# **Chapter ML:VI**

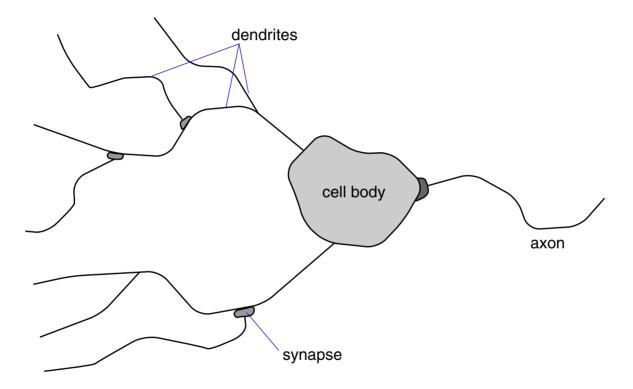
#### VI. Neural Networks

- Perceptron Learning
- □ Gradient Descent
- □ Multilayer Perceptron
- ☐ Radial Basis Functions

ML:VI-1 Neural Networks © STEIN 2005-2018

The Biological Model

#### Simplified model of a neuron:



ML:VI-2 Neural Networks ©STEIN 2005-2018

The Biological Model (continued)

#### Neuron characteristics:

- The numerous dendrites of a neuron serve as its input channels for electrical signals.
- At particular contact points between the dendrites, the so-called synapses, electrical signals can be initiated.
- A synapse can initiate signals of different strengths, where the strength is encoded by the frequency of a pulse train.
- □ The cell body of a neuron accumulates the incoming signals.
- If a particular stimulus threshold is exceeded, the cell body generates a signal, which is output via the axon.
- □ The processing of the signals is unidirectional. (from left to right in the figure)

ML:VI-3 Neural Networks © STEIN 2005-2018

History

- 1943 Warren McCulloch and Walter Pitts present a model of the neuron.
- 1949 Donald Hebb postulates a new learning paradigm: reinforcement only for active neurons. (those neurons that are involved in a decision process)
- 1958 Frank Rosenblatt develops the perceptron model.
- 1962 Rosenblatt proves the perceptron convergence theorem.
- 1969 Marvin Minsky and Seymour Papert publish a book on the limitations of the perceptron model.

1970

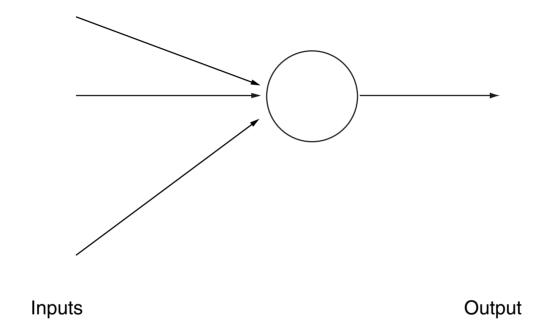
:

1985

1986 David Rumelhart and James McClelland present the multilayer perceptron.

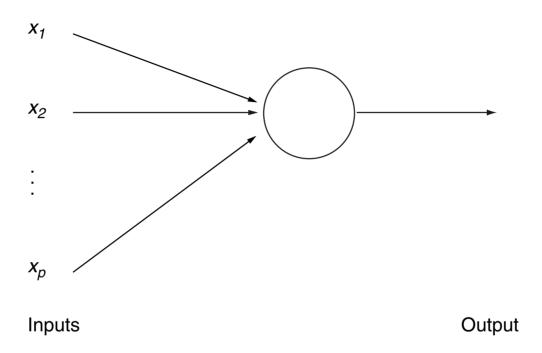
ML:VI-4 Neural Networks © STEIN 2005-2018

The Perceptron of Rosenblatt [1958]



ML:VI-5 Neural Networks ©STEIN 2005-2018

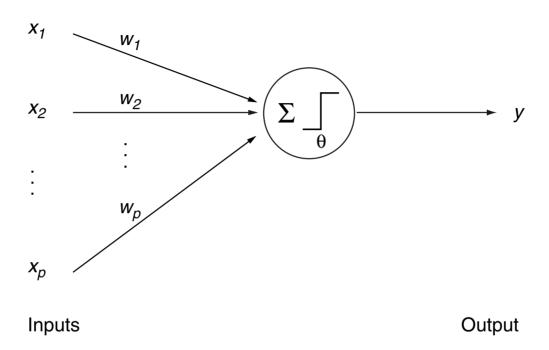
The Perceptron of Rosenblatt [1958]



$$x_j, w_j \in \mathbf{R}, \quad j = 1 \dots p$$

ML:VI-6 Neural Networks © STEIN 2005-2018

The Perceptron of Rosenblatt [1958]



$$x_j, w_j \in \mathbf{R}, \quad j = 1 \dots p$$

ML:VI-7 Neural Networks ©STEIN 2005-2018

#### Remarks:

☐ The perceptron is a "feed forward system".

ML:VI-8 Neural Networks ©STEIN 2005-2018

Specification of Classification Problems [ML Introduction]

#### Characterization of the model (model world):

- $\square$  X is a set of feature vectors, also called feature space.  $X \subseteq \mathbb{R}^p$
- $\Box$   $C = \{0, 1\}$  is a set of classes.  $C = \{-1, 1\}$  in the regression setting.
- $\neg c: X \to C$  is the ideal classifier for X. c is approximated by y (perceptron).
- $D = \{(\mathbf{x}_1, c(\mathbf{x}_1)), \dots, (\mathbf{x}_n, c(\mathbf{x}_n))\} \subseteq X \times C \text{ is a set of examples.}$

How could the hypothesis space H look like?

ML:VI-9 Neural Networks ©STEIN 2005-2018

Computation in the Perceptron [Regression]

If 
$$\sum_{j=1}^p w_j x_j \geq \theta$$
 then  $y(\mathbf{x}) = 1$ , and

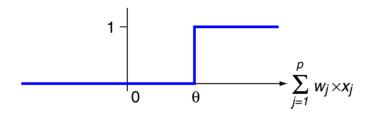
if 
$$\sum_{j=1}^p w_j x_j < \theta$$
 then  $y(\mathbf{x}) = 0$ .

ML:VI-10 Neural Networks ©STEIN 2005-2018

Computation in the Perceptron [Regression]

If 
$$\sum_{j=1}^p w_j x_j \geq \theta$$
 then  $y(\mathbf{x}) = 1$ , and

if 
$$\sum_{j=1}^p w_j x_j < \theta$$
 then  $y(\mathbf{x}) = 0$ .



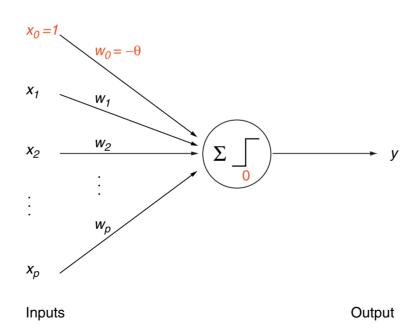
where 
$$\sum_{j=1}^{p} w_j x_j = \mathbf{w}^T \mathbf{x}$$
. (or other notations for the scalar product)

 $\rightarrow$  A hypothesis is determined by  $\theta, w_1, \dots, w_p$ .

ML:VI-11 Neural Networks ©STEIN 2005-2018

Computation in the Perceptron (continued)

$$y(\mathbf{x}) = \textit{heaviside}(\sum_{j=1}^p w_j x_j - \theta)$$
 
$$= \textit{heaviside}(\sum_{j=0}^p w_j x_j) \quad \text{with } w_0 = -\theta, \ x_0 = 1$$



 $\rightarrow$  A hypothesis is determined by  $\mathbf{w}_0, w_1, \dots, w_p$ .

