

# Künstliche Intelligenz

Prof. Dr. Sibylle Schwarz  
HTWK Leipzig, Fakultät IM  
Gustav-Freytag-Str. 42a, 04277 Leipzig  
Zimmer Z 411 (Zuse-Bau)  
<https://informatik.htwk-leipzig.de/schwarz>  
[sibylle.schwarz@htwk-leipzig.de](mailto:sibylle.schwarz@htwk-leipzig.de)

Sommersemester 2021

# Was ist Künstliche Intelligenz?

EU-Factsheet on Artificial Intelligence

(<https://digital-strategy.ec.europa.eu/en/library/factsheet-artificial-intelligence-europe>)

*Artificial intelligence (AI) refers to systems that show intelligent behaviour: by analysing their environment they can perform various tasks with some degree of autonomy to achieve specific goals.*

*Mobile phones, e-commerce tools, navigation systems and many other different sensors constantly gather **data** or **images**. AI, particularly **machine-learning** technologies, can learn from this torrent of data to **make predictions** and **create useful insights**.*

Aussage über das **derzeitige** (eingeschränkte) Verständnis von KI

# Können Maschinen denken?

Alan Turing 1950

Konkretisierung der Frage:  
Können Maschinen **denken**?

zur überprüfbareren Frage:  
Können Maschinen konstruiert werden, die einen  
**speziellen Test bestehen**?

# Imitation Game

Imitation Game (Alan Turing 1950):

- ▶ zwei verschlossene Räume,  
in einem befindet sich **Herr** A, im anderen **Frau** B
- ▶ eine Person C (Frager) stellt Fragen, A und B antworten
- ▶ Kommunikation über neutrales Medium,  
an welchem das Geschlecht nicht erkennbar ist,
- ▶ C soll herausfinden, in welchem der Räume Frau B ist
- ▶ Herr A versucht, C irrezuführen
- ▶ Frau B kooperiert mit C

Herr A besteht den Test, wenn ihn C für Frau B hält.

# Wie erkennt man Intelligenz: Turing-Test

Turing-Test 1950: verschiedene Versionen des Imitation Game

- ▶ A ist Machine statt Mann (B Person beliebigen Geschlechts)
- ▶ verschiedene Kooperationsverhalten von A und B

Vorschlag zur Bewertung natürlichsprachlicher  
Kommunikationsfähigkeiten

# Beginn koordinierter Forschung zur Künstlichen Intelligenz

John McCarthy

Programmiersprachen

Marvin Minsky

Kognitionswissenschaft

Claude Shannon

Informationstheorie

stellten 1955 die Vermutung auf, dass

„jeder Aspekt des Lernens oder jedes anderen Ausdrucks von Intelligenz prinzipiell so präzise beschrieben werden kann, dass sich eine Maschine konstruieren lässt, die ihn simuliert.“

# Begriff Künstliche Intelligenz

McCarthy formulierte das Ziel,

„herauszufinden, wie man Maschinen konstruiert, die

- ▶ natürliche Sprache benutzen,
- ▶ Abstraktionen und Begriffe entwickeln,
- ▶ Aufgaben lösen, die (bis dahin) nur Menschen lösen konnten,
- ▶ sich selbst verbessern.“

und prägte dafür den Begriff **Künstliche Intelligenz**.

# Beginn koordinierter Forschung zur Künstlichen Intelligenz

1956: erste Konferenz zur Künstlichen Intelligenz

Dartmouth Summer Research Project on Artificial Intelligence

Themen:

- ▶ Berechnungsmodelle in Computern
- ▶ Kommunikation mit Computern in natürlicher Sprache
- ▶ Berechenbarkeitstheorie
- ▶ Neuronale Netzwerke
- ▶ Selbst-Verbesserung
- ▶ Abstraktionen
- ▶ Zufälligkeit und Kreativität

# Forschung zur Künstlichen Intelligenz

Momentaufnahme 2006:

Dartmouth Artificial Intelligence Conference: The Next Fifty Years

Themen:

- ▶ Modelle des (menschlichen) Denkens
- ▶ Neuronale Netzwerke
- ▶ (Maschinelles) Lernen und Suchen
- ▶ Maschinelles Sehen
- ▶ Logisches Schließen
- ▶ Sprache und Kognition
- ▶ KI und Spiele
- ▶ Interaktion mit intelligenten Maschinen
- ▶ Ethische Fragen und zukünftige Möglichkeiten der KI

# KI-Erfolge – Auswahl

- ▶ 1945 frühe Schachprogramme (ohne Implementierung)
- ▶ 1955 Logic Theorist: automatischer Beweiser
- ▶ 1958 erster erfolgreicher Neurocomputer Mark I Perceptron
- ▶ 1961 General Problem Solver, z.B. zum Lösen von Rätseln und Intelligenztests
- ▶ 1972 erster mobiler Roboter
- ▶ ab ca. 1970 Beschränkung auf spezialisierte Expertensysteme
- ▶ 1976 MYCIN (Medizinisches Diagnosesystem)
- ▶ 1980 Dendral (Molekülstruktur aus Massenspektrogramm)
- ▶ 1982 XCON (Konfiguration von Computersystemen)
- ▶ ab ca. 1980 Expertensystem-Shells
- ▶ seit 1993 RoboCup Roboter-Fußball
- ▶ 1997 Deep Blue gewinnt gegen amtierenden Weltmeister
- ▶ 2011 Watson schlägt zwei Meister in Quizshow Jeopardy!
- ▶ 2016 AlphaGo schlägt Go-Meister
- ▶ ...

# Phasen in der KI-Geschichte

wechselnde Betonung **symbolischer** und **statistischer** KI-Gebiete

- ▶ ca. 1950-70  
symbolisch: Inferenz-Maschinen, Eliza, Lisp  
statistisch: KNN  
Robotik, Verarbeitung natürlicher Sprache
- ▶ ca. 1970-85 (symbolisch)  
Prolog, Expertensysteme (z.B. medizinisch)  
nichtmonotones Schließen
- ▶ ca. 1985-2000 (statistisch)  
maschinelles Lernen, KNN, evolutionäre Alg.,  
Schwarm-Intelligenz, (Fuzzy-Logik) ...  
autonome Fahrzeuge (Ernst Dickmanns)
- ▶ ca. 2000-2010 (symbolisch)  
Constraint-Programmierung  
SAT-Solver, Theorem-Prover  
Ontologien (Semantic Web), Beschreibungslogiken
- ▶ seit ca. 2010 (statistisch)  
Deep Learning, CNN

# Ansätze intelligenter Systeme

- ▶ Simulation menschlichen **Verhaltens**  
(Verständnis und eigenes Denken nicht notwendig)  
Modellierung von Kognition,  
statistische Verfahren, Training mit vielen Fällen  
Getroffene Entscheidungen werden nicht begründet.  
**schwache** künstliche Intelligenz
  
- ▶ Simulation des menschlichen **Denkens**  
(Verständnis und eigenes Denken notwendig)  
Denkmodelle, mentale Modelle als Grundlage  
logisches Schließen, Abstraktion  
Jede Entscheidungen kann nachvollziehbar begründet werden.  
**starke** künstliche Intelligenz

# Kritik am Turing-Test

Kritik:

**schwache KI genügt**, um den Turing-Test zu bestehen

1966: Maschinelle Psychotherapeutin Eliza besteht Turing-Test

Searle (1980) Chinese-Room-Argument:

eine (nicht chinesisch verstehende) Person B in einem verschlossenen Raum mit einem (riesigen) Regelbuch mit chinesischen Fragen und passenden Antworten.

- ▶ A stellt Fragen, B antwortet.
- ▶ B antwortet mit Hilfe des Buches immer passend, ohne die Frage verstanden zu haben.

These: (anscheinend) intelligentes Verhalten ist noch

**keine Intelligenz, wenn Verständnis fehlt** (Ansatz der starken KI)

außerdem: praktisch nicht umsetzbar (Datenmenge)

# Logische / regelbasierte KI-Methoden

**Wissensrepräsentation:** formale Beschreibung von Umwelt (Randbedingungen) und Problem

**Problemlöseverfahren:** zur Lösung vieler Probleme anwendbares Standardverfahren (z.B. logisches Schließen)

Beispiele:

- ▶ Entscheidungsbäume und -tabellen
- ▶ Regelsysteme, Logiken, logisches Schließen
- ▶ Constraint-Systeme und -Löser
- ▶ deklarative Programmierung (logisch, funktional)
- ▶ fallbasiertes Schließen (durch Analogien)
- ▶ Simulation

typische Anwendungen klassischer KI-Methoden:

- ▶ Entscheidungsunterstützung (z.B. Finanzwirtschaft)
- ▶ Diagnosesysteme (z.B. in Medizin, Technik)
- ▶ Bewegungs- und Ablaufplanung

# Statistische KI-Methoden

„Soft-Computing“ oft besser geeignet für Probleme

- ▶ die unvollständig beschrieben sind,
- ▶ die keine eindeutige Lösung haben,
- ▶ für die keine effizienten Lösungsverfahren bekannt sind, usw.

einige Ansätze:

- ▶ künstliche neuronale Netze
- ▶ evolutionäre Algorithmen
- ▶ Schwarm-Intelligenz, Ameisen-Algorithmen
- ▶ Fuzzy-Logiken, probabilistische Logiken

## Aktuelle Entwicklung

starker Fortschritt einiger KI-Methoden („Deep Learning“)  
in den letzten 10 Jahren aufgrund der Entwicklung bei

- ▶ Computertechnik: Parallelrechner, GPU (70% Einfluss)
- ▶ Speichermöglichkeit großer Datenmengen, Verfügbarkeit großer strukturierter und annotierter Datenmengen (20%)
- ▶ neue Typen künstlicher neuronaler Netze, neue Algorithmen (10%)

sowie starkem Medieninteresse an bestimmten Erfolgen, z.B.

- ▶ 1997 Deep Blue gewinnt gegen amtierenden Weltmeister
- ▶ 2011 Watson schlägt zwei Meister in Quizshow Jeopardy!
- ▶ 2012 erste Zulassung eines autonomen Fahrzeugs für den Test auf öffentlichen Straßen
- ▶ 2016 AlphaGo schlägt Go-Meister
- ▶ ...

führte zum aktuellen Aufblühen der KI-Euphorie

# Leistung aktueller (statistischer) KI-Systeme

nahe und teilweise über den menschlichen Fähigkeiten z.B. bei

- ▶ Erkennung von Objekten in Bildern
- ▶ Einordnung / Klassifikation von Objekten und Situationen
- ▶ Reaktion auf klar erkannte Situationen
- ▶ strategischen Spielen mit endlichem Zustandsraum  
z.B. Schach, Go

prinzipielle Herausforderungen:

- ▶ Zuverlässigkeit, Sicherheit
- ▶ Begründung, Erklärung

# Schwächen aktueller (statistischer) KI-Systeme

KI derzeit noch weit von menschlichen Fähigkeiten entfernt bzgl.

- ▶ Erkennung der eigenen Grenzen
- ▶ Intuition
- ▶ Aufstellen und Überprüfen sinnvoller Annahmen bei unvollständig vorhandener Information
- ▶ Lernen ohne vorheriges Training mit großen Mengen (manuell) annotierter Daten
- ▶ Übertragen von Wissen zwischen verschiedenen Anwendungsbereichen
- ▶ Kombination verschiedener Methoden
- ▶ Schließen bzgl. rechtlicher und moralischer Bezugssysteme, mentaler Modelle

# Einordnung in die Informatik

**Informatik** Wissenschaft von der Darstellung und Verarbeitung symbolischer Information durch Algorithmen

Einordnung in die Teilgebiete der Informatik:

**theoretisch** ▶ Sprachen zur Formulierung von Information und Algorithmen,  
▶ Berechenbarkeit durch Algorithmen,

Grundlagen, z.B. Logik, formale Sprachen

**technisch** ▶ maschinelle Darstellung von Information  
▶ Mittel zur Ausführung von Algorithmen

Parallelrechner, GPU, Anwendung z.B. in HW-Verifikation, technischer Diagnose

**praktisch** Entwurf und Implementierung von Algorithmen  
Grundlagen, z.B. Graph-Suchverfahren, Inferenzalgorithmen, Algorithmen zum Constraint-Lösen, Anwendung z.B. in SW-Verifikation

**angewandt** Anwendung von Algorithmen, z.B. Anwendung, z.B. KI, Spracherkennung, Bilderkennung, Suchmaschinen, autonome Agenten, Robotik

# Inhalt der Lehrveranstaltung

- ▶ Heuristische Suche / Spielbaumsuche
- ▶ Künstliche Neuronale Netze
- ▶ Unscharfes / probabilistisches Schließen
- ▶ Bayes-Netze
- ▶ Kausalität (Zusammenhang von Ursache und Wirkung)

# Literatur

Folien, Aufgaben, ... zur aktuellen Vorlesung unter

<https://informatik.htwk-leipzig.de/schwarz/lehre/ss21/kim>

Bücher:

- ▶ KI-Grundlagen:
  - ▶ Ingo Boersch, Jochen Heinsohn, Rolf Socher:  
Wissensverarbeitung (Spektrum, 2007)
  - ▶ Wolfgang Ertel:  
Grundkurs Künstliche Intelligenz (Springer, 2016)  
(elektronische Version in HTWK-Bibliothek)
  - ▶ Ronald Brachman, Hector Levesque:  
Knowledge Representation and Reasoning  
(Morgan Kaufmann 2004)
  - ▶ Stuart Russell, Peter Norvig:  
Künstliche Intelligenz (Pearson 2004)
- ▶ KNN:
  - ▶ Raúl Rojas: Neural Networks – A Systematic Introduction  
<https://page.mi.fu-berlin.de/rojas/neural/neuron.pdf>
- ▶ Kausalität:
  - ▶ Judea Pearl: The Book of Why  
<http://bayes.cs.ucla.edu/WHY/>

# Organisation

5 ECTS (Präsenzzeit 56 h, Vor- und Nachbereitungszeit 94 h)

- ▶ wöchentlich
  - ▶ eine Vorlesung
  - ▶ evtl. Übungsaufgaben
  - ▶ ein Seminar (gemeinsam, direkt vor der Vorlesung)
- ▶ Prüfungsvorleistung: Beleg
  - ▶ Referate zu passenden Fachartikeln
  - ▶ evtl. Präsentation der Lösung der Übungsaufgaben
- ▶ Prüfung: Klausur 90 min oder mündliche Prüfung (BBB)

# Was bisher geschah

- ▶ KI-Geschichte
- ▶ KI-Tests (Turing, Chinese Room)
- ▶ statistische / symbolische KI

# Daten, Wissen, Intelligenz

Umwelt		Reize, Eindrücke
Agent	Wahrnehmen, Beobachten	Daten
	Erkennen, Verstehen	Information
	Anwenden, Können	Wissen
	Lernen	Wissenserwerb (Intelligenz?)
	Reflektieren, Begründen, Erkennen der Grenzen, Verstehen	Intelligenz

# Beispiel: Daten, Information, Wissen, Intelligenz

**Daten** Darstellungsform (Syntax)  
Zeichenketten, Bilder, Ton, ... (z.B 39.7)

**Information** Bedeutung der Daten (Semantik)  
in einem bestimmten Kontext  
im Beispiel: Körpertemperatur= 39.7

**Wissen** Information mit einem Nutzen,  
trägt zur Lösung eines Problemes bei,  
Nutzen abhängig von vorhandenem Kontextwissen  
im Beispiel: Kontext Körpertemperatur > 39.0 ist Fieber,  
Fieber ist Symptom von COVID-19 (27%) oder  
Nebenwirkung einer Corona-Impfung oder ...  
Bei Verdacht auf COVID-19 testen,  
bei kürzlicher Impfung beobachten, sonst ...

**Wissenserwerb** selbständige Informationsgewinnung (auch zum Kontext)  
im Beispiel über (derzeit typische) Auslöser,  
Nebensymptome, Therapien

**Intelligenz** Diagnose und Auswahl aus Therapie-Alternativen speziell  
für die zu behandelnde Person durch Abwägung der zu  
erwartenden Wirkungen, ggf. Überweisung zu Spezialisten

# Explizites und implizites Wissen

## explizites Wissen

z.B. Fakten, Aussagen, Zusammenhänge, Verfahren  
ermöglicht maschinelle Verarbeitung

## implizites Wissen

z.B. Fähigkeiten wie Laufen, Autofahren,  
Schachspielen  
wird durch Training erworben,  
(ggf. mit Hilfe expliziten Wissens, z.B. Spielregeln)  
Nachbildung durch statistische Verfahren

Kommuniziert werden kann nur explizites Wissen.

