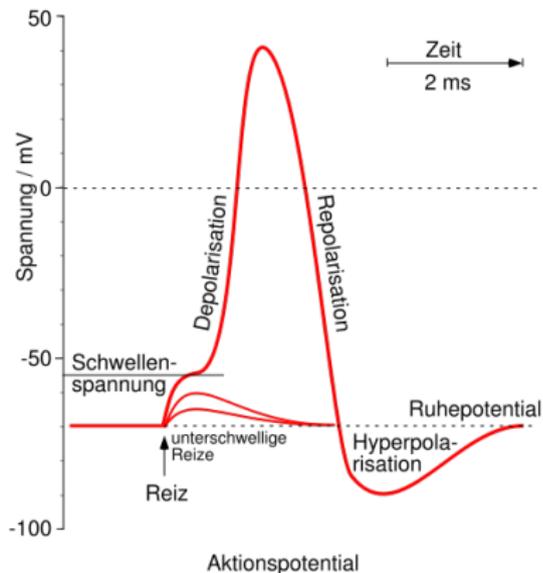
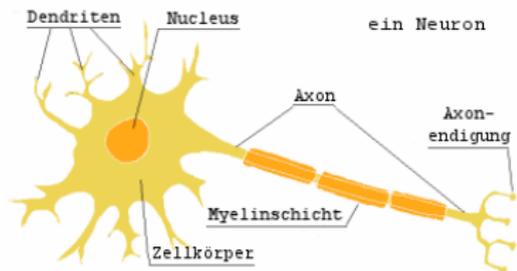


