

Last updated: Nov 26, 2012

# MULTILAYER PERCEPTRONS

# Outline

2

Multilayer Perceptrons

- Combining Linear Classifiers
- Learning Parameters

# Outline

3

Multilayer Perceptrons

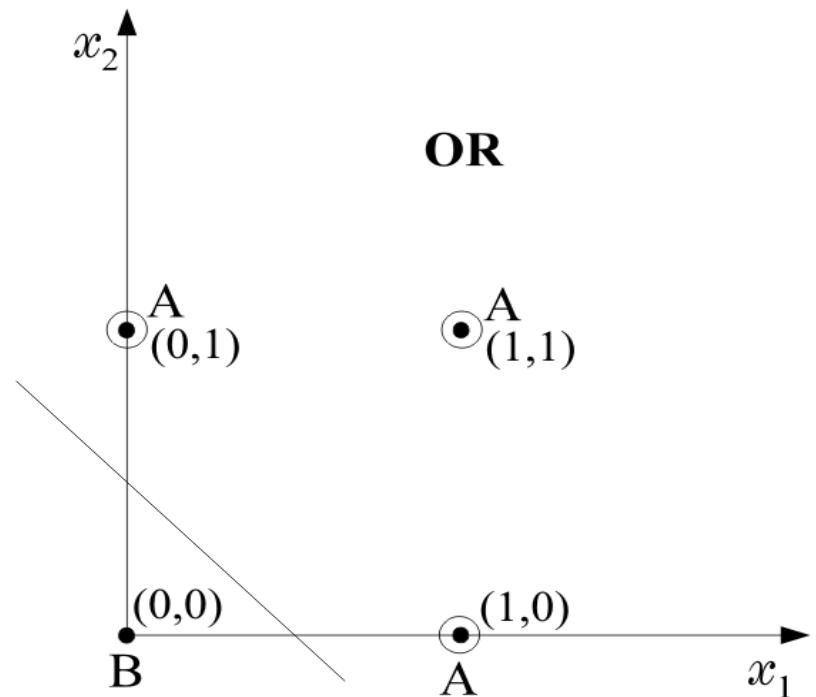
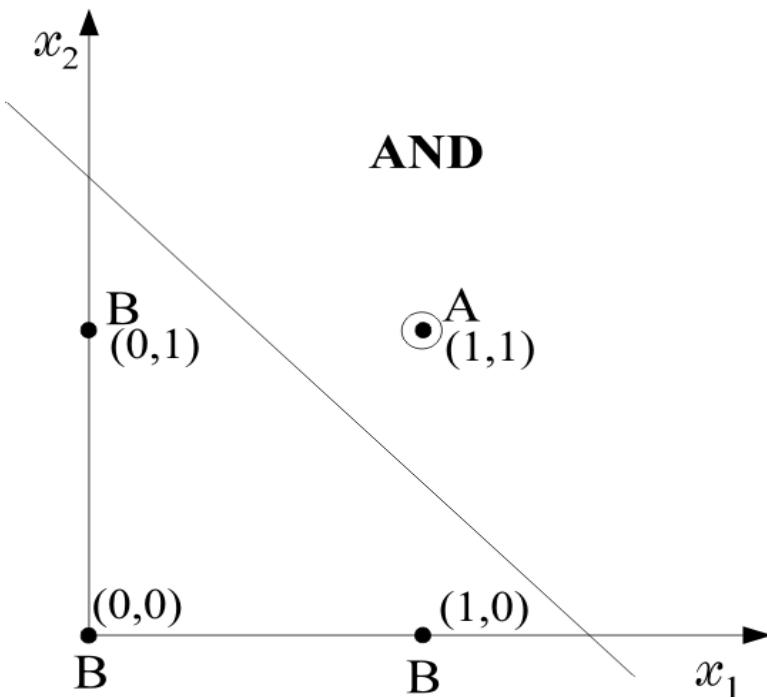
- **Combining Linear Classifiers**
- Learning Parameters

# Implementing Logical Relations

4

Multilayer Perceptrons

- AND and OR operations are linearly separable problems



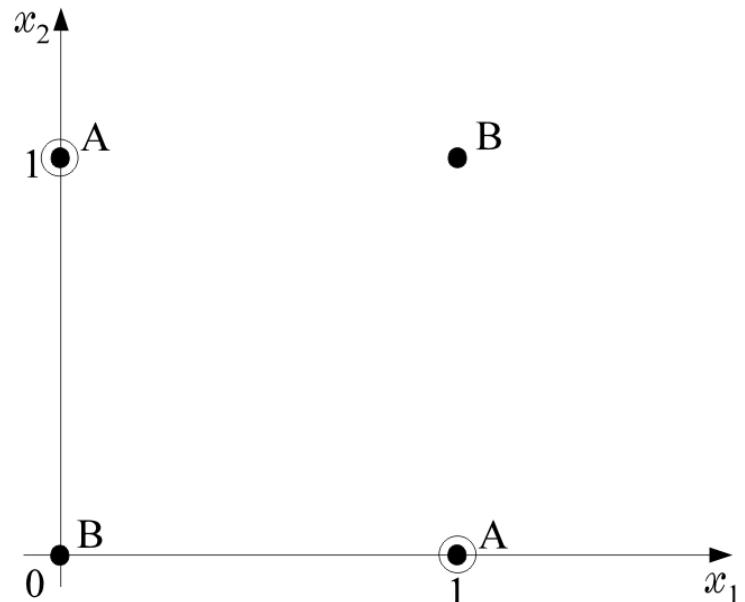
# The XOR Problem

5

Multilayer Perceptrons

- XOR is not linearly separable.

<b><math>x_1</math></b>	<b><math>x_2</math></b>	<b>XOR</b>	<b>Class</b>
0	0	0	B
0	1	1	A
1	0	1	A
1	1	0	B



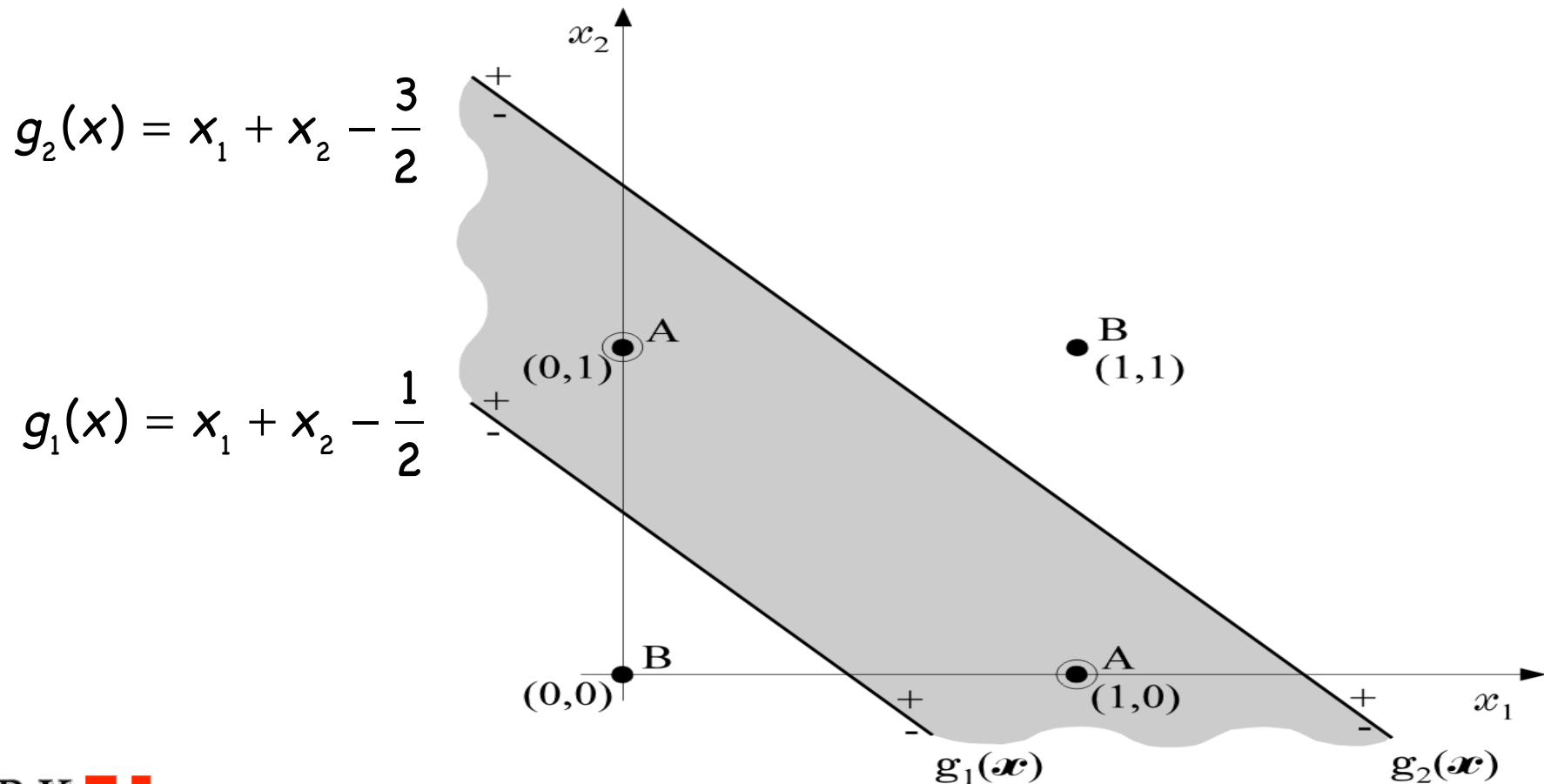
- How can we use linear classifiers to solve this problem?

# Combining two linear classifiers

6

Multilayer Perceptrons

- Idea: use a logical combination of two linear classifiers.