Transformation von implizitem in explizites Wissen notwendig

# Problemlösung durch Suche in Graphen – Beispiele

- ▶ Finden von Wegen in einem Graphen
  - ▶ Aufgabe:
    - ▶ gegeben: Graph  $G$  (Tafel)
    - ▶ gesucht: Weg (Pfad) in  $G$  von Knoten  $u$  zu Knoten  $v$
  - ▶ Lösungsidee: Suche im Graphen
- ▶ Münzenstapelspiel (für eine Person)
  - ▶ Aufgabe:
    - ▶ gegeben: Stapel von  $n$  Münzen
    - ▶ gesucht: Zugfolge durch erlaubte Züge (zwei Münzen von einem Stapel nehmen und auf beide Nachbarn verteilen) bis zu einer Situation, in der kein Zug möglich ist
  - ▶ Lösungsidee:
    - ▶ Modellierung als Zustandsübergangssystem
    - ▶ Suche im Graphen
- ▶ 3 Krüge
  - ▶ Aufgabe:
    - ▶ gegeben: 3 volle Krüge mit Volumen 4l, 7l, 9l,
    - ▶ gesucht: genau 6l in einem der 3 Krüge
  - ▶ Lösungsidee: Zustände als Knoten eines Suchbaumes

# Darstellung von Aufgabe und Lösung

Aufgabe:

- gegeben:
- ▶ Menge  $V$  von Zuständen (evtl. unendlich)  
oft beschrieben durch Eigenschaften
  - ▶ Startzustand  $s \in V$
  - ▶ Menge  $Z \subseteq V$  von Zielzuständen  
(oder Eigenschaften der Zielzustände)
  - ▶ mögliche Übergänge zwischen Zuständen  
Übergangsrelation  $E \subseteq V \times V$

**Lösung:** Folge von Zuständen (Weg von einem Start- zu einem Zielzustand) (Mitunter interessiert nur der erreichte Zielzustand.)

**Wissensrepräsentation:** als Graph  $G = (V, E)$

(Zustandsübergangssystem):

- ▶ Knotenmenge  $V$ : Zustände
- ▶ (gerichtete) Kanten: Zustandsübergänge

Entfaltung des Graphen zu einem Baum:

Pfade im Graphen = Knoten im Baum

# Problemlösen durch Suchen

- ▶ formale Darstellung des Problemes als Graph (z.B. Baum, DAG)
- ▶ formale Beschreibung der Lösung als Eigenschaft von
  - ▶ Pfaden im Graphen
  - ▶ Knoten im Baum

Möglichkeiten zum Problemlösen:

- ▶ Pfadsuche im Graphen
- ▶ Knotensuche im Baum

# Suche in Graphen

(schon bekannte) Verfahren zur Suche in Graphen (und Bäumen):

- ▶ Tiefensuche (depth-first search):  
Suche zuerst in Teilbäumen eines noch nicht besuchten Nachbarn des aktuellen Knotens
- ▶ Breitensuche (breadth-first search):  
Suche zuerst in Teilbäumen eines noch nicht besuchten Knotens mit der geringsten Tiefe

# Allgemeines Suchverfahren

- Daten:  $L_a$  Menge der noch zu expandierenden Knoten  
 $L_x$  Menge der expandierten Knoten  
 $s$  Startknoten  
 $\varphi$  Anforderungen an Lösung (Zielknoten)

Allgemeiner Suchalgorithmus:

1.  $L_a = \{s\}, L_x = \emptyset$
2. solange  $\neg L_a = \emptyset$ :
  - 2.1 Verschiebe einen auf **festgelegte Art** ausgewählten Knoten  $u$  aus  $L_a$  in  $L_x$
  - 2.2 Füge alle Nachbarn von  $u$ , die nicht in  $L_a \cup L_x$  enthalten sind, auf eine **festgelegte Art** in  $L_a$  ein  
(Abbruch falls ein Nachbar  $v$  von  $u$  die Bedingung  $\varphi$  erfüllt, also eine Lösung repräsentiert)

prominente Spezialfälle:

- Tiefensuche**   ▶ Verwaltung von  $L_a$  als **Stack**  
▶ Einfügen der Nachbarn an den **Anfang** der Liste  $L_a$   
▶ festgelegter Knoten wurde **zuletzt** in  $L_a$  eingefügt
- Breitensuche**   ▶ Verwaltung von  $L_a$  als **Queue**  
▶ Einfügen der Nachbarn an das **Ende** der Liste  $L_a$   
▶ festgelegter Knoten wurde **zuerst** in  $L_a$  eingefügt

## Schrittweise Vertiefung (iterative deepening)

beschränkte Tiefensuche:

1. festgelegte Tiefenbeschränkung  $m \in \mathbb{N}$
2. Tiefensuche auf allen Pfaden bis zur Tiefe  $m$

nicht vollständig, weiter entfernte Lösungen werden nicht gefunden

Schrittweise Vertiefung (iterative deepening):

Kombination aus Breiten- und Tiefensuche durch

Nacheinanderausführung der beschränkten Tiefensuche für alle  $m \in \mathbb{N}$ , solange keine Lösung gefunden wurde

vollständig, optimal

(asymptotischer) Zeit- und Platzbedarf wie Tiefensuche

# Gleiche-Kosten-Suche (kleinste bisherige Kosten)

(uniform-cost-search)

bei Zustandsübergängen mit verschiedenen **Kosten**

Ziel: Lösung (Pfad vom Start- zu einem Lösungsknoten) mit möglichst geringen Pfadkosten

(Pfadkosten = Summe der Kosten aller Übergänge auf dem Pfad)

**Bewertungsfunktion** für Knoten  $k : V \rightarrow \mathbb{R}_{\geq 0}$

$k(u)$  = minimale (bisher entdeckte) Pfadkosten vom Startknoten zu  $u$

Datenstruktur zur Verwaltung von  $L_a$ : Priority Queue

Priorität eines Knotens  $u$ :  $k(u)$

Beispiele:

- ▶ Breitensuche (Kosten = Tiefe des aktuellen Knotens  $u$ )
- ▶ kürzeste Wege (Kosten = minimale bisher bekannte Kosten vom Startknoten zum aktuellen Knoten  $u$ )

Dijkstra-Algorithmus

# Heuristische Suche – Motivation

Heuristik: Effizienzsteigerung durch Zusatzinformationen  
(z.B. Erfahrungswerte)

Anwendung bei

- ▶ Aufgaben mit mehreren Lösungen (z.B. Wege in Graphen)
- ▶ unterschiedliche Qualität der Lösungen  
(z.B. Länge des Weges)
- ▶ Suche nach **optimalen** Lösungen (z.B. kürzester Weg)
- ▶ falls vollständige Suche zu aufwendig

Ziele:

- ▶ Wahl einer geeigneten Such-Reihenfolge, unter welcher gute Lösungen zuerst gefunden werden
- ▶ Verwerfen von Knoten, die wahrscheinlich nicht zu einer Lösung führen  
(beabsichtigte Verletzung der Fairness-Eigenschaft)

# Schätzfunktionen

Ziel: sinnvolle Auswahl der in jedem Schritt zu expandierenden Knoten unter Verwendung von Zusatzinformationen

**Schätzfunktion** (heuristische Funktion)  $h : V \rightarrow \mathbb{R}_{\geq 0} \cup \{\infty\}$   
(oder in eine andere geordnete Menge)  
Schätzung der erwartete Restkosten vom Knoten  $u$   
bis zum Ziel

repräsentiert die Zusatzinformation

# Eigenschaften von Heuristiken

Schätzfunktion  $h : V \rightarrow \mathbb{R}_{\geq 0} \cup \{\infty\}$  heißt

**perfekt** (Schätzfunktion  $H(u)$ ), gdw.  $\forall u \in V : H(u) =$   
exakte Kosten einer optimalen Lösung durch  $u$   
( $H(u) = \infty$ , falls keine Lösung über  $u$  existiert)

**zielerkennend** gdw. für jeden Lösungsknoten  $u \in V$  gilt  $h(u) = 0$

**sicher** gdw. aus jedem Knoten  $u \in V$  mit  $h(u) = \infty$  ist  
kein Lösungsknoten erreichbar  
d.h.  $\forall u : (h(u) = \infty \rightarrow H(u) = \infty)$

**konsistent** gdw. für jeden Knoten  $u \in V$  und alle Folgeknoten  $v$   
von  $u$  gilt  $h(u) \leq w(u, v) + h(v)$   
( $w(u, v)$  Kosten des Übergangs von  $u$  nach  $v$ )

**nicht-überschätzend** gdw. für jeden Knoten  $u \in V$  gilt  
 $h(u) \leq H(u)$

Aus nicht-überschätzend folgt sicher und zielerkennend. (ÜA)

Aus zielerkennend und konsistent folgt nicht-überschätzend. (ÜA)

# Besten-Suche

(best-first-search)

Allgemeines Suchverfahren mit Bewertungsfunktion

$$f : V \rightarrow \mathbb{R}_{\geq 0} \cup \{\infty\}$$

mit folgender Strategie zur Auswahl der in jedem Schritt zu expandierenden Knoten:

- ▶ Knoten werden aufsteigend nach Bewertung  $f(u)$  expandiert,
- ▶ Expansion des Knotens  $u$  mit dem geringsten Wert  $f(u)$  zuerst
- ▶ Verwaltung von  $L_a$  als priority queue

Beispiel: Suche eines kürzesten Weges zwischen Orten A und B

- ▶ Bewertungsfunktion  $f(u)$ : bisherige Kosten bis zum Ort  $u$  (ohne Schätzfunktion, uniforme Kostensuche, Dijkstra)
- ▶ Bewertungsfunktion  $f(u)$ :  
Luftlinienentfernung des Ortes  $u$  von B (nur Schätzfunktion)

# Besten-Suche – Eigenschaften

zwei Methoden:

1. Knoten mit großen Werten **möglichst spät** expandieren
2. Knoten mit großen Werten **nicht** expandieren

- ▶ Bestensuche mit einer beliebigen Bewertungsfunktion ist nicht immer optimal.
- ▶ Bestensuche nach Methode 1 (fair) ist vollständig.
- ▶ Bestensuche nach Methode 2 ist nicht immer vollständig.

## Greedy-Suche (kleinste Restkosten)

Idee: Suche zuerst in Teilbäumen der noch nicht besuchten Knoten mit den geringsten (geschätzten) noch aufzuwendenden Kosten

Heuristische Funktion  $h : V \rightarrow \mathbb{R}_{\geq 0} \cup \{\infty\}$

$h(v)$  ist Abschätzung des von Knoten  $v$  aus den **noch notwendigen** Kosten zum Erreichen eines Zielzustandes

**Greedy-Suche:**

Besten-Suche mit Bewertungsfunktion  $f : V \rightarrow \mathbb{R}_{\geq 0} \cup \{\infty\}$ , wobei für jeden Knoten  $v \in V$  gilt

$$f(v) = h(v)$$

Eigenschaften der Greedy-Suche:

- ▶ optimal?
- ▶ vollständig?

# Bisherige Kosten

Kostenfunktion  $k : V \rightarrow \mathbb{R}_{\geq 0}$

$k(u)$  Kosten des besten (bisher bekannten) Pfades vom Startzustand zum Zustand  $u$

Kostenfunktion  $k : V \rightarrow \mathbb{R}_{\geq 0}$  heißt

**streng monoton wachsend**, falls für alle Knoten  $u$  und alle Nachfolger  $v$  von  $u$  gilt  $k(u) < k(v)$

Beispiele für Kostenfunktionen:

- ▶ Tiefe des Knotens im Suchbaum,
- ▶ maximale Entfernung vom Startknoten

## A\*-Suche (kleinste Gesamtkosten)

Idee: Suche zuerst in Teilbäumen der noch nicht besuchten Knoten mit dem **geringsten Wert der Schätzfunktion**

(Summe von bisherigen und geschätzten zukünftigen Kosten)

Funktionen

- ▶  $k : V \rightarrow \mathbb{R}_{\geq 0}$  – geringste bisher bekannte Kosten von einem Startzustand zu  $v$
- ▶  $h : V \rightarrow \mathbb{R}_{\geq 0}$  – geschätzte (geringste) Kosten von  $v$  zu einem Endzustand (Lösung)

**A\*-Suche:**

Besten-Suche mit Schätzfunktion  $f : V \rightarrow \mathbb{R}_{\geq 0}$ , wobei für jeden Knoten  $v \in V$  gilt

$$f(v) = k(v) + h(v)$$

**IDA\*-Suche:** Kombination von

- ▶ schrittweiser Vertiefung (iterative deepening)
- ▶ A\*-Suche

# Anwendungen

Planungsprobleme und kombinatorische Suchprobleme, z.B.

- ▶ Routenplanung
- ▶ TSP
- ▶ Verlegen von Leitungen
- ▶ Schaltkreis-Layout
- ▶ Scheduling
- ▶ Produktionsplanung
- ▶ Navigation (z.B. autonomer Fahrzeuge)

## Beispiel Schiebefax

- ▶ Zustände  $u \in \{0, \dots, 8\}^{3 \times 3}$ ,  $3 \times 3$ -Matrix mit Einträgen  $\{0, \dots, 8\}$  (jede Zahl genau einmal, 0 leeres Feld)
- ▶ Zulässige Züge: Verschieben des leeren Feldes auf ein Nachbarfeld d. h. Vertauschen von 0 und einem Wert in einem Nachbarfeld (gleicher Zeilen- oder Spaltenindex)
- ▶ Zielkonfiguration

1	2	3
8		4
7	6	5

- ▶ Aufgabeninstanz: gegebene Ausgangskonfiguration (Matrix), z.B.

8		3
2	1	4
7	6	5

- ▶ Lösung: Folge von zulässigen Zügen (Bewegung der Lücke 0) von der Ausgangs- zur Zielkonfiguration
- ▶ Bewertung der Lösung: Anzahl der Züge (Länge der Lösungsfolge)

# Schiebefax – Heuristische Funktionen

Heuristische Funktionen  $h_i : \{0, \dots, 8\}^{3 \times 3} \rightarrow \mathbb{N}$  mit

- $h_1$  Anzahl der Zahlen, die sich nicht an ihrer Zielposition befinden
- $h_2$  weitester Abstand einer Zahl zu ihrer Zielposition
- $h_3$  Summe der Manhattan-Abstände jeder Zahl zu ihrer Zielposition

Tafel: Bestensuche mit Bewertungsfunktionen  $f(u) = h_i(u)$

Qualität der Schätzfunktionen:

- ▶ gute Trennung verschiedener Zustände
- ▶ fair: zu jedem  $n \geq 0$  existieren nur endlich viele  $u \in V$  mit  $h(u) \leq n$

# Was bisher geschah

- ▶ Daten, Information, Wissen
- ▶ Wissensrepräsentation und -verarbeitung
- ▶ Wissensbasierte Systeme

## Wissensrepräsentation:

- ▶ Zustandsübergangssystem:  
Graph mit markierten Knoten  
(Zustände und deren Eigenschaften)
- ▶ Startzustand
- ▶ Eigenschaften der Zielzustände

**Lösung:** Pfad vom Start- zu einem Zielzustand

## Wissensverarbeitung: Suche im Graphen

- uninformiert:** Breiten-, Tiefen-, Gleiche-Kosten-Suche
- informiert:** Heuristik, Greedy-, A\*-Suche

# Zwei-Personen-Spiele

Zwei-Personen-(Brett)spiel:

- ▶ aktueller Spielzustand immer für beide Spieler sichtbar (vollständige Information)
- ▶ einer gewinnt, der andere verliert (Nullsummenspiel)

Wissensrepräsentation (Spielbaum):

- ▶ Menge von Zuständen (Min- und Max-Zustände)
- ▶ Startzustand
- ▶ Endzustände (ohne Fortsetzung)
- ▶ Nachfolgermenge  $S(v)$  = Menge von Zuständen (nach zulässigen Zügen)
- ▶ Bewertungsfunktion: Menge der Endzustände  $\rightarrow \mathbb{Z}$ 
  - ▶ positiv: Spieler (1, Max, beginnt) gewinnt
  - ▶ negativ: Gegner (0, Min) gewinnt

## Beispiel Nim (Variante)

- ▶  $n$  Münzen auf einem Stapel
- ▶ Spielzug: Teilen eines Stapels in zwei nichtleere Stapel ungleicher Größe
- ▶ Sobald ein Spieler keinen Zug mehr ausführen kann, hat er verloren (und der andere gewonnen).

