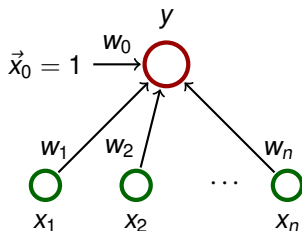


# ADALINE

## Architecture:



$\vec{w} = (w_0, w_1, \dots, w_n)$  and  $\vec{x} = (x_0, x_1, \dots, x_n)$  where  $x_0 = 1$ .

## Activity:

- ▶ inner potential:  $\xi = w_0 + \sum_{i=1}^n w_i x_i = \sum_{i=0}^n w_i x_i = \vec{w} \cdot \vec{x}$
- ▶ activation function:  $\sigma(\xi) = \xi$
- ▶ network function:  $y[\vec{w}](\vec{x}) = \sigma(\xi) = \vec{w} \cdot \vec{x}$

## Learning:

- ▶ Given a **training set**

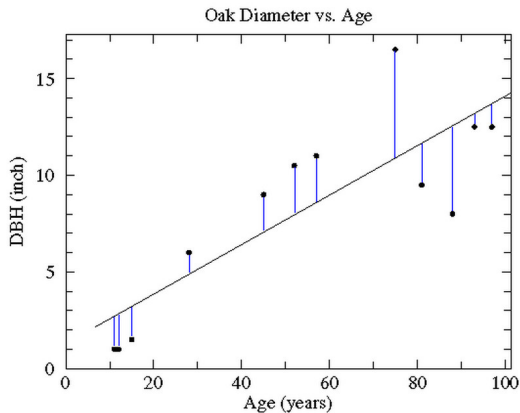
$$\mathcal{T} = \{(\vec{x}_1, d_1), (\vec{x}_2, d_2), \dots, (\vec{x}_p, d_p)\}$$

Here  $\vec{x}_k = (x_{k0}, x_{k1}, \dots, x_{kn}) \in \mathbb{R}^{n+1}$ ,  $x_{k0} = 1$ , is the  $k$ -th input, and  $d_k \in \mathbb{R}$  is the expected output.

Intuition: The network is supposed to compute an affine approximation of the function (some of) whose values are given in the training set.

# Oaks in Wisconsin

<b>Age (years)</b>	<b>DBH (inch)</b>
97	12.5
93	12.5
88	8.0
81	9.5
75	16.5
57	11.0
52	10.5
45	9.0
28	6.0
15	1.5
12	1.0
11	1.0

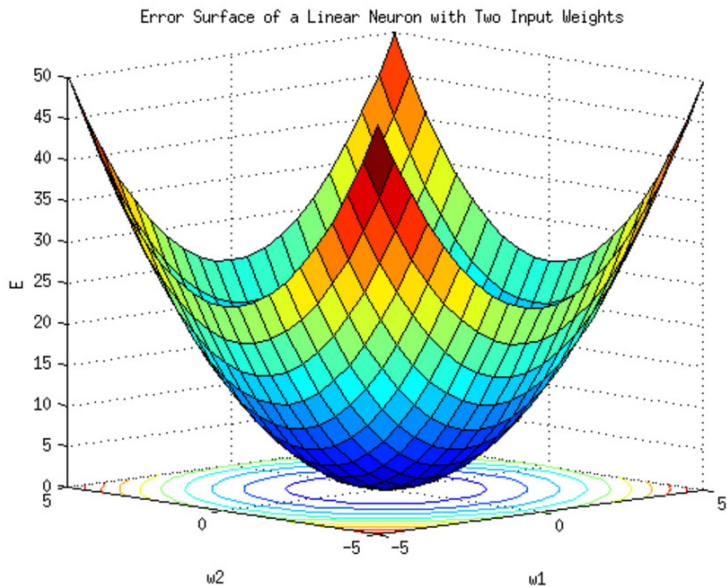


- ▶ **Error function:**

$$E(\vec{w}) = \frac{1}{2} \sum_{k=1}^p (\vec{w} \cdot \vec{x}_k - d_k)^2 = \frac{1}{2} \sum_{k=1}^p \left( \sum_{i=0}^n w_i x_{ki} - d_k \right)^2$$

- ▶ The goal is to find  $\vec{w}$  which minimizes  $E(\vec{w})$ .