# Combining two linear classifiers

7

Multilayer Perceptrons

Let  $f(x)$  be the unit step activation function:

$$f(x) = 0, \quad x < 0$$

$$f(x) = 1, \quad x \geq 0$$

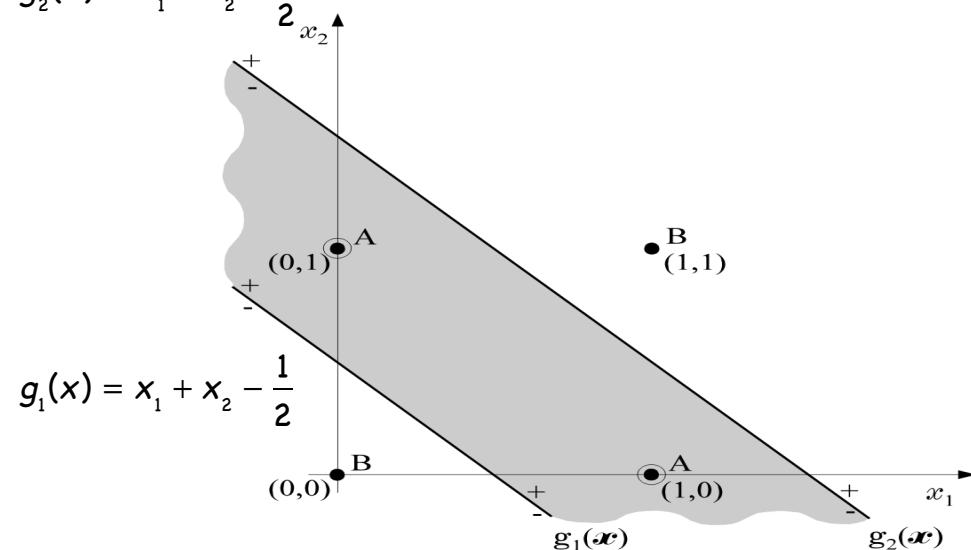
Observe that the classification problem is then solved by

$$f\left(y_1 - y_2 - \frac{1}{2}\right)$$

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

where

$$y_1 = f(g_1(x)) \text{ and } y_2 = f(g_2(x))$$



# Combining two linear classifiers

8

Multilayer Perceptrons

- This calculation can be implemented sequentially:
  1. Compute  $y_1$  and  $y_2$  from  $x_1$  and  $x_2$ .
  2. Compute the decision from  $y_1$  and  $y_2$ .
- Each layer in the sequence consists of one or more linear classifications.
- This is therefore a two-layer perceptron.

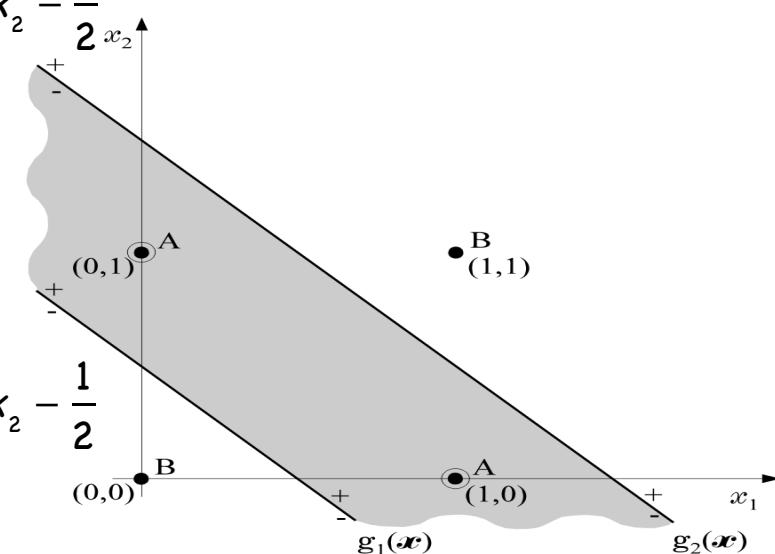
$$f\left(y_1 - y_2 - \frac{1}{2}\right)$$

where

$$y_1 = f(g_1(x)) \text{ and } y_2 = f(g_2(x))$$

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

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



# The Two-Layer Perceptron

9

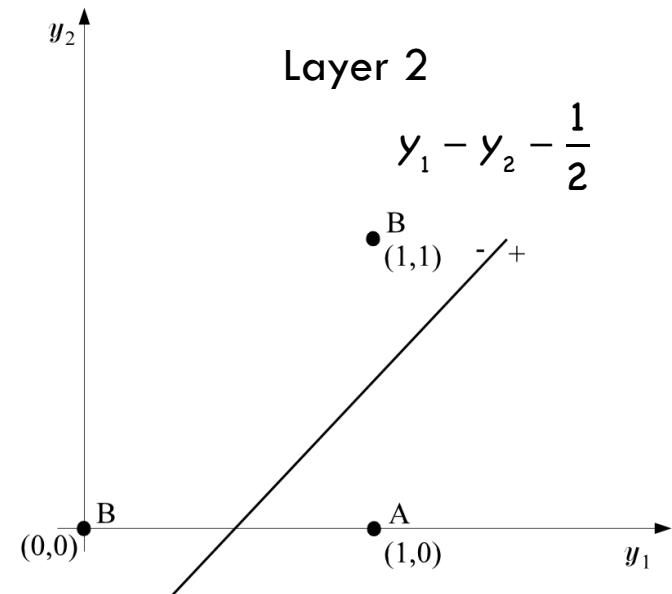
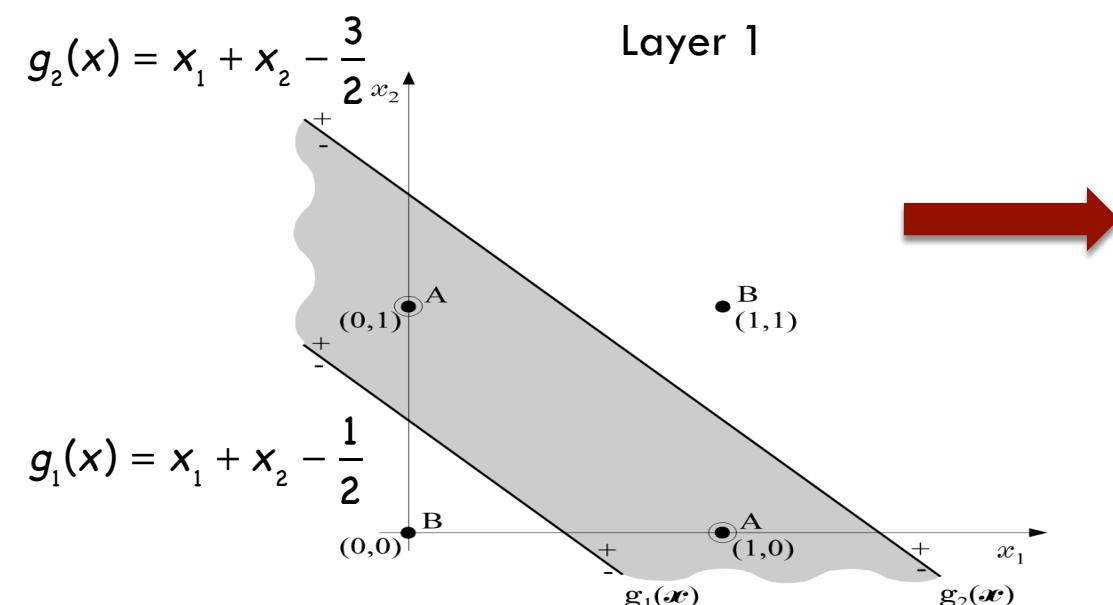
## Multilayer Perceptrons

Layer 1				Layer 2
$x_1$	$x_2$	$y_1$	$y_2$	
0	0	0(-)	0(-)	B(0)
0	1	1(+)	0(-)	A(1)
1	0	1(+)	0(-)	A(1)
1	1	1(+)	1(+)	B(0)

$$f\left(y_1 - y_2 - \frac{1}{2}\right)$$

where

$$y_1 = f(g_1(x)) \text{ and } y_2 = f(g_2(x))$$



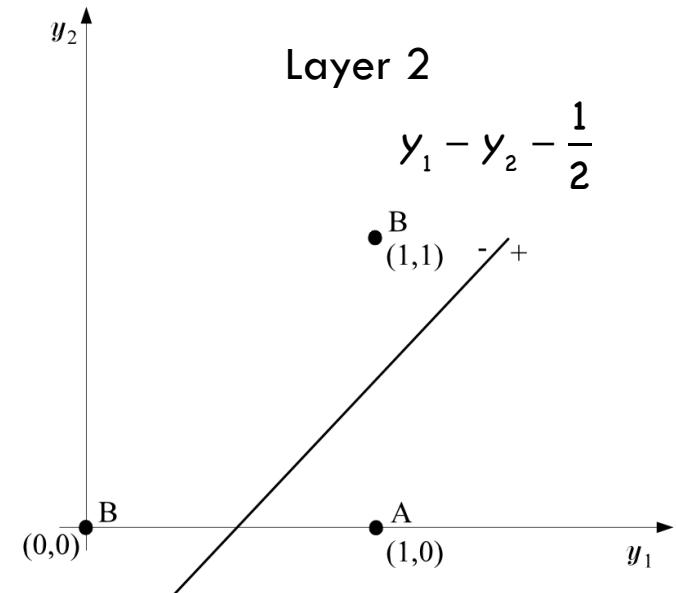
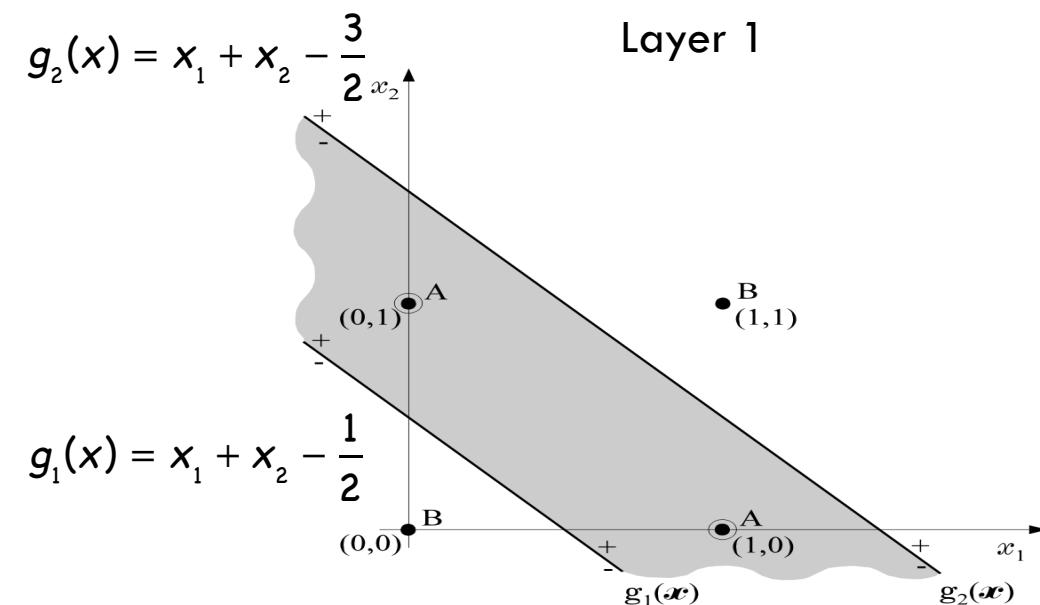
# The Two-Layer Perceptron

10

Multilayer Perceptrons

- The first layer performs a nonlinear mapping that makes the data linearly separable.

