

# Introduction to Deep Learning

## Lecture 2

### **Single Layer Perceptrons**

Slides by:

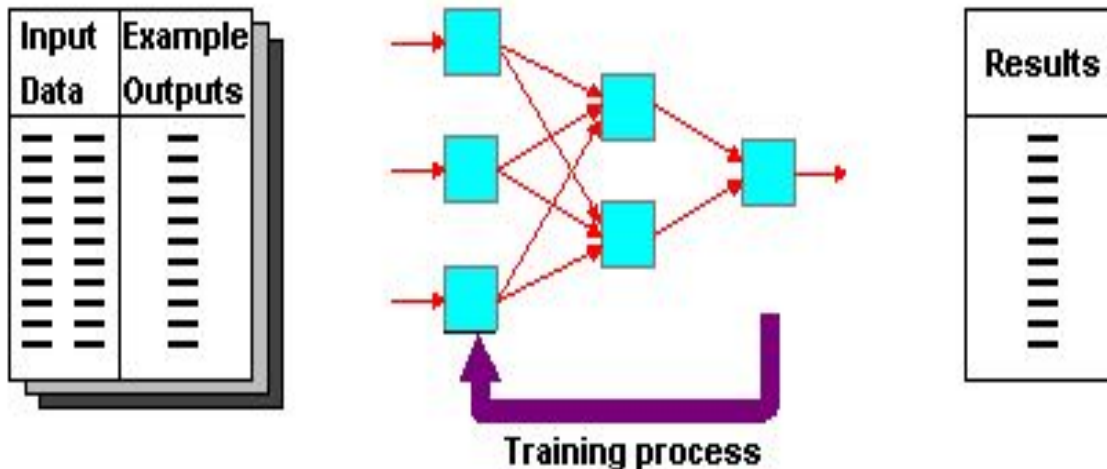
Dr. Wojtek Kowalczyk

Presenter:

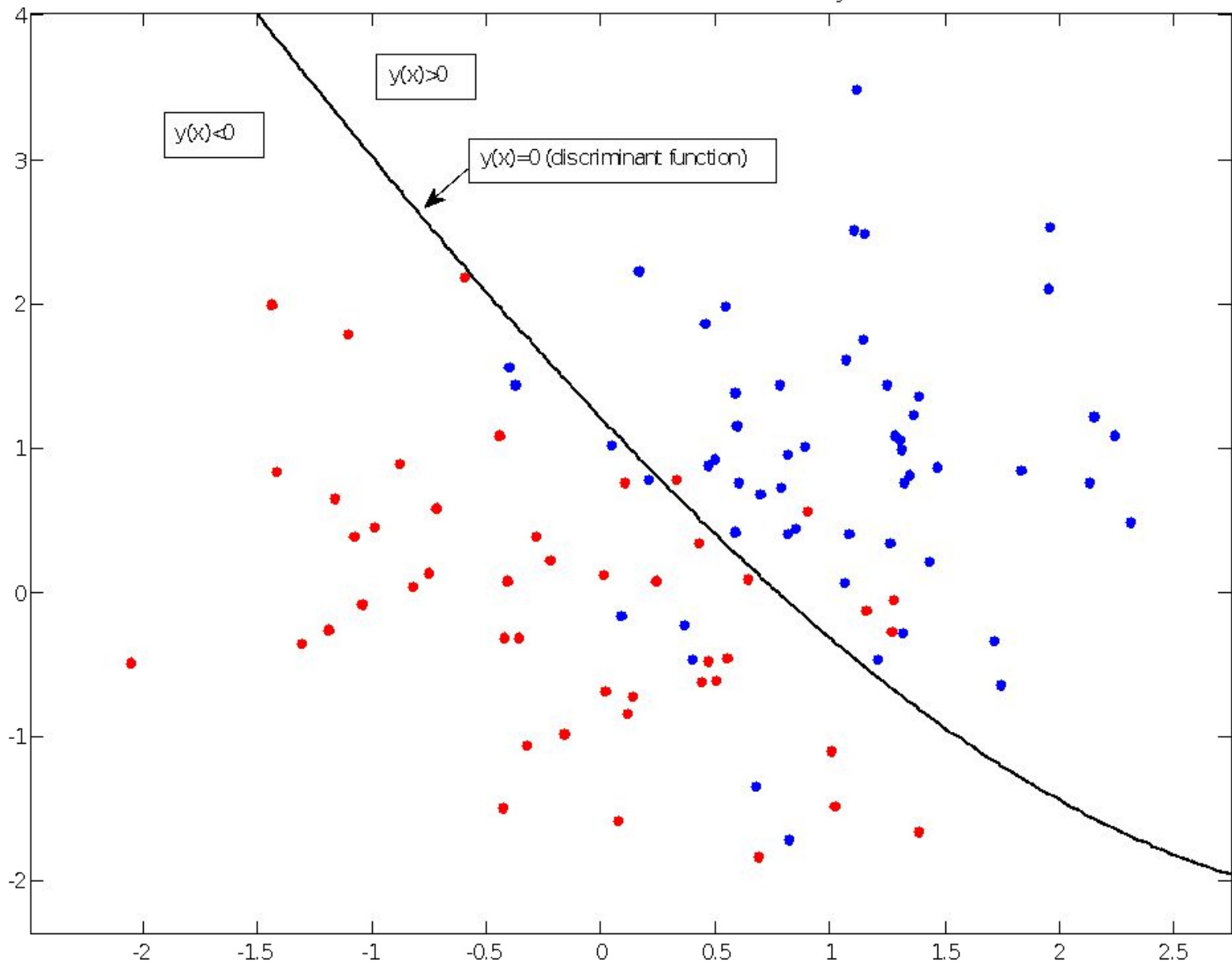
Msc. Andrius Bernatavicius

# Single Layer Perceptron

- Training and test data sets
- Training set: input & target
- Test set: only input



A Classification Problem and a Class Boundary



Neural Networks

Linear Models

# Single Layer Networks (Linear Models)

---

- Linear Discriminants
- Single Layer Perceptron
- Linear Separability and Cover's Theorem
- Learning Algorithms:
  - Perceptron Rule and Convergence Theorem
  - Least-Squares Method (Adaline)
  - Gradient Descent and Logistic Regression
- Generalized Linear Discriminants
- Multi-class perceptron learning algorithm