ML:VI-12 Neural Networks © STEIN 2005-2018

#### Remarks:

- If the weight vector is extended by  $w_0 = -\theta$ , and, if the feature vectors are extended by the constant feature  $x_0 = 1$ , the learning algorithm gets a canonical form. Implementations of neural networks introduce this extension often implicitly.
- Be careful with regard to the dimensionality of the weight vector: it is always denoted as where, regardless of the fact whether the  $w_0$ -dimension, with  $w_0 = -\theta$ , is included.
- ☐ The function *heaviside* is named after the mathematician Oliver Heaviside. [Heaviside: step function O. Heaviside]

ML:VI-13 Neural Networks © STEIN 2005-2018

 $return(\mathbf{w})$ 

Weight Adaptation [IGD Algorithm]

Algorithm: PT**Perceptron Training** DInput: Training examples  $(\mathbf{x}, c(\mathbf{x}))$  with  $|\mathbf{x}| = p + 1$ ,  $c(\mathbf{x}) \in \{0, 1\}$ . Learning rate, a small positive constant.  $\eta$ Internal: y(D)Set of y(x)-values computed from the elements x in D given some w. Output: Weight vector.  $\mathbf{w}$  $PT(D, \eta)$ initialize\_random\_weights( $\mathbf{w}$ ), t=02. REPEAT 3. t = t + 14.  $(\mathbf{x}, c(\mathbf{x})) = random\_select(D)$  $error = c(\mathbf{x}) - heaviside(\mathbf{w}^T\mathbf{x})$  //  $c(\mathbf{x}) \in \{0,1\}$ , heaviside  $\in \{0,1\}$ , error  $\in \{0,1,-1\}$ 5. 6.  $\Delta \mathbf{w} = \eta \cdot \textit{error} \cdot \mathbf{x}$ 7.  $\mathbf{w} = \mathbf{w} + \Delta \mathbf{w}$ **UNTIL**(convergence(D, y(D)) OR  $t > t_{max}$ )

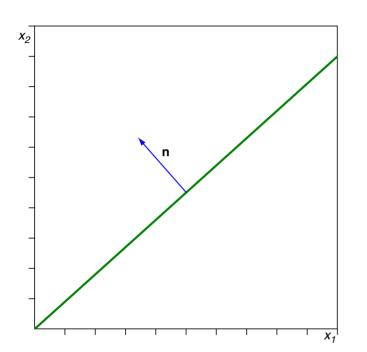
ML:VI-14 Neural Networks ©STEIN 2005-2018

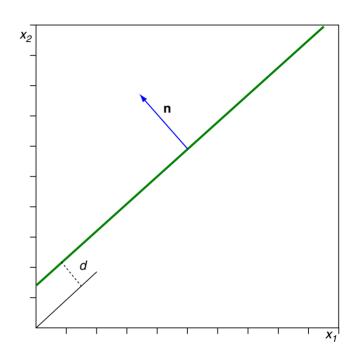
#### Remarks:

- $\Box$  The variable t denotes the time. At each point in time the learning algorithm gets an example presented and, as a consequence, may adapt the weight vector.
- The weight adaptation rule compares the true class  $c(\mathbf{x})$  (the ground truth) to the class computed by the perceptron. In case of a wrong classification of a feature vector  $\mathbf{x}$ , *error* is either -1 or +1, regardless of the exact numeric difference between  $c(\mathbf{x})$  and  $\mathbf{w}^T\mathbf{x}$ .
- $\supset y(D)$  is the set of  $y(\mathbf{x})$ -values given  $\mathbf{w}$  for the elements  $\mathbf{x}$  in D.

ML:VI-15 Neural Networks © STEIN 2005-2018

Weight Adaptation: Illustration in Input Space



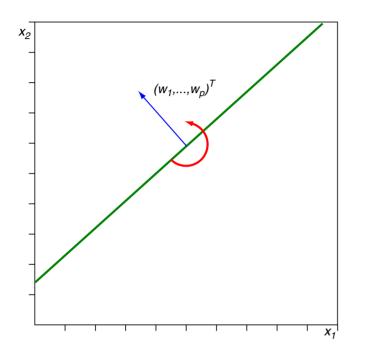


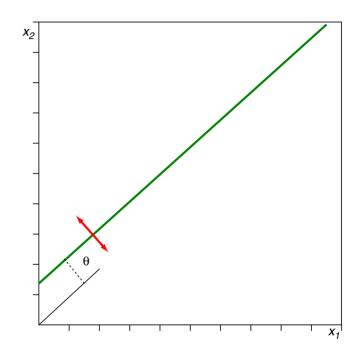
Definition of an (affine) hyperplane:  $L = \{\mathbf{x} \mid \mathbf{n}^T\mathbf{x} = d\}$  [Wikipedia]

- f n denotes a normal vector that is perpendicular to the hyperplane L.
- $\Box$  If  $||\mathbf{n}|| = 1$  then  $|\mathbf{n}^T \mathbf{x} d|$  gives the distance of any point  $\mathbf{x}$  to L.
- $\Box$  If  $sgn(\mathbf{n}^T\mathbf{x_1} d) = sgn(\mathbf{n}^T\mathbf{x_2} d)$ , then  $\mathbf{x_1}$  and  $\mathbf{x_2}$  lie on the same side of the hyperplane.

ML:VI-16 Neural Networks © STEIN 2005-2018

Weight Adaptation: Illustration in Input Space (continued)





Definition of an (affine) hyperplane:  $\mathbf{w}^T\mathbf{x} = 0 \iff \sum_{j=1}^p w_j x_j = \theta = -w_0$ 

(hyperplane definition as before, with notation taken from the classification problem setting)

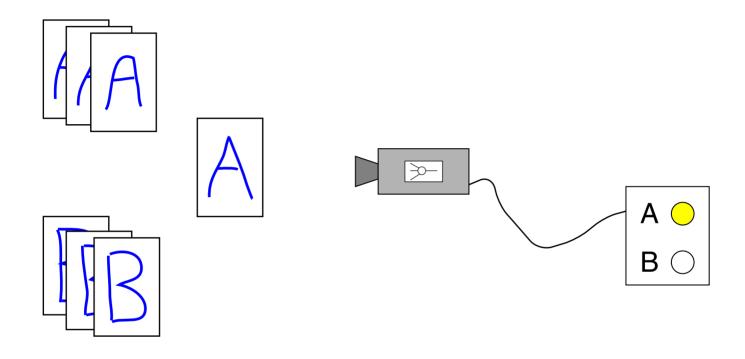
ML:VI-17 Neural Networks ©STEIN 2005-2018

#### Remarks:

- $\Box$  A perceptron defines a hyperplane that is perpendicular (= normal) to  $(w_1,\ldots,w_p)^T$ .
- $\Box$  The set of possible weight vectors  $\mathbf{w} = (w_0, w_1, \dots, w_p)^T$  form the hypothesis space H.
- □ Weight adaptation means learning, and the shown learning paradigm is supervised.
- $\Box$  For the weight adaptation in Line 6–7 of the *PT* Algorithm, note that if some  $x_j$  is zero,  $\Delta w_j$  will be zero as well. Keyword: Hebbian learning [Hebb 1949]
- Note that here (and in the following illustrations) the hyperplane movement is not the result of solving a regression problem in the (p+1)-dimensional input-output-space, where the residuals are to be minimized. Instead, the PT Algorithm takes each missclassified example as a trigger to correct the hyperplane's normal vector—without taking the effect on the other residuals into account.

ML:VI-18 Neural Networks ©STEIN 2005-2018

#### Example



- □ The examples are presented to the perceptron.
- □ The perceptron computes a value that is interpreted as class label.

ML:VI-19 Neural Networks © STEIN 2005-2018

Example (continued)

#### Encoding:

- The encoding of the examples is based on expressive features such as the number of line crossings, most acute angle, longest line, etc.
- $\Box$  The class label,  $c(\mathbf{x})$ , is encoded as a number. Examples from A are labeled with 1, examples from B are labeled with 0.

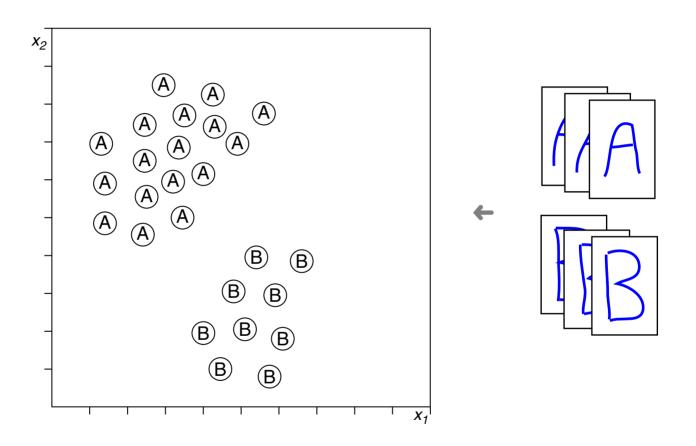
$$\begin{pmatrix} x_{1_1} \\ x_{1_2} \\ \vdots \\ x_{1_p} \end{pmatrix} \cdots \begin{pmatrix} x_{k_1} \\ x_{k_2} \\ \vdots \\ x_{k_p} \end{pmatrix} \cdots \begin{pmatrix} x_{l_1} \\ x_{l_2} \\ \vdots \\ x_{l_p} \end{pmatrix} \cdots \begin{pmatrix} x_{m_1} \\ x_{m_2} \\ \vdots \\ x_{m_p} \end{pmatrix}$$

$$Class  $A \simeq c(\mathbf{x}) = 1$ 

$$Class  $B \simeq c(\mathbf{x}) = 0$$$$$

ML:VI-20 Neural Networks © STEIN 2005-2018

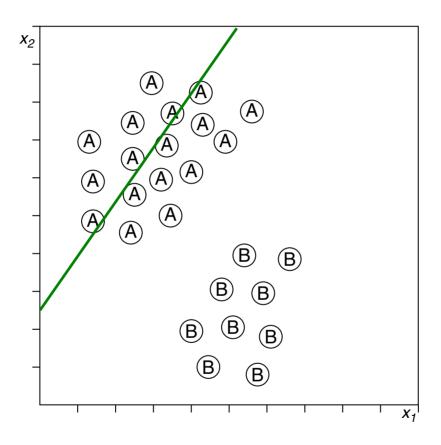
Example: Illustration in Input Space [PT Algorithm]



A possible configuration of encoded objects in the feature space X.