$$y_1 = f(g_1(x)) \text{ and } y_2 = f(g_2(x))$$

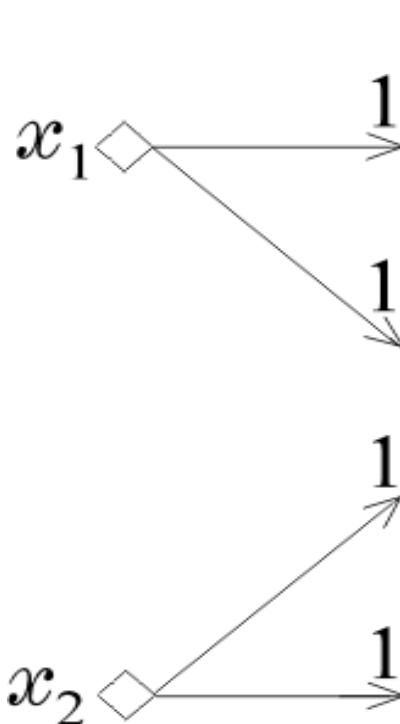


# The Two-Layer Perceptron Architecture

11

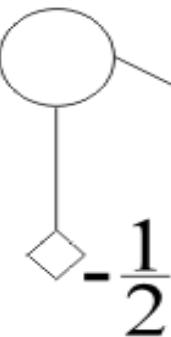
Multilayer Perceptrons

Input Layer

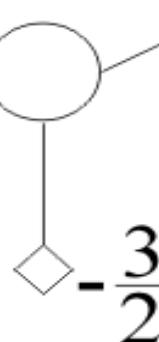


Hidden Layer

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



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



Output Layer

$$y_1 - y_2 - \frac{1}{2}$$



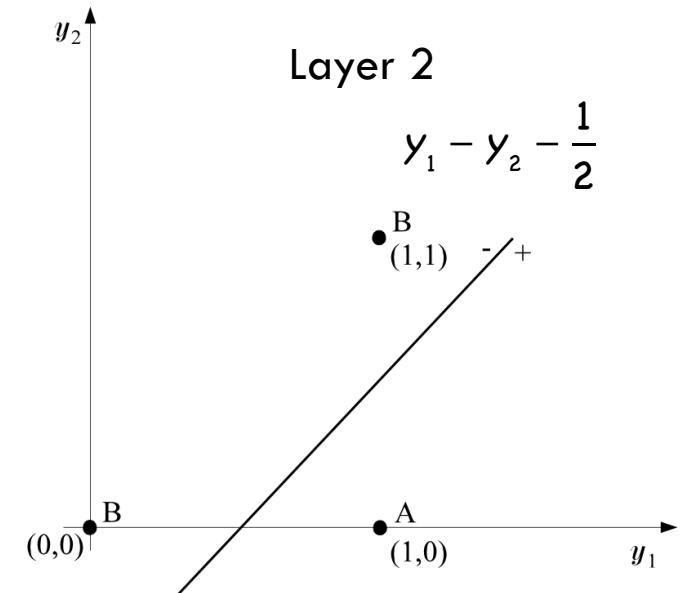
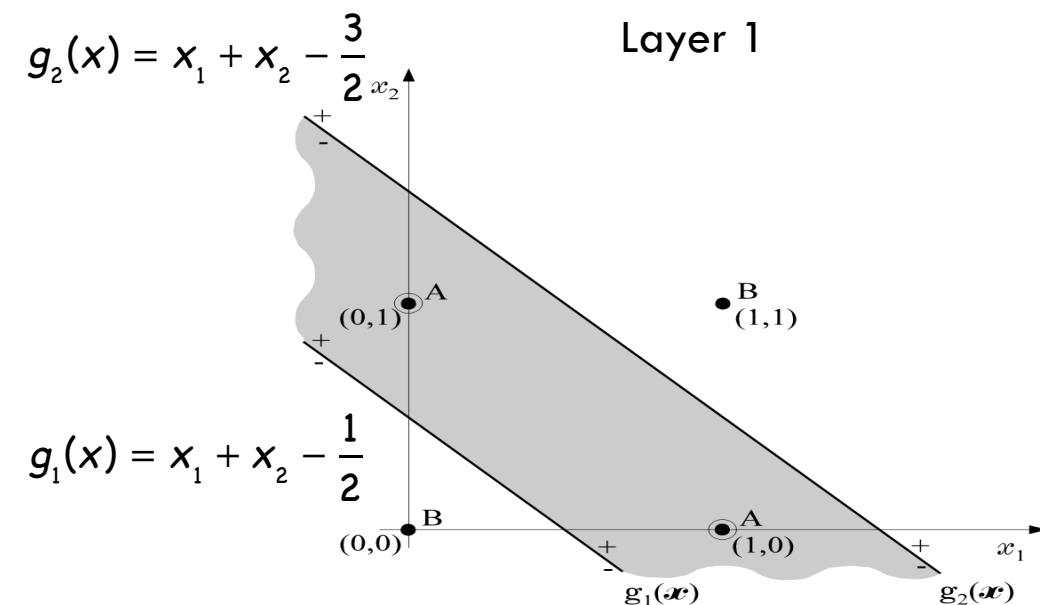
# The Two-Layer Perceptron

12

Multilayer Perceptrons

- Note that the hidden layer maps the plane onto the vertices of a unit square.

$$y_1 = f(g_1(x)) \text{ and } y_2 = f(g_2(x))$$

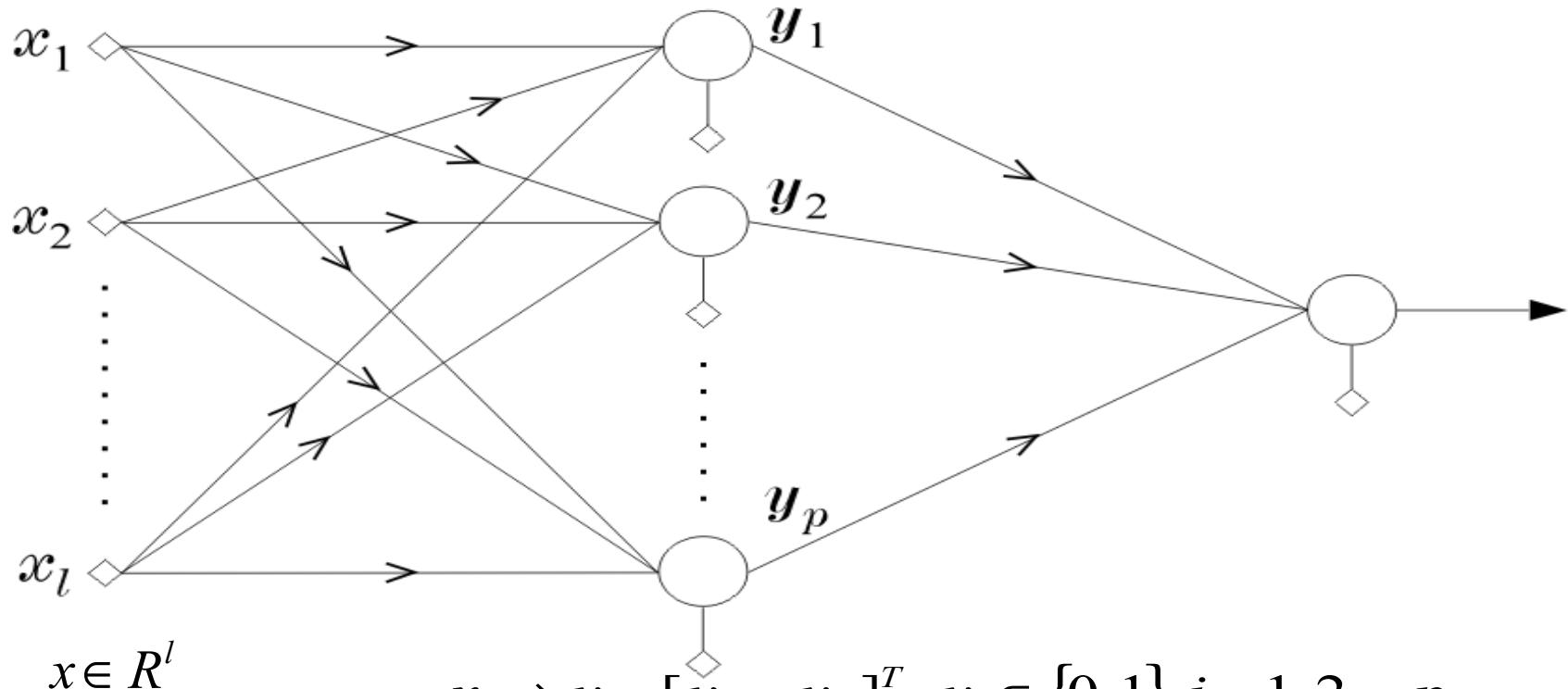


# Higher Dimensions

13

Multilayer Perceptrons

- Each hidden unit realizes a hyperplane discriminant function.
- The output of each hidden unit is 0 or 1 depending upon the location of the input vector relative to the hyperplane.

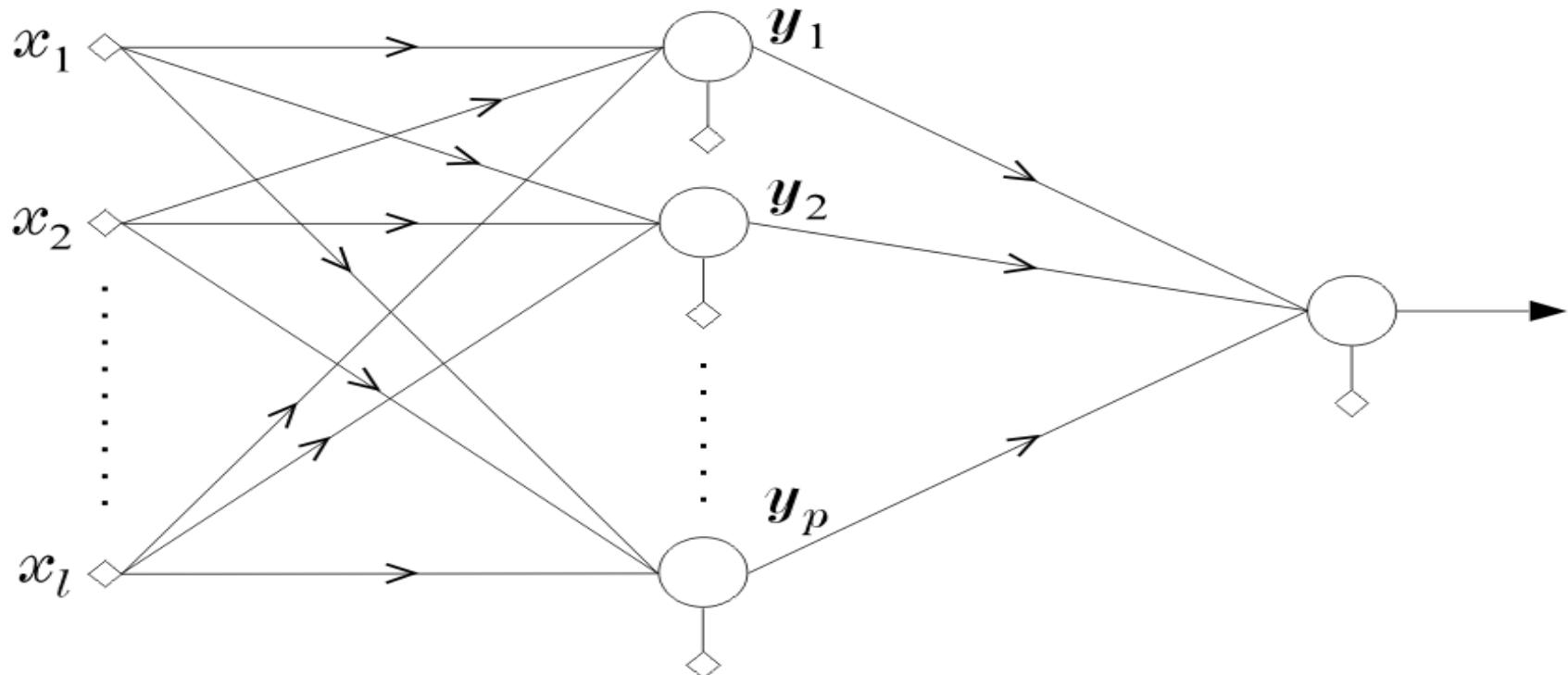


# Higher Dimensions

14

Multilayer Perceptrons

- Together, the hidden units map the input onto the vertices of a  $p$ -dimensional unit hypercube.