# Error function



# Gradient of the error function

Consider **gradient** of the error function:

$$\nabla E(\vec{w}) = \left( \frac{\partial E}{\partial w_0}(\vec{w}), \dots, \frac{\partial E}{\partial w_n}(\vec{w}) \right)$$

Intuition:  $\nabla E(\vec{w})$  is a vector in the **weight space** which points in the direction of the *steepest ascent* of the error function.

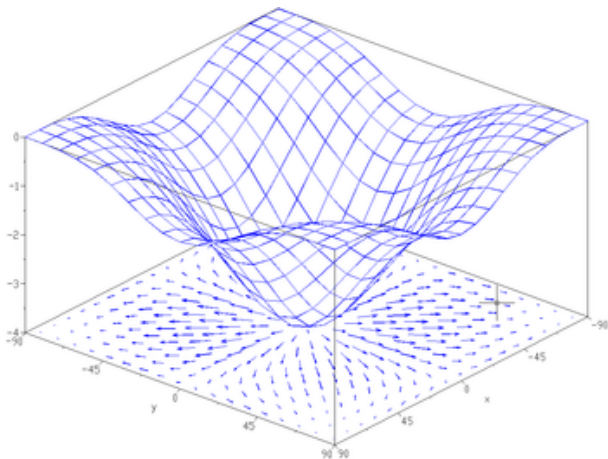
Note that the vectors  $\vec{x}_k$  are just parameters of the function  $E$ , and are thus fixed!

## Fact

If  $\nabla E(\vec{w}) = \vec{0} = (0, \dots, 0)$ , then  $\vec{w}$  is a global minimum of  $E$ .

For ADALINE, the error function  $E(\vec{w})$  is a convex paraboloid and thus has the unique global minimum.

# Gradient - illustration



Caution! This picture just illustrates the notion of gradient ... it is not the convex paraboloid  $E(\vec{w})$  !

## Gradient of the error function (ADALINE)

$$\begin{aligned}\frac{\partial E}{\partial w_\ell}(\vec{w}) &= \frac{1}{2} \sum_{k=1}^p \frac{\delta E}{\delta w_\ell} \left( \sum_{i=0}^n w_i x_{ki} - d_k \right)^2 \\ &= \frac{1}{2} \sum_{k=1}^p 2 \left( \sum_{i=0}^n w_i x_{ki} - d_k \right) \frac{\delta E}{\delta w_\ell} \left( \sum_{i=0}^n w_i x_{ki} - d_k \right) \\ &= \frac{1}{2} \sum_{k=1}^p 2 \left( \sum_{i=0}^n w_i x_{ki} - d_k \right) \left( \sum_{i=0}^n \left( \frac{\delta E}{\delta w_\ell} w_i x_{ki} \right) - \frac{\delta E}{\delta w_\ell} d_k \right) \\ &= \sum_{k=1}^p \left( \vec{w} \cdot \vec{x}_k - d_k \right) x_{k\ell}\end{aligned}$$

Thus

$$\nabla E(\vec{w}) = \left( \frac{\partial E}{\partial w_0}(\vec{w}), \dots, \frac{\partial E}{\partial w_n}(\vec{w}) \right) = \sum_{k=1}^p \left( \vec{w} \cdot \vec{x}_k - d_k \right) \vec{x}_k$$



# ADALINE - learning

## Batch algorithm (gradient descent):

**Idea:** In every step "move" the weights in the direction *opposite* to the gradient.

The algorithm computes a sequence of weight vectors

$\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \dots$

- ▶ weights in  $\vec{w}^{(0)}$  are randomly initialized to values close to 0
- ▶ in the step  $t + 1$ , weights  $\vec{w}^{(t+1)}$  are computed as follows:

$$\begin{aligned}\vec{w}^{(t+1)} &= \vec{w}^{(t)} - \varepsilon \cdot \nabla E(\vec{w}^{(t)}) \\ &= \vec{w}^{(t)} - \varepsilon \cdot \sum_{k=1}^p (\vec{w}^{(t)} \cdot \vec{x}_k - d_k) \cdot \vec{x}_k\end{aligned}$$

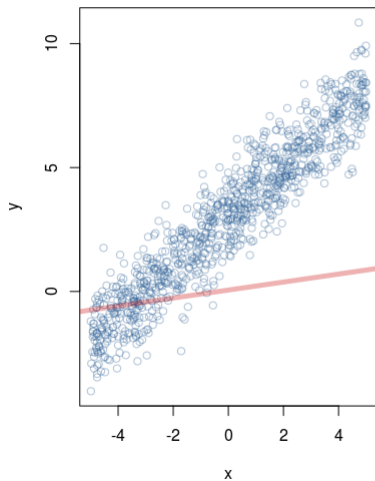
Here  $k = (t \bmod p) + 1$  and  $0 < \varepsilon \leq 1$  is a *learning rate*.