(eine mögliche) Modellierung als Zustandsübergangssystem:

**Zustände:**  $S : \mathbb{N} \rightarrow \mathbb{N}$  (Multimenge)

Münzanzahl  $\mapsto$  Anzahl der Stapel mit dieser Zahl an Münzen

**Startzustand:**  $S(n) = 1 \wedge \forall i \neq n : S(i) = 0$

**Endzustände:** kein Zug möglich

**Übergänge:** (erlaubte Züge) für  $x = x_1 + x_2 \wedge x_1 \neq x_2 \wedge x_1 x_2 \neq 0$ :

$S \rightarrow S'$  mit

$$S'(x) = S(x) - 1$$

$$\wedge S'(x_1) = S(x_1) + 1 \wedge S'(x_2) = S(x_2) + 1$$

$$\wedge \forall i \in \mathbb{N} \setminus \{x, x_1, x_2\} : S'(i) = S(i)$$

# Minimax-Werte in vollständigen Spielbäumen

- ▶ vollständiger Spielbaum  $B = (V, E)$
- ▶ Bewertung der Endzustände (Blätter im Spielbaum) bekannt
- ▶ Fortsetzung der Bewertungsfunktion von den Blättern auf alle Knoten im Spielbaum  $b : V \rightarrow \mathbb{Z}$

rekursive Berechnung (Minimax-Algorithmus) des Wertes eines Knotens  $v$  im Spielbaum:

$$m(v) = \begin{cases} b(v) & \text{falls } v \text{ Endzustand} \\ \max\{m(u) \mid u \in S(v)\} & \text{falls } v \text{ Max-Knoten} \\ \min\{m(u) \mid u \in S(v)\} & \text{falls } v \text{ Min-Knoten} \end{cases}$$

Beispiele (Tafel):

- ▶ Spielbaum,
- ▶ Nim mit  $n = 6$

Spielstrategie für Spieler 1 (Max):

Zug wählen, der zum Zustand mit höchstem Minimax-Wert führt

# Minimax-Werte mit Heuristik

bei unvollständigem Spielbaum: Kombination von

- ▶ heuristischer Knotenbewertung
- ▶ Berechnung der Minimax-Werte

Beispiele (Tafel): Tic-Tac-Toe

mit Schätzfunktion für den Spieler am Zug:

Differenz der Anzahlen der noch nicht blockierten Gewinntripel

auch dabei Spielstrategie für Spieler 1 (Max):

Zug wählen, der zum Zustand mit höchstem Minimax-Wert führt

## $\alpha$ - $\beta$ -Suche

Idee: Tiefensuche mit Verwaltung zusätzlicher Werte

$\alpha$  : bisher höchster Minimax-Wert an Max-Positionen

$\beta$  : bisher geringster Minimax-Wert an Min-Positionen

Bei Berechnung des Minimax-Wertes der Wurzel eines Teilbaumes  
Berechnungen für Enkel auslassen, sobald bekannt ist, dass sie  $\alpha$   
und  $\beta$  nicht verbessern können

$\alpha$ - $\beta$ -Pruning: Abtrennen jedes Kindes  $v$  eines

**min-Knotens**  $u$ , falls  $\beta(u) \leq \alpha(v)$

(min-Spieler kann durch Wahl eines zuvor  
untersuchten Kindes von  $u$  den geringeren  
Minimax-Wert  $\beta(u)$  erreichen als durch Wahl von  $v$ )

**max-Knotens**  $u$ , falls  $\alpha(u) \geq \beta(v)$

(max-Spieler kann durch Wahl eines zuvor  
untersuchten Kindes von  $u$  den höheren  
Minimax-Wert  $\alpha(u)$  erreichen als durch Wahl von  $v$ )

Beispiel (Tafel)

# Automatische Berechnung heuristischer Funktionen

Ziel: Bewertung von Spielzügen, d.h.

Bewertung von Knoten  $v$  (Spielsituationen) im Spielbaum

Beispiel: **Monte-Carlo-Baum-Suche** MCTS

Idee: Berechnung des Wertes für  $v$  aus simulierten Spielen

- ▶ simuliertes Spiel  $i$ : Folge von (zufälligen) Zügen bis Spielende
- ▶ Bewertung des Knotens  $v$  im Spielbaum durch
  - ▶  $n$  simulierte Spiele, beginnend in  $v$
  - ▶  $\forall i \in \{1, \dots, n\}$ :  $R_i =$  Ergebnis des Spieles  $i$
  - ▶ Berechnung des Wertes von  $v$  aus Ergebnissen  $\{R_1, \dots, R_n\}$   
z.B. Mittel, Gewinnwahrscheinlichkeit

übliche Modifikationen:

- ▶ Integration von spezifischem Wissen (z.B. Standard-Antworten, Eröffnungsbibliotheken) statt ausschließlich zufälliger Züge
- ▶ Spiele nur bis festgelegter Anzahl von Zügen
- ▶ Backpropagation: Update der Werte auf Pfad zur Wurzel
- ▶ Speichern und laufendes Anpassen schon berechneter Bewertungen von Spielsituationen über mehrere Spiele (Lernen)

Beispiele: Monte-Carlo Go (1993), AlphaGo (erfolgreich 2015/16)

# Was bisher geschah

Abgrenzung der Begriffe: Daten, Information, Wissen, Intelligenz  
Symbolische KI: Suchprobleme

Wissensrepräsentation (Darstellung des Kontextes):

Zustandsübergangssystem, Zielbedingung,  
(Verfahren zur Bestimmung der) Werte der Zielknoten

Lösung: Pfad zu einem Zielzustand im Zustandsübergangssystem  
Spielstrategie

Wissensverarbeitung (Lösungsverfahren): Pfadsuche (informiert,  
uninformiert), MiniMax-Werte,  $\alpha$ - $\beta$ -Suche

Heuristische Funktionen:

- ▶ notwendig für informierte Suche
- ▶ Eigenschaften
- ▶ mitunter automatische Berechnung möglich

Beispiel für Kombination symbolischer und statistischer Verfahren

- ▶ heuristische Spielbaum-Suche (symbolisch)
- ▶ automatische Bewertung der Knoten durch Simulation mehrerer Spiele (statistisch)

Monte-Carlo-Baum-Suche (MCTS)

# Entscheidungsunterstützung

Ziel der KI: intelligente **Entscheidungen** treffen oder vorschlagen  
(analog menschlichen Experten)

Entscheidung:

Auswahl einer aus mehreren Optionen abhängig von der (aktuellen) Situation, z.B.

- ▶ nächster zu expandierender Knoten im Suchbaum
- ▶ nächster Spielzug
- ▶ Einordnung von Objekten
- ▶ Diagnosen
- ▶ Kreditwürdigkeit
- ▶ Therapieansätze

# Bewertung von Objekten / Situationen

Ziel: Bewertung von **Objekten** (Fällen) anhand bestimmter **Merkmale** für Mengen  $O$  aller Objekte,  $W$  aller möglichen Werte

**Funktion**  $f : O \rightarrow W$ , z.B.

- ▶  $O$ : Knoten in Agenda,  $W$ : Priorität  $\in \mathbb{R}_{\geq 0}$
- ▶  $O$ : mögliche Spielzüge,  $W$ : Minimax-Wert  $\in \mathbb{R}$
- ▶  $O$ : Personen,  $W$ : Alter  $\in \mathbb{N}$
- ▶  $O$ : Personen,  $W$ : Geschlecht  $\in \{m, w, d\}$
- ▶  $O$ : Belegungen  $\beta : P \rightarrow \{0, 1\}$ ,  $W$ : Wahrheitswert  $\in \{0, 1\}$
- ▶  $O$ : digitale Bilder,  $W$ :  $\in 2^{\{\text{Katze, Hund, Maus}\}}$

Bewertung der Objekte anhand ihrer **Merkmale**, z.B.

- ▶ Position, bisherige und geschätzte zukünftige Kosten
- ▶ Eigenschaften des Spielzustandes, z.B. noch nicht blockierte Tripel
- ▶ Wahrheitswert  $\llbracket \varphi \rrbracket_{\beta} \in \{0, 1\}$  für gegebene Formel  $\varphi \in \text{AL}(P)$
- ▶ Anordnung der Pixel / Farbwerte im Bild (Matrix)

# Klassifikation

Ziel: Einteilung von **Objekten** (Fällen) anhand bestimmter **Merkmale** in **Klassen** (Auswahl einer Lösung für einen Fall aus einer Menge gegebener Alternativen)

**Klassifikation**: Bewertung / **Funktion**  $f : O \rightarrow W$

mit diskreter (meist endlicher) Menge  $W$  von **Klassen**

- M** Menge von Merkmalen  $m$  (Symptome, Attribute)  
jedes Merkmal  $m$  mit zugeordneter Menge  $V_m$  möglicher Werte  
z.B.  $M = \{a, g\}$  mit Alter  $a$ ,  $V_a = \mathbb{N}$ , Geschlecht  $g$ ,  
 $V_g = \{m, w, d\}$   
Merkmalsraum:  $\times_{m \in M} V_m$
- O**  $\subseteq M$  Menge aller Objekte (Fälle)  
Jedem Objekt  $o \in O$  sind seine Merkmalswerte  
 $m_o \in \times_{m \in M} V_m$  zugeordnet.  
z.B. Tom  $\in O$  mit  $m_{\text{Tom}} = (5, m)$
- K** (=  $W$ ) Menge aller Klassen (Diagnosen, Lösungen)  
jede Klasse repräsentiert eine Teilmenge des  
Merkmalsraumes  $\times_{m \in M} V_m$   
z.B. Baby, Kind, Mädchen

# Beispiel

**Merkmale** (Attribute) mit Menge möglicher Werte (Ausprägungen)

Beispiele:

- ▶ Name: Tom, Tina, Anna, Paul
- ▶ Geschlecht, Werte: w, m, d
- ▶ Alter: natürliche Zahl

**Objekte** definiert durch Zuordnung:

Merkmal  $\rightarrow$  Wert

Beispiele:

- ▶ Name  $\mapsto$  Tina, Alter  $\mapsto$  7, Geschlecht  $\mapsto$  w
- ▶ Name  $\mapsto$  Paul, Alter  $\mapsto$  87, Geschlecht  $\mapsto$  m

**Klassen** Kombination von Merkmalswerten

definiert Menge von Objekten

Beispiele:

- ▶ Rentner: Geschlecht = m und Alter  $\geq 65$
- ▶ Mädchen: Geschlecht = w und Alter  $\leq 14$
- ▶ Test vor Röntgen: Geschlecht = w und Alter  $\in \{15, 50\}$

# Klassifikationsprobleme

**gegeben:** Objekt  $o \in O$  mit Merkmalswerten  $m_o$   
Zuordnung  $K \rightarrow X_K$  mit  $X_K \subseteq M$  (Extension)  
jede Klasse definiert durch Merkmalswerte (oft  
Intervalle)

**gesucht:** Zuordnung Objekt  $o \in O$  zu Klasse  $K$  mit  $m_o \in X_K$

Beispiele:

- ▶ Klassifikation von
  - ▶ Objekte: Adler, Biber, Elefant, Fledermaus, Gorilla, Hecht, Pinguin, Specht
  - ▶ Merkmale: Federn, Schuppen, Fell, kann fliegen, kann schwimmen, legt Eier
  - ▶ Klassen: Säugetier, Fisch, Vogel

# Diagnose-Probleme

Spezialfall von Klassifikationsproblemen:

- ▶ Objekte: Fälle
- ▶ Merkmale: Fragen, Tests
- ▶ Merkmalswerte: Antworten
- ▶ Klassen: Diagnosen

Anwendungen, z.B. in

- ▶ Medizin:
  - ▶ Krankheitserkennung
  - ▶ Entscheidung für eine Therapie
  - ▶ Auswertung von Studien
- ▶ Technische Systeme:
  - ▶ Konfiguration
  - ▶ Feststellung von Störungen
  - ▶ Entscheidung für Vorgehen bei Behebung

# Statistische Verfahren

Einsatz zum Lösen von Problemen,

- ▶ die unvollständig beschrieben sind
- ▶ die keine eindeutige Lösung haben
- ▶ für die keine effizienten exakten Algorithmen bekannt sind

einige Ansätze:

- ▶ Fuzzy-Logik, probabilistische Logik
- ▶ **Künstliche neuronale Netze**
- ▶ Evolutionäre Algorithmen
- ▶ Schwarm-Intelligenz

**maschinelles Lernen:**

- ▶ Ansätze zur ständigen Selbstverbesserung von Verfahren zum Problemlösen
- ▶ derzeit teilweise erfolgreicher Einsatz auf den Teilgebieten Bewertung, Klassifikation

## (Natürliches und ) Maschinelles Lernen

(Schrittweise) Änderung eines Systems (Verfahrens zur Problemlösung), so dass es bei der zukünftigen Anwendung dasselbe oder ähnliche Probleme besser löst.

- ▶ Aufgaben (Problem): Menge von Eingaben
- ▶ Aufgabeninstanz: Eingabe
- ▶ Lösung der Instanz: Ausgabe
- ▶ Bewertung der Lösung: Zuordnung Lösung  $\rightarrow$  Güte

Schritte bei der Lösung von Aufgabeninstanzen mit Lerneffekt:  
Schüler (System) führt wiederholt aus:

1. verwendet ein Lösungsverfahren  $V$  für diese Aufgabe
2. bestimmt eine Lösung  $l$  der gegebenen Aufgabeninstanz
3. bestimmt (oder erfährt) eine Bewertung dieser Lösung  $l$
4. modifiziert das Lösungsverfahren  $V$  zu  $V'$ , um (in Zukunft) Lösungen mit besseren Bewertungen zu finden
5. wendet im nächsten Schritt zur Lösung dieser Aufgabe das Lösungsverfahren  $V'$  an

Lernen: Schritte 3 und 4

# Lernverfahren

## Lernen durch

- ▶ Auswendiglernen (gegebener Beispiele)
- ▶ Nachahmen
- ▶ Anleitung (Anweisungen)
- ▶ logische Ableitung neuer Lösungsverfahren
- ▶ Analogie (zu gegebenen Beispielen)  
anhand Ähnlichkeit
- ▶ Erfahrung (durch gegebene Beispiele)  
Fähigkeit zur Verallgemeinerung
- ▶ Probieren und Beobachten  
(Erzeugen eigener Beispiele)

## nach Art des Lernenden:

- ▶ natürliches Lernen
- ▶ maschinelles (künstliches) Lernen

# Lernen durch gegebene Beispiele

nach der zum Lernen verwendbaren Information:

überwachtes Lernen (supervised learning)

    korrigierendes Lernen (corrective learning)

    bestärkendes Lernen (reinforcement learning)

unüberwachtes Lernen (unsupervised learning)

gewünschte Eigenschaften des Löseverfahrens:

- ▶ Korrektheit  
der Lösungen für die gegebenen Beispiele
- ▶ Generalisierung  
„sinnvolle“ Lösungen für ähnliche Aufgaben

# Korrigierendes Lernen

**Trainingsmenge:** Menge von Paaren (Eingabe, Ausgabe)  
(partielle Funktion an Stützstellen)

**Lernziel:** (möglichst einfache) Funktion, die an den  
Stützstellen mit der Trainingsmenge übereinstimmt

**Rückmeldung:** Trainer sagt nach jedem Lernschritt die korrekte  
Ausgabe.

**Prinzip:** Lernen durch Nachahmen (mit Korrektur)

Anwendung z.B. bei

- ▶ Klassifizierung (Zuordnung von Objekten / Fällen zu Klassen,  
abhängig von den Merkmalen der Objekte)  
z.B. Zuordnung Sensorwerte → Alarmklasse  
Trainingsmenge ist  
Menge von Paaren (Objekteigenschaften, Klasse)
- ▶ Lernen von Funktionen: Trainingsmenge ist  
Menge von Paaren (Parameter, Funktionswert)

