
Solution: Three Layer Perceptron

- capable to classify vectors consisting of **ANY** union of polyhedral regions.
 - The idea is similar to the XOR problem.
 - Realizes more than one planes in the $\underline{y} \in R^p$ space.



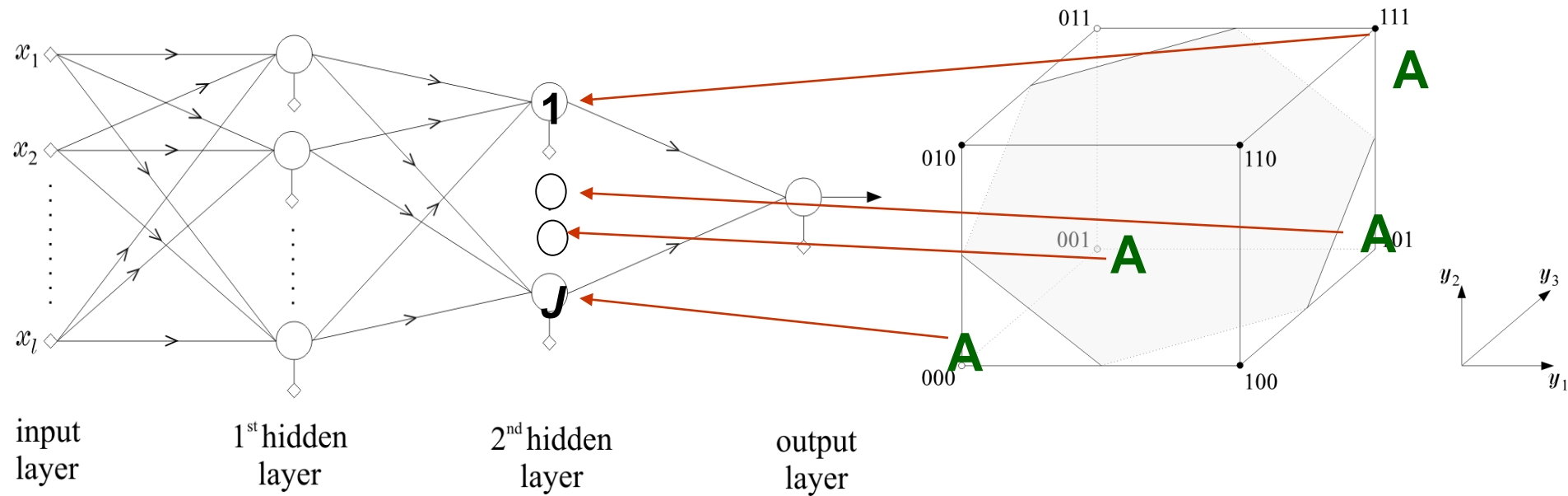
How does The Three Layer Perceptron Do It?

- Let, any J polyhedral regions constitutes vectors of class **A**.
- Learn a neuron in the 2nd hidden layer for each of J regions



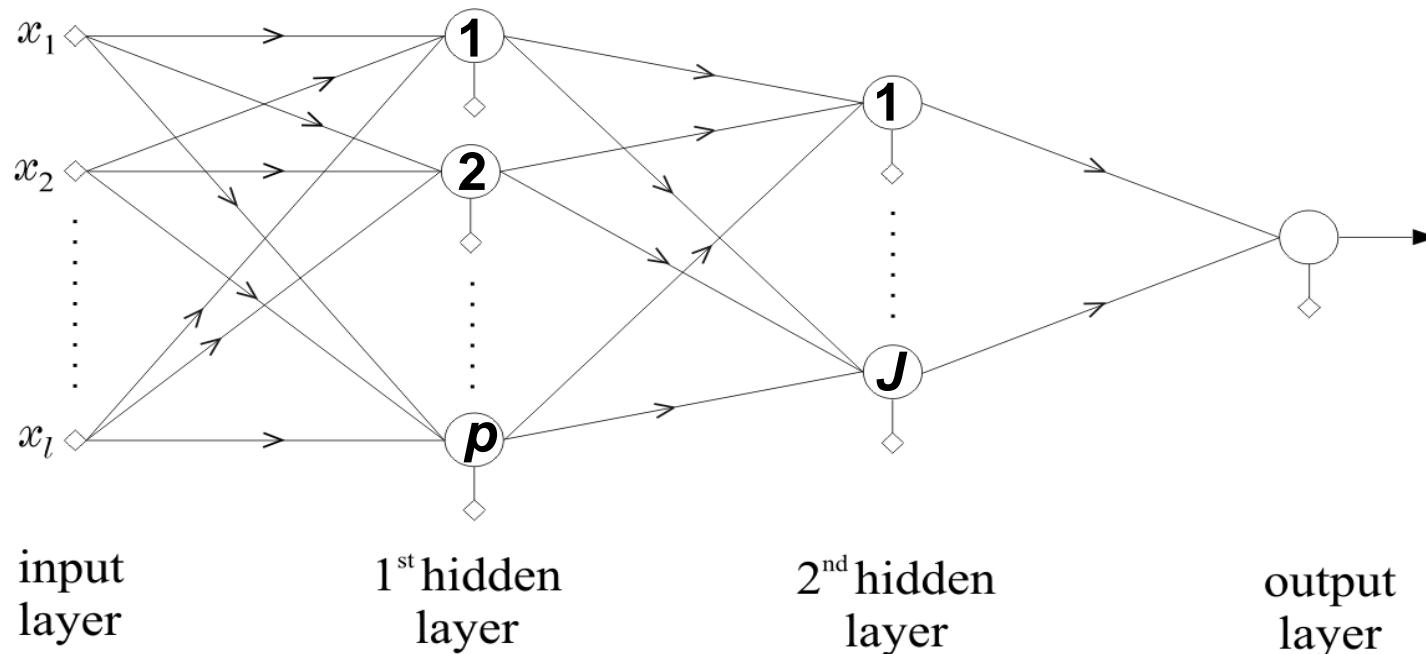
How does The Three Layer Perceptron Do It?

- Learn a neuron in the 2nd hidden layer for each of J regions



How does The Three Layer Perceptron Do It?

- For training vectors of a particular region of class **A**, only one of the 2nd-layer neuron produces 1, the rest of neurons produce 0.
- Now realize the output neuron as an OR gate.

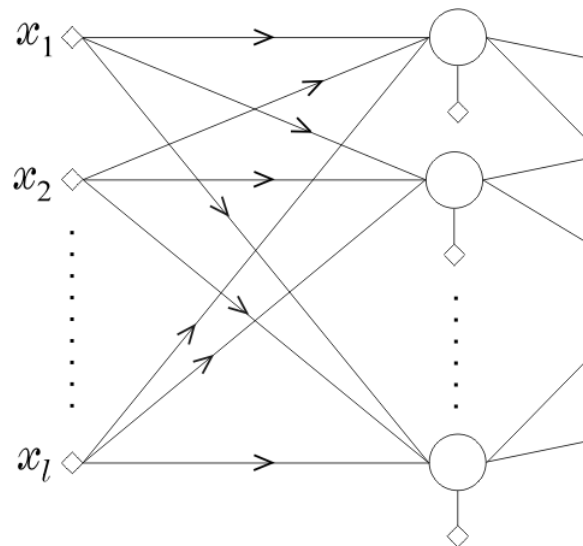


Training of a Multi Layer Perceptron (MLP)

- use rationale and develop a structure that **classifies correctly all the training patterns.**

OR

- choose a structure and compute the synaptic weights to **optimize a cost function.**

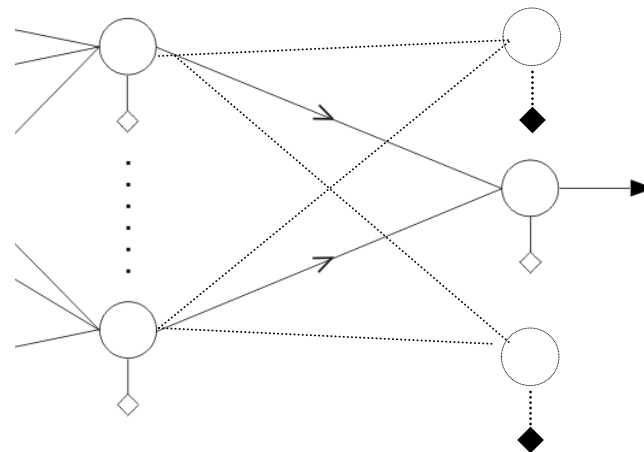


input
layer

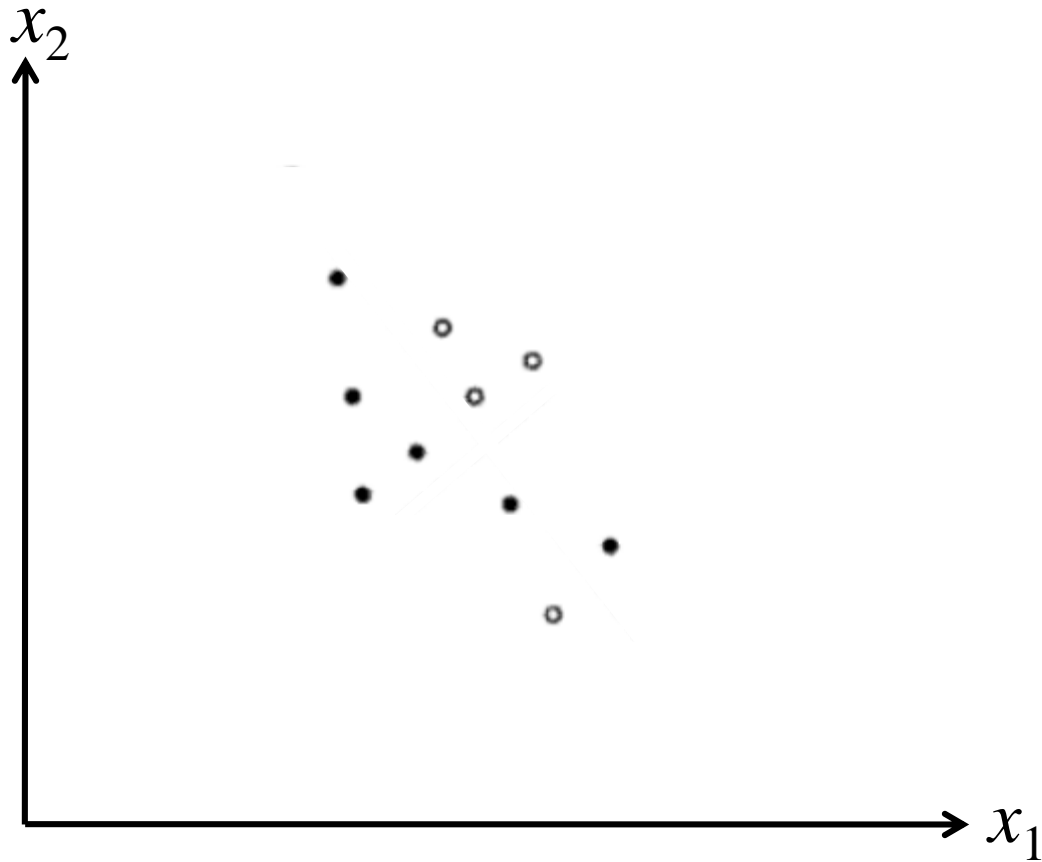
1st hidden
layer

$(L-1)^{\text{th}}$ hidden
layer

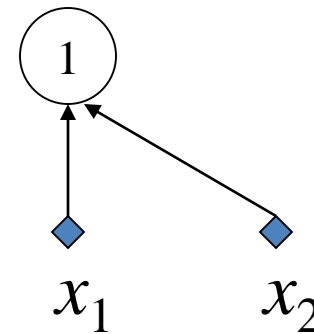
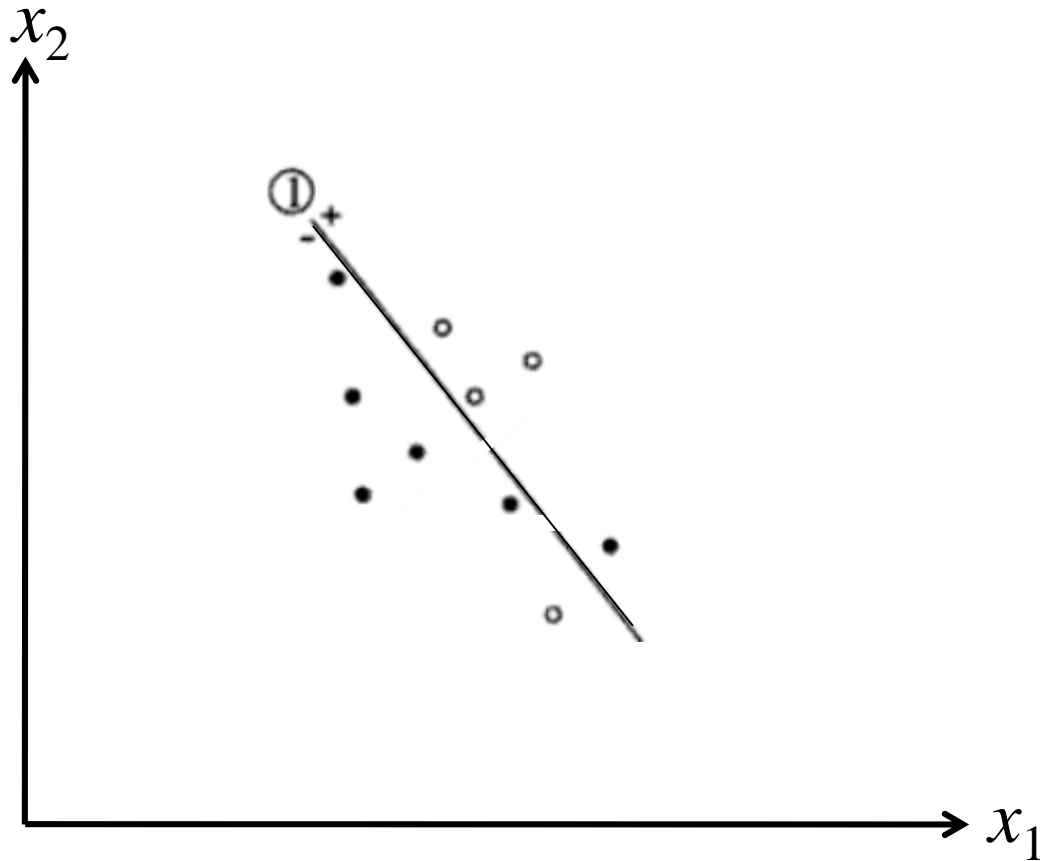
L^{th} or output
layer



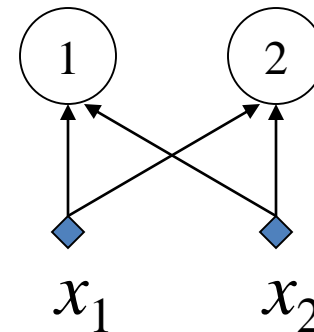
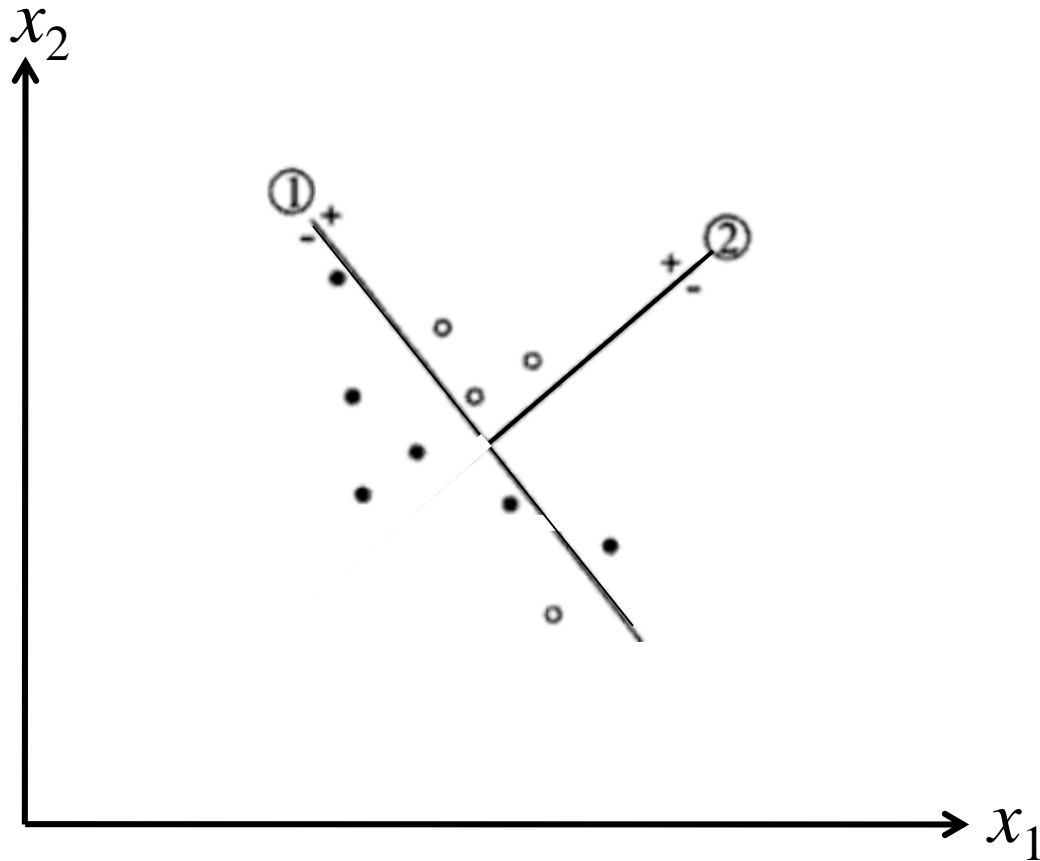
Algorithm Based on Exact Classification of Training Examples