## Proposition

For sufficiently small  $\varepsilon > 0$  the sequence  $\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \dots$  converges (componentwise) to the global minimum of  $E$  (i.e. to the vector  $\vec{w}$  satisfying  $\nabla E(\vec{w}) = \vec{0}$ ).

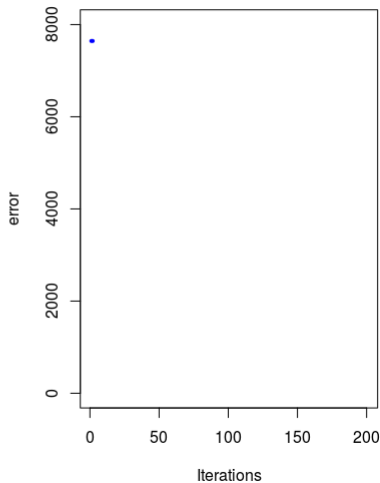
# ADALINE – Animation

Linear regression by gradient descent



Linear regression by gradient descent

Error function



Error function

# ADALINE - learning

## Online algorithm (Delta-rule, Widrow-Hoff rule):

- ▶ weights in  $\vec{w}^{(0)}$  initialized randomly close to 0
- ▶ in the step  $t + 1$ , weights  $\vec{w}^{(t+1)}$  are computed as follows:

$$\vec{w}^{(t+1)} = \vec{w}^{(t)} - \varepsilon(t) \cdot (\vec{w}^{(t)} \cdot \vec{x}_k - d_k) \cdot \vec{x}_k$$

Here  $k = t \bmod p + 1$  and  $0 < \varepsilon(t) \leq 1$  is a learning rate in the step  $t + 1$ .

Note that the algorithm does not work with the complete gradient but only with its part determined by the currently considered training example.

## Theorem (Widrow & Hoff)

*If  $\varepsilon(t) = \frac{1}{t}$ , then  $\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \dots$  converges to the global minimum of  $E$ .*

# ADALINE - classification

How to use the ADALINE for classification?

- ▶ The training set is

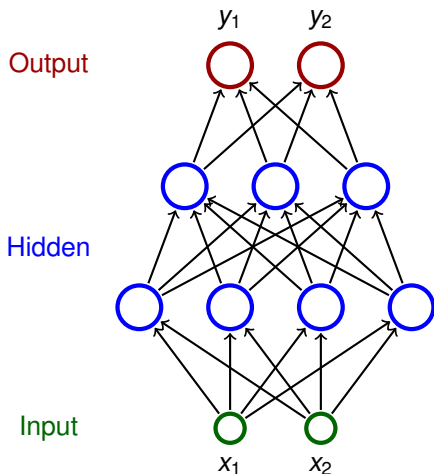
$$\mathcal{T} = \left\{ \left( \vec{x}_1, d_1 \right), \left( \vec{x}_2, d_2 \right), \dots, \left( \vec{x}_p, d_p \right) \right\}$$

where  $\vec{x}_k = (x_{k0}, x_{k1}, \dots, x_{kn}) \in \mathbb{R}^{n+1}$  and  $d_k \in \{1, -1\}$ .

Here  $d_k$  determines a class.

- ▶ Train the network using the ADALINE algorithm.
- ▶ We may expect the following:
  - ▶ if  $d_k = 1$ , then  $\vec{w} \cdot \vec{x}_k \geq 0$
  - ▶ if  $d_k = -1$ , then  $\vec{w} \cdot \vec{x}_k < 0$
- ▶ This does not have to be always true but if the training set is reasonably linearly separable, then the algorithm typically gives satisfactory results.