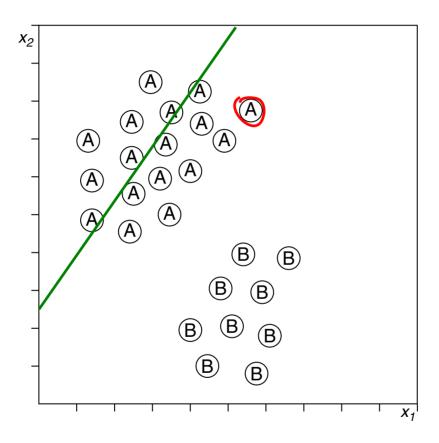
ML:VI-21 Neural Networks © STEIN 2005-2018

Example: Illustration in Input Space [PT Algorithm]



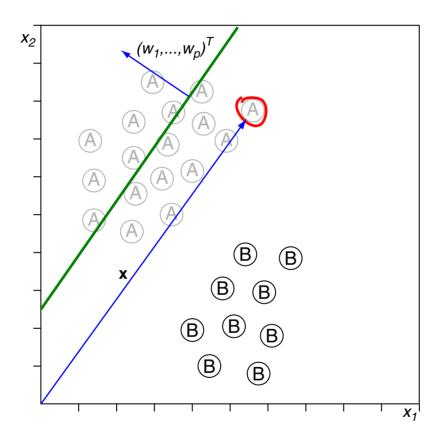
ML:VI-22 Neural Networks ©STEIN 2005-2018

Example: Illustration in Input Space [PT Algorithm]



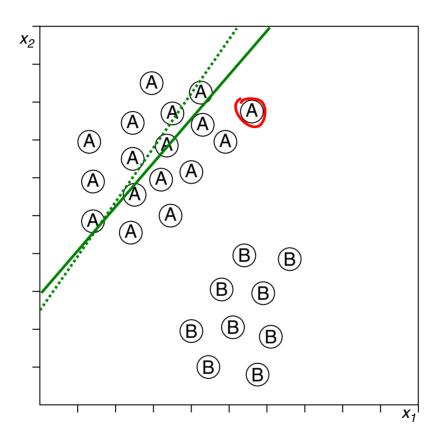
ML:VI-23 Neural Networks ©STEIN 2005-2018

Example: Illustration in Input Space [PT Algorithm]



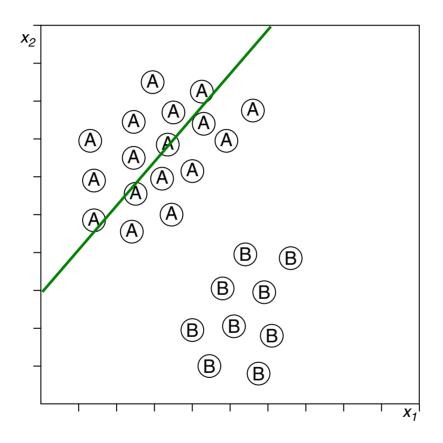
ML:VI-24 Neural Networks ©STEIN 2005-2018

Example: Illustration in Input Space [PT Algorithm]



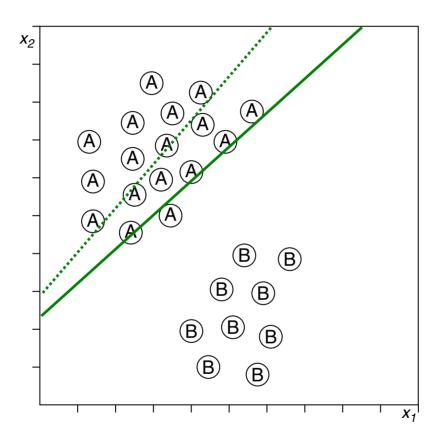
ML:VI-25 Neural Networks ©STEIN 2005-2018

Example: Illustration in Input Space [PT Algorithm]



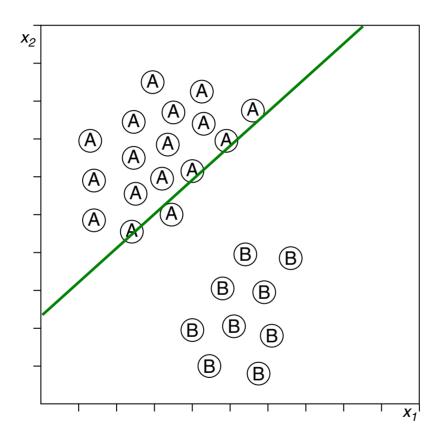
ML:VI-26 Neural Networks ©STEIN 2005-2018

Example: Illustration in Input Space [PT Algorithm]



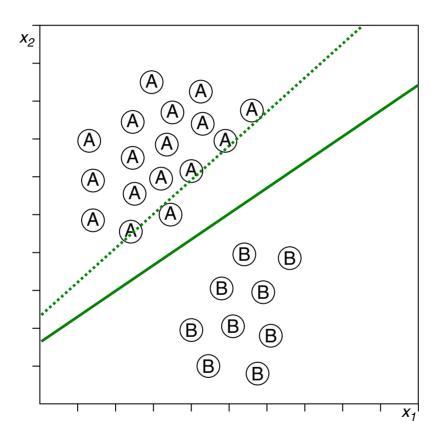
ML:VI-27 Neural Networks ©STEIN 2005-2018

Example: Illustration in Input Space [PT Algorithm]



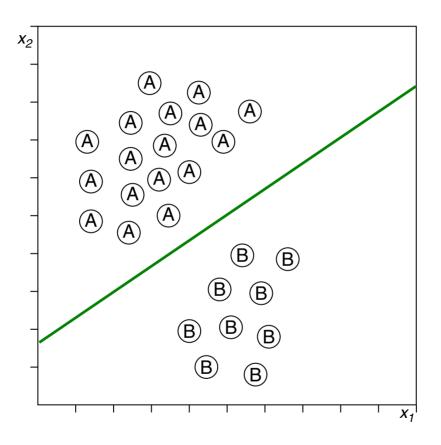
ML:VI-28 Neural Networks ©STEIN 2005-2018

Example: Illustration in Input Space [PT Algorithm]



ML:VI-29 Neural Networks ©STEIN 2005-2018

Example: Illustration in Input Space [PT Algorithm]



ML:VI-30 Neural Networks ©STEIN 2005-2018

Perceptron Convergence Theorem

#### Questions:

- 1. Which kind of learning tasks can be addressed with the functions in the hypothesis space *H*?
- 2. Can the PT Algorithm construct such a function for a given task?

ML:VI-31 Neural Networks © STEIN 2005-2018

Perceptron Convergence Theorem

#### Questions:

- 1. Which kind of learning tasks can be addressed with the functions in the hypothesis space *H*?
- 2. Can the *PT* Algorithm construct such a function for a given task?

#### Theorem 1 (Perceptron Convergence [Rosenblatt 1962])

Let  $X_0$  and  $X_1$  be two finite sets with vectors of the form  $\mathbf{x} = (1, x_1, \dots, x_p)^T$ , let  $X_1 \cap X_0 = \emptyset$ , and let  $\widehat{\mathbf{w}}$  define a separating hyperplane with respect to  $X_0$  and  $X_1$ . Moreover, let D be a set of examples of the form  $(\mathbf{x}, 0)$ ,  $\mathbf{x} \in X_0$  and  $(\mathbf{x}, 1)$ ,  $\mathbf{x} \in X_1$ . Then holds:

If the examples in D are processed with the PT Algorithm, the constructed weight vector  $\mathbf{w}$  will converge within a finite number of iterations.

ML:VI-32 Neural Networks © STEIN 2005-2018

Perceptron Convergence Theorem: Proof

#### **Preliminaries:**

The sets  $X_1$  and  $X_0$  are separated by a hyperplane  $\widehat{\mathbf{w}}$ . The proof requires that for all  $\mathbf{x} \in X_1$  the inequality  $\widehat{\mathbf{w}}^T \mathbf{x} > 0$  holds. This condition is always fulfilled, as the following consideration shows.

Let  $\mathbf{x}' \in X_1$  with  $\widehat{\mathbf{w}}^T \mathbf{x}' = 0$ . Since  $X_0$  is finite, the members  $\mathbf{x} \in X_0$  have a minimum positive distance  $\delta$  with regard to the hyperplane  $\widehat{\mathbf{w}}$ . Hence,  $\widehat{\mathbf{w}}$  can be moved by  $\frac{\delta}{2}$  towards  $X_0$ , resulting in a new hyperplane  $\widehat{\mathbf{w}}'$  that still fulfills  $(\widehat{\mathbf{w}}')^T \mathbf{x} < 0$  for all  $\mathbf{x} \in X_0$ , but that now also fulfills  $(\widehat{\mathbf{w}}')^T \mathbf{x} > 0$  for all  $\mathbf{x} \in X_1$ .

