

---

# Einführung in Web- und Data-Science

Differenzierbare Programmierung

Prof. Dr. Ralf Möller

**Universität zu Lübeck**

**Institut für Informationssysteme**

---

Differenzierbare Programmierung

# PERZEPTRONS



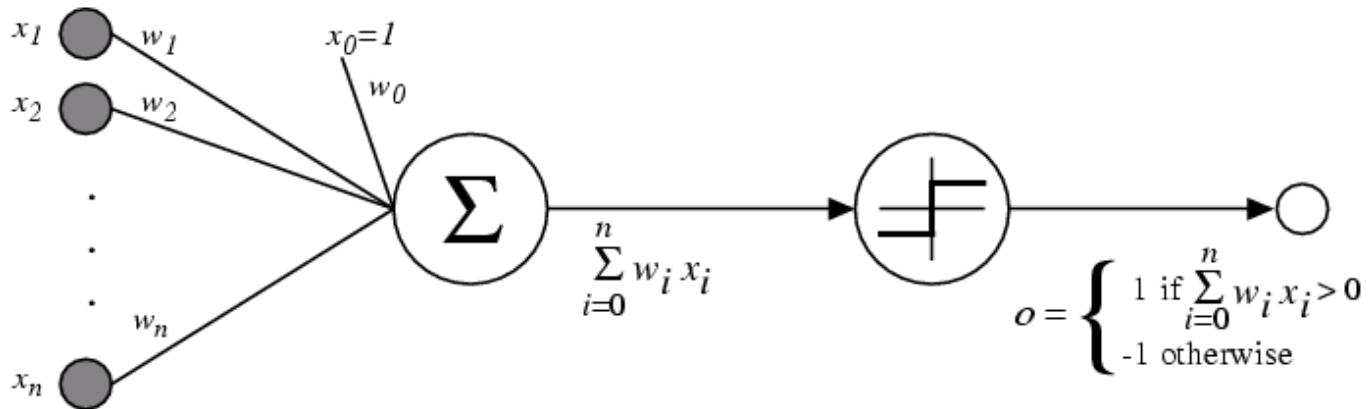
# Repräsentation von Funktionen ...

---

- ... durch Perzeptrons
  - Ein-Ebenen-Netzwerk (Linearer Klassifikator)
  - Mehrebenen-Netzwerke
  - Lernregel Fehlerrückführung (Backpropagation)

Frank Rosenblatt, *The Perceptron--a perceiving and recognizing automaton*. Report 85-460-1, Cornell Aeronautical Laboratory, **1957**

# Perzeptron



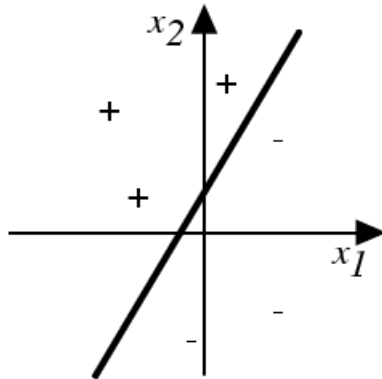
$$o(x_1, \dots, x_n) = \begin{cases} 1 & \text{if } w_0 + w_1 x_1 + \dots + w_n x_n > 0 \\ -1 & \text{sonst} \end{cases}$$

Manchmal einfacher geschrieben als (Ann.:  $x_0 = 1$ ):

$$o(\vec{x}) = \begin{cases} 1 & \text{if } \vec{w} \cdot \vec{x} > 0 \\ -1 & \text{sonst} \end{cases}$$



# Entscheidungslinie eines Perzeptrons



Repräsentiert lineare Funktion  $y = mx + b$

- Was machen die Gewichte?  
 $g(x_1, x_2) = AND(x_1, x_2)$ ?

Verallgemeinerung auf Entscheidungsebenen möglich

Vgl.: Warren McCulloch und William Pitts: *A logical calculus of the ideas immanent in nervous activity*. In: *Bulletin of Mathematical Biophysics*, Bd. 5, S. 115–133, 1943

# Trennlinien

---

- Wenn  $w_2x_2 + w_1x_1 + w_0 > 0$ , dann wird für  $(x_1, x_2)$  ein (+) vergeben.
- $w_2x_2 + w_1x_1 + w_0 = 0$  ist Trennlinie. Warum?
- Es wird die Gerade  $x_2 = -w_1/w_2 x_1 - w_0/w_2$  als Trennlinie verwendet.
- Wir betrachten  $y = mx + b$
- Bedingung lautet:  $y - mx - b > 0$
- Wähle  $m=1$  und  $b$  positiv. Dann betrachte  $(x, y) = (0, 2b)$ . Das liegt über der Geraden.
- $2b - 0 - b > 0$
- Also wird  $(0, 2b)$  mit 1 bewertet (+)

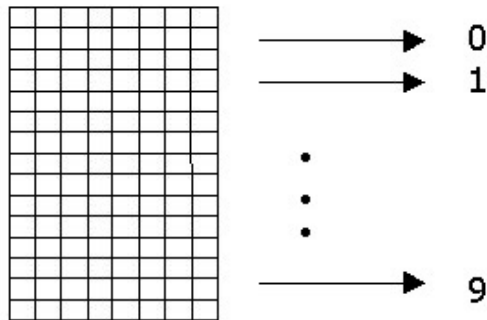
# Lassen sich alle Funktionen repräsentieren?

---

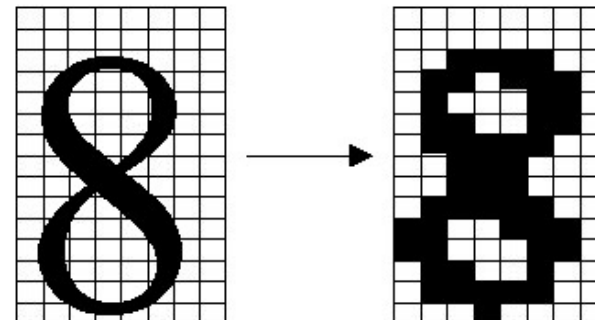
- Kann die Fehlerfunktion immer sinnvoll minimiert werden?
- XOR-Problem
  - Einführung weiterer Dimensionen
    - Erweiterung der Daten?
  - Besser: Einführung weiterer Ebenen im Netz

# Ein Anwendungsbeispiel

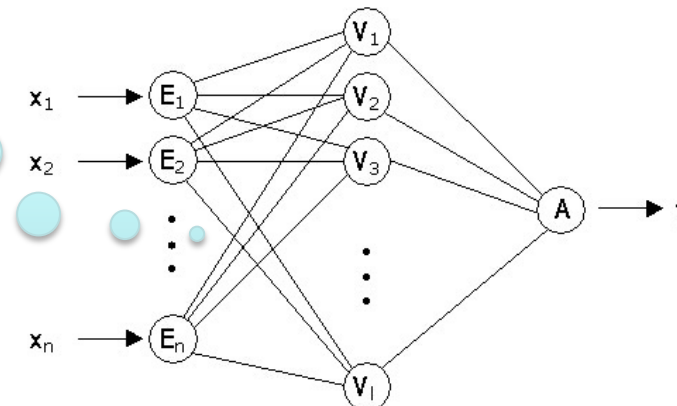
Das folgende Netz soll Ziffern von 0 bis 9 erkennen. Dafür wird zunächst das *Eingabefeld* in 8x15 Elemente aufgeteilt:



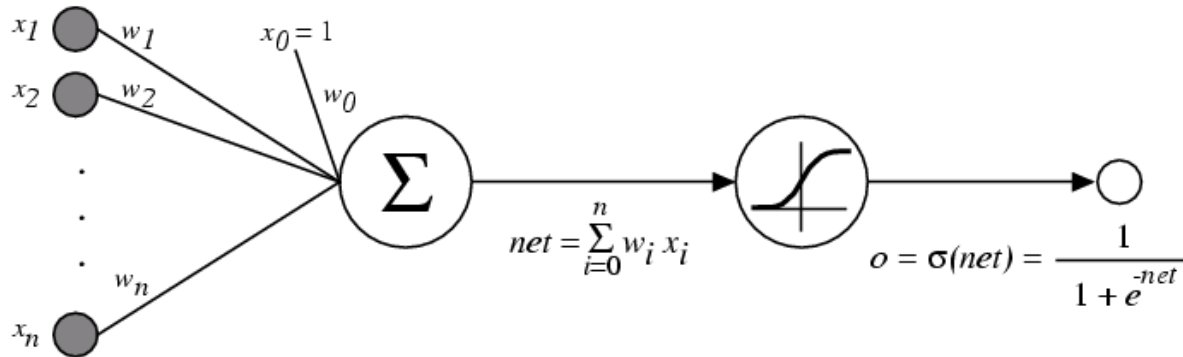
Die geschriebene Ziffer wird in eine Folge von Nullen und Einsen umgewandelt, wobei 0 für leere und 1 für übermalte Rasterpunkte steht:



Wie können wir die Parameter automatisch so bestimmen, dass geschriebene Ziffern erkannt und als Ausgabe  $y$  ausgegeben werden?



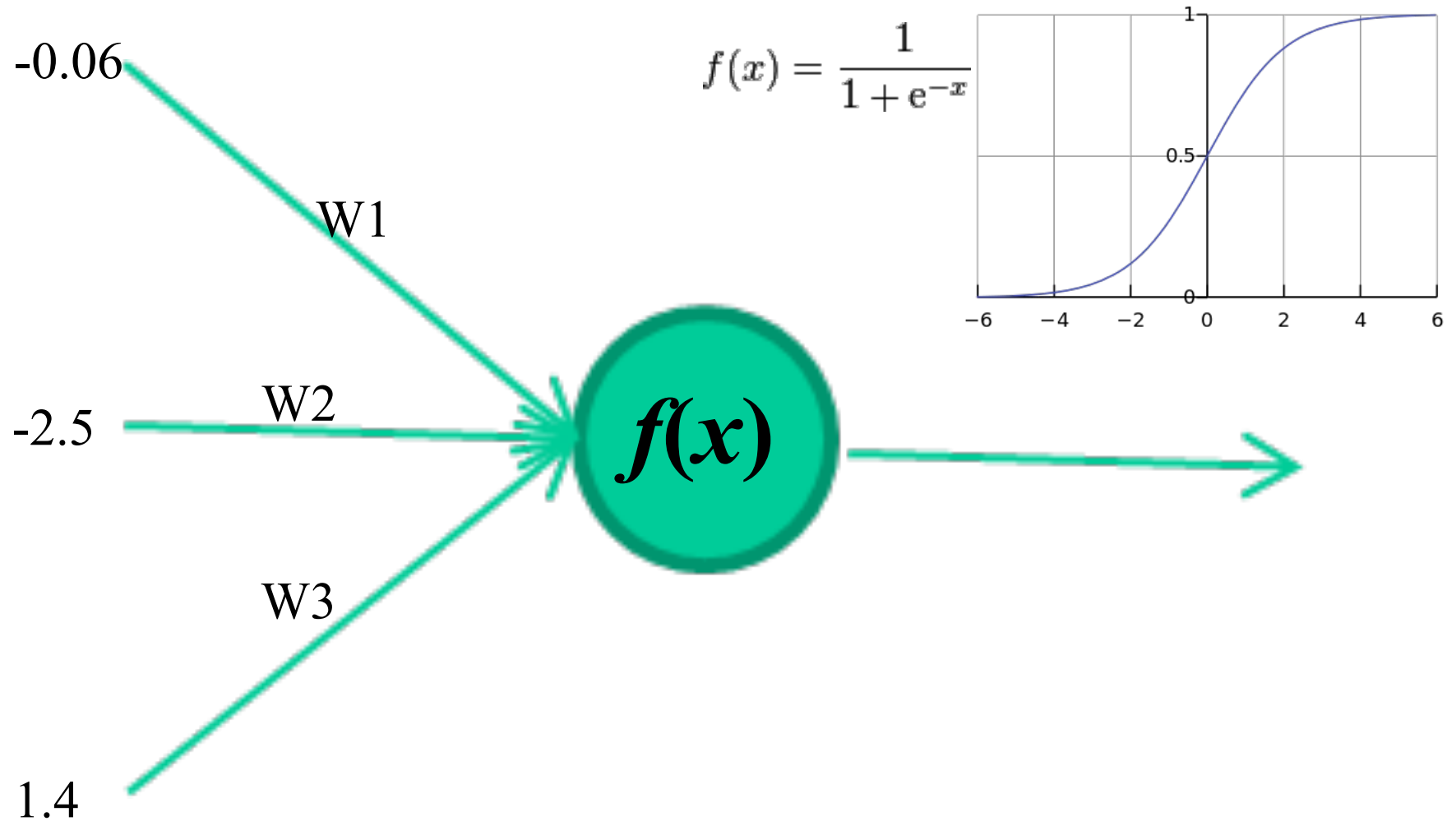
# Verwendung kontinuierlicher Funktionen



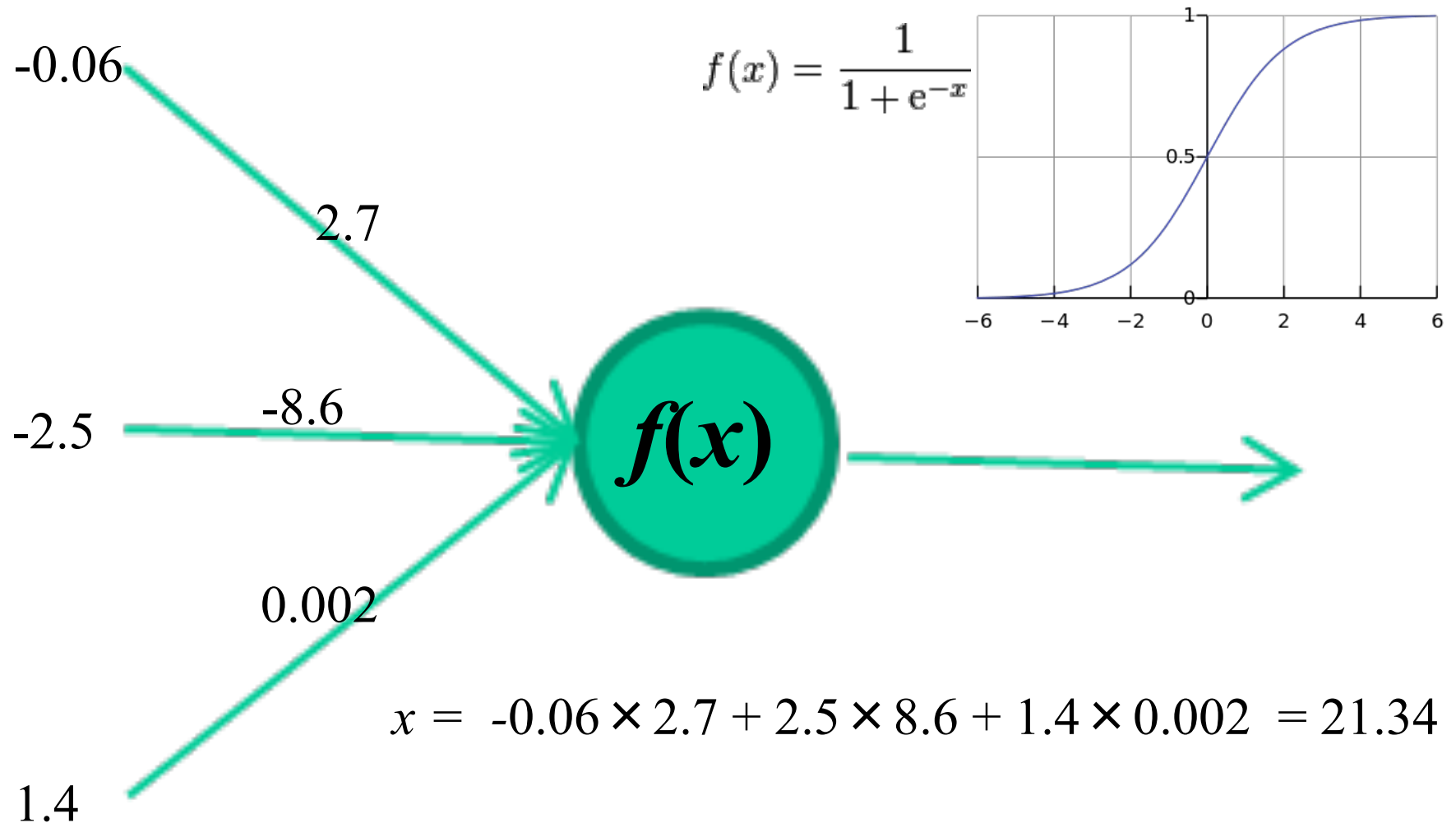
$\sigma(x)$  ist die Sigmoid-Funktion

$$\frac{1}{1 + e^{-x}}$$

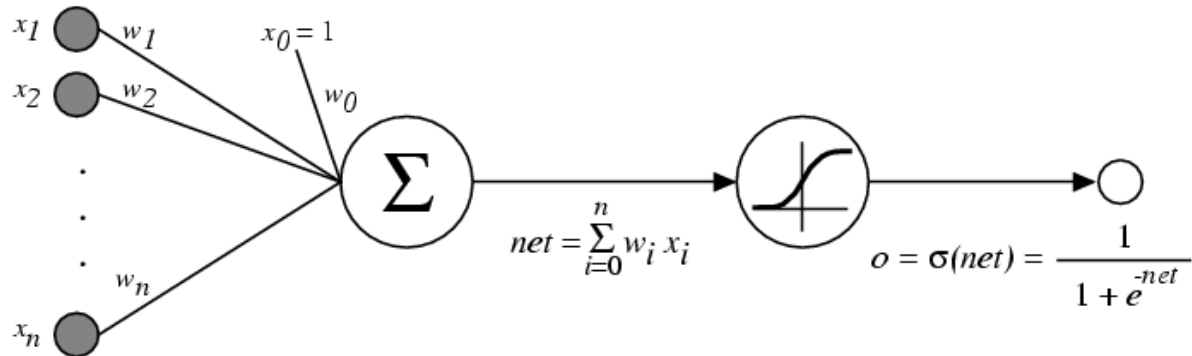
# Beispiel einer Sigmoid-Funktion



# Beispiel einer Sigmoid-Funktion



# Sigmoid-Einheit



$\sigma(x)$  ist die Sigmoid-Funktion (auch: logistische Funktion)

$$\frac{1}{1 + e^{-x}}$$

Eigenschaft:  $\frac{d\sigma(x)}{dx} = \sigma(x)(1 - \sigma(x))$  (Gradient)

- Wir können den Gradienten verwenden, um die Einheit anzupassen
- Wir kommen gleich darauf zurück



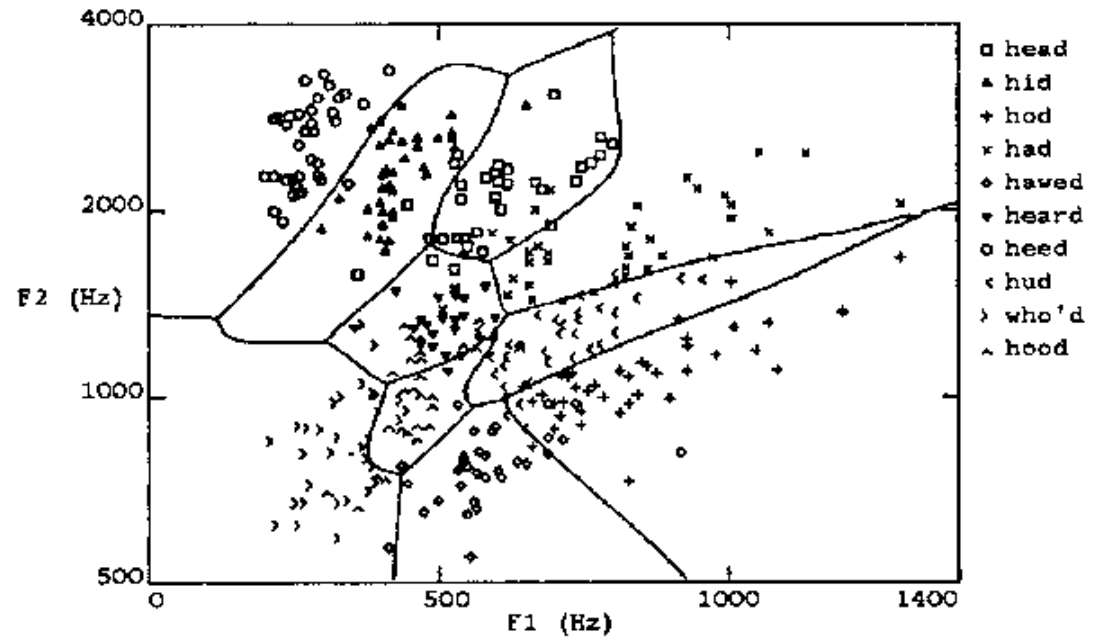
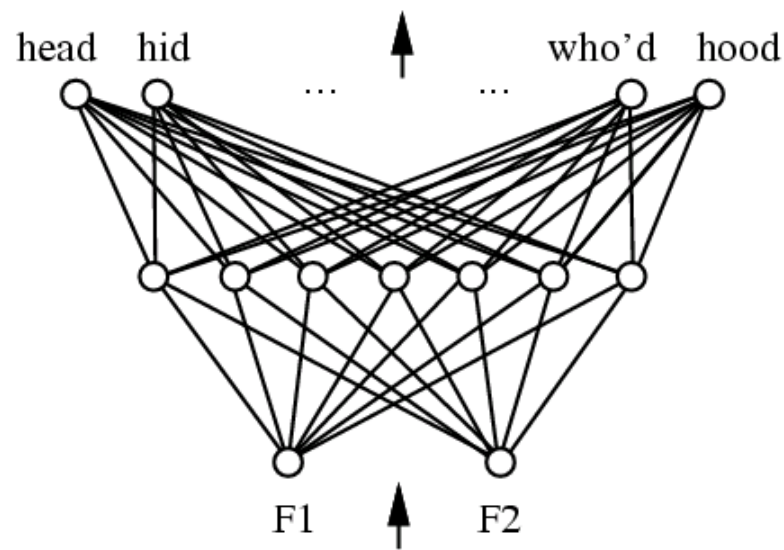
---