# Architecture – Multilayer Perceptron (MLP)



- ▶ Neurons partitioned into **layers**; one input layer, one output layer, possibly several hidden layers
- ▶ layers numbered from 0; the input layer has number 0
  - ▶ E.g. three-layer network has two hidden layers and one output layer
- ▶ Neurons in the  $i$ -th layer are connected with all neurons in the  $i + 1$ -st layer
- ▶ Architecture of a MLP is typically described by numbers of neurons in individual layers (e.g. 2-4-3-2)

## Notation:

- ▶ Denote
  - ▶  $X$  a set of *input* neurons
  - ▶  $Y$  a set of *output* neurons
  - ▶  $Z$  a set of *all* neurons ( $X, Y \subseteq Z$ )
- ▶ individual neurons denoted by indices  $i, j$  etc.
  - ▶  $\xi_j$  is the inner potential of the neuron  $j$  *after the computation stops*
  - ▶  $y_j$  is the output of the neuron  $j$  *after the computation stops*

(define  $y_0 = 1$  is the value of the formal unit input)

- ▶  $w_{ji}$  is the **weight of the connection from  $i$  to  $j$**   
(in particular,  $w_{j0}$  is the weight of the connection from the formal unit input, i.e.  $w_{j0} = -b_j$  where  $b_j$  is the bias of the neuron  $j$ )
- ▶  $j_{\leftarrow}$  is a set of all  $i$  such that  $j$  is adjacent from  $i$   
(i.e. there is an arc **to**  $j$  from  $i$ )
- ▶  $j_{\rightarrow}$  is a set of all  $i$  such that  $j$  is adjacent to  $i$   
(i.e. there is an arc **from**  $j$  to  $i$ )

## Activity:

- ▶ inner potential of neuron  $j$ :

$$\xi_j = \sum_{i \in j_{\leftarrow}} w_{ji} y_i$$

- ▶ activation function  $\sigma_j$  for neuron  $j$  (arbitrary differentiable) [ e.g. logistic sigmoid  $\sigma_j(\xi) = \frac{1}{1+e^{-\lambda_j \xi}}$  ]
- ▶ State of non-input neuron  $j \in Z \setminus X$  after the computation stops:

$$y_j = \sigma_j(\xi_j)$$

( $y_j$  depends on the configuration  $\vec{w}$  and the input  $\vec{x}$ , so we sometimes write  $y_j(\vec{w}, \vec{x})$ )

- ▶ The network computes a function  $\mathbb{R}^{|X|}$  to  $\mathbb{R}^{|Y|}$ . Layer-wise computation: First, all input neurons are assigned values of the input. In the  $\ell$ -th step, all neurons of the  $\ell$ -th layer are evaluated.

## Learning:

- ▶ Given a **training set**  $\mathcal{T}$  of the form

$$\left\{ \left( \vec{x}_k, \vec{d}_k \right) \mid k = 1, \dots, p \right\}$$

Here, every  $\vec{x}_k \in \mathbb{R}^{|X|}$  is an *input vector* and every  $\vec{d}_k \in \mathbb{R}^{|Y|}$  is the desired network output. For every  $j \in Y$ , denote by  $d_{kj}$  the desired output of the neuron  $j$  for a given network input  $\vec{x}_k$  (the vector  $\vec{d}_k$  can be written as  $(d_{kj})_{j \in Y}$ ).

- ▶ **Error function:**

$$E(\vec{w}) = \sum_{k=1}^p E_k(\vec{w})$$

where

$$E_k(\vec{w}) = \frac{1}{2} \sum_{j \in Y} \left( y_j(\vec{w}, \vec{x}_k) - d_{kj} \right)^2$$



# MLP – learning algorithm

## Batch algorithm (gradient descent):

The algorithm computes a sequence of weight vectors  $\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \dots$

- ▶ weights in  $\vec{w}^{(0)}$  are randomly initialized to values close to 0
- ▶ in the step  $t + 1$  (here  $t = 0, 1, 2 \dots$ ), weights  $\vec{w}^{(t+1)}$  are computed as follows:

$$w_{ji}^{(t+1)} = w_{ji}^{(t)} + \Delta w_{ji}^{(t)}$$

where

$$\Delta w_{ji}^{(t)} = -\varepsilon(t) \cdot \frac{\partial E}{\partial w_{ji}}(\vec{w}^{(t)})$$

is a *weight update* of  $w_{ji}$  in step  $t + 1$  and  $0 < \varepsilon(t) \leq 1$  is a *learning rate* in step  $t + 1$ .