ML:VI-33 Neural Networks © STEIN 2005-2018

Perceptron Convergence Theorem: Proof

#### **Preliminaries:**

- The sets  $X_1$  and  $X_0$  are separated by a hyperplane  $\widehat{\mathbf{w}}$ . The proof requires that for all  $\mathbf{x} \in X_1$  the inequality  $\widehat{\mathbf{w}}^T \mathbf{x} > 0$  holds. This condition is always fulfilled, as the following consideration shows.
  - Let  $\mathbf{x}' \in X_1$  with  $\widehat{\mathbf{w}}^T \mathbf{x}' = 0$ . Since  $X_0$  is finite, the members  $\mathbf{x} \in X_0$  have a minimum positive distance  $\delta$  with regard to the hyperplane  $\widehat{\mathbf{w}}$ . Hence,  $\widehat{\mathbf{w}}$  can be moved by  $\frac{\delta}{2}$  towards  $X_0$ , resulting in a new hyperplane  $\widehat{\mathbf{w}}'$  that still fulfills  $(\widehat{\mathbf{w}}')^T \mathbf{x} < 0$  for all  $\mathbf{x} \in X_0$ , but that now also fulfills  $(\widehat{\mathbf{w}}')^T \mathbf{x} > 0$  for all  $\mathbf{x} \in X_1$ .
- □ By defining  $X' = X_1 \cup \{-\mathbf{x} \mid \mathbf{x} \in X_0\}$ , the searched w fulfills  $\mathbf{w}^T \mathbf{x} > 0$  for all  $\mathbf{x} \in X'$ . Then, with  $c(\mathbf{x}) = 1$  for all  $\mathbf{x} \in X'$ , *error*  $\in \{0, 1\}$  (instead of  $\{0, 1, -1\}$ ). [PT Algorithm, Line 5]

ML:VI-34 Neural Networks © STEIN 2005-2018

Perceptron Convergence Theorem: Proof

#### **Preliminaries:**

- The sets  $X_1$  and  $X_0$  are separated by a hyperplane  $\widehat{\mathbf{w}}$ . The proof requires that for all  $\mathbf{x} \in X_1$  the inequality  $\widehat{\mathbf{w}}^T \mathbf{x} > 0$  holds. This condition is always fulfilled, as the following consideration shows.
  - Let  $\mathbf{x}' \in X_1$  with  $\widehat{\mathbf{w}}^T \mathbf{x}' = 0$ . Since  $X_0$  is finite, the members  $\mathbf{x} \in X_0$  have a minimum positive distance  $\delta$  with regard to the hyperplane  $\widehat{\mathbf{w}}$ . Hence,  $\widehat{\mathbf{w}}$  can be moved by  $\frac{\delta}{2}$  towards  $X_0$ , resulting in a new hyperplane  $\widehat{\mathbf{w}}'$  that still fulfills  $(\widehat{\mathbf{w}}')^T \mathbf{x} < 0$  for all  $\mathbf{x} \in X_0$ , but that now also fulfills  $(\widehat{\mathbf{w}}')^T \mathbf{x} > 0$  for all  $\mathbf{x} \in X_1$ .
- By defining  $X' = X_1 \cup \{-\mathbf{x} \mid \mathbf{x} \in X_0\}$ , the searched  $\mathbf{w}$  fulfills  $\mathbf{w}^T \mathbf{x} > 0$  for all  $\mathbf{x} \in X'$ . Then, with  $c(\mathbf{x}) = 1$  for all  $\mathbf{x} \in X'$ ,  $error \in \{0, 1\}$  (instead of  $\{0, 1, -1\}$ ). [PT Algorithm, Line 5]
- □ The *PT* Algorithm performs a number of iterations, where  $\mathbf{w}(t)$  denotes the weight vector for iteration t, which form the basis for the weight vector  $\mathbf{w}(t+1)$ .  $\mathbf{x}(t) \in X'$  denotes the feature vector chosen in round t. The first (and randomly chosen) weight vector is denoted as  $\mathbf{w}(0)$ .

ML:VI-35 Neural Networks © STEIN 2005-2018

Perceptron Convergence Theorem: Proof

#### **Preliminaries:**

- The sets  $X_1$  and  $X_0$  are separated by a hyperplane  $\widehat{\mathbf{w}}$ . The proof requires that for all  $\mathbf{x} \in X_1$  the inequality  $\widehat{\mathbf{w}}^T \mathbf{x} > 0$  holds. This condition is always fulfilled, as the following consideration shows.
  - Let  $\mathbf{x}' \in X_1$  with  $\widehat{\mathbf{w}}^T \mathbf{x}' = 0$ . Since  $X_0$  is finite, the members  $\mathbf{x} \in X_0$  have a minimum positive distance  $\delta$  with regard to the hyperplane  $\widehat{\mathbf{w}}$ . Hence,  $\widehat{\mathbf{w}}$  can be moved by  $\frac{\delta}{2}$  towards  $X_0$ , resulting in a new hyperplane  $\widehat{\mathbf{w}}'$  that still fulfills  $(\widehat{\mathbf{w}}')^T \mathbf{x} < 0$  for all  $\mathbf{x} \in X_0$ , but that now also fulfills  $(\widehat{\mathbf{w}}')^T \mathbf{x} > 0$  for all  $\mathbf{x} \in X_1$ .
- By defining  $X' = X_1 \cup \{-\mathbf{x} \mid \mathbf{x} \in X_0\}$ , the searched  $\mathbf{w}$  fulfills  $\mathbf{w}^T \mathbf{x} > 0$  for all  $\mathbf{x} \in X'$ . Then, with  $c(\mathbf{x}) = 1$  for all  $\mathbf{x} \in X'$ ,  $error \in \{0, 1\}$  (instead of  $\{0, 1, -1\}$ ). [PT Algorithm, Line 5]
- □ The *PT* Algorithm performs a number of iterations, where  $\mathbf{w}(t)$  denotes the weight vector for iteration t, which form the basis for the weight vector  $\mathbf{w}(t+1)$ .  $\mathbf{x}(t) \in X'$  denotes the feature vector chosen in round t. The first (and randomly chosen) weight vector is denoted as  $\mathbf{w}(0)$ .
- □ Recall the Cauchy-Schwarz inequality:  $||\mathbf{a}||^2 \cdot ||\mathbf{b}||^2 \ge (\mathbf{a}^T \mathbf{b})^2$ , where  $||\mathbf{x}|| := \sqrt{\mathbf{x}^T \mathbf{x}}$  denotes the Euclidean norm.

ML:VI-36 Neural Networks © STEIN 2005-2018

Perceptron Convergence Theorem: Proof (continued)

### Line of argument:

- (a) We state a lower bound for how much  $||\mathbf{w}||$  must change from its initial value after n iterations (to become a separating hyperplane). The derivation of this lower bound exploits the presupposed linear separability of  $X_0$  and  $X_1$ .
- (b) We state an upper bound for how much  $||\mathbf{w}||$  can change from its initial value after n iterations. The derivation of this upper bound exploits the finiteness of  $X_0$  and  $X_1$ , which in turn guarantees the existence of an upper bound for the norm of the maximum feature vector.
- (c) We observe that the lower bound grows quadratically in n, whereas the upper bound grows linearly. From the relation "lower bound < upper bound" we derive a finite upper bound for n.

ML:VI-37 Neural Networks © STEIN 2005-2018

Perceptron Convergence Theorem: Proof (continued)

1. The <u>PT</u> Algorithm computes in iteration t the scalar product  $\mathbf{w}(t)^T \mathbf{x}(t)$ . If classified correctly,  $\mathbf{w}(t)^T \mathbf{x}(t) > 0$  and  $\mathbf{w}$  is unchanged. Otherwise,  $\mathbf{w}(t+1) = \mathbf{w}(t) + \eta \cdot \mathbf{x}(t)$  [Line 5-7].