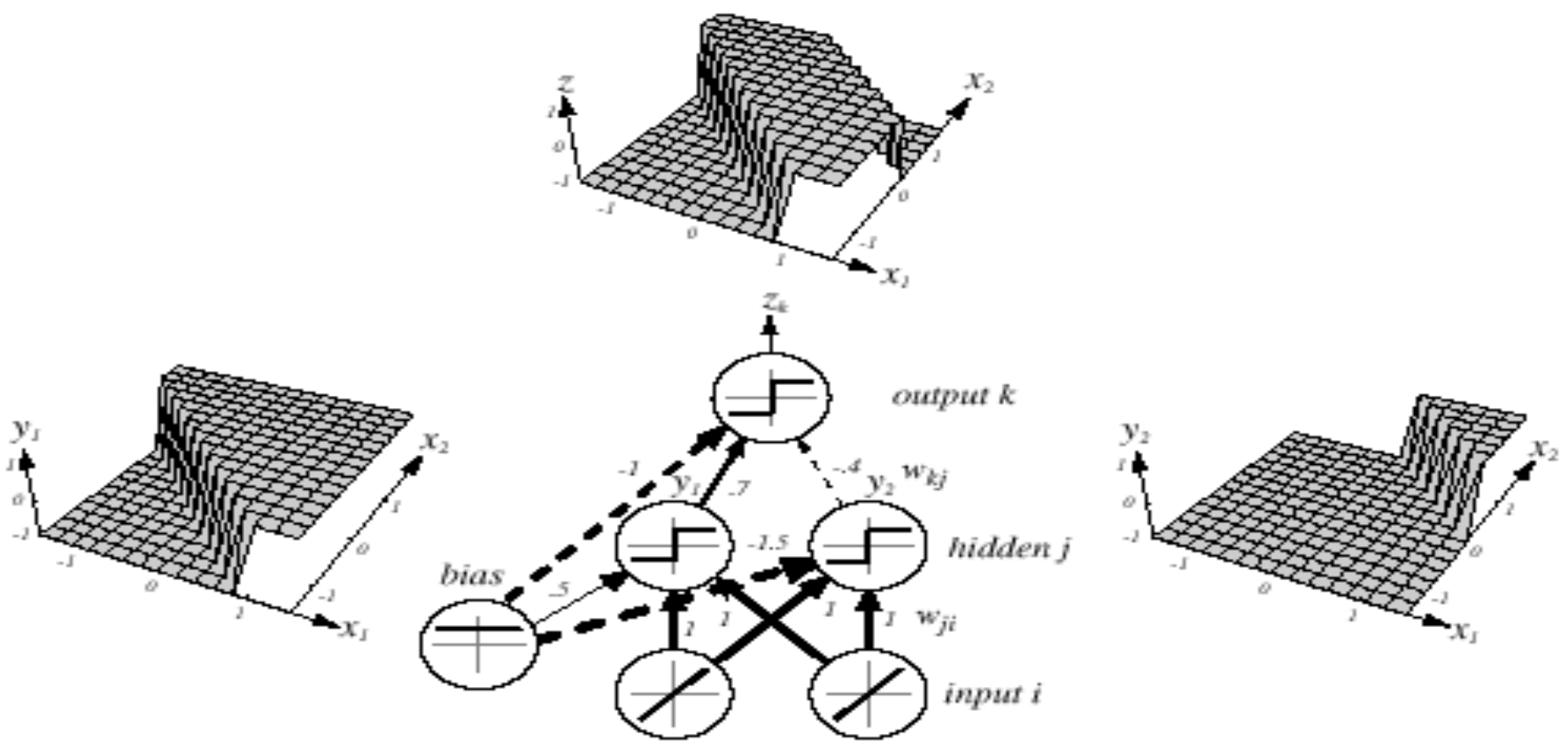
Differenzierbare Programmierung

# MEHR SCHICHTEN PERZEPTRONS

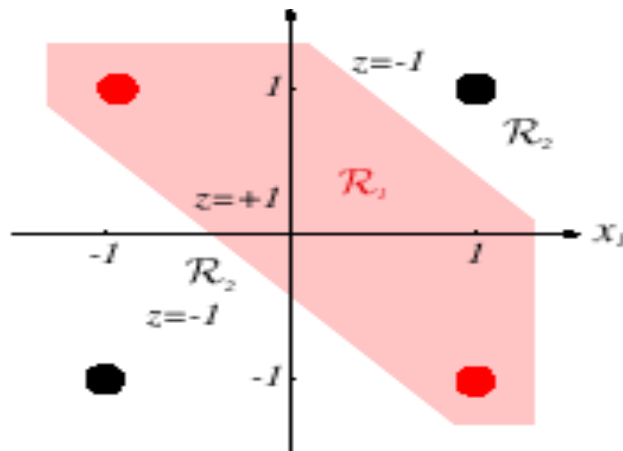


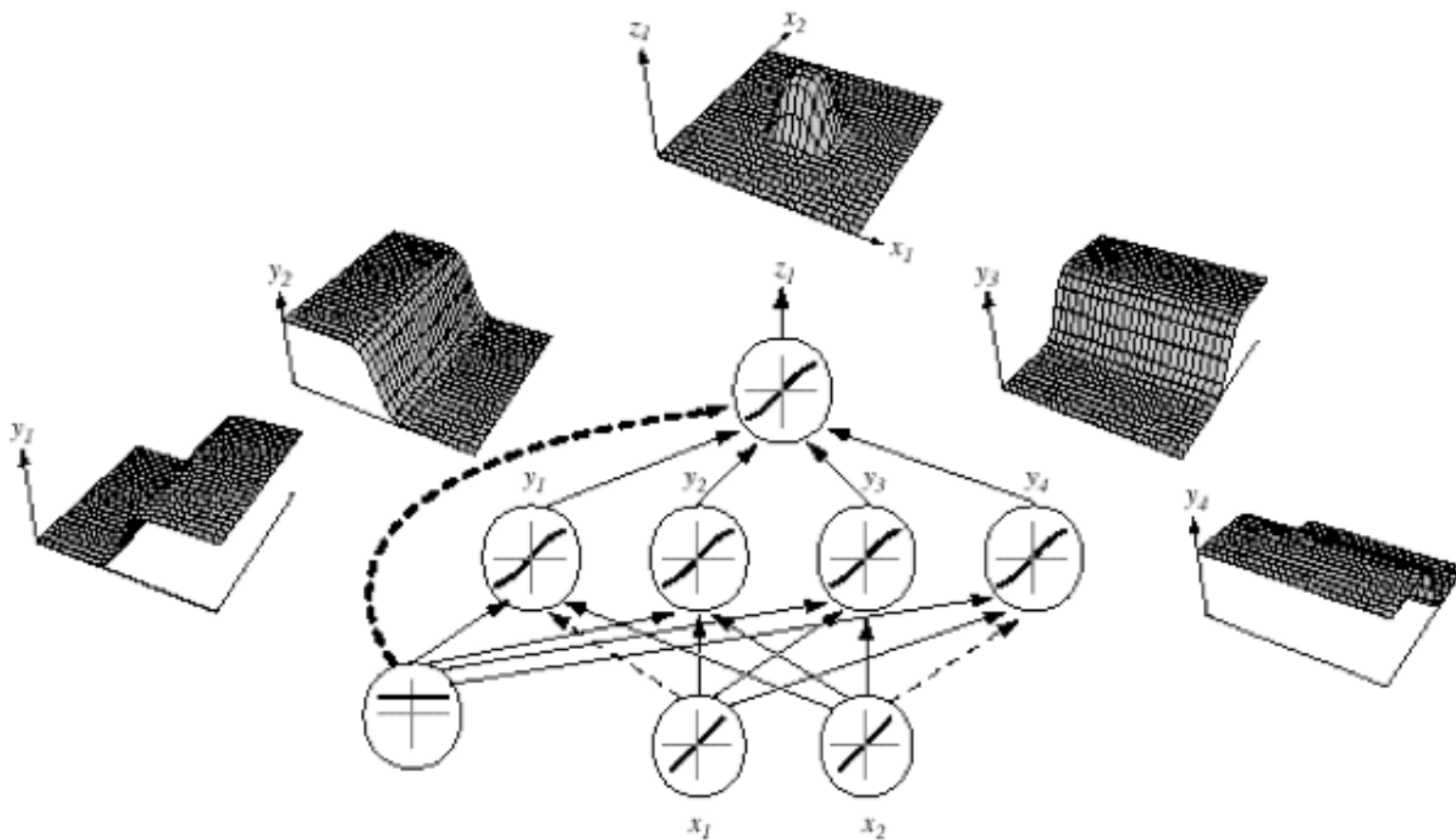
# Mehr-Ebenen Netze von Sigmoid-Einheiten



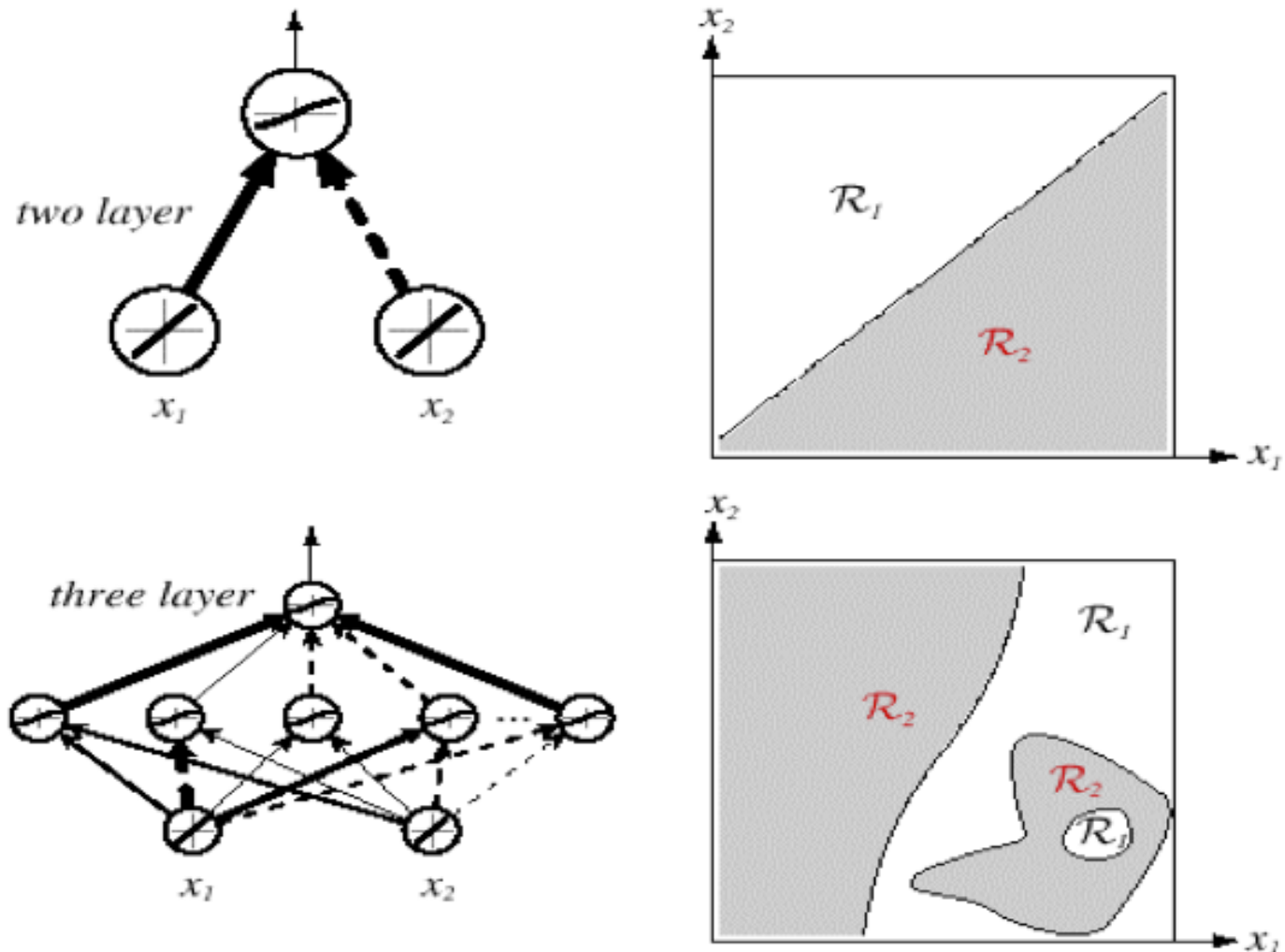


$$Z = x_1 \text{ XOR } x_2 = (x_1 \text{ OR } x_2) \text{ AND NOT}(x_1 \text{ AND } x_2)$$





**FIGURE 6.2.** A 2-4-1 network (with bias) along with the response functions at different units; each hidden output unit has sigmoidal activation function  $f(\cdot)$ . In the case shown, the hidden unit outputs are paired in opposition thereby producing a “bump” at the output unit. Given a sufficiently large number of hidden units, any continuous function from input to output can be approximated arbitrarily well by such a network. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

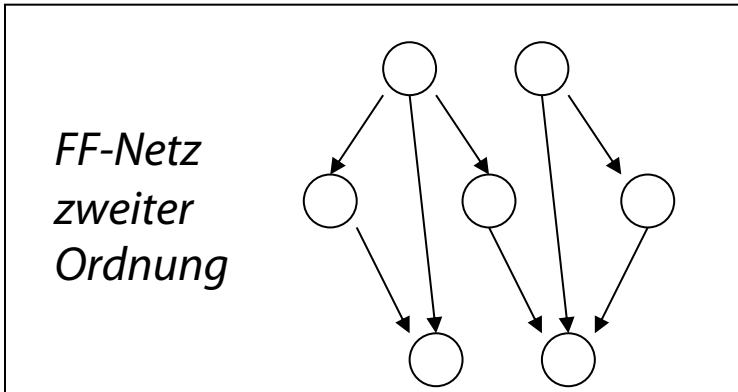
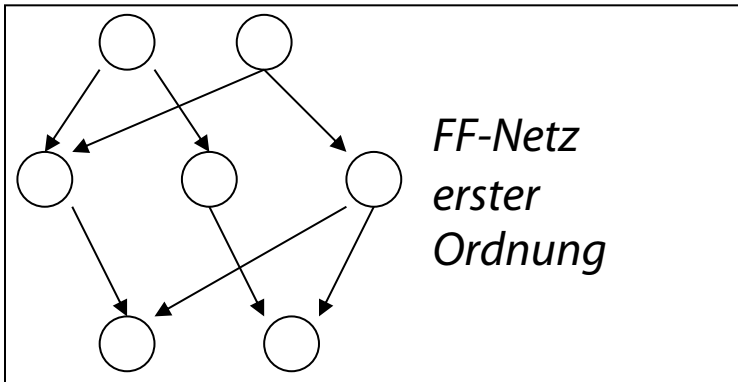


**FIGURE 6.3.** Whereas a two-layer network classifier can only implement a linear decision boundary, given an adequate number of hidden units, three-, four- and higher-layer networks can implement arbitrary decision boundaries. The decision regions need not be convex or simply connected. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

# Netztopologien

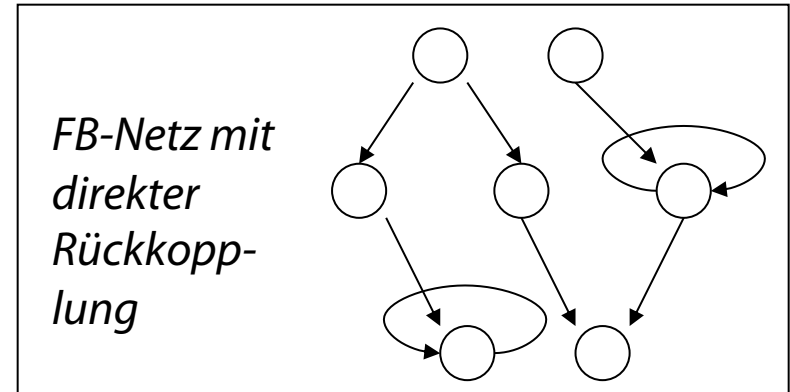
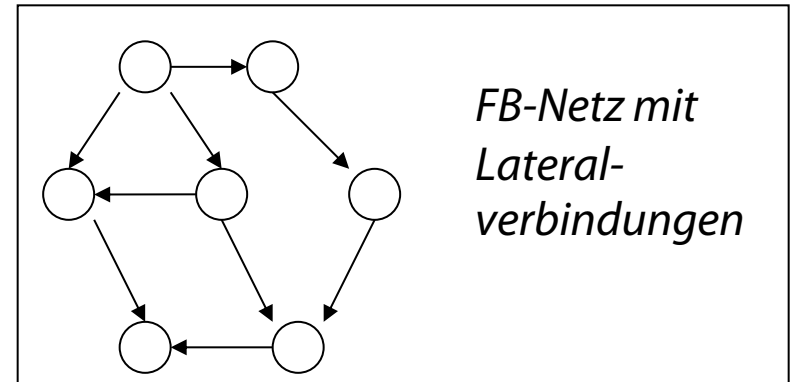
## FeedForward-Netze:

Gerichtete Verbindungen nur von niedrigen zu höheren Schichten



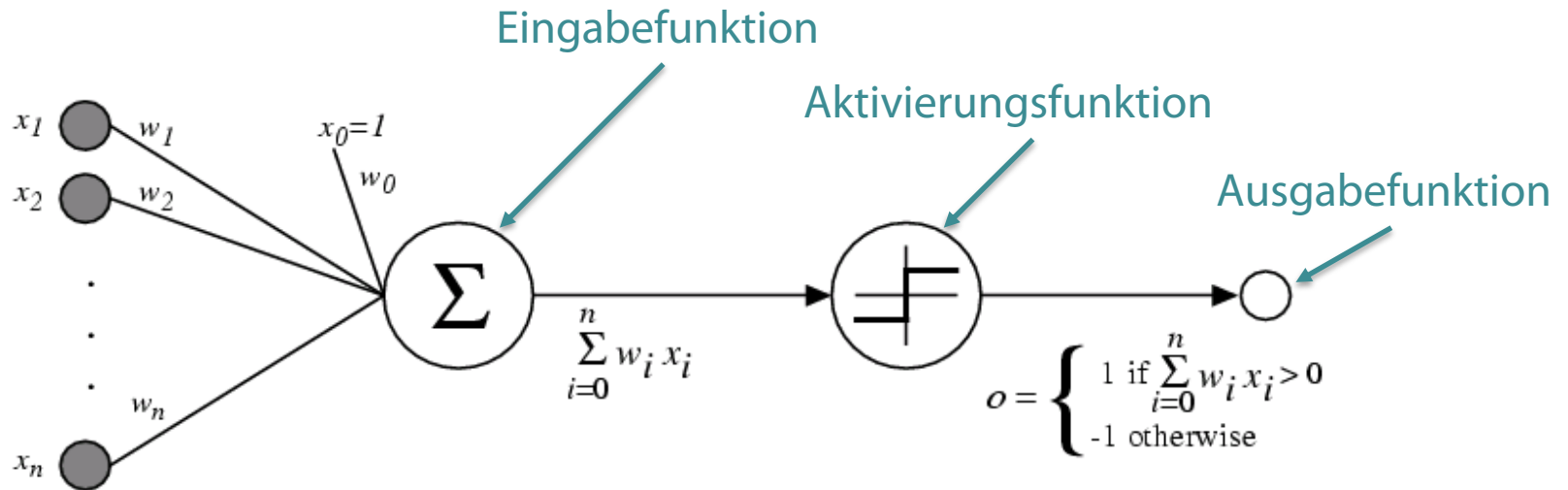
## FeedBack-Netze (rekurrente Netze):

Verbindungen zwischen allen Schichten möglich (Abrollen)



# Netzwerk von Perzeptrons

- Perzeptron:



- Was muss ich vor dem Lernen/Trainieren bestimmen?
  - Netzwerk-Topologie?
  - Perzeptron bezogen?

---

Differenzierbare Programmierung

# **BEGRIFFSBESTIMMUNG**





# Die Eingabefunktion

Die Eingabe- oder Propagierungsfunktion berechnet aus dem Eingabevektor  $\vec{x} = (x_1, \dots, x_m)$  und dem Gewichtsvektor  $\vec{w}_k = (w_{1,k}, \dots, w_{m,k})$  den Nettoinput des  $k$ -ten Knotens. Es gibt folgende Inputfunktionen:

- Summe:  $net_k = f_{in}(\vec{w}_k, \vec{x}) = \sum_{i=1}^m w_{i,k} \cdot x_i$
- Maximalwert:  $net_k = f_{in}(\vec{w}_k, \vec{x}) = \max_i (w_{i,k} \cdot x_i)$
- Produkt:  $net_k = f_{in}(\vec{w}_k, \vec{x}) = \prod_{i=1}^m w_{i,k} \cdot x_i$
- Minimalwert:  $net_k = f_{in}(\vec{w}_k, \vec{x}) = \min_i (w_{i,k} \cdot x_i)$

# Die Aktivierungsfunktion

Mit der Aktivierungsfunktion (auch: Transferfunktion) wird aus dem Nettoinput  $net_k$  der Aktivierungszustand  $a_k$  eines Knotens berechnet.