# Discriminant Functions

- A *discriminant function for classes  $C_1$  and  $C_2$*  is any function  $y(\mathbf{x})$  such that an input vector  $\mathbf{x}$  is assigned to class  $C_1$  if  $y(\mathbf{x}) > 0$  (and to  $C_2$  if  $y(\mathbf{x}) < 0$ ) [what to do with ties?]
- In case of  $N$  classes, we need  $N$  discriminant functions  $y_1(\mathbf{x}), y_2(\mathbf{x}), \dots, y_N(\mathbf{x})$ , such the  $k$ -th function  $y_k(\mathbf{x})$  discriminates class  $C_k$  from all other classes, for  $k=1, \dots, N$ , i.e.,

$\mathbf{x}$  is assigned to class  $C_k$  iff  $y_k(\mathbf{x}) > y_n(\mathbf{x})$  for all  $k \neq n$

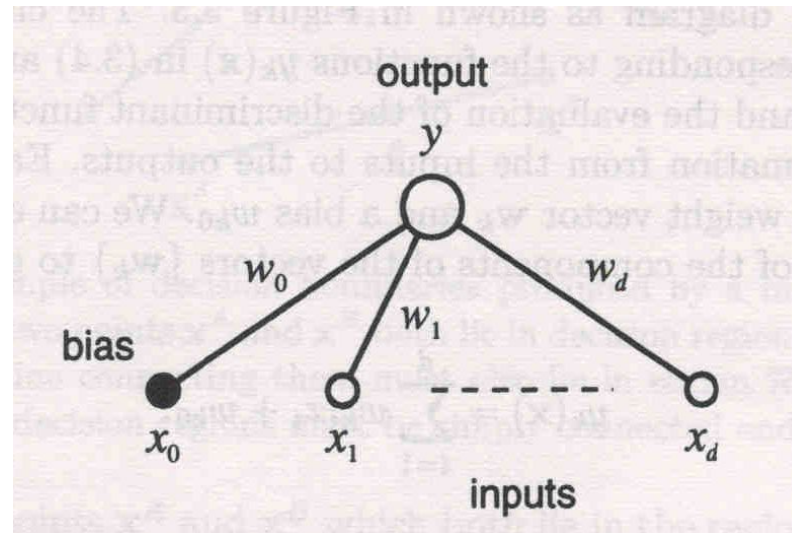
- In case of two classes ( $C_1$  and  $C_2$ ) we often use  $y(\mathbf{x})=y_1(\mathbf{x})-y_2(\mathbf{x})$  as a discriminant; then the sign of  $y(\mathbf{x})$  decides on the class:  
if  $y(\mathbf{x}) > 0$  then  $\mathbf{x}$  in  $C_1$  else  $\mathbf{x}$  in  $C_2$

# Linear Discriminant Functions

A *linear discriminant function in n-dimensional space* is a function  $y(\mathbf{x})$  in the form:

$$y(\mathbf{x}) = w_0 + w_1x_1 + \dots + w_nx_n = w_0 + \mathbf{w}^T\mathbf{x}$$

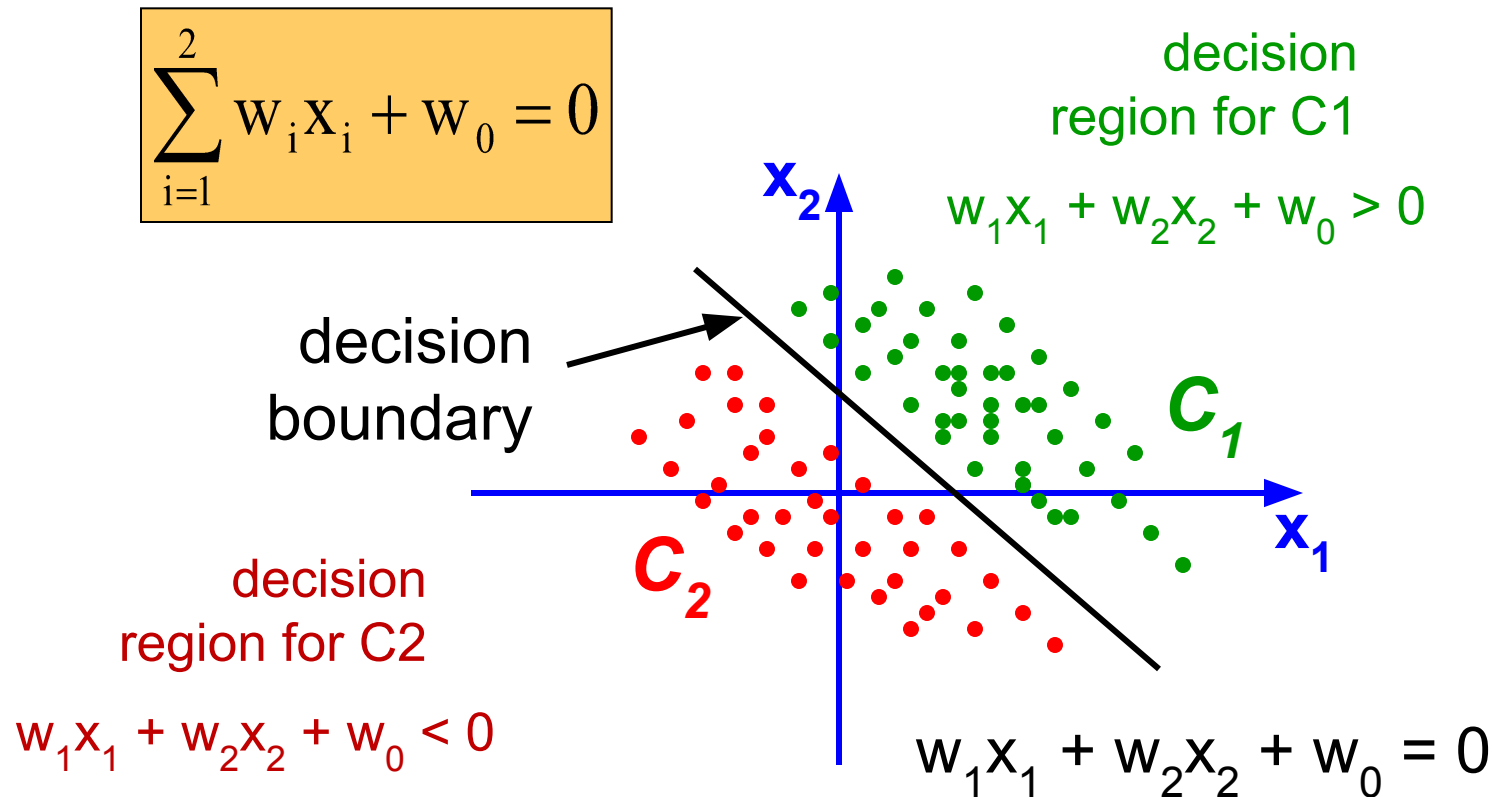
Neuron:



Linear discriminant function defines a **point** (in  $\mathbb{R}^1$ ), a **line** (in  $\mathbb{R}^2$ ), a **plane** (in  $\mathbb{R}^3$ ), a **hyperplane** (in  $\mathbb{R}^n$ ), that splits the space into “positive” and “negative” half-spaces.