## Bestärkendes Lernen (reinforcement learning)

**Trainingsmenge:** Menge von Paaren (Eingabe, Erfolg  $\in \{\text{ja, nein}\}$ )

**Lernziel:** (möglichst einfache) Funktion, die den Stützstellen korrekte Werte zuordnet

**Rückmeldung:** Trainer sagt nach jedem Lernschritt, ob die Ausgabe korrekt war.

**Idee:** Lernen durch Probieren

- ▶ **Klassifizierung:** Trainingsmenge ist Menge von Objekten (mit ihren Eigenschaften)  
Bewertung der Lösung: ja, falls Zuordnung zur korrekten Klasse, sonst nein
- ▶ **Lernen von Plänen** (Anlagestrategien, Bewegungsabläufe usw.)  
z.B. Steuern eines autonomen Fahrzeuges  
Trainingsmenge: Strecke( $n$ ),  
Folge von Paaren (Sensordaten, Steuersignale)  
Bewertung der Lösung: ja, falls Plan zum Erfolg geführt hat  
(z.B. Fahrzeug fährt  $> n$  km ohne Eingriff) , sonst nein

# Unüberwachtes Lernen

Trainingsmenge: Menge von Eingaben

- Lernziel:
- ▶ Gruppierung ähnliche Muster
  - ▶ oft auch topologisch sinnvolle Anordnung

Idee: Lernen ohne Trainer (ohne Rückmeldung)

- ▶ Entdecken von Strukturen
- ▶ Selbstorganisation von Objekten zu Gruppen (mit gemeinsamen Merkmalen, typische Vertreter)
- ▶ topologierhaltende Abbildungen (z.B. Körperteile → Gehirnregionen)
- ▶ Assoziation (z.B. in Schrifterkennung)

# Neuronale Netze

Neuron – Nerv (griechisch)

Modellierung und Simulation der Strukturen und Mechanismen im Nervensystem von Lebewesen

Biologisches Vorbild	Mathematisches Modell
Nervenzellen (Neuronen)	künstliche Neuronen
Struktur (eines Teiles) eines Nervensystems	künstliche neuronale Netze (KNN) unterschiedlicher Struktur
Aktivierung von Neuronen, Reizübertragung	künstlichen Neuronen zugeordnete Funktionen
Anpassung (Lernen)	Änderungen verschiedener Parameter des KNN

# Natürliche Neuronen

ZNS besteht aus miteinander verbundenen Nervenzellen (Neuronen)

Struktur eines Neurons:

- ▶ Zellkörper
- ▶ Dendriten
- ▶ Synapsen (verstärkende, hemmende)
- ▶ Axon

## Natürliche Neuronen – Funktionsweise

Informationsübertragung durch elektrochemische Vorgänge:

- ▶ aktivierte Zelle setzt an Synapsen Neurotransmitter frei,
- ▶ Neurotransmitter ändern die Durchlässigkeit der Zellmembran für Ionen an den Dendriten der empfangenden Zelle,
- ▶ Potential innerhalb der empfangenden Zelle ändert sich durch diffundierende Ionen,
- ▶ überschreitet die Summe der an allen Synapsen entstandenen Potentiale (Gesamtpotential) der Zelle einen Schwellwert, entsteht ein Aktionspotential (Zelle feuert),
- ▶ Aktionspotential (Spannungsspitze) durchquert das Axon (Nervenfasern) zu den Synapsen zu Nachbarzellen,
- ▶ aktivierte Zelle setzt an Synapsen Neurotransmitter frei, usw.

Stärke der Information durch Häufigkeit der Spannungsspitzen (Frequenzmodulation).

# Eigenschaften natürlicher neuronaler Netze

- ▶ geringe Taktrate  $10^{-3}$  s
- ▶ parallele Arbeit sehr vieler ( $10^{11}$ ) Neuronen
- ▶ Neuronen sehr stark miteinander vernetzt (ca. 10 000 Nachbarn)
- ▶ Verarbeitungseinheit = Speicher

## Vorteile:

- ▶ hohe Arbeitsgeschwindigkeit durch Parallelität,
- ▶ Funktionsfähigkeit auch nach Ausfall von Teilen des Netzes,
- ▶ Lernfähigkeit,
- ▶ Möglichkeit zur Generalisierung

Ziel: Nutzung dieser Vorteile zum Problemlösen durch Wissensrepräsentation als künstliche neuronale Netze

# Natürliche Neuronen – Lernen

Speicherung von Informationen durch Anpassung der Durchlässigkeit (Leitfähigkeit) der Synapsen

- ▶ **Regel von Hebb** (1949):  
Synapsen zwischen gleichzeitig aktiven Zellen werden immer durchlässiger (Reizschwelle wird verringert),  
Verbindung an dieser Synapse wird stärker
- ▶ lange nicht benutzte Synapsen verlieren mit der Zeit ihre Durchlässigkeit  
Verbindung an dieser Synapse wird schwächer.

# Anwendungen künstlicher neuronaler Netze

## Anwendungsgebiete:

- ▶ Bildverarbeitung, z.B.
  - ▶ Objekterkennung
  - ▶ Szenenerkennung
  - ▶ Schrifterkennung
  - ▶ Kantenerkennung
- ▶ Medizin, z.B. Auswertung von Bildern, Langzeit-EKGs
- ▶ automatische Spracherkennung
- ▶ Sicherheit, z.B. Biometrische Identifizierung
- ▶ Wirtschaft, z.B. Aktienprognosen, Kreditrisikoabschätzung
- ▶ Robotik, z.B. Lernen von Bewegungsabläufen
- ▶ Steuerung autonomer Fahrzeuge

# Geschichte künstlicher neuronaler Netze

- ▶ 1943, Warren McCulloch, Walter Pitts:  
A logical calculus of the ideas immanent in nervous activity
- ▶ 1949, Donald O. Hebb: Lernmodell  
The organization of behaviour
- ▶ 1957 Frank Rosenblatt: Perzeptron (1 Schicht)  
erster Neurocomputer MARK 1  
(Ziffernerkennung in  $20 \times 20$ -Bildsensor)
- ▶ 1969, Marvin Minsky, Seymour Papert: Perceptrons
- ▶ 1971 Perzeptron mit 8 Schichten
- ▶ 1974 Backpropagation (Erfindung)
- ▶ 1982, Teuvo Kohonen: selbstorganisierende Karten
- ▶ 1982, John Hopfield: Hopfield-Netze
- ▶ 1985, Backpropagation (Anwendung)
- ▶ 1997, long short-term memory (Erfindung)
- ▶ 2000, Begriff Deep Learning für KNN, Faltungsnetze (CNN)
- ▶ 2006, long short-term memory (Anwendung)
- ▶ 2009, verstärkt Training mit GPUs
- ▶ 2015, AlphaGo, AlphaZero, ...

# Was bisher geschah

Abgrenzung der Begriffe:

- ▶ Daten, Information, Wissen, Lernen, Intelligenz
- ▶ Symbolische / statistische KI

KI – grundsätzliches Ziel:

- ▶ sinnvolle **Entscheidungen** treffen / unterstützen  
z.B. Spielzug, -strategie, Diagnose, Therapie

KI – grundlegende Methode:

- ▶ **Bewertung / Klassifikation**  
von Optionen (Objekte, Diagnosen, Aktionen)  
anhand ihrer **Merkmale**

Verfahren (bisher):

- ▶ WH: Suchprobleme (überwiegend symbolisch)  
uninformiert / informiert mit / ohne Gegenspieler
- ▶ Maschinelles Lernen (statistisch)  
überwacht (korrigierend / reinforcement) / unüberwacht
- ▶ Künstliche neuronale Netze (statistisch)  
biologische Grundlagen und Modell

# Künstliche Neuronen:

## McCulloch-Pitts-Neuron ohne Hemmung

einfaches abstraktes Neuronenmodell von  
McCulloch und Pitts, 1943

Aufbau eines künstlichen Neurons  $u$  (Tafel)

Eingabe:	$x = (x_1, \dots, x_{m_u}) \in \{0, 1\}^{m_u}$	(ankommende Reize)
Schwellwert:	$\theta_u \in \mathbb{R}$	(Reizschwelle)
Ausgabe:	$f(x_1, \dots, x_{m_u}) \in \{0, 1\}$	(weitergegebener Reiz)

Parameter eines McCulloch-Pitts-Neurons  $u$  ohne Hemmung:

- ▶  $m_u$ : Anzahl der (erregenden) Eingänge
- ▶  $\theta_u$ : Schwellwert

# McCulloch-Pitts-Neuron ohne Hemmung: Funktionen

**Eingangsfunktion** des Neurons  $u$ :  $I_u: \{0, 1\}^{m_u} \rightarrow \mathbb{R}$  mit

$$I_u(x_1, \dots, x_{m_u}) = \sum_{i=1}^{m_u} x_i$$

(Summe aller erregenden Eingänge des Neurons  $u$ )

**Aktivierungsfunktion** des Neurons  $u$

(abhängig vom Schwellwert  $\theta_u$ ):  $A_u: \mathbb{R} \times \mathbb{R} \rightarrow \{0, 1\}$  mit

$$A_u(\theta_u, v) = \begin{cases} 1 & \text{falls } v \geq \theta_u \\ 0 & \text{sonst} \end{cases}$$

(Stufenfunktion mit Stufe bei  $\theta_u$ )

**Ausgabefunktion** des Neurons  $u$ :  $O_u: \{0, 1\} \rightarrow \{0, 1\}$  mit

$$O_u(v) = v$$

(Identität)

# McCulloch-Pitts-Neuron ohne Hemmung: Berechnung

vom Neuron  $u$  berechnete Funktion:  $f_u: \{0, 1\}^{m_u} \rightarrow \{0, 1\}$  mit

$$\begin{aligned} f_u(x_1, \dots, x_{m_u}) &= O_u(A_u(\theta_u, I_u(x_1, \dots, x_{m_u}))) \\ &= \begin{cases} 1 & \text{falls } \sum_{i=1}^{m_u} x_i \geq \theta_u \\ 0 & \text{sonst} \end{cases} \end{aligned}$$

$m_u$ -stellige Boolesche Funktion

# McCulloch-Pitts-Neuron ohne Hemmung: Beispiele

elementare Boolesche Funktionen  $\vee, \wedge$

mehrstellige  $\vee, \wedge$

Existiert zu jeder Booleschen Funktion  $f : \{0, 1\}^n \rightarrow \{0, 1\}$  ein McCulloch-Pitts-Neuron ohne Hemmung, welches  $f$  berechnet?

Nein, nur **monotone** Boolesche Funktionen,  
z.B.  $\neg$  nicht

Warum?

## Geometrische Interpretation

Jedes McCulloch-Pitts-Neuron  $u$  mit  $m_u$  Eingängen teilt die Menge  $\{0, 1\}^{m_u}$  in zwei Teilmengen:

$$\begin{aligned} f_u^{-1}(1) &= \{(x_1, \dots, x_{m_u}) \in \{0, 1\}^{m_u} \mid f(x_1, \dots, x_{m_u}) = 1\} \\ &= \{(x_1, \dots, x_{m_u}) \in \{0, 1\}^{m_u} \mid \sum_{i=1}^{m_u} x_i \geq \theta_u\} \end{aligned}$$

und

$$\begin{aligned} f_u^{-1}(0) &= \{(x_1, \dots, x_{m_u}) \in \{0, 1\}^{m_u} \mid f(x_1, \dots, x_{m_u}) = 0\} \\ &= \{(x_1, \dots, x_{m_u}) \in \{0, 1\}^{m_u} \mid \sum_{i=1}^{m_u} x_i < \theta_u\} \end{aligned}$$

geometrische Interpretation als Teilräume des  $R^m$

**Grenze** zwischen beiden Bereichen:

$(m_u - 1)$ -dimensionaler Teilraum  $\sum_{i=1}^{m_u} x_i = \theta$   
parallele Schnitte (abhängig von  $\theta$ )

# Geometrische Interpretation: Beispiele

Beispiele:

- ▶ Neuron  $u$  mit  $m_u = 2$  Eingängen und Schwellwert  $\theta_u = 1$

$$f_u(x_1, x_2) = \begin{cases} 1 & \text{falls } x_1 + x_2 \geq 1 \\ 0 & \text{sonst} \end{cases}$$

Bereich der  $x_1, x_2$ -Ebene mit  $f_u(x_1, x_2) = 1$  ist die Halbebene mit  $x_2 \geq 1 - x_1$ .

$x_2 = g(x_1) = 1 - x_1$  ist eine **lineare Trennfunktion** zwischen den Halbebenen mit  $f_u(x_1, x_2) = 0$  und  $f_u(x_1, x_2) = 1$ .

- ▶ Neuron  $v$  mit  $m_v = 3$  Eingängen und  $\theta_v = 1$

# Linear trennbare Funktionen

Zwei **Mengen**  $A, B \subseteq \mathbb{R}^n$  heißen genau dann **linear trennbar**, wenn eine lineare Funktion  $g : \mathbb{R}^n \rightarrow \mathbb{R}$  mit

$g(x_1, \dots, x_n) = a_0 + \sum_{i=1}^n a_i x_i$  existiert, so dass

- ▶ für alle  $(x_1, \dots, x_n) \in A$  gilt  $g(x_1, \dots, x_n) > 0$
- ▶ für alle  $(x_1, \dots, x_n) \in B$  gilt  $g(x_1, \dots, x_n) < 0$

(eindeutig beschreiben durch  $n + 1$ -Tupel  $(a_0, a_1, \dots, a_n)$  )

Eine **Boolesche Funktion**  $f : \{0, 1\}^n \rightarrow \{0, 1\}$  heißt genau dann **linear trennbar**, wenn die Mengen  $f^{-1}(0)$  und  $f^{-1}(1)$  linear trennbar sind.

Beispiele:  $\vee, \wedge, \neg x_1, x_1 \rightarrow x_2, x_1 \wedge \neg x_2$

Die Boolesche Funktion XOR ist nicht linear trennbar.