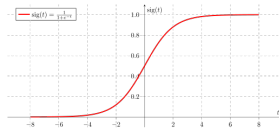
Folgende Aktivierungsfunktionen sind gebräuchlich:

- Lineare Aktivierungsfkt.:  $a_k = f_{act}(net_k) = c_k \cdot net_k$

Schwellenwert,  
häufig 0

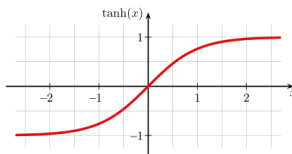
- Binäre Schwellenwertfkt.:  $a_k = f_{act}(net_k) = \begin{cases} 1, & \text{falls } net_k \geq \theta_k, \\ 0, & \text{sonst.} \end{cases}$

- Fermi-Fkt. (logistische Fkt.):  $a_k = f_{act}(net_k) = \frac{1}{1 + e^{-\frac{net_k}{T}}}$



Spezialfall: Sigmoid:  $T=1$

- Tangens hyperbolicus:  $a_k = f_{act}(net_k) = \frac{e^{net_k} - e^{-net_k}}{e^{net_k} + e^{-net_k}} = \frac{1 + \tanh(net_k)}{2}$

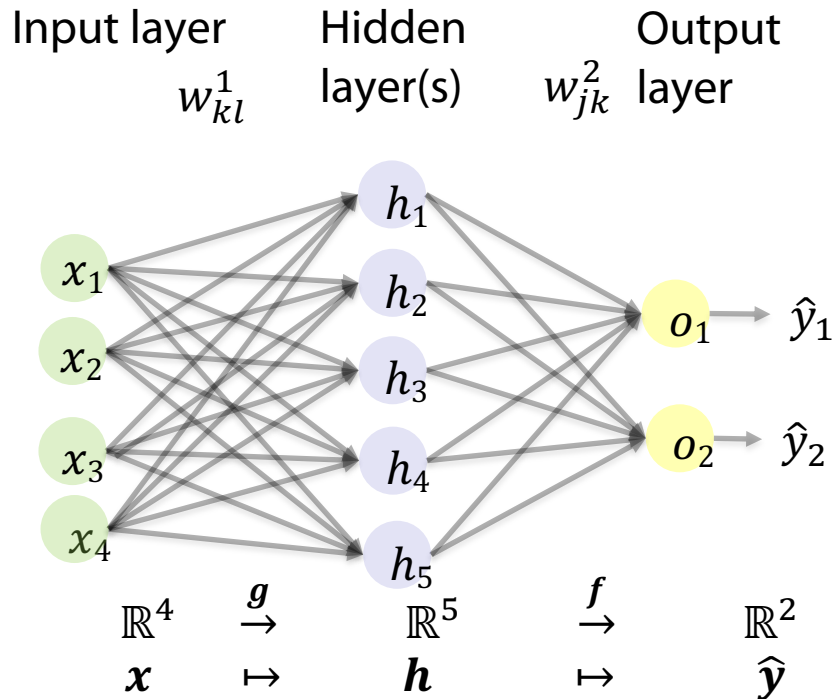


# Die Ausgabefunktion

- Die Ausgabefunktion berechnet aus der Aktivierung  $a_k$  den Wert  $o_k$ , der als Ausgabe an die nächste Schicht weitergegeben wird.
- In den meisten Fällen ist die Ausgabefunktion die **Identität**, d.h.  $o_k = a_k$ .
- Für binäre Ausgaben wird manchmal auch eine Schwellenwertfunktion verwendet:

$$o_k = f_{out}(a_k) = \begin{cases} 1, & \text{falls } a_k \geq \theta_k, \\ 0, & \text{sonst.} \end{cases}$$

# Beispiel



$i$ : Ebene  $i$   
 $b^i$ : Verzerrung für  $i$   
 $W^i$ : Gewichtsmatrix für  $i$   
 $\sigma^i$ : Aktivierungsfkt. für  $i$

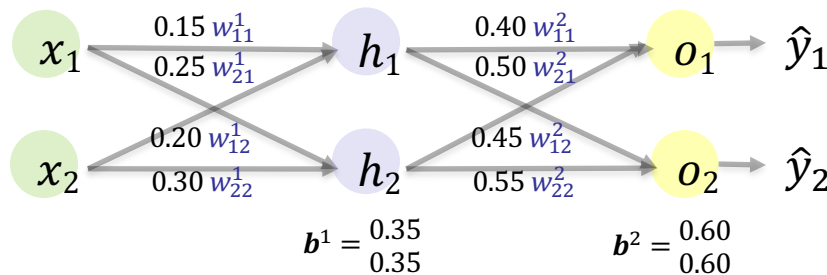
$out^i$ :  $\sigma^i(W^i out^{i-1} + b^i)$   
 Aktivierung in  $i$

$net^i$ :  $W^i out^{i-1} + b^i$   
 Lineare Ausgabe in Ebene  $i$

$$\hat{y} = f(g(x; W^1, b^1); W^2, b^2)$$

$$= \sigma^2(W^2 \sigma^1(W^1 x + b^1) + b^2)$$

# Bespiel: Berechnung



- $net_{h_1} = w_{11}^1 \cdot x_1 + w_{12}^1 \cdot x_2 + b_1^1$
- $net_{h_1} = 0.15 \cdot 0.05 + 0.20 \cdot 0.10 + 0.35 = 0.3775$
- $out_{h_1} = \frac{1}{1+e^{-net_{h_1}}} = \frac{1}{1+e^{-0.3775}} = 0.593269992$
- $out_{h_2} = 0.596884378$
- $\hat{y}_1 = out_{o_1} = 0.75136507$
- $\hat{y}_2 = out_{o_2} = 0.772928465$
- $target_{o_1} = 0.01$
- $target_{o_2} = 0.99$

Aktivierungsfkt:

$$\sigma_1, \sigma_2 = \frac{1}{1 + e^{-x}}$$

Eingabe:

$$x = \begin{pmatrix} 0.05 \\ 0.10 \end{pmatrix}$$

---

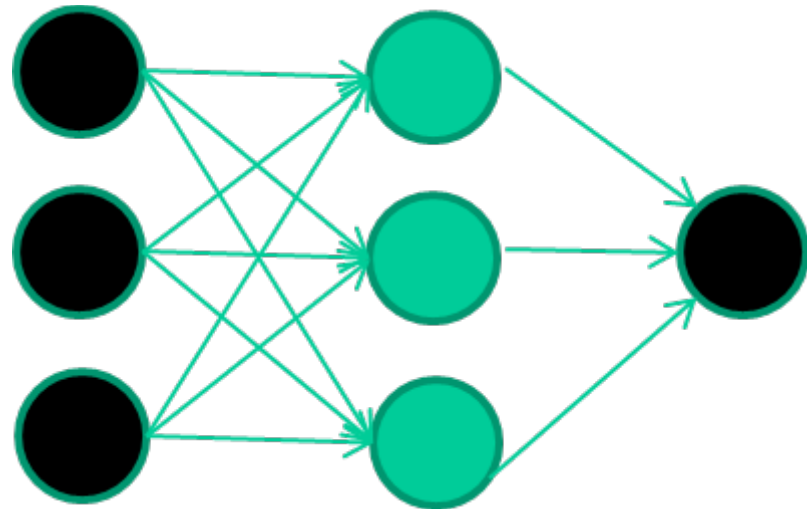
Differenzierbare Programmierung

# **LERNEN VON GEWICHTEN (IDEE)**



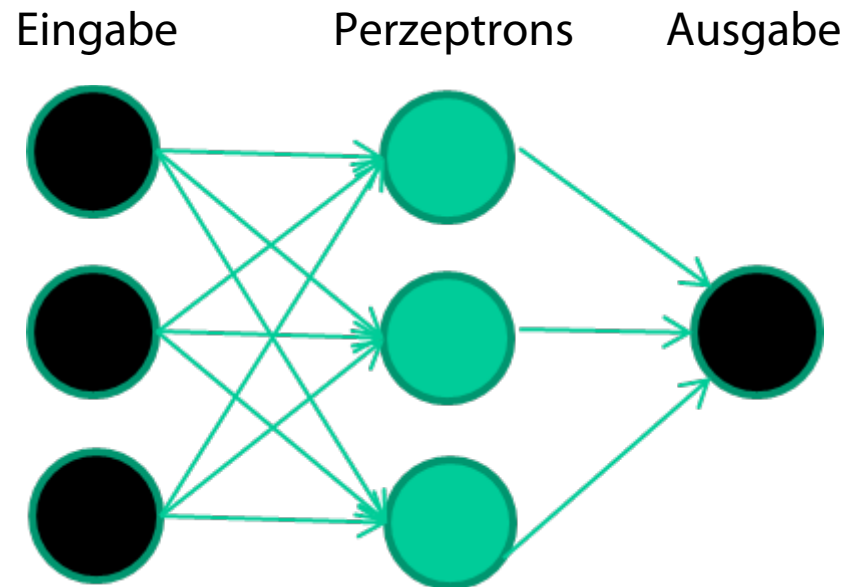
# Einstellen der Gewichte mit Trainingsdaten...

<i>Fields</i>	<i>class</i>
1.4 2.7 1.9	0
3.8 3.4 3.2	0
6.4 2.8 1.7	1
4.1 0.1 0.2	0
etc ...	



# Anlernen des Netzwerks

<i>Fields</i>	<i>class</i>
1.4 2.7 1.9	0
3.8 3.4 3.2	0
6.4 2.8 1.7	1
4.1 0.1 0.2	0
etc ...	

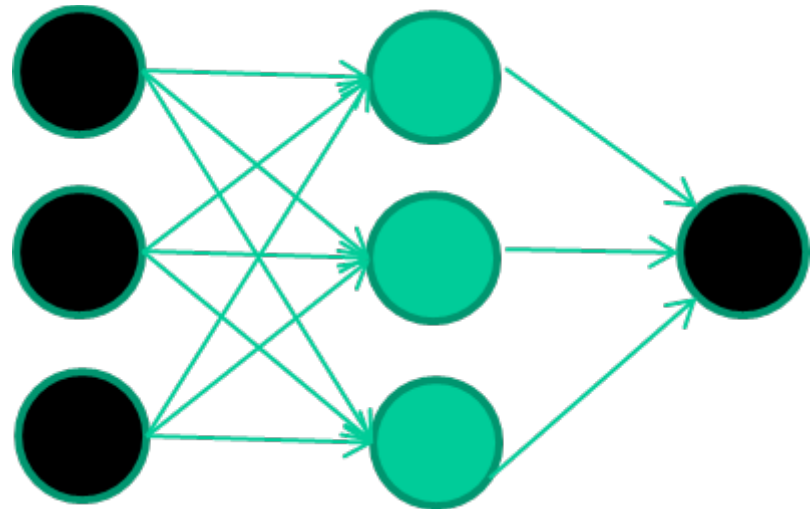




# Trainingsdaten

<i>Fields</i>	<i>class</i>
1.4 2.7 1.9	0
3.8 3.4 3.2	0
6.4 2.8 1.7	1
4.1 0.1 0.2	0
etc ...	

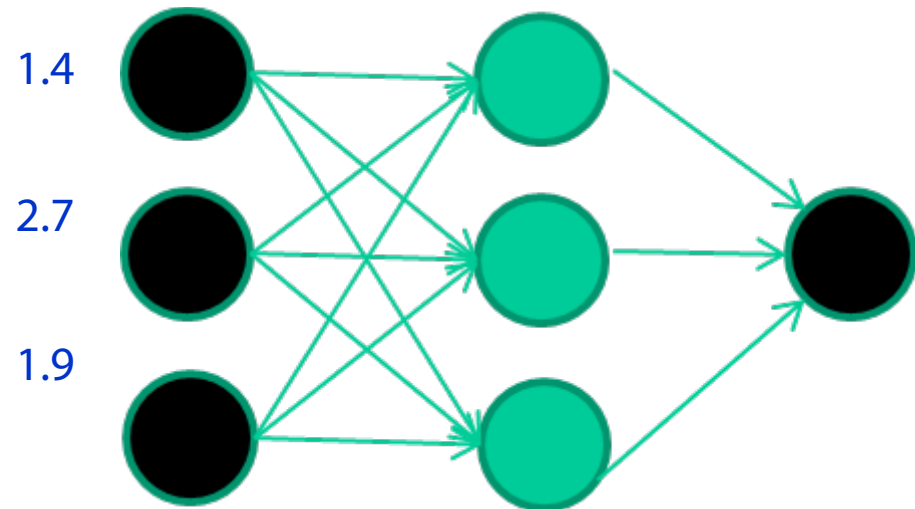
Initialisierung mit zufälligen Gewichten



# Trainingsdaten

<i>Fields</i>	<i>class</i>
1.4 2.7 1.9	0
3.8 3.4 3.2	0
6.4 2.8 1.7	1
4.1 0.1 0.2	0
etc ...	

## Präsentierung eines Trainingsdatensatzes



# Trainingsdaten

*Fields*                      *class*

1.4 2.7 1.9                      0

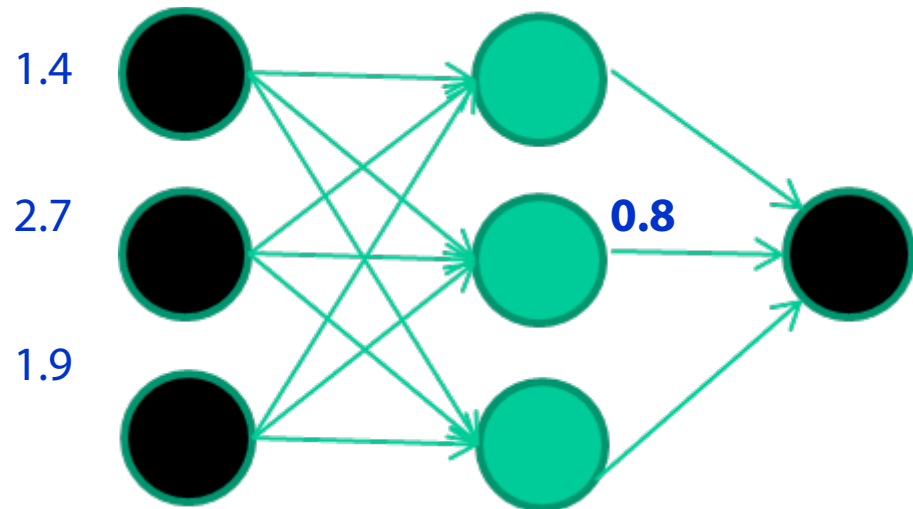
3.8 3.4 3.2                      0

6.4 2.8 1.7                      1

4.1 0.1 0.2                      0

etc ...

Durchpropagierung zur Ausgabe



# Trainingsdaten

*Fields* *class*

1.4 2.7 1.9 0

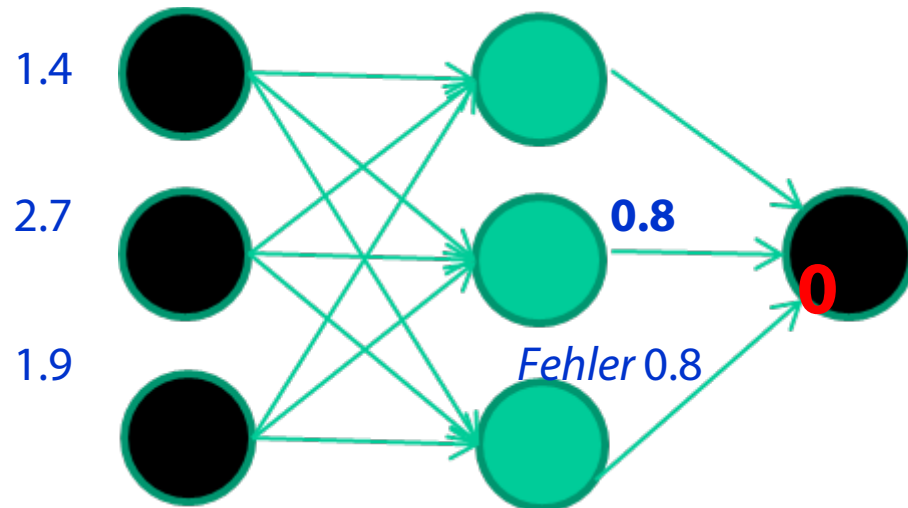
3.8 3.4 3.2 0

6.4 2.8 1.7 1

4.1 0.1 0.2 0

etc ...

Vergleich mit der Zielausgabe



# Trainingsdaten

*Fields* *class*

1.4 2.7 1.9 0

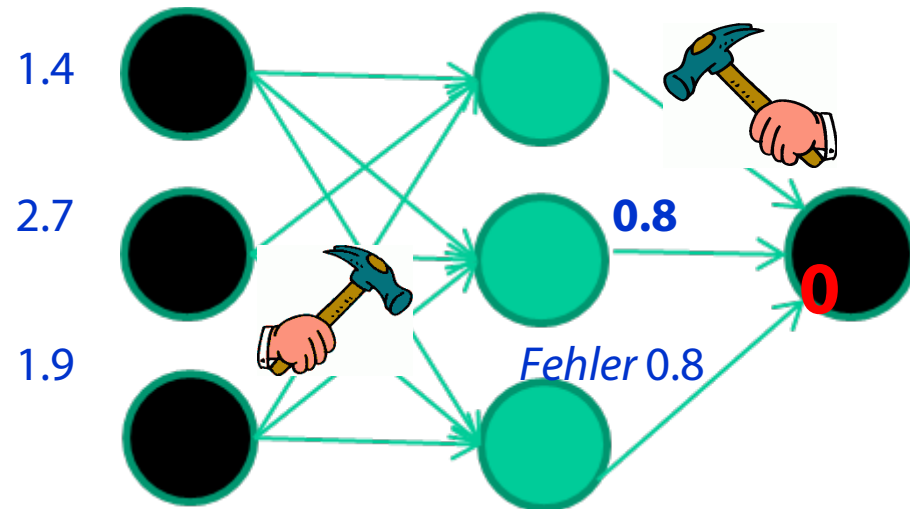
3.8 3.4 3.2 0

6.4 2.8 1.7 1

4.1 0.1 0.2 0

etc ...

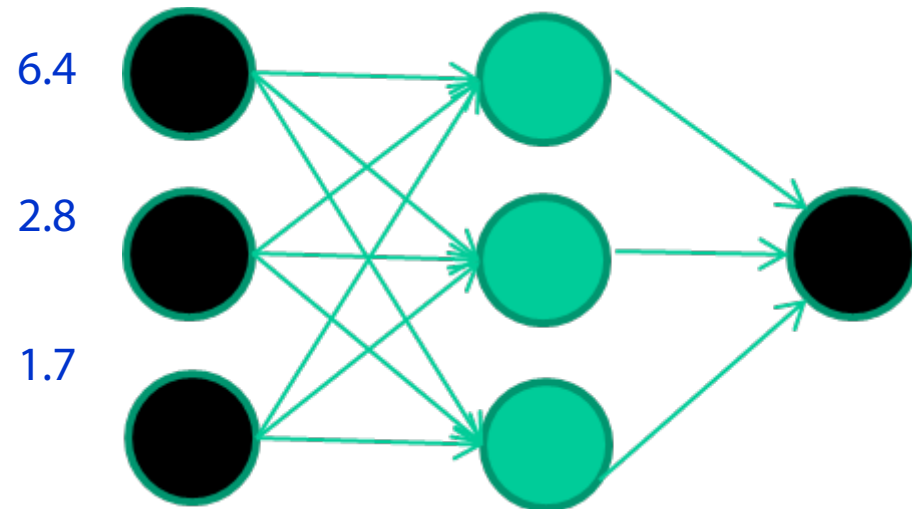
Anpassen der Gewichte gemäß Fehler



# Trainingsdaten

<i>Fields</i>	<i>class</i>
1.4 2.7 1.9	0
3.8 3.4 3.2	0
6.4 2.8 1.7	1
4.1 0.1 0.2	0
etc ...	