# Geometric View in 2D

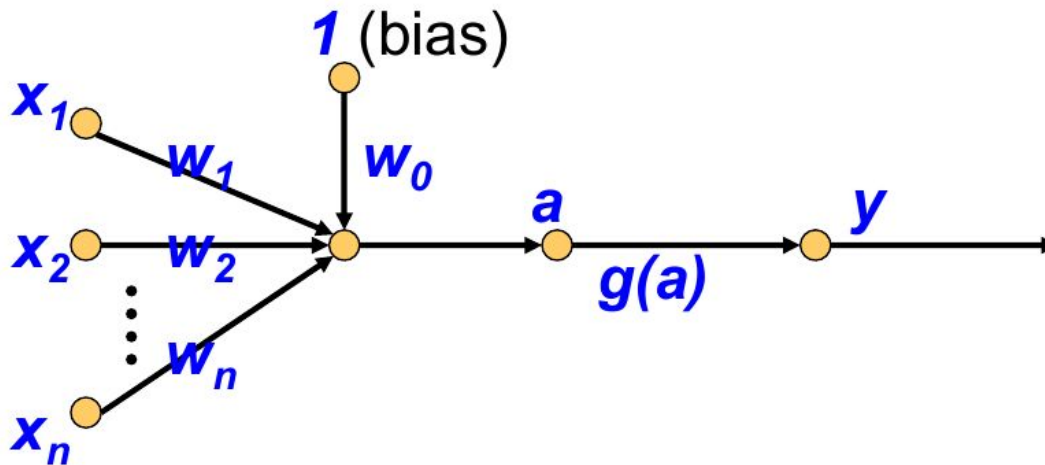
The equation below describes a (hyper-)plane in the input space consisting of real valued 2D vectors. The plane splits the input space into two regions, each of them describing one class.



# Perceptron: McCulloch-Pitts Model

The (McCulloch-Pitts) **perceptron** is a single node NN (or a single layer NN) with a non-linear function  **$g(\mathbf{a})$** :

$$a = w_0 + \sum w_i x_i; \quad g(a) = \begin{cases} +1 & \text{if } a \geq 0 \\ -1 & \text{if } a < 0 \end{cases}$$





# Perceptron Training

---

- How can we train a perceptron for a classification task?
- We try to find suitable values for the **weights** in such a way that the training examples are correctly classified.
- Geometrically, we try to find a **hyper-plane** that separates the examples of the two classes.
- Two classes  $C_1$  and  $C_2$  are *linearly separable* if there exists a hyperplane that separates them.

# Perceptron learning algorithm

initialize  $\mathbf{w}$  randomly;

**while** (there are misclassified training examples)

    Select a misclassified example  $(\mathbf{x}, d)$

$$\mathbf{w}_{\text{new}} = \mathbf{w}_{\text{old}} + \eta d \mathbf{x};$$

**end-while**;

$\eta > 0$  is a learning rate parameter (step size);

## Motivation:

**If**  $\mathbf{x}$  misclassified and  $d=1 \Rightarrow \mathbf{w}\mathbf{x}$  should be bigger,

**If**  $\mathbf{x}$  misclassified and  $d=-1 \Rightarrow \mathbf{w}\mathbf{x}$  should be smaller:

$$(\mathbf{w} + \eta d \mathbf{x})\mathbf{x} = \mathbf{w}\mathbf{x} + \eta d \mathbf{x}^2$$

**This rule does exactly what we want ( $\mathbf{x}^2$  is  $>0$ ) !!!**

## Conventions:

- $\mathbf{w}$  is a row vector;
- $\mathbf{x}$  is a column vector
- $\mathbf{x}^2$  is  $\mathbf{x}^T * \mathbf{x}$

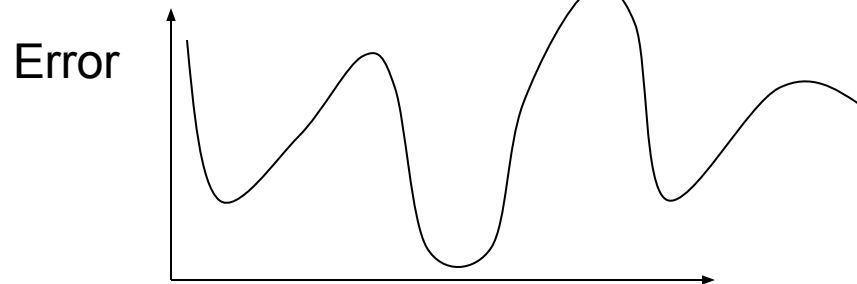
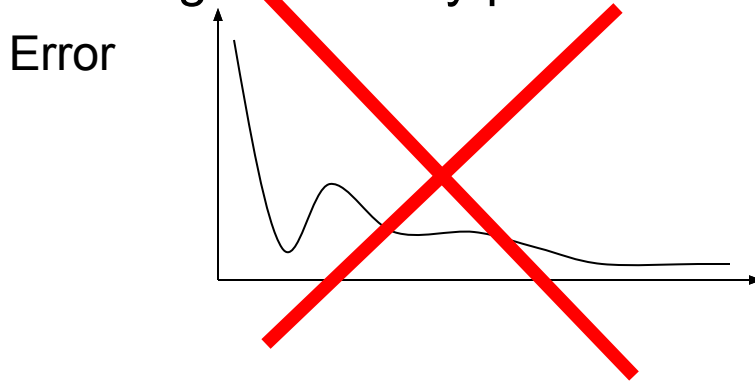
# Main properties

## 1) Perceptron Convergence Theorem:

If the classes  $C_1$ ,  $C_2$  are linearly separable (that is, there exists a hyper-plane that separates them) then the perceptron algorithm applied to  $C_1 \cup C_2$  terminates successfully after a finite number of iterations. [*the value of the learning rate  $\eta$  is not essential*]

## 2) Bad Behavior:

If the classes  $C_1$ ,  $C_2$  are **not** linearly separable then the perceptron algorithm may produce worse and worse results:



# Improvement: (Naive) Pocket Algorithm

0. Start with a random set of weights;  
put them in a 'pocket'
1. Select at random a pattern
2. If it is misclassified then
  - 2A. apply the perceptron update rule
  - 2B. check whether new weights are better than those kept in the pocket; if so put new weights to the pocket (and remove the old ones)*
3. goto 1.