# Einführung in künstliche neuronale Netze

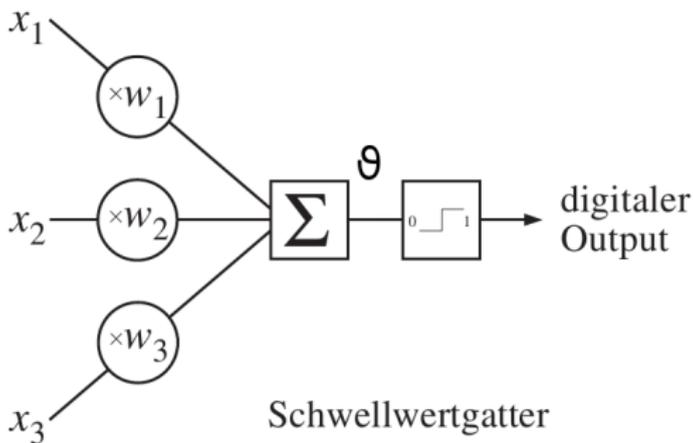
Stefan Fürtinger

Woche der Modellierung mit Mathematik  
5.–11. Februar 2012

# Ein biologisches Neron



# Ein künstliches Neuron

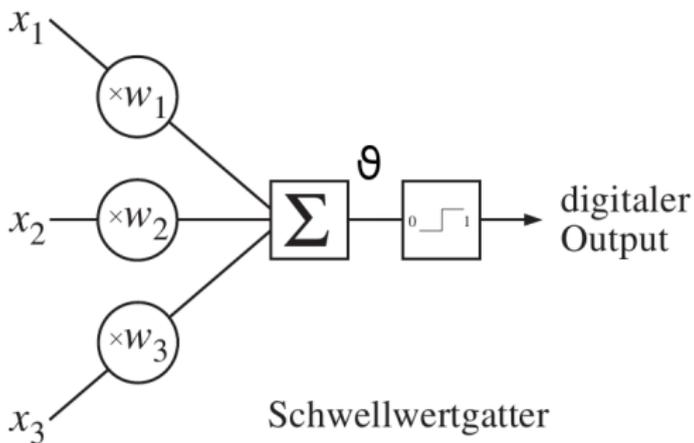


Warren McCulloch



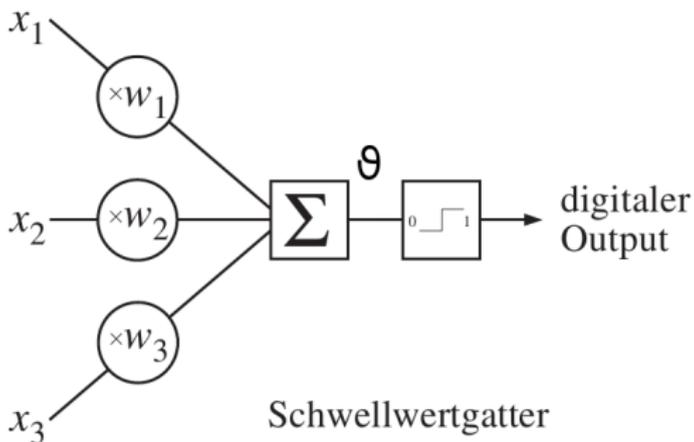
Walter Pitts

# Ein künstliches Neuron



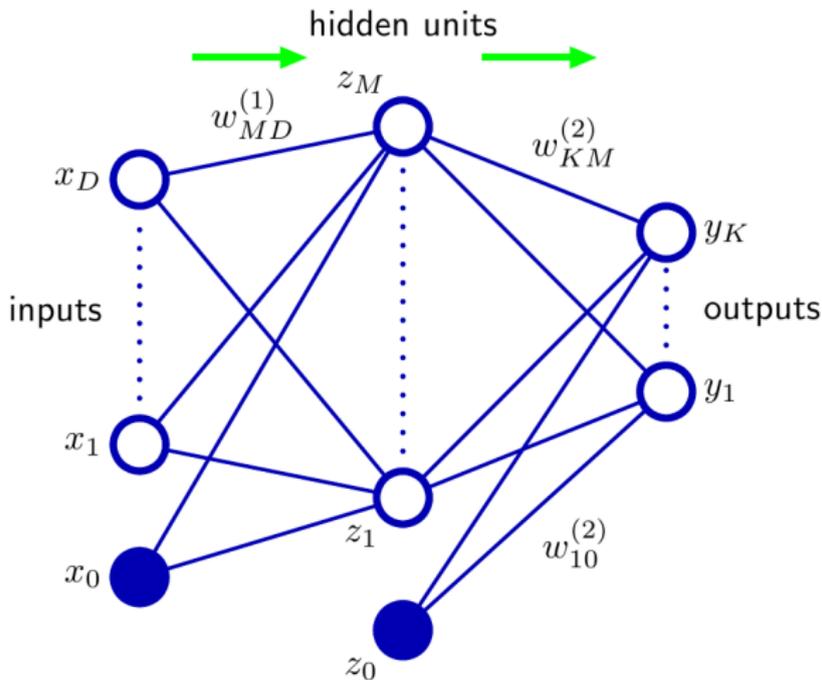
Einführung...