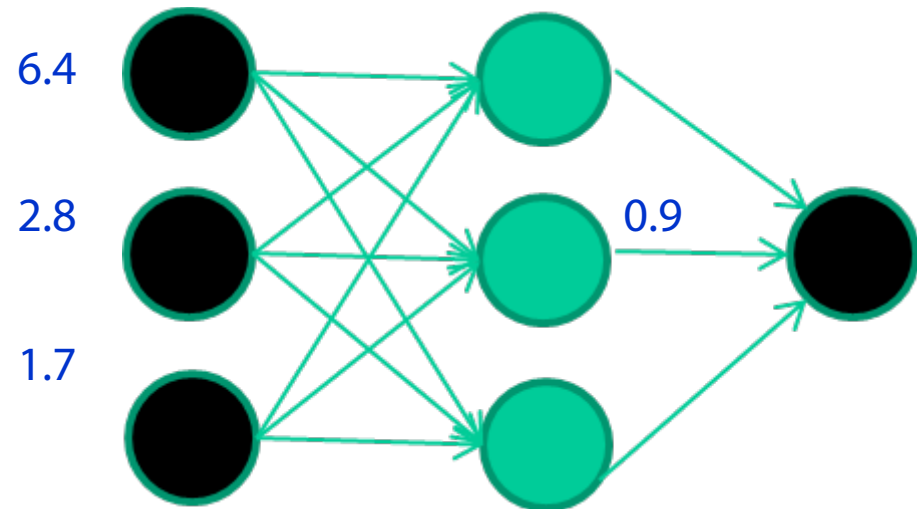
## Präsentierung eines Trainingsdatensatzes



# Trainingsdaten

<i>Fields</i>	<i>class</i>
1.4 2.7 1.9	0
3.8 3.4 3.2	0
6.4 2.8 1.7	1
4.1 0.1 0.2	0
etc ...	

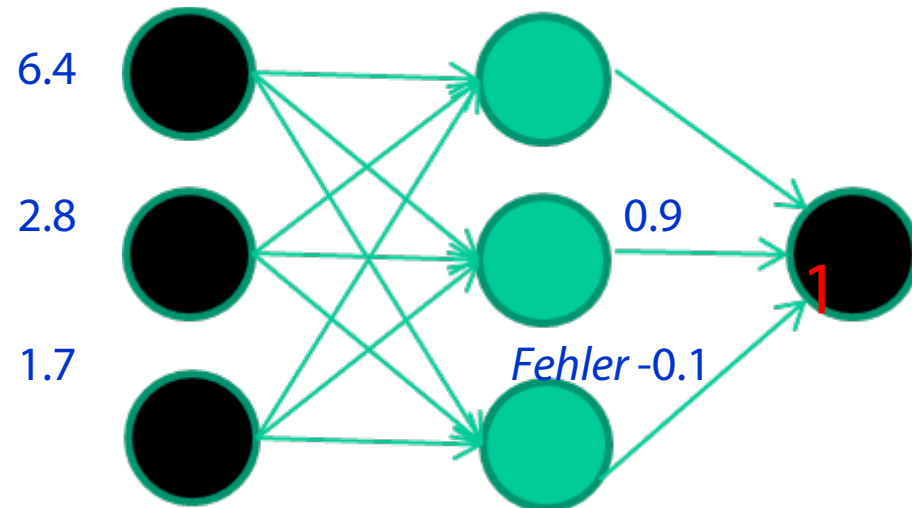
## Durchpropagierung zur Ausgabe



# Trainingsdaten

<i>Fields</i>	<i>class</i>
1.4 2.7 1.9	0
3.8 3.4 3.2	0
6.4 2.8 1.7	1
4.1 0.1 0.2	0
etc ...	

## Vergleich mit der Zielausgabe

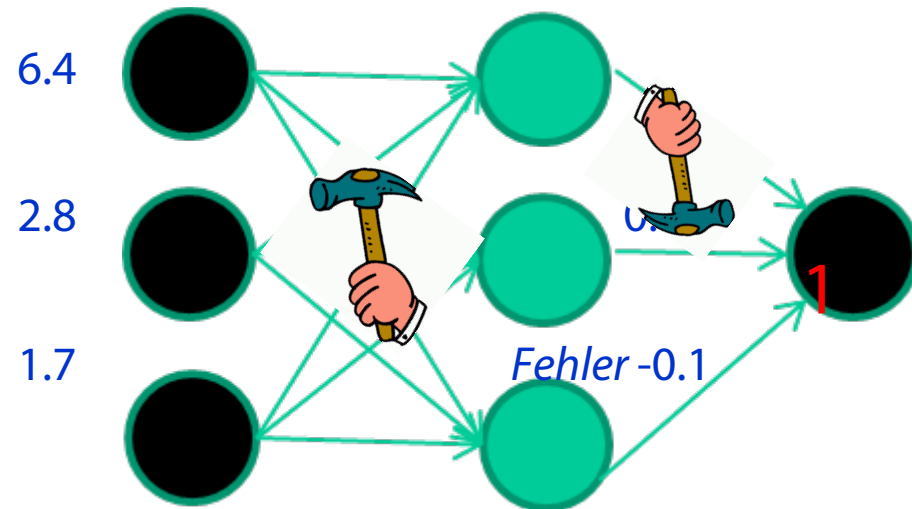




# Trainingsdaten

<i>Fields</i>	<i>class</i>
1.4 2.7 1.9	0
3.8 3.4 3.2	0
6.4 2.8 1.7	1
4.1 0.1 0.2	0
etc ...	

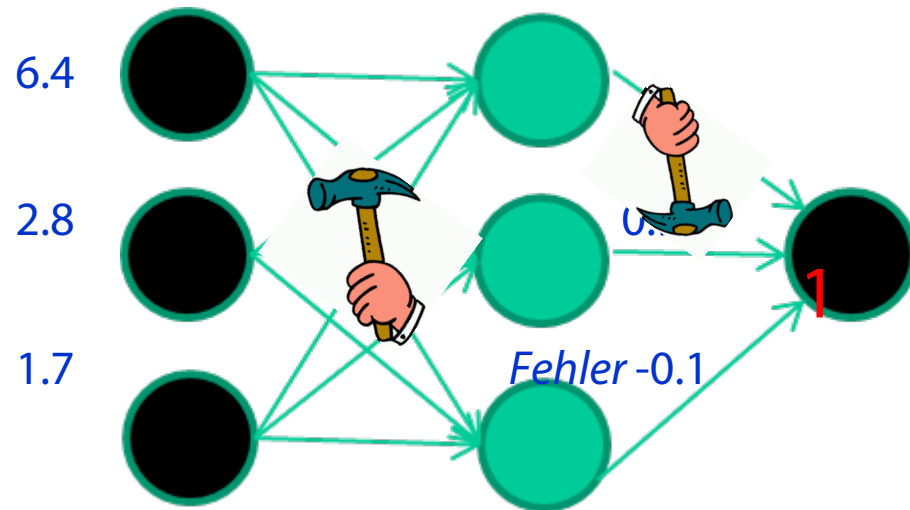
## Anpassen der Gewichte gemäß Fehler



# Trainingsdaten

<i>Fields</i>	<i>class</i>
1.4 2.7 1.9	0
3.8 3.4 3.2	0
6.4 2.8 1.7	1
4.1 0.1 0.2	0
etc ...	

Und so weiter ....



Wiederhole tausend-, vielleicht millionenmal – jedesmal mit einer zufälligen Trainingsinstanz und einer kleinen Anpassung der Gewichte

*Verfahren zur Gewichtsanzpassung müssen Fehler minimieren*

---

Differenzierbare Programmierung

# **LERNEN VON GEWICHTEN (EINE EBENE)**

# Eine Ebene: Perzeptron-Lernregel (Delta-Lernregel)

---

$$w_i \leftarrow w_i + \Delta w_i$$

wobei

$$\Delta w_i \leftarrow \eta(t - o)x_i$$

und

- $t = c(\vec{x})$  der Zielwert ist
- $o$  ist die Ausgabe des Perzeptrons
- $\eta$  ist eine kleine Konstante (z.B. 0,1): die Lernrate

Gewichte häufig nur nach Verarbeitung eines ganzen Datensatzes  $D$  angepasst

# Begründung für die Delta-Regel

---

- Idee: minimiere den quadratischen Fehler
  - $D$  Trainingsmenge
  - $t_d$  Wert für  $d \in D$
  - $o_d$  Ausgabe für  $d$

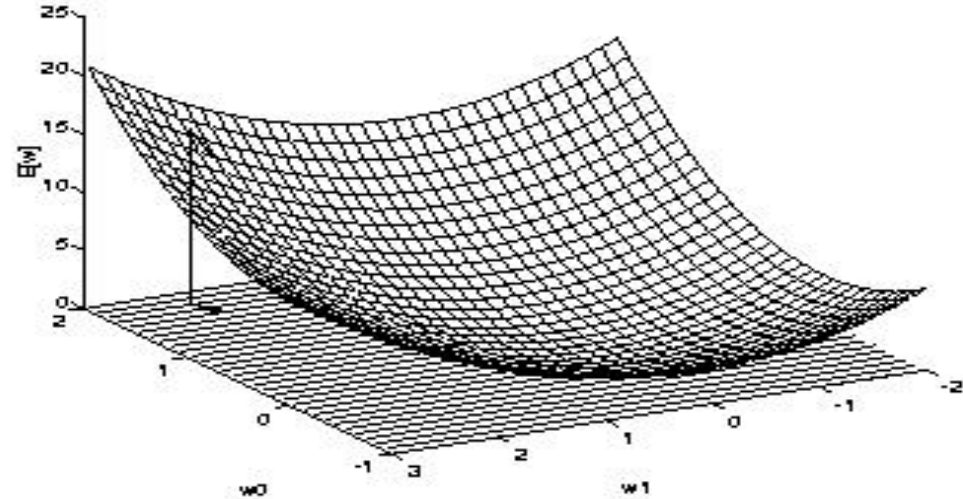
$$E[\vec{w}] \equiv \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

- Minimumbestimmung mit 1. Ableitung

# Absteigender Gradient

- Gradient

$$\nabla E[\vec{w}] \equiv \left[ \frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$$

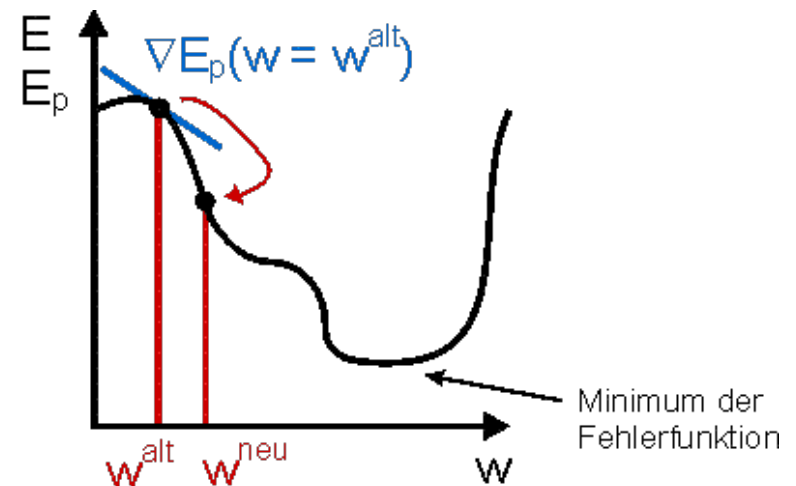


- Lernregel

$$\Delta \vec{w} = -\eta \nabla E[\vec{w}]$$

- i.e.

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$



# Gradient

$$\begin{aligned}\frac{\partial E}{\partial w_i} &= \frac{\partial}{\partial w_i} \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2 \\ &= \frac{1}{2} \sum_{d \in D} \frac{\partial}{\partial w_i} (t_d - o_d)^2 \\ &= \frac{1}{2} \sum_{d \in D} 2(t_d - o_d) \frac{\partial}{\partial w_i} (t_d - o_d) \\ &= \sum_{d \in D} (t_d - o_d) \frac{\partial}{\partial w_i} (t_d - \vec{w} \cdot \vec{x}_d) \\ &= \sum_{d \in D} (t_d - o_d) (-x_{i,d}) = - \sum_{d \in D} (t_d - o_d) x_{i,d}\end{aligned}$$

# Absteigender Gradient (Forts.)

- Gradient

$$\nabla E[\vec{w}] \equiv \left[ \frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$$

$$\frac{\partial E}{\partial w_i} = - \sum_{d \in D} (t_d - o_d) x_{i,d}$$

- Lernregel

$$\Delta \vec{w} = -\eta \nabla E[\vec{w}]$$

- i.e.

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$

Iterativ pro Datenpunkt  $d$

$$\Delta w_i = \eta (t - o) x_i$$

$$\Delta w_i = \eta \sum_{d \in D} (t_d - o_d) x_{i,d}$$



# Algorithmus für *eine* Gewichtsanzpassung

- Jedes Trainingsbeispiel sei ein Paar  $\langle \vec{x}, t \rangle$ 
  - $\vec{x}$  ist ein Inputvektor (Vektor mit Feature-Werten)
  - $t$  ist der Zielwert (Klassenlabel/Ausgabewert)
  - $\eta$  ist die Lernrate
- Initialisiere jedes  $w_i$  mit einem beliebigen, kleinen Wert
- Bis die Abbruchbedingung erfüllt ist (z.B. Anzahl an verarbeiteten Trainingsbeispielen):
  - Initialisiere jedes  $\Delta w_i$  mit 0
  - Für jedes Trainingsbeispiel  $\langle \vec{x}, t \rangle$ 
    - Berechne  $o_t$
    - Für jedes Gewicht  $w_i$ : 
$$\Delta w_i \leftarrow \Delta w_i + \eta(t - o_t)x_i$$
  - Für jedes  $w_i$ : 
$$w_i \leftarrow w_i + \Delta w_i$$

# Perzeptron-Lernregel

---

Man kann zeigen, dass der Vorgang konvergiert, ...

- ... wenn die Daten linear separierbar sind
- ... und  $\eta$  genügend klein gewählt wird

Schon früher untersucht:

D. Hebb: *The organization of behavior. A neuropsychological theory.*  
Erlbaum Books, Mahwah, N.J., **1949**

Netzwerke daher von manchen  
als künstliche neuronale Netze bezeichnet

Später für mehrschichtige Netze erweitert (Deep Learning):  
Fehlerrückführung durch mehrere Ebenen (Backpropagation)

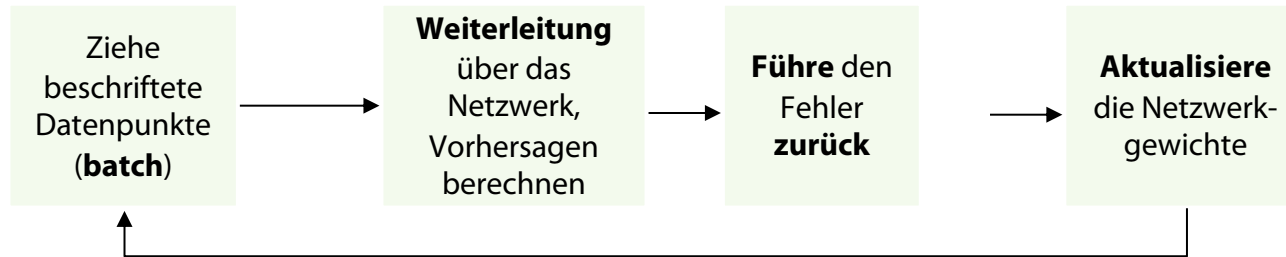
---

Differenzierbare Programmierung

# **LERNEN VON GEWICHTEN (MEHRE EBENEN)**



# Backpropagating Idee



## Backpropagation idee

- Erzeuge ein **Fehlersignal**, das die Differenz zwischen Prognosen und Zielwerten misst
- Verwenden Sie das Fehlersignal, um die Gewichte zu ändern und genauere Vorhersagen rückwärts zu erhalten
- Zugrunde liegende Mathematik: Kettenregel

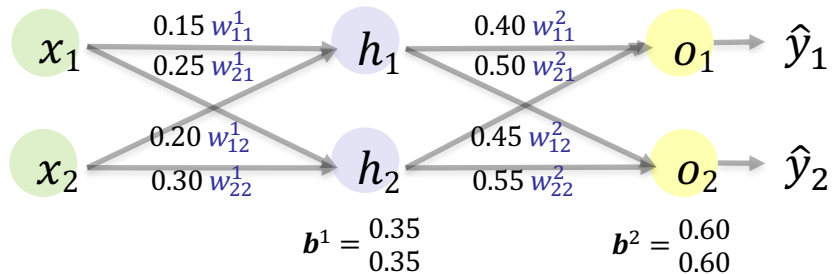
Kettenregel (1-dim)

$$\frac{dh}{dx} = \frac{df}{dg} \frac{dg}{dx} \quad (\text{für } h(x) = f(g(x)))$$

Idee: D. Hebb: *The organization of behavior. A neuropsychological theory.* Erlbaum Books, Mahwah, N.J., 1949

Für neuronale Netze: D. E. Rumelhart, G. E. Hinton, and R. J. Williams. *Learning representations by back-propagating errors*, 1986

# Beispiel: Backward Pass für $E_{total}$



- $\hat{y}_1 = out_{o_1} = 0.75136507$   
vs  $target_{o_1} = 0.01$
- $\hat{y}_2 = out_{o_2} = 0.772928465$   
vs  $target_{o_2} = 0.99$

Aktivierungsfkt.:

$$\sigma_1, \sigma_2 = \frac{1}{1 + e^{-x}}$$

Fehlerfkt.:

$E_{total}$

$$= \sum \frac{1}{2} (target - out)^2$$

Eingabe:

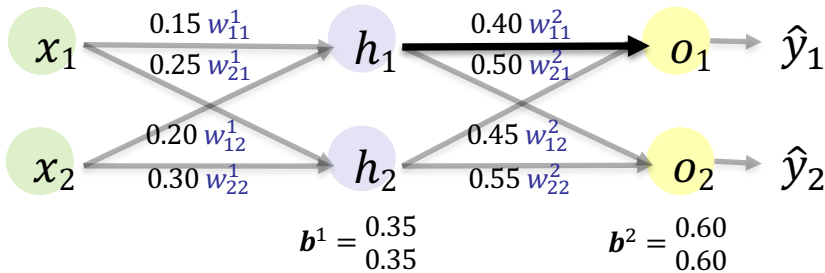
$$x = \begin{matrix} 0.05 \\ 0.10 \end{matrix}$$

Zielausgabe:

$$target = \begin{matrix} 0.01 \\ 0.99 \end{matrix}$$

- $E_{o_1} = \frac{1}{2} (target_{o_1} - out_{o_1})^2$   
 $= \frac{1}{2} (0.01 - 0.75136507)^2 = 0.27481108$
- $E_{o_2} = 0.023560026$
- $E_{total} = E_{o_1} + E_{o_2} = 0.298371109$

# Beispiel: Backward Pass für $w_{11}^2$



- $\frac{\partial E_{total}}{\partial E_{w_{11}^2}} = \frac{\partial E_{total}}{\partial out_{o_1}} \cdot \frac{\partial out_{o_1}}{\partial net_{o_1}} \cdot \frac{\partial net_{o_1}}{\partial E_{w_{11}^2}}$
- $\frac{\partial E_{total}}{\partial out_{o_1}} = 2 \cdot \frac{1}{2} (target_{o_1} - out_{o_1})^{2-1} \cdot (-1) + 0$   
 $= out_{o_1} - target_{o_1} = 0.75136507 - 0.01$   
 $= 0.74136507$

Aktivierungsfkt

$$\sigma_1, \sigma_2 = \frac{1}{1 + e^{-x}}$$

Eingabe:

$$x = \begin{matrix} 0.05 \\ 0.10 \end{matrix}$$

Zielausgabe:

$$target = \begin{matrix} 0.01 \\ 0.99 \end{matrix}$$

Fehlerfkt:

$$E_{total}$$

$$= \sum \frac{1}{2} (target - out)^2$$

- $\frac{\partial out_{o_1}}{\partial net_{o_1}} = out_{o_1} \cdot (1 - out_{o_1}) = 0.186815602$
- $\frac{\partial net_{o_1}}{\partial E_{w_{11}^2}} = 1 \cdot out_{h_1}$
- $\frac{\partial E_{total}}{\partial E_{w_{11}^2}} = 0.74136507 \cdot 0.186815602 \cdot out_{h_1}$

Welcher Fehler sollte als nächstes propagiert werden?