ML:VI-38 Neural Networks ©STEIN 2005-2018

Perceptron Convergence Theorem: Proof (continued)

- 1. The <u>PT</u> Algorithm computes in iteration t the scalar product  $\mathbf{w}(t)^T \mathbf{x}(t)$ . If classified correctly,  $\mathbf{w}(t)^T \mathbf{x}(t) > 0$  and  $\mathbf{w}$  is unchanged. Otherwise,  $\mathbf{w}(t+1) = \mathbf{w}(t) + \eta \cdot \mathbf{x}(t)$  [Line 5-7].
- 2. A sequence of n incorrectly classified feature vectors, (x(t)), along with the weight adaptation,  $\mathbf{w}(t+1) = \mathbf{w}(t) + \eta \cdot \mathbf{x}(t)$ , results in the series  $\mathbf{w}(n)$ :

$$\mathbf{w}(1) = \mathbf{w}(0) + \eta \cdot \mathbf{x}(0)$$

$$\mathbf{w}(2) = \mathbf{w}(1) + \eta \cdot \mathbf{x}(1) = \mathbf{w}(0) + \eta \cdot \mathbf{x}(0) + \eta \cdot \mathbf{x}(1)$$

$$\vdots$$

$$\mathbf{w}(n) = \mathbf{w}(0) + \eta \cdot \mathbf{x}(0) + \dots + \eta \cdot \mathbf{x}(n-1)$$

ML:VI-39 Neural Networks ©STEIN 2005-2018

Perceptron Convergence Theorem: Proof (continued)

- 1. The <u>PT</u> Algorithm computes in iteration t the scalar product  $\mathbf{w}(t)^T \mathbf{x}(t)$ . If classified correctly,  $\mathbf{w}(t)^T \mathbf{x}(t) > 0$  and  $\mathbf{w}$  is unchanged. Otherwise,  $\mathbf{w}(t+1) = \mathbf{w}(t) + \eta \cdot \mathbf{x}(t)$  [Line 5-7].
- 2. A sequence of n incorrectly classified feature vectors, (x(t)), along with the weight adaptation,  $\mathbf{w}(t+1) = \mathbf{w}(t) + \eta \cdot \mathbf{x}(t)$ , results in the series  $\mathbf{w}(n)$ :

$$\mathbf{w}(1) = \mathbf{w}(0) + \eta \cdot \mathbf{x}(0)$$

$$\mathbf{w}(2) = \mathbf{w}(1) + \eta \cdot \mathbf{x}(1) = \mathbf{w}(0) + \eta \cdot \mathbf{x}(0) + \eta \cdot \mathbf{x}(1)$$

$$\vdots$$

$$\mathbf{w}(n) = \mathbf{w}(0) + \eta \cdot \mathbf{x}(0) + \ldots + \eta \cdot \mathbf{x}(n-1)$$

3. The hyperplane defined by  $\widehat{\mathbf{w}}$  separates  $X_1$  and  $X_0$ :  $\forall \mathbf{x} \in X' : \widehat{\mathbf{w}}^T \mathbf{x} > 0$  Let  $\delta := \min_{\mathbf{x} \in X'} \widehat{\mathbf{w}}^T \mathbf{x}$ . Observe that  $\delta > 0$  holds.

ML:VI-40 Neural Networks ©STEIN 2005-2018

### Perceptron Convergence Theorem: Proof (continued)

- 1. The <u>PT</u> Algorithm computes in iteration t the scalar product  $\mathbf{w}(t)^T \mathbf{x}(t)$ . If classified correctly,  $\mathbf{w}(t)^T \mathbf{x}(t) > 0$  and  $\mathbf{w}$  is unchanged. Otherwise,  $\mathbf{w}(t+1) = \mathbf{w}(t) + \eta \cdot \mathbf{x}(t)$  [Line 5-7].
- 2. A sequence of n incorrectly classified feature vectors, (x(t)), along with the weight adaptation,  $\mathbf{w}(t+1) = \mathbf{w}(t) + \eta \cdot \mathbf{x}(t)$ , results in the series  $\mathbf{w}(n)$ :

$$\mathbf{w}(1) = \mathbf{w}(0) + \eta \cdot \mathbf{x}(0)$$

$$\mathbf{w}(2) = \mathbf{w}(1) + \eta \cdot \mathbf{x}(1) = \mathbf{w}(0) + \eta \cdot \mathbf{x}(0) + \eta \cdot \mathbf{x}(1)$$

$$\vdots$$

$$\mathbf{w}(n) = \mathbf{w}(0) + \eta \cdot \mathbf{x}(0) + \dots + \eta \cdot \mathbf{x}(n-1)$$

- 3. The hyperplane defined by  $\widehat{\mathbf{w}}$  separates  $X_1$  and  $X_0$ :  $\forall \mathbf{x} \in X' : \widehat{\mathbf{w}}^T \mathbf{x} > 0$  Let  $\delta := \min_{\mathbf{x} \in X'} \widehat{\mathbf{w}}^T \mathbf{x}$ . Observe that  $\delta > 0$  holds.
- 4. Analyze the scalar product of  $\mathbf{w}(n)$  and  $\hat{\mathbf{w}}$ :

$$\widehat{\mathbf{w}}^T \mathbf{w}(n) = \widehat{\mathbf{w}}^T \mathbf{w}(0) + \eta \cdot \widehat{\mathbf{w}}^T \mathbf{x}(0) + \ldots + \eta \cdot \widehat{\mathbf{w}}^T \mathbf{x}(n-1)$$

$$\Rightarrow \widehat{\mathbf{w}}^T \mathbf{w}(n) \geq \widehat{\mathbf{w}}^T \mathbf{w}(0) + n\eta \delta \geq 0 \quad (\text{for } n \geq n_0 \text{ with sufficiently large } n_0 \in \mathbf{N})$$

$$\Rightarrow (\widehat{\mathbf{w}}^T \mathbf{w}(n))^2 \geq (\widehat{\mathbf{w}}^T \mathbf{w}(0) + n\eta \delta)^2$$

ML:VI-41 Neural Networks ©STEIN 2005-2018

### Perceptron Convergence Theorem: Proof (continued)

- 1. The <u>PT</u> Algorithm computes in iteration t the scalar product  $\mathbf{w}(t)^T \mathbf{x}(t)$ . If classified correctly,  $\mathbf{w}(t)^T \mathbf{x}(t) > 0$  and  $\mathbf{w}$  is unchanged. Otherwise,  $\mathbf{w}(t+1) = \mathbf{w}(t) + \eta \cdot \mathbf{x}(t)$  [Line 5-7].
- 2. A sequence of n incorrectly classified feature vectors, (x(t)), along with the weight adaptation,  $\mathbf{w}(t+1) = \mathbf{w}(t) + \eta \cdot \mathbf{x}(t)$ , results in the series  $\mathbf{w}(n)$ :

$$\mathbf{w}(1) = \mathbf{w}(0) + \eta \cdot \mathbf{x}(0)$$

$$\mathbf{w}(2) = \mathbf{w}(1) + \eta \cdot \mathbf{x}(1) = \mathbf{w}(0) + \eta \cdot \mathbf{x}(0) + \eta \cdot \mathbf{x}(1)$$

$$\vdots$$

$$\mathbf{w}(n) = \mathbf{w}(0) + \eta \cdot \mathbf{x}(0) + \dots + \eta \cdot \mathbf{x}(n-1)$$

- 3. The hyperplane defined by  $\widehat{\mathbf{w}}$  separates  $X_1$  and  $X_0$ :  $\forall \mathbf{x} \in X' : \widehat{\mathbf{w}}^T \mathbf{x} > 0$ Let  $\delta := \min_{\mathbf{x} \in X'} \widehat{\mathbf{w}}^T \mathbf{x}$ . Observe that  $\delta > 0$  holds.
- 4. Analyze the scalar product of  $\mathbf{w}(n)$  and  $\hat{\mathbf{w}}$ :

$$\widehat{\mathbf{w}}^T \mathbf{w}(n) = \widehat{\mathbf{w}}^T \mathbf{w}(0) + \eta \cdot \widehat{\mathbf{w}}^T \mathbf{x}(0) + \ldots + \eta \cdot \widehat{\mathbf{w}}^T \mathbf{x}(n-1)$$

$$\Rightarrow \widehat{\mathbf{w}}^T \mathbf{w}(n) \geq \widehat{\mathbf{w}}^T \mathbf{w}(0) + n\eta \delta \geq 0 \quad (\text{for } n \geq n_0 \text{ with sufficiently large } n_0 \in \mathbf{N})$$

$$\Rightarrow (\widehat{\mathbf{w}}^T \mathbf{w}(n))^2 \geq (\widehat{\mathbf{w}}^T \mathbf{w}(0) + n\eta \delta)^2$$

5. Apply the Cauchy-Schwarz inequality:

$$||\widehat{\mathbf{w}}||^2 \cdot ||\mathbf{w}(n)||^2 \geq (\widehat{\mathbf{w}}^T \mathbf{w}(0) + n\eta \delta)^2 \quad \Rightarrow \quad ||\mathbf{w}(n)||^2 \geq \frac{(\widehat{\mathbf{w}}^T \mathbf{w}(0) + n\eta \delta)^2}{||\widehat{\mathbf{w}}||^2}$$

ML:VI-42 Neural Networks ©STEIN 2005-2018

Perceptron Convergence Theorem: Proof (continued)

6. Consider again the weight adaptation  $\mathbf{w}(t+1) = \mathbf{w}(t) + \eta \cdot \mathbf{x}(t)$ :

$$\begin{aligned} ||\mathbf{w}(t+1)||^2 &= ||\mathbf{w}(t) + \eta \cdot \mathbf{x}(t)||^2 \\ &= (\mathbf{w}(t) + \eta \cdot \mathbf{x}(t))^T (\mathbf{w}(t) + \eta \cdot \mathbf{x}(t)) \\ &= \mathbf{w}(t)^T \mathbf{w}(t) + \eta^2 \cdot \mathbf{x}(t)^T \mathbf{x}(t) + 2\eta \cdot \mathbf{w}(t)^T \mathbf{x}(t) \\ &\leq ||\mathbf{w}(t)||^2 + ||\eta \cdot \mathbf{x}(t)||^2 \quad (\text{since } \mathbf{w}(t)^T \mathbf{x}(t) < 0) \end{aligned}$$

ML:VI-43 Neural Networks ©STEIN 2005-2018

Perceptron Convergence Theorem: Proof (continued)

6. Consider again the weight adaptation  $\mathbf{w}(t+1) = \mathbf{w}(t) + \eta \cdot \mathbf{x}(t)$ :

$$\begin{aligned} ||\mathbf{w}(t+1)||^2 &= ||\mathbf{w}(t) + \eta \cdot \mathbf{x}(t)||^2 \\ &= (\mathbf{w}(t) + \eta \cdot \mathbf{x}(t))^T (\mathbf{w}(t) + \eta \cdot \mathbf{x}(t)) \\ &= \mathbf{w}(t)^T \mathbf{w}(t) + \eta^2 \cdot \mathbf{x}(t)^T \mathbf{x}(t) + 2\eta \cdot \mathbf{w}(t)^T \mathbf{x}(t) \\ &\leq ||\mathbf{w}(t)||^2 + ||\eta \cdot \mathbf{x}(t)||^2 \quad (\text{since } \mathbf{w}(t)^T \mathbf{x}(t) < 0) \end{aligned}$$

7. Consider the series  $\mathbf{w}(n)$  from Step 2:

$$||\mathbf{w}(n)||^{2} \leq ||\mathbf{w}(n-1)||^{2} + ||\eta \cdot \mathbf{x}(n-1)||^{2}$$

$$\leq ||\mathbf{w}(n-2)||^{2} + ||\eta \cdot \mathbf{x}(n-2)||^{2} + ||\eta \cdot \mathbf{x}(n-1)||^{2}$$

$$\leq ||\mathbf{w}(0)||^{2} + ||\eta \cdot \mathbf{x}(0)||^{2} + \dots + ||\eta \cdot \mathbf{x}(n-1)||^{2}$$

$$= ||\mathbf{w}(0)||^{2} + \sum_{j=0}^{n-1} ||\eta \cdot \mathbf{x}(j)||^{2}$$

ML:VI-44 Neural Networks ©STEIN 2005-2018

Perceptron Convergence Theorem: Proof (continued)

6. Consider again the weight adaptation  $\mathbf{w}(t+1) = \mathbf{w}(t) + \eta \cdot \mathbf{x}(t)$ :

$$\begin{aligned} ||\mathbf{w}(t+1)||^2 &= ||\mathbf{w}(t) + \eta \cdot \mathbf{x}(t)||^2 \\ &= (\mathbf{w}(t) + \eta \cdot \mathbf{x}(t))^T (\mathbf{w}(t) + \eta \cdot \mathbf{x}(t)) \\ &= \mathbf{w}(t)^T \mathbf{w}(t) + \eta^2 \cdot \mathbf{x}(t)^T \mathbf{x}(t) + 2\eta \cdot \mathbf{w}(t)^T \mathbf{x}(t) \\ &\leq ||\mathbf{w}(t)||^2 + ||\eta \cdot \mathbf{x}(t)||^2 \quad (\text{since } \mathbf{w}(t)^T \mathbf{x}(t) < 0) \end{aligned}$$

7. Consider the series w(n) from Step 2:

$$||\mathbf{w}(n)||^{2} \leq ||\mathbf{w}(n-1)||^{2} + ||\eta \cdot \mathbf{x}(n-1)||^{2}$$

$$\leq ||\mathbf{w}(n-2)||^{2} + ||\eta \cdot \mathbf{x}(n-2)||^{2} + ||\eta \cdot \mathbf{x}(n-1)||^{2}$$

$$\leq ||\mathbf{w}(0)||^{2} + ||\eta \cdot \mathbf{x}(0)||^{2} + \dots + ||\eta \cdot \mathbf{x}(n-1)||^{2}$$

$$= ||\mathbf{w}(0)||^{2} + \sum_{j=0}^{n-1} ||\eta \cdot \mathbf{x}(j)||^{2}$$

8. With  $\varepsilon := \max_{\mathbf{x} \in X'} ||\mathbf{x}||^2$  follows  $||\mathbf{w}(n)||^2 \le ||\mathbf{w}(0)||^2 + n\eta^2 \varepsilon$ 

ML:VI-45 Neural Networks ©STEIN 2005-2018

Perceptron Convergence Theorem: Proof (continued)

9. Both inequalities (see Step 5 and Step 8) must be fulfilled:

$$\begin{aligned} ||\mathbf{w}(n)||^2 & \geq \frac{(\widehat{\mathbf{w}}^T \mathbf{w}(0) + n\eta \delta)^2}{||\widehat{\mathbf{w}}||^2} \quad \text{and} \quad ||\mathbf{w}(n)||^2 & \leq ||\mathbf{w}(0)||^2 + n\eta^2 \varepsilon \\ & \Rightarrow \quad \frac{(\widehat{\mathbf{w}}^T \mathbf{w}(0) + n\eta \delta)^2}{||\widehat{\mathbf{w}}||^2} & \leq \quad ||\mathbf{w}(n)||^2 & \leq ||\mathbf{w}(0)||^2 + n\eta^2 \varepsilon \\ & \Rightarrow \quad \frac{(\widehat{\mathbf{w}}^T \mathbf{w}(0) + n\eta \delta)^2}{||\widehat{\mathbf{w}}||^2} & \leq \quad ||\mathbf{w}(0)||^2 + n\eta^2 \varepsilon \end{aligned}$$

$$\mathbf{Set} \ \mathbf{w}(0) = \mathbf{0} : \qquad \Rightarrow \qquad \frac{n^2 \eta^2 \delta^2}{||\widehat{\mathbf{w}}||^2} & \leq \quad n\eta^2 \varepsilon \\ \Leftrightarrow \qquad n & \leq \quad \frac{\varepsilon}{\delta^2} \cdot ||\widehat{\mathbf{w}}||^2 \end{aligned}$$

ML:VI-46 Neural Networks ©STEIN 2005-2018

Perceptron Convergence Theorem: Proof (continued)

9. Both inequalities (see Step 5 and Step 8) must be fulfilled:

$$\begin{aligned} ||\mathbf{w}(n)||^2 & \geq \frac{(\widehat{\mathbf{w}}^T \mathbf{w}(0) + n\eta\delta)^2}{||\widehat{\mathbf{w}}||^2} \quad \text{and} \quad ||\mathbf{w}(n)||^2 & \leq ||\mathbf{w}(0)||^2 + n\eta^2 \varepsilon \\ \Rightarrow \quad \frac{(\widehat{\mathbf{w}}^T \mathbf{w}(0) + n\eta\delta)^2}{||\widehat{\mathbf{w}}||^2} & \leq \quad ||\mathbf{w}(n)||^2 & \leq ||\mathbf{w}(0)||^2 + n\eta^2 \varepsilon \\ \Rightarrow \quad \frac{(\widehat{\mathbf{w}}^T \mathbf{w}(0) + n\eta\delta)^2}{||\widehat{\mathbf{w}}||^2} & \leq \quad ||\mathbf{w}(0)||^2 + n\eta^2 \varepsilon \end{aligned}$$
 Set  $\mathbf{w}(0) = \mathbf{0}$ :  $\Rightarrow \quad \frac{n^2 \eta^2 \delta^2}{||\widehat{\mathbf{w}}||^2} & \leq \quad n\eta^2 \varepsilon$   $\Leftrightarrow \quad n \quad \leq \quad \frac{\varepsilon}{\delta^2} \cdot ||\widehat{\mathbf{w}}||^2$ 

→ The *PT* Algorithm terminates within a finite number of iterations.

Observe: 
$$\frac{(\widehat{\mathbf{w}}^T\mathbf{w}(0) + n\eta\delta)^2}{||\widehat{\mathbf{w}}||^2} \ \in \ \Theta(n^2) \quad \text{ and } \quad ||\mathbf{w}(0)||^2 + n\eta^2\varepsilon \ \in \ \Theta(n)$$

ML:VI-47 Neural Networks ©STEIN 2005-2018

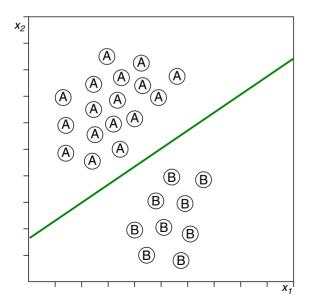
Perceptron Convergence Theorem: Discussion

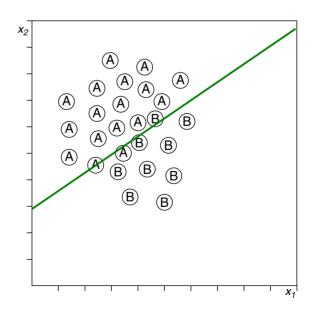
- □ If a separating hyperplane between  $X_0$  and  $X_1$  exists, the <u>PT Algorithm</u> will converge. If no such hyperplane exists, convergence cannot be guaranteed.
- A separating hyperplane can be found in polynomial time with linear programming. The PT Algorithm, however, may require an exponential number of iterations.

ML:VI-48 Neural Networks ©STEIN 2005-2018

## Perceptron Convergence Theorem: Discussion

- □ If a separating hyperplane between  $X_0$  and  $X_1$  exists, the <u>PT Algorithm</u> will converge. If no such hyperplane exists, convergence cannot be guaranteed.
- A separating hyperplane can be found in polynomial time with linear programming. The PT Algorithm, however, may require an exponential number of iterations.
- Classification problems with noise (right-hand side) are problematic:





ML:VI-49 Neural Networks ©STEIN 2005-2018

### Classification Error

Gradient descent considers the true error (better: the hyperplane distance) and will converge even if  $X_1$  and  $X_0$  cannot be separated by a hyperplane. However, this convergence process is of an asymptotic nature and no finite iteration bound can be stated.

Gradient descent applies the so-called delta rule, which will be derived in the following. The delta rule forms the basis of the backpropagation algorithm.

ML:VI-50 Neural Networks © STEIN 2005-2018

### Classification Error

Gradient descent considers the true error (better: the hyperplane distance) and will converge even if  $X_1$  and  $X_0$  cannot be separated by a hyperplane. However, this convergence process is of an asymptotic nature and no finite iteration bound can be stated.

Gradient descent applies the so-called delta rule, which will be derived in the following. The delta rule forms the basis of the backpropagation algorithm.

Consider the linear perceptron *without* a threshold function:

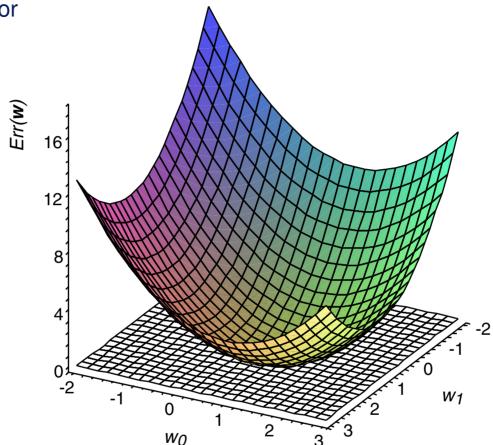
$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} = \sum_{j=0}^p w_j x_j$$
 [Heaviside]

The classification error  $Err(\mathbf{w})$  of a weight vector (= hypothesis)  $\mathbf{w}$  with regard to D can be defined as follows:

$$Err(\mathbf{w}) = \frac{1}{2} \sum_{(\mathbf{x}, c(\mathbf{x})) \in D} (c(\mathbf{x}) - y(\mathbf{x}))^2$$
 [Singleton error]

ML:VI-51 Neural Networks ©STEIN 2005-2018

Classification Error



The gradient of  $Err(\mathbf{w})$ ,  $\nabla Err(\mathbf{w})$ , defines the steepest ascent or descent:

$$\nabla \textit{Err}(\mathbf{w}) = \left(\frac{\partial \textit{Err}(\mathbf{w})}{\partial w_0}, \frac{\partial \textit{Err}(\mathbf{w})}{\partial w_1}, \cdots, \frac{\partial \textit{Err}(\mathbf{w})}{\partial w_p}\right)$$

ML:VI-52 Neural Networks © STEIN 2005-2018

## Weight Adaptation

$$\mathbf{w} \leftarrow \mathbf{w} + \Delta \mathbf{w}$$
 where  $\Delta \mathbf{w} = -\eta \nabla \textit{Err}(\mathbf{w})$ 

Componentwise (j = 0, ..., p) weight adaptation [PT Algorithm]:

$$w_j \leftarrow w_j + \Delta w_j$$
 where  $\Delta w_j = -\eta \frac{\partial}{\partial w_j} \textit{Err}(\mathbf{w})$ 

ML:VI-53 Neural Networks ©STEIN 2005-2018

## Weight Adaptation

$$\mathbf{w} \leftarrow \mathbf{w} + \Delta \mathbf{w}$$
 where  $\Delta \mathbf{w} = -\eta \nabla \textit{Err}(\mathbf{w})$ 

Componentwise (j = 0, ..., p) weight adaptation [PT Algorithm]:

$$w_j \leftarrow w_j + \Delta w_j$$
 where  $\Delta w_j = -\eta \frac{\partial}{\partial w_j} \textit{Err}(\mathbf{w})$ 

$$\frac{\partial}{\partial w_j} \textit{Err}(\mathbf{w}) \ = \ \frac{\partial}{\partial w_j} \, \frac{1}{2} \sum_{\substack{(\mathbf{x}, c(\mathbf{x})) \in D}} (c(\mathbf{x}) - y(\mathbf{x}))^2 = \frac{1}{2} \sum_{\substack{(\mathbf{x}, c(\mathbf{x})) \in D}} \frac{\partial}{\partial w_j} (c(\mathbf{x}) - y(\mathbf{x}))^2$$

ML:VI-54 Neural Networks ©STEIN 2005-2018

## Weight Adaptation

$$\mathbf{w} \leftarrow \mathbf{w} + \Delta \mathbf{w}$$
 where  $\Delta \mathbf{w} = -\eta \nabla \textit{Err}(\mathbf{w})$ 

Componentwise (j = 0, ..., p) weight adaptation [PT Algorithm]:

$$w_j \leftarrow w_j + \Delta w_j$$
 where  $\Delta w_j = -\eta \frac{\partial}{\partial w_j} \textit{Err}(\mathbf{w})$ 

$$\begin{split} \frac{\partial}{\partial w_j} \textit{Err}(\mathbf{w}) &= \frac{\partial}{\partial w_j} \frac{1}{2} \sum_{(\mathbf{x}, c(\mathbf{x})) \in D} (c(\mathbf{x}) - y(\mathbf{x}))^2 = \frac{1}{2} \sum_{(\mathbf{x}, c(\mathbf{x})) \in D} \frac{\partial}{\partial w_j} (c(\mathbf{x}) - y(\mathbf{x}))^2 \\ &= \frac{1}{2} \sum_{(\mathbf{x}, c(\mathbf{x})) \in D} 2(c(\mathbf{x}) - y(\mathbf{x})) \cdot \frac{\partial}{\partial w_j} (c(\mathbf{x}) - y(\mathbf{x})) \end{split}$$

ML:VI-55 Neural Networks ©STEIN 2005-2018

## Weight Adaptation

$$\mathbf{w} \leftarrow \mathbf{w} + \Delta \mathbf{w}$$
 where  $\Delta \mathbf{w} = -\eta \nabla \textit{Err}(\mathbf{w})$ 

Componentwise (j = 0, ..., p) weight adaptation [PT Algorithm]:

$$w_j \leftarrow w_j + \Delta w_j$$
 where  $\Delta w_j = -\eta \frac{\partial}{\partial w_j} \textit{Err}(\mathbf{w})$ 

$$\begin{split} \frac{\partial}{\partial w_j} \textit{Err}(\mathbf{w}) &= \frac{\partial}{\partial w_j} \frac{1}{2} \sum_{(\mathbf{x}, c(\mathbf{x})) \in D} (c(\mathbf{x}) - y(\mathbf{x}))^2 = \frac{1}{2} \sum_{(\mathbf{x}, c(\mathbf{x})) \in D} \frac{\partial}{\partial w_j} (c(\mathbf{x}) - y(\mathbf{x}))^2 \\ &= \frac{1}{2} \sum_{(\mathbf{x}, c(\mathbf{x})) \in D} 2(c(\mathbf{x}) - y(\mathbf{x})) \cdot \frac{\partial}{\partial w_j} (c(\mathbf{x}) - y(\mathbf{x})) \\ &= \sum_{(\mathbf{x}, c(\mathbf{x})) \in D} (c(\mathbf{x}) - \mathbf{w}^T \mathbf{x}) \cdot \frac{\partial}{\partial w_j} (c(\mathbf{x}) - \mathbf{w}^T \mathbf{x}) \end{split}$$

ML:VI-56 Neural Networks ©STEIN 2005-2018

## Weight Adaptation

$$\mathbf{w} \leftarrow \mathbf{w} + \Delta \mathbf{w}$$
 where  $\Delta \mathbf{w} = -\eta \nabla \textit{Err}(\mathbf{w})$ 