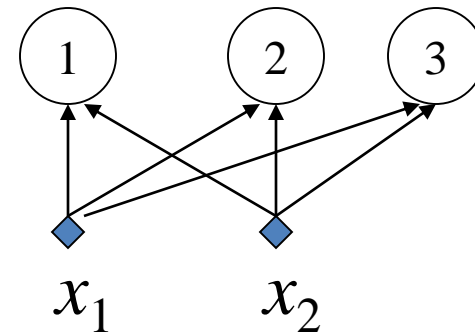
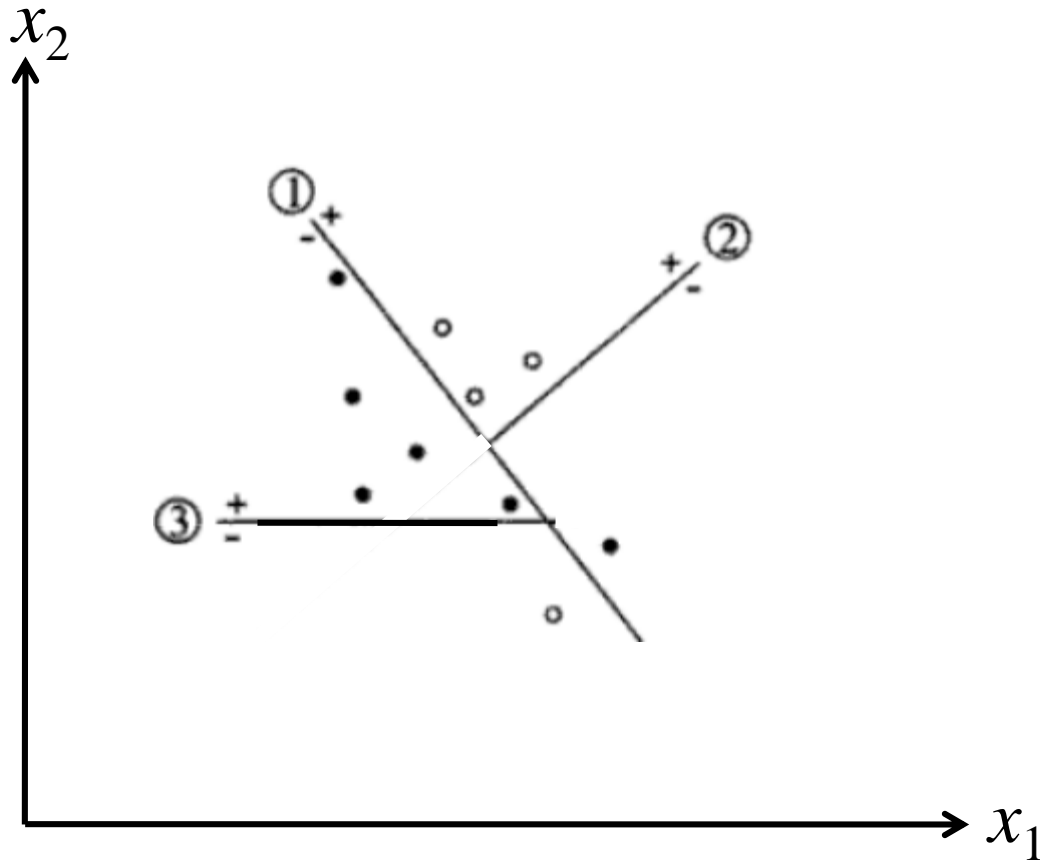
Algorithm Based on Exact Classification of Training Examples



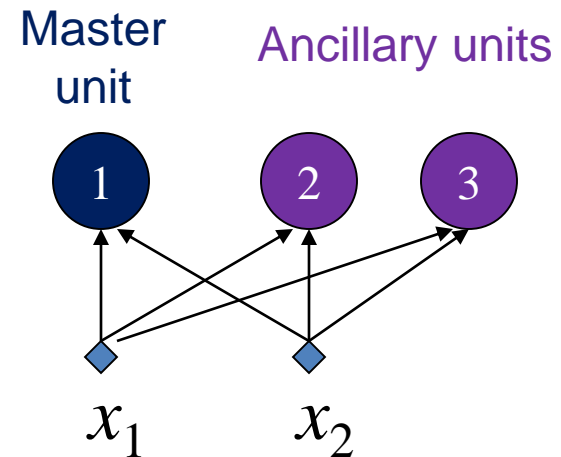
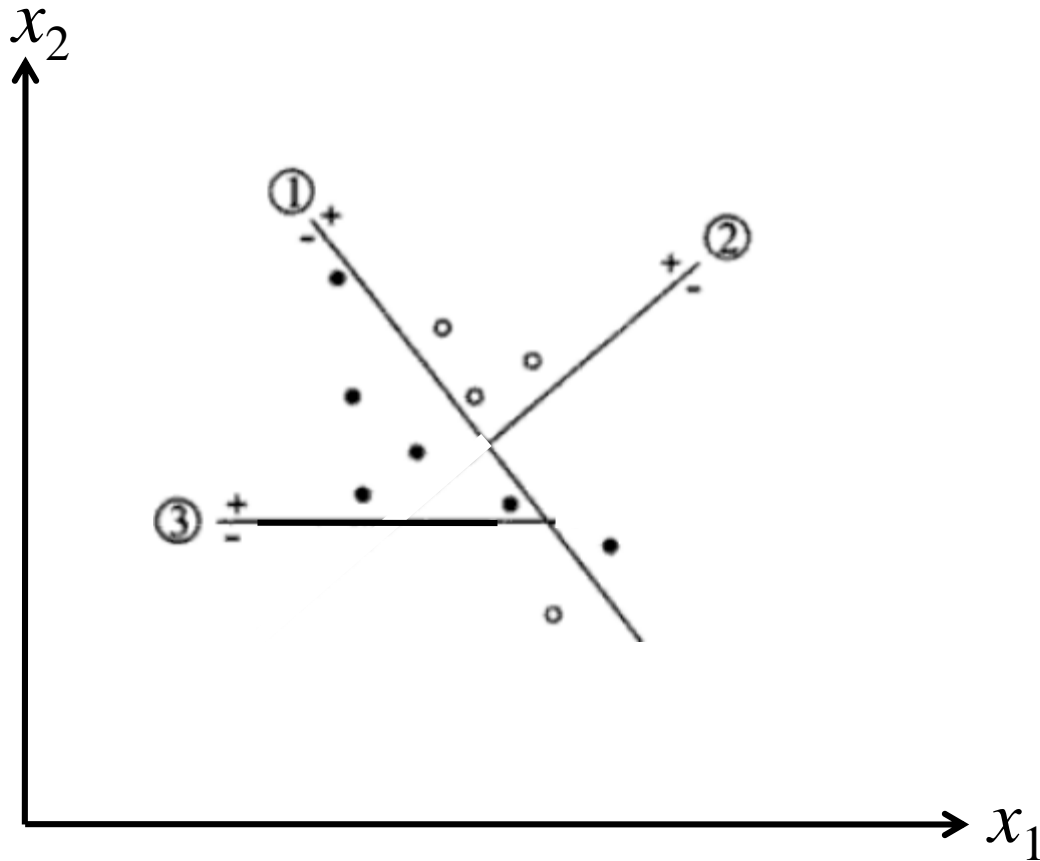
Algorithm Based on Exact Classification of Training Examples



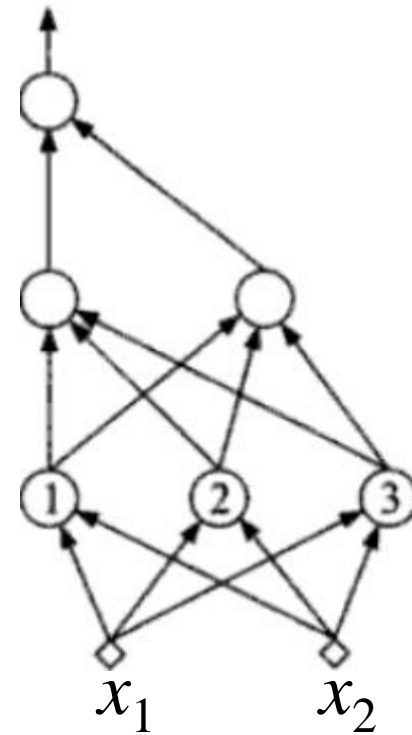
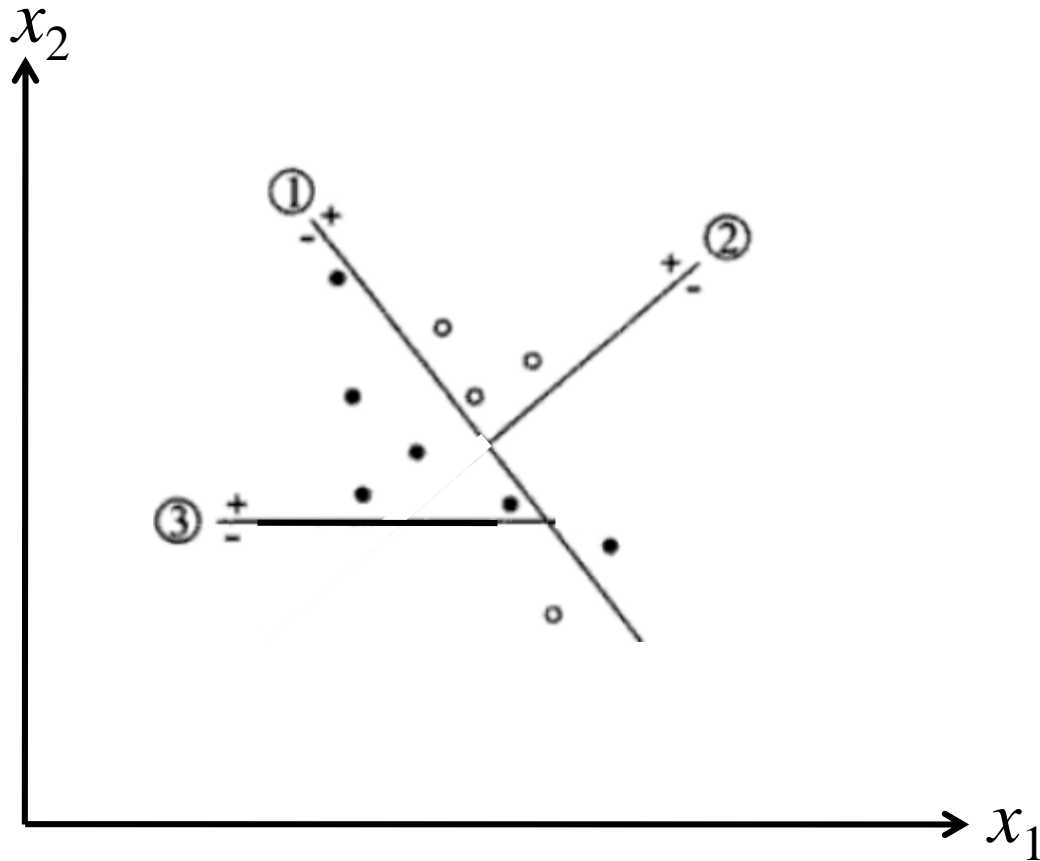
Algorithm Based on Exact Classification of Training Examples



Algorithm Based on Exact Classification of Training Examples



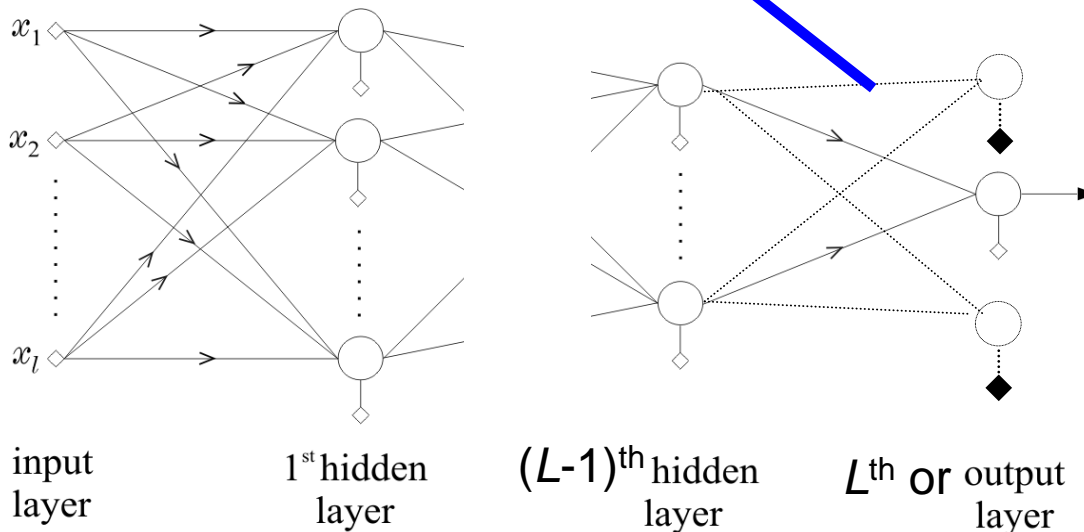
Algorithm Based on Exact Classification of Training Examples



Iterative update of Synaptic weights: The Backpropagation Algorithm

- computes the weights iteratively, subject to a cost function is optimized

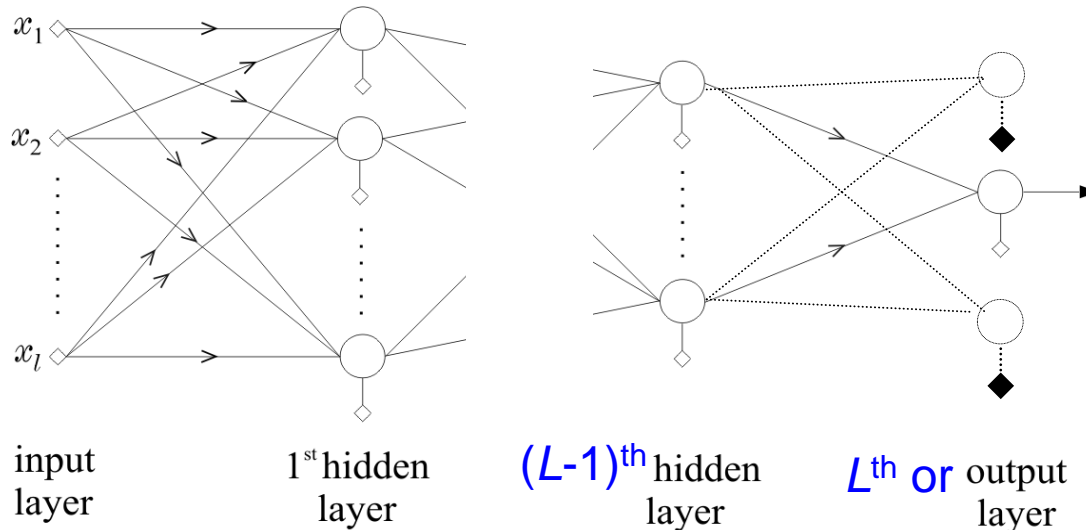
$$w'_j(\text{new}) = w'_j(\text{old}) + \Delta w'_j$$



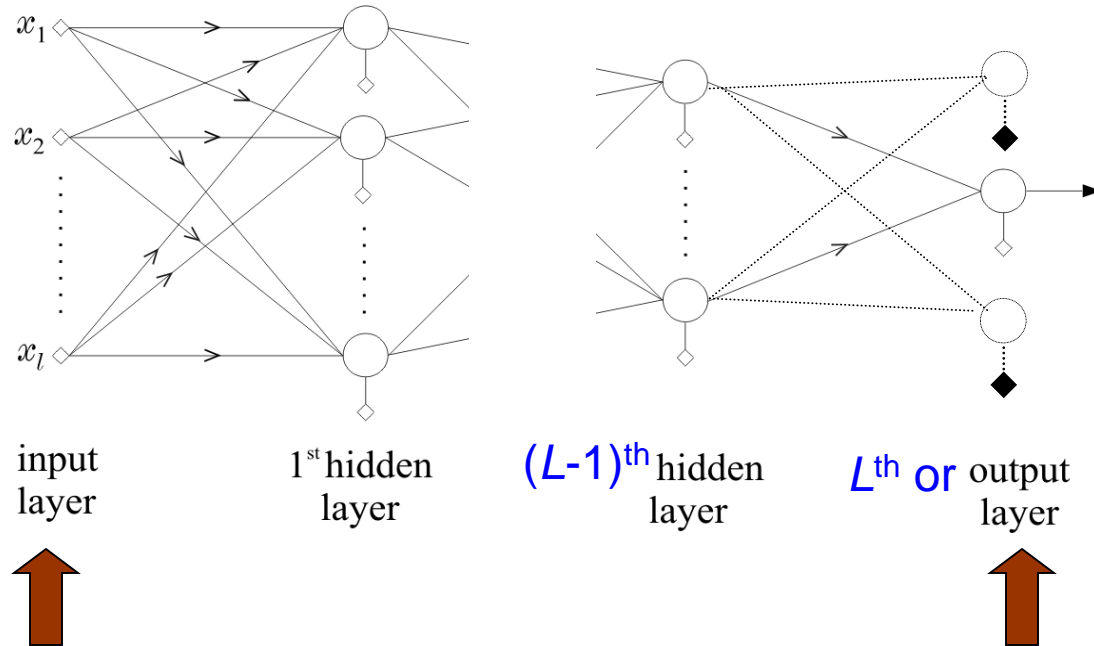
Iterative update of Synaptic weights: The Backpropagation Algorithm