$\delta_{o_1}$   $out_{h_1}$

- $new\_w_{11}^2 = w_{11}^2 - \frac{\partial E_{total}}{\partial E_{w_{11}^2}} = 0.317832959$

# Backpropagation Algorithmus (Eine Runde)

1. Eingabe: Initialisieren des Eingabevektors  $x = \mathit{out}^0$
2. Feedforward: Für  $i = 1, 2, \dots, M$   
$$\mathit{net}^i = \mathbf{W}^{(i)} \mathit{out}^{i-1} + \mathbf{b}_i \text{ and } \mathit{out}^i = \sigma_i(\mathit{net}^i)$$

3. Berechnen Sie den Fehler für die letzte Ebene

$$\delta^M = \nabla_{\hat{y}} E \odot \sigma'(\mathit{net}^M)$$

4. Fehler durch propagieren: Für  $i = M-1, M-2, \dots,$

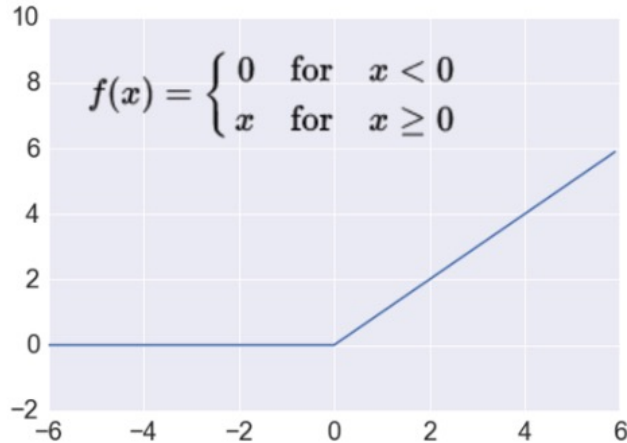
$$\delta^i = (\mathbf{w}^{i+1})^\top \delta^{i+1} \odot \sigma'(\mathit{net}^i)$$

5. Berechnen der Gradienten:

$$\frac{\partial E}{\partial w_{jk}^i} = \mathit{out}_k^{i-1} \delta_j^i \quad \text{und} \quad \frac{\partial E}{\partial b_j^i} = \delta_j^i$$

Hadamard product  $\odot$ :  $\begin{pmatrix} a \\ b \end{pmatrix} \odot \begin{pmatrix} c \\ d \end{pmatrix} = \begin{pmatrix} ac \\ bd \end{pmatrix}$

# Aktivierung: ReLU (rectified linear unit)



<http://adilmoujahid.com/images/activation.png>

Nimmt eine reellwertige Zahl und legt einen Schwellenwert bei Null fest  $f(x) = \max(0, x)$

$$\mathbb{R}^n \rightarrow \mathbb{R}_+^n$$

Die meisten Deep Netzwerke verwenden heutzutage ReLU

Trainiert viel **schneller**

- beschleunigt die Konvergenz der SGD
- aufgrund linearer, nicht beschränkter Form

Weniger aufwändige Berechnungen

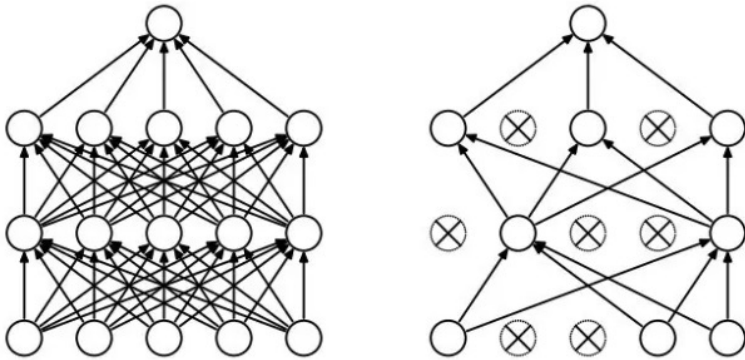
- im Vergleich zu Sigmoid/Tanh (Exponentiale etc.)
- Implementierung durch einfaches Schwellenwert einer Matrix bei Null

**Ausdrucksstarker**

Beugt dem **gradient vanishing problem** vor



# Regularization



## Dropout

- Zufälliges Ablegen von Einheiten (zusammen mit ihren Verbindungen) während des Trainings
- Jeder Knoten wird mit fester Wahrscheinlichkeit  $p$  beibehalten, unabhängig von anderen Knoten
- **Hyperparameter**  $p$  zu wählen (abgestimmt)

Srivastava, Nitish, et al. "[Dropout: a simple way to prevent neural networks from overfitting.](#)" *Journal of machine learning research* (2014)

## L2 = Gewichtszerrfall (weight decay)

- Regularisierungsterm, der große Gewichte bestraft, werden hinzugefügt
- Der Gewichtszerrfallwert bestimmt, wie dominant die Regularisierung während der Gradientenberechnung ist
- Großer Gewichtszerrfallwert  $\rightarrow$  Große Strafe für große Gewichte

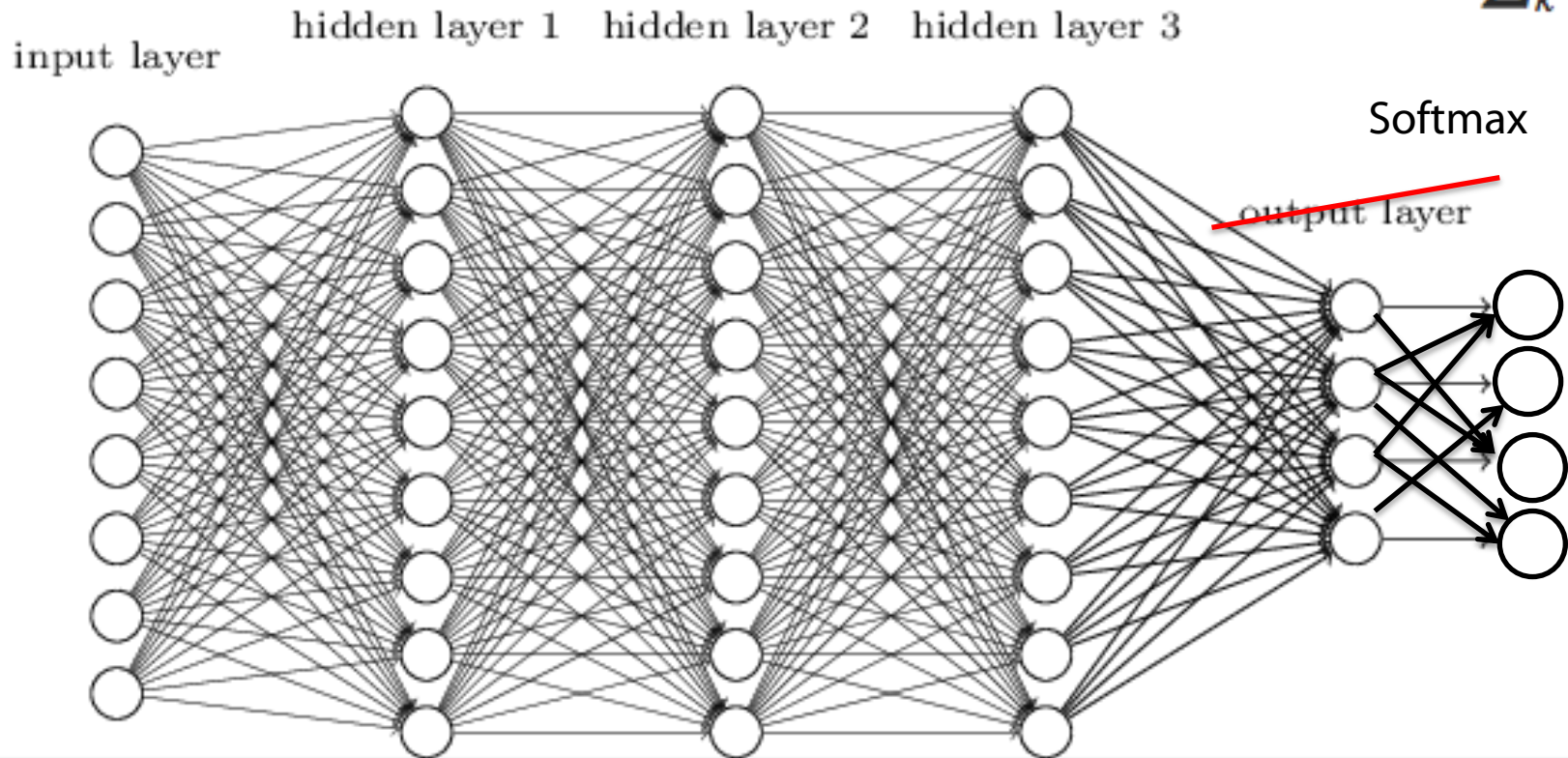
$$J_{reg}(\theta) = J(\theta) + \lambda \sum_k \theta_k^2$$

## Early-stopping

- Verwenden eines Validierungsfehlers, um zu entscheiden, wann das Training beendet werden soll
- Stoppe, wenn sich die überwachte Menge nach  $n$  aufeinanderfolgenden Epochen nicht verbessert hat
- $n$  wird auch patience (Geduld?) genannt

# Softmax Ausgabeschicht

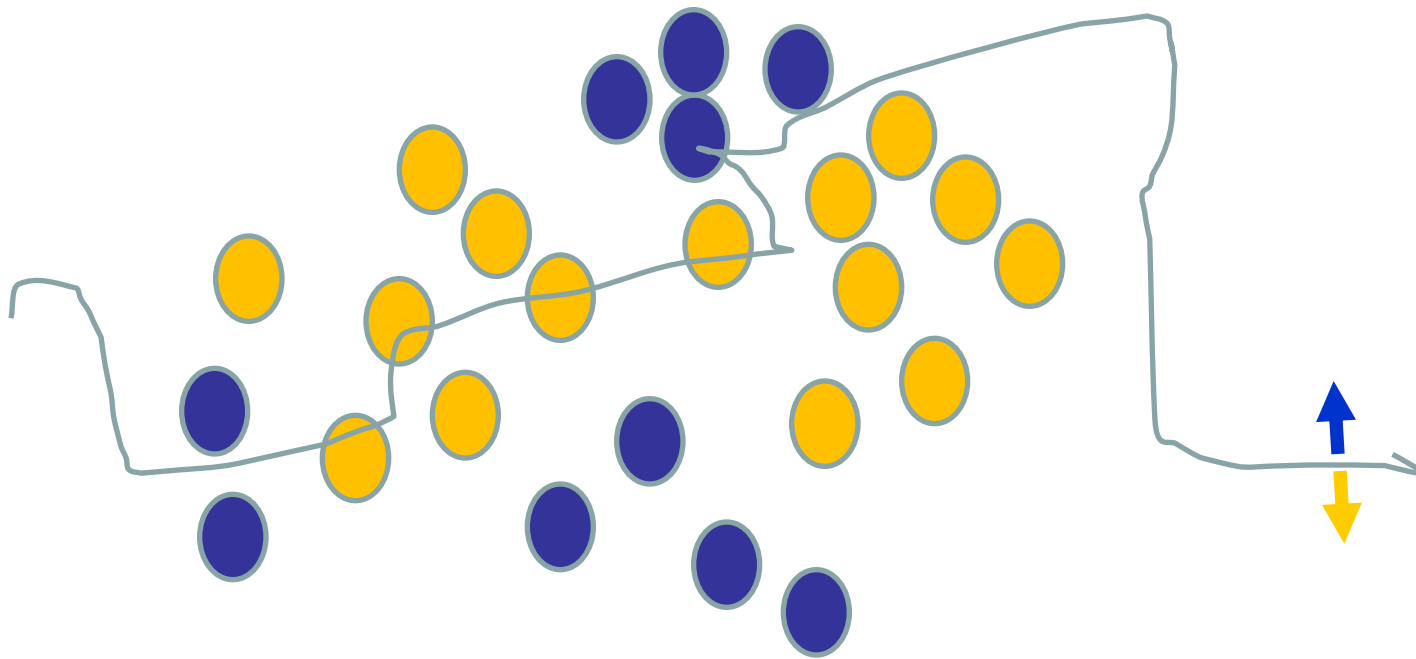
$$a_j^L = \frac{e^{z_j^L}}{\sum_k e^{z_k^L}}$$



Netzwerk gibt eine Wahrscheinlichkeitsverteilung aus!

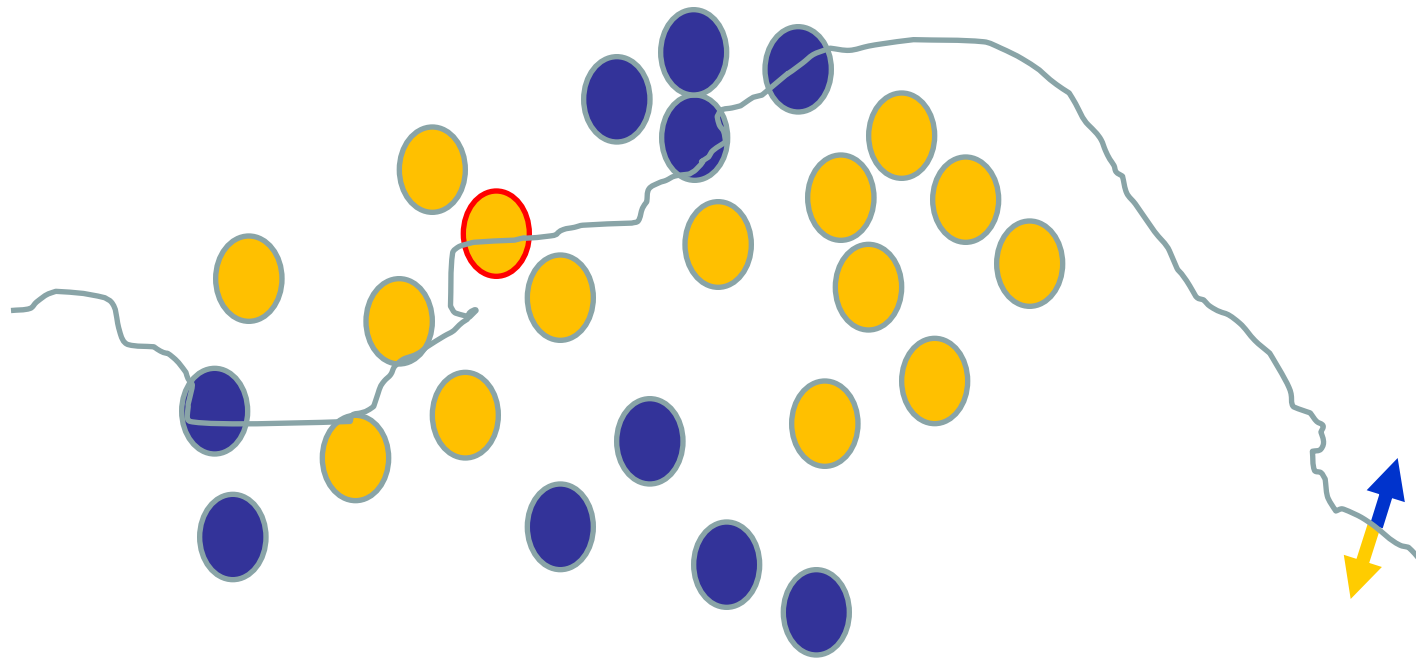
# Perspektive der Entscheidungsgrenzenanpassung

## Zufällige Initialgewichte



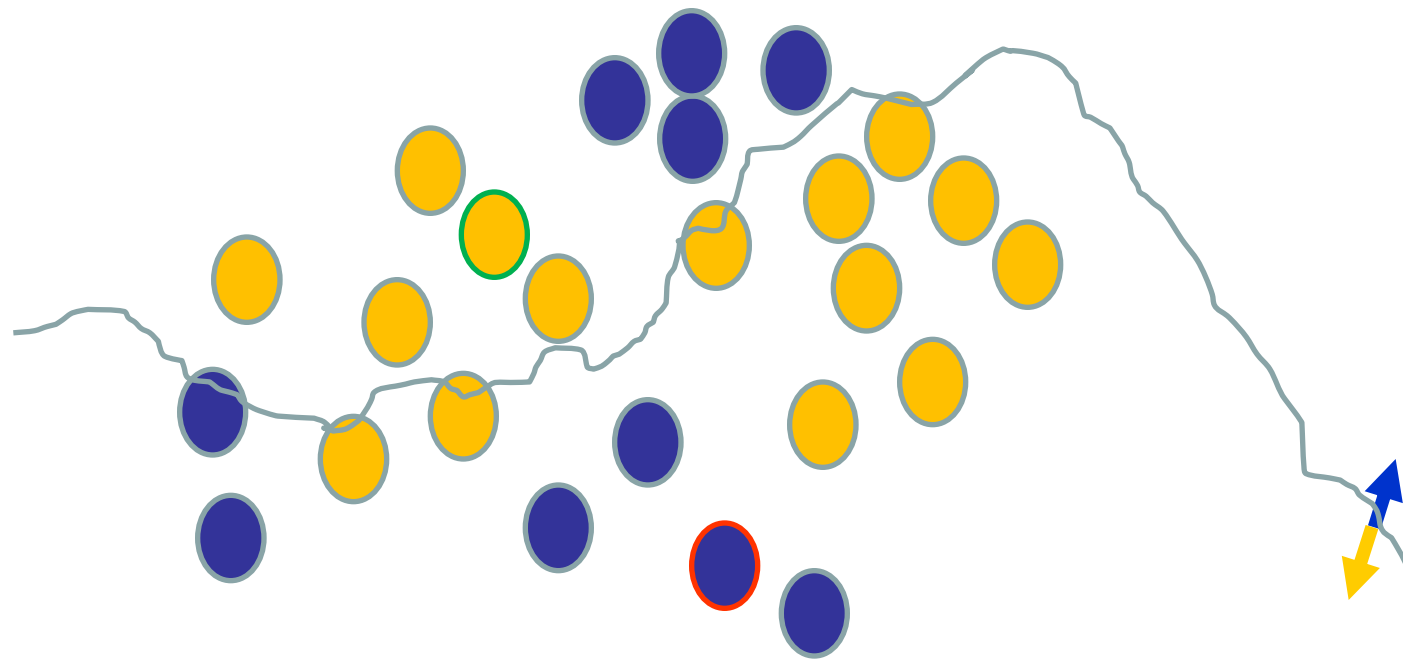
# Perspektive der Entscheidungsgrenzenanpassung

Verwenden einer Trainingsinstanz / Anpassung der Gewichte



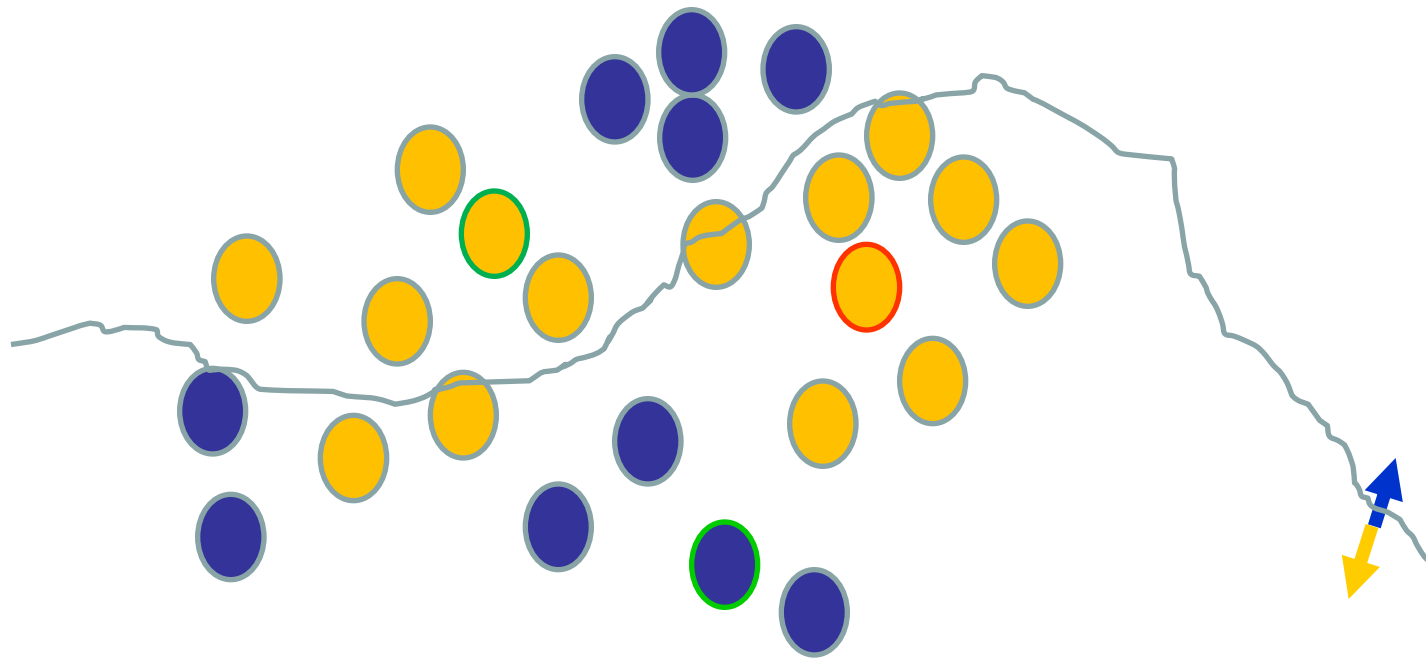
# Perspektive der Entscheidungsgrenzenanpassung

Verwenden einer Trainingsinstanz / Anpassung der Gewichte



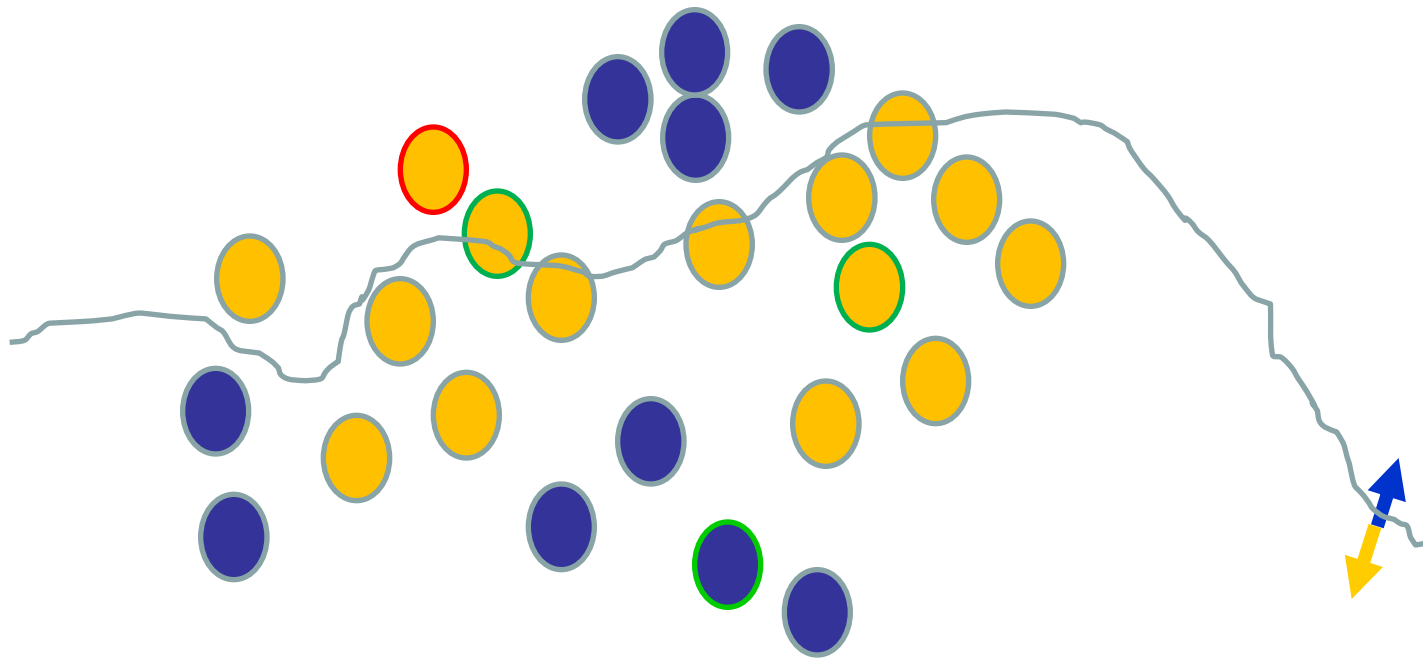
# Perspektive der Entscheidungsgrenzenanpassung