# McCulloch-Pitts-Neuron mit Hemmung

McCulloch-Pitts-Neuron  $u$  mit Hemmung:

Eingabewerte:  $x = (x_1, \dots, x_{m_u}) \in \{0, 1\}^{m_u}$       erregend  
 $y = (y_1, \dots, y_{m'_u}) \in \{0, 1\}^{m'_u}$       hemmend

Schwellwert:  $\theta_u \in \mathbb{R}$

Ausgabe:  $f(x_1, \dots, x_{m_u}, y_1, \dots, y_{m'_u}) \in \{0, 1\}$

Parameter eines McCulloch-Pitts-Neurons  $u$  (mit Hemmung):

- ▶  $m_u$ : Anzahl der erregenden Eingänge
- ▶  $m'_u$ : Anzahl der hemmenden Eingänge
- ▶  $\theta_u$ : Schwellwert

## Funktionen bei hemmenden Eingängen

**Eingangsfunktion** des Neurons  $u$ :  $I_u : \{0, 1\}^{m_u+m'_u} \rightarrow \mathbb{R} \times \mathbb{R}$

$$I_u(x_1, \dots, x_{m_u}, y_1, \dots, y_{m'_u}) = \left( \sum_{i=1}^{m_u} x_i, \sum_{i=1}^{m'_u} y_i \right)$$

(Summe aller erregenden Eingänge des Neurons  $u$ ,  
Summe aller hemmenden Eingänge des Neurons  $u$ )

**Aktivierungsfunktion** des Neurons  $u$  (abhängig von  $\theta_u$ ):

$A_u : \mathbb{R} \times (\mathbb{R} \times \mathbb{R}) \rightarrow \{0, 1\}$

$$A_u(\theta_u, (x, y)) = \begin{cases} 1 & \text{falls } x \geq \theta_u \text{ und } y \leq 0 \\ 0 & \text{sonst} \end{cases}$$

(Stufenfunktion)

**Ausgabefunktion** des Neurons  $u$ :  $O_u : \{0, 1\} \rightarrow \{0, 1\}$  mit

$$O_u(v) = v$$

(Identität)

## Berechnung bei hemmenden Eingängen

Gesamtfunktion des Neurons  $u$

$$f_u(x_1, \dots, x_{m_u}, y_1, \dots, y_{m'_u}) = O_u(A_u(\theta_u, I_u(x_1, \dots, x_{m_u}, y_1, \dots, y_{m'_u})))$$

Jedes McCulloch-Pitts-Neuron  $u$  mit  $m_u$  erregenden Eingängen,  $m'_u$  hemmenden Eingängen und Schwellwert  $\theta_u$  repräsentiert die Boolesche Funktion  $f_u : \{0, 1\}^{m_u+m'_u} \rightarrow \{0, 1\}$ :

$$f_u(x_1, \dots, x_{m_u}, y_1, \dots, y_{m'_u}) = \begin{cases} 1 & \text{falls } \sum_{i=1}^{m_u} x_i \geq \theta_u \\ & \text{und } \sum_{i=1}^{m'_u} y_i \leq 0 \\ 0 & \text{sonst} \end{cases}$$

Beispiele mit Hemmung:

- ▶ elementare Boolesche Funktion:  $\neg$
- ▶ komplexere Boolesche Funktionen, z.B.

$$x_1 \wedge \neg x_2$$

$$\neg x_1 \wedge x_2 \wedge x_3,$$

$$\neg(x_1 \vee \neg x_2 \vee \neg x_3)$$

# McCulloch-Pitts-Netze

McCulloch-Pitts-Netz:

gerichteter Graph mit

- ▶ McCulloch-Pitts-Neuronen als Ecken und
- ▶ gerichteten Kanten zwischen Neuronen  
zwei Arten: erregend, hemmend

Berechnung der Neuronen-Funktionen  
(entsprechend Struktur des Netzes):

- ▶ parallel
- ▶ sequentiell
- ▶ rekursiv

# McCulloch-Pitts-Netze

## Ein-Schicht-McCulloch-Pitts-Netz

parallele Schaltung mehrerer

McCulloch-Pitts-Neuronen

repräsentiert Boolesche Funktionen mit mehreren  
Ausgaben

Beispiel: Parallelschaltung von  $x_1 \wedge \neg x_2$  und  $\neg x_1 \wedge x_2$

## Mehr-Schicht-McCulloch-Pitts-Netz

parallele und sequentielle Schaltung mehrerer

McCulloch-Pitts-Neuronen

Beispiel: XOR

Analogie zu logischen Schaltkreisen

Jede Boolesche Funktion  $f : \{0, 1\}^n \rightarrow \{0, 1\}$  lässt sich durch ein  
McCulloch-Pitts-Netz berechnen.

McCulloch-Pitts-Netz mit zwei Schichten genügt  
(analog DNF, CNF in Aussagenlogik)

## Modifikationen von McCulloch-Pitts-Neuronen

- ▶ Durch Vervielfachung eines Einganges erhöht sich seine Wirkung (sein Gewicht).
- ▶ Vervielfachung (absolut) hemmender Eingänge ändert die berechnete Funktion nicht.
- ▶ relative Hemmung:  
hemmende Eingänge verhindern das Feuern der Zelle nicht völlig, sondern erschweren es (erhöhen den Schwellwert, negatives Gewicht).
- ▶ Absolute Hemmung lässt sich durch relative Hemmung mit großer Schwellwerterhöhung (auf Anzahl aller erregenden Eingänge +1) simulieren.
- ▶ Durch Einführung von Gewichten wird Trennung in hemmende und erregende Eingänge überflüssig.

# Parameter künstlicher Neuronen

verschiedene künstliche Neuronenmodelle unterscheiden sich in:

- ▶ Anzahl Typen der Ein- und Ausgabewerte,
- ▶ zulässige Gewichte an den Eingangskanten,
- ▶ Eingabe-, Ausgabe- und Aktivierungsfunktion

Jedes Neuron mit  $m$  Eingängen repräsentiert eine Funktion von  $m$  Eingabewerten

# Schwellwertneuronen

Idee: gewichtete Eingänge

- ▶ zur Modellierung der Stärke der synaptischen Bindung
- ▶ ermöglichen Lernen durch Änderung der Gewichte

Mathematisches Modell:

Schwellwertneuron (Perzeptron)

Eingabewerte:  $x = (x_1, \dots, x_m) \in \{0, 1\}^m$

Eingangsgewichte:  $w = (w_1, \dots, w_m) \in \mathbb{R}^m$

Schwellwert:  $\theta \in \mathbb{R}$

Ausgabe:  $a(x_1, \dots, x_m) \in \{0, 1\}$       Aktivität

Parameter eines Schwellwertneurons  $u$ :

- ▶  $m_u$ : Anzahl der (erregenden) Eingänge
- ▶  $(w_1, \dots, w_{m_u}) \in \mathbb{R}^{m_u}$ : Eingangsgewichte
- ▶  $\theta_u$ : Schwellwert

## Schwellwertneuronen: Funktionen

**Eingangsfunktion** des Neurons  $u$  (abhängig von  $(w_1, \dots, w_{m_u})$ ):

$I_u: \mathbb{R}^{m_u} \times \{0, 1\}^{m_u} \rightarrow \mathbb{R}$  mit

$$I_u(w_1, \dots, w_{m_u}, x_1, \dots, x_{m_u}) = \sum_{i=1}^{m_u} w_i x_i$$

(gewichtete Summe aller Eingänge des Neurons  $u$ )

**Aktivierungsfunktion** des Neurons  $u$  (abhängig von  $\theta_u$ ):

$A_u: \mathbb{R} \times \mathbb{R} \rightarrow \{0, 1\}$  mit

$$A_u(\theta_u, v) = \begin{cases} 1 & \text{falls } v \geq \theta_u \\ 0 & \text{sonst} \end{cases}$$

(Stufenfunktion)

**Ausgabefunktion** des Neurons  $u$ :  $O_u: \{0, 1\} \rightarrow \{0, 1\}$  mit

$$O_u(v) = v$$

(Identität)

## Schwellwertneuronen: Berechnung

vom Neuron  $u$  berechnete Funktion:  $f_u: \{0, 1\}^{m_u} \rightarrow \{0, 1\}$  mit

$$\begin{aligned} f_u(x_1, \dots, x_{m_u}) &= O_u(A_u(\theta_u, I_u(w_1, \dots, w_{m_u}, x_1, \dots, x_{m_u}))) \\ &= \begin{cases} 1 & \text{falls } \langle w, x \rangle \geq \theta_u \\ 0 & \text{sonst} \end{cases} \end{aligned}$$

Wiederholung:

$\sum_{i=1}^n w_i x_i = \langle w, x \rangle$  Skalarprodukt

der Vektoren  $w = (w_1, \dots, w_n)$  und  $x = (x_1, \dots, x_n)$

Jedes Schwellwertneuron  $u$  mit  $m_u$  Eingängen repräsentiert eine Boolesche Funktion  $f_u: \{0, 1\}^{m_u} \rightarrow \{0, 1\}$

Auch mit Schwellwertneuronen lassen sich nur linear trennbare Boolesche Funktionen berechnen (XOR nicht).

Beispiele:  $\vee, \wedge, \rightarrow, ((x_1 \wedge (x_3 \vee \neg x_2)) \vee (\neg x_2 \wedge x_3))$

## Schwellwertneuronen: geometrische Interpretation

Jedes Schwellwertneuron  $u$  mit  $m_u$  Eingängen teilt die Menge  $\{0, 1\}^{m_u}$  der **Eingabevektoren** (Punkte im  $\mathbb{R}^{m_u}$ ) in zwei Teilmengen (Teilräume des  $\mathbb{R}^{m_u}$ ):

$$\begin{aligned}f_u^{-1}(1) &= \{(x_1, \dots, x_{m_u}) \in \{0, 1\}^{m_u} \mid f(x_1, \dots, x_{m_u}) = 1\} \\ &= \{(x_1, \dots, x_{m_u}) \in \{0, 1\}^{m_u} \mid \langle w, x \rangle \geq \theta_u\}\end{aligned}$$

und

$$\begin{aligned}f_u^{-1}(0) &= \{(x_1, \dots, x_{m_u}) \in \{0, 1\}^{m_u} \mid f(x_1, \dots, x_{m_u}) = 0\} \\ &= \{(x_1, \dots, x_{m_u}) \in \{0, 1\}^{m_u} \mid \langle w, x \rangle < \theta_u\}\end{aligned}$$

Grenze: durch  $\langle w, x \rangle = \theta_u$  beschriebene  $(m_u - 1)$ -dimensionale Hyperebene (Teilraum)  
(parallele Schnitte)

## Schwellwert als Gewicht (Bias-Neuronen)

Neuron mit Schwellwert  $\theta$

Hinzufügen eines zusätzlichen Eingangs  $x_0$  (bias neuron)  
mit Wert  $x_0 = 1$  (konstant)

Gewicht des Einganges  $x_0$ :  $w_0 = -\theta$

$$\begin{aligned} \sum_{i=1}^n w_i x_i \geq \theta & \quad \text{gdw.} \quad \sum_{i=1}^n w_i x_i - \theta \geq 0 \\ & \quad \text{gdw.} \quad \sum_{i=0}^n w_i x_i \geq 0 \end{aligned}$$

# Überwachtes Lernen einzelner Schwellwertneuronenn

**Aufgabe:** Konstruktion eines Schwellwertneurons zur Berechnung einer Booleschen Funktion  
 $f : \{0, 1\}^m \rightarrow \{0, 1\}$

**Trainingsmenge:** Menge  $T$  von Paaren  $(x, t)$  aus

- ▶ Eingabevektoren  $x \in \{0, 1\}^m$  und
- ▶ Funktionswerten  $t = f(x) \in \{0, 1\}$

(Werte der Funktion  $f$  an Stützstellen)

**Struktur des Schwellwertneurons:** Schwellwertneuron mit  $m + 1$  Eingängen (bias  $x_0$ )  
und Eingangsgewichten  $(w_0, \dots, w_m) \in \mathbb{R}^{m+1}$

**Idee:** automatisches Lernen der Funktion durch (wiederholte) Änderung der Gewichte

**Lernziel:** Gewichte  $(w'_0, \dots, w'_m) \in \mathbb{R}^{m+1}$ , so dass das Schwellwertneuron die Funktion  $f$  berechnet  
(Korrektheit an Stützstellen)

# $\Delta$ -Regel

Idee: Lernen aus Fehlern (und deren Korrektur)

Delta-Regel:

$$\forall i \in \{0, \dots, m\} : w_i' = w_i + \Delta w_i \quad \text{mit} \quad \Delta w_i = \eta x_i (t - y)$$

- ▶ Trainingswert  $t$
- ▶ vom Netz berechneter Wert  $y$
- ▶ **Lernrate**  $\eta \in \mathbb{R}$  (Grad der Verstärkung der Verbindung)

korrigierendes Lernen,  
(falls  $x_i$  aktiv und  $y \neq t$ )

Beispiel:  $\neg, \wedge, \rightarrow$

## $\Delta$ -Lernverfahren für Schwellwertneuronen

- ▶ Beginn mit **zufälligen Eingangsgewichten**  $(w_0, \dots, w_n) \in \mathbb{R}^m$  (Schwellwert als Gewicht),
- ▶ die folgenden Schritte so oft wiederholen, bis der Fehler verschwindet (oder hinreichend klein ist):
  1. Bestimmung der Schwellwertneuron-**Ausgabe**  $y$  für Trainingspaar  $(x, t)$
  2. Bestimmung des **Fehlers**  $t - y$  der tatsächlichen zur gewünschten Ausgabe vom Trainingsziel  $t$  (als Funktion  $e(w_0, \dots, w_m)$  von den aktuellen Gewichten  $w_0, \dots, w_m$ ),
  3. Bestimmung geeigneter **Gewichtsänderungen**  $\Delta w_i$
  4. Zuordnung der **neuen Gewichte**  $w'_i = w_i + \Delta w_i$  zur Verringerung des (zukünftigen) Fehlers ( $e(w'_0, \dots, w'_n) < e(w_0, \dots, w_n)$ )

# Online-Lernen und Batch-Lernen

Lernen durch schrittweise

1. Berechnung des Fehlers
2. Berechnung der notwendigen Gewichtsänderungen
3. Änderung der Gewichte

Verfahren nach Zeitpunkt der Gewichtsänderung:

**Online-Lernen** Berechnung von Fehler und Gewichtsänderungen für jedes Trainingsmuster, Änderung der Gewichte sofort für jedes Trainingpaar

**Batch-Lernen** (Lernen in Epochen)

Epoche: Berechnung für jedes Paar der Trainingsmenge

Berechnung von Fehler und Gewichtsänderungen für die gesamte Trainingsmenge (z.B. Summe über alle Trainingpaare)

Änderung der Gewichte erst nach einer ganzen Epoche

# Konvergenz des Lernverfahrens

Konvergenzsatz:

Für jede Trainingsmenge

$$\mathcal{T} \subseteq \{(x^{(i)}, t^{(i)}) \mid \forall i \in \{1, \dots, n\} : x^{(i)} \in \{0, 1\}^m \wedge t^{(i)} \in \{0, 1\}\},$$

für welche die Mengen

$$\mathcal{T}_0 = \{x \mid (x, 0) \in \mathcal{T}\} \text{ und } \mathcal{T}_1 = \{x \mid (x, 1) \in \mathcal{T}\}$$

linear trennbar sind,

terminieren sowohl Online- als auch Batch-Lernen eines Schwellwertneurons (passender Struktur) nach endlich vielen Schritten.

Die vom so trainierten Schwellwertneuron berechnete Funktion trennt die Mengen  $\mathcal{T}_0$  und  $\mathcal{T}_1$  voneinander.

# Netze aus Schwellwertneuronen

## Ein-Schicht-Schwellwertneuronen-Netz

parallele Schaltung mehrerer Schwellwertneuronen  
repräsentiert Boolesche Funktionen mit mehreren  
Ausgaben

Beispiel: Parallelschaltung von  $x_1 \wedge x_2$  und  $\neg x_1 \wedge \neg x_2$

## Mehr-Schicht-Schwellwertneuronen-Netz

parallele und sequentielle Schaltung mehrerer  
Schwellwertneuronen

Jede Boolesche Funktion  $f : \{0, 1\}^n \rightarrow \{0, 1\}$  lässt sich durch ein  
Schwellwertneuronen-Netz berechnen.

Schwellwertneuronen-Netz mit zwei Schichten genügt  
(analog DNF, CNF in Aussagenlogik)

# Netze aus Schwellwertneuronen

## Ein-Schicht-Schwellwertneuronen-Netz

parallele Schaltung mehrerer Schwellwertneuronen  
repräsentiert Boolesche Funktionen mit mehreren  
Ausgaben

Beispiel: Parallelschaltung von  $x_1 \wedge x_2$  und  $\neg x_1 \wedge \neg x_2$

## Mehr-Schicht-Schwellwertneuronen-Netz

parallele und sequentielle Schaltung mehrerer  
Schwellwertneuronen

Jede Boolesche Funktion  $f : \{0, 1\}^n \rightarrow \{0, 1\}$  lässt sich durch ein  
Schwellwertneuronen-Netz berechnen.