# Ein künstliches Neuron



<http://www.neuralesnetz.de/applets/NettoInput.html>

# Ein neuronales Netz



Terminologie eines Netzes

# Was bisher geschah...

1943 McCulloch-Pitts Neuron Das erste mathematische Modell  
eines Neurons. Heute bekannt als **Perceptron**

## Was bisher geschah...

1943 McCulloch-Pitts Neuron Das erste mathematische Modell  
eines Neurons. Heute bekannt als **Perceptron**

1951 Edmond und Minsky bauen den ersten Neurocomputer

## Was bisher geschah...

1943 McCulloch-Pitts Neuron Das erste mathematische Modell eines Neurons. Heute bekannt als **Perceptron**

1951 Edmond und Minsky bauen den ersten Neurocomputer

**Problem:** Wie **lernt** der Computer? Vorschläge?

## Was bisher geschah...

1943 McCulloch-Pitts Neuron Das erste mathematische Modell eines Neurons. Heute bekannt als **Perceptron**

1951 Edmond und Minsky bauen den ersten Neurocomputer

**Problem:** Wie **lernt** der Computer? Vorschläge?



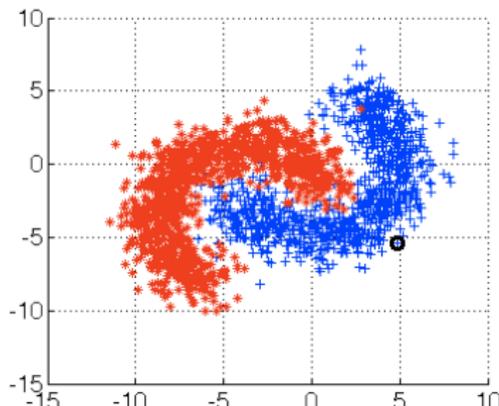
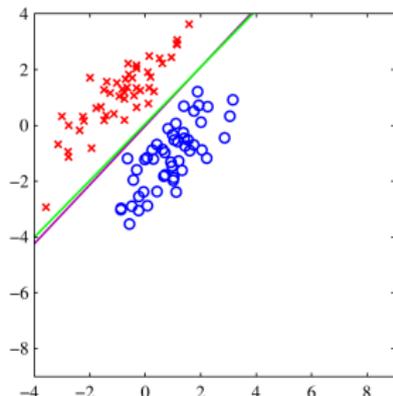
**Anpassung der Gewichte**

1959 Empory und Rosenblatt das Perceptron kann lernen →  
enthusiastische Forschung



# Was sonst noch geschah...

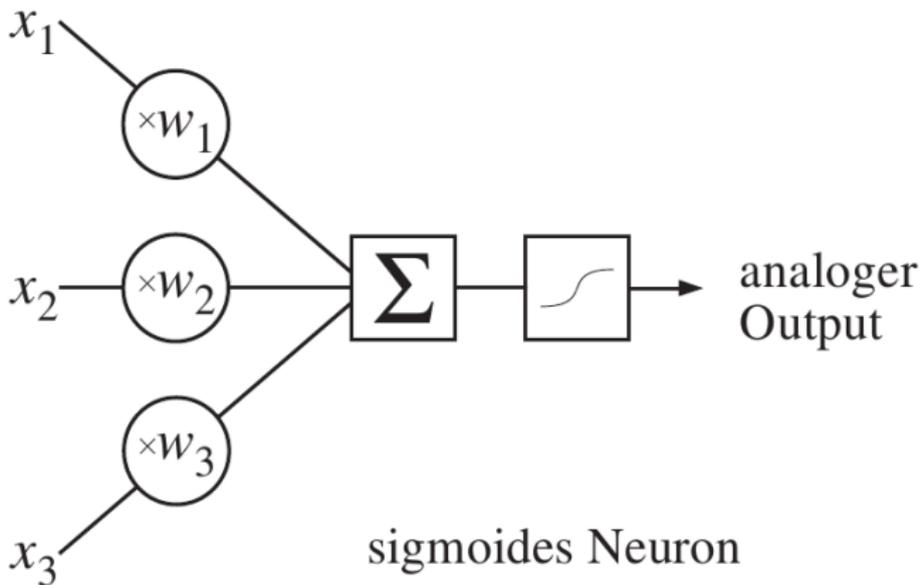
1969 Minsky und Papert Das Buch "Perceptrons" erscheint und dämpft die Euphorie...



## Satz

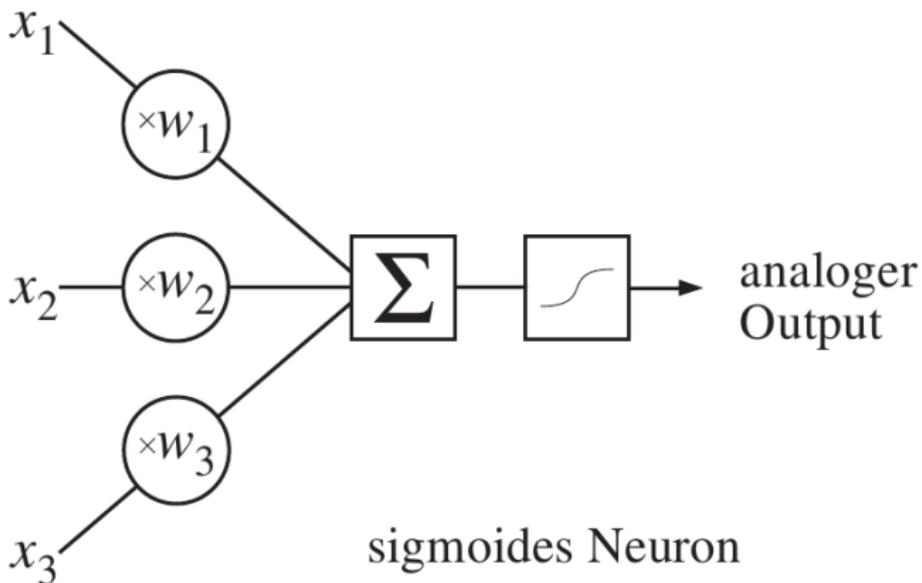
*Ein einschichtiges Netzwerk kann nur Lösungen linear separabler Probleme lernen*