Verwenden einer Trainingsinstanz / Anpassung der Gewichte



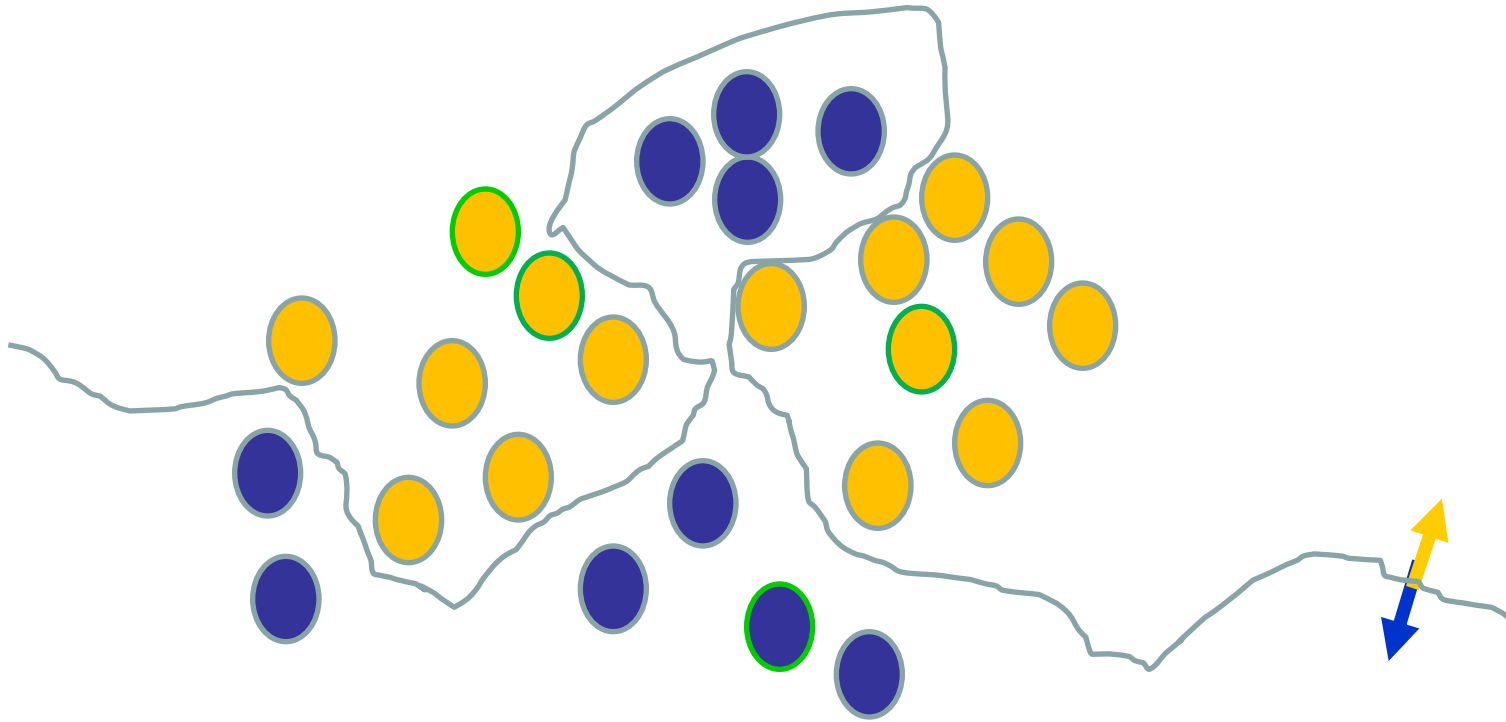
# Perspektive der Entscheidungsgrenzenanpassung

Verwenden einer Trainingsinstanz / Anpassung der Gewichte



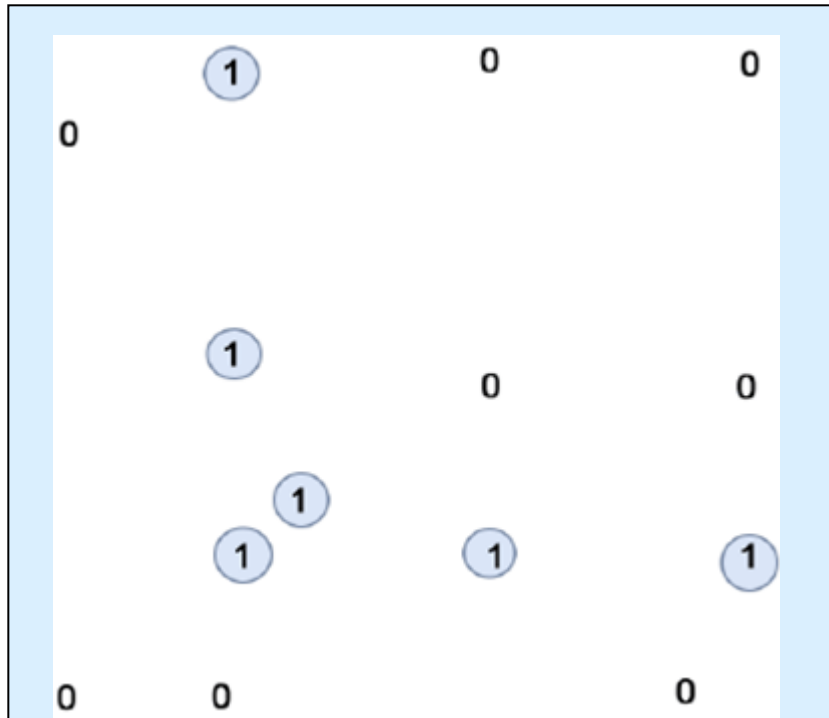
# Perspektive der Entscheidungsgrenzenanpassung

Und schließlich....

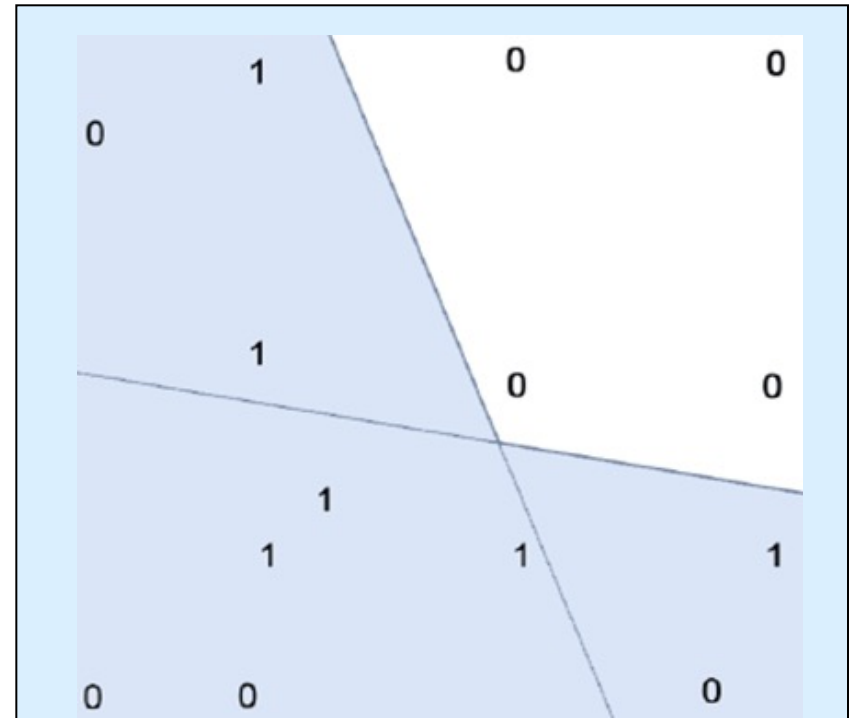




# Fehler bei einer ungeeigneten Trainingsmenge



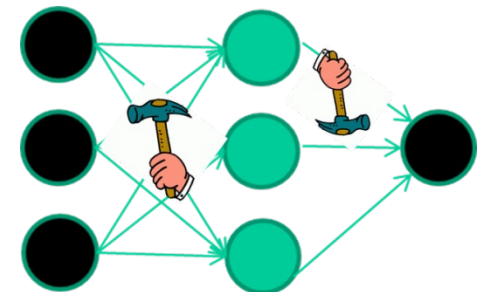
Zu viele Trainingsbeispiele:  
das Netz hat „auswendig gelernt“



Zu wenig Trainingsbeispiele:  
keine richtige Klassifikation

# Lernverfahren: Resümee

- Epoche: einmalige Verwendung des gesamten Datensatzes
  - Manchmal zu groß
- Batch: Zerlegung eines Datensatzes in Teilmengen
  - Für jede Teilmenge: Passe Gewichte an
  - Anzahl der Teilmengen: #Iterationen
- Tausende von kleinen Anpassungen, jede macht das Netz besser für die letzte Eingabe (aber vielleicht schlechter für frühere Eingaben)
- Verwende mehrere Epochen
- Wieviele Epochen? Batches?
  - Ausprobieren
- Durch verdammt gutes Glück kommt oft eine Funktion heraus, die gut genug in Anwendungen funktioniert



---

Differenzierbare Programmierung

# APPLICATIONS



# Word-Word Associations in Document Retrieval

---

Recap bag-of-words approaches

- Client profiles, TF-IDF

**Words are not independent** of each other

Need to represent some aspects of **word semantics**

# Point(wise) Mutual Information: PMI

- **Measure of association** used in information theory and statistics

$$\text{pmi}(x; y) \equiv \log \frac{p(x, y)}{p(x)p(y)} = \log \frac{p(x|y)}{p(x)} = \log \frac{p(y|x)}{p(y)}$$

- Positive PMI:  $\text{PPMI}(x, y) = \max(\text{pmi}(x, y), 0)$
- Quantifies the discrepancy between the **probability of their coincidence** given their **joint distribution** and their **individual distributions**, assuming independence
- **Finding collocations and associations between words**
- **Countings** of occurrences and co-occurrences of words in a text corpus can be used to **approximate the probabilities**  $p(x)$  or  $p(y)$  and  $p(x,y)$  respectively

# PMI – Example

word 1	word 2	count word 1	count word 2	count of co-occurrences	PMI
puerto	rico	1938	1311	1159	10.0349081703
hong	kong	2438	2694	2205	9.72831972408
los	angeles	3501	2808	2791	9.56067615065
carbon	dioxide	4265	1353	1032	9.09852946116
prize	laureate	5131	1676	1210	8.85870710982
san	francisco	5237	2477	1779	8.83305176711
nobel	prize	4098	5131	2498	8.68948811416
ice	hockey	5607	3002	1933	8.6555759741
star	trek	8264	1594	1489	8.63974676575
car	driver	5578	2749	1384	8.41470768304
it	the	283891	3293296	3347	-1.72037278119
are	of	234458	1761436	1019	-2.09254205335
this	the	199882	3293296	1211	-2.38612756961
is	of	565679	1761436	1562	-2.54614706831
and	of	1375396	1761436	2949	-2.79911817902
a	and	984442	1375396	1457	-2.92239510038
in	and	1187652	1375396	1537	-3.05660070757
to	and	1025659	1375396	1286	-3.08825363041
to	in	1025659	1187652	1066	-3.12911348956
of	and	1761436	1375396	1190	-3.70663100173

- Counts of pairs of words getting the **most and the least PMI scores** in the first 50 millions of words in **Wikipedia** (dump of October 2015)
- Filtering by 1,000 or more co-occurrences.
- The frequency of each count can be obtained by dividing its value by 50,000,952. (Note: natural log is used to calculate the PMI values in this example, instead of log base 2)

# PMI – Co-occurrence Matrix

Add-2 Smoothed Count( $w, context$ )

	computer	data	pinch	result	sugar
apricot	2	2	3	2	3
pineapple	2	2	3	2	3
digital	4	3	2	3	2
information	3	8	2	6	2

PPMI( $w, context$ )

	computer	data	pinch	result	sugar
apricot	-	-	2.25	-	2.25
pineapple	-	-	2.25	-	2.25
digital	1.66	0.00	-	0.00	-
information	0.00	0.57	-	0.47	-



# Embedding Approaches to Word Semantics

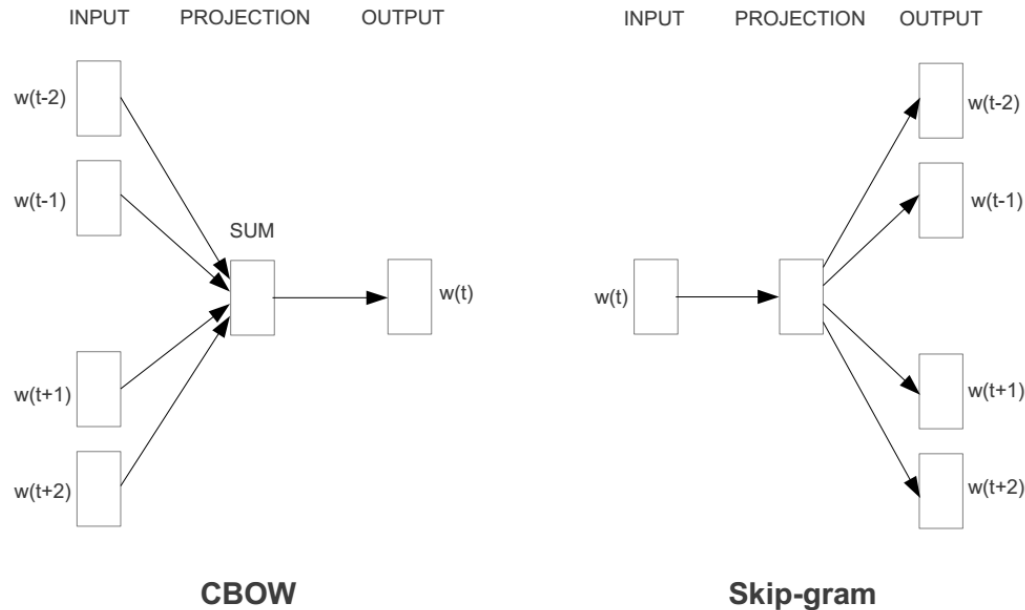
---

- Represent each word with a low-dimensional vector
- Word similarity = vector similarity
- Key idea: Predict surrounding words of every word



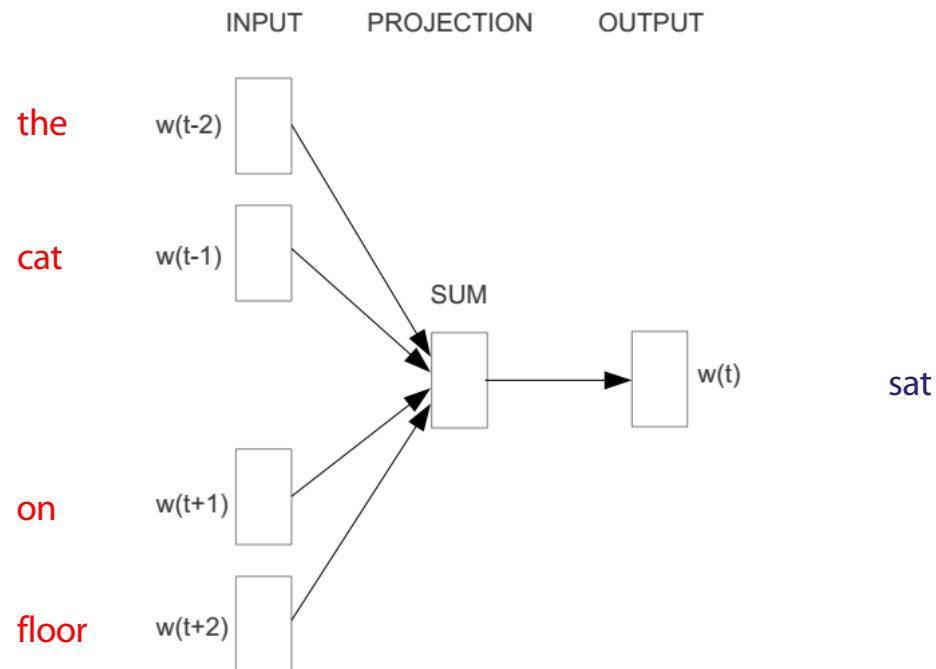
# Represent the meaning of **words** – word2vec

- 2 basic structural models:
  - **Continuous Bag of Words (CBOW)**: use a window of words to predict the middle word
  - **Skip-gram (SG)**: use a word to predict the surrounding ones in window.

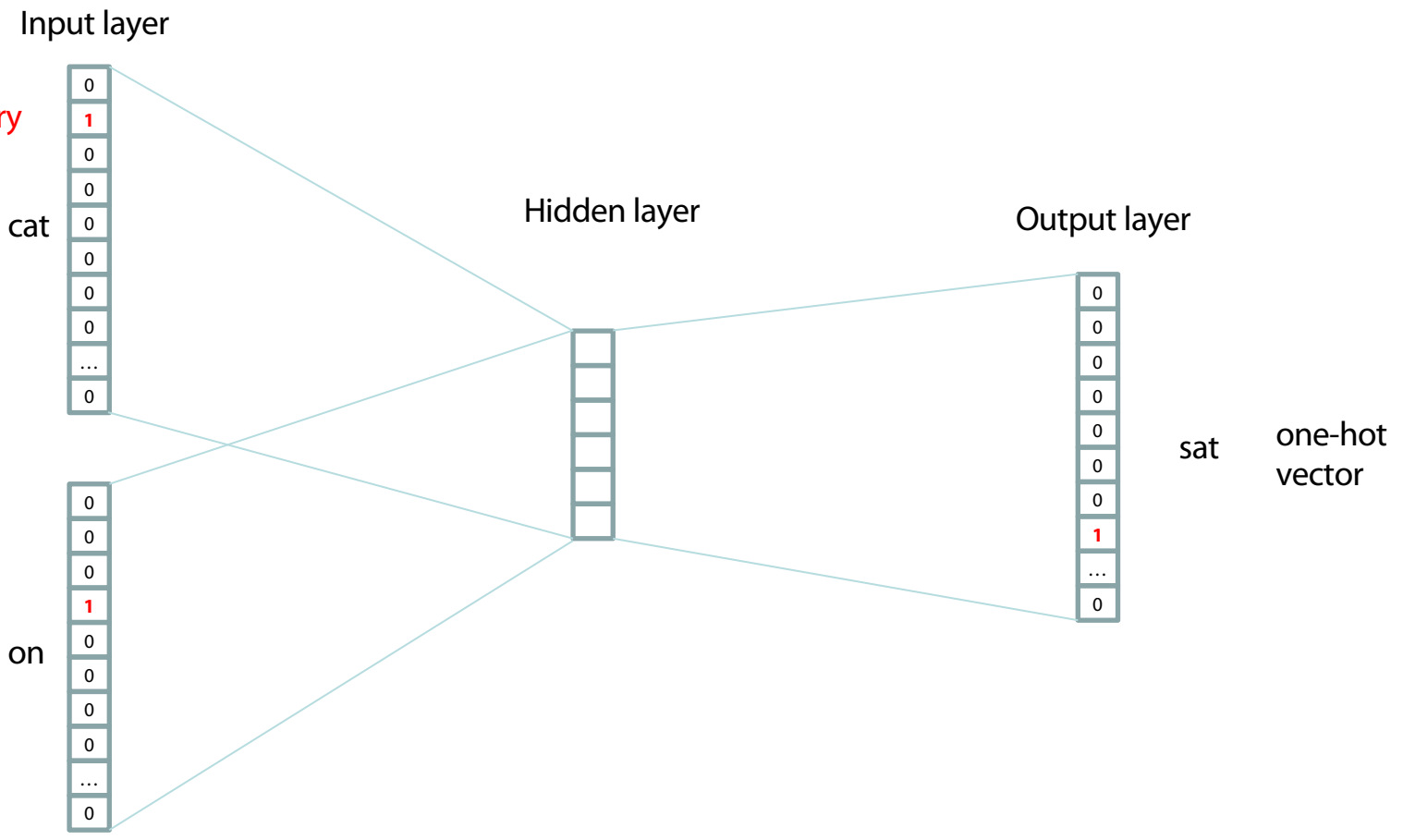


# Word2vec – Continuous Bag of Word

- E.g. “The cat <sat> on floor”
  - Window size = 2



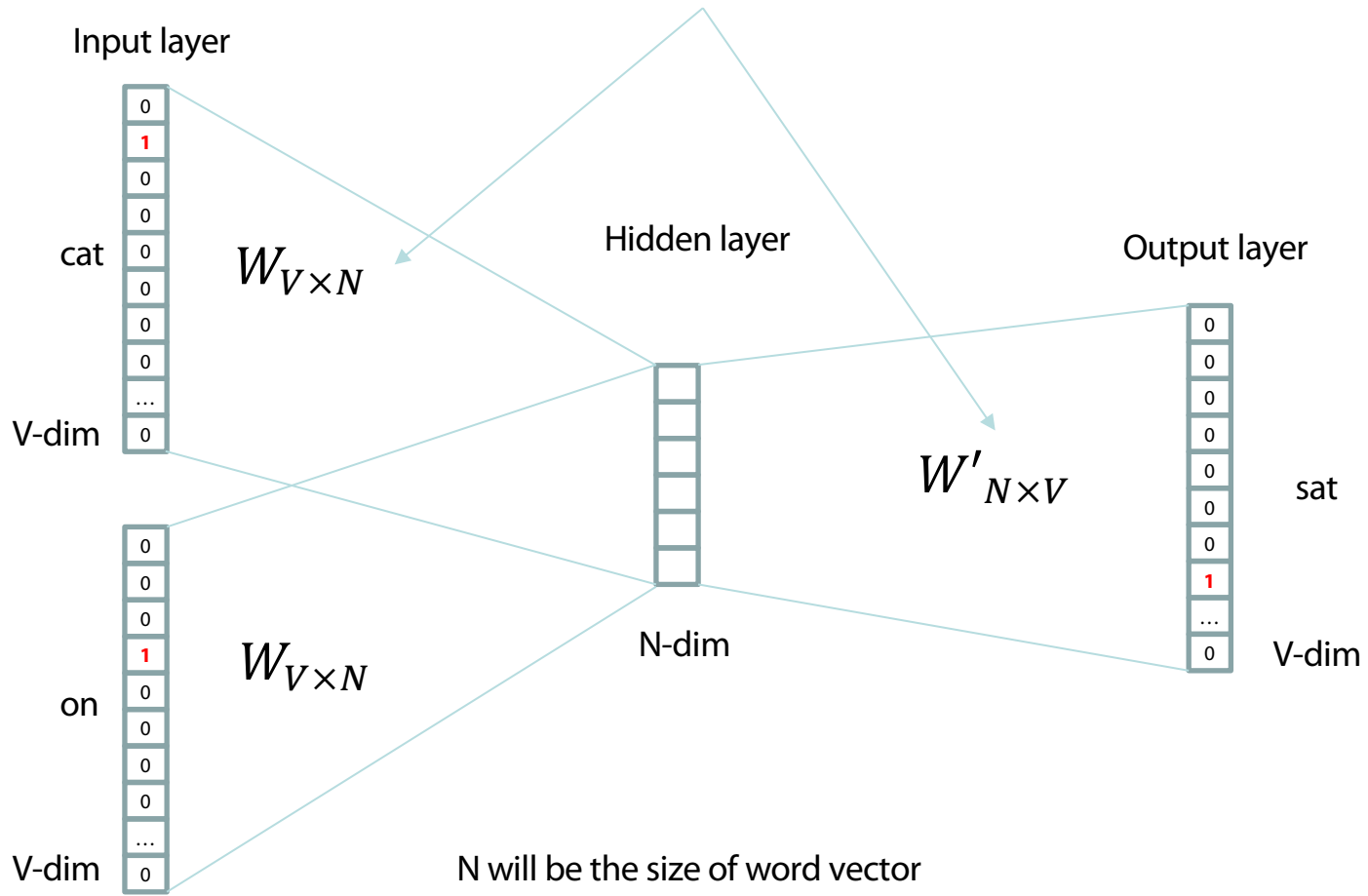
Index of cat in vocabulary



one-hot vector

sat one-hot vector

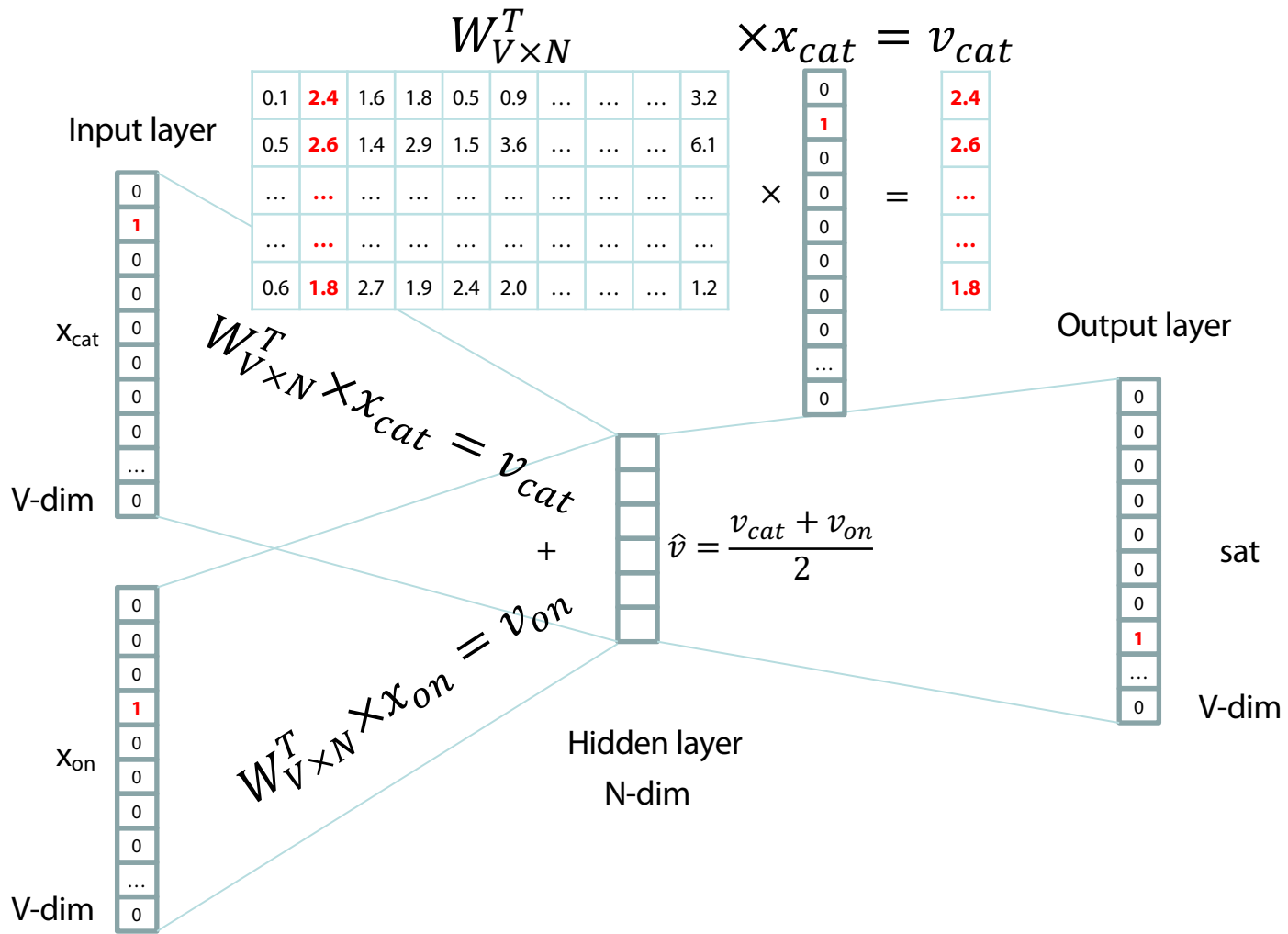
We must learn  $W$  and  $W'$

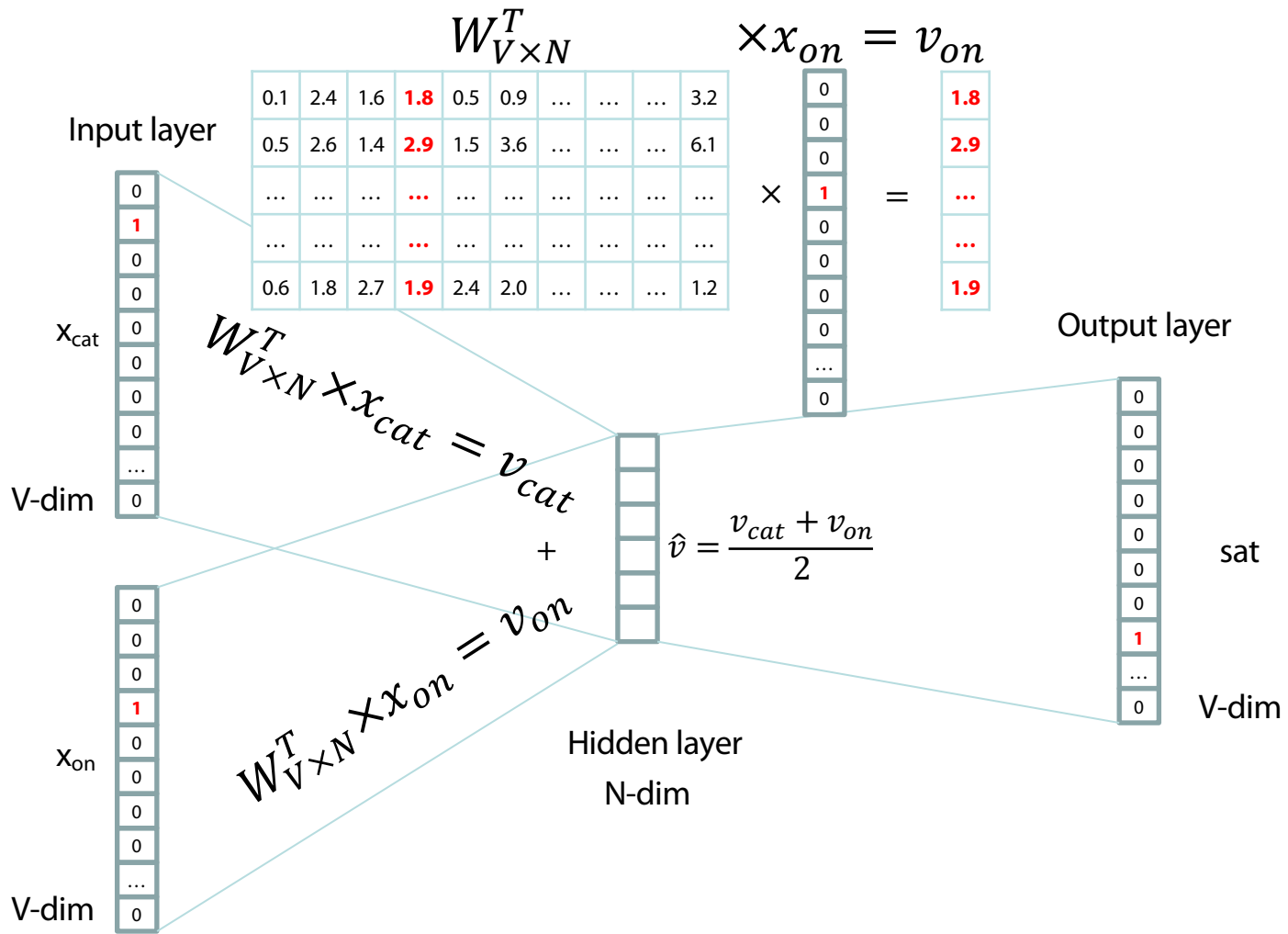


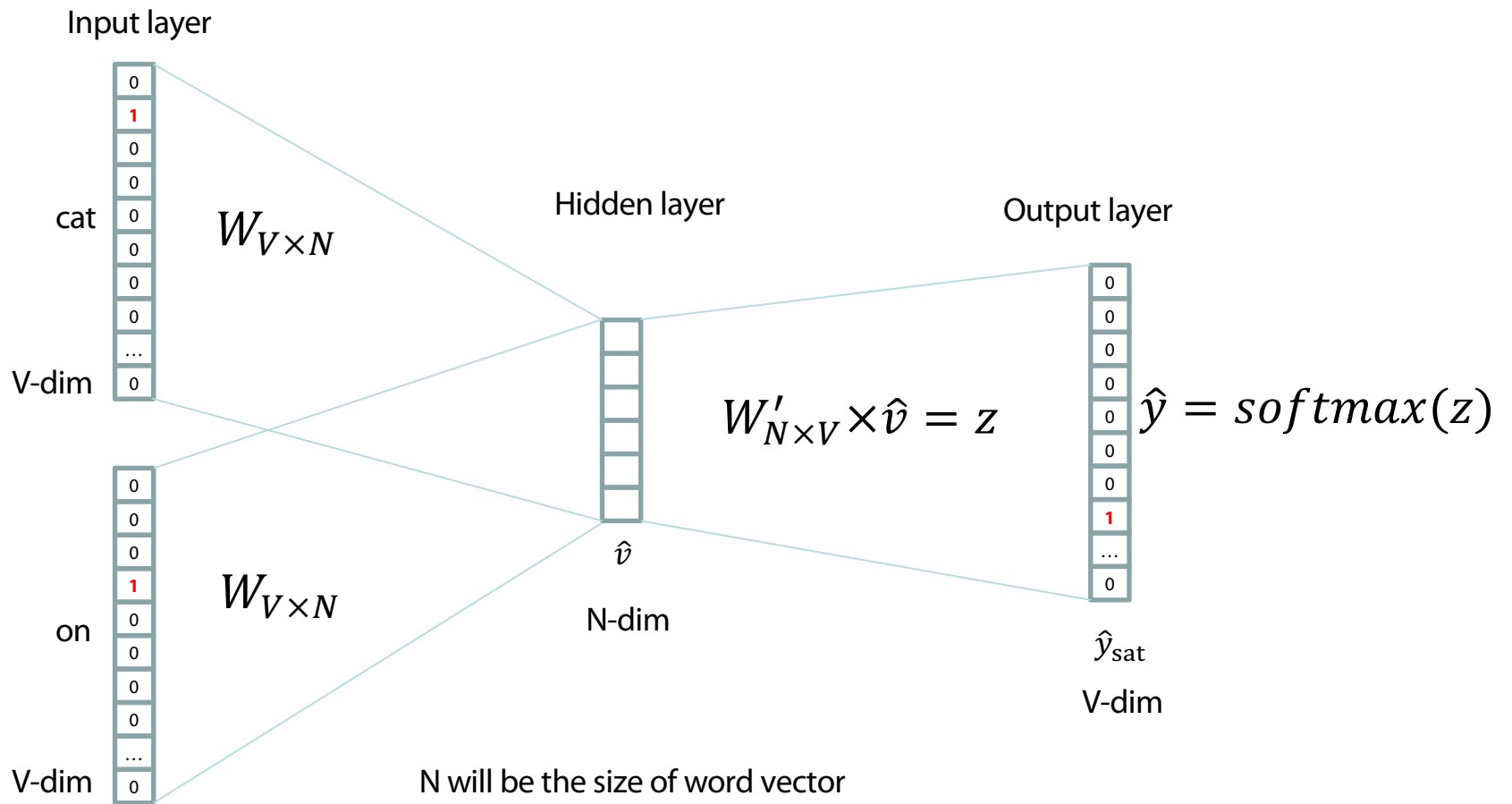
# Deep Learning

---

- Hidden layer represents feature space
  - Making explicit features in the data...
  - ... that are relevant for a certain task
- Determine features automatically
  - Learning suitable mappings into feature space
- Deep learning also known as representation learning









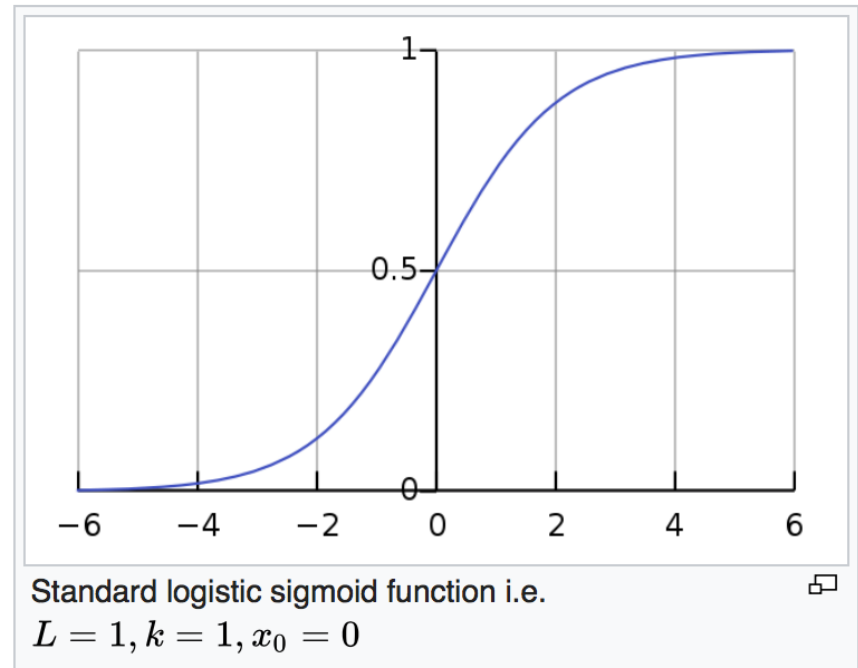
# Logistic function

A **logistic function** or **logistic curve** is a common "S" shape (**sigmoid curve**), with equation:

$$f(x) = \frac{L}{1 + e^{-k(x-x_0)}}$$

where

- $e$  = the **natural logarithm** base (also known as **Euler's number**),
- $x_0$  = the  $x$ -value of the sigmoid's midpoint,
- $L$  = the curve's maximum value, and
- $k$  = the steepness of the curve.<sup>[1]</sup>

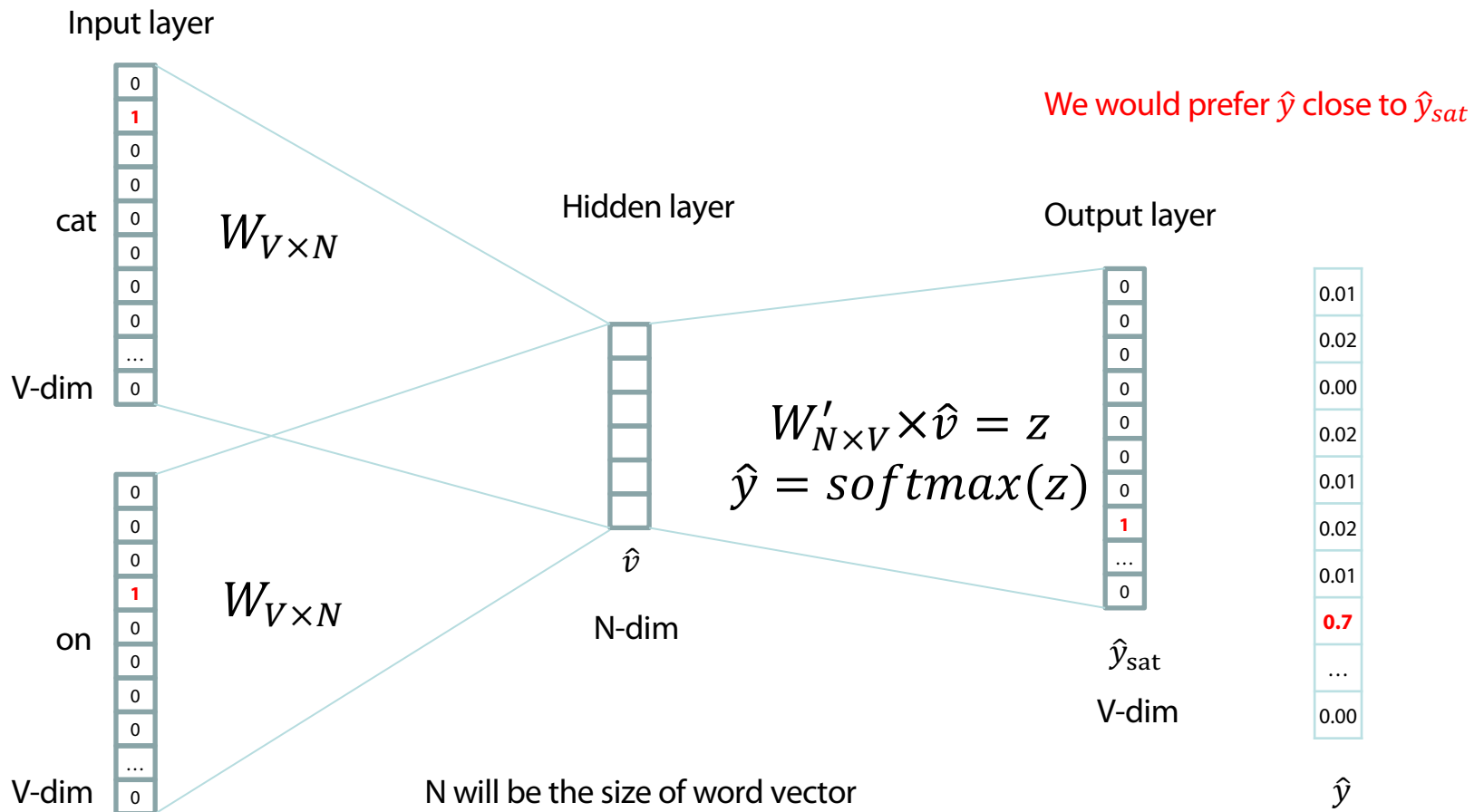


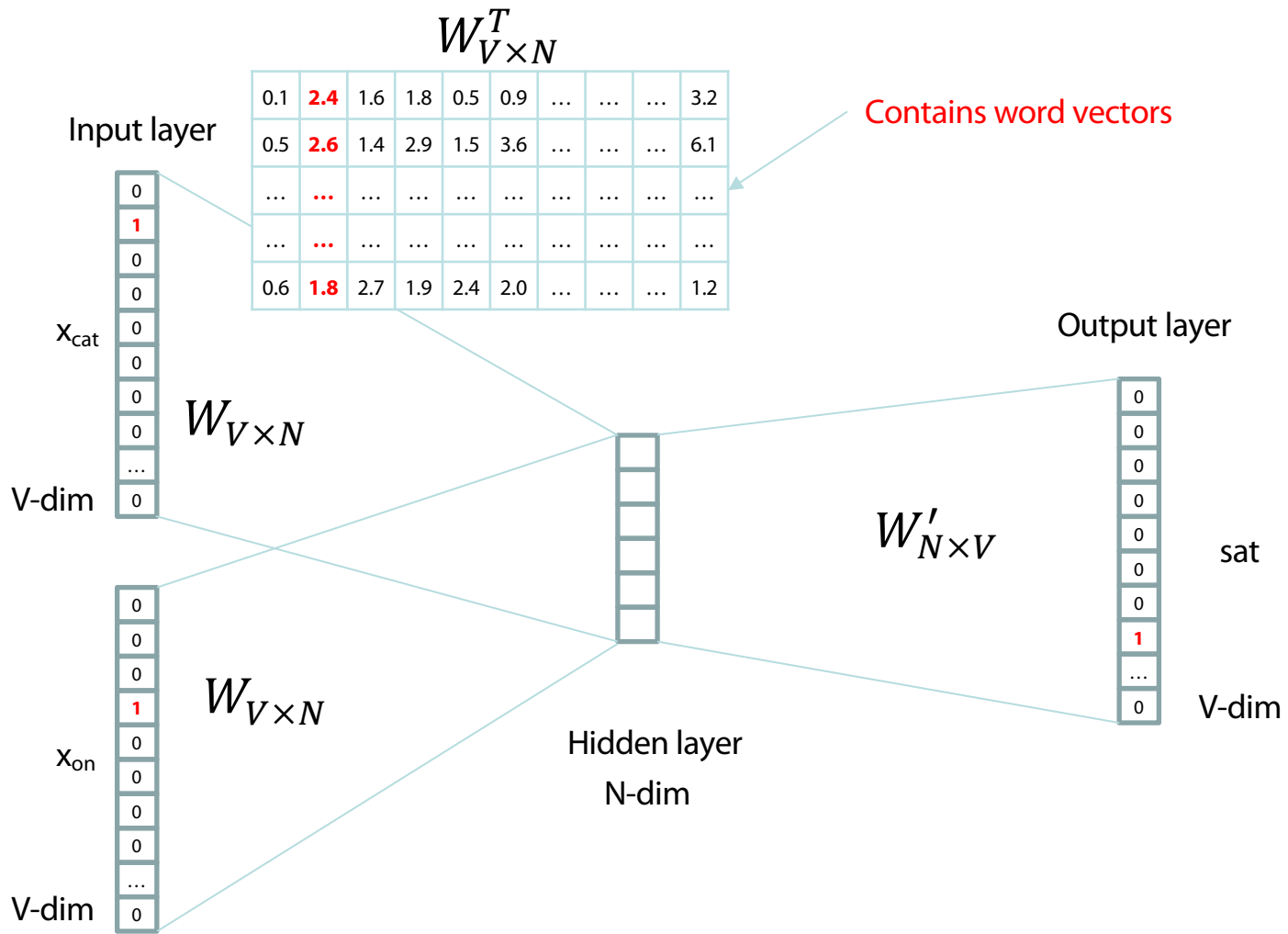
# softmax( $\mathbf{z}$ )

The **softmax function**, or **normalized exponential function**, is a generalization of the **logistic function** that "squashes" a  $K$ -dimensional vector  $\mathbf{z}$  of arbitrary real values to a  $K$ -dimensional vector  $\sigma(\mathbf{z})$  of real values in the range  $[0, 1]$  that add up to 1. The function is given by

$$\sigma : \mathbb{R}^K \rightarrow [0, 1]^K$$
$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} \quad \text{for } j = 1, \dots, K.$$

In **probability theory**, the output of the softmax function can be used to represent a **categorical distribution** – that is, a **probability distribution** over  $K$  different possible outcomes.





Consider either  $W$  or  $W'$  as the word's representation.

# Word Analogies

Test for linear relationships, examined by Mikolov et al. (2014)

a:b :: c:?



$$d = \arg \max_x \frac{(w_b - w_a + w_c)^T w_x}{\|w_b - w_a + w_c\| \|w_x\|}$$

man:woman :: king:?

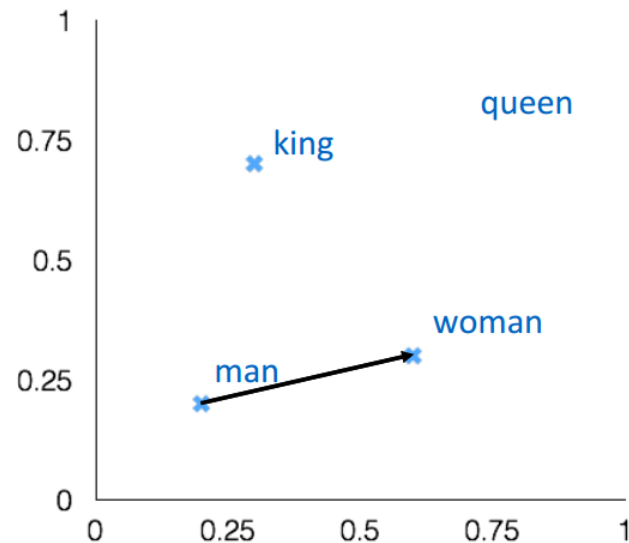
+ king [ 0.30 0.70 ]

- man [ 0.20 0.20 ]

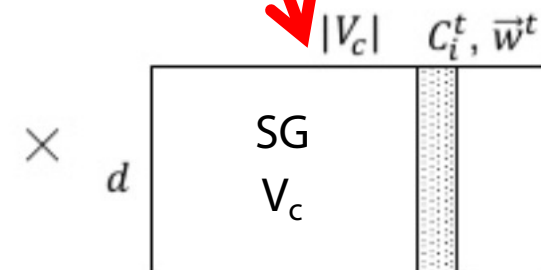
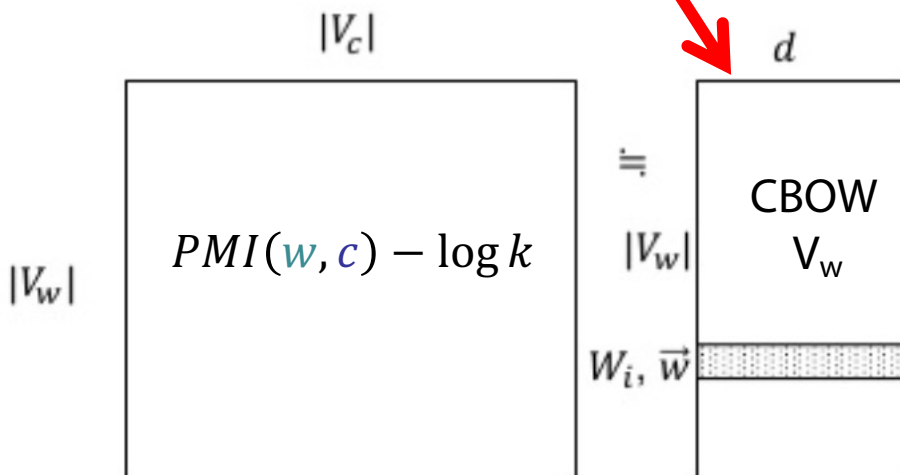
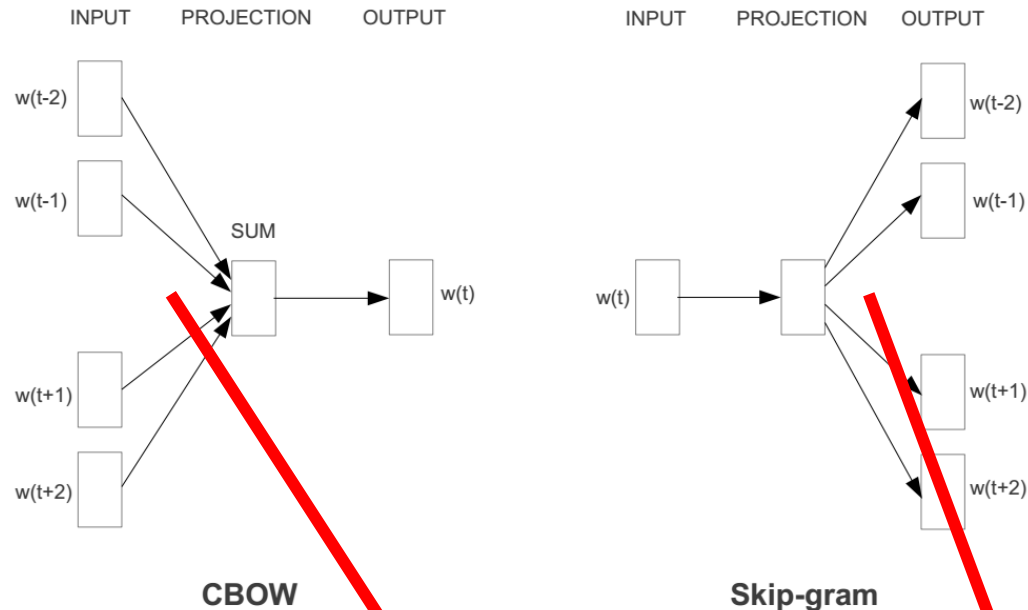
+ woman [ 0.60 0.30 ]

---

queen [ 0.70 0.80 ]



# The Picture: CBOW and Skip-Gram (SG)



“Neural Word Embeddings as Implicit Matrix Factorization”  
Levy & Goldberg, NeurIPS 2014



# Convolution

Input image



Convolution  
Kernel

$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$

Feature map



# Convolutional Neural Networks (CNNs)

Main CNN idea for text:

**Compute vectors for n-grams** and group them afterwards

Example: "this takes too long" compute vectors for:

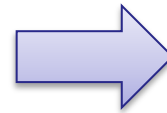
This takes, takes too, too long, this takes too, takes too long, this takes too long

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

Input matrix

1	0	1
0	1	0
1	0	1

Convolutional  
3x3 filter



1 <sub>x=-1</sub>	1 <sub>x=0</sub>	1 <sub>x=1</sub>	0	0
0 <sub>x=0</sub>	1 <sub>x=1</sub>	1 <sub>x=0</sub>	1	0
0 <sub>x=1</sub>	0 <sub>x=0</sub>	1 <sub>x=1</sub>	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

Convolved  
Feature

[http://deeplearning.stanford.edu/wiki/index.php/Feature\\_extraction\\_using\\_convolution](http://deeplearning.stanford.edu/wiki/index.php/Feature_extraction_using_convolution)



# Convolutional Neural Networks (CNNs)

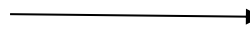
Main CNN idea for text:

Compute vectors for n-grams and **group them afterwards**

Feature Map

6	4	8	5
5	4	5	8
3	6	7	7
7	9	7	2

max pool  
2x2 filters  
and stride 2

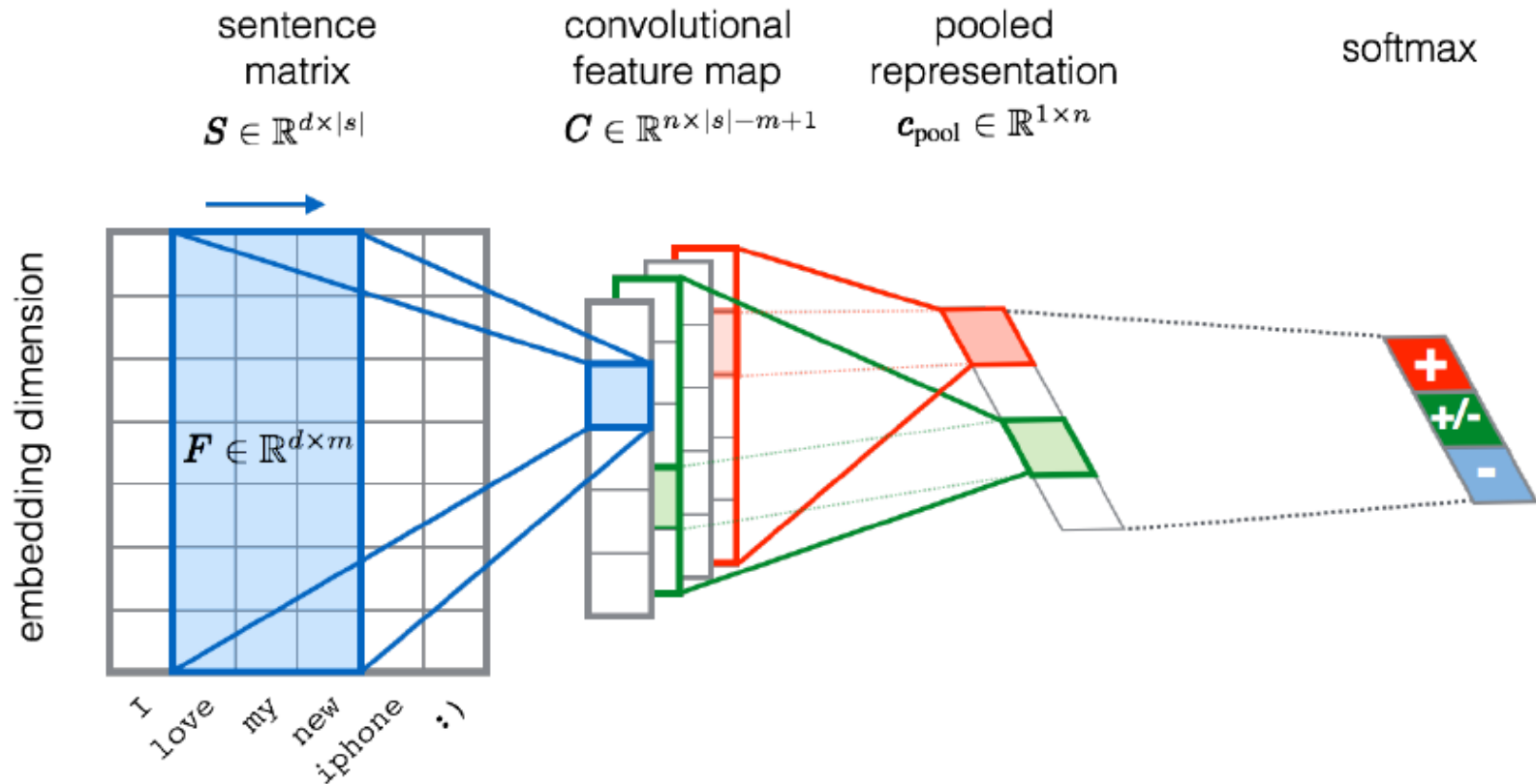


Max-Pooling


Dimensionsreduktion

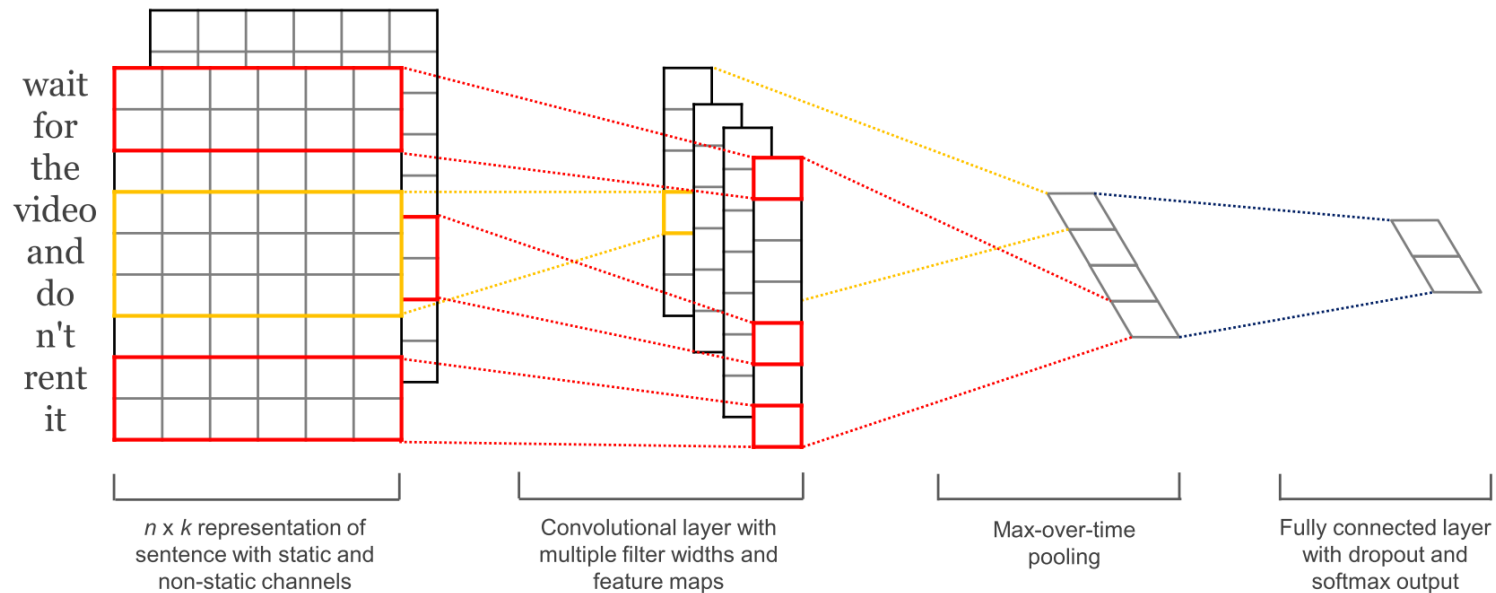
<https://shafeentejani.github.io/assets/images/pooling.gif>

# CNN for text classification



Severyn, Aliaksei, and Alessandro Moschitti. "UNITN: Training Deep Convolutional Neural Network for Twitter Sentiment Classification." *SemEval@NAACL-HLT*. 2015.

# CNN with multiple filters

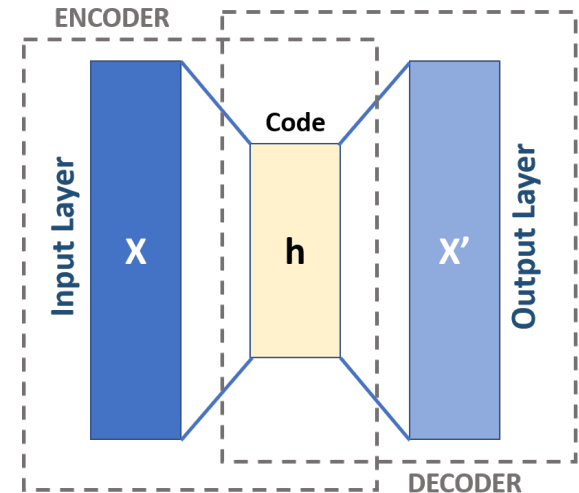


Kim, Y. "Convolutional Neural Networks for Sentence Classification", EMNLP (2014)

sliding over 3, 4 or 5 words at a time

# Differential Programming: Basic Idea

- Example: Computer Vision
  - Estimate position of light source
  - Use standard learning approach



- Develop render in appropriate programming language (e.g., Julia)
- Differentiate renderer (e.g., Zygote)
  - Use differentiated render for backpropagation

# Geschichtlicher Überblick

Anfänge

- **1943:** McCulloch und Pitts beschreiben und definieren eine Art erster neuronaler Netzwerke.
- **1949:** Formulierung der Hebb'schen Lernregel (nach Hebb)
- **1957:** Entwicklung des Perzeptrons durch Rosenblatt

Ernüchterung

- **1969:** Minsky und Papert untersuchen das Perzeptron mathematisch und zeigen dessen Grenzen, etwa beim XOR-Problem, auf.

Renaissance

- **1982:** Beschreibung der ersten selbstorganisierenden Netze (nach *biologischem Vorbild*) durch van der Malsburg und Kohonen und eines richtungweisenden Artikels von Hopfield, indem die ersten rückgekoppelten Netze (Hopfield-Netze, nach *physikalischen Vorbild*) beschrieben werden
- **1986:** Das Lernverfahren Backpropagation für mehrschichtige Perzeptrons wird entwickelt.

Boom

- **Ab 2000:** Deep Learning (Hinton, LeCun, Bengio, Ng, et al.)
- **Ab 2020:** Differentiable Programming (Lecun et al.)