Assume:

- Multiple layers
- more than one neurons in each layer
- any number of classes



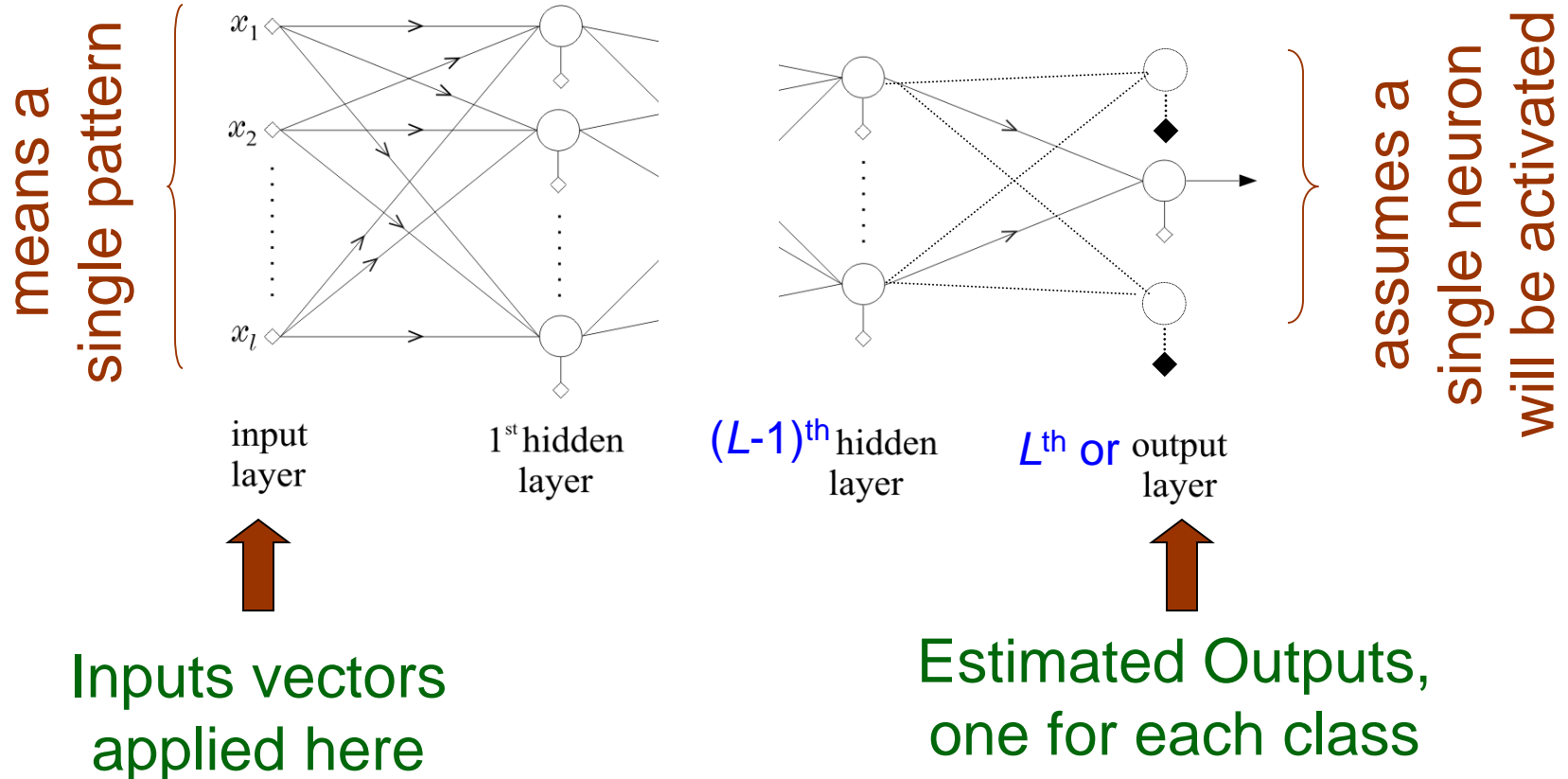
Iterative update of Synaptic weights: The Backpropagation Algorithm



Inputs vectors
applied here

Estimated Outputs,
one for each class

Iterative update of Synaptic weights: The Backpropagation Algorithm



Iterative update of Synaptic weights: The Backpropagation Algorithm

Let:

- 4 classes

Training Sample#	Class#
Sample#1	1
Sample#2	3
Sample#3	2
Sample#4	4
Sample#5	2

Iterative update of Synaptic weights: The Backpropagation Algorithm

Let:

- 4 classes

Training Sample#	Class#	Class Vector
Sample#1	1	1 0 0 0
Sample#2	3	0 0 1 0
Sample#3	2	0 1 0 0
Sample#4	4	0 0 0 1
Sample#5	2	0 1 0 0

The Backpropagation Algorithm

- Recall the perceptron algorithm:

- We update with this

$$\underline{w}(\text{new}) = \underline{w}(\text{old}) + \Delta \underline{w}$$

The Backpropagation Algorithm

- Recall the perceptron algorithm:
 - We update with this $\underline{w}(\text{new}) = \underline{w}(\text{old}) + \Delta \underline{w}$
- Backpropagation updates multiple nodes for a number of layers:

$$\mathbf{w}'_j(\text{new}) = \mathbf{w}'_j(\text{old}) + \Delta \mathbf{w}'_j$$

The Backpropagation Algorithm

- Recall the perceptron algorithm:
 - We update with this $\underline{w}(\text{new}) = \underline{w}(\text{old}) + \Delta \underline{w}$
- Backpropagation updates multiple nodes for a number of layers:

$$\mathbf{w}'_j(\text{new}) = \mathbf{w}'_j(\text{old}) + \Delta \mathbf{w}'_j$$

r-th layer

j-th neuron

The Backpropagation Algorithm

- Another difference is the activation function:
- Perceptron algorithm uses unit activation function:

$$f(x) = \begin{cases} 1 & x > 0 \\ 0 & x < 0 \end{cases}$$

- This function is not differentiable at $x=0$.

The Backpropagation Algorithm

- Another difference is the activation function:
- Perceptron algorithm uses unit activation function:

$$f(x) = \begin{cases} 1 & x > 0 \\ 0 & x < 0 \end{cases}$$

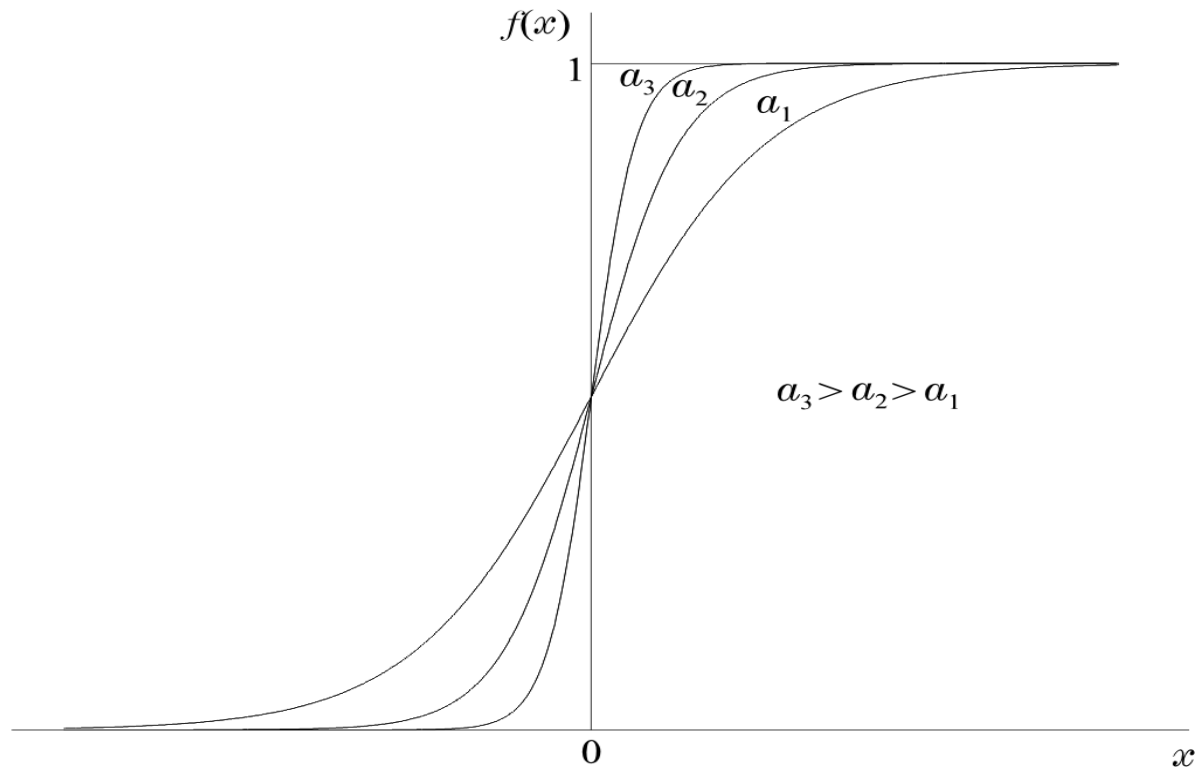
– This function is not differentiable at $x=0$.

- Backpropagation uses logistic function:

$$f(x) = \frac{1}{1 + \exp(-ax)}$$

Logistic function

The Logistic function



$$f(x) = \frac{1}{1 + \exp(-ax)}$$

The Backpropagation Algorithm

- Similar to perceptron algorithm: *Backpropagation also iteratively updates weights*

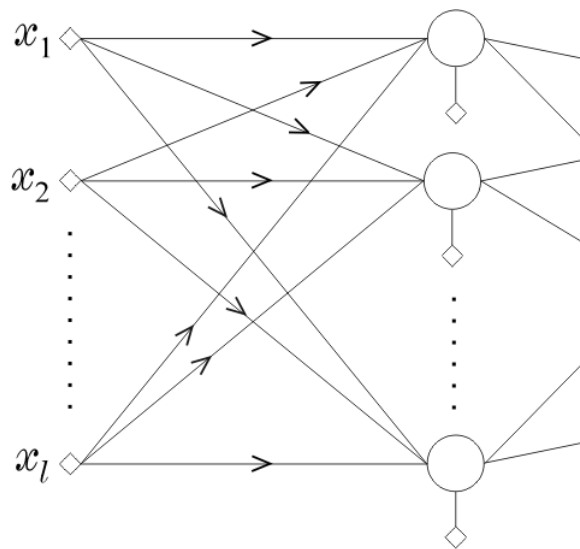
$$\mathbf{w}_j^r(\text{new}) = \mathbf{w}_j^r(\text{old}) + \Delta \mathbf{w}_j^r$$

where, $\Delta \mathbf{w}_j^r = -\mu \frac{\partial J}{\partial \mathbf{w}_j^r}$

and $J = \sum_{i=1}^N \mathcal{E}(i)$

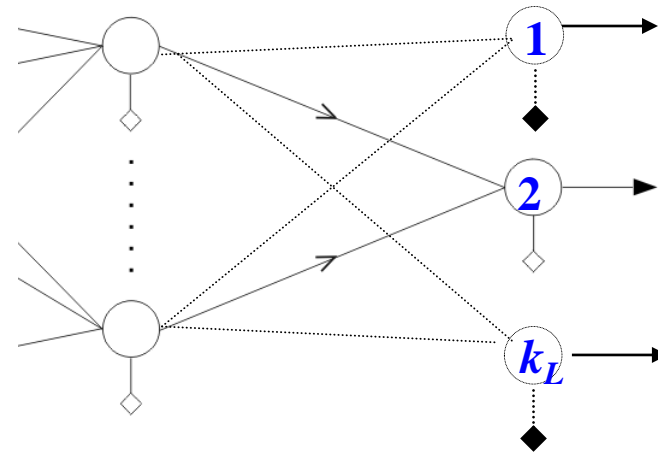
Define the Terms

- L layers of neurons
- k_r neurons in r^{th} layer
- k_0 nodes in the input layer = input feature dimension = l
- k_L nodes in the output layer = output class dimension



input
layer

1st hidden
layer

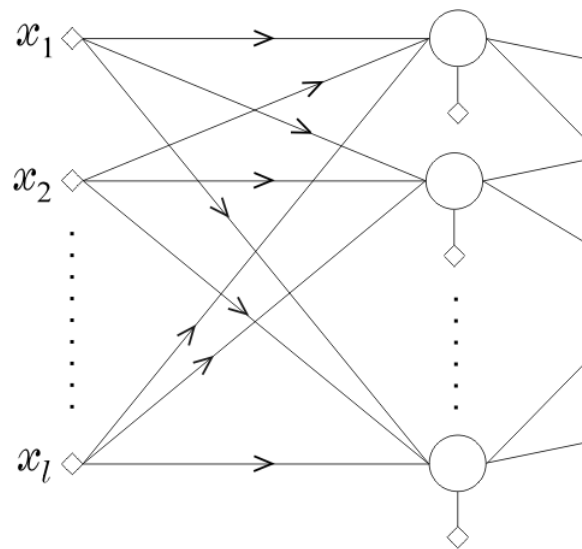


$(L-1)^{\text{th}}$ hidden
layer

L^{th} or
output
layer

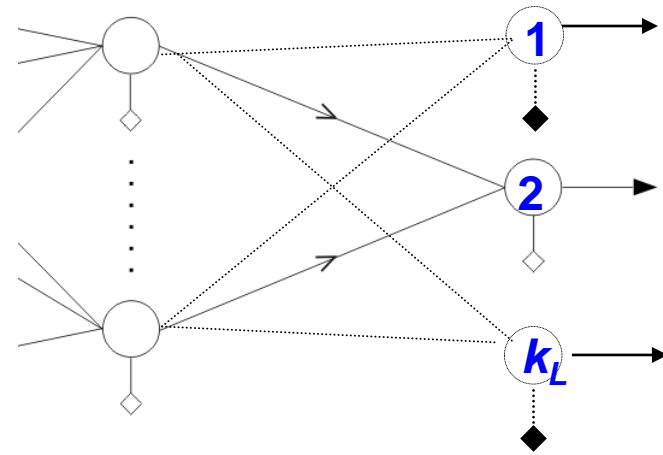
Define the Terms

- *Remember:* The number of classes is more than 2, it is K_L .
- Class value of a sample is no longer a single variable, rather it is a vector of k_L dimension.