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

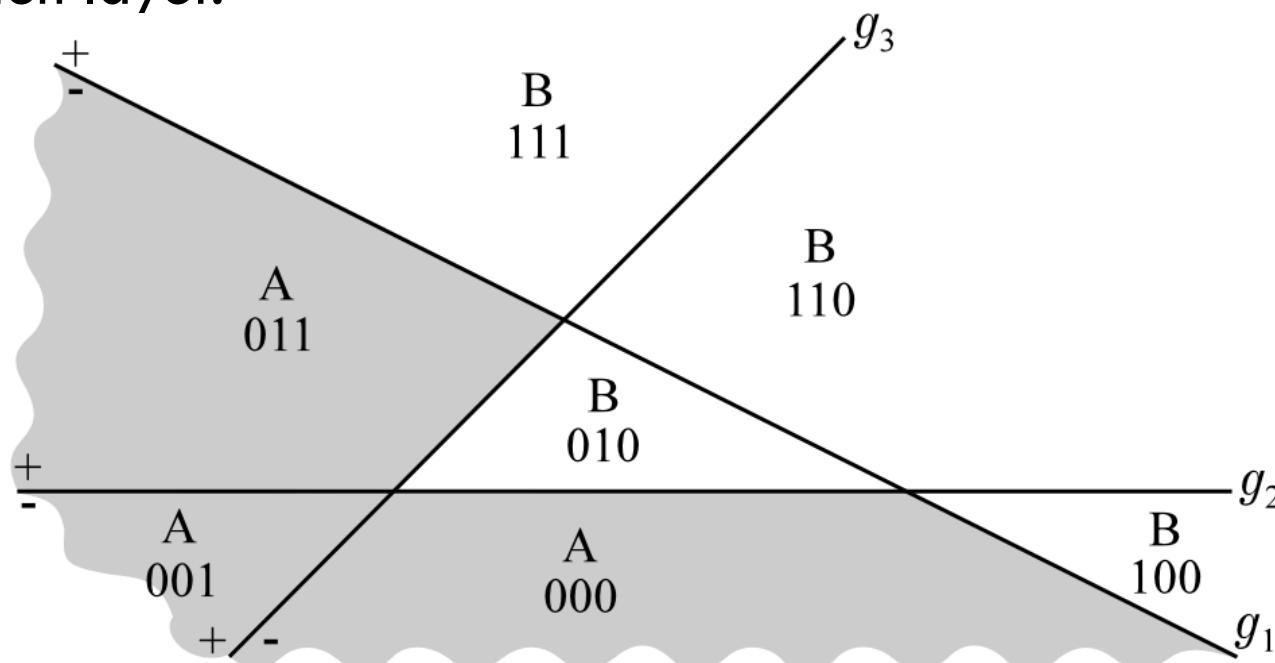
$$\underline{x} \rightarrow \underline{y} = [y_1, \dots, y_p]^T, y_i \in \{0, 1\} \quad i = 1, 2, \dots, p$$

# Two-Layer Perceptron

15

Multilayer Perceptrons

- These  $p$  hyperplanes partition the  $l$ -dimensional input space into polyhedral regions
- Each region corresponds to a different vertex of the  $p$ -dimensional hypercube represented by the outputs of the hidden layer.

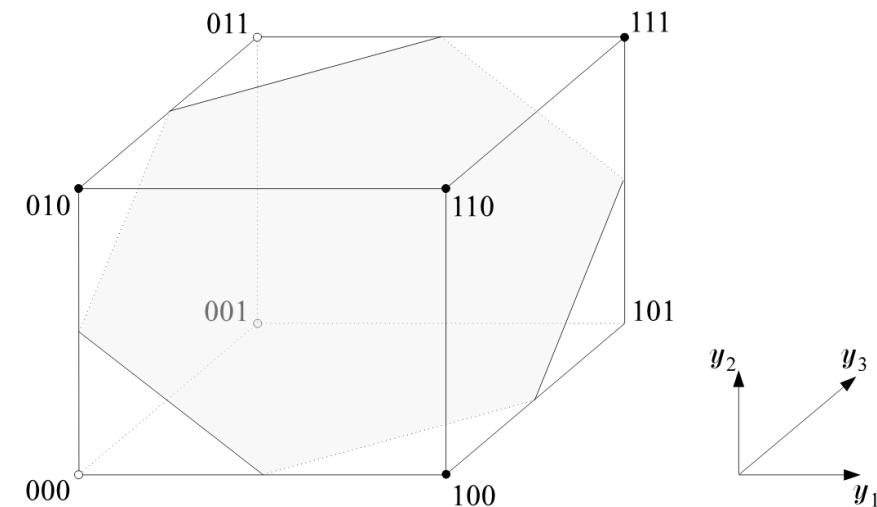
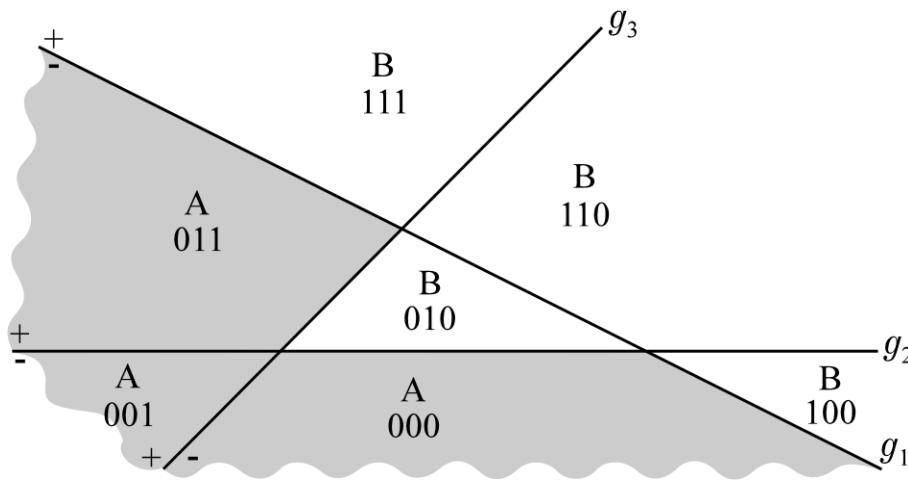


# Two-Layer Perceptron

16

Multilayer Perceptrons

- In this example, the vertex  $(0, 0, 1)$  corresponds to the region of the input space where:
  - $g_1(x) < 0$
  - $g_2(x) < 0$
  - $g_3(x) > 0$



# Limitations of a Two-Layer Perceptron

17

Multilayer Perceptrons

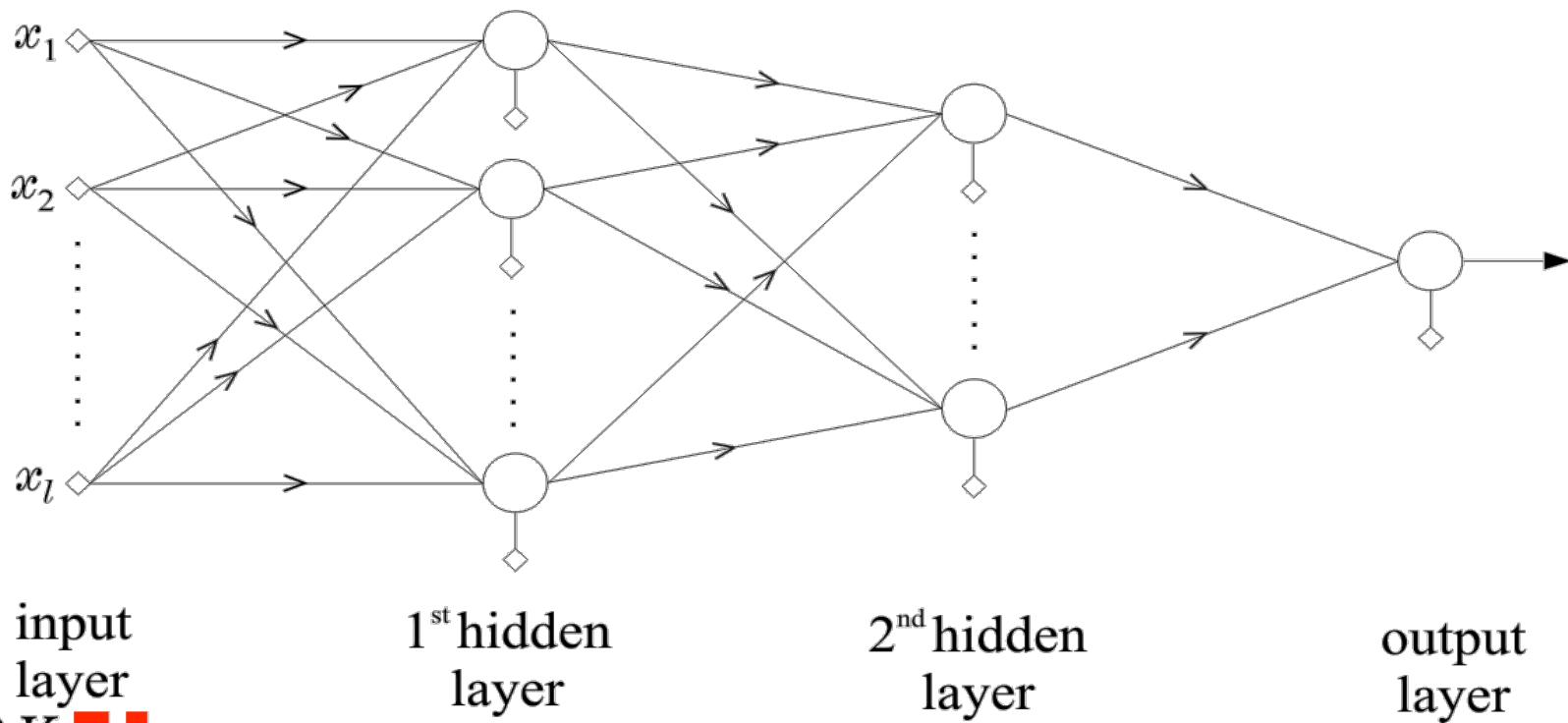
- The output neuron realizes a hyperplane in the transformed space that partitions the p vertices into two sets.
- Thus, the two layer perceptron has the capability to classify vectors into **classes that consist of unions of polyhedral regions.**
- But **NOT ANY** union. It depends on the relative position of the corresponding vertices.
- How can we solve this problem?