Note that  $\frac{\partial E}{\partial w_{ji}}(\vec{w}^{(t)})$  is a component of the gradient  $\nabla E$ , i.e. the weight update can be written as  $\vec{w}^{(t+1)} = \vec{w}^{(t)} - \varepsilon(t) \cdot \nabla E(\vec{w}^{(t)})$ .

# MLP – error function gradient

For every  $w_{ji}$  we have

$$\frac{\partial E}{\partial w_{ji}} = \sum_{k=1}^p \frac{\partial E_k}{\partial w_{ji}}$$

where for every  $k = 1, \dots, p$  holds

$$\frac{\partial E_k}{\partial w_{ji}} = \frac{\partial E_k}{\partial y_j} \cdot \sigma'_j(\xi_j) \cdot y_i$$

and for every  $j \in Z \setminus X$  we get

$$\frac{\partial E_k}{\partial y_j} = y_j - d_{kj} \quad \text{for } j \in Y$$

$$\frac{\partial E_k}{\partial y_j} = \sum_{r \in j^{\rightarrow}} \frac{\partial E_k}{\partial y_r} \cdot \sigma'_r(\xi_r) \cdot w_{rj} \quad \text{for } j \in Z \setminus (Y \cup X)$$

(Here all  $y_j$  are in fact  $y_j(\vec{w}, \vec{x}_k)$ ).

## MLP – error function gradient

- ▶ If  $\sigma_j(\xi) = \frac{1}{1+e^{-\lambda_j \xi}}$  for all  $j \in Z$ , then

$$\sigma'_j(\xi_j) = \lambda_j y_j (1 - y_j)$$

and thus for all  $j \in Z \setminus X$ :

$$\frac{\partial E_k}{\partial y_j} = y_j - d_{kj} \quad \text{for } j \in Y$$

$$\frac{\partial E_k}{\partial y_j} = \sum_{r \in \vec{j}} \frac{\partial E_k}{\partial y_r} \cdot \lambda_r y_r (1 - y_r) \cdot w_{rj} \quad \text{for } j \in Z \setminus (Y \cup X)$$

- ▶ If  $\sigma_j(\xi) = a \cdot \tanh(b \cdot \xi_j)$  for all  $j \in Z$ , then

$$\sigma'_j(\xi_j) = \frac{b}{a} (a - y_j)(a + y_j)$$

# MLP – computing the gradient

Compute  $\frac{\partial E}{\partial w_{ji}} = \sum_{k=1}^p \frac{\partial E_k}{\partial w_{ji}}$  as follows:

Initialize  $\mathcal{E}_{ji} := 0$

(By the end of the computation:  $\mathcal{E}_{ji} = \frac{\partial E}{\partial w_{ji}}$ )

For every  $k = 1, \dots, p$  do:

- 1. forward pass:** compute  $y_j = y_j(\vec{w}, \vec{x}_k)$  for all  $j \in Z$
- 2. backward pass:** compute  $\frac{\partial E_k}{\partial y_j}$  for all  $j \in Z$  using *backpropagation* (see the next slide!)
- 3.** compute  $\frac{\partial E_k}{\partial w_{ji}}$  for all  $w_{ji}$  using

$$\frac{\partial E_k}{\partial w_{ji}} := \frac{\partial E_k}{\partial y_j} \cdot \sigma'_j(\xi_j) \cdot y_i$$

- 4.**  $\mathcal{E}_{ji} := \mathcal{E}_{ji} + \frac{\partial E_k}{\partial w_{ji}}$

The resulting  $\mathcal{E}_{ji}$  equals  $\frac{\partial E}{\partial w_{ji}}$ .

# MLP – backpropagation

Compute  $\frac{\partial E_k}{\partial y_j}$  for all  $j \in Z$  as follows:

- ▶ if  $j \in Y$ , then  $\frac{\partial E_k}{\partial y_j} = y_j - d_{kj}$
- ▶ if  $j \in Z \setminus Y \cup X$ , then assuming that  $j$  is in the  $\ell$ -th layer and assuming that  $\frac{\partial E_k}{\partial y_r}$  has already been computed for all neurons in the  $\ell + 1$ -st layer, compute

$$\frac{\partial E_k}{\partial y_j} = \sum_{r \in j^{\rightarrow}} \frac{\partial E_k}{\partial y_r} \cdot \sigma'_r(\xi_r) \cdot w_{rj}$$

(This works because all neurons of  $r \in j^{\rightarrow}$  belong to the  $\ell + 1$ -st layer.)