input
layer

1st hidden
layer

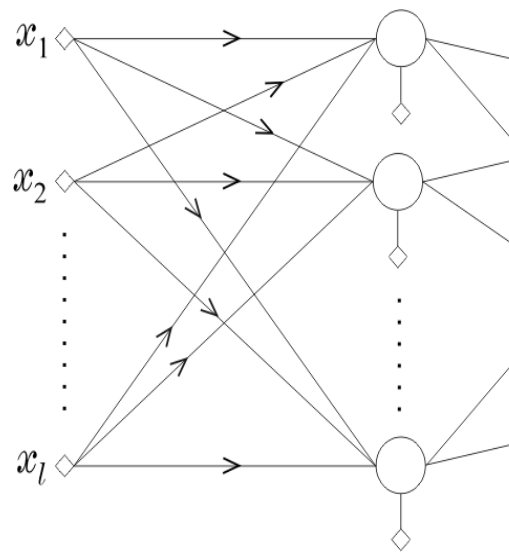


$(L-1)^{\text{th}}$ hidden
layer

L^{th} or
output
layer

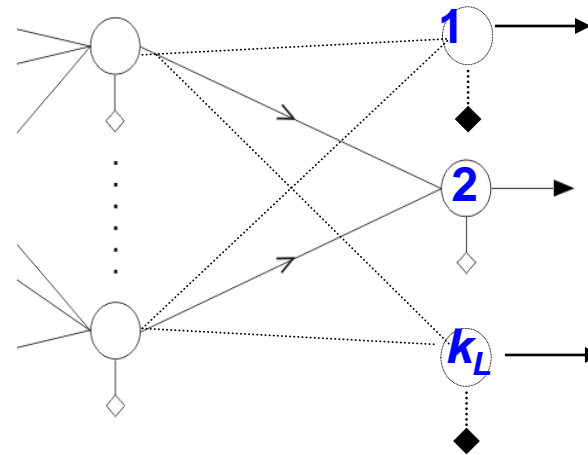
Define the Terms

- N training samples: $(\mathbf{x}(i), \mathbf{y}(i))$, for $i = 1, 2, 3, \dots, N$
- Features of i th training sample: $\mathbf{x}(i) = [x_1(i), \dots, x_{k_0}(i)]^T$
- Class of i th training sample: $\mathbf{y}(i) = [y_1(i), \dots, y_{k_L}(i)]^T$



input
layer

1st hidden
layer

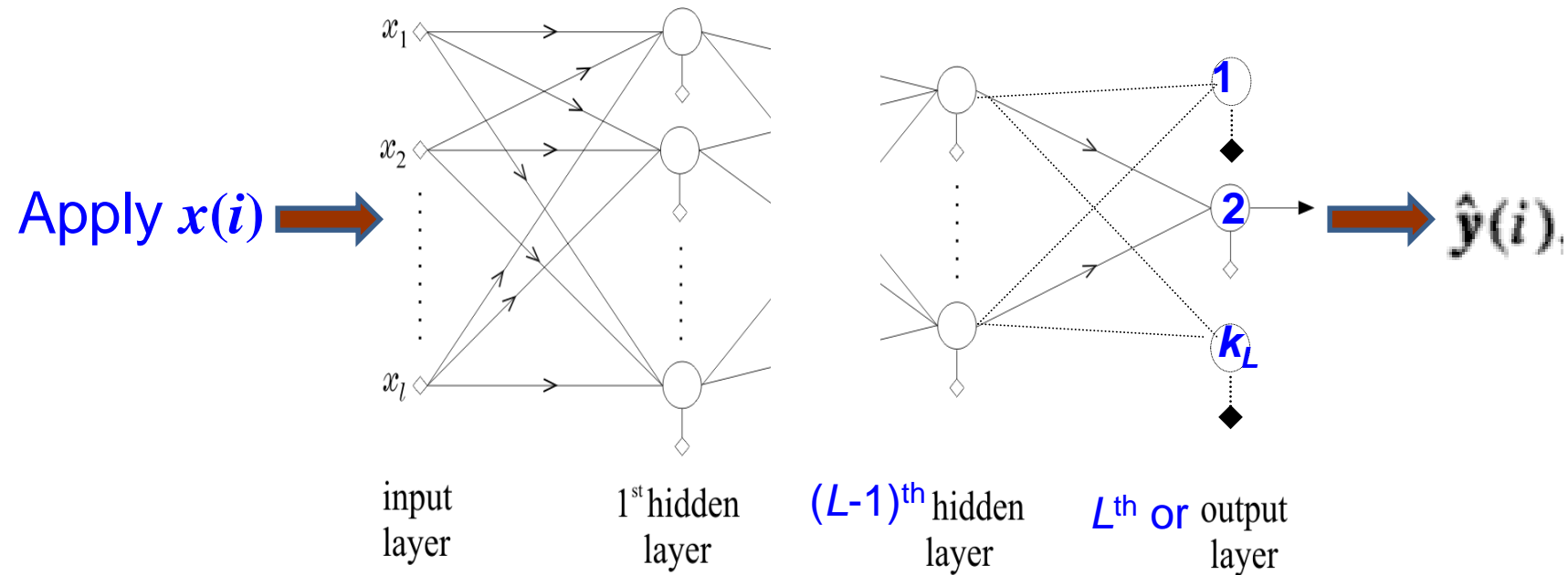


$(L-1)$ th hidden
layer

L th or
output
layer

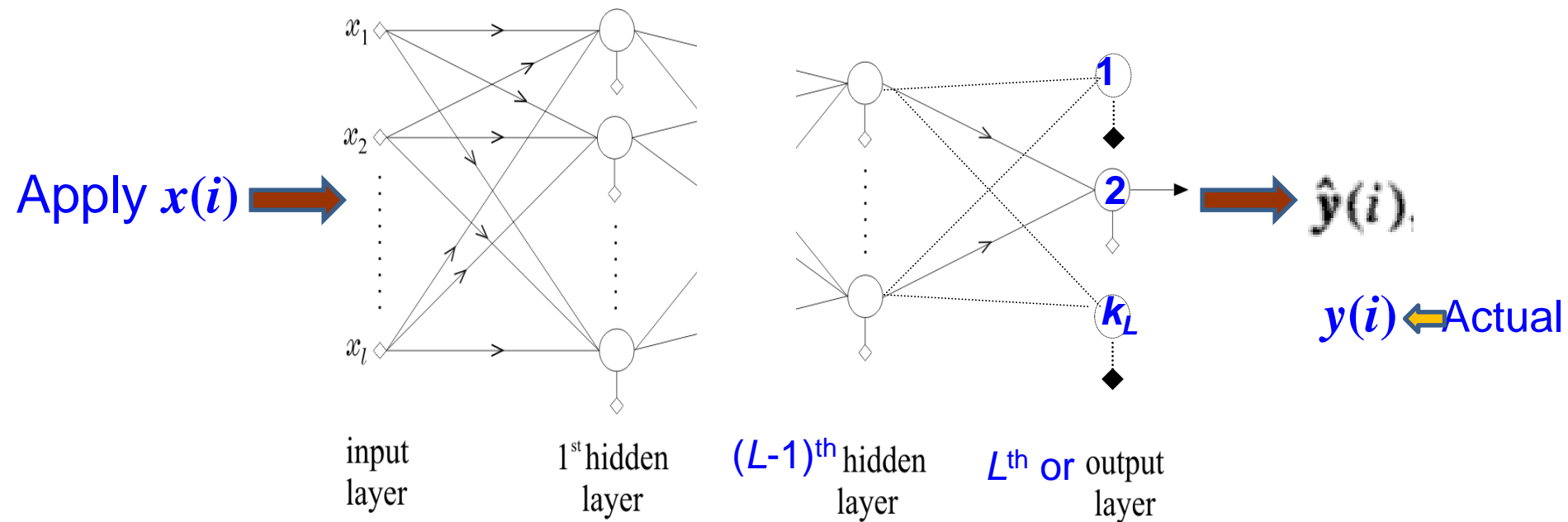
Define the Terms

- During training, apply i^{th} training vector $x(i)$, and output is $\hat{y}(i)$, instead of $y(i)$



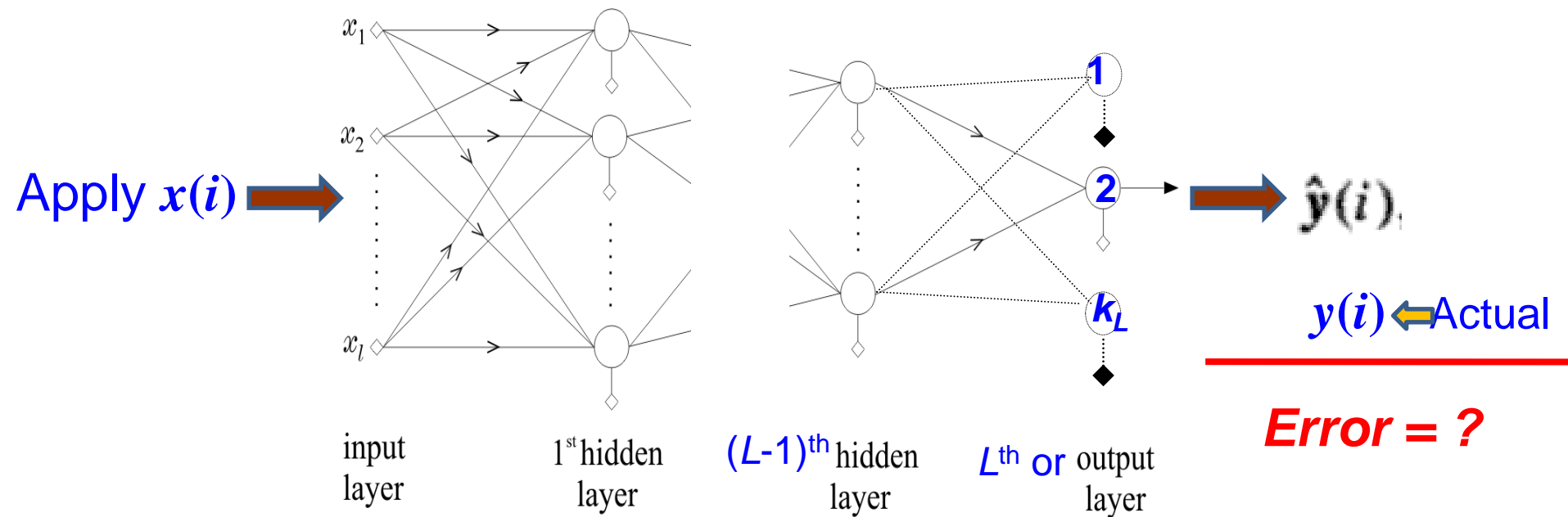
Define the Terms

- During training, apply i^{th} training vector $\mathbf{x}(i)$, and output is $\hat{\mathbf{y}}(i)$, instead of $\mathbf{y}(i)$



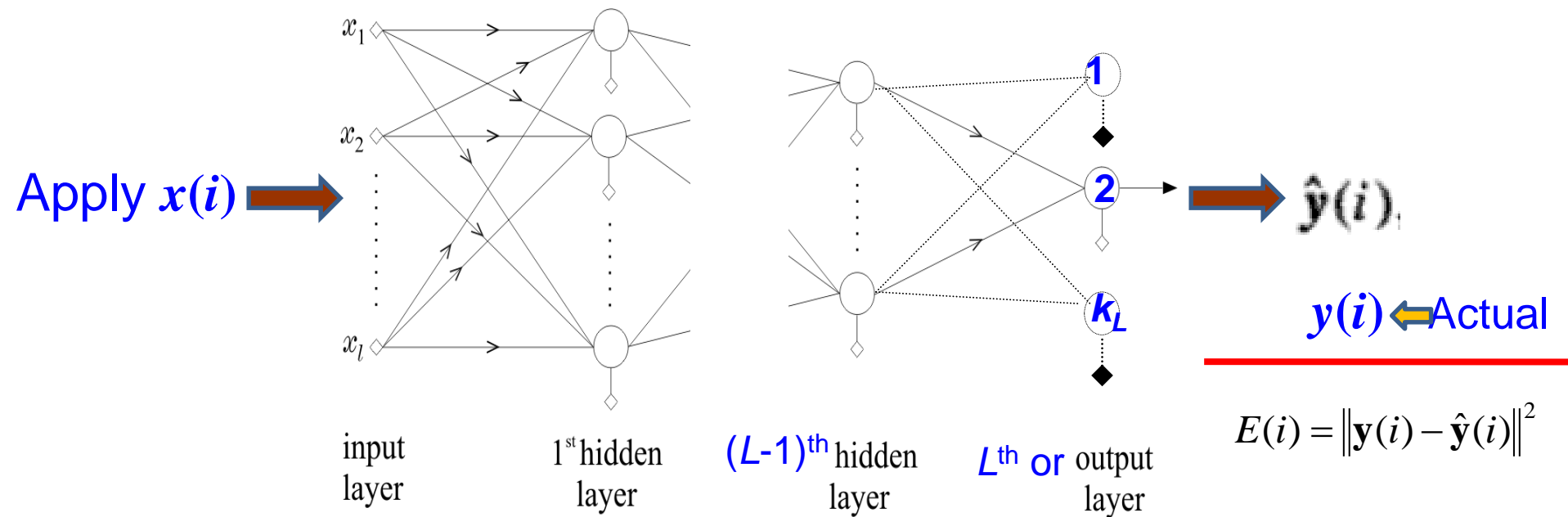
Define the Terms

- During training, apply i^{th} training vector $\mathbf{x}(i)$, and output is $\hat{\mathbf{y}}(i)$, instead of $\mathbf{y}(i)$



Define the Terms

- During training, apply i^{th} training vector $\mathbf{x}(i)$, and output is $\hat{\mathbf{y}}(i)$, instead of $\mathbf{y}(i)$



Define the Terms

- During training, apply i^{th} training vector $\mathbf{x}(i)$, and output is $\hat{\mathbf{y}}(i)$, instead of $\mathbf{y}(i)$
- Error for i^{th} vector:

$$\mathcal{E}(i) = \frac{1}{2} \sum_{m=1}^{k_L} e_m^2(i) \equiv \frac{1}{2} \sum_{m=1}^{k_L} (y_m(i) - \hat{y}_m(i))^2, \quad i = 1, 2, \dots, N$$

- Total Error:

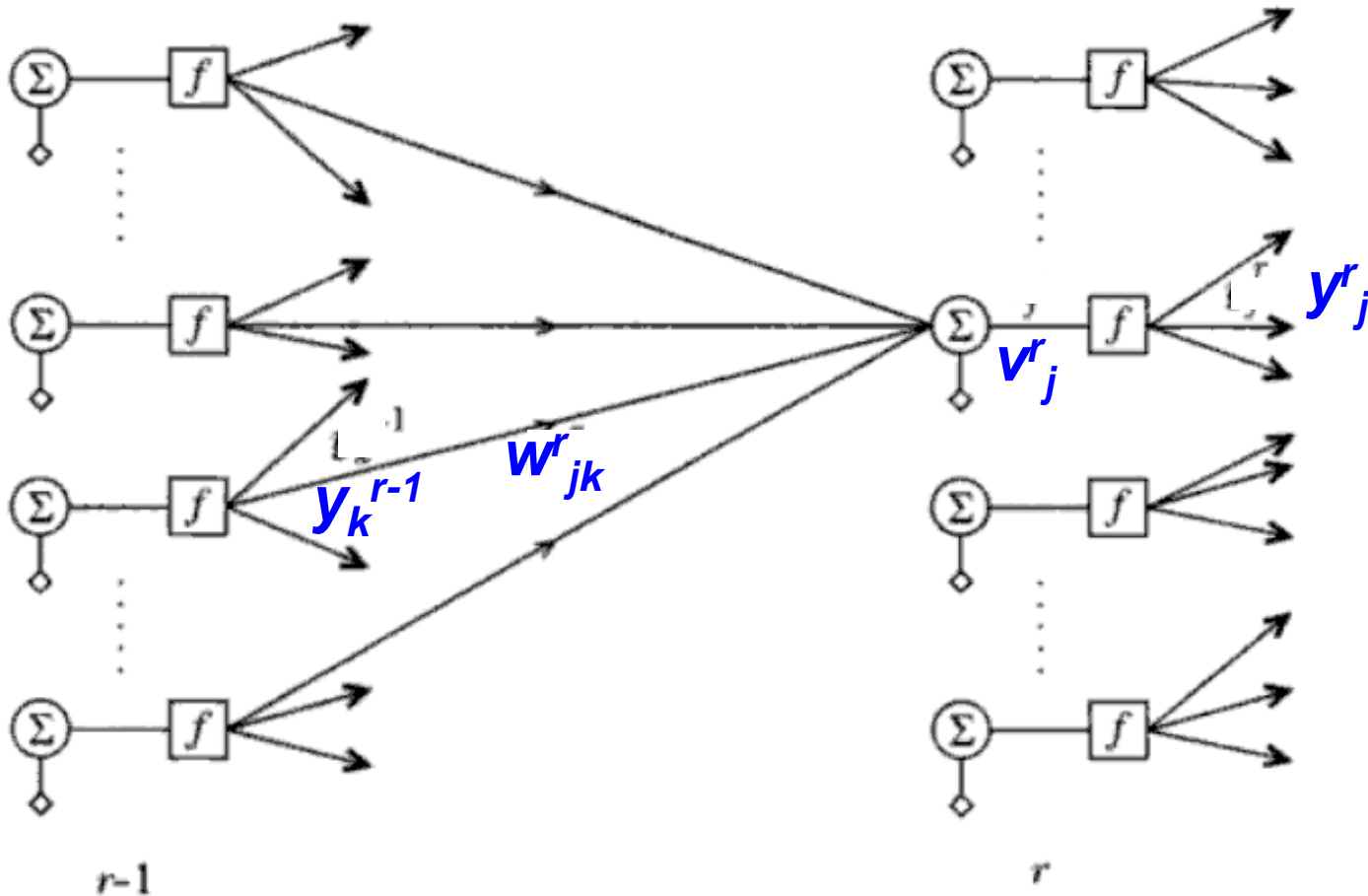
$$J = \sum_{i=1}^N \mathcal{E}(i)$$

Define the Terms

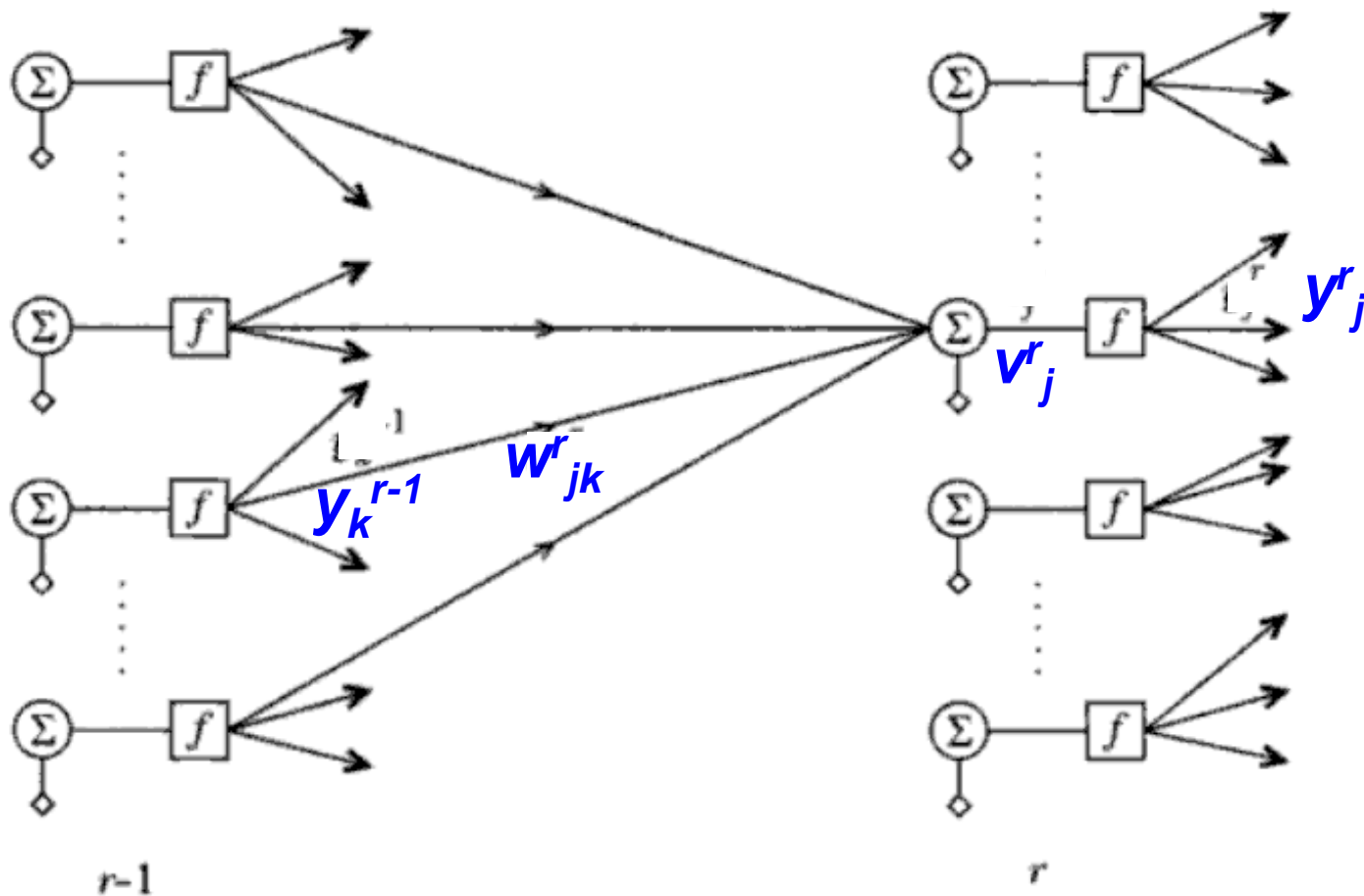
- We need to calculate $\Delta \mathbf{w}_j^r = -\mu \frac{\partial J}{\partial \mathbf{w}_j^r} = -\mu \sum_{i=1}^N \frac{\partial E(i)}{\partial \mathbf{w}_j^r(i)}$

Define the Terms

- We need to calculate $\Delta \mathbf{w}_j^r = -\mu \frac{\partial J}{\partial \mathbf{w}_j^r} = -\mu \sum_{i=1}^N \frac{\partial E(i)}{\partial \mathbf{w}_j^r(i)}$
- J depends on \mathbf{w}_j^r and passes through \mathbf{v}_j^r



Define the Terms



$$v_j^r(i) = \sum_{k=1}^{k_{r-1}} w_{jk}^r y_k^{r-1}(i) + w_{jo}^r \equiv \sum_{k=0}^{k_{r-1}} w_{jk}^r y_k^{r-1}(i)$$

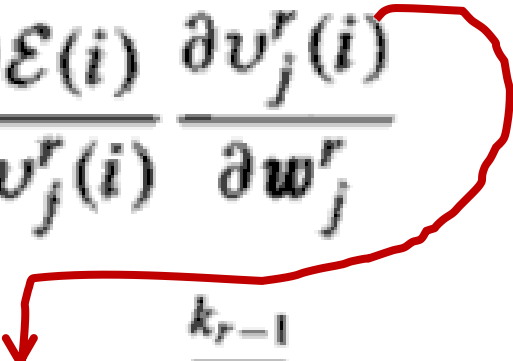
Define the Terms

$$\frac{\partial \mathcal{E}(i)}{\partial \mathbf{w}_j'} = \frac{\partial \mathcal{E}(i)}{\partial v_j'(i)} \frac{\partial v_j'(i)}{\partial \mathbf{w}_j'}$$

Define the Terms

$$\frac{\partial \mathcal{E}(i)}{\partial \mathbf{w}_j^r} = \frac{\partial \mathcal{E}(i)}{\partial v_j^r(i)} \boxed{\frac{\partial v_j^r(i)}{\partial \mathbf{w}_j^r}}$$

Define the Terms

$$\frac{\partial \mathcal{E}(i)}{\partial \mathbf{w}_j^r} = \frac{\partial \mathcal{E}(i)}{\partial v_j^r(i)} \frac{\partial v_j^r(i)}{\partial \mathbf{w}_j^r}$$


Recall, $v_j^r(i) = \sum_{k=1}^{k_r-1} w_{jk}^r y_k^{r-1}(i) + w_{jo}^r \equiv \sum_{k=0}^{k_r-1} w_{jk}^r y_k^{r-1}(i)$

Define the Terms

$$\frac{\partial \mathcal{E}(i)}{\partial \mathbf{w}_j^r} = \frac{\partial \mathcal{E}(i)}{\partial v_j^r(i)} \frac{\partial v_j^r(i)}{\partial \mathbf{w}_j^r}$$

Recall, $v_j^r(i) = \sum_{k=1}^{k_r-1} w_{jk}^r y_k^{r-1}(i) + w_{j0}^r \equiv \sum_{k=0}^{k_r-1} w_{jk}^r y_k^{r-1}(i)$

Therefore, $\frac{\partial}{\partial \mathbf{w}_j^r} v_j^r(i) \equiv \begin{bmatrix} \frac{\partial}{\partial w_{j0}^r} v_j^r(i) \\ \vdots \\ \frac{\partial}{\partial w_{jk_{r-1}}^r} v_j^r(i) \end{bmatrix}$

Define the Terms

$$\frac{\partial \mathcal{E}(i)}{\partial \mathbf{w}_j^r} = \frac{\partial \mathcal{E}(i)}{\partial v_j^r(i)} \frac{\partial v_j^r(i)}{\partial \mathbf{w}_j^r}$$

Recall, $v_j^r(i) = \sum_{k=1}^{k_{r-1}} w_{jk}^r y_k^{r-1}(i) + w_{j0}^r \equiv \sum_{k=0}^{k_{r-1}} w_{jk}^r y_k^{r-1}(i)$

Therefore, $\frac{\partial}{\partial \mathbf{w}_j^r} v_j^r(i) \equiv \begin{bmatrix} \frac{\partial}{\partial w_{j0}^r} v_j^r(i) \\ \vdots \\ \frac{\partial}{\partial w_{jk_{r-1}}^r} v_j^r(i) \end{bmatrix} = \begin{bmatrix} +1 \\ y_1^{r-1}(i) \\ \vdots \\ y_{k_{r-1}}^{r-1}(i) \end{bmatrix}$

Define the Terms

$$\frac{\partial \mathcal{E}(i)}{\partial \mathbf{w}_j^r} = \frac{\partial \mathcal{E}(i)}{\partial v_j^r(i)} \frac{\partial v_j^r(i)}{\partial \mathbf{w}_j^r}$$

Recall, $v_j^r(i) = \sum_{k=1}^{k_r-1} w_{jk}^r y_k^{r-1}(i) + w_{j0}^r \equiv \sum_{k=0}^{k_r-1} w_{jk}^r y_k^{r-1}(i)$

Therefore, $\frac{\partial}{\partial \mathbf{w}_j^r} v_j^r(i) \equiv \begin{bmatrix} \frac{\partial}{\partial w_{j0}^r} v_j^r(i) \\ \vdots \\ \frac{\partial}{\partial w_{jk_{r-1}}^r} v_j^r(i) \end{bmatrix} = \mathbf{y}^{r-1}(i)$

Define the Terms

The diagram illustrates the chain rule for the derivative of the error $\mathcal{E}(i)$ with respect to the weight w_j^r . The main equation is:

$$\frac{\partial \mathcal{E}(i)}{\partial w_j^r} = \frac{\partial \mathcal{E}(i)}{\partial v_j^r(i)} \frac{\partial v_j^r(i)}{\partial w_j^r}$$

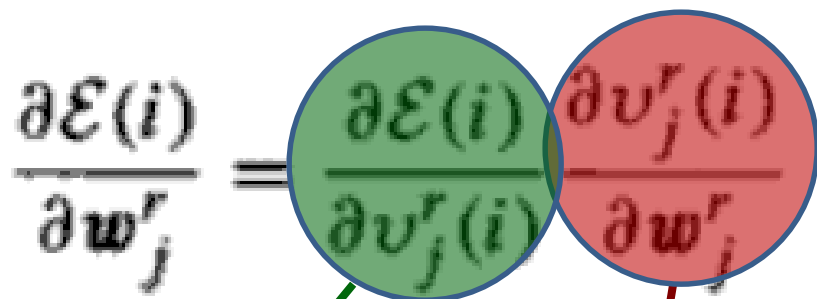
The first term, $\frac{\partial \mathcal{E}(i)}{\partial v_j^r(i)}$, is highlighted in a green circle. A green arrow points from this term to its definition:

$$\frac{\partial \mathcal{E}(i)}{\partial v_j^r(i)} \equiv \delta_j^r(i)$$

The second term, $\frac{\partial v_j^r(i)}{\partial w_j^r}$, is highlighted in a red circle. A red arrow points from this term to its definition:

$$\frac{\partial}{\partial w_j^r} v_j^r(i) \equiv \begin{bmatrix} \frac{\partial}{\partial w_{j0}^r} v_j^r(i) \\ \vdots \\ \frac{\partial}{\partial w_{jk_{r-1}}^r} v_j^r(i) \end{bmatrix} = \mathbf{y}^{r-1}(i)$$

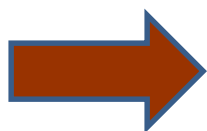
Define the Terms

$$\frac{\partial \mathcal{E}(i)}{\partial \mathbf{w}_j^r} = \frac{\partial \mathcal{E}(i)}{\partial v_j^r(i)} \frac{\partial v_j^r(i)}{\partial \mathbf{w}_j^r}$$


$$\frac{\partial \mathcal{E}(i)}{\partial v_j^r(i)} \equiv \delta_j^r(i)$$

$$\frac{\partial}{\partial \mathbf{w}_j^r} v_j^r(i) \equiv \begin{bmatrix} \frac{\partial}{\partial w_{j0}^r} v_j^r(i) \\ \vdots \\ \frac{\partial}{\partial w_{jk_{r-1}}^r} v_j^r(i) \end{bmatrix} = \mathbf{y}^{r-1}(i)$$

$$\Delta \mathbf{w}_j^r = -\mu \sum_{i=1}^N \frac{\partial E(i)}{\partial \mathbf{w}_j^r(i)}$$



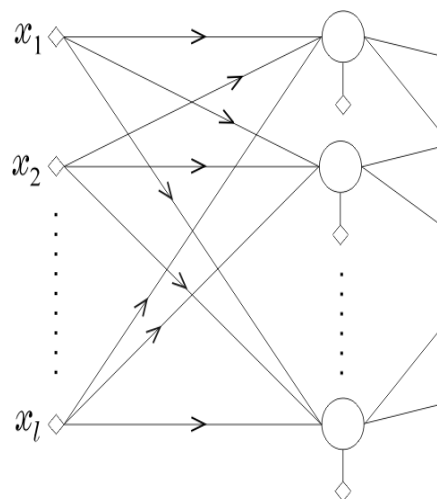
$$\Delta \mathbf{w}_j^r = -\mu \sum_{i=1}^N \delta_j^r(i) \mathbf{y}^{r-1}(i)$$

Define the Terms

$$\Delta \mathbf{w}_j^r = -\mu \sum_{i=1}^N \delta_j^r(i) \mathbf{y}^{r-1}(i)$$

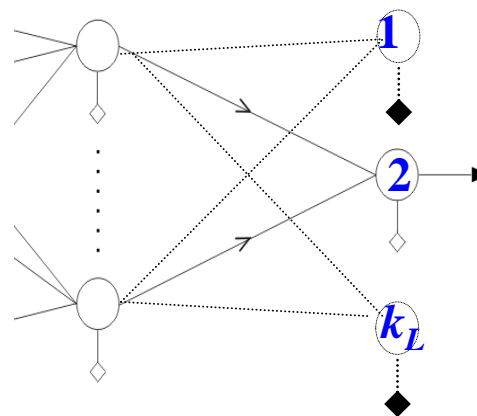
Define the Terms

$$\Delta \mathbf{w}_j^r = -\mu \sum_{i=1}^N \delta_j^r(i) \mathbf{y}^{r-1}(i)$$



input
layer

1st hidden
layer

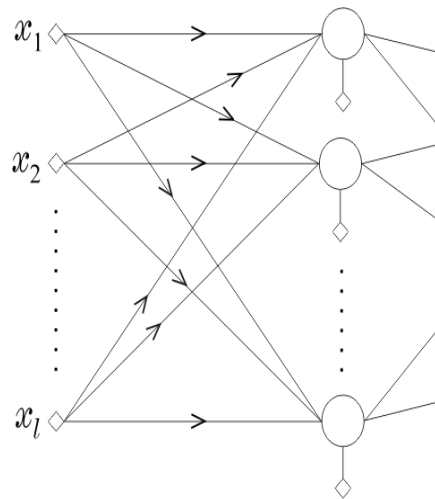


$(L-1)^{\text{th}}$ hidden
layer

L^{th} or
output
layer

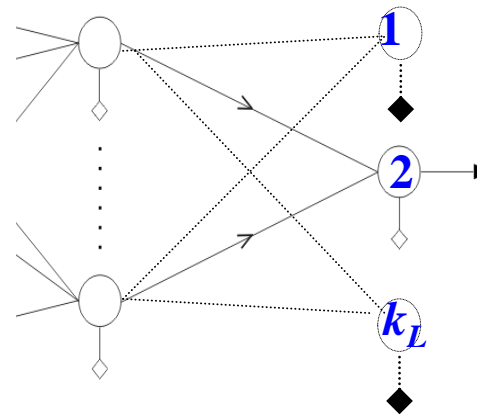
Define the Terms

$$\Delta \mathbf{w}_j^r = -\mu \sum_{i=1}^N \delta_j^r(i) \mathbf{y}^{r-1}(i)$$



input
layer

1st hidden
layer



$(L-1)^{\text{th}}$ hidden
layer

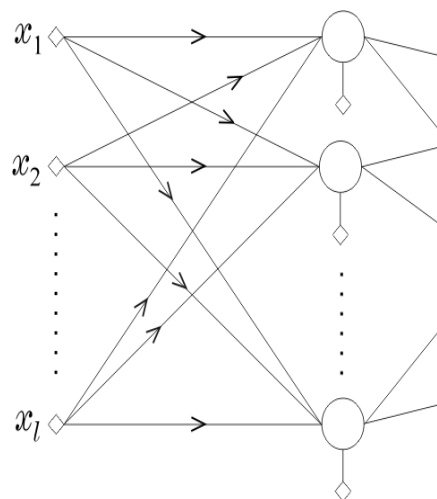
L^{th} or
output
layer



1. Find δ at
Layer L

Define the Terms

$$\Delta \mathbf{w}_j^r = -\mu \sum_{i=1}^N \delta_j^r(i) \mathbf{y}^{r-1}(i)$$



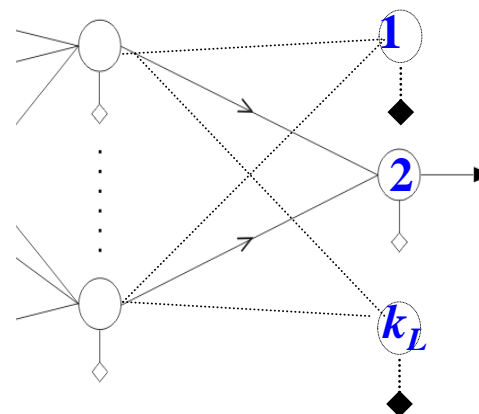
input
layer

1st hidden
layer

($L-1$)th hidden
layer

L^{th} or
output
layer

so on . . .