# The Three-Layer Perceptron

18

Multilayer Perceptrons

- Suppose that Class A consists of the union of K polyhedra in the input space.
- Use K neurons in the 2<sup>nd</sup> hidden layer.
- Train each to classify one Class A vertex as positive, the rest negative.
- Now use an output neuron that implements the OR function.

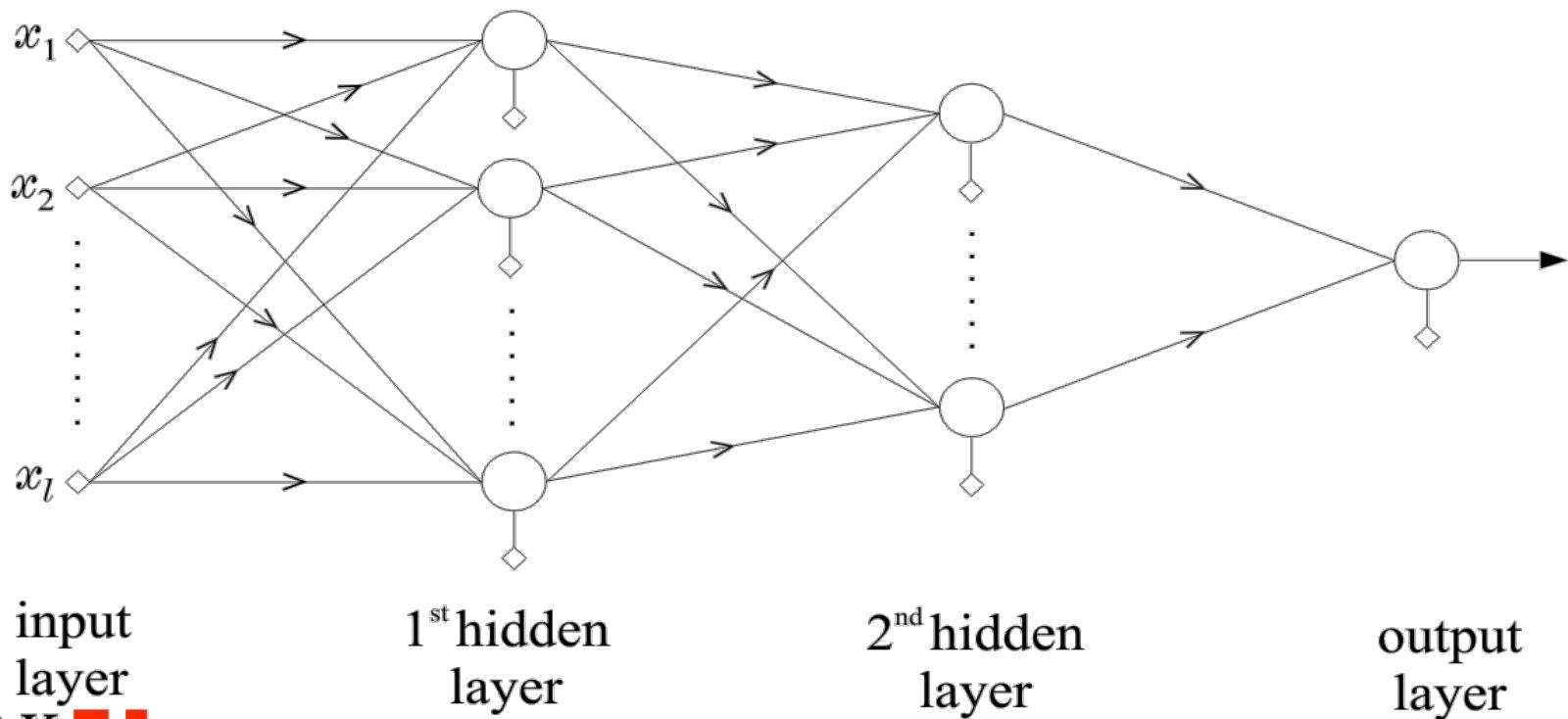


# The Three-Layer Perceptron

19

Multilayer Perceptrons

- Thus the three-layer perceptron can separate classes resulting from any union of polyhedral regions in the input space.

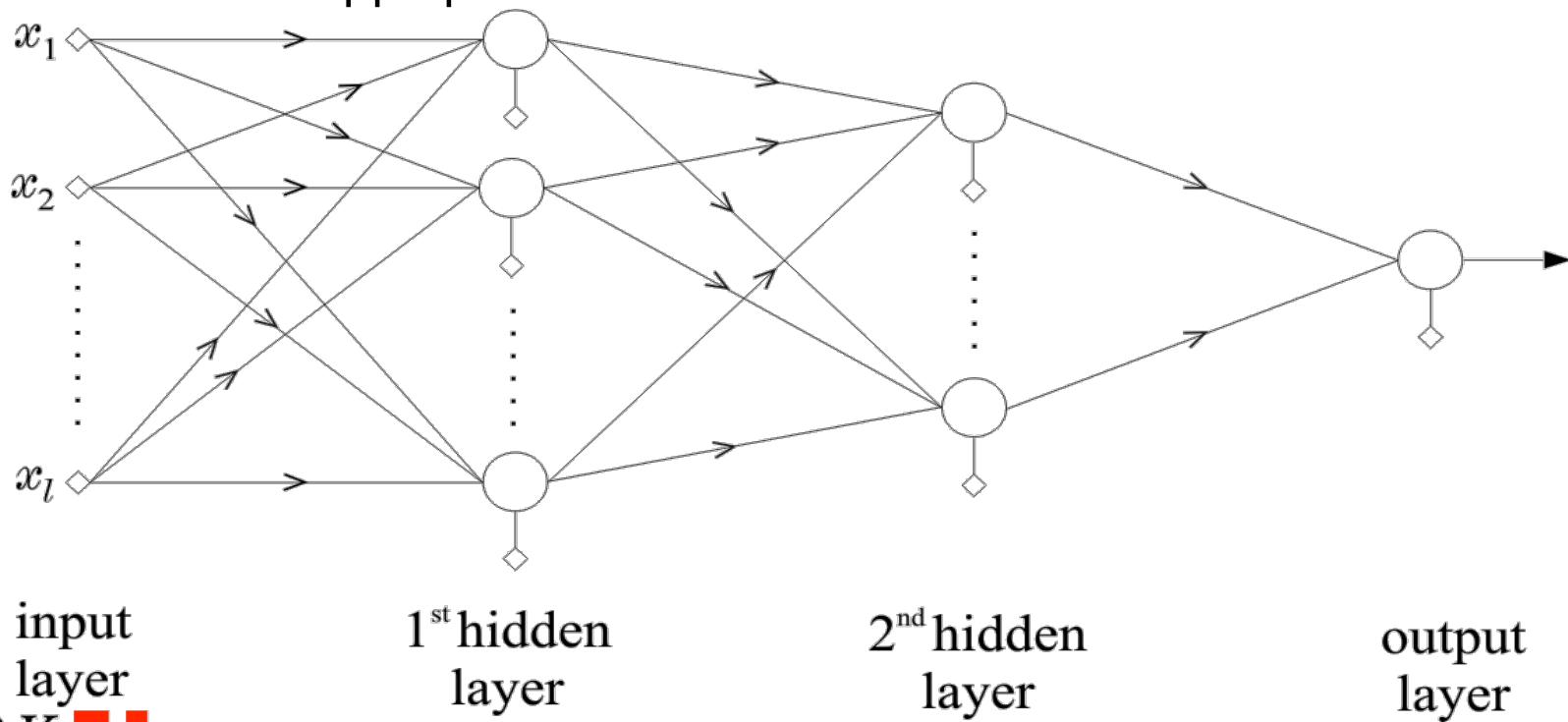


# The Three-Layer Perceptron

20

Multilayer Perceptrons

- The first layer of the network forms the **hyperplanes** in the input space.
- The second layer of the network forms the **polyhedral regions** of the input space
- The third layer forms the appropriate **unions of these regions** and maps each to the appropriate class.



# Outline

21

Multilayer Perceptrons

- Combining Linear Classifiers
- Learning Parameters

# Training Data

22

Multilayer Perceptrons

- The training data consist of  $N$  input-output pairs:

$$(\mathbf{y}(i), \mathbf{x}(i)), \quad i \in 1, \dots, N$$

where

$$\mathbf{y}(i) = \left[ y_1(i), \dots, y_{k_L}(i) \right]^t$$

and

$$\mathbf{x}(i) = \left[ x_1(i), \dots, x_{k_0}(i) \right]^t$$

# Choosing an Activation Function

23

Multilayer Perceptrons

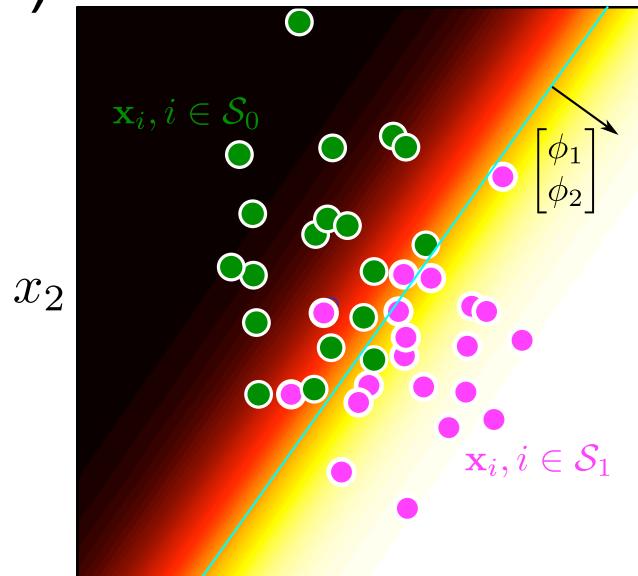
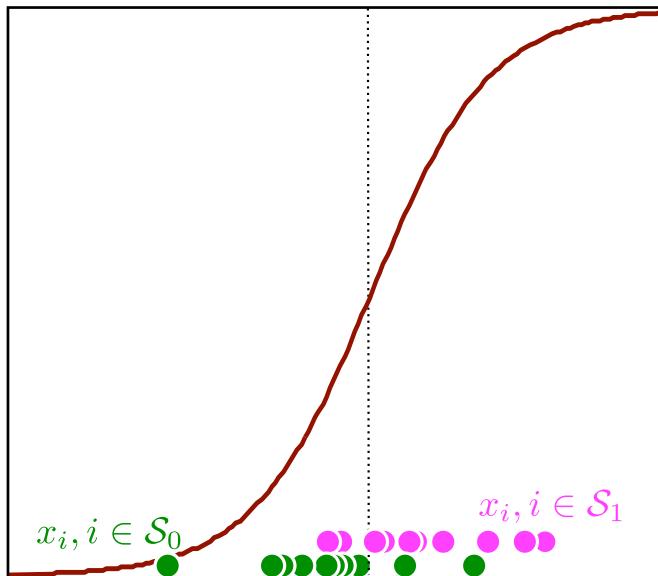
- The unit step activation function means that the error rate of the network is a discontinuous function of the weights.
- This makes it difficult to learn optimal weights by minimizing the error.
- To fix this problem, we need to use a smooth activation function.
- A popular choice is the sigmoid function we used for logistic regression:

# Smooth Activation Function

24

Multilayer Perceptrons

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



$$\mathbf{w}^t \phi$$

$$x_1$$

# Output: Two Classes

25

Multilayer Perceptrons

- For a binary classification problem, there is a single output node with activation function given by

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

- Since the output is constrained to lie between 0 and 1, it can be interpreted as the probability of the input vector belonging to Class 1.