Schwellwertneuronen-Netz mit zwei Schichten genügt  
(analog DNF, CNF in Aussagenlogik)

# Feed-Forward-Netze (FFN)

- ▶  $V = \bigcup_{k=1}^n V_k$  mit  $\forall i < j \in \{1, \dots, n\} : V_i \cap V_j = \emptyset$   
Zerlegung der Menge der Neuronen in  $n$  disjunkte **Schichten**
- ▶ Menge der Eingangsneuronen:  $V_1$  (je ein Eingang)
- ▶ Menge der Ausgangsneuronen:  $V_n$  (je ein Ausgang)
- ▶ Neuronen aller anderen Schichten heißen versteckte Neuronen
- ▶  $E \subseteq \bigcup_{k=1}^{n-1} V_k \times V_{k+1}$   
nur vorwärtsgerichtete Kanten zwischen benachbarten Schichten
- ▶ Gewichte bilden  $m \times m$ -Matrix (mit  $m = \text{Anzahl aller Neuronen}$ )
- ▶ für FFN besteht die Gewichtsmatrix aus unabhängigen Blöcken  
Blöcke sind die Gewichtsmatrizen zwischen den Schichten

FFN als Berechnungsmodell:

- ▶ parallele Berechnung (in den Neuronen einer Schicht)
- ▶ sequentielle Berechnung (in miteinander verbundenen Neuronen benachbarter Schichten)

# Perzeptron (historisch)

1958 Frank Rosenblatt, Idee: Modell der Netzhaut (Retina)

Aufbau des Perzeptrons:

1. Schicht (Eingabeschicht) : Menge  $S$  von Stimulus-Zellen  
(Verteilung)
2. Schicht (Mittelschicht) : Menge  $A$  von Assoziations-Zellen  
(Vorverarbeitung)
3. Schicht (Perzeptron-Schicht) : Menge  $R$  von Response-Zellen  
Muster-Assoziator aus Schwellwertneuronen  
(eigentliche Verarbeitung)

Verbindungen:

- ▶ zufällig zwischen Neuronen der Eingabeschicht und Neuronen der Mittelschicht  
feste Gewichte (zufällig)
- ▶ von jedem Neuron der Mittelschicht zu jedem Neuron der Ausgabeschicht  
trainierbare Gewichte

Jedes Ausgabeneuron teilt die Eingabemuster in zwei Klassen  
(akzeptierte und nicht-akzeptierte)

# Ein-Schicht-FFN

- ▶ Abstraktion von der Eingabeschicht im historischen Perzeptron-Modell
- ▶ nur Perzeptron-Schicht (Muster-Assoziator)
- ▶ Parallele Berechnung mehrerer künstlicher Neuronen (hier Schwellwertneuronen)

Eingänge:  $(x_1, \dots, x_m) \in \{0, 1\}^m$

Ausgänge:  $(y_1, \dots, y_n) \in \{0, 1\}^n$

Gewichtsmatrix  $W \in \mathbb{R}^{m \times n}$

Gesamtberechnung des Ein-Schicht-FFN  $f : \{0, 1\}^m \rightarrow \{0, 1\}^n$  des Neurons mit gewichteter Summe als Aktivierungsfunktion:

$f(x_1, \dots, x_m) = (y_1, \dots, y_n)$  mit  $\forall k \in \{1, \dots, n\}$ :

$$y_k = \begin{cases} 1 & \text{falls } \sum_{i=1}^m x_i w_{ij} \geq 0 \\ 0 & \text{sonst} \end{cases}$$

(Matrixmultiplikation)

# Ein-Schicht-FFN: Training mit $\Delta$ -Regel

überwachtes Lernen

Trainingsmenge: Menge von Paaren  $(x, t)$  aus

- ▶ Eingabevektoren  $x \in \{0, 1\}^m$  und
- ▶ gewünschten Ausgabevektoren  $t \in \{0, 1\}^n$

Lernen mit Delta-Regel für Ein-Schicht-FFN:

- ▶ Beginn mit zufälligen Eingangsgewichten  $w_{ij} \in \mathbb{R}$ ,
- ▶ für jede Eingabe der Trainingsmenge  $(x, t)$ :
  1. Netz berechnet die Ausgabe  $y = xW$ ,
  2. Zuordnung neuer Gewichte  $w'_{ij}$  durch Delta-Regel:

$$w'_{ij} = w_{ij} + \Delta(w_{ij}) \quad \text{mit} \quad \Delta(w_{ij}) = \eta x_i (t_j - y_j)$$

- ▶ wiederholen, bis der Fehler klein genug ist.

Das Lernverfahren mit Delta-Regel konvergiert für

- ▶ jede linear trennbare Boolesche Funktion  $f$  und
- ▶ hinreichend kleine Lernquote  $\eta$

in endliche vielen Schritten zu einem Ein-Schicht-FFN, welche die Funktion  $f$  berechnet.

# Künstliche Neuronen mit reellen Ein- und Ausgängen

Parameter:

Eingänge:  $x_1, \dots, x_m \in \mathbb{R}^m$

Eingangsgewichte  $w_1, \dots, w_m \in \mathbb{R}^m$

Ausgang:  $f(\langle x, w \rangle) \in \mathbb{R}$

- ▶ Eingangsfunktion  $I : \mathbb{R}^m \rightarrow \mathbb{R}$
- ▶ Aktivierungsfunktion  $A : \mathbb{R} \rightarrow \mathbb{R}$
- ▶ Ausgabefunktion  $O : \mathbb{R} \rightarrow \mathbb{R}$

Gesamtberechnung  $f : \mathbb{R}^m \rightarrow \mathbb{R}$  des Neurons:

$$f(x_1, \dots, x_m) = O(A(I(x_1, \dots, x_m)))$$

# Klassifikation durch Ein-Schicht-FFN

Klassifikation:

Zerlegung einer Menge  $M$  von Werten in (paarweise disjunkte) Klassen  $\{C_1, \dots, C_n\}$ , welche die Wertemenge vollständig überdecken

$$\bigcup_{i=1}^n C_i = M \quad (\forall i \neq j : C_i \cap C_j = \emptyset)$$

Klassifikation des  $\mathbb{R}^m$  durch KNN:

- ▶ Eingänge  $(x_1, \dots, x_m) \in \mathbb{R}^m$
- ▶ Ausgänge  $(y_1, \dots, y_n) \in \{0, 1\}^n$   
für jede Klasse  $C_i$  ein Ausgabeneuron  $y_i$   
Ausgang  $y_i = 1$  gdw. Eingabe  $(x_1, \dots, x_m) \in C_i$

überwachtes Training des Ein-Schicht-FFN:

- ▶ zufällige Startgewichte
- ▶ schrittweise Modifikation der Gewichte zur Verringerung des Fehlers

Ein-Schicht-FFN erkennt nur linear trennbare Klassen

Problem: Wie trainiert man Mehrschicht-FFN?

## Auswahl durch Mehrschicht-FFN – Beispiel

Beispiel: Auswahl aller Punkte im Einheitsquadrat

$$y = \begin{cases} 1 & \text{falls } 0 \leq x_1 \leq 1 \wedge 0 \leq x_2 \leq 1 \\ 0 & \text{sonst} \end{cases}$$

durch das 2-Schicht-FFN mit

- ▶ Eingängen  $x_1, x_2$  und  $x_0$  (bias)
- ▶ Ausgang  $y$
- ▶ versteckten Neuronen  $z_1, \dots, z_4$  und  $z_0$  (bias)
- ▶ Gewichte der ersten Schicht (zwischen  $(x_0, x_1, x_2)$  und  $(z_1, \dots, z_4)$ ):

$$W_1 = \begin{pmatrix} 1 & 0 & 1 & 0 \\ 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & -1 \end{pmatrix}$$

$z_1$  feuert gdw.  $x_1 \leq 1$ ,  $z_2$  feuert gdw.  $x_1 \geq 0$

$z_3$  feuert gdw.  $x_2 \leq 1$ ,  $z_4$  feuert gdw.  $x_2 \geq 0$

- ▶ Gewichte der zweiten Schicht (zwischen  $(z_0, \dots, z_4)$  und  $y$ ):

$$W_2 = (-7/2, 1, 1, 1, 1)^T$$

# Gesamtmatrix des FFN – Beispiel

	$x_0$	$x_1$	$x_2$	$z_0$	$z_1$	$z_2$	$z_3$	$z_4$	$y$
$x_0$	0	0	0	0	1	0	1	0	0
$x_1$	0	0	0	0	1	-1	0	0	0
$x_2$	0	0	0	0	0	0	1	-1	0
$z_0$	0	0	0	0	0	0	0	0	$-7/2$
$z_1$	0	0	0	0	0	0	0	0	1
$z_2$	0	0	0	0	0	0	0	0	1
$z_3$	0	0	0	0	0	0	0	0	1
$z_4$	0	0	0	0	0	0	0	0	1
$y$	0	0	0	0	0	0	0	0	0

## Mehr-Schicht-FFN mit linearer Aktivierung

Netzeingänge:  $(x_1, \dots, x_{k_0}) \in \mathbb{R}^m$

Netzausgänge:  $(y_1, \dots, y_{k_l}) \in \mathbb{R}^n$

Neuronen ( $l$  Schichten):  $(z_1^0, \dots, z_{k_0}^0) \in \mathbb{R}^{k_1}$  (Eingabeneuronen)

$\vdots$  (versteckte Neuronen)

$(z_1^l, \dots, z_{k_l}^l) \in \mathbb{R}^{k_l}$  (Ausgabeneuronen)

Gewichtsmatrizen  $W^{(j)} \in \mathbb{R}^{k_j \times k_{j+1}}$  für jedes  $j \in \{0, \dots, l-1\}$

lineare Aktivierungsfunktion  $I: \mathbb{R} \rightarrow \mathbb{R}$  mit  $I(x) = mx$

Ausgabe des Neurons  $z_i^j$  in Schicht  $j$ :

$$f(z_1^{j-1}, \dots, z_{k_{j-1}}^{j-1}) = O(A(I(x_1, \dots, x_{k_{j-1}}))) = m \left( \sum_{l=1}^{k_{j-1}} w_{li}^{(j)} z_l^{(j-1)} \right)$$

Netzausgabe:

$$f(x_1, \dots, x_m) = m'(x_1, \dots, x_m) W^{(0)} \dots W^{(l-1)} = m'(x_1, \dots, x_m) W$$

mit  $W = W^{(0)} \dots W^{(l-1)}$  (Matrixmultiplikation)

Jede Funktion, die von einem Mehr-Schicht-FFN mit linearer Aktivierung berechnet wird, kann also auch durch ein Ein-Schicht-FFN mit linearer Aktivierung berechnet werden.

# Approximation von Funktionen

gegeben: Menge von Trainingspaaren  $\{(x^{(1)}, t^{(1)}), \dots, (x^{(k)}, t^{(k)})\}$   
 $k$  Stützstellen und Werte an diesen Stützstellen  
(z.B. Messwerte)

Ziel:

Konstruktion eines KNN zur Approximation dieser Funktion durch

- ▶ lineare Funktionen
- ▶ Stufenfunktionen
- ▶ komplexere Funktionen

# Quadratischer Fehler

Approximation einer Menge von Trainingspaaren  
(Funktionswerte an Stützstellen)  
durch Funktion gegebenen Typs (z.B. linear)

- ▶ Trainingsmenge liefert Stützstellen:

$$(x_{k1}, \dots, x_{kn}, t_k)_{k \in \{1, \dots, m\}}$$

- ▶ approximierende Funktion  $f : \mathbb{R}^n \rightarrow \mathbb{R}$
- ▶ Fehler an der Stützstelle  $(x_{k1}, \dots, x_{kn})$ :

$$t_k - f(x_{k1}, \dots, x_{kn})$$

- ▶ quadratischer Fehler an der Stützstelle  $(x_{k1}, \dots, x_{kn})$ :

$$E_k = (t_k - f(x_{k1}, \dots, x_{kn}))^2$$

- ▶ quadratischer Gesamtfehler (Summe über alle Trainingspaare / Stützstellen):

$$E = \sum_{k=1}^m (t_k - f(x_{k1}, \dots, x_{kn}))^2$$

Trainingsziel: Minimierung des quadratischen Fehlers

## Beispiel

Bestimmung der Parameter  $m, n$  einer Geraden  $y = f(x) = mx + n$  aus einer Menge gegebener (ungenauer) Trainingspaare  $(x, t)$ , z.B.:

$$\{(1, 10), (2, 7), (4, 5), (5, 1)\}$$

(ganz einfaches) Ein-Schicht-FFN:

- ▶ ein Eingang  $x_1$ , ein Bias-Neuron  $x_0$
- ▶ ein Ausgangsneuron  $y$
- ▶ Gewichte:  $w_0 = n, w_1 = m$

Funktionen des Ausgabeneurons  $y$ :

- ▶ Eingangsfunktion  $I$ : gewichtete Summe  $nx_0 + mx_1 = mx_1 + n$
- ▶ Aktivierungsfunktion  $A$ : Identität (linear)
- ▶ Ausgangsfunktion  $O$ : Identität

Dieses Netz berechnet die Funktion

$$f(x) = O(A(I(x_1))) = I(x_1) = mx_1 + n$$

Ermittlung der Parameter  $m, n$  durch Training des Netzes ( $\Delta$ -Regel)

# Methode der kleinsten Quadrate

direkte Berechnung mit Hilfe der partiellen Ableitungen nach  $m$  und  $n$

$$E = \sum_{k=1}^l (t_k - f(x_k))^2 = \sum_{k=1}^l (t_k - mx_k - n)^2$$

partielle Ableitungen nach  $m$  und  $n$ :

$$\begin{aligned} \frac{\partial E}{\partial m} &= \sum_{k=1}^l (-2) (t_k - mx_k - n) x_k \\ &= -2 \left( \sum_{k=1}^l t_k x_k - m \sum_{k=1}^l x_k^2 - n \sum_{k=1}^l x_k \right) \end{aligned}$$

$$\begin{aligned} \frac{\partial E}{\partial n} &= \sum_{k=1}^l (-2) (t_k - mx_k - n) \\ &= -2 \left( \sum_{k=1}^l t_k - m \sum_{k=1}^l x_k - nl \right) \end{aligned}$$

## Bestimmung der Parameter

Im Minimum von  $f$  sind alle partiellen Ableitungen 0.  
Das ergibt ein lineares Gleichungssystem für  $m$  und  $n$ :

$$\begin{aligned}\sum_{k=1}^l t_k x_k - m \sum_{k=1}^l x_k^2 - n \sum_{k=1}^l x_k &= 0 \\ \sum_{k=1}^l t_k - m \sum_{k=1}^l x_k - ln &= 0\end{aligned}$$

mit den Lösungen

$$\begin{aligned}n &= \frac{\sum_{k=1}^l t_k - m \sum_{k=1}^l x_k}{l} \\ m &= \frac{l \sum_{k=1}^l t_k x_k - \left(\sum_{k=1}^l t_k\right) \left(\sum_{k=1}^l x_k\right)}{\sum_{k=1}^l x_k^2 - \left(\sum_{k=1}^l x_k\right)^2}\end{aligned}$$

im Beispiel  $m = -2, n = 47/4$

# Berechnung der Gewichts-Verschiebungen

Ziel: Minimierung des Fehlers durch schrittweise Verschiebung des Gewichtsvektors

Methode: Gradientenabstiegsverfahren

Verschiebung des Gewichtsvektors in Richtung des steilsten Abstieges (entgegen dem steilsten Anstieg) der Fehlerfunktion (als Funktion der Gewichte)

steilster Anstieg: Gradient (partielle Ableitungen)

Gradientenabstiegsverfahren führt oft, aber nicht immer zu einem geeigneten (globalen) Minimum der Fehlerkurve, endet mitunter in lokalem Minimum

Voraussetzung: Fehlerfunktion ist **differenzierbar**

zur Anwendung in KNN:  $f_u$  differenzierbar gdw.

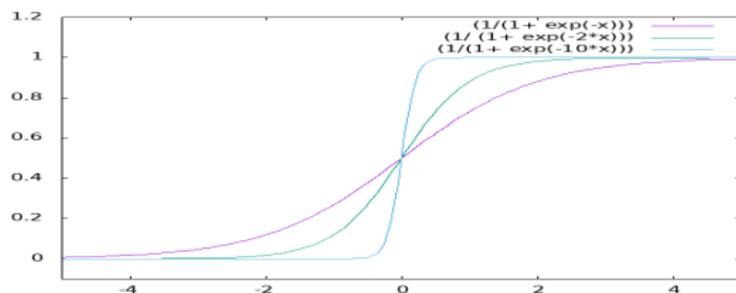
- ▶ Eingangsfunktion  $I_u$  differenzierbar: gewichtete Summe (ok)
- ▶ Ausgangsfunktion  $O_u$  differenzierbar : linear (ok)
- ▶ Aktivierungsfunktion  $A_u$  **differenzierbar**

# Sigmoide Aktivierungsfunktion

differenzierbare Approximation der Stufenfunktion:

sigmoide Funktion

$$f(x) = \frac{1}{1 + e^{-x}} \quad \text{mit Parameter } c > 0: \quad f(x) = \frac{1}{1 + e^{-cx}}$$



- + überall differenzierbar  
Ableitung im Punkt  $x$ :

$$s'(x) = \left( \frac{1}{1 + e^{-x}} \right)' = \frac{e^{-x}}{(1 + e^{-x})^2} = \frac{1}{1 + e^{-x}} \left( 1 - \frac{1}{1 + e^{-x}} \right) = s(x)(1 - s(x))$$

in jedem Punkt eindeutige Abstiegsrichtung

- erreicht die Werte 0 und 1 nie,  
Toleranzbereiche notwendig, so entstehen Ungenauigkeiten

# Lineare Aktivierungsfunktionen und ReLU

lineare Aktivierung:  $\forall x \in \mathbb{R} : A(x) = mx + n$

- + einfach (schnell) zu berechnen
  - überall differenzierbar
  - Ableitung: konstant  $m$
  - in jedem Punkt  $x > 0$  eindeutige Abstiegsrichtung
- Gesamtfunktion  $f_u$  bleibt linear

ReLU(Rectified Linear Units):  $\forall x \in \mathbb{R} : A(x) = \max(0, x)$

- + einfach (schnell) zu berechnen
  - fast überall differenzierbar
  - Ableitung: Stufenfunktion, 0 bei  $x < 0$ , 1 bei  $x > 0$ ,
  - in jedem Punkt  $x > 0$  eindeutige Abstiegsrichtung
- Problem: Ableitung nicht definiert bei  $x = 0$   
(aber: praktisch nicht relevant)

# Backpropagation in FFN

(Bryson, Ho 1969, Rummelhard, McClelland 1986)

Ziel: Geeignete Modifikation aller Gewichte im FFN zur Verringerung des Gesamtfehlers

Idee:

- ▶ Betrachte jedes Gewicht  $w_{uv}$  als Eingangsgewicht des Teilnetzes zwischen Neuron  $v$  und Netz-Ausgängen
- ▶ Netzeingabe in dieses Teilnetz ist  $w_{uv}o_u$  mit Netzausgabe  $o_u$  des Neurons  $u$
- ▶ partielle Ableitung  $\frac{\partial E}{\partial w_{uv}} = o_u \delta_v$  mit Fehleranteil  $\delta_v = o_v(1 - o_v) \sum_p w_{vp} \delta_p$ , wobei  $p$  über alle direkten Nachfolger von  $v$  läuft

# Backpropagation-Training

in jedem Schritt 2 Durchläufe des FFN:

**Vorwärts-Schritt:** Berechnung der Netzausgabe

Speichern der Netzausgabe  $o_u$  in jedem Neuron  $u$

Speichern der Ableitung der Netzausgabe  $o_u(1 - o_u)$   
in jedem Neuron  $u$

**Rückwärts-Schritt:** Berechnung des Fehleranteils  $\delta_u$  jedes Neurons  
aus den Fehleranteilen aller Nachfolger-Neuronen

$$\delta_u = o_u(1 - o_u) \sum_p w_{vu} \delta_p,$$

Speichern der Fehleranteile  $\delta_u$  in jedem Neuron  $u$

danach Anpassung aller Gewichte um  $\Delta w_{uv} = -\eta o_u \delta_v$

# Backpropagation-Lernen allgemein

- ▶ Instanziierung aller Gewichte mit kleinen zufälligen Werten
- ▶ BP-Verfahren für eine Epoche:
  - ▶ BP-Verfahren für jedes Trainingsmuster  $(x, t)$ :
    - ▶ Vorwärtsschritt (Ausgabe-Berechnung):  
für jede Schicht  $s$  (Beginn bei Eingabeschicht):  
Berechnung der Vektoren  $z^{(s)} = I(y^{(s-1)})$  und  
 $y^{(s)} = A(z^{(s)}) = A(I(y^{(s-1)}))$  für jedes Neuron der Schicht  $s$
    - ▶ Rückwärtsschritt (Gewichtsdifferenzen):  
für die Ausgabeschicht  $k$ :  
Berechnung des Vektors  $d^{(k)} = (t - y^{(k)})y^{(k)}(1 - y^{(k)})$   
für jede Schicht  $s$  (Beginn bei letzter versteckter Schicht  
 $k - 1$ ):  
Berechnung des Vektors  $d_j^{(s)} = y_j^s(1 - y_j^{(s)}) \sum_{m=1}^{n^{(s+1)}} d_m^{(s+1)} w_{mj}$   
für jedes Neuron  $j$  der Schicht  $s$
    - ▶ Aktualisierung aller Gewichte:  $w_{ij}^{(s)} := w_{ij}^{(s)} + \eta d_j^{(s)} y_i^{(s)}$   
danach weiter mit nächstem Trainingsmuster  $(x', t')$
  - danach weiter mit nächster Epoche
- ▶ Ende, falls erreichte Änderung des Fehlers klein (unter einer Schranke)

# Backpropagation-Lernen mit Trägheit

zur Vermeidung von

- ▶ Oszillationen in „Schluchten“ und
- ▶ Abbremsen auf Plateaus

$$w_{ij}^{(s)} := (1 + \alpha)w_{ij}^{(s)} + \eta d_j^{(s)} y_i^{(s)}$$

mit Trägheit  $\alpha$

# Anwendung von FFN mit Backpropagation

KNN zur Muster-Klassifikation

Klassifikation von Eingabemustern, z.B.

- ▶ optische Zeichenerkennung  
(z.B. Buchstaben, abstrahiert von Schriftart)
- ▶ Erkennung akustischer Signale (z.B. Stimmen)
- ▶ englische Ausspracheregeln (NETTALK)
- ▶ Datenkompression (Eingabe = Ausgabe, Code in der versteckten Schicht)
- ▶ Vertrauenswürdigkeit von Bankkunden (Risikoklassen)
- ▶ Vorhersage (Wetter, Aktienkurse)
- ▶ bisher: Boolesche Funktionen  
(Klassifikation von Eingabevektoren nach Ausgabe-Wahrheitswerten)

# Qualität von BP-Netzen

gute Generalisierung:

KNN klassifiziert die meisten neuen Eingabemuster einer Testdatenmenge (nicht aus der Trainingsmenge) richtig  
abstrahiert von kleinen Abweichungen  
abhängig von

- ▶ Netzarchitektur (nicht zu viele versteckte Neuronen)
- ▶ Auswahl der Trainingsmenge

Problem:

übertrainierte Netze kennen die Trainingsmenge „auswendig“  
(Overfitting)

# Prominente Aktivierungsfunktionen

- ▶ Stufenfunktion  $A(x) = (x \geq 0)$ ,
- ▶ sigmoide Funktion  $A(x) = \frac{1}{1+e^{-x}}$ ,  $A'(x) = \tanh x$   
(differenzierbare Approximation der Stufenfunktion)
- ▶ lineare Funktion  $A(x) = ax + b$
- ▶ ReLU  $A(x) = \max(0, x)$ ,
- ▶ analytische Funktion  $A(x) = \log(1 + e^x)$ ,  
(differenzierbare ReLU-Approximation, softplus)
- ▶ Radiale Basisfunktion (mehrere Eingaben)  $A(x) = r(d(x, w))$  mit  
Eingabe(-vektor)  $x$ , Gewicht(-svektor)  $w$  (Zentrum), Abstand  $d$ ,  
radiale Funktion  $r : \mathbb{R}_{\geq 0} \rightarrow [0, 1]$  mit
  - ▶  $r(0) = 1$
  - ▶ monoton fallend:  $\forall x, y \in \mathbb{R}_{\geq 0} : (x < y) \rightarrow (r(x) < r(y))$z.B.  $d$  Manhattan-Metrik,  $r(x) = \max(0, 1 - x)$
- ▶ Softmax (mehrere Ein- und Ausgaben)  $A(x_i) = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$   
zur Klassifizierung in mehrere Klassen  
(Wahrscheinlichkeitsverteilung)

# Hybride Netze

Idee:

- ▶ Netze Kombination verschiedenartiger Neuronen
- ▶ schichtweise verschiedene Aktivierungsfunktionen
- ▶ Schichten enthalten (derzeit) meist Neuronen derselben Art
- ▶ Trennung von Neuronen in Folge mehrerer  
eigene Neuronen für Eingangs-, Aktivierungs-,  
Ausgangsfunktion
- ▶ verschiedene Lernverfahren je Schicht

## Beobachtungen im visuellen System:

- ▶ sendet **vorverarbeitete** Signale an Gehirn
- ▶ **heterogenes** Netz: verschiedene Neurone haben verschiedene Wirkungen (Funktionen)
- ▶ Neuronen derselben Schicht haben dieselbe Funktion
- ▶ Verbindung benachbarter Neuronen  
horizontale Zellen berechnen Mittelwert (der Helligkeit)  
wirken hemmend auf Signale nahe beim Mittelwert
- ▶ ähnlich **Faltung** in digitaler Bildverarbeitung:  
Funktionswert eines Pixels hängt von Werten benachbarter Pixel ab

## Faltung (1D)

stetig:

Für  $f, g : \mathbb{R} \rightarrow \mathbb{R}$  ist die Faltung von  $f$  und  $g$  definiert durch

$$\forall p \in \mathbb{R} : (f * g)(p) = \int_{-\infty}^{\infty} f(x)g(p-x)dx$$

diskret:

Für  $f, g : \mathbb{Z} \rightarrow \mathbb{Z}$  ist die Faltung (Convolution) von  $f$  und  $g$  definiert durch

$$\forall p \in \mathbb{Z} : (f * g)(p) = \sum_{x=-\infty}^{\infty} f(x)g(p-x)$$

diskret mit endlichem Bereich  $\text{pos} = [m, \dots, n] \subseteq \mathbb{Z}$  für  $f$ :

Für  $f : \text{pos} \rightarrow \mathbb{Z}, g : \mathbb{Z} \rightarrow \mathbb{Z}$  ist die **Faltung** (Convolution) von  $f$  und  $g$  definiert durch

$$\forall p \in \mathbb{Z} : (f * g)(p) = \sum_{x \in \text{pos}} f(x)g(p-x)$$

## Diskrete Faltung (1D) – Beispiel

$f : \mathbb{Z} \rightarrow \mathbb{Z}, g : \{1, 2\} \rightarrow \mathbb{Z}$  mit

$$f(x) = \begin{cases} 3 & \text{falls } x = 1 \\ 1 & \text{falls } x = 2 \end{cases} \quad g(x) = \begin{cases} 2 & \text{falls } x = 1 \\ 1 & \text{falls } x = 2 \\ 0 & \text{sonst} \end{cases}$$

häufig: Darstellung von  $f$  (Faltungskern)

als Vektor mit endlichem Definitionsbereich  $f = (f(1), f(2)) = (3, 1)$

$$(f * g)(p) = \begin{cases} 6 & \text{falls } p = 2 \\ 5 & \text{falls } p = 3 \\ 1 & \text{falls } p = 4 \\ 0 & \text{sonst} \end{cases}$$

Berechnung (des relevanten Teils):

$$(f * g)(1) = f(1)g(0) + f(2)g(-1) = 3 \cdot 0 + 1 \cdot 0 = 0$$

$$(f * g)(2) = f(1)g(1) + f(2)g(0) = 3 \cdot 2 + 1 \cdot 0 = 6$$

$$(f * g)(3) = f(1)g(2) + f(2)g(1) = 3 \cdot 1 + 1 \cdot 2 = 5$$

$$(f * g)(4) = f(1)g(3) + f(2)g(2) = 3 \cdot 0 + 1 \cdot 1 = 1$$

$$(f * g)(5) = f(1)g(4) + f(2)g(3) = 3 \cdot 0 + 1 \cdot 0 = 0$$

## Faltung (2D)

Für  $f, g : \mathbb{R}^2 \rightarrow \mathbb{R}$  ist die **Faltung** (Convolution) von  $f$  und  $g$  definiert durch

$$\forall (p_x, p_y) \in \mathbb{R}^2 : (f * g)(p_x, p_y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) g(p_x - x, p_y - y) dx dy$$

diskreter Fall (Bilder):

Für  $f, g : \{0, \dots, m-1\} \times \{0, \dots, m-1\} \rightarrow \mathbb{R}$  ist die **Faltung** (Convolution) von  $f$  und  $g$  definiert durch

$\forall (p_x, p_y) \in \{0, \dots, m-1\} \times \{0, \dots, m-1\} :$

$$(f * g)(p_x, p_y) = \sum_{x=0}^{m-1} \sum_{y=0}^{n-1} f(x, y) g(p_x - x, p_y - y)$$

Erweiterung auf höhere Dimensionen analog

# Multiskalen-Strategien

(alte) Idee aus der digitalen Bildverarbeitung / -erkennung

Ziel: Bildanalyse in verschiedenen Größen (Auflösungsstufen)

- ▶ wenig aufwendige Untersuchung grober Strukturen
- ▶ feinere Untersuchung feiner Strukturen im Bild

Idee: Multiskalenraum (ähnlich menschlicher Wahrnehmung)

enthält Originalbild  $B_0$  und weniger detaillierte Versionen  $B_k$

Erzeugung der Multiskalen-Bilder (Pyramiden) ( $B_0, B_1, B_2, \dots$ )

durch wiederholte Ausführung der Schrittfolge

1. Glättung (durch Tiefpass-Filter, z.B. Gauß-Filter)
2. **reduce**: Komprimierung durch geringere Abtastrate, z.B. Gauß-Pyramide: Löschen jeder zweiten Zeile und Spalte
3. **expand**: Umkehrung durch Interpolation (nicht verlustfrei) zur Vergleichbarkeit der Bilder verschiedener Auflösungen

Laplace-Pyramide (redundanzarme Darstellung):

Schichten enthalten Differenz

zwischen Bild  $B_k$  und  $\text{expand}(\text{reduce } B_k)$

# Bild-Pyramiden durch Faltungen

zur Erkennung von Features:

- ▶ Flächen gleicher Farbe
- ▶ Kanten
- ▶ Formen
- ▶ Texturen, ...

Bilder enthalten Informationen auf verschiedenen Ebenen  
(Auflösungsstufen),

kleinteilige Beobachtung lenkt evtl. von wesentlichen Merkmalen ab

Umsetzung durch Multiskalen-Bilder (Pyramiden)

entstehen durch mehrfache Wiederholung von

- ▶ Glättung (durch geeignete Filter, oft **Faltung**)
- ▶ **Komprimierung** durch geringere Abtastrate,  
z.B. Gauß-Pyramide: Löschen jeder zweiten Zeile und Spalte

CNN-Idee:

Umsetzung durch Aufeinanderfolgen verschiedener Schichten in  
(feed-forward)-KNN

# Neocognitron

Fukushima (1975):

Cognitron: A Self-Organizing Multilayered Neural Network Model

(1983) Neocognitron: A Neural Network Model for a Mechanism of Visual Pattern Recognition

Motivation: Erkennung handschriftlicher Ziffern

Aufbau Neocognitron:

- ▶ Eingabe-Schicht
  - ▶ vier (oder mehr) versteckte Stufen aus je zwei Schichten:
    1. Transformation in 12 Bilder (Ebenen)  
Feature-Extraktion (Faltungen mit je einem  $3 \times 3$ -Kern)  
Filterkerne durch Eingangsgewichte definiert (weight sharing)  
Gewichte durch Trainingsmuster gelernt
    2. Kombination mehrerer transformierter Bilder  
z.B. punktweise gewichtete Summe, Max  
Gewichte nicht trainiert
  - ▶ Ausgabe nach letzter Kombinations-Schicht  
(Klassifikation)
  - ▶ inkrementelles Lernen stufenweise von Ein- zu Ausgabeschicht
- mehrere Varianten mit überwachtem und unüberwachtem Lernen

# KNN zur Klassifikation

z.B. von handgeschriebenen Ziffern

- ▶ Linearer Klassifikator (keine versteckte Schicht)
- ▶ Ein-Schicht-FFN (eine versteckte Schicht)
- ▶ Mehr-Schicht-FFN (mehrere versteckte Schichten)
- ▶ CNN LeNet-1 (1989)
- ▶ CNN LeNet-5
- ▶ ...