2. Find δ at
Layer $L-1$

1. Find δ at
Layer L

Define the Terms

- Calculate $\delta_j^r(i) = \frac{\partial E(i)}{\partial v_j^r(i)}$

Define the Terms

- Calculate $\delta_j^r(i) = \frac{\partial E(i)}{\partial v_j^r(i)}$
- For $r = L$

$$\delta_j^L(i) = \frac{\partial E(i)}{\partial v_j^L(i)}$$

Define the Terms

- Calculate $\delta_j^r(i) = \frac{\partial E(i)}{\partial v_j^r(i)}$
- For $r = L$

$$\delta_j^L(i) = \frac{\partial E(i)}{\partial v_j^L(i)}$$

$$E(i) = \frac{1}{2} \sum_{m=1}^{k_L} e_m^2(i) \equiv \frac{1}{2} \sum_{m=1}^{k_L} (f(v_m^L(i)) - y_m(i))^2$$

Define the Terms

- Calculate $\delta_j^r(i) = \frac{\partial E(i)}{\partial v_j^r(i)}$
- For $r = L$

$$\delta_j^L(i) = \frac{\partial E(i)}{\partial v_j^L(i)}$$

$$E(i) = \frac{1}{2} \sum_{m=1}^{k_L} e_m^2(i) \equiv \frac{1}{2} \sum_{m=1}^{k_L} (f(v_m^L(i)) - y_m(i))^2$$

$$\delta_j^L(i) = e_j(i) f'(v_j^L(i))$$

Define the Terms

- For $r < L$
- Calculate $\delta_j^{r-1}(i)$ from $\delta_j^r(i)$

Define the Terms

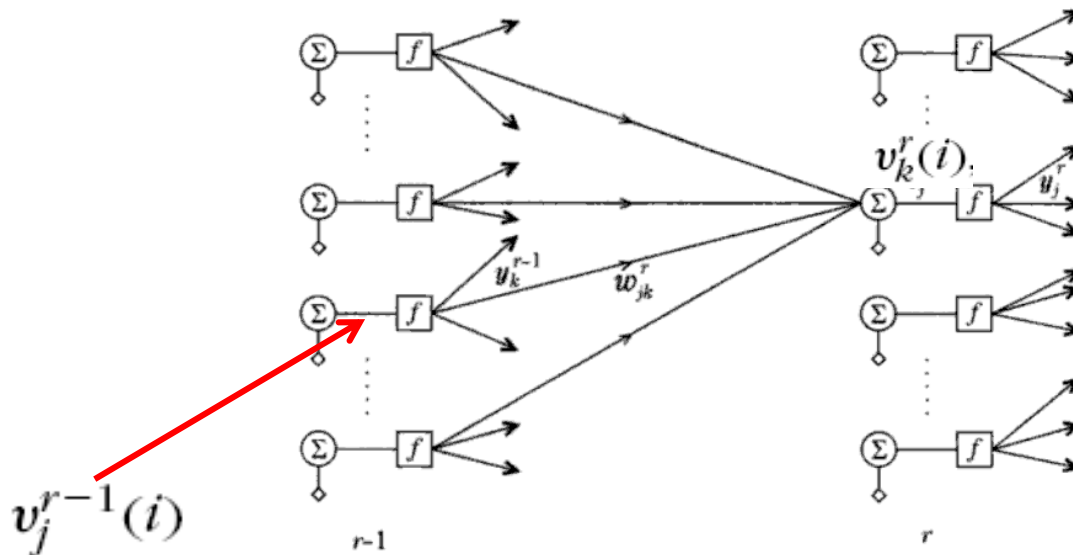
- For $r < L$
- Calculate $\delta_j^{r-1}(i)$ from $\delta_j^r(i)$

- We know,

$$\delta_j^{r-1}(i) = \frac{\partial E(i)}{\partial v_j^{r-1}(i)}$$

Define the Terms

- We need to calculate,
$$\delta_j^{r-1}(i) = \frac{\partial E(i)}{\partial v_j^{r-1}(i)}$$
- However, $v_j^{r-1}(i)$ influences all $v_k^r(i)$, for $k = 1, 2, 3, \dots, k_r$



Define the Terms

- We need to calculate,

$$\delta_j^{r-1}(i) = \frac{\partial E(i)}{\partial v_j^{r-1}(i)}$$

- However, $v_j^{r-1}(i)$ influences all $v_k^r(i)$, for $k = 1, 2, 3, \dots, k_r$

- Therefore,

$$\frac{\partial E(i)}{\partial v_j^{r-1}(i)} = \sum_{k=1}^{k_r} \frac{\partial E(i)}{\partial v_k^r(i)} \frac{\partial v_k^r(i)}{\partial v_j^{r-1}(i)}$$

Define the Terms

- For $r < L$
- Calculate

$$\delta_j^{r-1}(i) = \frac{\partial E(i)}{\partial v_j^{r-1}(i)}$$

$$\frac{\partial E(i)}{\partial v_j^{r-1}(i)} = \sum_{k=1}^{k_r} \frac{\partial E(i)}{\partial v_k^r(i)} \frac{\partial v_k^r(i)}{\partial v_j^{r-1}(i)}$$

Define the Terms

- For $r < L$
- Calculate

$$\delta_j^{r-1}(i) = \frac{\partial E(i)}{\partial v_j^{r-1}(i)}$$

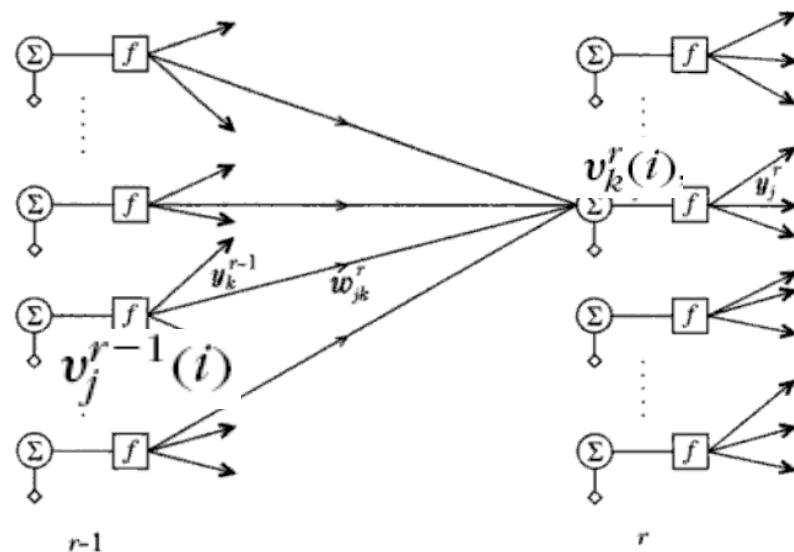
$$\frac{\partial E(i)}{\partial v_j^{r-1}(i)} = \sum_{k=1}^{k_r} \frac{\partial E(i)}{\partial v_k^r(i)} \frac{\partial v_k^r(i)}{\partial v_j^{r-1}(i)}$$

$$\delta_j^{r-1}(i) = \sum_{k=1}^{k_r} \delta_k^r(i) \frac{\partial v_k^r(i)}{\partial v_j^{r-1}(i)}$$

Define the Terms

$$\delta_j^{r-1}(i) = \sum_{k=1}^{k_r} \delta_k^r(i) \frac{\partial v_k^r(i)}{\partial v_j^{r-1}(i)}$$

$$\frac{\partial v_k^r(i)}{\partial v_j^{r-1}(i)} = \frac{\partial \left[\sum_{m=0}^{k_{r-1}} w_{km}^r y_m^{r-1}(i) \right]}{\partial v_j^{r-1}(i)}$$



where, $y_m^{r-1}(i) = f(v_m^{r-1}(i))$

Define the Terms

$$\delta_j^{r-1}(i) = \sum_{k=1}^{k_r} \delta_k^r(i) \frac{\partial v_k^r(i)}{\partial v_j^{r-1}(i)}$$

$$\frac{\partial v_k^r(i)}{\partial v_j^{r-1}(i)} = \frac{\partial \left[\sum_{m=0}^{k_{r-1}} w_{km}^r y_m^{r-1}(i) \right]}{\partial v_j^{r-1}(i)} \quad \text{where, } y_m^{r-1}(i) = f(v_m^{r-1}(i))$$

then,

$$\frac{\partial v_k^r(i)}{\partial v_j^{r-1}(i)} = w_{kj}^r f'(v_j^{r-1}(i))$$

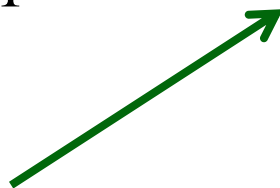
Define the Terms

$$\delta_j^{r-1}(i) = \sum_{k=1}^{k_r} \delta_k^r(i) \frac{\partial v_k^r(i)}{\partial v_j^{r-1}(i)}$$

$$\frac{\partial v_k^r(i)}{\partial v_j^{r-1}(i)} = w_{kj}^r f'(v_j^{r-1}(i))$$

Define the Terms

$$\delta_j^{r-1}(i) = \sum_{k=1}^{k_r} \delta_k^r(i) \frac{\partial v_k^r(i)}{\partial v_j^{r-1}(i)}$$


$$\frac{\partial v_k^r(i)}{\partial v_j^{r-1}(i)} = w_{kj}^r f'(v_j^{r-1}(i))$$

Define the Terms

$$\delta_j^{r-1}(i) = \sum_{k=1}^{k_r} \delta_k^r(i) \frac{\partial v_k^r(i)}{\partial v_j^{r-1}(i)}$$

$$\frac{\partial v_k^r(i)}{\partial v_j^{r-1}(i)} = w_{kj}^r f'(v_j^{r-1}(i))$$

$$\delta_j^{r-1}(i) = \left[\sum_{k=1}^{k_r} \delta_k^r(i) w_{kj}^r \right] f'(v_j^{r-1}(i))$$

Define the Terms

$$\delta_j^{r-1}(i) = \left[\sum_{k=1}^{k_r} \delta_k^r(i) w_{kj}^r \right] f'(v_j^{r-1}(i))$$

$$\delta_j^{r-1}(i) = e_j^{r-1}(i) f'(v_j^{r-1}(i))$$

where,
$$e_j^{r-1}(i) = \sum_{k=1}^{k_r} \delta_k^r(i) w_{kj}^r$$

Define the Terms

- Only remaining is the derivative of the logistic function:

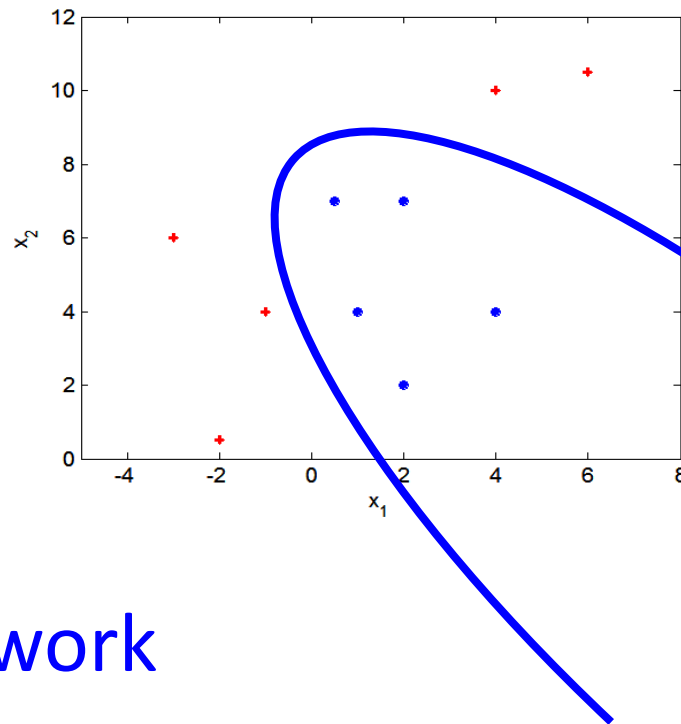
$$f'(x) = \alpha f(x)(1 - f(x))$$

The Algorithm

- Initialization:
 - Start with small random weights
- Forward Computations: $v_j^r(i), y_j^r(i) = f(\tilde{v}_j^r(i)),$
- Backward Computation: $\delta_j^L(i)$ and $\delta_j^{r-1}(i)$
- Update weight: $w_j^r(\text{new}) = w_j^r(\text{old}) + \Delta w_j^r$

$$\Delta w_j^r = -\mu \sum_{i=1}^N \delta_j^r(i) y^{r-1}(i)$$

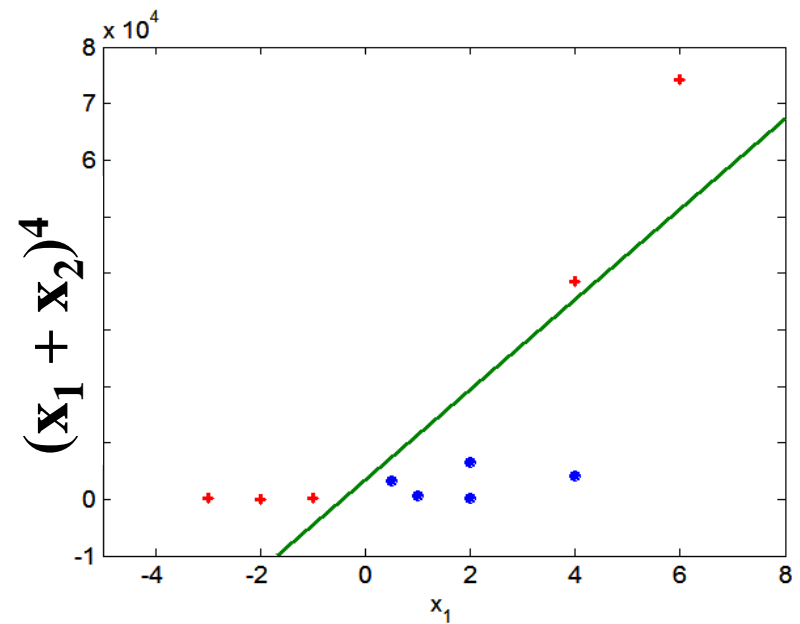
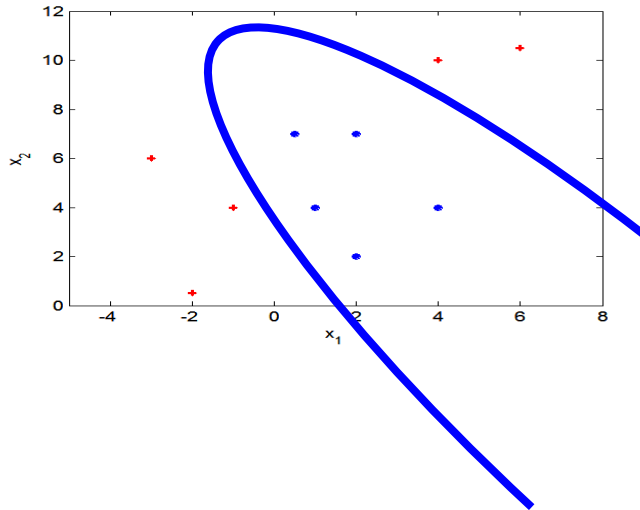
Some Nonlinear Classifiers