# Complexity of the batch algorithm

Computation of  $\frac{\partial E}{\partial w_{ji}}(\vec{w}^{(t-1)})$  stops in time linear in the size of the network plus the size of the training set.

(assuming unit cost of operations including computation of  $\sigma'_r(\xi_r)$  for given  $\xi_r$ )

**Proof sketch:** The algorithm does the following  $p$  times:

1. forward pass, i.e. computes  $y_j(\vec{w}, \vec{x}_k)$
2. backpropagation, i.e. computes  $\frac{\partial E_k}{\partial y_j}$
3. computes  $\frac{\partial E_k}{\partial w_{ji}}$  and adds it to  $\mathcal{E}_{ji}$  (a constant time operation in the unit cost framework)

The steps 1. - 3. take linear time.

Note that the speed of convergence of the gradient descent cannot be estimated ...

# MLP – learning algorithm

## Online algorithm:

The algorithm computes a sequence of weight vectors  $\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \dots$

- ▶ weights in  $\vec{w}^{(0)}$  are randomly initialized to values close to 0
- ▶ in the step  $t + 1$  (here  $t = 0, 1, 2 \dots$ ), weights  $\vec{w}^{(t+1)}$  are computed as follows:

$$w_{ji}^{(t+1)} = w_{ji}^{(t)} + \Delta w_{ji}^{(t)}$$

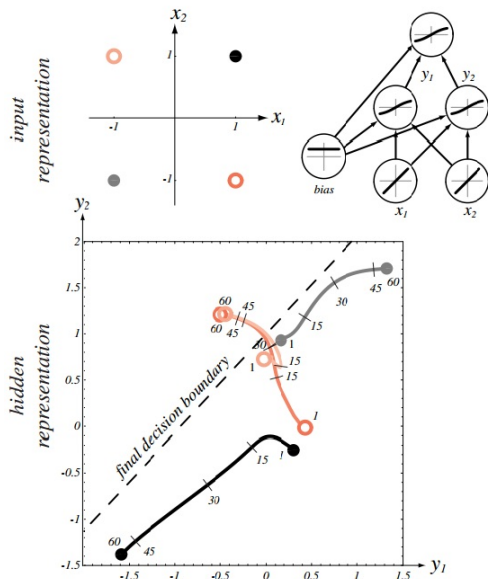
where

$$\Delta w_{ji}^{(t)} = -\varepsilon(t) \cdot \frac{\partial E_k}{\partial w_{ji}}(w_{ji}^{(t)})$$

is the *weight update* of  $w_{ji}$  in the step  $t + 1$  and  $0 < \varepsilon(t) \leq 1$   
is the *learning rate* in the step  $t + 1$ .

There are other variants determined by selection of the training examples used for the error computation (more on this later).

# Illustration of the gradient descent – XOR

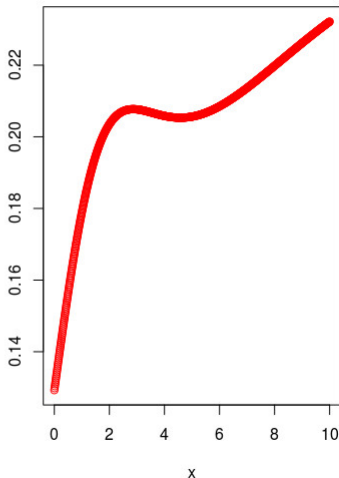
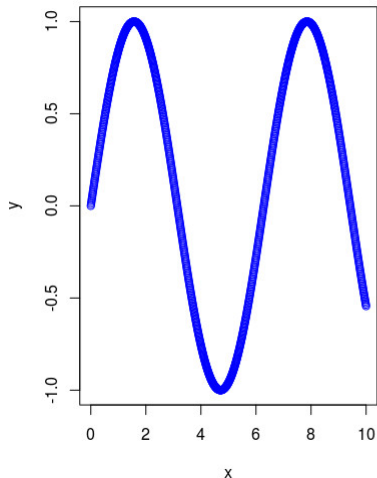


Source: Pattern Classification (2nd Edition); Richard O. Duda, Peter E. Hart, David G. Stork



# Animation (sin(x)), network 1-5-1)

One iteration:



10 iterations: