



Grundlagen der Künstlichen Intelligenz

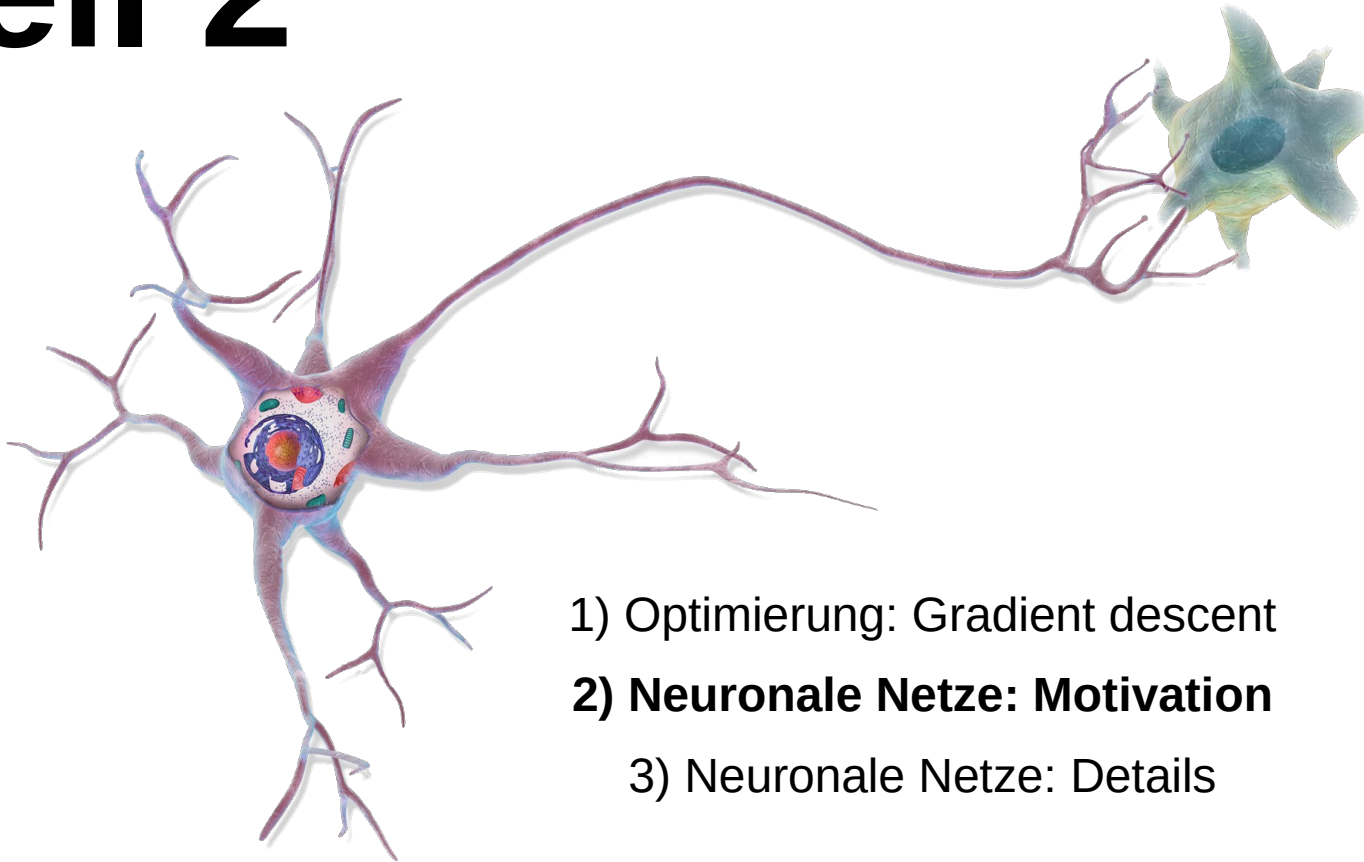
Wintersemester 25/26

Vorlesung 5

Neuronale Netze, Backpropagation Teil 2

Prof. Dr. Pascal Friederich
T.T.-Prof. Dr. Peer Nowack

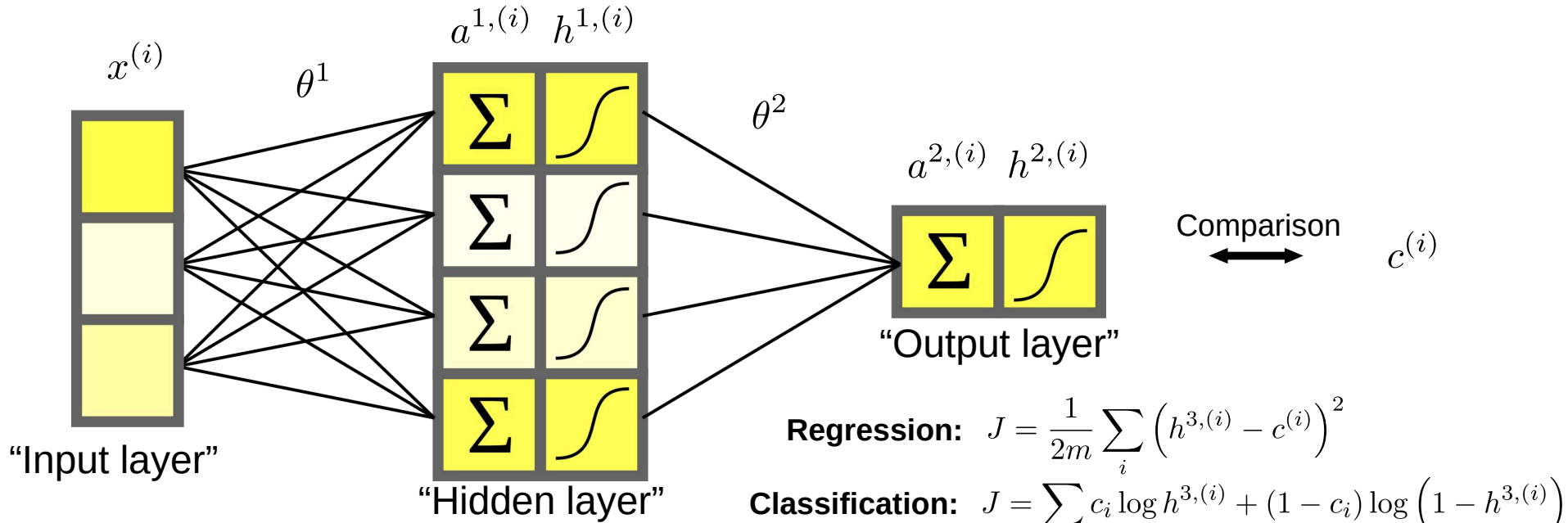
Teil 2



- 1) Optimierung: Gradient descent
- 2) Neuronale Netze: Motivation**
- 3) Neuronale Netze: Details

Neuronale Netze: Backpropagation

$$x^{(i)} \xrightarrow{\theta^1} a^{1,(i)} = (\theta^1)^T x^{(i)} \longrightarrow h^{1,(i)} = \sigma(a^{1,(i)}) \xrightarrow{\theta^2} a^{2,(i)} = (\theta^2)^T h^{1,(i)} \longrightarrow h^{2,(i)} = \sigma(a^{2,(i)})$$



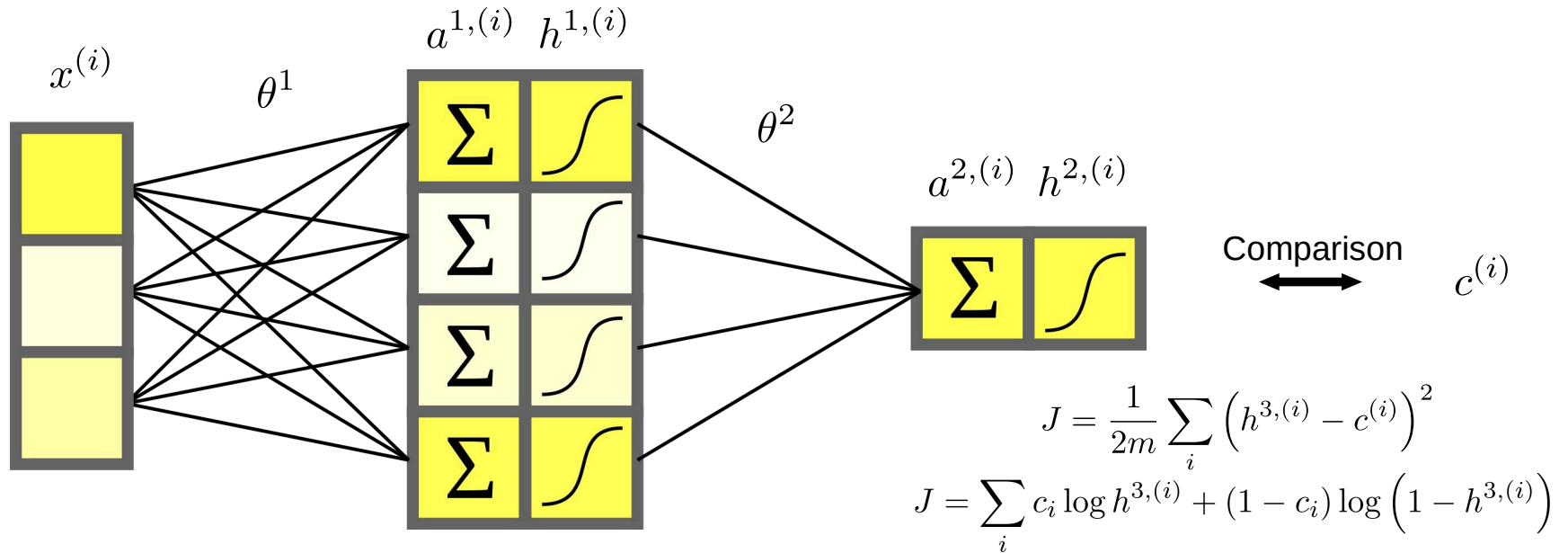
Regression: $J = \frac{1}{2m} \sum_i (h^{3,(i)} - c^{(i)})^2$

Classification: $J = \sum_i c_i \log h^{3,(i)} + (1 - c_i) \log (1 - h^{3,(i)})$

$$\theta^1 \rightarrow \theta^1 - \alpha \frac{\partial J}{\partial \theta^1}$$

$$\theta^2 \rightarrow \theta^2 - \alpha \frac{\partial J}{\partial \theta^2}$$

Neuronale Netze: Backpropagation



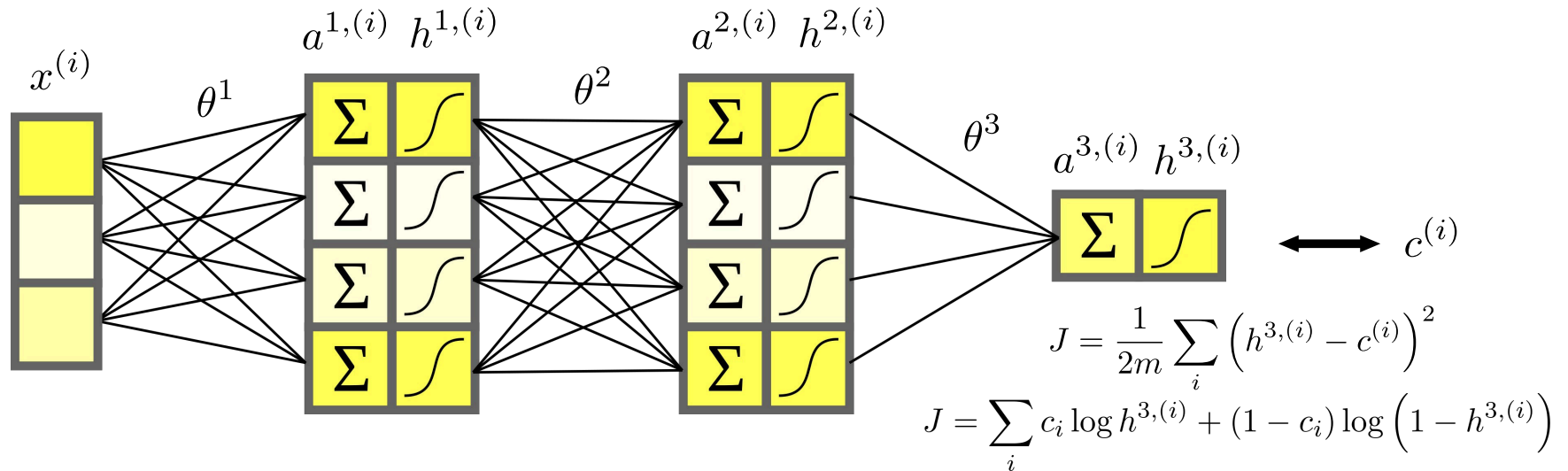
$$\frac{\partial J}{\partial \theta^1} = \frac{\partial J}{\partial h^2} \frac{\partial h^2}{\partial a^2} \frac{\partial a^2}{\partial h^1} \frac{\partial h^1}{\partial a^1} \frac{\partial a^1}{\partial \theta^1}$$

$$\theta^1 \rightarrow \theta^1 - \alpha \frac{\partial J}{\partial \theta^1}$$

$$\frac{\partial J}{\partial \theta^2} = \frac{\partial J}{\partial h^2} \frac{\partial h^2}{\partial a^2} \frac{\partial a^2}{\partial \theta^2}$$

$$\theta^2 \rightarrow \theta^2 - \alpha \frac{\partial J}{\partial \theta^2}$$

Neuronale Netze: Backpropagation



$$\frac{\partial J}{\partial \theta^1} = \frac{\partial J}{\partial h^3} \frac{\partial h^3}{\partial a^3} \frac{\partial a^3}{\partial h^2} \frac{\partial h^2}{\partial a^2} \frac{\partial a^2}{\partial h^1} \frac{\partial h^1}{\partial a^1} \frac{\partial a^1}{\partial \theta^1}$$

$$\theta^1 \rightarrow \theta^1 - \alpha \frac{\partial J}{\partial \theta^1}$$

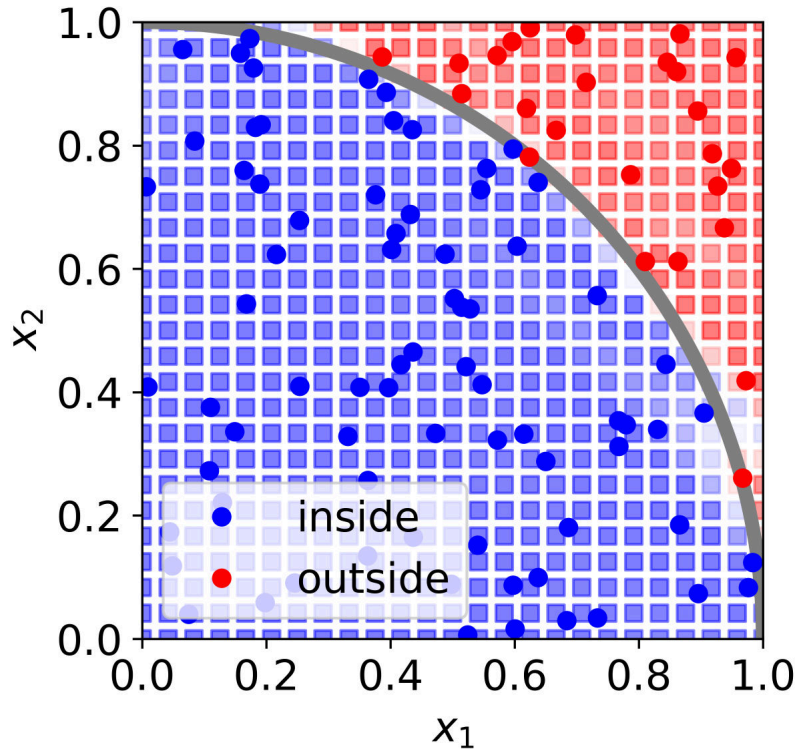
$$\frac{\partial J}{\partial \theta^2} = \frac{\partial J}{\partial h^3} \frac{\partial h^3}{\partial a^3} \frac{\partial a^3}{\partial h^2} \frac{\partial h^2}{\partial a^2} \frac{\partial a^2}{\partial \theta^2}$$

$$\theta^2 \rightarrow \theta^2 - \alpha \frac{\partial J}{\partial \theta^2}$$

$$\frac{\partial J}{\partial \theta^3} = \frac{\partial J}{\partial h^3} \frac{\partial h^3}{\partial a^3} \frac{\partial a^3}{\partial \theta^3}$$

$$\theta^3 \rightarrow \theta^3 - \alpha \frac{\partial J}{\partial \theta^3}$$

Neuronale Netze: Ergebnisse



Genauigkeit
Training: 100 %
Test: >98 %

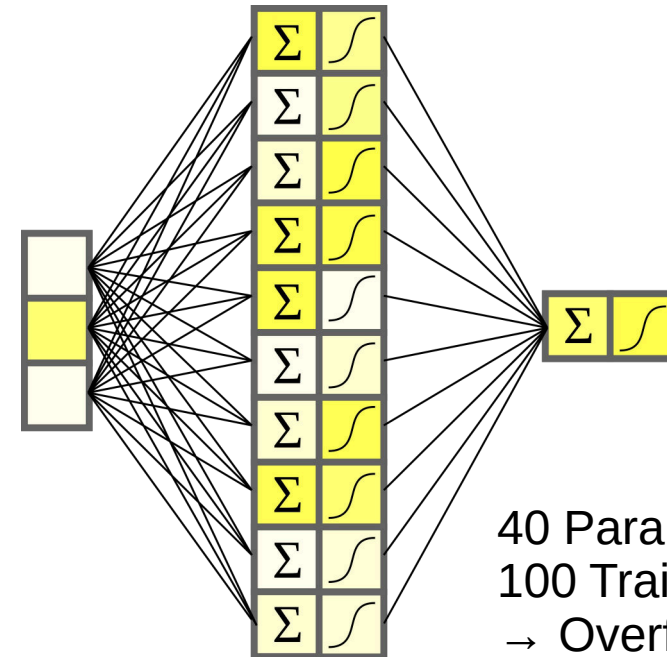
Neuronales Netz

$$a^1(x) = (\theta^1)^T x$$

$$h^1(x) = \sigma(a^1(x))$$

$$a^2(x) = (\theta^2)^T h^1(x)$$

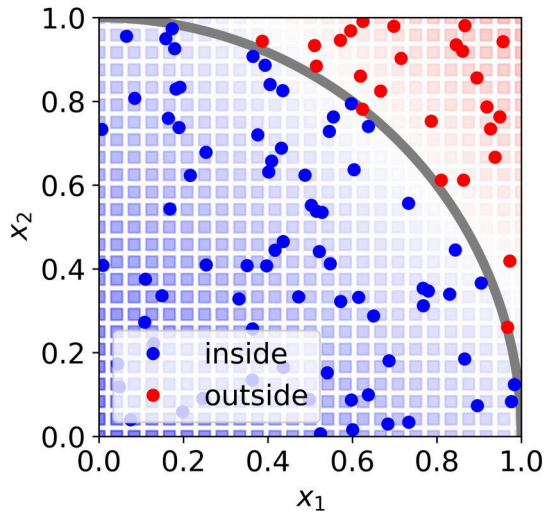
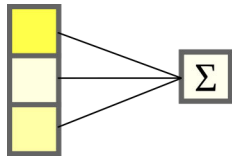
$$f(x) = h^2(x) = \sigma(a^2(x))$$



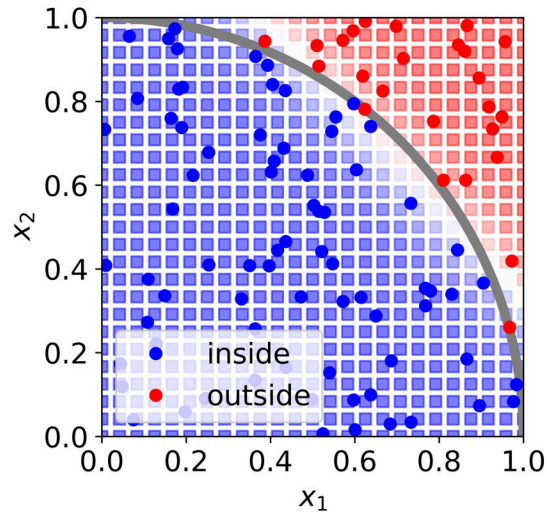
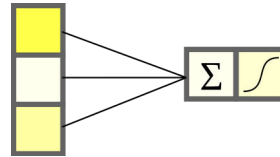
40 Parameter
100 Trainingspunkte
→ Overfitting?

Vergleich mit letzter Vorlesung: Lineare und logistische Regression

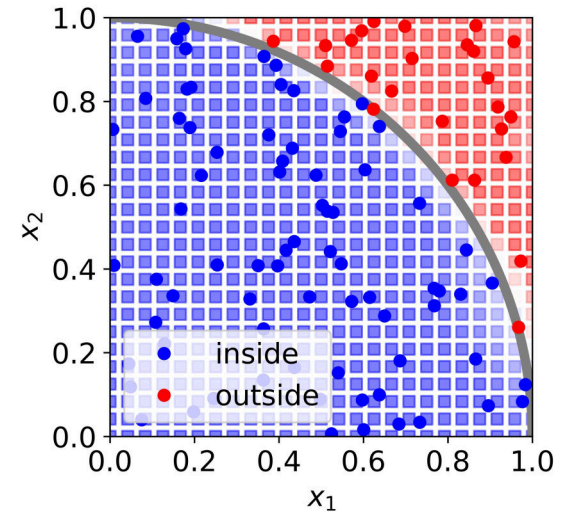
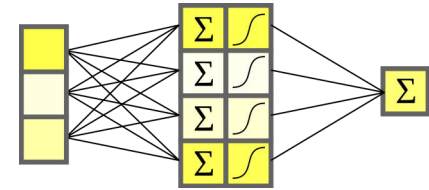
Lineare Regression



Logistische Regression

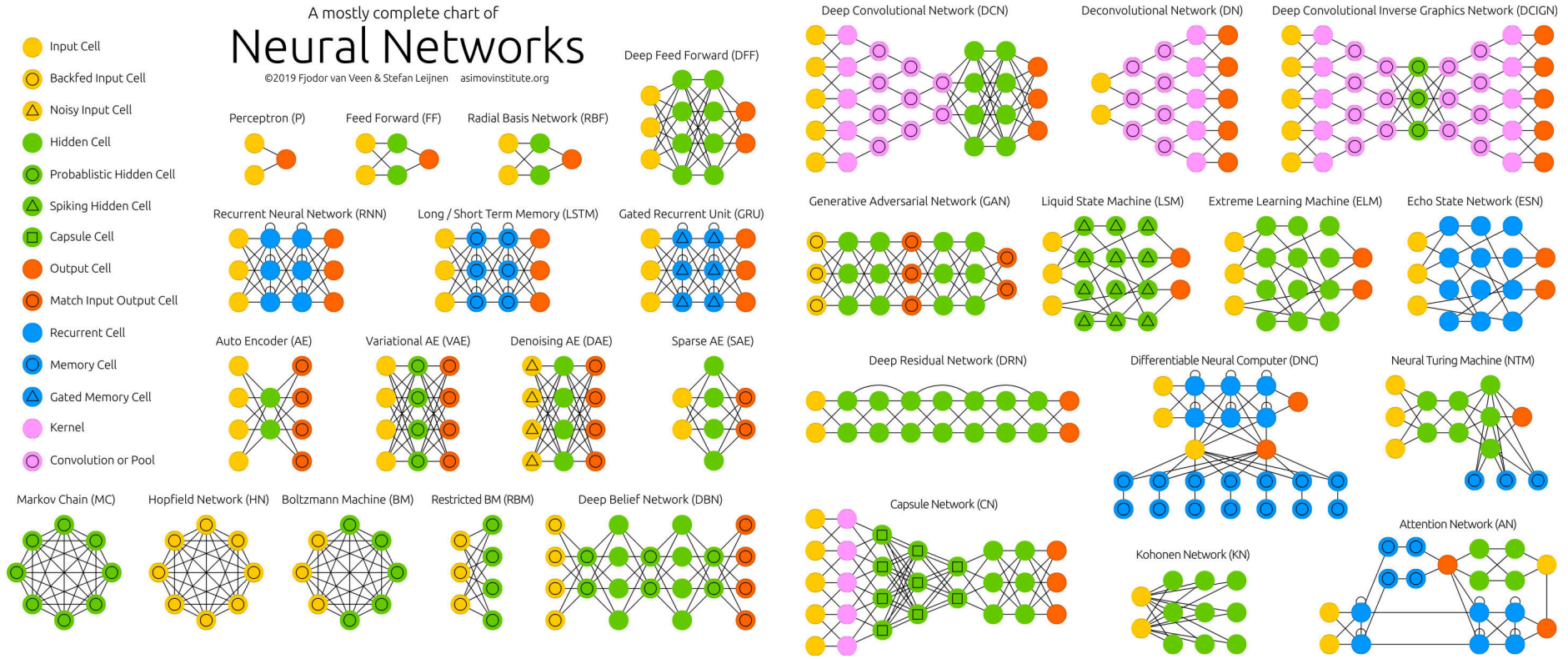


Neuronales Netz



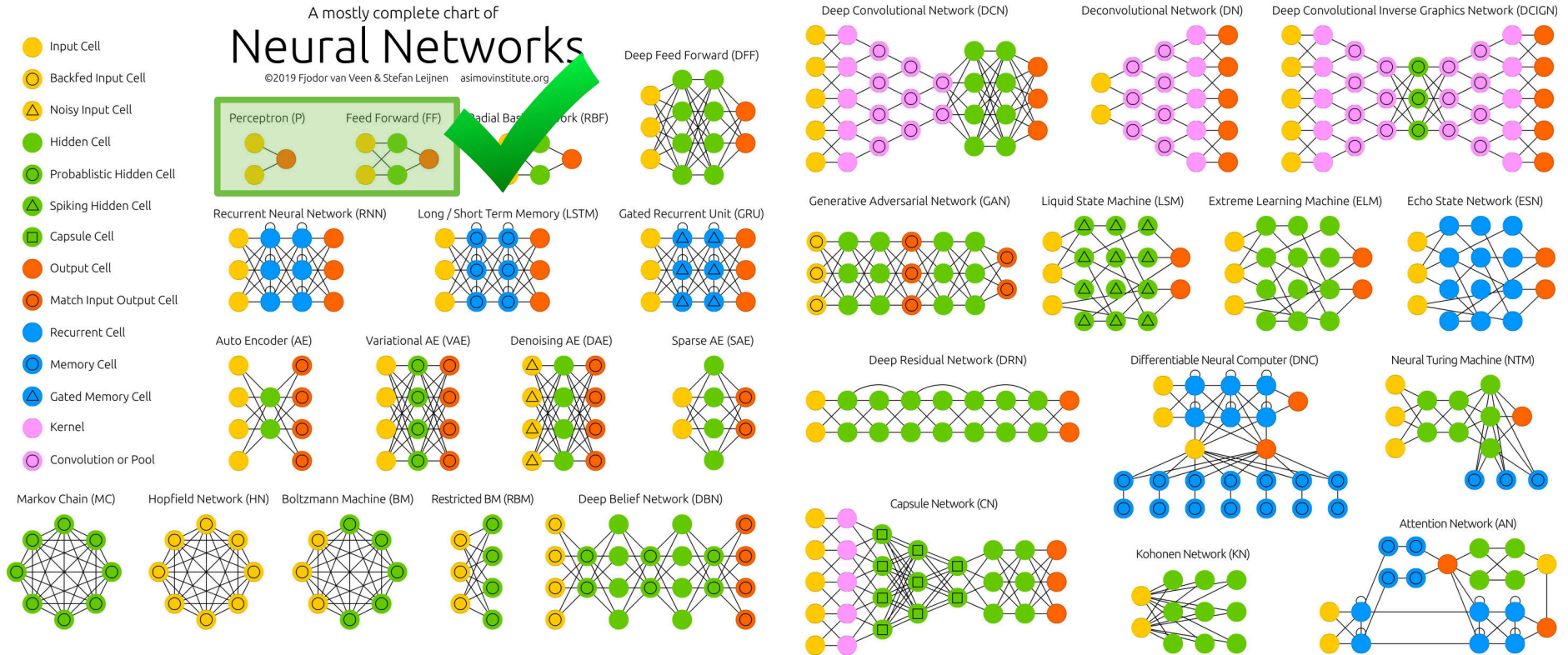
Und jetzt?

A “mostly” complete chart of neural networks (*Asimov institute*)

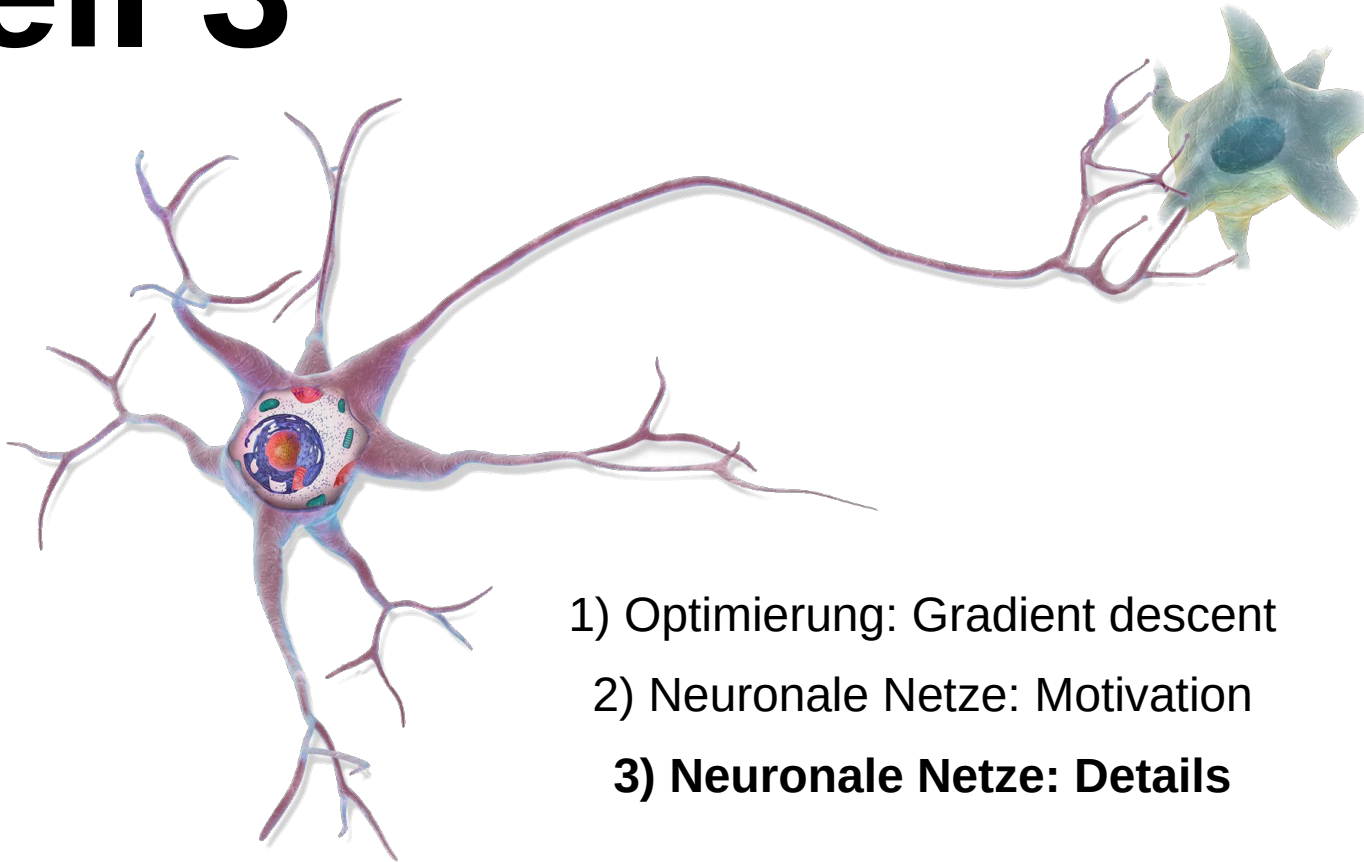


Und jetzt?

A “mostly” complete chart of neural networks (*Asimov institute*)



Teil 3



- 1) Optimierung: Gradient descent
- 2) Neuronale Netze: Motivation
- 3) Neuronale Netze: Details**

Warum benötigt man die Aktivierungsfunktion?

Neuronales Netz mit einem hidden layer

$$\begin{array}{ccc} h = f^{(1)}(x; \mathbf{W}, c) & & y = f^{(2)}(h; w, b) \\ & \searrow & \swarrow \\ & f(x; \mathbf{W}, c, w, b) = f^{(2)}(f^{(1)}(x)) & \end{array}$$

Beide Funktionen sind linear:
(ohne Bias)

$$f^{(1)}(x) = \mathbf{W}^\top x$$

$$f^{(2)}(h) = h^\top w$$



$$f(x) = x^\top \mathbf{W} w$$

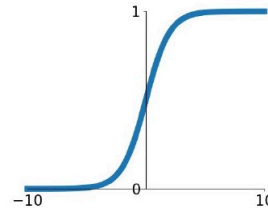
Das ist äquivalent zu einfachem linearen Modell: $f(x) = x^\top w'$ mit $w' = \mathbf{W} w$

Aktivierungsfunktionen

Viele verschiedene Aktivierungsfunktionen können verwendet werden

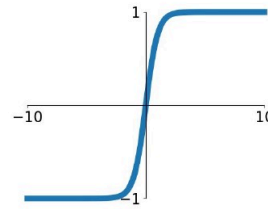
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



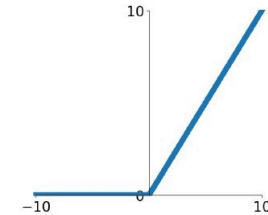
tanh

$$\tanh(x)$$



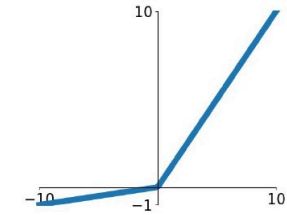
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

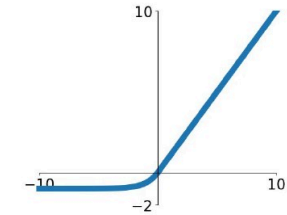


Maxout

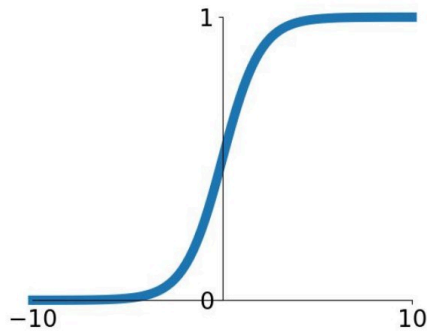
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



Aktivierungsfunktionen: Sigmoid



Sigmoid

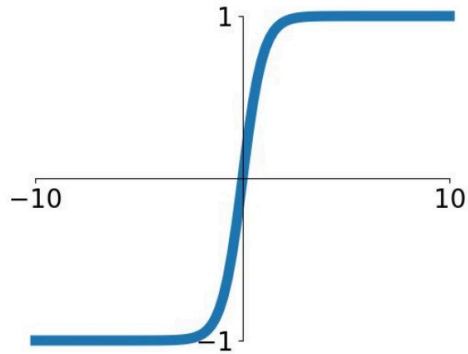
Berechnet: $f(a) = \sigma(a) = \frac{1}{1 + \exp(-a)}$

- Wertebereich [0,1]
- Historisch, Interpretation als saturierende Aktivierungsrate der Neuronen

Probleme

- Wenn Input groß/klein: Flach → Kein Gradient
- Output nicht um 0 zentriert: Macht Initialisierung schwer
- $\exp()$ ist aufwändig zu berechnen

Aktivierungsfunktionen: Tangens hyperbolicus (Tanh)



tanh(x)

Berechnet: $h(a) = \tanh(a)$

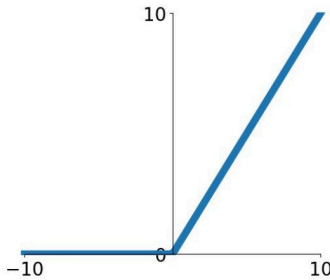
- Um 0 zentriert

Probleme

- Wenn Input groß/klein: Flach → Kein Gradient

Aktivierungsfunktionen: Tangens hyperbolicus (Tanh)

Berechnet: $h(a) = \max(0, a)$



ReLU
(Rectified Linear Unit)

- Sättigt nicht bei großem Input
- Sehr effizient zu berechnen
- Konvergiert typischerweise schneller als sigmoid/tanh

Probleme

- Wenn Input klein: Flach → Kein Gradient
- Output nicht um 0 zentriert: Macht Initialisierung schwer

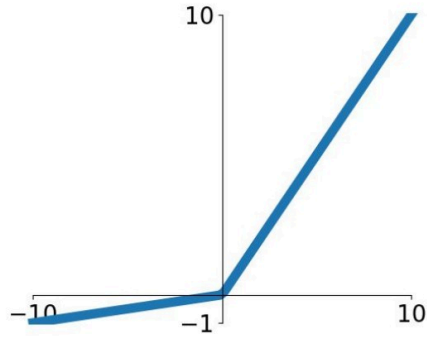
Nebenbemerkung: ReLU

- **ReLU** ist ein Akronym für „**R**ectified **L**inear **U**nit“
- Der Name geht aber zurück auf **Relu** Patrascu, Sysadmin an der Universität von Toronto



Image: Twitter, Relu Patrascu, @ikoflexer

Aktivierungsfunktionen: Leaky ReLU



Leaky ReLU

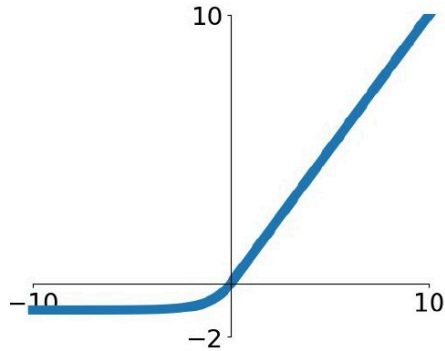
Berechnet: $h(a) = \max(0.1a, a)$

- Sättigt nicht → Hat immer einen Gradienten
- Sehr effizient zu berechnen
- Konvergiert typischerweise schneller als sigmoid/tanh

Erweiterung: Parametric Rectifier (PReLU)

- Berechnet: $h(a) = \max(\alpha a, a)$
- Parameter α muss/kann gelernt werden

Aktivierungsfunktionen: Exponential Linear Unit



Berechnet:
$$h(x) = \begin{cases} a & \text{if } a > 0 \\ \alpha(\exp(a) - 1) & \text{if } a \leq 0 \end{cases}$$

- Alle Vorteile von ReLU
- Outputs im Mittel näher bei 0
- Benötigt die Berechnung von $\exp()$

Exponential Linear Units (ELU)

Aktivierungsfunktionen: Auswahl

Praktisches Vorgehen

- Benutze **ReLU**, aber justiere die Lernrate und die Initialisierung der Parameter
- Teste **Leaky ReLU** und **ELU**
- Teste **tanh**, aber erwarte nicht zu viel
- Benutze keine **sigmoid** Funktion (außer als Output für Klassifizierungsprobleme)

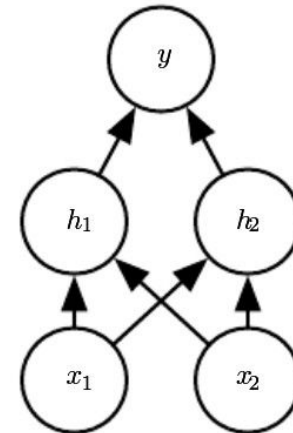
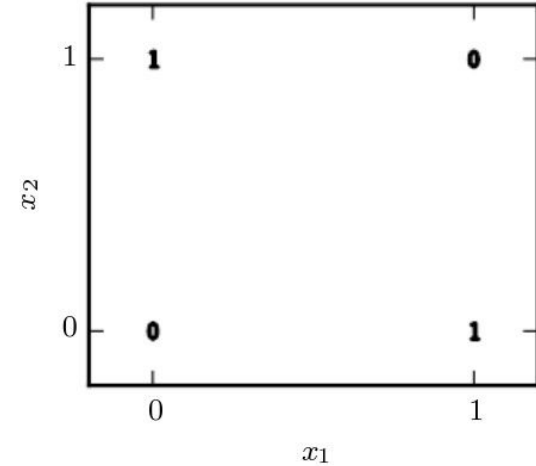
Fazit

- Aktivierungsfunktion als Hyperparameter: Nächste Woche mehr!

Beispiel: XOR

- Sehr einfaches Problem
- Vier Input Vektoren, vier Labels (0, 1)
- **Nicht lösbar** mit **linearem Modell**
- **Lösbar mit**
 - Neuronalem Netz mit einem hidden layer mit zwei Neuronen
 - Aktivierung: step (nur hier als Beispiel)

$$f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \text{step} \left(\mathbf{W}^\top \mathbf{x} + \mathbf{c} \right) + b$$



Beispiel: XOR

Eine mögliche Lösung: $f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \text{step}(\mathbf{W}^\top \mathbf{x} + \mathbf{c}) + b$

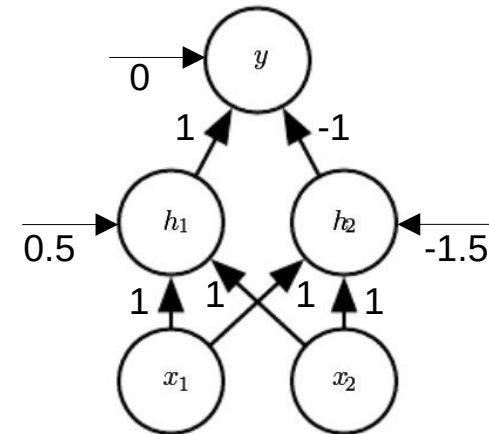
Schritt für Schritt:

$$\mathbf{X} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} \quad \mathbf{XW} = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix} \quad \mathbf{XW} + \mathbf{c} = \begin{bmatrix} -0.5 & -1.5 \\ 0.5 & -0.5 \\ 0.5 & -0.5 \\ 1.5 & 0.5 \end{bmatrix}$$

Für jede Zeile

$$\text{step}(\mathbf{XW} + \mathbf{c}) = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} \quad \mathbf{w}^\top \text{step}(\mathbf{XW} + \mathbf{c}) + b = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

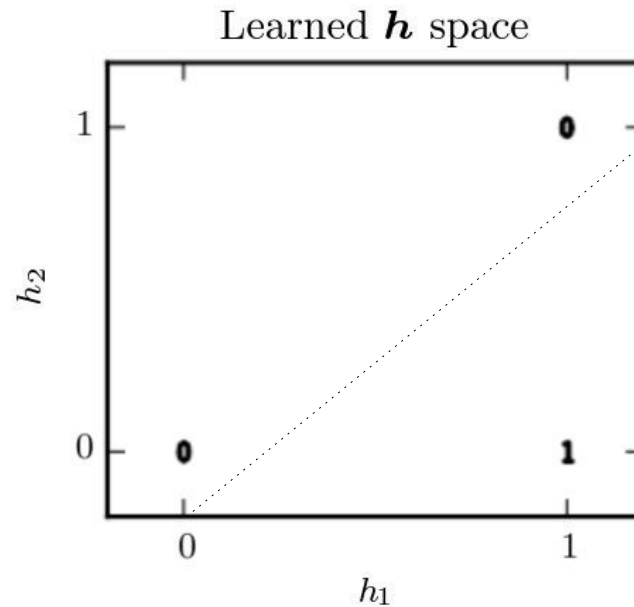
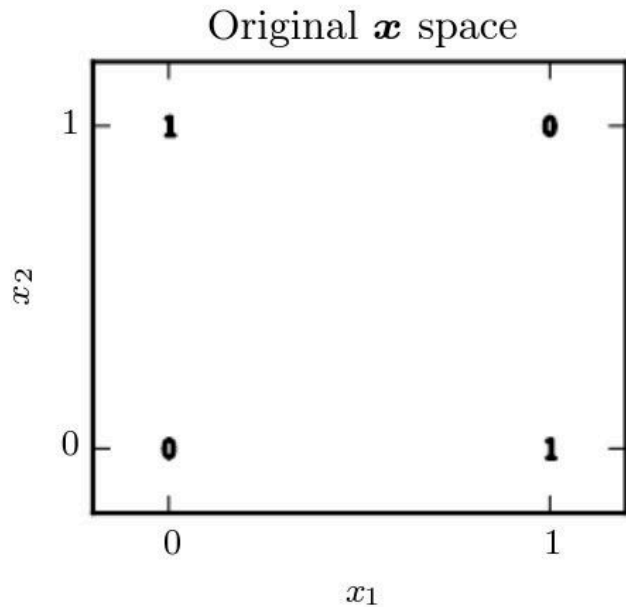
$$\mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$
$$\mathbf{c} = \begin{bmatrix} -0.5 \\ -1.5 \end{bmatrix}$$
$$\mathbf{w} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$
$$b = 0$$



Interpretation 1: Lernen einer linear separierbaren Repräsentation

$$f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \text{step}(\mathbf{W}^\top \mathbf{x} + \mathbf{c}) + b = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

Erklärung



$$\max\{0, \mathbf{XW} + \mathbf{c}\} = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}$$

Interpretation 1: Lernen der Logik

Aufgabe

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

Erklärung

$$h_1 = (x_1 \text{ OR } x_2)$$

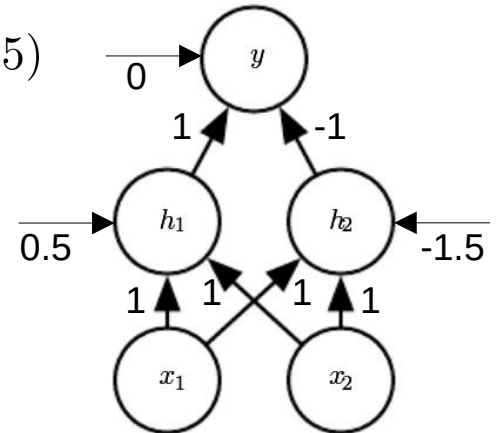
$$h_1 = \text{step}(1 \cdot x_1 + 1 \cdot x_2 - 0.5)$$

$$h_2 = (x_1 \text{ AND } x_2)$$

$$h_2 = \text{step}(1 \cdot x_1 + 1 \cdot x_2 - 1.5)$$

$$y = h_1 \text{ AND NOT } h_2$$

$$y = (h_1 - h_2)$$



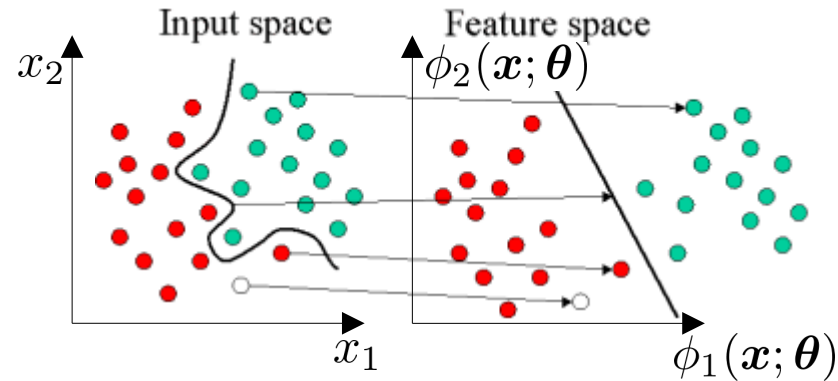
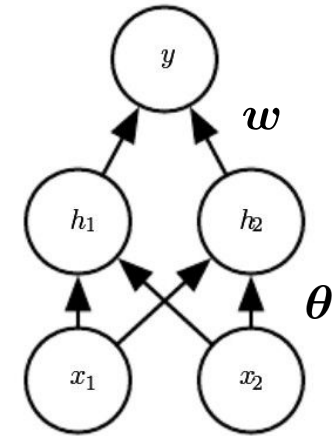
Allgemeiner: Representation Learning

Interpretation von neuronalen Netzen

- Kombination aus:
 - Nicht-lineare Transformation des Inputs (Gewichte θ)
 - Lineares Modell (Gewichte w)

$$y = f(\mathbf{x}; \boldsymbol{\theta}, \mathbf{w}) = \phi(\mathbf{x}; \boldsymbol{\theta})^\top \mathbf{w}$$

\swarrow
 h



Warum mehr als ein Layer?

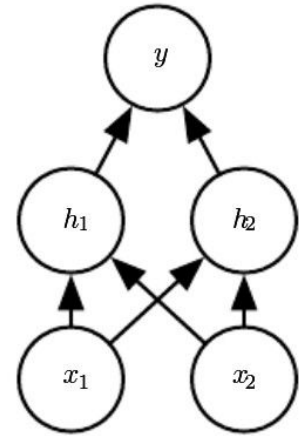
Aufbau des neuronalen Netzes

- **Modular:** Typischerweise in Form von aufeinanderfolgenden Layers
 - Freiheitsgrade: Anzahl, Art und Größe der Layers
- **Universal Approximation Theorem** (Hornik 1989, Cybenko 1989)
 - Neuronales Netz mit **einem hidden layer** kann alle Funktionen approximieren
 - (Natürlich mit Einschränkungen, die hier zu weit gehen)
 - Praktisch sind dafür aber sehr viele Neuronen notwendig
- **Tiefere neuronale Netze** können mit weniger Parametern bessere Modelle liefern
 - Besser: weniger Parameter, kleinere Validierungs- und Testfehler
 - Aber: eventuell schwieriger zu optimieren/trainieren

Welche Gewichte nutzen wir?

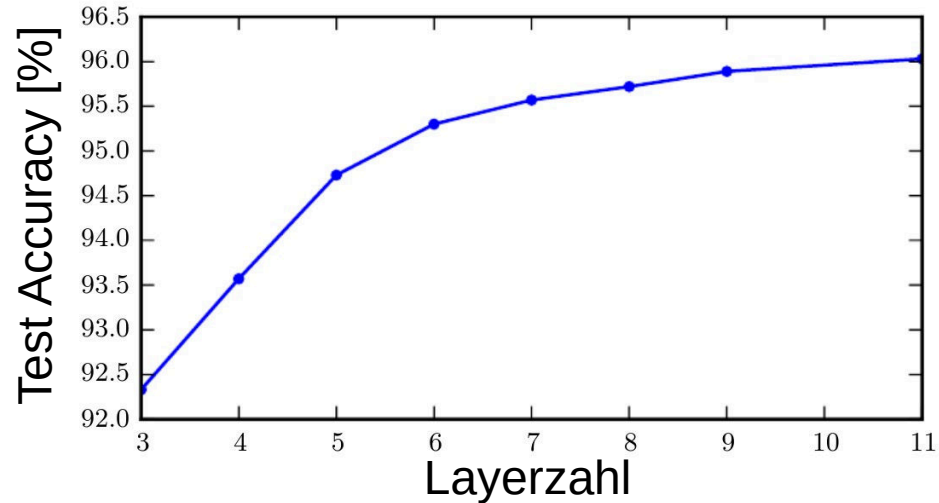
Initialisierung der Gewichte vor dem Training

- **Typisch:** Zufällige Gewichte
- Aber warum?
- Angenommen, wir setzen alle Gewichte w auf den gleichen Wert
- → Symmetrie in allen Gewichten und Informationsflüssen
- → wir werden immer wieder gleiche Gradienten erhalten
- → Initialisieren mit zufälligen (nicht-identischen) Gewichten



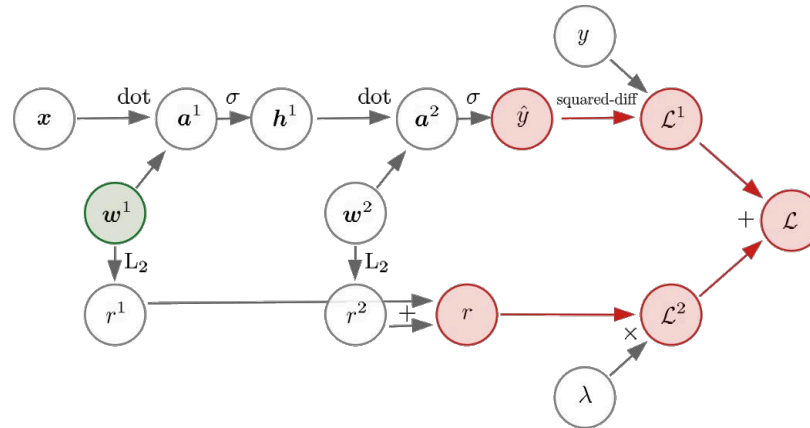
Tiefe neuronale Netze

Klassifizierung von fotografierten Hausnummern



<https://arxiv.org/pdf/1312.6082.pdf>

Teil 4



$$\nabla_{w^1} \mathcal{L} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}^1} \nabla_{w^1} \mathcal{L}^1 + \frac{\partial \mathcal{L}}{\partial \mathcal{L}^2} \nabla_{w^1} \mathcal{L}^2 = \dots$$

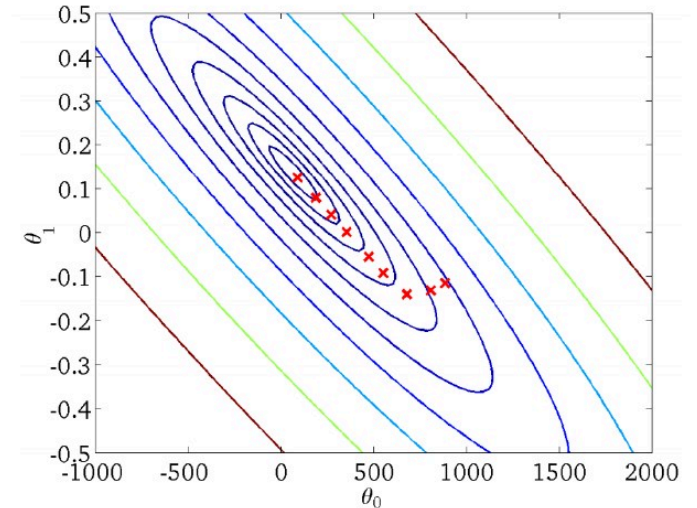
- 1) Optimierung: Gradient descent
- 2) Neuronale Netze: Motivation
- 3) Neuronale Netze: Details
- 4) Neuronale Netze: Back-propagation**

Gradient descent

- **Bisher:** Training von **logistischer Regression** mit **gradient descent**
- **Jetzt:** Neuronale Netze, fast identischer Algorithmus
- **Unterschiede:**
 - Sehr viel mehr Parameter
 - Nicht konvex: Viele lokale Optima
 - Gradient descent funktioniert hier häufig nicht „einfach so“
→ NN-Winter von ~2000 bis 2012
- **In den letzten 5-10 Jahren:**
 - Mehr Compute-Power
 - Mehr Daten
 - Einige Tricks

Gradient descent

- **Grundidee aus Teil 1**
 - Berechne Gradienten
 - Bewege die Parameter in die entgegengesetzte Richtung
- **Für neuronale Netze**
 - Grundsätzlich das selbe, aber:
 - ... höher-dimensional
 - ... schwieriger zu visualisieren



Loss Funktion

- **Ziel: Minimierung der Loss Funktion / Kostenfunktion**
 - Gegeben: Daten \mathcal{D} (M Datenpunkte \mathbf{x}_i)
 - Gesucht: Optimale Gewichte θ
 - → Minimiere Loss Funktion (inkl. Regularisierung)

$$\mathcal{L}(\theta, \mathcal{D}) = \sum_{i=1}^M l(\mathbf{x}_i, \theta) + \lambda \text{ penalty}(\theta)$$
$$\theta = \left\{ \mathbf{W}^{(1)}, \dots, \mathbf{W}^{(L)}, \mathbf{b}^{(1)}, \dots, \mathbf{b}^{(L)} \right\}$$

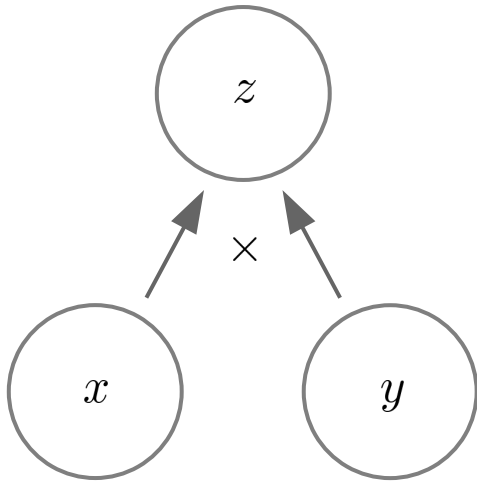
- **Vorgehen: Berechnung der partiellen Ableitungen**

- Gewichte \mathbf{W} : $\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(l)}}$ → Rekursive Anwendung der **Kettenregel**

- Bias \mathbf{b} : $\frac{\partial \mathcal{L}}{\partial \mathbf{b}^{(l)}}$ $\frac{d}{dw} L(h(\dots(w))) = \frac{dL}{dh} \frac{dh}{d \dots} \frac{d \dots}{d \dots} \frac{d \dots}{dw}$

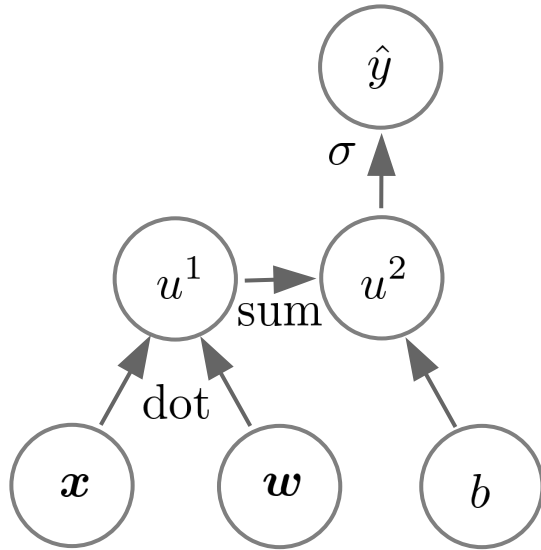
Wie kann man das automatisieren und implementieren?

→ Backpropagation durch Darstellung des computational graphs

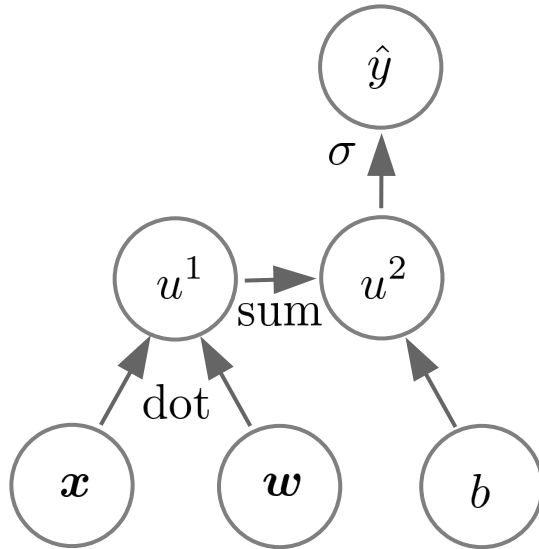


$$z = xy$$

Computational graph



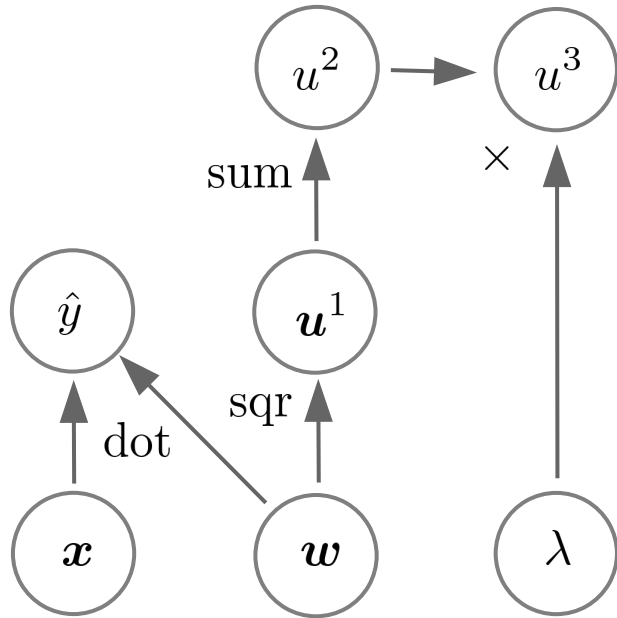
Computational graph



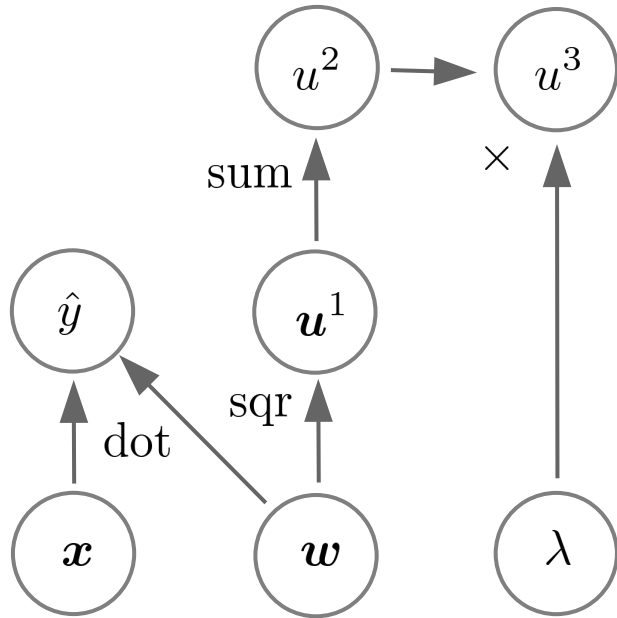
$$\hat{y} = \sigma(\mathbf{x}^\top \mathbf{w} + b)$$

Logistische Regression

Computational graph



Computational graph



Lineare Regression $x^\top w$
mit Regularisierung $\lambda \sum_i w_i^2$

Backpropagation aus computational graph

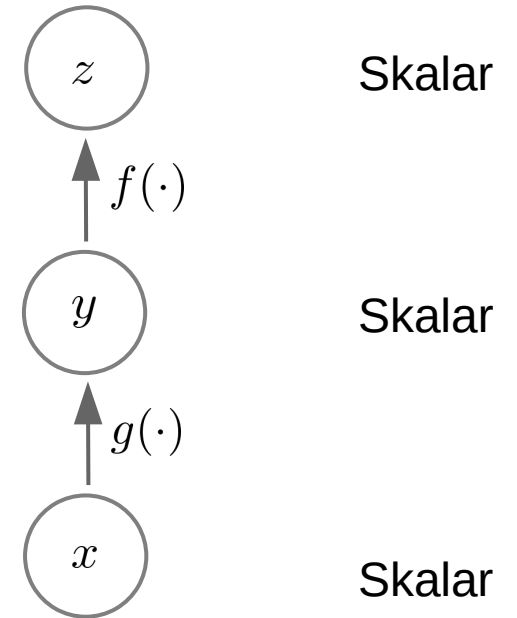
Einfaches Beispiel

$$x \in \mathbb{R}, y \in \mathbb{R}$$

$$y = g(x)$$

$$z = f(y) = f(g(x))$$

$$\longrightarrow \frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

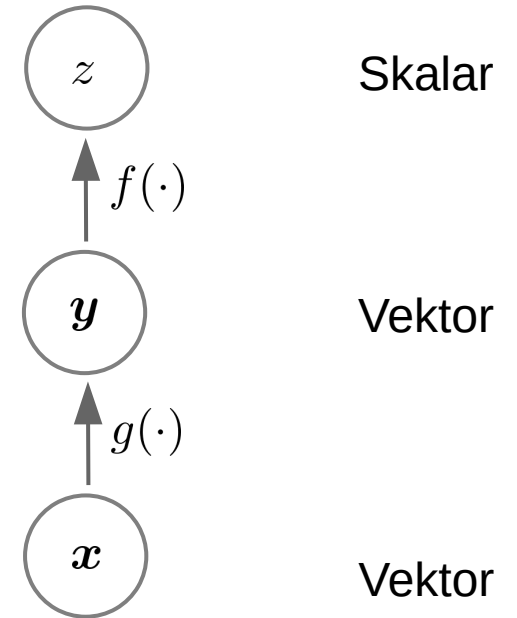


Backpropagation aus computational graph

Etwas komplizierter

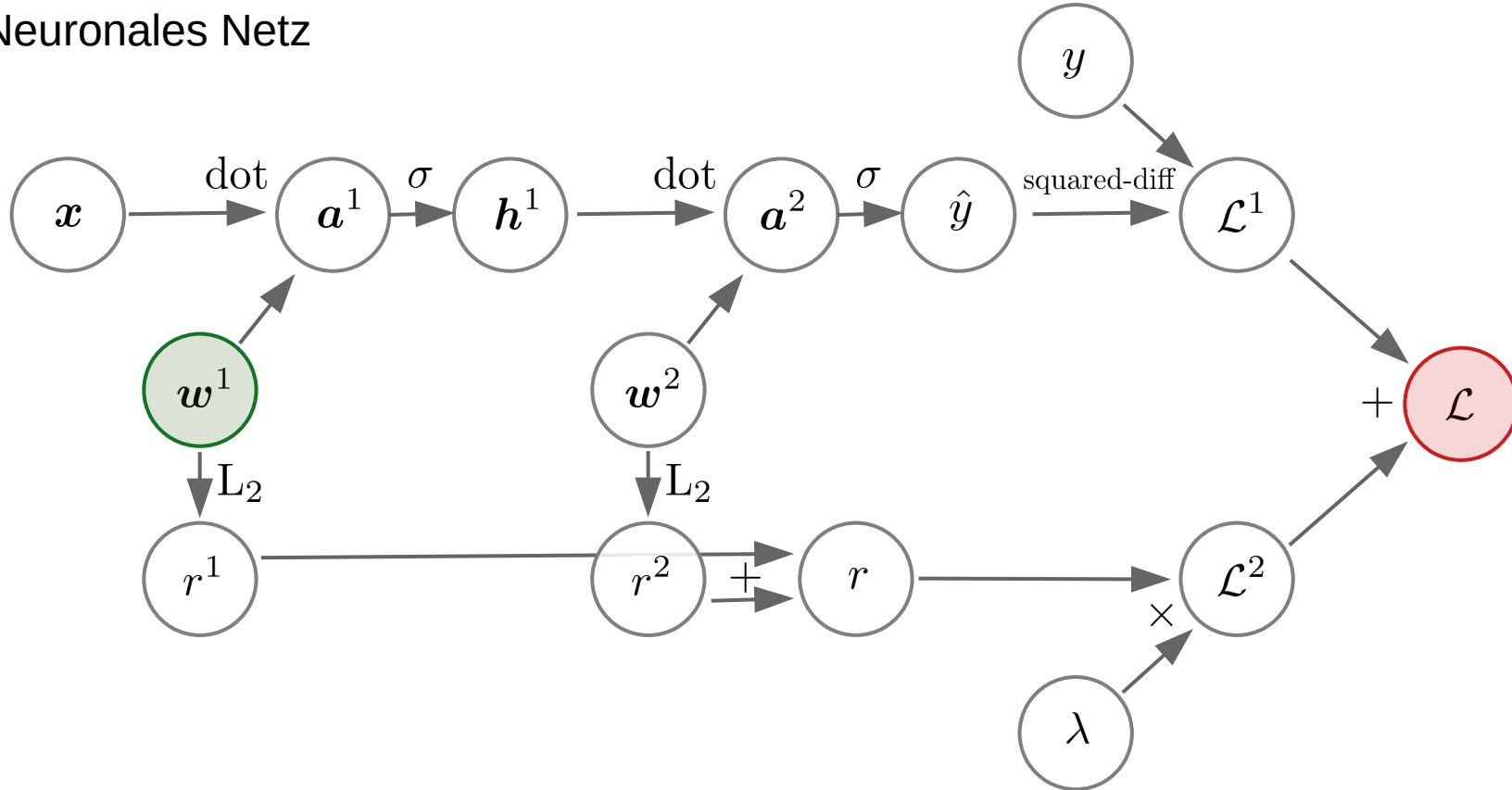
$$\begin{aligned} \mathbf{x} &\in \mathbb{R}^m, \mathbf{y} \in \mathbb{R}^n \\ \mathbf{y} &= g(\mathbf{x}) \\ z &= f(\mathbf{y}) \\ \longrightarrow \frac{\partial z}{\partial x_i} &= \sum_j \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i} \\ \longrightarrow \nabla_{\mathbf{x}} z &= \left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)^\top \nabla_{\mathbf{y}} z \end{aligned}$$

Jacobi Matrix Gradient



Backpropagation aus computational graph

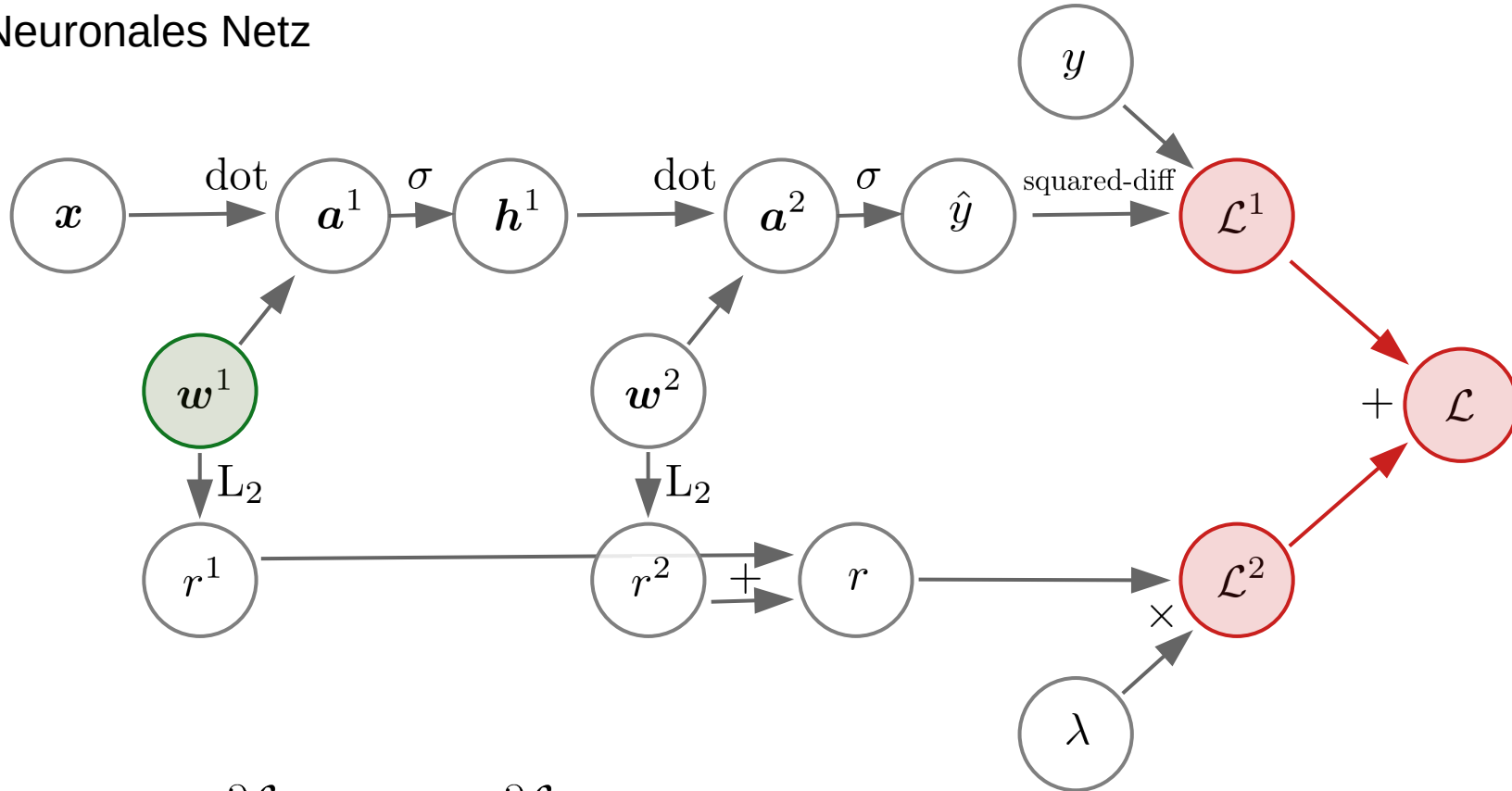
Neuronales Netz



$$\nabla_{w^1} \mathcal{L} = \dots$$

Backpropagation aus computational graph

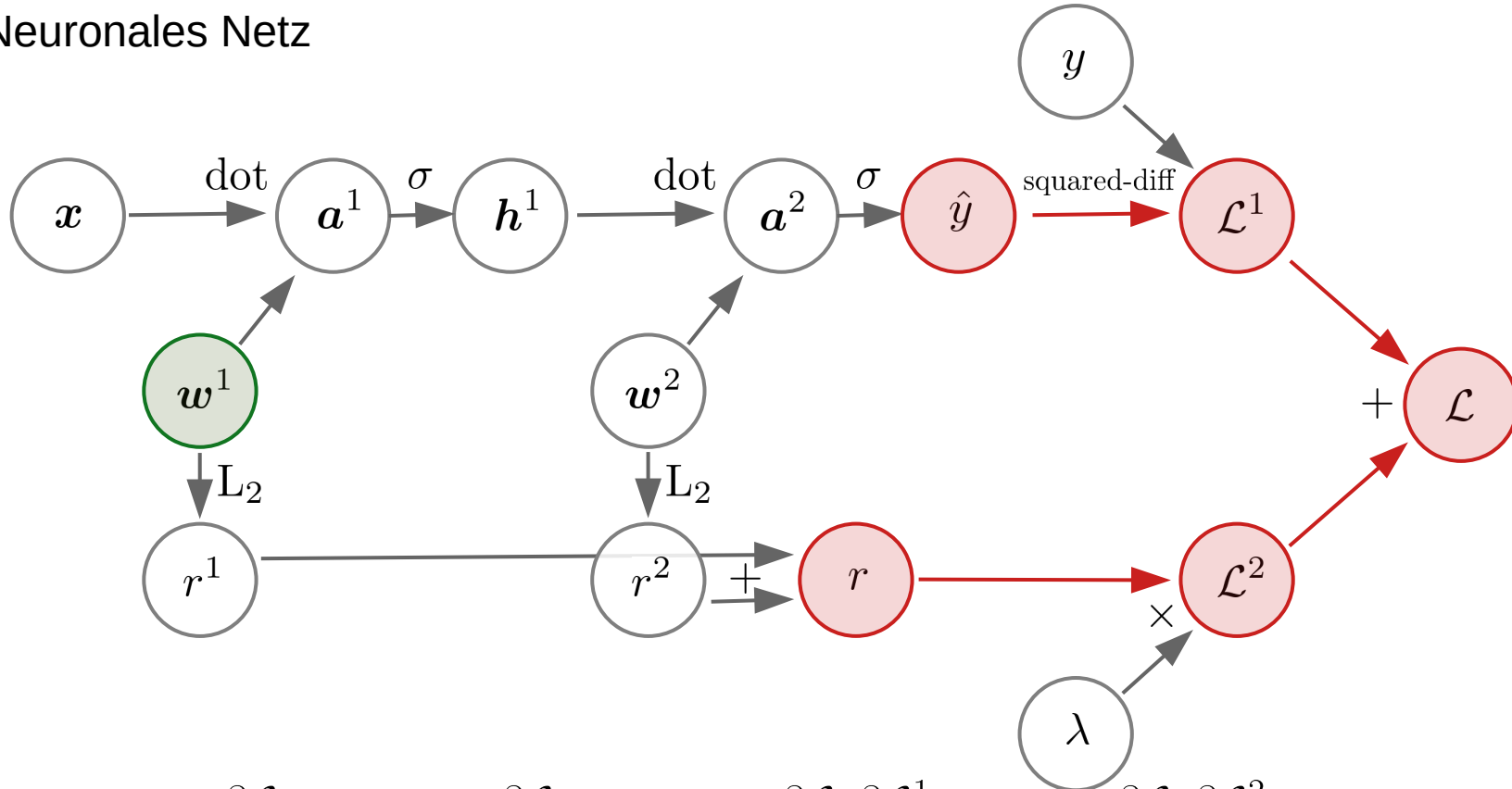
Neuronales Netz



$$\nabla_{w^1} \mathcal{L} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}^1} \nabla_{w^1} \mathcal{L}^1 + \frac{\partial \mathcal{L}}{\partial \mathcal{L}^2} \nabla_{w^1} \mathcal{L}^2 = \dots$$

Backpropagation aus computational graph

Neuronales Netz



$$\nabla_{w^1} \mathcal{L} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}^1} \nabla_{w^1} \mathcal{L}^1 + \frac{\partial \mathcal{L}}{\partial \mathcal{L}^2} \nabla_{w^1} \mathcal{L}^2 = \frac{\partial \mathcal{L}}{\partial \mathcal{L}^1} \frac{\partial \mathcal{L}^1}{\partial \hat{y}} \nabla_{w^1} \hat{y} + \frac{\partial \mathcal{L}}{\partial \mathcal{L}^2} \frac{\partial \mathcal{L}^2}{\partial r} \nabla_{w^1} r = \dots$$