Componentwise (j = 0, ..., p) weight adaptation [PT Algorithm]:

 $(\mathbf{x},c(\mathbf{x}))\in D$ 

$$\begin{split} \frac{\partial}{\partial w_j} \textit{Err}(\mathbf{w}) &= \frac{\partial}{\partial w_j} \frac{1}{2} \sum_{(\mathbf{x}, c(\mathbf{x})) \in D} (c(\mathbf{x}) - y(\mathbf{x}))^2 = \frac{1}{2} \sum_{(\mathbf{x}, c(\mathbf{x})) \in D} \frac{\partial}{\partial w_j} (c(\mathbf{x}) - y(\mathbf{x}))^2 \\ &= \frac{1}{2} \sum_{(\mathbf{x}, c(\mathbf{x})) \in D} 2(c(\mathbf{x}) - y(\mathbf{x})) \cdot \frac{\partial}{\partial w_j} (c(\mathbf{x}) - y(\mathbf{x})) \\ &= \sum_{(\mathbf{x}, c(\mathbf{x})) \in D} (c(\mathbf{x}) - \mathbf{w}^T \mathbf{x}) \cdot \frac{\partial}{\partial w_j} (c(\mathbf{x}) - \mathbf{w}^T \mathbf{x}) \\ &= \sum_{(\mathbf{x}, c(\mathbf{x})) \in D} (c(\mathbf{x}) - \mathbf{w}^T \mathbf{x}) (-x_j) \end{split}$$

ML:VI-57 Neural Networks

Weight Adaptation: Batch Gradient Descent [IGD Algorithm]

Algorithm: BGD Batch Gradient Descent Input: DTraining examples  $(\mathbf{x}, c(\mathbf{x}))$  with  $|\mathbf{x}| = p + 1$ ,  $c(\mathbf{x}) \in \{0, 1\}$ .  $(c(\mathbf{x}) \in \{-1, 1\})$ Learning rate, a small positive constant.  $\eta$ Internal: Set of y(x)-values computed from the elements x in D given some w. y(D)Output: Weight vector.  $\mathbf{w}$  $BGD(D, \eta)$ initialize\_random\_weights( $\mathbf{w}$ ), t=0REPEAT 3. t = t + 14.  $\Delta \mathbf{w} = 0$ 5. FOREACH  $(\mathbf{x}, c(\mathbf{x})) \in D$  DO  $error = c(\mathbf{x}) - \mathbf{w}^T \mathbf{x}$ 6.  $\Delta \mathbf{w} = \Delta \mathbf{w} + \eta \cdot \text{error} \cdot \mathbf{x}$ 7. 8. ENDDO 9.  $\mathbf{w} = \mathbf{w} + \Delta \mathbf{w}$ 

**UNTIL**(convergence(D, y(D)) OR  $t > t_{max}$ )

ML:VI-58 Neural Networks

 $return(\mathbf{w})$ 

10.

11.

#### Remarks:

- $\triangle \mathbf{w} \sim -\nabla \textit{Err}(\mathbf{w})$  (and not to "+") to descend to the minimum.
- Each BGD iteration "REPEAT ... UNTIL" corresponds to finding the direction of steepest error descent as  $\nabla \textit{Err}(\mathbf{w_t}) = \sum_{(\mathbf{x}, c(\mathbf{x})) \in D} \left( c(\mathbf{x}) \mathbf{w_t}^T \mathbf{x} \right) \cdot \mathbf{x}$   $\frac{\partial}{\partial \mathbf{w}} \sum_{\mathbf{x}, c(\mathbf{x}) \in D} \left( c(\mathbf{x}) \mathbf{w_t}^T \mathbf{x} \right)^2$  and updating  $\mathbf{w_t}$  by taking a step of length  $\eta$  in this direction.
- Using a constant step size  $\eta$  is can severely impair the speed of convergence. When taking the optimal step size  $\eta_t := \operatorname{argmin}_{\eta} \operatorname{Err}(\mathbf{w}_t - \eta \cdot \nabla \operatorname{Err}(\mathbf{w}_t))$  at each iteration t, it can be shown that gradient descent (merely) has a linear rate of convergence. [Meza 2010]
- As criterion for the *convergence* function may serve the global error, either quantified as the sum of the squared residuals,  $Err(\mathbf{w}_t)$ , or as the norm of the error gradient,  $||\nabla Err(\mathbf{w}_t)||$ , which are compared to some small positive bound  $\varepsilon$ .

ML:VI-59 Neural Networks ©STEIN 2005-2018

Weight Adaptation: Delta Rule

The weight adaptation in the BGD Algorithm is set-based: before modifying a weight component in w, the total error of all examples (the "batch") is computed.

Weight adaptation with regard to a *single* example  $(\mathbf{x}, c(\mathbf{x})) \in D$ :

$$\Delta \mathbf{w} = \eta \cdot (c(\mathbf{x}) - \mathbf{w}^T \mathbf{x}) \cdot \mathbf{x}$$

This adaptation rule is known under different names:

- □ delta rule
- Widrow-Hoff rule
- adaline rule
- least mean squares (LMS) rule

The classification error  $\textit{Err}_d(\mathbf{w})$  of a weight vector (= hypothesis)  $\mathbf{w}$  with regard to a *single* example  $d \in D$ ,  $d = (\mathbf{x}, c(\mathbf{x}))$ , is given as:

$$\textit{Err}_d(\mathbf{w}) = \frac{1}{2}(c(\mathbf{x}) - \mathbf{w}^T \mathbf{x})^2$$
 [Batch error]

ML:VI-60 Neural Networks © STEIN 2005-2018

IGD

Algorithm:

Weight Adaptation: Incremental Gradient Descent [Algorithms: LMS BGD PT]

Incremental Gradient Descent

Input: Training examples  $(\mathbf{x}, c(\mathbf{x}))$  with  $|\mathbf{x}| = p + 1$ ,  $c(\mathbf{x}) \in \{0, 1\}$ .  $(c(\mathbf{x}) \in \{-1, 1\})$ DLearning rate, a small positive constant.  $\eta$ Internal: y(D) Set of y(x)-values computed from the elements x in D given some w. Output: Weight vector.  $\mathbf{w}$  $IGD(D, \eta)$ initialize\_random\_weights( $\mathbf{w}$ ), t=02. REPEAT 3. t = t + 14. FOREACH  $(\mathbf{x}, c(\mathbf{x})) \in D$  DO 5.  $error = c(\mathbf{x}) - \mathbf{w}^T \mathbf{x}$ 6.  $\Delta \mathbf{w} = \eta \cdot \text{error} \cdot \mathbf{x}$ 7.  $\mathbf{w} = \mathbf{w} + \Delta \mathbf{w}$ 8. ENDDO **UNTIL**(convergence(D, y(D)) OR  $t > t_{max}$ )  $return(\mathbf{w})$ 10.

ML:VI-61 Neural Networks ©STEIN 2005-2018

#### Remarks:

- The classification error *Err* of incremental gradient descent is specific for each training example  $d \in D$ ,  $d = (\mathbf{x}, c(\mathbf{x}))$ :  $\textit{Err}_d(\mathbf{w}) = \frac{1}{2}(c(\mathbf{x}) \mathbf{w}^T\mathbf{x})^2$
- The sequence of incremental weight adaptations approximates the gradient descent of the batch approach. If  $\eta$  is chosen sufficiently small, this approximation can happen at arbitrary accuracy.
- ☐ The computation of the total error of batch gradient descent enables larger weight adaptation increments.
- Compared to batch gradient descent, the example-based weight adaptation of incremental gradient descent can better avoid getting stuck in a local minimum of the error function.
- □ Incremental gradient descent is also called *stochastic* gradient descent.

ML:VI-62 Neural Networks © STEIN 2005-2018

### Remarks (continued):

- □ When, as is done here, the residual sum squares, RSS, is chosen as error (loss) function, the incremental gradient descent algorithm [IGD] corresponds to the least mean squares algorithm [LMS].
- The incremental gradient descend algorithm [IGD] looks similar to the perceptron training algorithm [PT], since these algorithms differ only in the error computation (Line 5) where the latter applies the Heaviside function. However, this subtle syntactic difference is a significant conceptual difference, entailing a number of consequences:
  - Gradient descent is a regression approach and exploits the residua, which are provided by an error function of choice, and whose differential is evaluated to control the hyperplane movement.
  - The PT algorithm is not based on residuals (in the (p+1)-dimensional input-output-space) but refers to the input space only, where it simply evaluates the side of the hyperplane as a binary feature (correct side or not).
  - Provided linear separability, the PT algorithm will converge within a finite number of iterations, which, however, cannot be guaranteed for gradient descent.
  - Gradient descent will converge even if the data is not linearly separable.
  - Data sets can be constructed whose classes are linearly separable, but where gradient descent will not determine a hyperplane that classifies all examples correctly (whereas the *PT* Algorithm of course does).

ML:VI-63 Neural Networks © STEIN 2005-2018