# Output: $K > 2$ Classes

26

Multilayer Perceptrons

- For a  $K$ -class problem, we use  $K$  outputs, and the softmax function given by

$$y_k = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

- Since the outputs are constrained to lie between 0 and 1, and sum to 1,  $y_k$  can be interpreted as the probability that the input vector belongs to Class  $K$ .

# Non-Convex

27

Multilayer Perceptrons

- Now each layer of our multi-layer perceptron is a logistic regressor.
- Recall that optimizing the weights in logistic regression results in a convex optimization problem.
- Unfortunately the cascading of logistic regressors in the multi-layer perceptron makes the problem non-convex.
- This makes it difficult to determine an exact solution.
- Instead, we typically use **gradient descent** to find a locally optimal solution to the weights.
- The specific learning algorithm is called the **backpropagation** algorithm.

# Nonlinear Classification and Regression: Outline

28

## Multilayer Perceptrons

- Multi-Layer Perceptrons
  - **The Back-Propagation Learning Algorithm**
- Generalized Linear Models
  - Radial Basis Function Networks
  - Sparse Kernel Machines
    - Nonlinear SVMs and the Kernel Trick
    - Relevance Vector Machines

# The Backpropagation Algorithm

Paul J. Werbos. Beyond Regression: New Tools for Prediction and Analysis in the Behavioral Sciences. PhD thesis, Harvard University, 1974

Rumelhart, David E.; Hinton, Geoffrey E., Williams, Ronald J. (8 October 1986). "Learning representations by back-propagating errors". *Nature* **323** (6088): 533–536.



Werbos



Rumelhart



Hinton

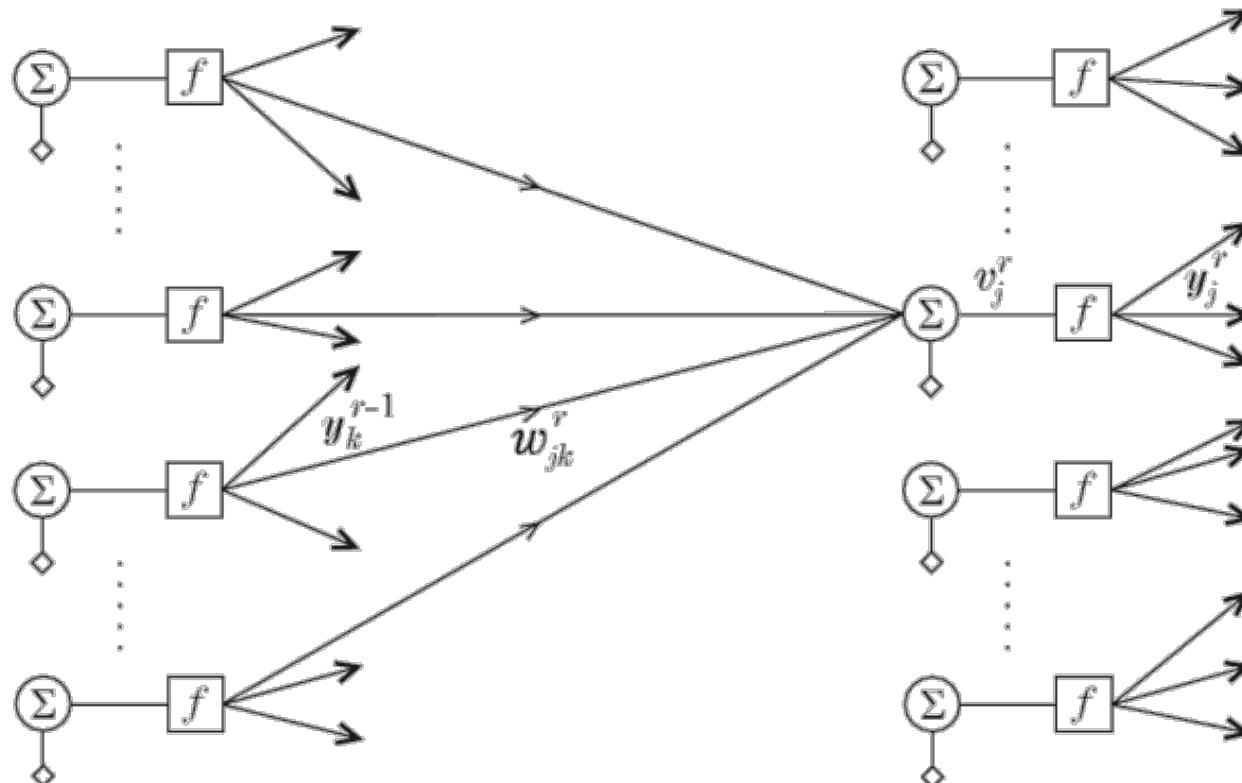
# Notation

30

## Multilayer Perceptrons

- Assume a network with  $L$  layers

- $k_0$  nodes in the input layer.
- $k_r$  nodes in the  $r$ 'th layer.



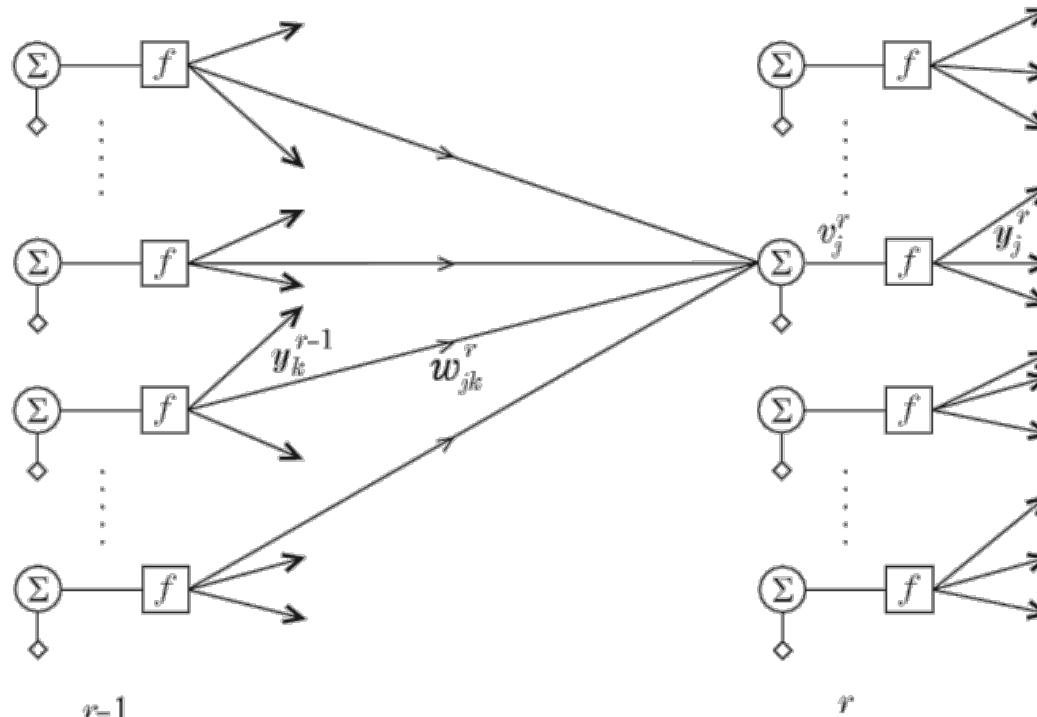
# Notation

31

Multilayer Perceptrons

Let  $y_k^{r-1}$  be the output of the  $k$ th neuron of layer  $r - 1$ .

Let  $w_{jk}^r$  be the weight of the synapse on the  $j$ th neuron of layer  $r$  from the  $k$ th neuron of layer  $r - 1$ .

