

Faculty of Mathematics, Physics and Informatics  
Comenius University Bratislava



# Neural Networks

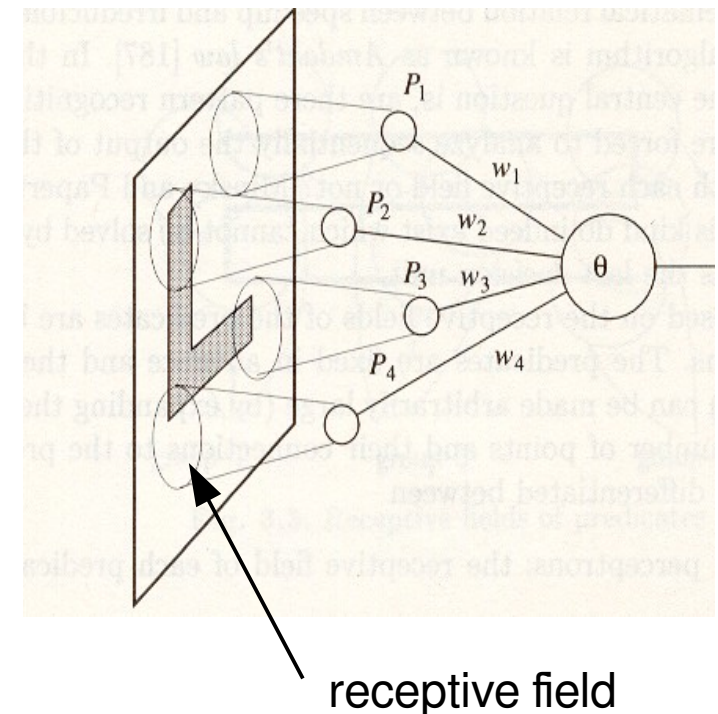
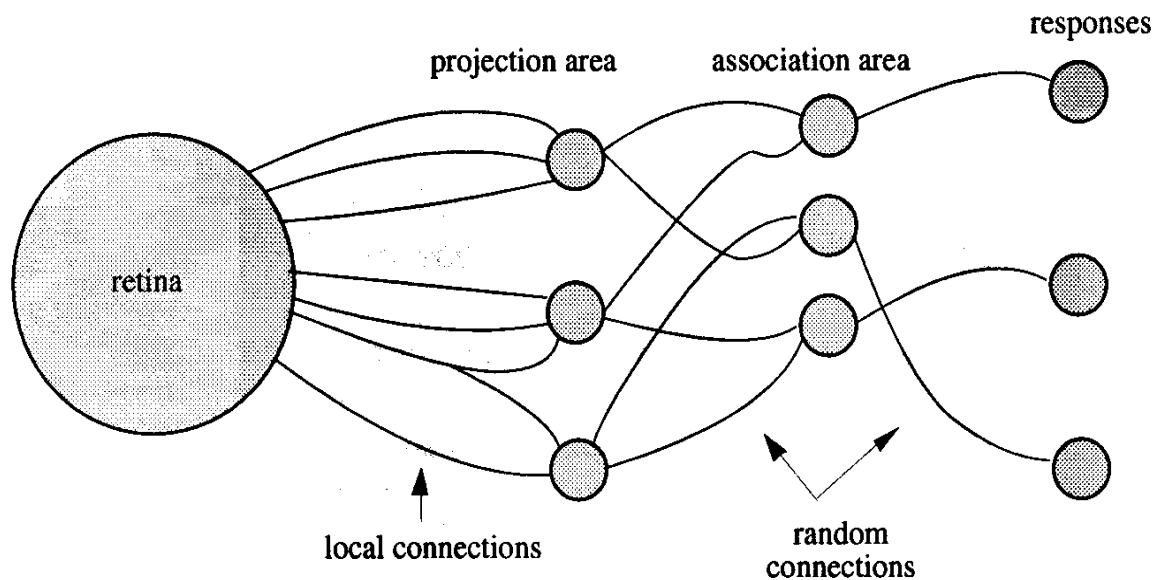
## Lecture 2

### Simple perceptron

# Classical perceptron

In 1958, F. Rosenblatt (American psychologist) proposed perceptron, a more general computational model (than McCulloch-Pitts' TL units) with free parameters, stochastic connectivity and threshold element.

In 1950, Hubel & Wiesel “decoded” the structure of retina and receptive fields.

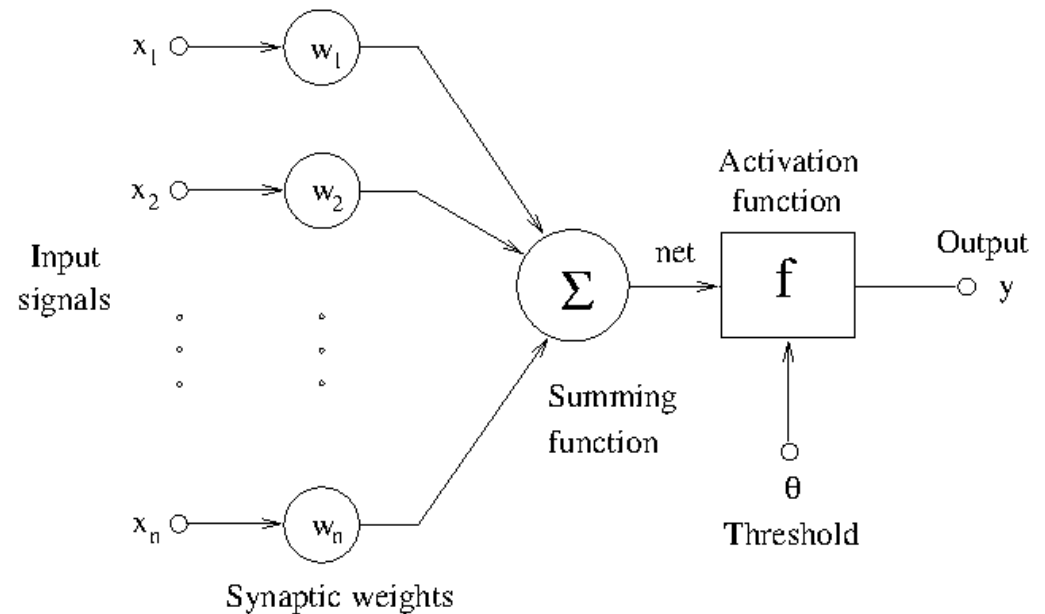


# Discrete perceptron

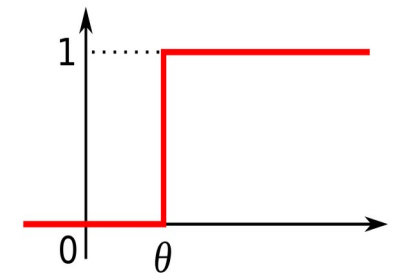
- Inputs  $\mathbf{x}$  , weights  $\mathbf{w}$ , output  $y$
- Function:

$$y = f\left(\sum_{j=1}^n w_j x_j - \theta\right)$$

$$y = f\left(\sum_{j=1}^{n+1} w_j x_j\right) \quad \begin{array}{l} x_{n+1} = -1 \\ w_{n+1} = \theta \end{array}$$



- $f$  = threshold (step) function
- Supervised learning** – with teacher  $d$
- e.g.  $d = 1$  for  $\forall \mathbf{x} \in C1$ ,  $d = 0$  for  $\forall \mathbf{x} \in C2$
- Learning: if  $\mathbf{w}^T \mathbf{x} \leq 0$  but  $\mathbf{x} \in C1$ , then  $\mathbf{w}(t+1) = \mathbf{w}(t) + \mathbf{x}$   
if  $\mathbf{w}^T \mathbf{x} > 0$  but  $\mathbf{x} \in C2$ , then  $\mathbf{w}(t+1) = \mathbf{w}(t) - \mathbf{x}$



Or:  $w_j(t+1) = w_j(t) + \alpha (d - y) x_j$        $\alpha$  = learning rate

# Perceptron algorithm

*Given:* input-target  $\{\mathbf{x}^{(p)}, d^{(p)}\}$  pairs, unipolar perceptron

*Initialization:* randomize weights, shuffle the pairs, set learning rate, set  $E = 0$ .

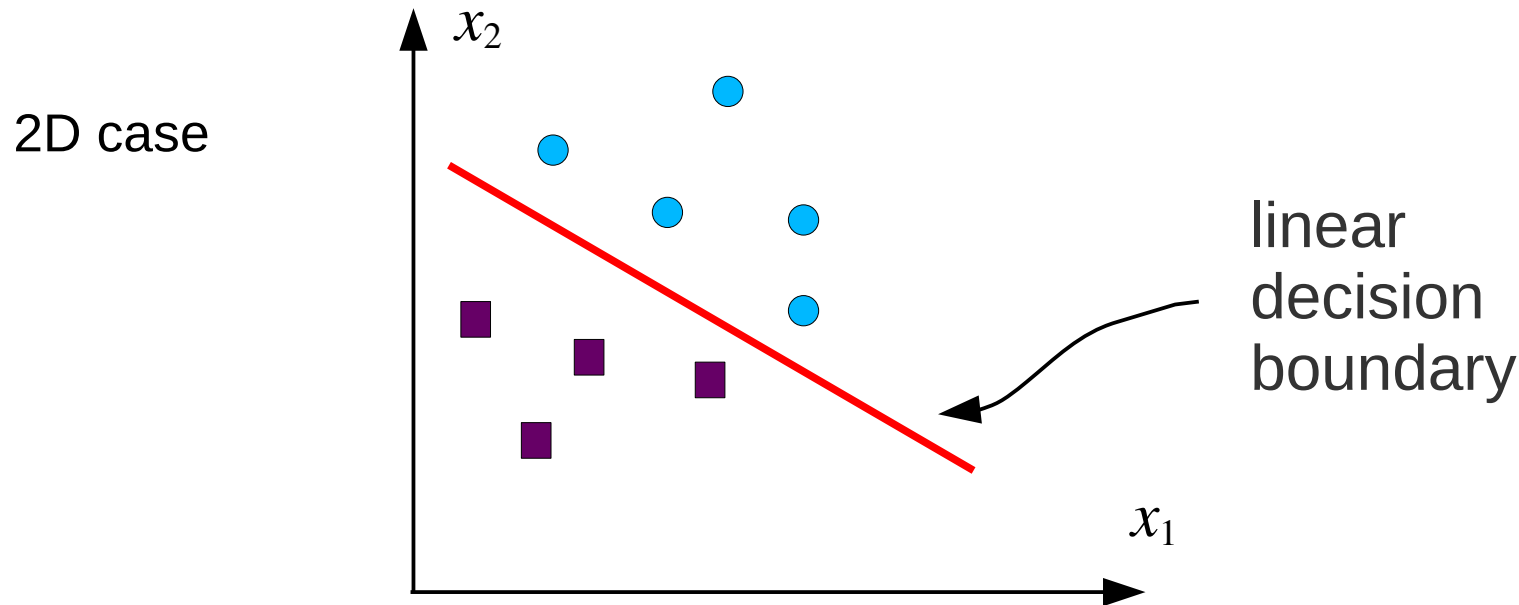
*Training:*

1. choose next input  $\mathbf{x}$ , compute output  $y$
2. evaluate error function,  $e(t) = \frac{1}{2} (d^{(p)} - y^{(p)})^2$ ,  $E = E + e(t)$
3. if  $e(t) > 0$ , adjust weights using perceptron rule
4. if not all inputs used, then goto 1, else goto 5
5. if  $E == 0$  (all inputs in the set classified correctly), then end  
else reshuffle the pairs,  $E = 0$ , go to 1.

# Perceptron classification capacity

$$w_1x_1 + w_2x_2 + \dots + w_nx_n = \theta$$

linear separability of two classes

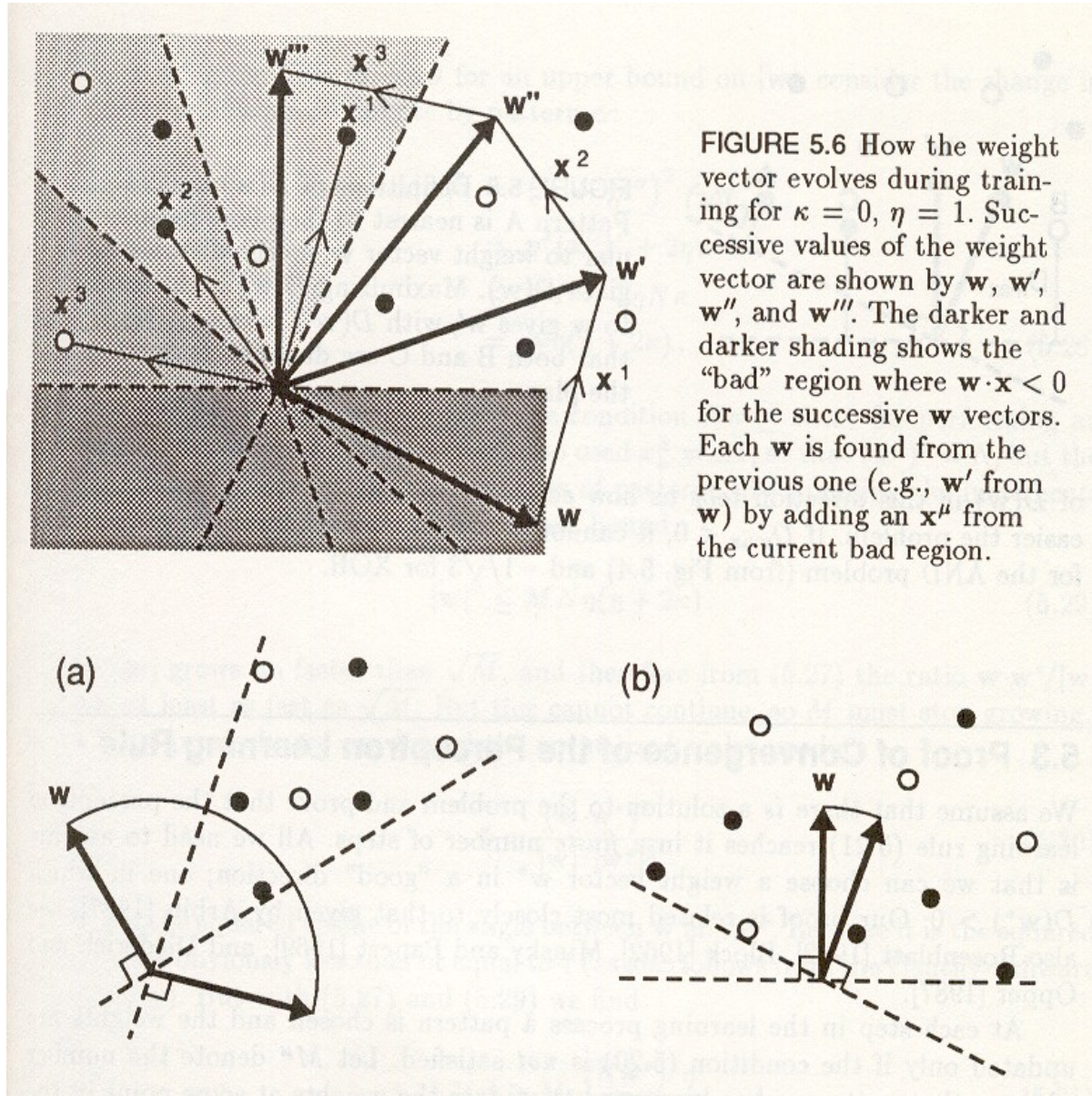


**Fixed-increment convergence theorem** (Rosenblatt, 1962): “Let the classes *A* and *B* be finite and linearly separable, then perceptron learning algorithm converges (updates its weight vector) in a finite number of steps.”



# Finding a solution

$w^T x > 0$  for C1  
 $w^T x \leq 0$  for C2



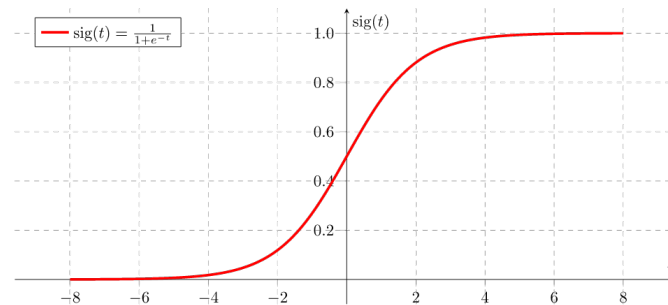
(Hertz et al, 1990)

# Continuous perceptron

- **Nonlinear unit** with activation function:  $y = f(net) = 1 / (1 + e^{-net})$

- Has nice properties:

- boundedness
- monotonicity
- differentiability

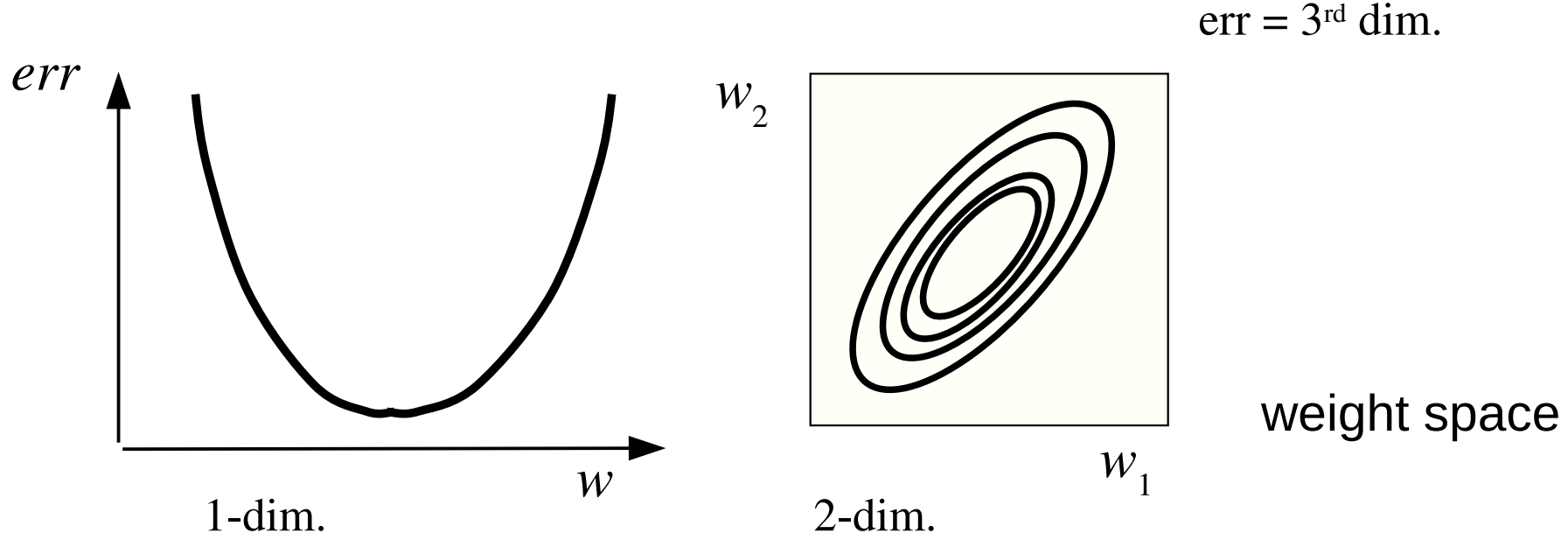


(unipolar)  
sigmoid

- Error function (e.g. quadratic):  $E(\mathbf{w}) = \sum_p e^{(p)} = \frac{1}{2} \sum_p (d^{(p)} - y^{(p)})^2$  over dataset  $p$ , also called **loss function** (objective function)
- We want to minimize the error function: necessary conditions  $e(\mathbf{w}^*) \leq e(\mathbf{w})$  and  $\nabla e(\mathbf{w}^*) = 0$ , gradient operator  $\nabla = [\partial/\partial w_1, \partial/\partial w_2, \dots]^T$ . Minimizing  $E(\mathbf{w})$  leads to
- (stochastic, online) **gradient descent learning**:  
$$w_j(t+1) = w_j(t) + \alpha (d^{(p)} - y^{(p)}) f'(net) x_j = w_j(t) + \alpha \delta^{(p)} x_j^{(p)}$$

# Error surface for a continuous perceptron

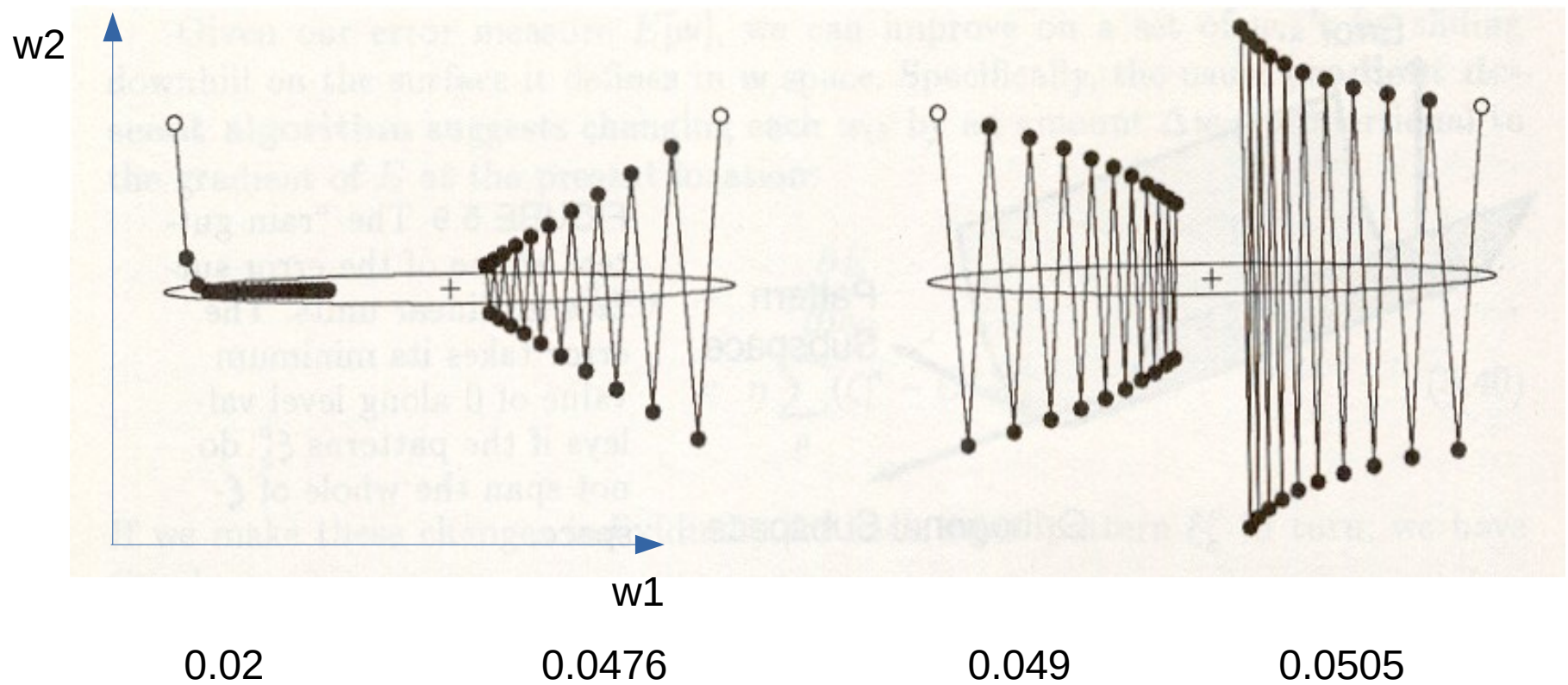
- Assume 1 neuron, linear or with a sigmoid function
- The output error  $e = f(w_1, w_2, \dots, w_n)$ , assume quadratic error
- For a linear neuron with  $n$  inputs, we have a convex function (quadratic bowl); vertical cross-sections are parabolas; horizontal cross-sections are ellipses.





# Effect of learning rate

$$E = w_1^2 + 20 w_2^2$$



(Hertz et al, 1990)

# Linear neuron as a least-squares filter

- Consider:  $y = \mathbf{w}^T \mathbf{x} = \mathbf{x}^T \mathbf{w}$ , input-target pairs  $\{\mathbf{x}^{(p)}, d^{(p)}\}$ ,  $p = 1, \dots, N$
- Collect inputs  $\mathbf{X} = [\mathbf{x}^{(1)} \mathbf{x}^{(2)} \dots \mathbf{x}^{(N)}]^T$  ( $N \times n$  matrix)
- Let  $\mathbf{e} = [e^{(1)} e^{(2)} \dots e^{(N)}]^T$  then output error  $\mathbf{e} = \mathbf{d} - \mathbf{X} \cdot \mathbf{w}$
- Gauss-Newton method:  $E(\mathbf{w}) = 1/2 \sum_p (d^{(p)} - y^{(p)})^2$ , compute  $\nabla \mathbf{e}$  ( $n \times N$ )
- $j_{pk} = \partial e^{(p)} / \partial w_k \Rightarrow$  Jacobian  $\mathbf{J}(t) = [j_{pk}]$  is ( $N \times n$ )  $\mathbf{J}(t) = -\mathbf{X}(t) = [\nabla \mathbf{e}^T]$
- $\mathbf{e}'(\mathbf{w})^{(N)} = \mathbf{e}(\mathbf{w}) + \mathbf{J}^{(N)} \cdot (\mathbf{w} - \mathbf{w}^{(N)})$ . [linearity assumption of error f.]
- Substitute [ $N \equiv t$ ]  $\mathbf{w}(t+1) = \arg \min_{\mathbf{w}} \{ 1/2 \|\mathbf{e}'(t, \mathbf{w})\|^2 \} \dots$
- Update  $\mathbf{w}(t+1) = \mathbf{w}(t) - (\mathbf{J}^T(t) \mathbf{J}(t))^{-1} \mathbf{J}(t) \mathbf{e}(t) = \mathbf{w}(t) + (\mathbf{X}^T(t) \mathbf{X}(t))^{-1} \mathbf{X}(t) [\mathbf{d}(t) - \mathbf{X}(t) \mathbf{w}(t)] = [\mathbf{X}^T(t) \mathbf{X}(t)]^{-1} \mathbf{X}(t) \mathbf{d}(t) \Rightarrow \mathbf{w}(t+1) = \mathbf{X}^+(t) \mathbf{d}(t)$
- “The weight vector  $\mathbf{w}(t+1)$  solves the least-squares problem in an observation interval until time  $t$ .”
- neuron as a linear regressor

# Alternative loss function: cross entropy

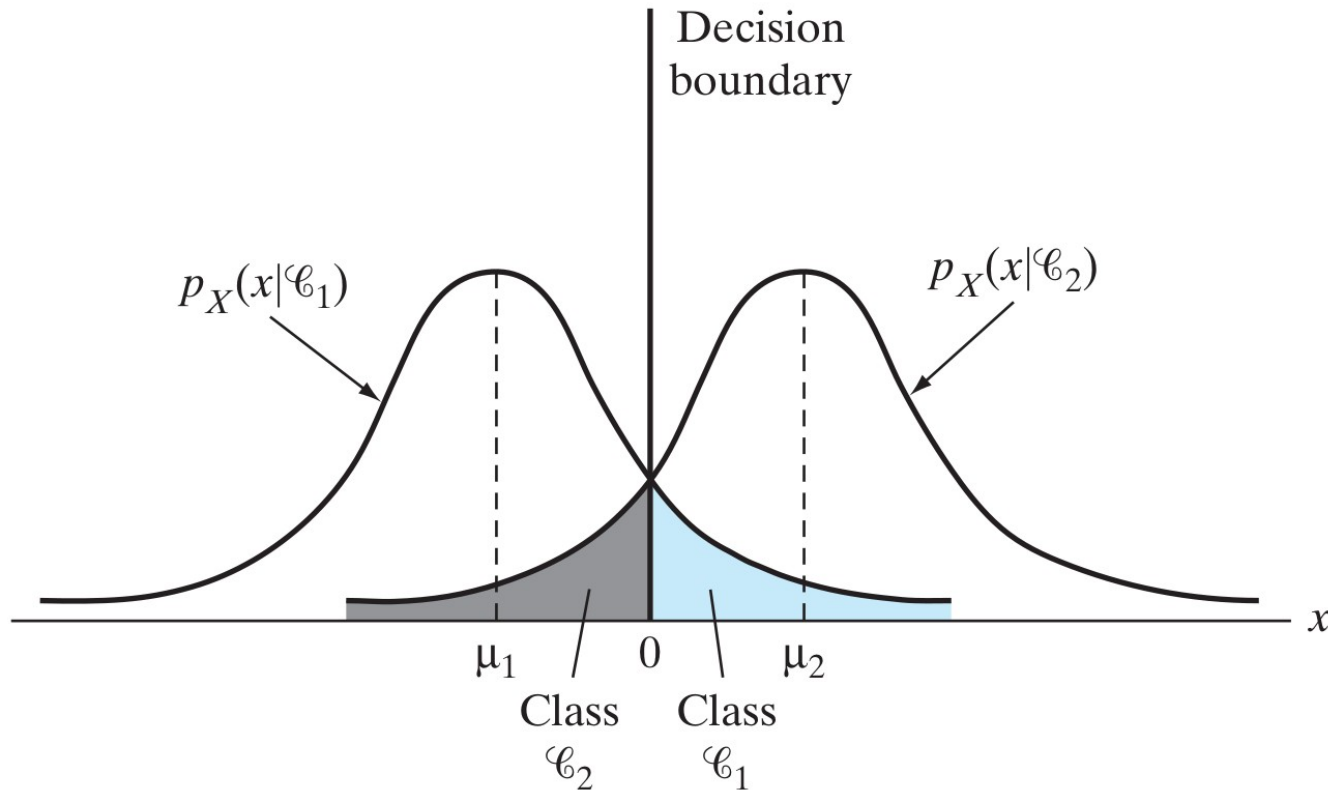
- Useful for classification, leads to probability (uncertainty)
- Error function – **cross-entropy** (for one output):  
[relative entropy b/w empirical probability distribution  $(d^{(p)}, 1 - d^{(p)})$  and output distribution  $(y, 1 - y)$ ]

$$E_{CE}(\mathbf{w}) = \sum_p E^{(p)} = - \sum_p [d^{(p)} \ln y^{(p)} + (1 - d^{(p)}) \ln (1 - y^{(p)})]$$

- → negative log-likelihood (in probabilistic models)
- minimization of  $E^{(p)}$  results in a learning rule:
$$w_j(t+1) = w_j(t) + \alpha (d^{(p)} - y^{(p)}) x_j^{(p)}$$
- *Note:* In case of 2 classes one can use logistic unit (for  $C_1: d = 1; C_2: d = 0$ ),
- then the output  $y$  can be interpreted as  $P(C_1 | \mathbf{x}) = 1 - P(C_2 | \mathbf{x})$

# Bayes classifier for two classes

- (linear) Bayes classifier for a 1D Gaussian environment
- for convenience, decision border is at  $x = 0$



# Perceptron link to Bayes classifier

## Assumptions:

- random vector  $X$ , two classes  $C_1: E[X] = \mathbf{m}_1$ ,  $C_2: E[X] = \mathbf{m}_2$
- covariance matrix  $\mathbf{C} = E[(X - \mathbf{m}_1)(X - \mathbf{m}_1)^T] = E[(X - \mathbf{m}_2)(X - \mathbf{m}_2)^T]$

We can express conditional probability density function:

$$f(\mathbf{x} | C_i) = [(2\pi)^{m/2} \det(\mathbf{C})^{1/2}]^{-1} \exp[-1/2(\mathbf{x} - \mathbf{m}_i)^T \mathbf{C}^{-1}(\mathbf{x} - \mathbf{m}_i)]$$

$\mathbf{x}$  – observation vector,  $i = \{1, 2\}$

- the 2 classes are equiprobable, i.e.  $p_1 = p_2$  (a priori probs)
- (mis)classifications carry the same cost, i.e.  $\omega_{12} = \omega_{21}$ ,  $\omega_{11} = \omega_{22} = 0$

**Bayes classifier:** “If  $p_1(\omega_{21} - \omega_{11}) f(\mathbf{x} | C_1) > p_2(\omega_{12} - \omega_{22}) f(\mathbf{x} | C_2)$ , assign the observation vector  $\mathbf{x}$  to  $C_1$ . Otherwise, assign it to  $C_2$ .”

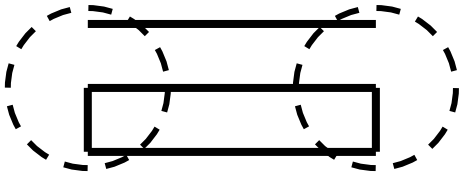
## Bayes classifier (ctd)

- Define likelihood ratio  $\Lambda(\mathbf{x}) = f(\mathbf{x} | C_1) / f(\mathbf{x} | C_2)$  and threshold  $\xi = [p_2(\omega_{12} - \omega_{22})] / [p_1(\omega_{21} - \omega_{11})]$ . Then:
- $\log \Lambda(\mathbf{x}) = -\frac{1}{2} (\mathbf{x} - \mathbf{m}_1)^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{m}_1) + \frac{1}{2} (\mathbf{x} - \mathbf{m}_2)^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{m}_2)$
- $\log \xi = 0$
- *Then*: we get a linear Bayes classifier  $y = \mathbf{w}^T \mathbf{x} + b$  where  $y = \log \Lambda(\mathbf{x})$ ,  $\mathbf{w} = \mathbf{C}^{-1}(\mathbf{m}_1 - \mathbf{m}_2)$ ,  $b = \frac{1}{2} (\mathbf{m}_2^T \mathbf{C}^{-1} \mathbf{m}_2 - \mathbf{m}_1^T \mathbf{C}^{-1} \mathbf{m}_1) \rightarrow$   
**log-likelihood test**: If  $y > 0$ , then  $\mathbf{x} \in C_1$ , else  $C_2$ .
- Differences b/w Perceptron (P) and Bayes classifier (BC):
  - P assumes linear separability, BC does not
  - P convergence algorithm is non-parametric, unlike BC
  - P convergence algorithm is adaptive and simple, unlike BC.



# Perceptron limits - XOR

1.



2.



3.



4.



- Consider a perceptron classifying shapes as connected or disconnected and taking inputs from shape ends (shown as dashed circles for pattern 1)
- The problem arises because a single layer of processing local knowledge cannot be combined into global knowledge
- No feature-weighting machine (such as a simple perceptron) can do this type of separation, because information about the relation between the bits of evidence is lost (proven by Minsky & Papert, 1969)
- This problem caused the loss of interest in connectionism (in 1970s), since many real problems are not linearly separable.

# Summary

- binary and continuous perceptron
- single perceptron can linearly separate two classes:
- perceptron as a detector (of half input space)
- optimization: gradient descent learning
- link to adaptive filtering – error correction learning
- two types of error (loss) functions
- link to statistics: probabilistic Bayes classifier
- limitations of a simple perceptron
- a single neuron model, with any activation function, is linear in its parameters.