# Ein Schritt zurück...



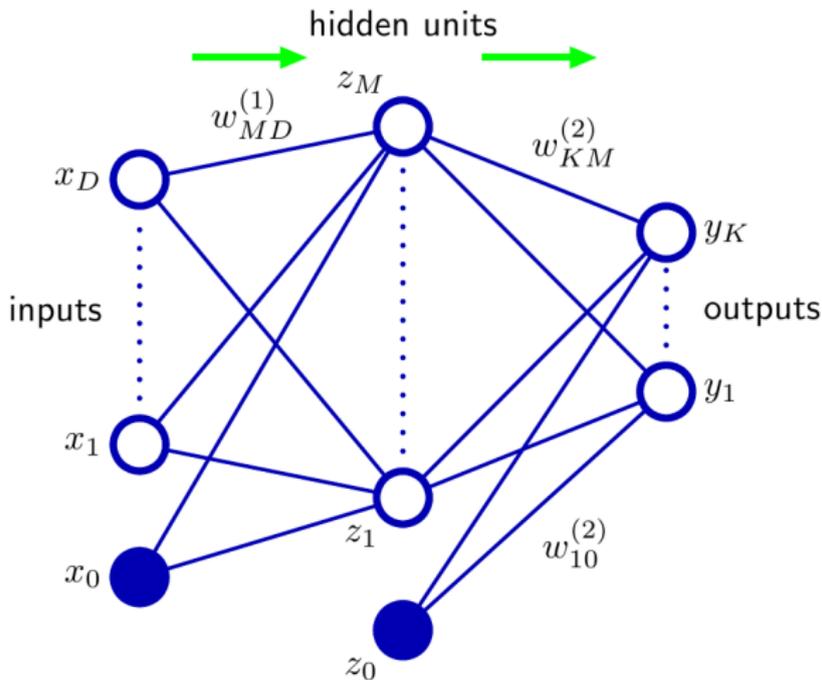
Eine neue Übergangsfunktion

# Ein Schritt zurück...



<http://www.neuronalesnetz.de/applets/ActivityFunction.html>

# ...ein Schritt vorwärts



Formale Beschreibung eines Netzwerks

# Ein Zweischichtiges Netzwerk

Input Variablen

$$x_1, \dots, x_D \quad \text{mit Bias } x_0$$

Der Netinput am  $j$ -ten Neuron des Hidden Layer ist

$$a_j = \sum_{i=1}^D w_{ji}^{(1)} x_i + w_{j0}^{(1)}$$

Der Aktivitätslevel des  $j$ -ten Neurons des Hidden Layer ist

$$z_j = h(a_j), \quad j = 1, \dots, M$$

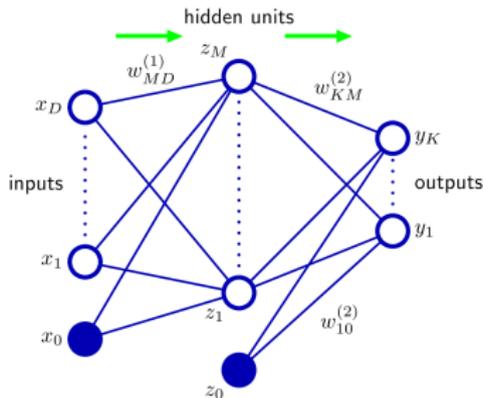
Der Aktivitätslevel am  $k$ -ten Neuron des Output Layer ist

$$b_k = \sum_{j=1}^M w_{kj}^{(2)} z_j + w_{k0}^{(2)}$$

Der Output des  $k$ -ten Neurons des Output Layer ist

$$y_k = g(b_k), \quad k = 1, \dots, K$$

# Formale Beschreibung des Netzwerks



Fasst man Gewichte und Inputs zusammen

$$\mathbf{w} = (w_{10}^{(1)}, \dots, w_{M0}^{(1)}, w_{11}^{(1)}, \dots, w_{M1}^{(1)}, \dots, w_{10}^{(2)}, \dots, w_{K0}^{(2)}, \dots)$$