Testing if new weights are better (less errors) is very expensive !!!

# Gallant's Pocket Algorithm

**main idea:** measure the quality of weights by counting the number of consecutive correct classifications ('current\_run')

0. Start with a random set of weights; put them in the 'pocket';

*best\_run* := 0;     *current\_run* := 0;

1. Select a pattern at random

2. If it is misclassified then

    apply the update rule; *current\_run* := 0;

else

*current\_run*++;

    if *current\_run* > *best\_run* then

        put new weights to the pocket; *best\_run* := *current\_run*

3. Goto 1

# Gallant's Pocket Algorithm with ratchet

0. Start with a random setting of weights; put it in the 'pocket'
1. *best\_run:=0; current\_run:=0;*
2. Select at random a pattern
3. **IF** it is correctly classified **THEN**  
*current\_run:=current\_run+1*  
  
**IF** *best\_run < current\_run* **AND** current weights  
are better than those kept in the pocket **THEN**  
*pocket:=current weights; best\_run=current\_run;*  
**ELSE**  
*current\_run=0;*  
update weights



# Pocket convergence theorem

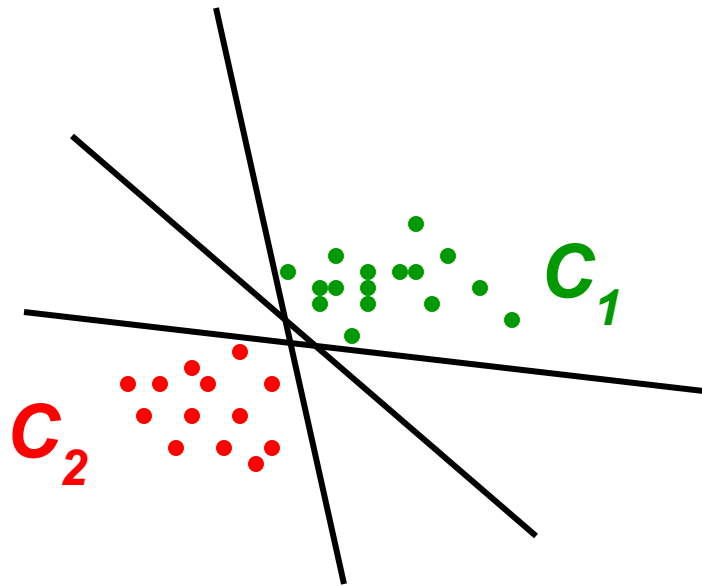
## **Pocket Convergence Theorem:**

The pocket algorithm converges with probability 1 to optimal weights (even if sets are not separable!)

## **Practice:**

- 1) The convergence rate is quite high
- 2) Pocket algorithm is better than the Perceptron algorithm
- 3) Both algorithms have very limited use
- 4) There may be many separating hyperplanes ...

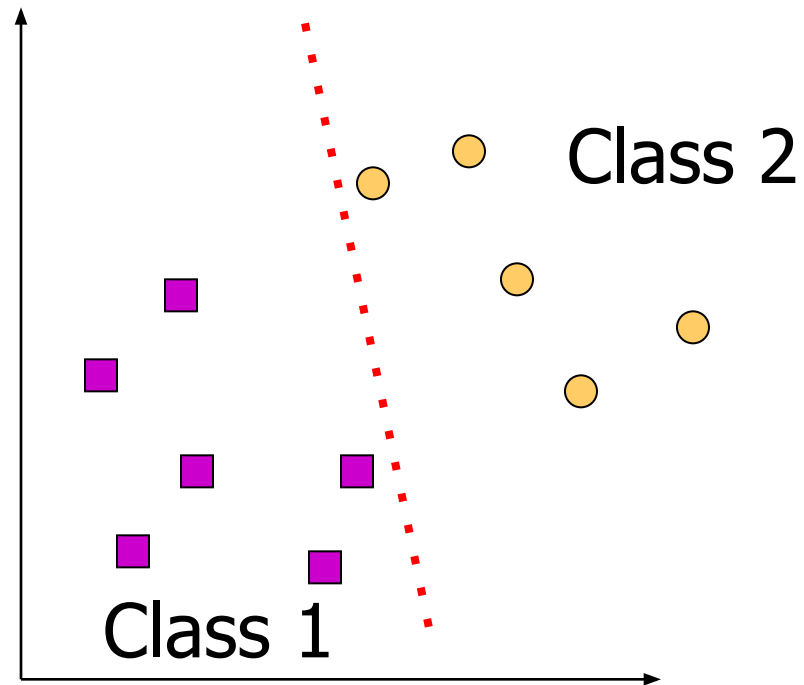
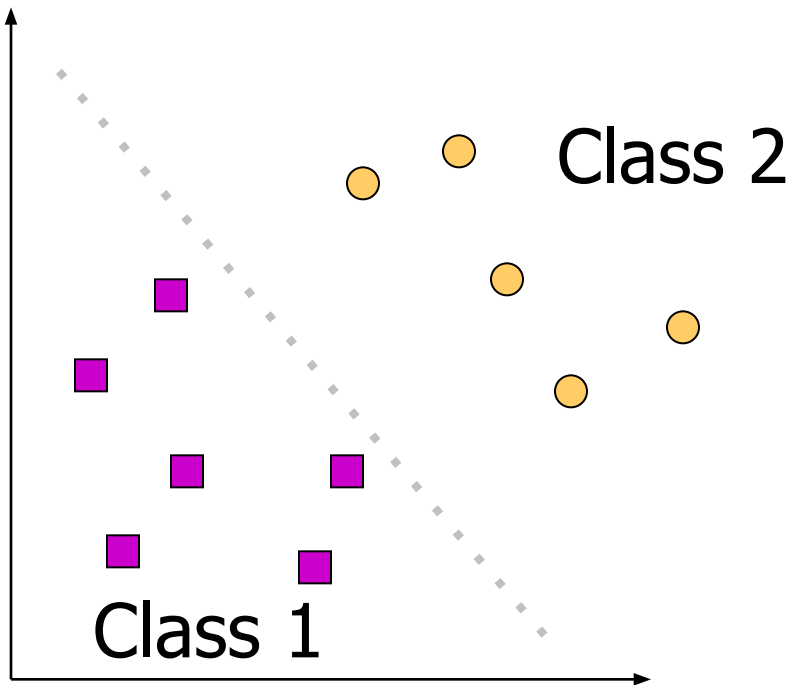
# Which separating hyperplane is best?