- Neural Network
- Decision Tree
- Non-linear SVM

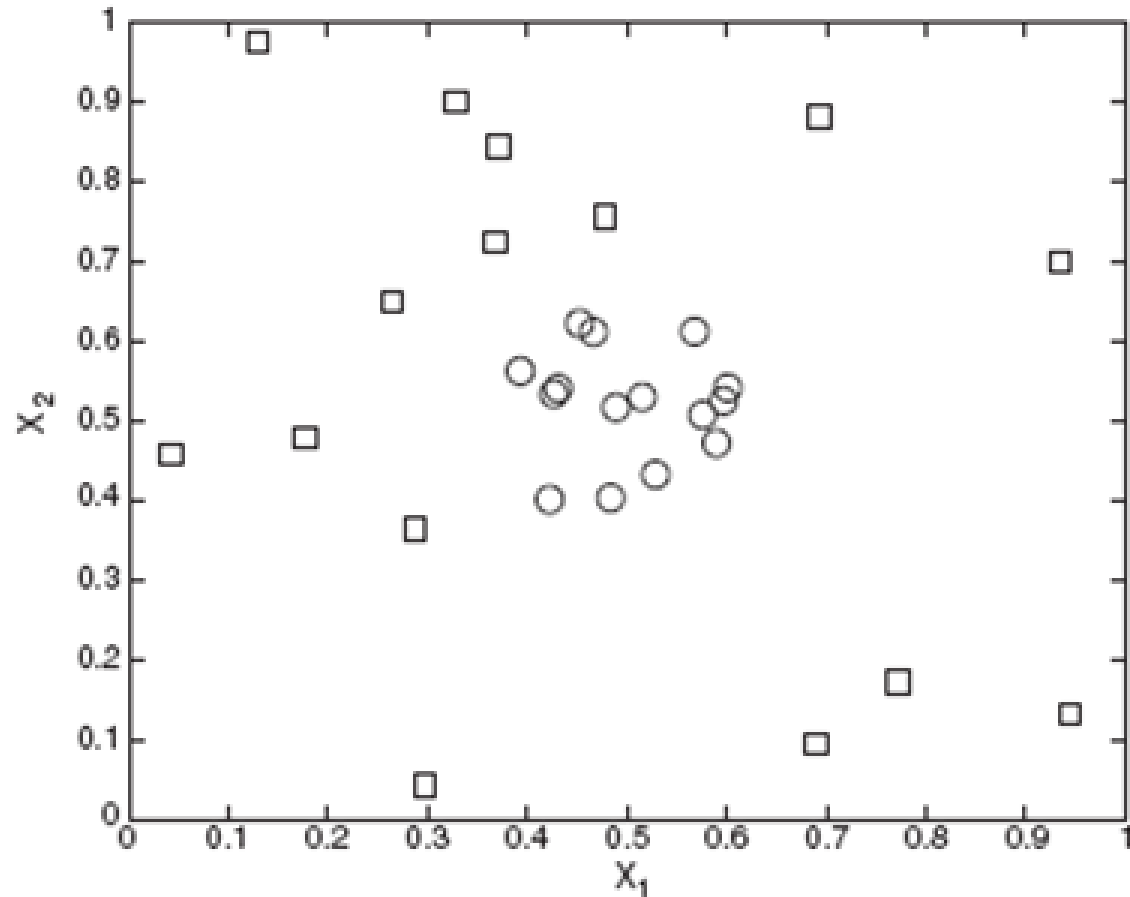
Nonlinear SVM

- Transform data into higher dimensional space



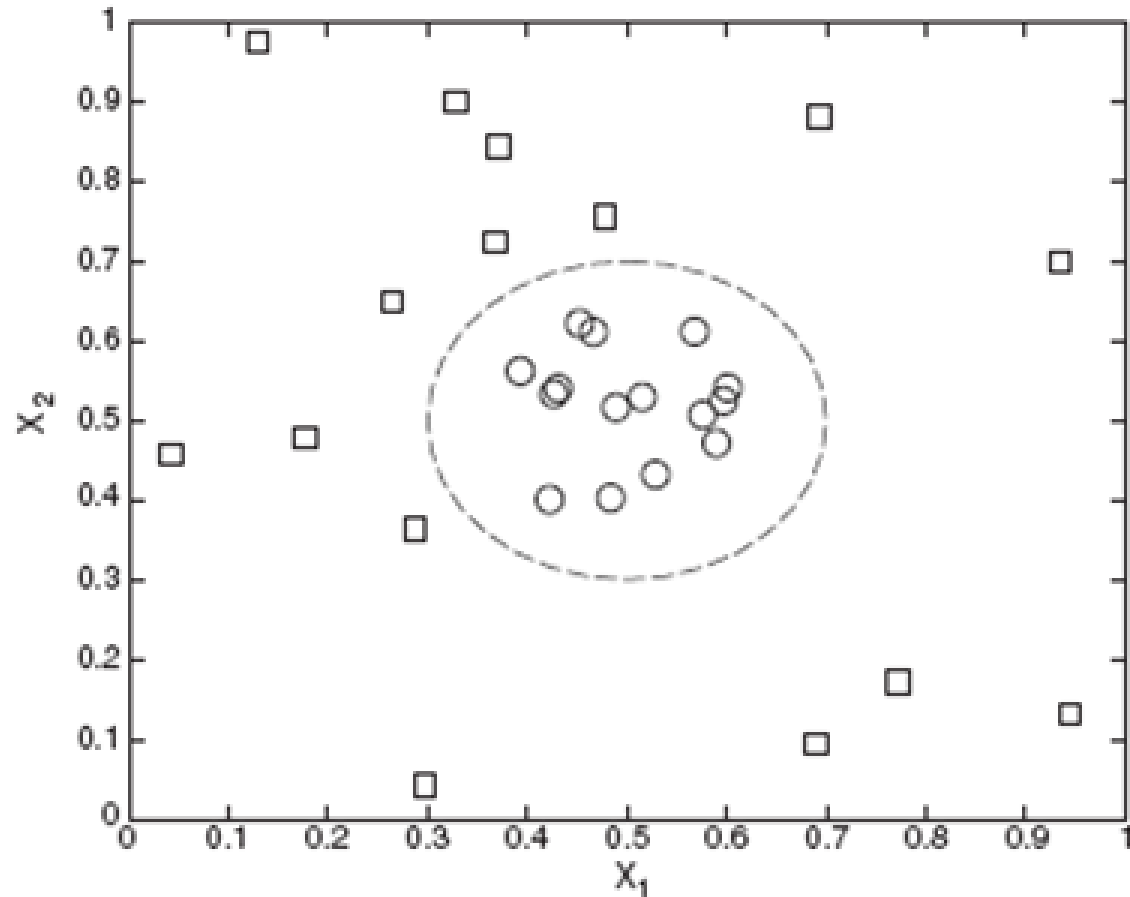
Nonlinear SVM

- Another Example



Nonlinear SVM

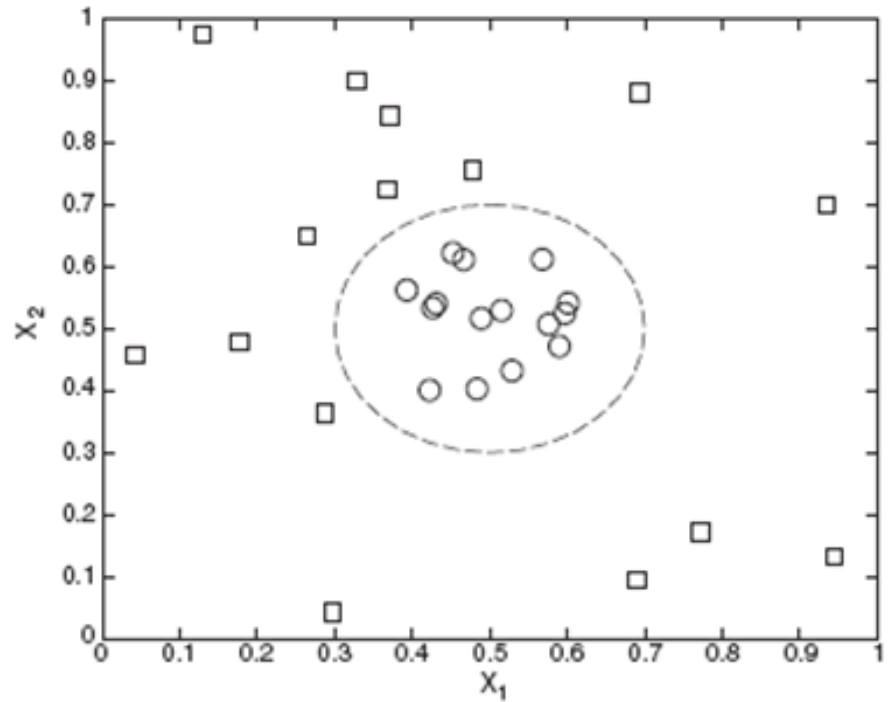
- Another Example



Nonlinear SVM

- Decision boundary:

$$\sqrt{(x_1 - 0.5)^2 + (x_2 - 0.5)^2} = 0.2$$



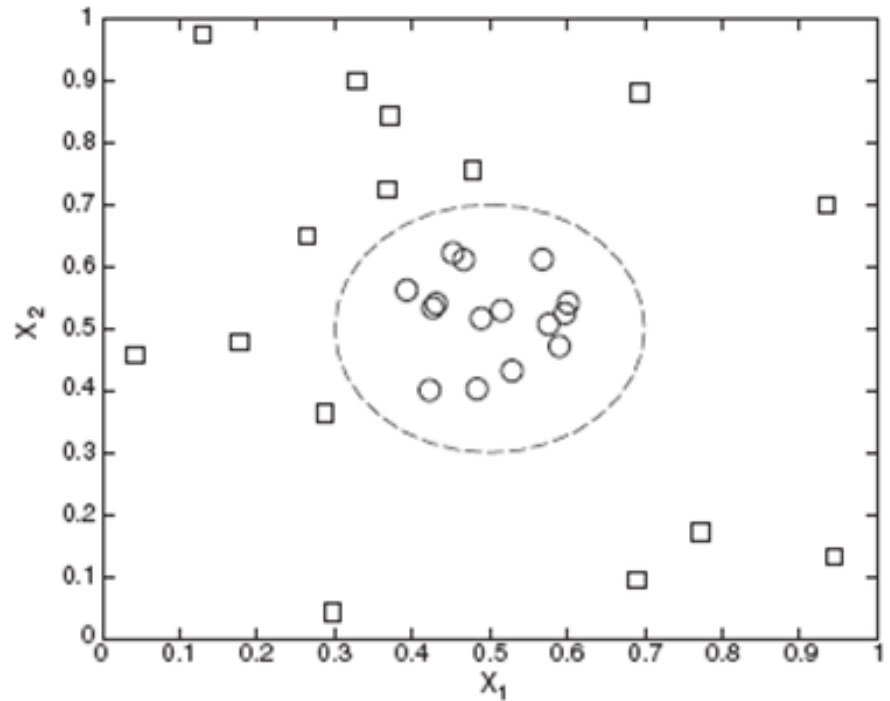
Nonlinear SVM

- Decision boundary:

$$\sqrt{(x_1 - 0.5)^2 + (x_2 - 0.5)^2} = 0.2$$

- 2 classes are defined as

$$y(x_1, x_2) = \begin{cases} 1, & \text{if } \sqrt{(x_1 - 0.5)^2 + (x_2 - 0.5)^2} > 0.2 \\ -1, & \text{otherwise} \end{cases}$$



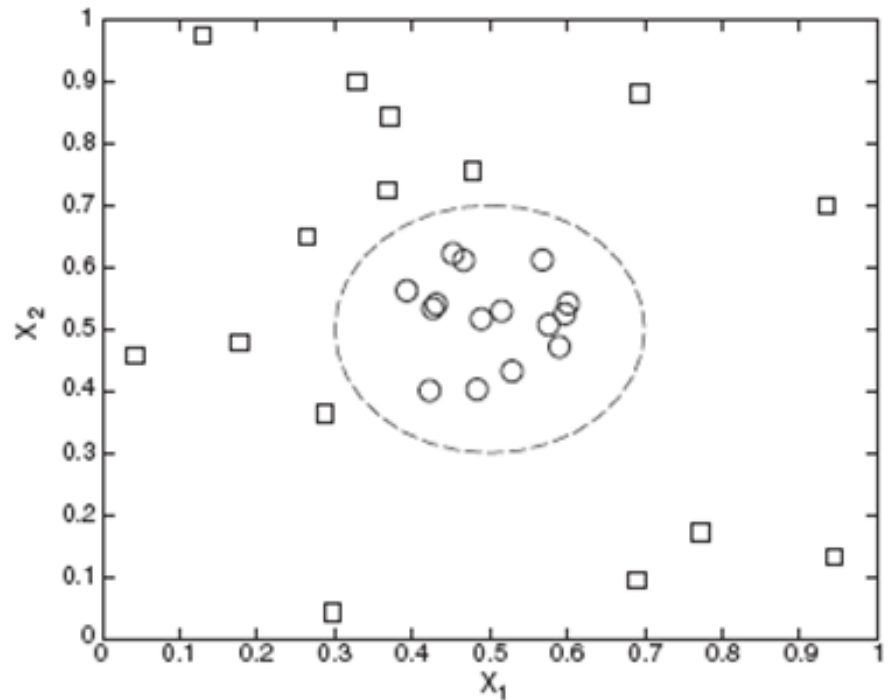
Nonlinear SVM

- Decision boundary:

$$\sqrt{(x_1 - 0.5)^2 + (x_2 - 0.5)^2} = 0.2$$

can be written as

$$x_1^2 - x_1 + x_2^2 - x_2 = -0.46$$



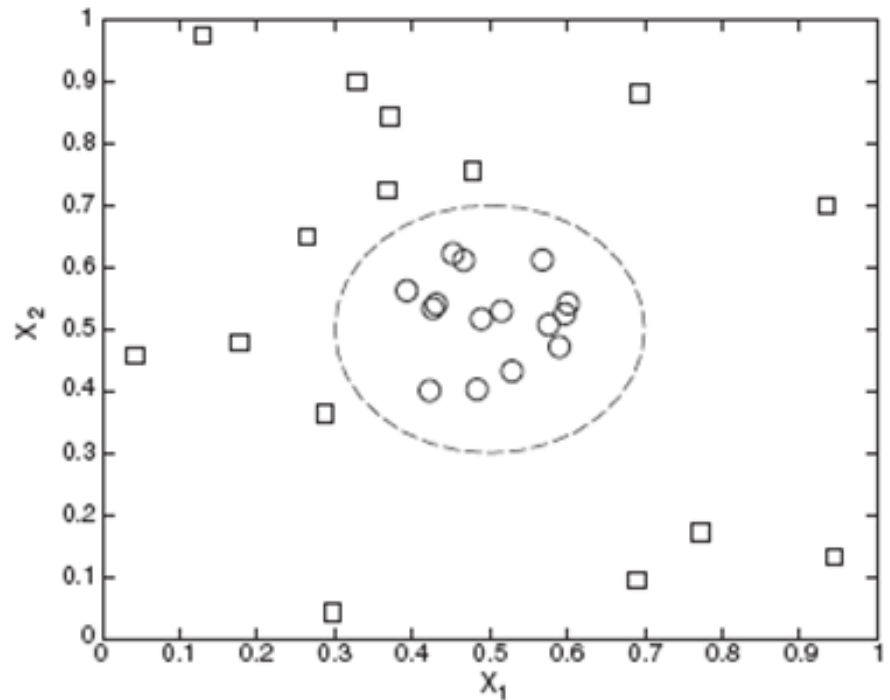
Nonlinear SVM

- Decision boundary:

$$\sqrt{(x_1 - 0.5)^2 + (x_2 - 0.5)^2} = 0.2$$

can be written as

$$\underbrace{x_1^2 - x_1}_{y_1} + \underbrace{x_2^2 - x_2}_{y_2} = -0.46$$



Nonlinear SVM

- Decision boundary:

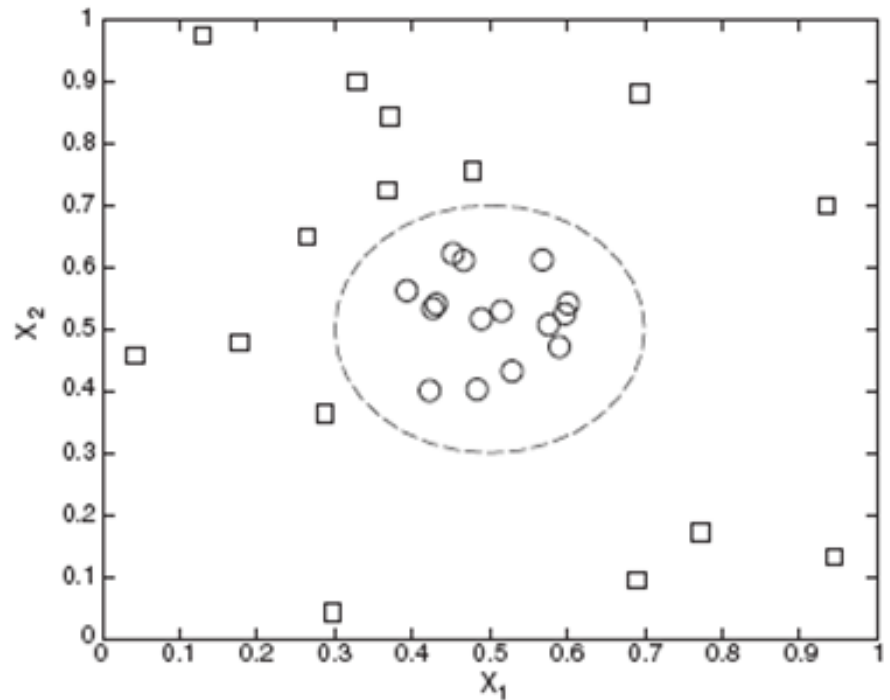
$$\sqrt{(x_1 - 0.5)^2 + (x_2 - 0.5)^2} = 0.2$$

can be written as

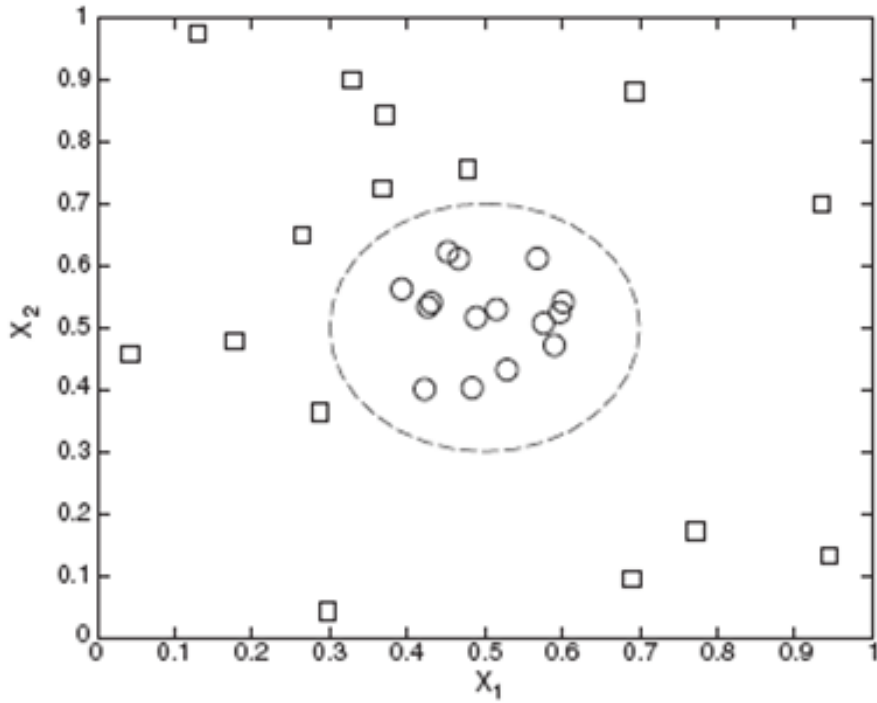
$$\underbrace{x_1^2 - x_1}_{y_1} + \underbrace{x_2^2 - x_2}_{y_2} = -0.46$$



$$y_1 + y_2 = -0.46$$

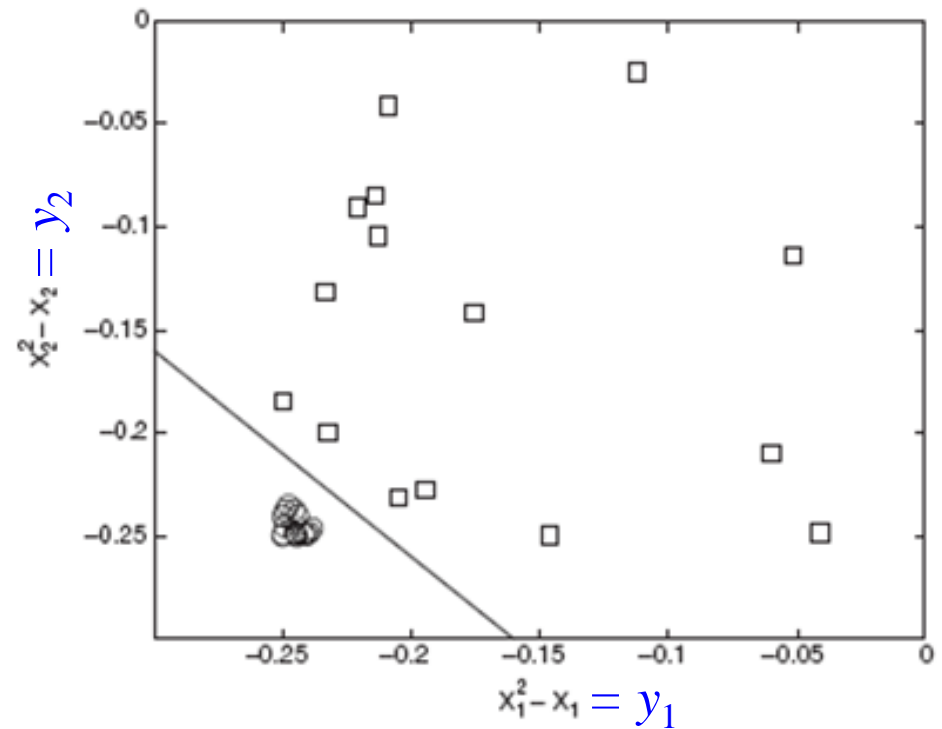
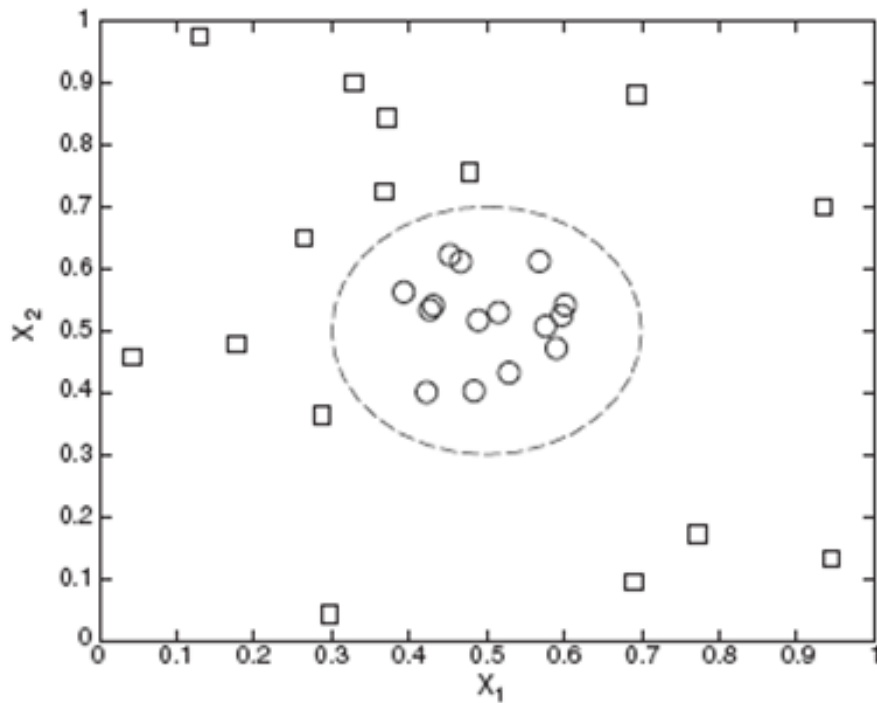


Nonlinear SVM



$$\sqrt{(x_1 - 0.5)^2 + (x_2 - 0.5)^2} = 0.2$$

Nonlinear SVM



$$\sqrt{(x_1 - 0.5)^2 + (x_2 - 0.5)^2} = 0.2$$

$$x_1^2 - x_1 + x_2^2 - x_2 = -0.46$$

OR,

$$y_1 + y_2 = -0.46$$

Nonlinear SVM

- We need a transformation like this

$$\Phi : (x_1, x_2) \rightarrow (x_1^2 - x_1, x_2^2 - x_2)$$

Nonlinear SVM

- We need a transformation like this

$$\Phi : (x_1, x_2) \rightarrow (x_1^2 - x_1, x_2^2 - x_2)$$

OR,

$$\Phi : (x_1, x_2) \rightarrow (y_1, y_2)$$

Nonlinear SVM

- We need a transformation like this

$$\Phi : (x_1, x_2) \rightarrow (x_1^2 - x_1, x_2^2 - x_2)$$

OR,

$$\Phi : (x_1, x_2) \rightarrow (y_1, y_2)$$

OR,

$$\Phi : \mathbf{x} \rightarrow \mathbf{y}$$

Nonlinear SVM

- We need a transformation like this

$$\Phi : (x_1, x_2) \rightarrow (x_1^2 - x_1, x_2^2 - x_2)$$

- OR, more generally:

$$\Phi : (x_1, x_2) \rightarrow (x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1x_2, 1)$$

Nonlinear SVM

- With the transform

$$\Phi : (x_1, x_2) \rightarrow (x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1x_2, 1)$$

The equation of the classifier will be of the form:

$$w_5x_1^2 + w_4x_2^2 + w_3\sqrt{2}x_1 + w_2\sqrt{2}x_2 + w_1\sqrt{2}x_1x_2 + w_0 = 0$$

Nonlinear SVM

- With the transform

$$\Phi : (x_1, x_2) \rightarrow (x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1x_2, 1)$$

The equation of the classifier will be of the form:

$$w_5x_1^2 + w_4x_2^2 + w_3\sqrt{2}x_1 + w_2\sqrt{2}x_2 + w_1\sqrt{2}x_1x_2 + w_01 = 0$$

$$w_5y_5 + w_4y_4 + w_3y_3 + w_2y_2 + w_1y_1 + w_01 = 0$$

Nonlinear SVM

$$\Phi : (x_1, x_2) \rightarrow (x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1x_2, 1) \rightarrow (y_1, y_2, \dots, y_5)$$

$$w_5x_1^2 + w_4x_2^2 + w_3\sqrt{2}x_1 + w_2\sqrt{2}x_2 + w_1\sqrt{2}x_1x_2 + w_01$$

- The main idea: *linearly separability* increases as the feature dimension increases

Formulation of a Nonlinear SVM

- With the new feature vectors $\Phi(\vec{x})$, replace all \mathbf{x} with $\Phi(\vec{x})$ in linear SVM:

- Minimize $L(w) = \frac{\|\vec{w}\|^2}{2}$

- Subject to $y_i (\vec{w} \bullet \Phi(\vec{x}_i) + b) \geq 1$

- The Dual function is:

$$L_D = \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)$$

Nonlinear SVM

- Once we get the solution of λ 's, find w and b using following equations:

$$\vec{w} = \sum_{i=1}^N \lambda_i y_i \Phi(\vec{x}_i)$$

$$\lambda_i \{ y_i (\sum_j \lambda_j y_j \Phi(\vec{x}_j) \cdot \Phi(\vec{x}_i) + b) - 1 \} = 0$$

Nonlinear SVM

- The new object \mathbf{z} is classified as:

$$\begin{aligned} f(\vec{z}) &= \text{sign}(\vec{w} \cdot \Phi(\vec{z}) + b) \\ &= \text{sign}\left(\sum_i \lambda_i y_i \Phi(\vec{x}_i) \cdot \Phi(\vec{z}) + b\right) \end{aligned}$$

Issues in the Design Nonlinear SVM

- The mapping function is often unclear
- The dimensionality increases, leading to high computation

Issues in the Design Nonlinear SVM

- The mapping function is often unclear
- The dimensionality increases, leading to high computation

Solution?

Issues in the Design Nonlinear SVM

- Note the equations:

$$L_D = \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)$$

$$\lambda_i \{ y_i (\sum_j \lambda_j y_j \Phi(\vec{x}_j) \cdot \Phi(\vec{x}_i) + b) - 1 \} = 0$$

$$f(\vec{z}) = \text{sign}(\sum_i \lambda_i y_i \Phi(\vec{x}_i) \cdot \Phi(\vec{z}) + b)$$

Issues in the Design Nonlinear SVM

- Note the equations:

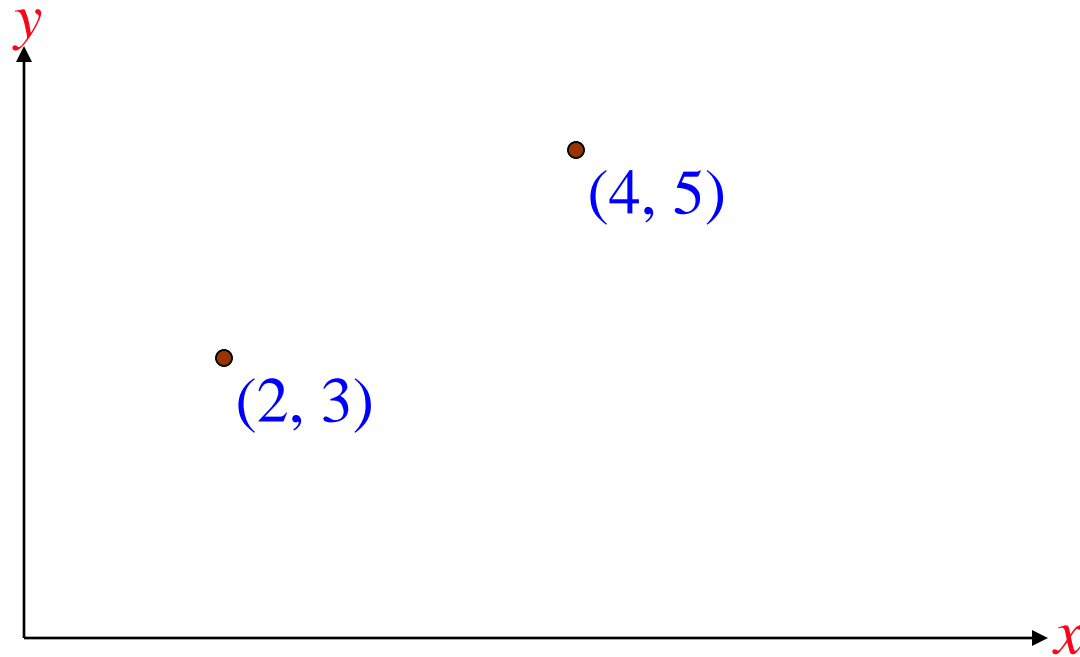
$$L_D = \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)$$

$$\lambda_i \{ y_i (\sum_j \lambda_j y_j \Phi(\vec{x}_j) \cdot \Phi(\vec{x}_i) + b) - 1 \} = 0$$

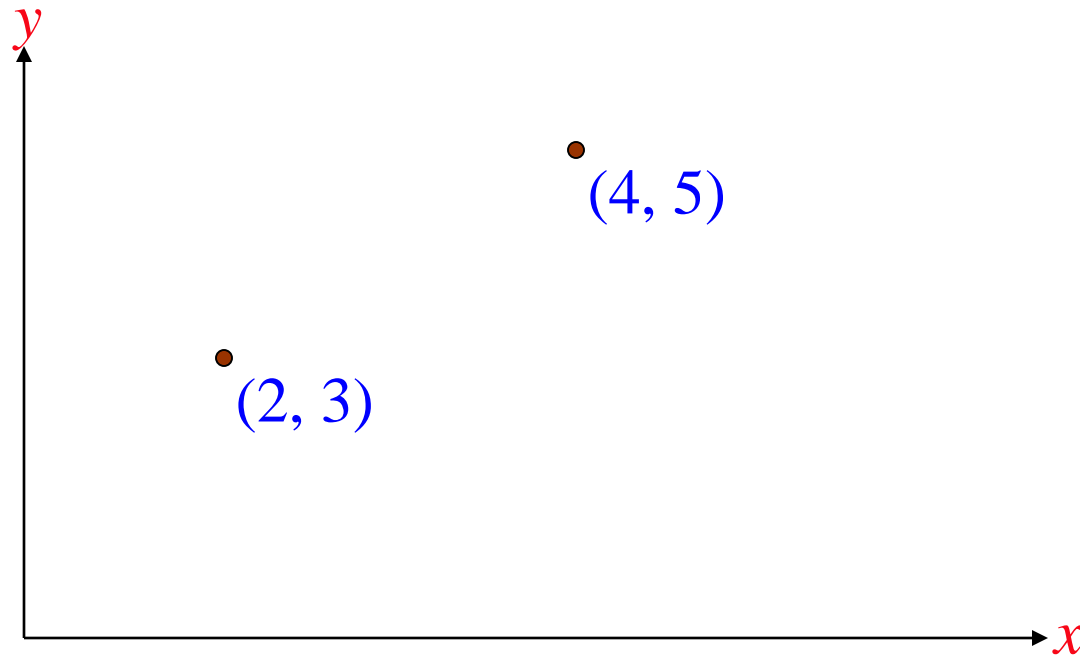
$$f(\vec{z}) = \text{sign}(\sum_i \lambda_i y_i \Phi(\vec{x}_i) \cdot \Phi(\vec{z}) + b)$$

- The dot product $\Phi(\vec{x}_j) \cdot \Phi(\vec{x}_i)$ is a similarity measurement

Similarity/Distance Measurement



Similarity/Distance Measurement



$$\text{distance} = \left[(2-4)^2 + (3-5)^2 \right]^{1/2}$$

$$\text{Cosine Similarity} = \frac{2 \cdot 4 + 3 \cdot 5}{\sqrt{(2^2 + 3^2)} \sqrt{(4^2 + 5^2)}}$$

Issues in the Design Nonlinear SVM

- We can calculate the term $\Phi(\vec{x}_j) \cdot \Phi(\vec{x}_i)$ in the original space using **Kernel trick**

Issues in the Design Nonlinear SVM

- We can calculate the term $\Phi(\vec{x}_j) \cdot \Phi(\vec{x}_i)$ in the original space using **Kernel trick**

Example:

$$\begin{aligned}\Phi(\vec{u}) \cdot \Phi(\vec{v}) &= (u_1^2, u_2^2, \sqrt{2}u_1, \sqrt{2}u_2, \sqrt{2}u_1u_2, 1) \\ &\quad \cdot (v_1^2, v_2^2, \sqrt{2}v_1, \sqrt{2}v_2, \sqrt{2}v_1v_2, 1) \\ &\cdot \\ &\cdot \\ &= (\vec{u} \cdot \vec{v} + 1)^2\end{aligned}$$

Issues in the Design Nonlinear SVM

- Kernel trick

$$\Phi(\vec{u}) \cdot \Phi(\vec{v}) = (\vec{u} \cdot \vec{v} + 1)^2$$

Issues in the Design Nonlinear SVM

- Kernel trick

$$\Phi(\vec{u}) \cdot \Phi(\vec{v}) = (\vec{u} \cdot \vec{v} + 1)^2$$

$$K(\vec{u}, \vec{v}) = (\vec{u} \cdot \vec{v} + 1)^2$$

Issues in the Design Nonlinear SVM

- Some Kernel Functions are:

$$K(\vec{u}, \vec{v}) = (\vec{u} \cdot \vec{v} + 1)^p$$

$$K(\vec{u}, \vec{v}) = e^{-\|\vec{u} - \vec{v}\|^2 / (2\sigma^2)}$$