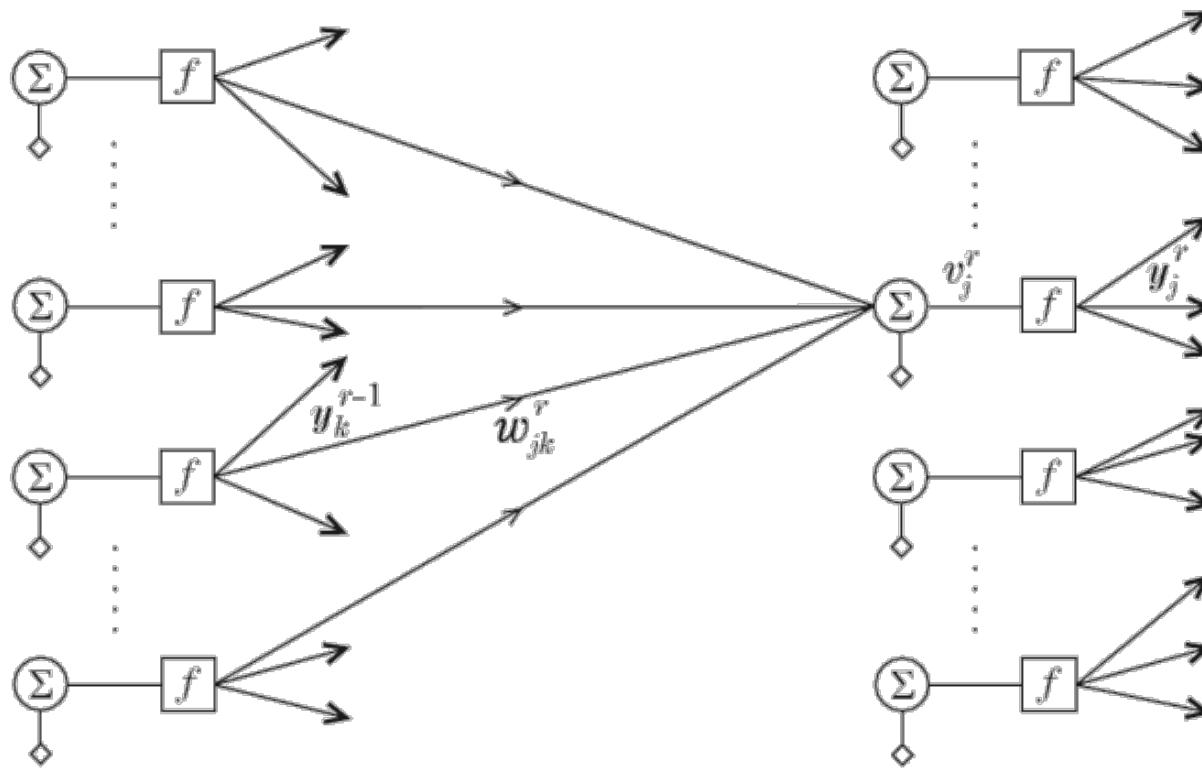


# Input

32

Multilayer Perceptrons

$$y_k^0(i) = x_k(i), \ k = 1, \dots, k_0$$

 $r-1$  $r$

# Notation

33

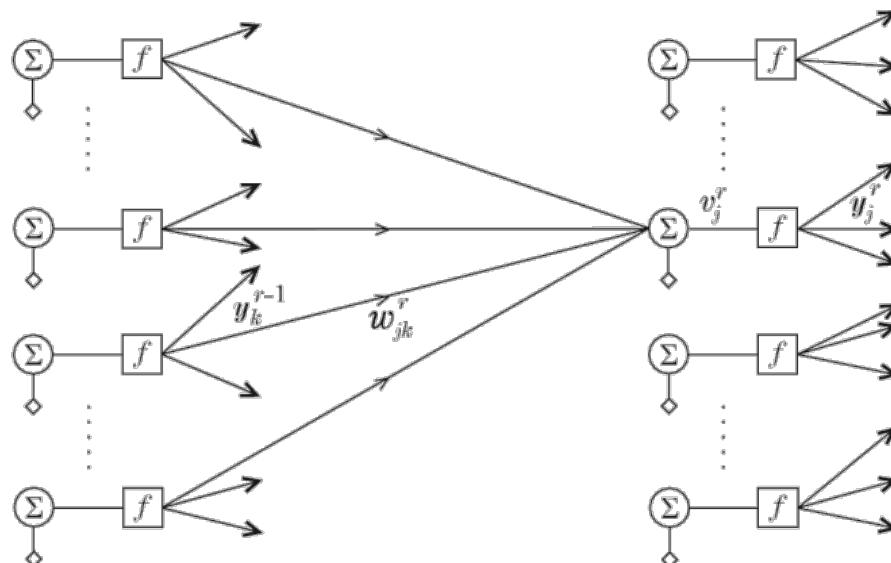
Multilayer Perceptrons

Let  $v_j^r$  be the total input to the  $j$ th neuron of layer  $r$ :

$$v_j^r(i) = (\mathbf{w}_j^r)^t \mathbf{y}^{r-1}(i) = \sum_{k=0}^{k_{r-1}} w_{jk}^r y_k^{r-1}(i)$$

where we define  $y_0^r(i) = +1$  to incorporate the bias term.

Then  $y_j^r(i) = f(v_j^r(i)) = f\left(\sum_{k=0}^{k_{r-1}} w_{jk}^r y_k^{r-1}(i)\right)$



# Cost Function

34

Multilayer Perceptrons

- A common cost function is the squared error:

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

where  $\varepsilon(i) \triangleq \frac{1}{2} \sum_{m=1}^{k_L} (e_m(i))^2 = \frac{1}{2} \sum_{m=1}^{k_L} (y_m(i) - \hat{y}_m(i))^2$

and

$\hat{y}_m(i) = y_k^r(i)$  is the output of the network.

# Cost Function

35

Multilayer Perceptrons

- To summarize, the error for input  $i$  is given by

$$\varepsilon(i) = \frac{1}{2} \sum_{m=1}^{k_L} (e_m(i))^2 = \frac{1}{2} \sum_{m=1}^{k_L} (\hat{y}_m(i) - y_m(i))^2$$

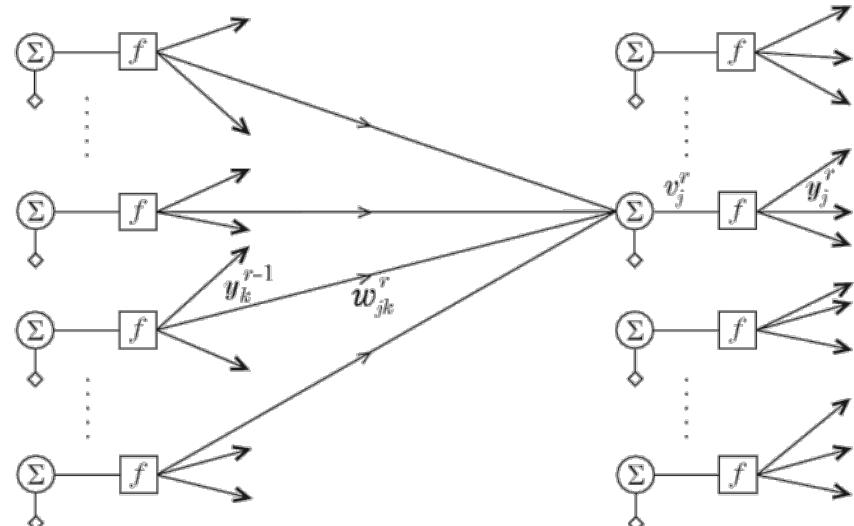
where  $\hat{y}_m(i) = y_k^r(i)$  is the output of the output layer

and each layer is related to the previous layer through

$$y_j^r(i) = f(v_j^r(i))$$

and

$$v_j^r(i) = (\mathbf{w}_j^r)^t \mathbf{y}^{r-1}(i)$$



# Gradient Descent

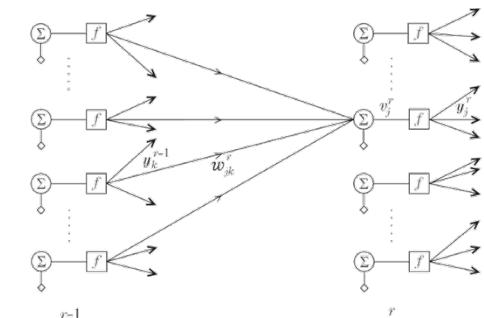
36

Multilayer Perceptrons

$$\varepsilon(i) = \frac{1}{2} \sum_{m=1}^{k_L} (e_m(i))^2 = \frac{1}{2} \sum_{m=1}^{k_L} (\hat{y}_m(i) - y_m(i))^2$$

- Gradient descent starts with an initial guess at the weights over all layers of the network.
- We then use these weights to compute the network output  $\hat{y}(i)$  for each input vector  $x(i)$  in the training data.
- This allows us to calculate the error  $\varepsilon(i)$  for each of these inputs.
- Then, in order to minimize this error, we incrementally update the weights in the negative gradient direction:

$$w_j^r(\text{new}) = w_j^r(\text{old}) - \mu \frac{\partial J}{\partial w_j^r} = w_j^r(\text{old}) - \mu \sum_{i=1}^N \frac{\partial \varepsilon(i)}{\partial w_j^r}$$



# Gradient Descent

37

Multilayer Perceptrons

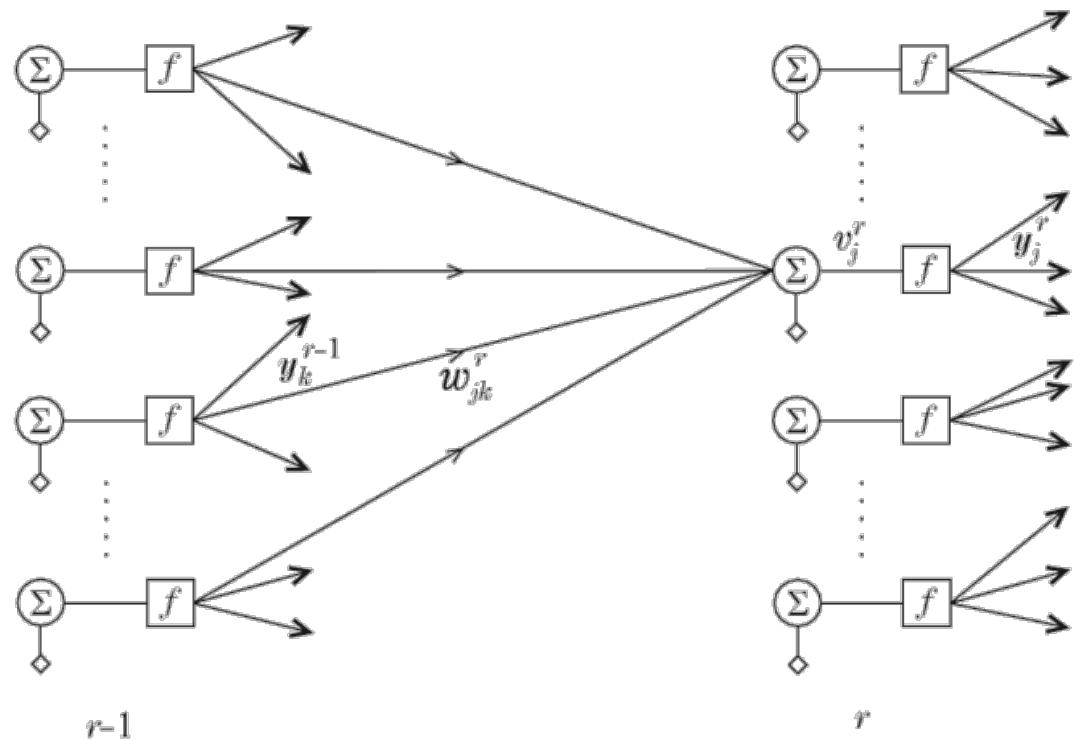
□ Since  $v_j^r(i) = (\mathbf{w}_j^r)^t \mathbf{y}^{r-1}(i)$ ,

the influence of the  $j$ th weight of the  $r$ th layer on the error can be expressed as:

$$\begin{aligned}\frac{\partial \varepsilon(i)}{\partial \mathbf{w}_j^r} &= \frac{\partial \varepsilon(i)}{\partial v_j^r(i)} \frac{\partial v_j^r(i)}{\partial \mathbf{w}_j^r} \\ &= \delta_j^r(i) \mathbf{y}^{r-1}(i)\end{aligned}$$