We are interested in accuracy on the test set!



# Examples of Bad Decision Boundaries



**Idea:**

define a “continuous error measure” and try to minimize it !

## Cover's Theorem (1965):

What is the chance that a randomly labeled set of  $N$  points (in general position) in  $d$ -dimensional space, is linearly separable?

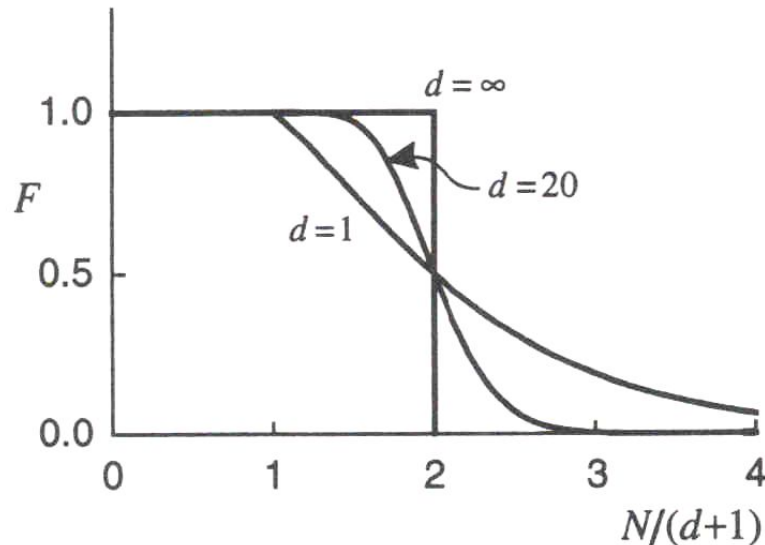


Figure 3.7. Plot of the fraction  $F(N, d)$  of the dichotomies of  $N$  data points in  $d$  dimensions which are linearly separable, as a function of  $N/(d + 1)$ , for various values of  $d$ .

$$F(N, d) = \begin{cases} 1 & \text{when } N \leq d + 1 \\ \frac{1}{2^{N-1}} \sum_{i=0}^d \binom{N-1}{i} & \text{when } N \geq d + 1 \end{cases} \quad (3.30)$$

# Cover's Theorem in highly dimensional spaces:

1) if the number of points in  $d$ -dimensional space is smaller than  $2^d$  then they are almost always linearly separable

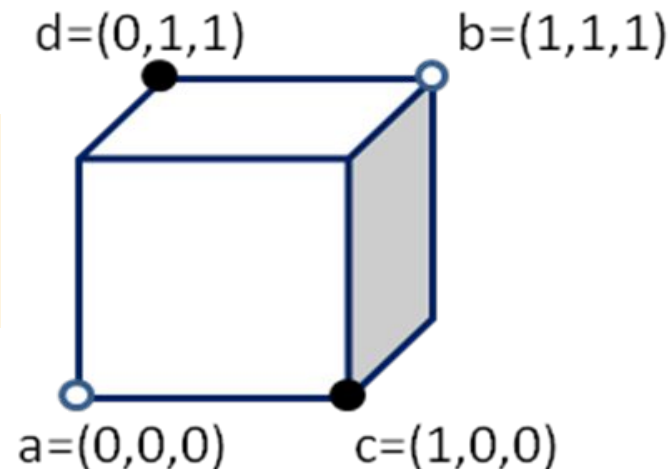
2) if the number of points in  $d$ -dimensional space is bigger than  $2^d$  then they are almost always not linearly separable

A quick check:

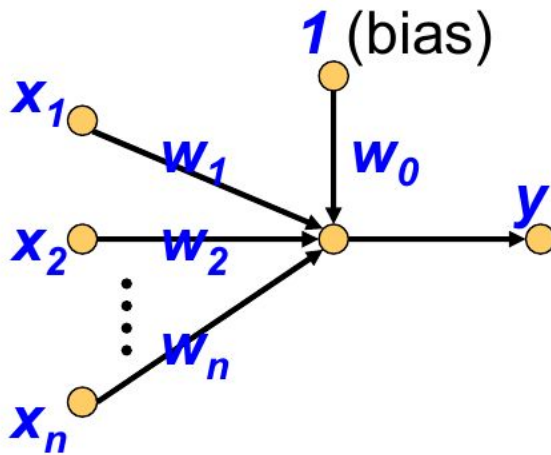
*Are points  $\{a, b\}$  linearly separable from  $\{c, d\}$ ?*

*How does it relate to the Cover's theorem?*

“general position in  $d$  dimensions”  
=  
“not in less than  $d$  dimensions”



# Adaline: Adaptive Linear Element



$$y = w_0 + \sum w_i x_i$$

- Activation function  $g(x)=x$  (identity)
- The desired outputs are -1's or 1's, the actual outputs (y's) are "real numbers"
- Main idea: minimize the squared error:

$$Error = \sum_{Examples} (y_i - d_i)^2$$

# Error function

on pattern  $i$  :  $E(i) = (y_i - d_i)^2$

on all patterns:  $E = \sum (y_i - d_i)^2$

But

$$y_i = w_0 + w_1 x_1 + \dots + w_k x_k,$$

so

$$E = \sum ((w_0 + w_1 x_1 + \dots + w_k x_k) - d_i)^2$$

thus

$E$  is a function of  $w_0, w_1, \dots, w_k$ .

How can we find the minimum of  $E(w_0, w_1, \dots, w_k)$  ?

***By the gradient descent algorithm ...***

# Gradient Descent Algorithm

How to find a minimum of a function  $f(x,y)$  ?