$$\mathbf{x} = (x_0, \dots, x_D)$$

so erhält man

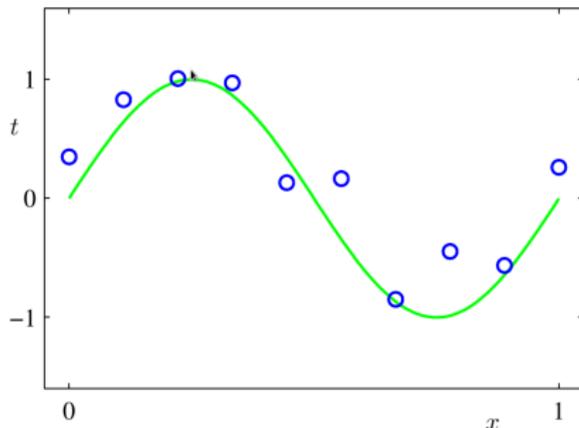
$$y_k(\mathbf{x}, \mathbf{w}) = g(b_k) = g \left( \sum_{j=0}^M w_{kj}^{(2)} h \left( \sum_{i=0}^D w_{ji}^{(1)} x_i \right) \right)$$

# Was tun mit $g$ und $h$ ?

Die Übertragungsfunktionen  $g$  und  $h$  sind stark **problemabhängig**

# Was tun mit $g$ und $h$ ?

Die Übertragungsfunktionen  $g$  und  $h$  sind stark **problemabhängig**



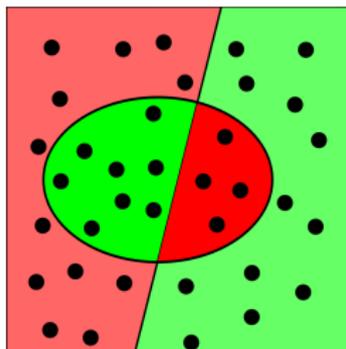
Regressionsproblem

Linear

$$h(a) = ca, \quad c \in \mathbb{R}$$

# Was tun mit $g$ und $h$ ?

Die Übertragungsfunktionen  $g$  und  $h$  sind stark **problemabhängig**



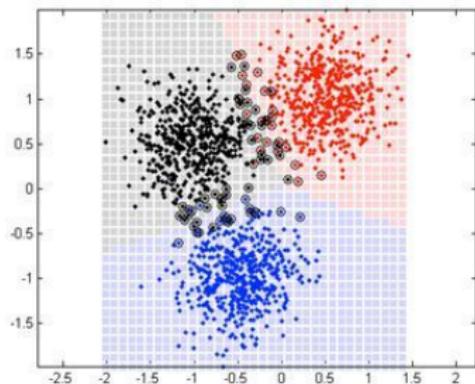
(Multiple) Binäre Klassifizierung

Sigmoid

$$h(a) = \frac{1}{1 + e^{-a}}$$

# Was tun mit $g$ und $h$ ?

Die Übertragungsfunktionen  $g$  und  $h$  sind stark **problemabhängig**



Mehrklassen Klassifizierung

Softmax

$$h(a_1, \dots, a_K) = \frac{e^{a_k}}{\sum_{j=1}^K e^{a_j}}$$

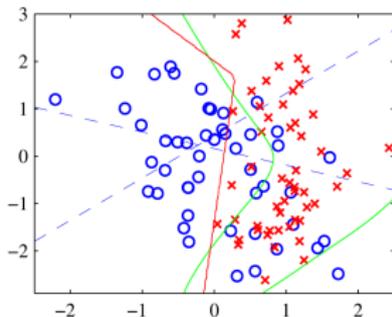
# Der nächste Schritt

Beginn 1980er Beweis des folgenden Resultats gelingt

## Satz

*Ein zweischichtiges neuronales Netzwerk mit linearem Output Neuron und sigmoidalen Hidden Neuronen kann jede beliebige stetige Funktion auf einem kompakten Gebiet mit beliebiger Genauigkeit approximieren, sofern die Anzahl der Neuronen im Hidden Layer nur groß genug ist.*

Damit ging's nun Schlag auf Schlag...