where

$$\delta_j^r(i) \triangleq \frac{\partial \varepsilon(i)}{\partial v_j^r(i)}$$



$r-1$

$r$

# Gradient Descent

38

## Multilayer Perceptrons

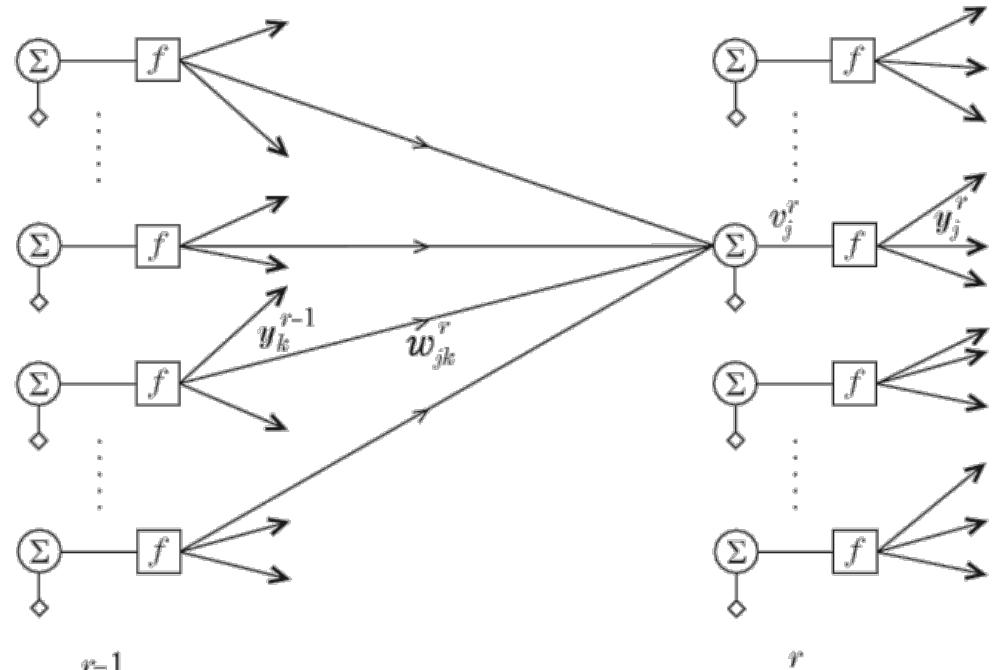
$$\frac{\partial \varepsilon(i)}{\partial w_j^r} = \delta_j^r(i) y^{r-1}(i),$$

where

$$\delta_j^r(i) \triangleq \frac{\partial \varepsilon(i)}{\partial v_j^r(i)}$$

For an intermediate layer  $r$ ,  
we cannot compute  $\delta_j^r(i)$  directly.

However,  $\delta_j^r(i)$  can be computed inductively,  
starting from the output layer.



# Backpropagation: The Output Layer

39

Multilayer Perceptrons

$$\frac{\partial \varepsilon(i)}{\partial \mathbf{w}_j^r} = \delta_j^r(i) \mathbf{y}^{r-1}(i), \text{ where } \delta_j^r(i) \triangleq \frac{\partial \varepsilon(i)}{\partial v_j^r(i)}$$

$$\text{and } \varepsilon(i) = \frac{1}{2} \sum_{m=1}^{k_L} (e_m(i))^2 = \frac{1}{2} \sum_{m=1}^{k_L} (\hat{y}_m(i) - y_m(i))^2$$

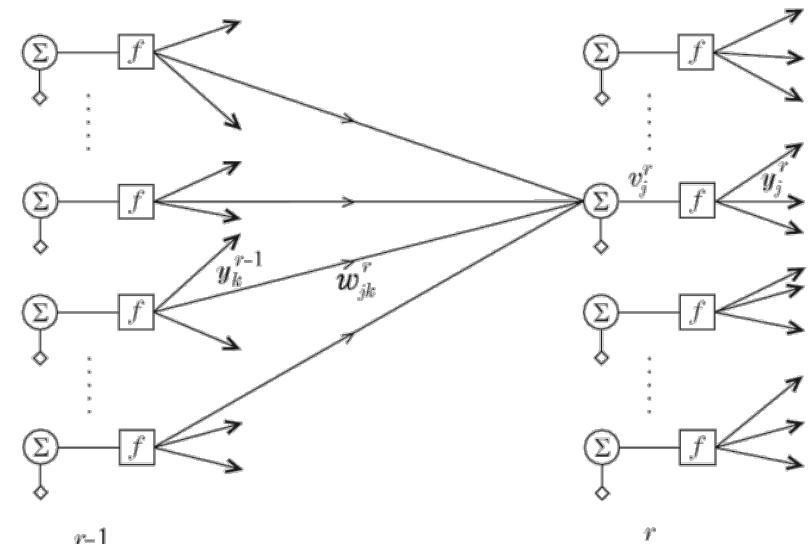
Recall that  $\hat{y}_m(i) = y_j^L(i) = f(v_j^L(i))$

Thus at the output layer we have

$$\delta_j^L(i) = \frac{\partial \varepsilon(i)}{\partial v_j^L(i)} = \frac{\partial \varepsilon(i)}{\partial e_j^L(i)} \frac{\partial e_j^L(i)}{\partial v_j^L(i)} = e_j^L(i) f'(v_j^L(i))$$

$$f(a) = \frac{1}{1 + \exp(-a)} \rightarrow f'(a) = f(a)(1 - f(a))$$

$$\boxed{\delta_j^L(i) = e_j^L(i) f(v_j^L(i)) (1 - f(v_j^L(i)))}$$



# Backpropagation: Hidden Layers

40

Multilayer Perceptrons

- Observe that the dependence of the error on the total input to a neuron in a previous layer can be expressed in terms of the dependence on the total input of neurons in the following layer:

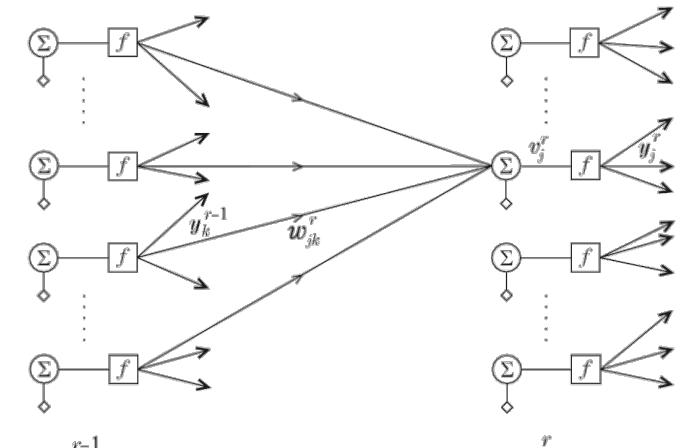
$$\delta_j^{r-1}(i) = \frac{\partial \varepsilon(i)}{\partial v_j^{r-1}(i)} = \sum_{k=1}^k \frac{\partial \varepsilon(i)}{\partial v_k^r(i)} \frac{\partial v_k^r(i)}{\partial v_j^{r-1}(i)} = \sum_{k=1}^k \delta_k^r(i) \frac{\partial v_k^r(i)}{\partial v_j^{r-1}(i)}$$

where  $v_k^r(i) = \sum_{m=0}^{k-1} w_{km}^r y_m^r(i) = \sum_{m=0}^{k-1} w_{km}^r f(v_m^r(i))$

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

and so  $\delta_j^{r-1}(i) = \frac{\partial \varepsilon(i)}{\partial v_j^{r-1}(i)} = f'(v_j^{r-1}(i)) \sum_{k=1}^k \delta_k^r(i) w_{kj}^r = f(v_j^L(i)) (1 - f(v_j^L(i))) \sum_{k=1}^k \delta_k^r(i) w_{kj}^r$

Thus once the  $\delta_k^r(i)$  are determined they can be propagated backward to calculate  $\delta_j^{r-1}(i)$  using this inductive formula.



# Backpropagation: Summary of Algorithm

41

Multilayer Perceptrons

Repeat until convergence

## 1. Initialization

- Initialize all weights with small random values

## 2. Forward Pass

- For each input vector, run the network in the forward direction, calculating:

$$v_j^r(i) = (\mathbf{w}_j^r)^T \mathbf{y}^{r-1}(i); \quad \mathbf{y}_j^r(i) = f(v_j^r(i))$$

$$\text{and finally } \varepsilon(i) = \frac{1}{2} \sum_{m=1}^{k_L} (e_m(i))^2 = \frac{1}{2} \sum_{m=1}^{k_L} (\hat{y}_m(i) - y_m(i))^2$$

## 3. Backward Pass

- Starting with the output layer, use our inductive formula to compute the  $\delta_j^{r-1}(i)$ :

- Output Layer (Base Case):  $\delta_j^L(i) = e_j^L(i) f'(v_j^L(i))$

- Hidden Layers (Inductive Case):  $\delta_j^{r-1}(i) = f'(v_j^{r-1}(i)) \sum_{k=1}^{k_r} \delta_k^r(i) w_{kj}^r$

## 4. Update Weights

$$\mathbf{w}_j^r(\text{new}) = \mathbf{w}_j^r(\text{old}) - \mu \sum_{i=1}^N \frac{\partial \varepsilon(i)}{\partial \mathbf{w}_j^r} \quad \text{where } \frac{\partial \varepsilon(i)}{\partial \mathbf{w}_j^r} = \delta_j^r(i) \mathbf{y}^{r-1}(i)$$

# Batch vs Online Learning

42

Multilayer Perceptrons

- As described, on each iteration backprop updates the weights based upon all of the training data.  
This is called **batch learning**.

$$\mathbf{w}_j^r(\text{new}) = \mathbf{w}_j^r(\text{old}) - \mu \sum_{i=1}^N \frac{\partial \varepsilon(i)}{\partial \mathbf{w}_j^r} \quad \text{where } \frac{\partial \varepsilon(i)}{\partial \mathbf{w}_j^r} = \delta_j^r(i) \mathbf{y}^{r-1}(i)$$

- An alternative is to update the weights after each training input has been processed by the network, based only upon the error for that input. This is called **online learning**.

$$\mathbf{w}_j^r(\text{new}) = \mathbf{w}_j^r(\text{old}) - \mu \frac{\partial \varepsilon(i)}{\partial \mathbf{w}_j^r} \quad \text{where } \frac{\partial \varepsilon(i)}{\partial \mathbf{w}_j^r} = \delta_j^r(i) \mathbf{y}^{r-1}(i)$$

# Batch vs Online Learning

43

Multilayer Perceptrons

- One advantage of batch learning is that averaging over all inputs when updating the weights should lead to smoother convergence.
- On the other hand, the randomness associated with online learning might help to prevent convergence toward a local minimum.
- Changing the order of presentation of the inputs from epoch to epoch may also improve results.

# Remarks

44

Multilayer Perceptrons

- Local Minima
  - The objective function is in general non-convex, and so the solution may not be globally optimal.
- Stopping Criterion
  - Typically stop when the change in weights or the change in the error function falls below a threshold.
- Learning Rate
  - The speed and reliability of convergence depends on the learning rate  $\mu$ .