1. Start with an arbitrary point  $(x_0, y_0)$
2. Find a direction in which  $f$  is decreasing most rapidly

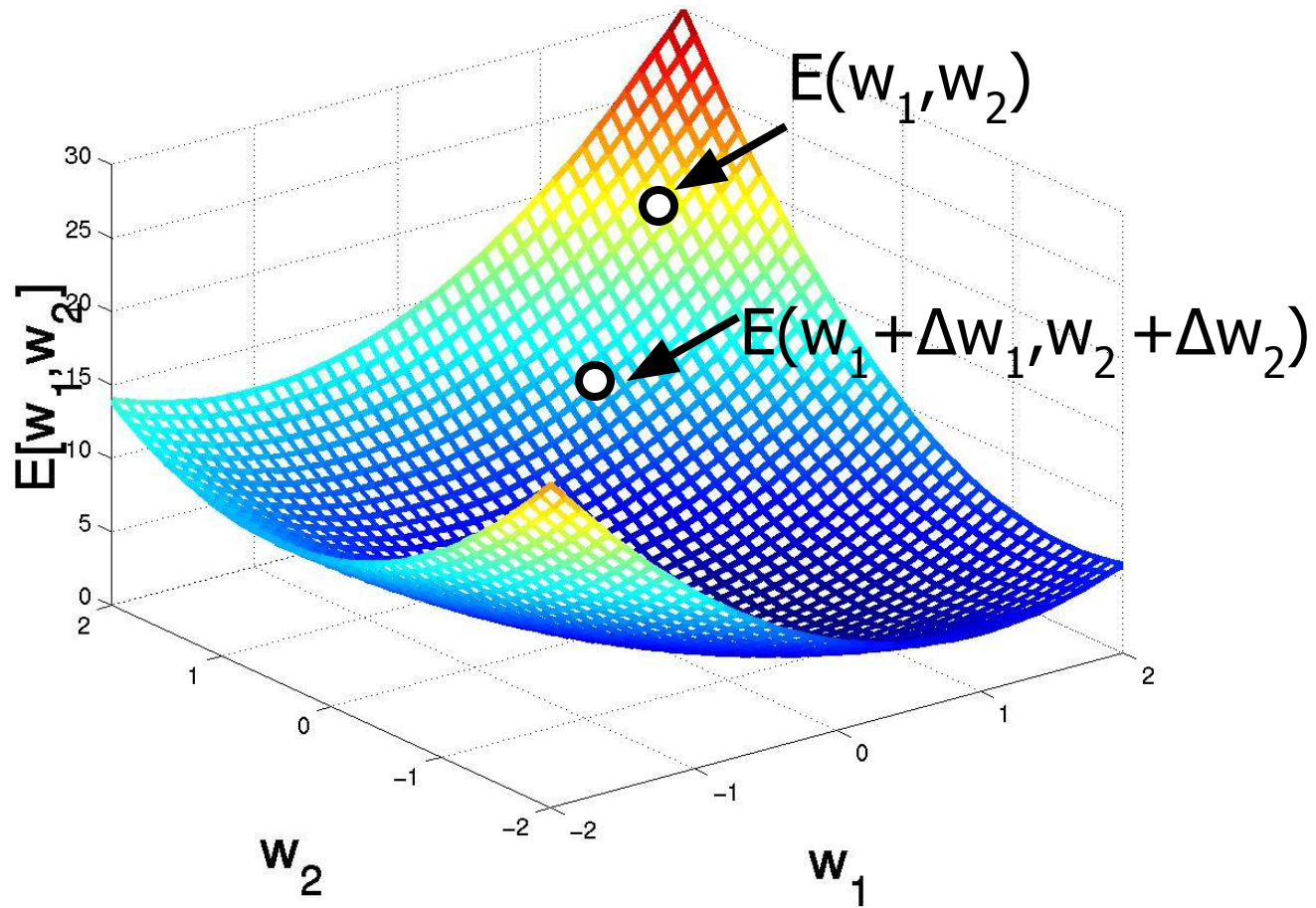
$$-\left[ \frac{\partial f(x_0, y_0)}{\partial x}, \frac{\partial f(x_0, y_0)}{\partial y} \right]$$

3. Make a small step in this direction

$$(x_0, y_0) = (x_0, y_0) - \eta \left[ \frac{\partial f(x_0, y_0)}{\partial x}, \frac{\partial f(x_0, y_0)}{\partial y} \right]$$

4. Repeat the whole process

# Gradient Descent



# Adaline: Gradient Descent

- Find  $w_i$ 's that minimize the squared error

$$E(w_0, \dots, w_m) = \sum (y-d)^2$$

- Gradient:

$$\nabla E[w] = [\partial E / \partial w_0, \dots, \partial E / \partial w_m] ; \Delta w = -\eta \nabla E[w]$$

$$\partial E / \partial w_i = \partial / \partial w_i \sum (y-d)^2 = \partial / \partial w_i \sum (\sum_j w_j x_j - d)^2 = 2 \sum (y-d)(x_i)$$

(summation over all examples)

- The weights should be updated by:  $w_i = w_i + \eta \sum (y-d)x_i$
- Summation over all cases? Split the summation by examples!  
I.e., for each training example,  $w_i = w_i + \eta (y-d)x_i$



# Adaline Training Algorithm

1. start with a random set of weights  $w$
2. select a pattern  $x$  (e.g., at random)
3. update weights:

$$w := w + \eta * (y - d) * x, \text{ (Adaline rule)}$$

where  $d$  - desired output on input  $x$  and  $y = wx$

4. goto 2

**The step size  $\eta$  should be relatively small**

# Adaline and Linear Regression

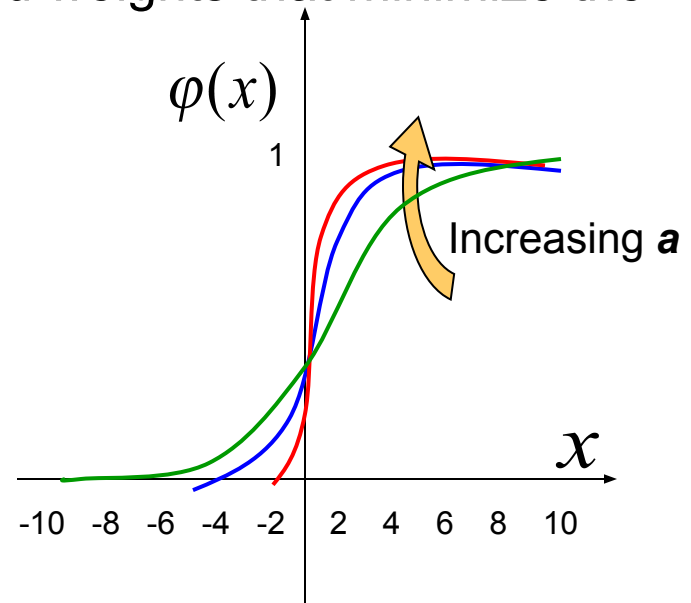
- Perceptron can be used for classification problems only
- Adaline doesn't require that outputs are -1's or 1's - arbitrary values are allowed
- Therefore Adaline can be used for solving Linear Regression Problems
- There is a direct (fast) algorithm for Linear Regression!
- **But:** Adaline requires no memory: it "learns on-the-fly"; it's biologically justified (?)