Beliebig komplexe Klassifizierungsprobleme und mehr...

# Lernen lernen

1986 Rumelhart, Hinton und Williams machen das Konzept des **Backpropagation** Ansatzes populär → Lernen in mehrschichtigen Netzwerken möglich

Wie lernt ein Netzwerk?

# Lernen lernen

1986 Rumelhart, Hinton und Williams machen das Konzept des **Backpropagation** Ansatzes populär → Lernen in mehrschichtigen Netzwerken möglich

Wie lernt ein Netzwerk?



**Supervised** Learning

# Lernen lernen

1986 Rumelhart, Hinton und Williams machen das Konzept des **Backpropagation** Ansatzes populär → Lernen in mehrschichtigen Netzwerken möglich

Wie lernt ein Netzwerk?



**Supervised** Learning

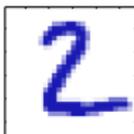
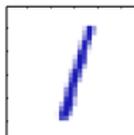


**Unsupervised** Learning

# Supervised Learning

Erfordert ein **Training Set** bestehend aus

Samples



Targets

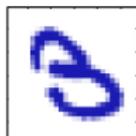
3

1

2

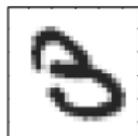
# Formale Schreibweise

Datenvorverarbeitung



# Formale Schreibweise

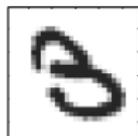
Datenvorverarbeitung



Konvertierung in ein Graustufenbild

# Formale Schreibweise

Datenvorverarbeitung

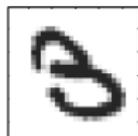


8x8 Pixel 8-Bit Graustufenbild als Vektor

$$\mathbf{x} = (0, 1, 6, 15, 12, 1, 0, 0, 0, 7, 16, 6, 6, 10, \dots) \in \mathbb{R}^{64}$$

# Formale Schreibweise

Datenvorverarbeitung



8x8 Pixel 8-Bit Graustufenbild als Vektor

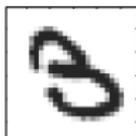
$$\mathbf{x} = (0, 1, 6, 15, 12, 1, 0, 0, 0, 7, 16, 6, 6, 10, \dots) \in \mathbb{R}^{64}$$

Alle  $N$  Samples

$$\{\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^N\}$$

# Formale Schreibweise

Datenvorverarbeitung



8x8 Pixel 8-Bit Graustufenbild als Vektor

$$\mathbf{x} = (0, 1, 6, 15, 12, 1, 0, 0, 0, 7, 16, 6, 6, 10, \dots) \in \mathbb{R}^{64}$$

Alle  $N$  Samples

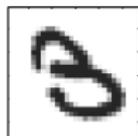
$$\{\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^N\}$$

zusammen mit den Targets  $\{1, 2, 3\}$  in einem Targetvektor

$$\mathbf{t} = (3, 1, 2, \dots) \in \mathbb{R}^N$$

# Formale Schreibweise

Datenvorverarbeitung



8x8 Pixel 8-Bit Graustufenbild als Vektor

$$\mathbf{x} = (0, 1, 6, 15, 12, 1, 0, 0, 0, 7, 16, 6, 6, 10, \dots) \in \mathbb{R}^{64}$$

Alle  $N$  Samples

$$\{\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^N\}$$

zusammen mit den Targets  $\{1, 2, 3\}$  in einem Targetvektor

$$\mathbf{t} = (3, 1, 2, \dots) \in \mathbb{R}^N$$

ergibt das **Training Set**

# Was sagt das Netz?

- 1 **Initialisierung** des Netzwerks mit Zufallsgewichten  $w^0$

# Was sagt das Netz?

- 1 **Initialisierung** des Netzwerks mit Zufallsgewichten  $w^0$
- 2 **Fehlerbestimmung**: Auswertung einer **Errorfunktion** z.B.

$$E(w^0) = |y(x^1, w^0) - t_1|^2 + \dots + |y(x^N, w^0) - t_N|^2 = \sum_{n=1}^N |y(x^n, w^0) - t_n|^2$$

→ Abweichung des Netzwerkoutputs von den gewünschten Targets

# Was sagt das Netz?

- 1 **Initialisierung** des Netzwerks mit Zufallsgewichten  $w^0$
- 2 **Fehlerbestimmung**: Auswertung einer **Errorfunktion** z.B.

$$E(w^0) = |y(x^1, w^0) - t_1|^2 + \dots + |y(x^N, w^0) - t_N|^2 = \sum_{n=1}^N |y(x^n, w^0) - t_n|^2$$

→ Abweichung des Netzwerkoutputs von den gewünschten Targets

Ist  $E(w^0) = 0$  → Glückstreffer: *perfekte* Netzantwort

# Was sagt das Netz?

- 1 **Initialisierung** des Netzwerks mit Zufallsgewichten  $\mathbf{w}^0$
- 2 **Fehlerbestimmung**: Auswertung einer **Errorfunktion** z.B.

$$E(\mathbf{w}^0) = |y(\mathbf{x}^1, \mathbf{w}^0) - t_1|^2 + \dots + |y(\mathbf{x}^N, \mathbf{w}^0) - t_N|^2 = \sum_{n=1}^N |y(\mathbf{x}^n, \mathbf{w}^0) - t_n|^2$$

→ Abweichung des Netzwerkoutputs von den gewünschten  
Targets

Ist  $E(\mathbf{w}^0) = 0$  → Glückstreffer: *perfekte* Netzantwort

Ist  $E(\mathbf{w}^0) > 0$  → Finde  $\mathbf{w}^1$  sodass  $E(\mathbf{w}^1) < E(\mathbf{w}^0)$

Hoffnung: nach  $K$  Schritten gilt

$$E(\mathbf{w}^0) > E(\mathbf{w}^1) > \dots > E(\mathbf{w}^K) = \min_{\mathbf{w}} E(\mathbf{w})$$

# Was sagt das Netz?

- 1 **Initialisierung** des Netzwerks mit Zufallsgewichten  $\mathbf{w}^0$
- 2 **Fehlerbestimmung**: Auswertung einer **Errorfunktion** z.B.

$$E(\mathbf{w}^0) = |y(\mathbf{x}^1, \mathbf{w}^0) - t_1|^2 + \dots + |y(\mathbf{x}^N, \mathbf{w}^0) - t_N|^2 = \sum_{n=1}^N |y(\mathbf{x}^n, \mathbf{w}^0) - t_n|^2$$

→ Abweichung des Netzwerkoutputs von den gewünschten  
Targets

Ist  $E(\mathbf{w}^0) = 0$  → Glückstreffer: *perfekte* Netzantwort

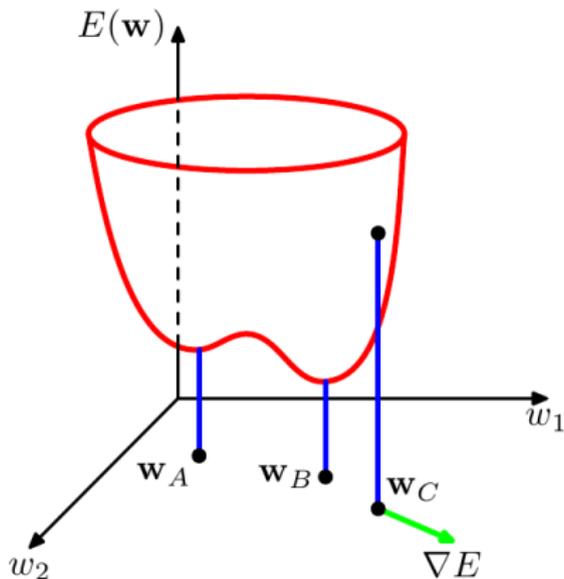
Ist  $E(\mathbf{w}^0) > 0$  → Finde  $\mathbf{w}^1$  sodass  $E(\mathbf{w}^1) < E(\mathbf{w}^0)$

Hoffnung: nach  $K$  Schritten gilt

$$E(\mathbf{w}^0) > E(\mathbf{w}^1) > \dots > E(\mathbf{w}^K) = \min_{\mathbf{w}} E(\mathbf{w})$$

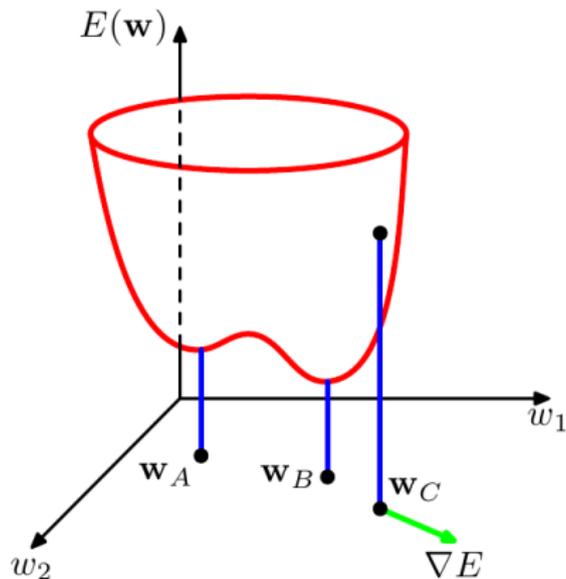
**Optimierungsproblem**

# Ein Minimierungsproblem



Gradientenabstieg

# Ein Minimierungsproblem



$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E(\mathbf{w}^{\tau})$$

# Backpropagation

## Algorithm

*Initialisierung:* Zufallsgewichte  $\mathbf{w}^0$ , setze  $\tau = 0$ ,  $tol = 1e - 3$

1: *while*  $|\nabla E(\mathbf{w}^\tau)| > tol$  *do*

6: *end while*

# Backpropagation

## Algorithm

*Initialisierung:* Zufallsgewichte  $\mathbf{w}^0$ , setze  $\tau = 0$ ,  $tol = 1e - 3$

1: *while*  $|\nabla E(\mathbf{w}^\tau)| > tol$  *do*

2:     *Forward-Pass:* Berechnung des Netzoutputs  $y(\mathbf{x}^k, \mathbf{w}^\tau)$  für jedes Sample  $\mathbf{x}^k$

6: *end while*

# Backpropagation

## Algorithm

*Initialisierung:* Zufallsgewichte  $\mathbf{w}^0$ , setze  $\tau = 0$ ,  $tol = 1e - 3$

1: *while*  $|\nabla E(\mathbf{w}^\tau)| > tol$  *do*

2:     *Forward-Pass:* Berechnung des Netzoutputs  $y(\mathbf{x}^k, \mathbf{w}^\tau)$  für jedes Sample  $\mathbf{x}^k$

3:     *Fehlerbestimmung:* Auswertung von  $E(\mathbf{w}^\tau)$

6: *end while*

# Backpropagation

## Algorithm

*Initialisierung:* Zufallsgewichte  $\mathbf{w}^0$ , setze  $\tau = 0$ ,  $tol = 1e - 3$

1: *while*  $|\nabla E(\mathbf{w}^\tau)| > tol$  *do*

2:     *Forward-Pass:* Berechnung des Netzoutputs  $y(\mathbf{x}^k, \mathbf{w}^\tau)$  für jedes Sample  $\mathbf{x}^k$

3:     *Fehlerbestimmung:* Auswertung von  $E(\mathbf{w}^\tau)$

4:     *Backward Pass:* Berechnung der neuen Gewichte

$$w_{ji}^{\tau+1} = w_{ji}^\tau - \eta \frac{\partial E(\mathbf{w}^\tau)}{\partial w_{ji}}$$

6: *end while*

# Backpropagation

## Algorithm

*Initialisierung:* Zufallsgewichte  $\mathbf{w}^0$ , setze  $\tau = 0$ ,  $tol = 1e - 3$

1: *while*  $|\nabla E(\mathbf{w}^\tau)| > tol$  *do*

2:     *Forward-Pass:* Berechnung des Netzoutputs  $y(\mathbf{x}^k, \mathbf{w}^\tau)$  für jedes Sample  $\mathbf{x}^k$

3:     *Fehlerbestimmung:* Auswertung von  $E(\mathbf{w}^\tau)$

4:     *Backward Pass:* Berechnung der neuen Gewichte

$$w_{ji}^{\tau+1} = w_{ji}^\tau - \eta \frac{\partial E(\mathbf{w}^\tau)}{\partial w_{ji}}$$

5:     *Update:*  $\tau \leftarrow \tau + 1$

6: *end while*

# An die Maschinen...