Neuronale Netze: Computational cost

- **Forward propagation**

$$a = wx + b$$

- Matrix-Vektor Produkt
- Ca. eine Multiplikation und Addition pro Gewicht

- **Back propagation**

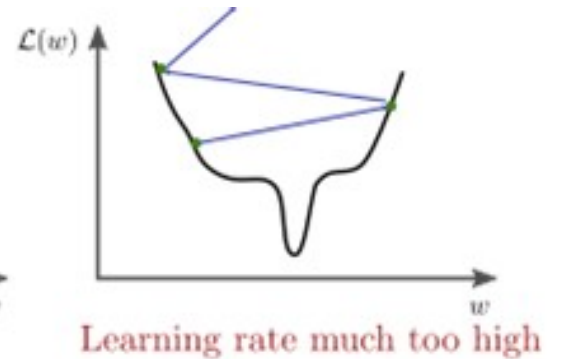
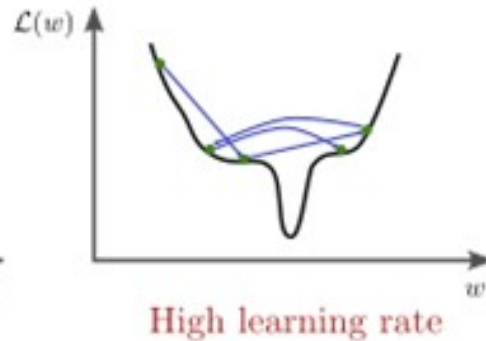
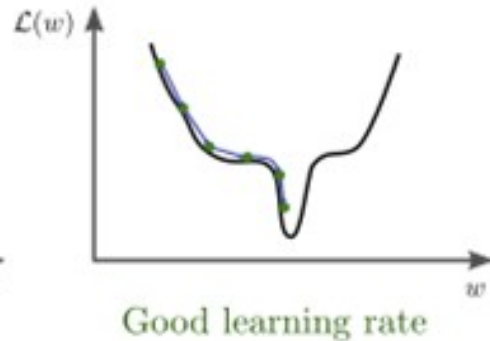
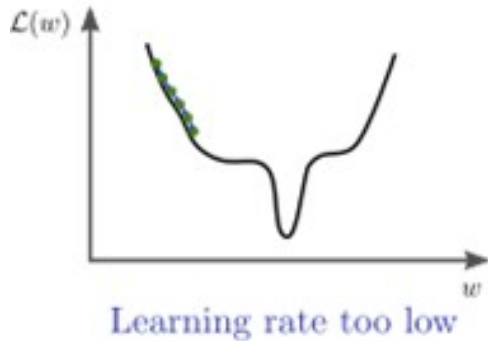
- Matrix-Vektor Produkt und äußeres Produkt
- Ca. zwei Multiplikationen und Additionen pro Gewicht

- **Daraus folgt**

- Kosten sind linear in der Zahl der Schichten
- Kosten sind quadratisch in der Zahl der Neuronen pro Schicht

Training: Gradient Descent plus Tricks

- **Gradient Descent**
 - Lernrate muss „gut“ gewählt werden



Training: Gradient Descent plus Tricks

- **Verbesserung: Momentum (Impuls/„Schwung“)**

- Gradient wie vorher
- Jetzt: Berechne den gleitenden Mittelwert des Gradienten

$$g_t = \frac{1}{b} \sum_{i \in I_t} \nabla_{\theta} \mathcal{L}(x_i; \theta_t)$$

Minibatch

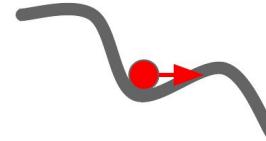
$$m_0 = \mathbf{0}, \quad m_{t+1} = \gamma m_t + (1 - \gamma) g_t$$

- Parameter Update: $\theta_{t+1} = \theta_t - \eta m_{t+1}$

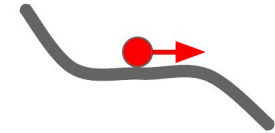
$\gamma = 0$: Gradient Descent

$\gamma > 0$: Momentum

Local Minima



Saddle points



- Interpretation: Gradient wirkt sich auf „Geschwindigkeit“ der Anpassung aus, und damit nur indirekt auf den Update-Schritt selbst

Training: Gradient descent plus Tricks

- **Weitere Tricks**

- **Gradient Normalization**

- Idee: Größere Schritte in flachen Gebieten, kleinere Schritte in steileren Gebieten

Gradient: $\mathbf{g}_k = \nabla_{\theta} \mathcal{L}(\theta_k)$ Running average gradient norm: $\mathbf{v}_{t+1,i} = \gamma \mathbf{v}_{t,i} + (1 - \gamma) \mathbf{g}_{t,i}^2$

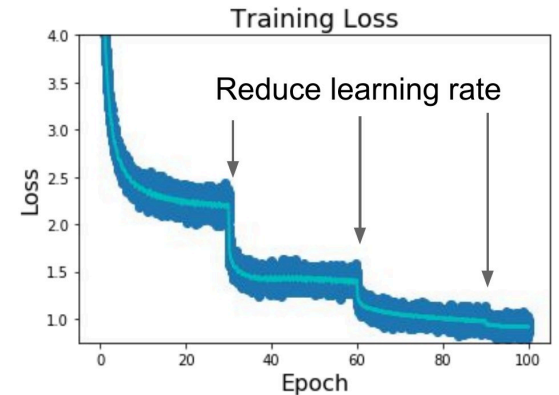
Update: $\theta_{t+1,i} = \theta_{t,i} - \frac{\alpha}{\sqrt{\mathbf{v}_{t+1,i} + \epsilon}} \mathbf{g}_{t,i}$

- **Learning rate decay**

- Idee: Passe die Lernrate nach k Epochen an

- **Adaptive Momentum (Adam)**

- Idee: Verbinde Momentum mit adaptiver Lernrate und Gradient Normalisierung



Sollten wir immer Backpropagation nutzen?

- Neuronale Netze: Grob inspiriert durch Biologie
- Aber: Im Gehirn gibt es (vermutlich) keinen Mechanismus, der Fehler ausrechnet und damit eine Backpropagation erlaubt
- Es gibt Alternativen für Backpropagation (oft durch Biologie motiviert)

PMLR Proceedings of Machine Learning Research
Volume 139 | JMLR | MLOSS | FAQ | Submission Form

arXiv > cs > arXiv:1901.09049

nature machine intelligence
Explore content | About the journal | Publish with us | Subscribe

nature > nature machine intelligence > news & views > article

News & Views | Published: 16 March 2020

ARTIFICIAL NEURAL NETWORKS

An alternative to backpropagation through time

Luca Manneschi & Eleni Vasilaki

Nature Machine Intelligence 2, 155–156 (2020) | Cite this article

1640 Accesses | 8 Citations | 3 Altmetric | Metrics

nature
Explore content | About the journal | Publish with us

nature > articles > article

Maass

Article | Open access | Published: 07 August 2024

Fully forward mode training for optical neural networks

Zhiwei Xue, Tiankuang Zhou, Zhihao Xu, Shaoliang Yu, Qionghai Dai & Lu Fang

Nature 632, 280–286 (2024) | Cite this article

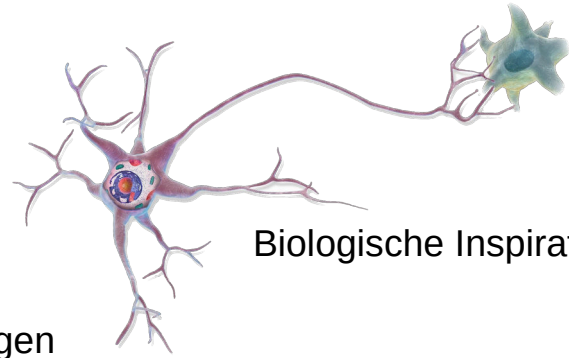
28k Accesses | 65 Altmetric | Metrics

Training Recurrent Neural Networks via Forward Propagation Through Time

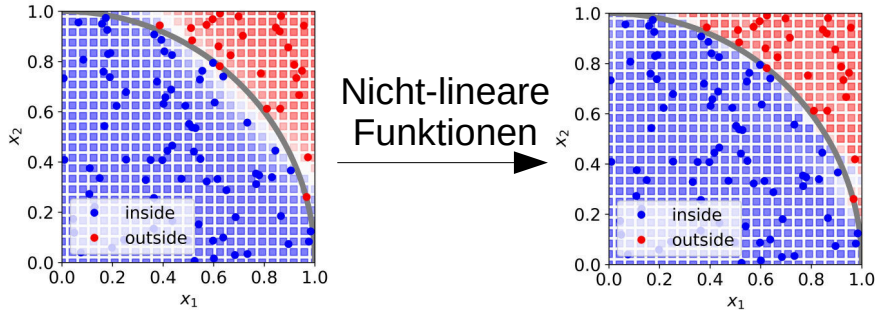
Anil Kag, Venkatesh Saligrama Proceedings of the 38th International Conference on Machine Learning, PMLR 139:5189–5200, 2021.

- Aber: Bisher ist Backpropagation **der** dominante Algorithmus fürs Training

Zusammenfassung

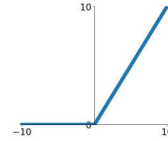


Biologische Inspiration



Nicht-lineare Funktionen

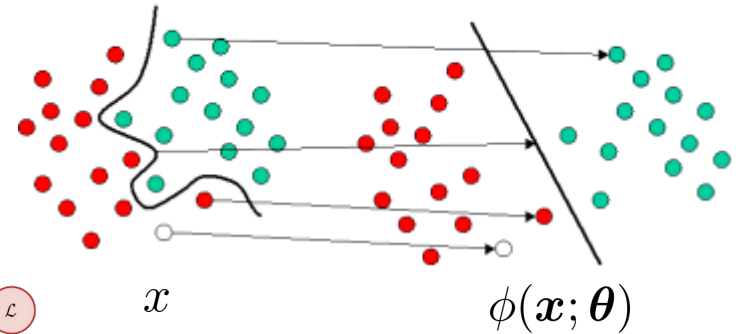
Aktivierungen



ReLU (Rectified Linear Unit)

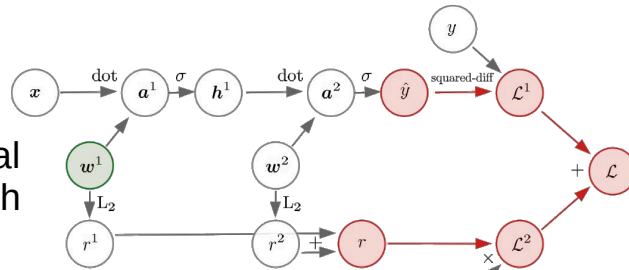
Representation learning

Input space Feature space



- 1) Optimierung: Gradient descent
- 2) Neuronale Netze: Motivation
- 3) Neuronale Netze: Details
- 4) Back-propagation: Details

Computational graph



Kettenregel

$$\nabla_{w^1} \mathcal{L} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}^1} \nabla_{w^1} \mathcal{L}^1 + \frac{\partial \mathcal{L}}{\partial \mathcal{L}^2} \nabla_{w^1} \mathcal{L}^2 = \dots$$