# “Smooth Perceptron” => Logistic Regression

- Main Idea:  
Replace the sign function by its “smooth approximation” and use the steepest descent algorithm to find weights that minimize the error (as with ADALINE)

$$\varphi(x) = \frac{1}{1+e^{-ax}} \text{ with } a > 0$$

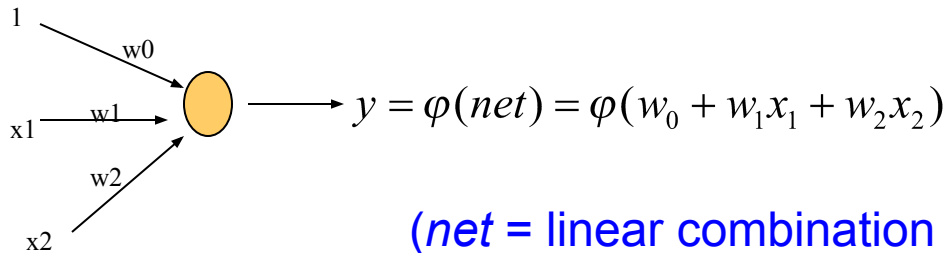
where  $x = w_0 + w_1x_1 + w_2x_2$



- The function to be optimized is a bit more complicated than in ADALINE case
- FORTUNATELY: the update rules are simple !

# Derivation of update rules for simple net

- Derive the Delta rule for the following network



(*net* = linear combination of inputs)

$$E(w_0, w_1, w_2) = \frac{1}{2} (y - d)^2 = \frac{1}{2} (\varphi(w_0 + w_1x_1 + w_2x_2) - d)^2$$

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i} \text{ for } i \text{ in } \{0,1,2\}$$

- We need to find  $\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \frac{\partial E}{\partial w_2}$

# Derivation of Delta Rule

$$\begin{aligned}\frac{\partial E(w_0, w_1, w_2)}{\partial w_1} &= \frac{1}{2} \frac{\partial (\varphi(w_0 + w_1 x_1 + w_2 x_2) - d)^2}{\partial w_1} \\&= \frac{1}{2} 2(\varphi(w_0 + w_1 x_1 + w_2 x_2) - d) \frac{\partial (\varphi(w_0 + w_1 x_1 + w_2 x_2) - d)}{\partial w_1} \\&= (\varphi(net) - d) \varphi'(net) \frac{\partial (w_0 + w_1 x_1 + w_2 x_2)}{\partial w_1} \\&= (\varphi(net) - d) \varphi'(net) x_1\end{aligned}$$

- From similar calculations we get:

$$\frac{\partial E(w_0, w_1, w_2)}{\partial w_2} = (\varphi(net) - d) \varphi'(net) x_2$$

- and  $\frac{\partial E(w_0, w_1, w_2)}{\partial w_0} = (\varphi(net) - d) \varphi'(net)$

## Concluding:

---

$$\Delta w_0 = \eta(d - \varphi(net))\varphi'(net)$$

$$\Delta w_1 = \eta(d - \varphi(net))\varphi'(net)x_1$$

$$\Delta w_2 = \eta(d - \varphi(net))\varphi'(net)x_2$$

$$\Delta \mathbf{w} = \eta(d - \varphi(net))\varphi'(net)\mathbf{x}$$

- It's good to know that for the logistic sigmoid function:  
 $\phi(x) = 1/(1 + \exp(-x))$  we have:  $\phi'(x) = \phi(x)(1 - \phi(x)) = \text{output}(1 - \text{output})$
- *Adaline: Linear Regression*
- *“A Neuron”: “Logistic Regression” (what is the error function?)*

# Summary of learning rules:

Perceptron learning rule:

$$\Delta \mathbf{w} = \boldsymbol{\eta} * \mathbf{x} * (\text{out} - \mathbf{d}) \text{ (for misclassified, output -1/1)}$$

$\mathbf{x}$  = input vector  
 $\text{out}$  = output of the network;  
 $\text{out} = f(\text{net})$ , where:  
 $\text{net}$  = linear comb. of inputs

Adaline learning rule:

$$\Delta \mathbf{w} = \boldsymbol{\eta} * \mathbf{x} * (\text{out} - \mathbf{d})$$

Perceptron:  $f(\text{net}) = \text{sign}(\text{net})$

Adaline :  $f(\text{net}) = \text{net}$

“A Neuron” learning rule:

$$\Delta \mathbf{w} = \boldsymbol{\eta} * \mathbf{x} * (\text{out}(\mathbf{x}) - \mathbf{d}) * \text{out}'(\mathbf{x})$$

“Neuron”:  $f(\text{net}) = 1/(1+\exp(-\text{net}))$

# Perceptron for multi-class problems

▮ So far, we were considering binary classification problems: how to separate two sets of points with a (linear) model? So we were looking for a single (linear) discriminant function...

▮ Multi-class classification problems: we want to separate  $c > 2$  sets of points

## **Linear Separability for multi-class problems:**

There exist  **$c$  linear discriminant functions**  $y_1(x), \dots, y_c(x)$  such that each  $\mathbf{x}$  is assigned to class  $C_k$  if and only if  $y_k(x) > y_j(x)$  for all  $j \neq k$

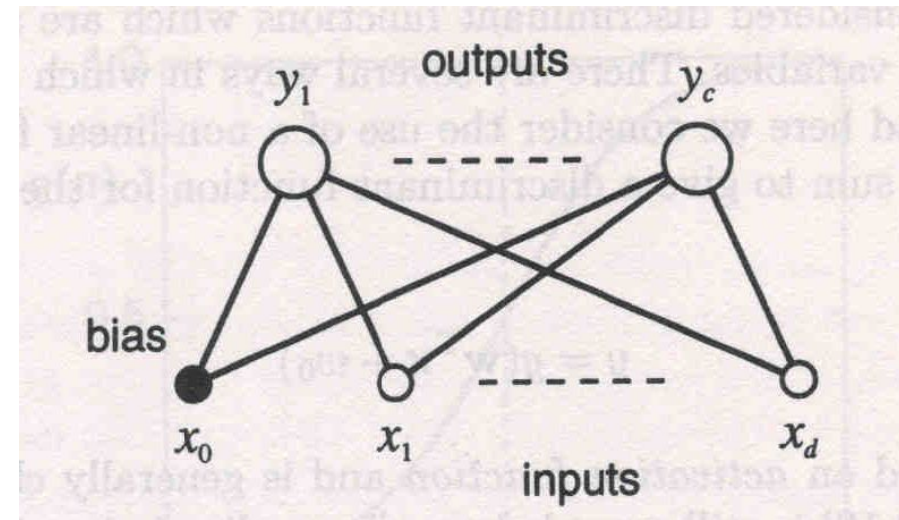
All algorithms discussed so far can be generalized to handle multi-class classification problems. In some cases the generalization is easy, in others not.

## **Generalized Perceptron convergence theorem:**

*If the  $c$  sets of points are linearly separable then the generalized perceptron algorithm terminates after a finite number of iterations, separating all classes.*



# Generalized Perceptron Algorithm (Duda et al.)



initialize weights  $\mathbf{w}$  at random

**while** (there are misclassified training examples)

    Select a misclassified example  $(\mathbf{x}, c_i)$

    Then *some nodes* are activated more than *the node  $c_i$*

        1) update weights of *these nodes* by  $-\mathbf{x}$ :  $\mathbf{w} = \mathbf{w} - \mathbf{x}$ ;

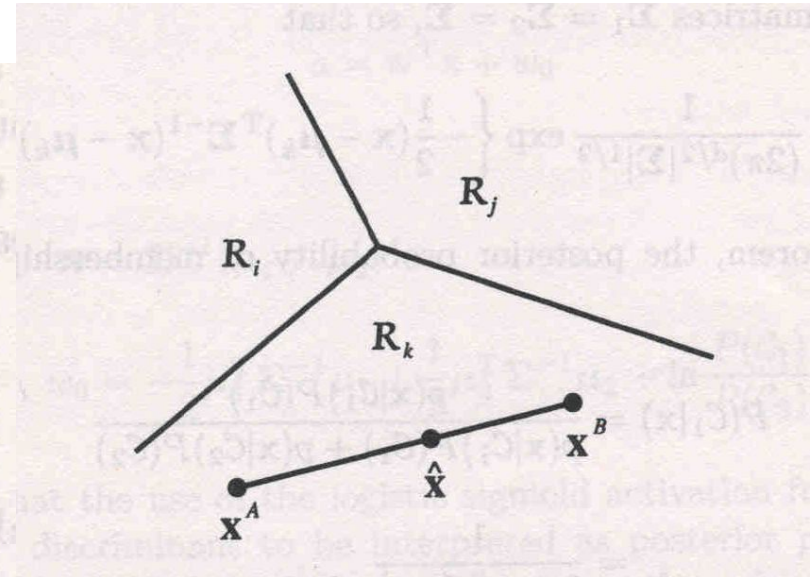
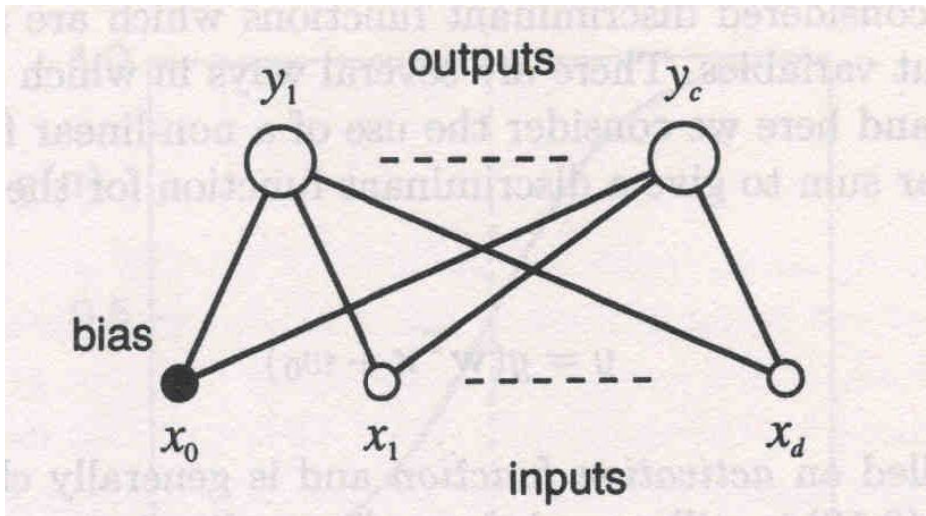
        2) update weights of the node  $c_i$  by  $\mathbf{x}$ :  $\mathbf{w} = \mathbf{w} + \mathbf{x}$ ;

        3) leave weights of all other nodes unchanged

**end-while;**

# Perceptron for multi-class problems

A network of  $c$  perceptrons that share the same input vector represent  $c$  linear discriminant functions and can be used for solving multi-class classification problems.



Note that  $y_i - y_j$  is a linear function which separates class  $i$  from  $j$ , for all  $i, j$ .  
So decision regions are intersections of half-spaces!

Decision regions are always **convex**: for any two points  $x^A$  and  $x^B$  from the same region, the whole line interval between  $x^A$  and  $x^B$  is also in this region.

# Generalized Linear Discriminants

- Allowing non-linear basis functions as inputs for linear models dramatically increases the power of these models: in theory such networks can model any boundary (or function), provided we use an appropriate set of basis functions



**Linear Models (Single Layer Networks) are very important!**

- An important class of basis functions are “radial functions” that measure the distance of  $x$  to specific “reference points”  $x_k$ . This leads to the concept of **Radial Basis Function Networks (RBF-networks)**
- RBF networks will be discussed later ...

# To Remember

- Discriminant functions, linear discriminants, linear separability, Cover's theorem (the plot and the interpretation!)
- The perceptron learning algorithm and key properties: convergence and bad behaviour on non-separable data sets
- The pocket algorithm (also with ratchet) of Gallant
- Adaline and (logistic) Perceptron; Incremental vs Batch Learning
- Derivation of learning rules for Adaline and Perceptron
- Multi-class linear separability
- The generalized perceptron algorithm
- The concept of SVM; quadratic optimization criterion
- Generalized Linear Discriminant