# Convolutional Neural Networks

z.B. Alex Krizhevsky, . . . , 2012:

ImageNet Classification with Deep Convolutional Neural Networks

prinzipieller Aufbau:

- ▶ Eingabe-Schicht
- ▶ Versteckte Stufen aus je mehreren Schichten
  - ▶ Faltungs-Schicht (Feature-Maps)  
alle Gewichte gleich
  - ▶ evtl. ReLU-Schicht (nichtlinear)
  - ▶ gelegentlich Subsampling-Schicht (Pooling)
- mehrfache Wiederholung (deep), evtl. in verschiedenen Reihenfolgen
- ▶ evtl. klassische Schichten mit vollständigen Verbindungen zwischen benachbarten Schichten
- ▶ Ausgabe-Schicht

nach und nach Entwicklung komplexerer Konstruktionen, z.B.

- ▶ AlexNet (Dropout-Schichten)
- ▶ GoogLeNet (Inception)
- ▶ ResNet (skip connections)

# CNN-Schichten

aktuelle CNNs bestehen im Wesentlichen aus mehrfacher Wiederholung folgender Schichten:

- ▶ **Faltung** (convolutional)
- ▶ **Komprimierung** / Auswahl  
(pooling, meist Durchschnitt oder Max)
- ▶ vollständig verbunden (fully connected, FC)
- ▶ Normalisierung (batch normalization, BN)
- ▶ Softmax (Wahrscheinlichkeitsverteilung)

CNNs unterscheiden sich in

- ▶ Reihenfolge der Schichten
- ▶ Anzahl der Schichten / Wiederholungen
- ▶ spezielle topologische Merkmale,  
z.B. Vorwärtskanten, die Schichten überspringen

## CNN – prominente Beispiele

- ▶ VGG16, ImageNet-Datenbank (seit 2010):
  - ▶  $< 15 \cdot 10^6$  annotierte Bilder
  - ▶  $< 10 \cdot 10^3$  Klassen

ImageNet Large-Scale Visual Recognition Challenge (ILSVRC, seit 2010), Gewinner 2014: VGG16

- ▶ Nachcolorierung: Eingabe Grauwertbild, Ausgabe Farbbild  
<https://richzhang.github.io/colorization> (2016)
- ▶ SegNet (2015), Aufgabe: semantische Segmentierung (Jeder Bildposition wird Bedeutung zugeordnet)  
Eingabe und Ausgabe: Bild derselben Größe,  
Farben des Ausgabebildes sind Klassen  
Idee: Verknüpfung VGG16 mit „gespiegelter“ Version

# CNN-Lernen

Überwachtes Lernen durch Backpropagation  
(angepasste Verfahren für jeden Schicht-/ Neuronentyp)

Trainieren von CNN ist i.A. sehr aufwendig (Zeit, Daten)

deshalb häufig auch Nutzung vortrainierter Netze

Idee:

- ▶ Ausgangspunkt: Netz, welches schon auf allgemeine Daten trainiert wurde  
z.B. Klassifikation, Segmentierung
- ▶ Training auf spezielle Aufgabe anhand spezieller Datensätze  
z.B. andere bzw. feinere Klassen

# Was bisher geschah

## Maschinelles Lernen

- ▶ überwacht: korrigierend / bestärkend
- ▶ unüberwacht
- ▶ Anwendungen

## Künstliche Neuronen

- ▶ Biologisches und mathematisches Neuronenmodell
- ▶ Eingang- / Aktivierungs- / Ausgangsfunktion mit Beispielen
- ▶ McCulloch-Pitts-Neuron ohne / mit Hemmung, Schwellwertelement, allgemeines Neuronenmodell

## Künstliche Neuronale Netze: Feed-Forward-Netze

- ▶ Ein-Schicht-FFN (ohne versteckte Neuronen)  
Lernen mit  $\Delta$ -Regel
- ▶ Mehr-Schicht-FFN / DNN (mit versteckten Schichten)  
Lernen mit Backpropagation (Gradientenabstieg)
- ▶ Spezialfall CNN: hybride Netze aus Schichten von Neuronen verschiedener Typen, u.A. Faltungsschichten

## Rekurrente Netze: Motivation

Ziel: Nachnutzung von Informationen aus vorangegangenen Schritten (Gedächtnis), z.B. zur

- ▶ Repräsentation zeitlicher Folgen von Mustern
- ▶ Zeitreihenanalyse und -voraussage
- ▶ Erkennung von Sätzen (Grammatik)
- ▶ Verarbeitung von Mustern variabler Längen (betrachtet als Sequenzen)

mögliche Ansätze:

- ▶ gleitendes Zeitfenster:  
FFN mit  $n$  Eingabeneuronen  
Eingabemuster enthält Informationen aus  $n$  vorangegangenen Schritten  
Nachteil: beschränkte Breite des Zeitfensters
  - ▶ Erkennen „entfernter“ Abhängigkeiten schwierig
  - ▶ viele Eingabeneuronen nötig
- ▶ rekurrente KNN

# Allgemeine KNN

Netzstruktur (Topologie):

gerichteter Graph  $G = (V, E)$  mit

- ▶ endliche Menge  $V = \{v_1, \dots, v_n\}$  von Knoten (Neuronen)  
evtl. einige als Eingabe- bzw. Ausgabeneuronen  
gekennzeichnet (nicht notwendig)
- ▶ Menge  $E \subseteq V \times V$  von (gewichteten) Kanten

eine Gewichtsmatrix  $\mathbb{R}^{V \times V}$  für alle möglichen Verbindungen  
zwischen Neuronen

# Beispiel

- ▶ zwei McCulloch-Pitts-Neuronen  $u, v$
- ▶ Eingang  $x \in \{0, 1\}$
- ▶ Ausgang  $y \in \{0, 1\}$
- ▶ erregende Kanten:  $(x, u), (x, v), (u, u), (u, v), (v, y)$
- ▶ hemmende Kanten  $(v, v), (v, u)$
- ▶ Schwellwerte  $\theta_u = 1, \theta_v = 2$

# Zustand rekurrenter Netze

Zustand eines neuronalen Netzes (zeitveränderlich)

Aktivierung aller Neuronen:

Zuordnung  $S$ : Neuron  $\rightarrow \mathbb{R}$

Übersetzung in Zustandsübergangssysteme

(endliche Automaten mit Ausgabe in Knoten, Moore)

Satz: Zu jedem NFA existiert ein rekurrentes Netz aus McCulloch-Pitts-Neuronen, welches dieselben Zustandsübergänge simuliert.

# Mathematisches Modell: Rekursion

Wiederholung: KNN als Berechnungsmodell

FFN als Berechnungsmodell:

- ▶ parallele Berechnung (in den Neuronen einer Schicht)
- ▶ sequentielle Berechnung (in miteinander verbundenen Neuronen benachbarter Schichten)  
Nacheinanderausführung von Funktionen

rekurrentes Netz als Berechnungsmodell:

- ▶ mehrmalige Nacheinanderausführung einer Funktion (ohne Abbruchbedingung)  
Berechnung einer rekursiven Funktion (Fixpunkt)

## „Entwerrung“ rekurrenter Netze

Idee:

- ▶ Verarbeitung von Eingaben zu Ausgaben eines Neurons kostet einen Zeitschritt
- ▶ für jeden Zeitschritt eine Kopie aller Neuronen und Kanten dazwischen,
- ▶ Ersetzung der Rückwärtskanten durch Vorwärtskanten zur nächsten Kopie.

In diesem expandierten Netz ist Lernen der Vorwärtskanten durch Backpropagation-Verfahren möglich:

- ▶ Durchlauf jeder Netz-Kopie ist ein Zeitschritt,
- ▶ Lernen durch Backpropagation des entwirrten KNN (Backpropagation through time)

# Jordan-Netze

Idee (1986): Nachnutzung der **Netzausgaben**

Netz-Topologie:

- ▶ Feed-Forward-Netz mit trainierbaren Vorwärtskanten,
- ▶ für jedes **Ausgabeneuron** ein zusätzliches **Kontextneuron** in der Eingabeschicht  
(zur Speicherung der Netzausgaben)  
Aktivierungsfunktion: Identität
- ▶ zusätzliche Verbindungen von jedem Neuron der Ausgabeschicht zu seinem Kontextneuron mit festen Gewichten  $\lambda$  (meist  $\lambda = 1$ ),  
Speicherung der Ausgaben
- ▶ evtl. direkte Verbindungen von jedem Kontextneuron zu sich selbst mit festem Gewicht  $\gamma$   
(zur weiteren Speicherung der Netzausgaben)
- ▶ zusätzliche Verbindungen von jedem Kontextneuron zu jedem Neuron der ersten versteckten Schicht mit trainierbaren Gewichten,  
(zur Verwendung der gespeicherten Ausgabe im Folgeschritt)

# Elman-Netze

Idee (1990): Nachnutzung der Aktivierung der **versteckten Neuronen**

Netz-Topologie:

- ▶ Feed-Forward-Netz (z.B. SRN 3-Schicht-FFN)
- ▶ für jedes **versteckte Neuron** ein zusätzliches **Kontextneuron** in der vorigen Schicht  
(zur Speicherung der Aktivierung)  
Aktivierungsfunktion: Identität
- ▶ zusätzliche Verbindungen von jedem versteckten Neuron zu seinem Kontextneuron mit festem Gewicht 1  
Speicherung der Aktivierung aller versteckten Neuronen
- ▶ zusätzliche Verbindungen von jedem Kontextneuron zu jedem Neuron der Schicht des Originalneurons mit trainierbaren Gewichten,  
(zur Verwendung der gespeicherten Aktivierung im Folgeschritt)

# Long Short-Term Memory (LSTM)

Idee (1997): Nachnutzung der Aktivierung der **versteckten Neuronen**

- ▶ Knoten: LSTM unit (memory block + gate units) besteht aus
  - ▶ Kern: CEC (constant error carousel)  
lineare Aktivierung und Rückkopplung zu sich selbst  
speichert Zustand, Rückkopplung im nächsten Schritt
  - ▶ zusätzliche gate units vor und nach CEC steuern die Fähigkeit, Eingaben zu verarbeiten und speichern (input gates) bzw. Ausgaben weiterzugeben (output gates)
- ▶ Netz-Topologie: 3-Schicht-FFN, LSTM units in versteckter Schicht

prominente Modifikationen:

- ▶ forget gate steuert (multipliziert, i.A. schwächt) Gewicht der rückkoppelnden Kante
- ▶ Peephole Connections (gewichtete Kanten von CEC zu allen gates) bei schwacher Aktivierung: Isolation des CEC

Anwendungen:

- ▶ Zeitreihenanalyse, -vorhersage
- ▶ Handschrifterkennung, Spracherkennung, Übersetzung

# Assoziativspeicher

Ziel: Musterassoziation

(Training mit endlich vielen Musterpaaren)

Generalisierung:

- ▶ Aus Zuordnung: Muster  $x \rightarrow$  Muster  $y$  folgt für jedes zu  $x$  **ähnliche** Muster  $x'$  die Zuordnung: Muster  $x' \rightarrow$  Muster  $y$ .
- ▶ Ziel: sinnvolle Zuordnung „verrauschter“ oder unvollständiger Eingabemuster

Netztypen:

**heteroassoziativ** Eingabemuster  $x \in \mathbb{R}^m$ ,  
Ausgabemuster  $y \in \mathbb{R}^n$

**Mustererkennung** Spezialfall heteroassoziativer Netze  
assoziiert Muster mit Identifikator, z.B. für Klasse

**autoassoziativ** Spezialfall heteroassoziativer Netze  
Ein- und Ausgabemuster  $x \in \mathbb{R}^m$  (prinzipiell) gleich

# Heteroassoziativer Speicher

**Topologie:** vollständig verbundenes Ein-Schicht-FFN,  
Eingänge  $x \in \mathbb{R}^m$ , Ausgänge  $y \in \mathbb{R}^n$ ,  
alles Schwellwertneuronen, Gewichtsmatrix  $W \in \mathbb{R}^{m \times n}$

**Aktivierung:** Signum = Stufenfunktion

$$A(x) = \text{sgn}(x) \begin{cases} -1 & \text{falls } x < 0 \\ 0 & \text{falls } x = 0 \\ 1 & \text{sonst} \end{cases}$$

**Berechnung der Gewichte:** Methode der kleinsten Quadrate =  
Minimierung von

Berechnung der Eingangsfunktion analog Ein-Schicht-FFN:

$$I(x_1, \dots, x_m) = (x_1, \dots, x_m) \begin{pmatrix} w_{11} & \dots & w_{1n} \\ \vdots & \ddots & \vdots \\ w_{m1} & \dots & w_{mn} \end{pmatrix}$$

dasselbe kürzer:  $I(x) = xW$

# Heteroassoziativer Speicher: Beispiel

Berechnung für 3-2-Netz mit Gewichtsmatrix  $W$ :

$$W = \begin{pmatrix} -1 & 1 \\ 1 & -1 \\ -1 & 1 \end{pmatrix}$$

Eingabe  $x = (1, -1, 1)$

Berechnung:

$$\text{sgn}(xW) = \text{sgn} \left( (1, -1, 1) \begin{pmatrix} -1 & 1 \\ 1 & -1 \\ -1 & 1 \end{pmatrix} \right) = \text{sgn}(-3, 3) \mapsto (-1, 1)$$

Ausgabe  $y = (-1, 1)$

# Heteroassoziativer Speicher: Training

Idee: Lernen aus gleichzeitiger Aktivität (Hebb)

biologisches Vorbild: Synapsen zwischen gleichzeitig aktiven Neuronen ( $x_i$  und  $y$ ) werden verstärkt (synaptische Plastizität)

Trainingsmenge  $\{(x^{(i)}, y^{(i)}) \mid i \in \{1, \dots, k\}\}$  (Bipolarvektoren)

Ziel Gewichtsmatrix  $W$ , so dass für jedes Trainingspaar  $(x^{(i)}, y^{(i)})$  gilt: mit  $\text{sgn}(x^{(i)} W) = y^{(i)}$  (komponentenweise)

Startgewichte alle  $w_{ij} = 0$

Lernregel von Hebb  $\Delta w_{ij} = \eta x_i y_j$ , hier mit Lernrate  $\eta = 1$

Training : Gewichtsbestimmung  $\Delta w_{kl} = x_k y_l$   
je Trainingspaar  $(x, y)$  einmal  
( $W$  ist Korrelationsmatrix von  $x$  und  $y$ )

Alternative zum Training: direkte Berechnung der Gewichte (z.B. mit Methode der kleinsten Quadrate)

# Bidirektionaler Assoziativspeicher (BAM)

(heteroassoziativer Speicher von Musterpaaren)

Netz-Topologie:

- ▶ Eingangsschicht, Ausgangsschicht, keine versteckten Neuronen
- ▶ Eingänge  $x \in \{-1, 1\}^m$
- ▶ Ausgänge  $y \in \{-1, 1\}^n$
- ▶ vollständige **symmetrische** Verbindungen zwischen jedem Eingangs- und jedem Ausgangsneuron  
(ungerichteter vollständiger bipartiter Graph  $K_{m,n}$ )
- ▶ Gewichte an jeder Kante (für beide Richtungen gleich)  
Gewichte in Gewichtsmatrix  $W \in R^{m \times n}$
- ▶ Eingangsfunktion: gewichtete Summe
- ▶ Aktivierung: Signum
- ▶ Ausgangsfunktion: Identität

# BAM: Berechnung und Training

(wie im heteroassoziativen Speicher)

überwachtes Lernen

Trainingsmenge  $\{(x^{(i)}, y^{(i)}) \mid i \in \{1, \dots, k\}\}$

**Eingabe** : Startzustand (initiale Zustände der Eingangsneuronen  $x \in \{-1, 1\}^m$ )

**Lernregel** von Hebb:  $\Delta w_{ij} = \eta x_i y_j$  mit  $\eta = 1$

**Berechnung** : Folge von Schritten (Zustandsübergängen),

- ▶ synchron oder
- ▶ abwechselnd

für beide Neuronenschichten:

Aktualisierung: Neuberechnung der Aktivierung der anderen Schicht

wiederholen, bis stabiler Zustand erreicht (Fixpunkt) oder Abbruch ausgelöst wird

**Ausgabe** : Zustand der Ausgangsneuronen des stabilen Netzes

## BAM: Beispiel

ein Trainingspaar  $(x, y)$  mit  
 $x = (1, -1, 1)$  und  $y = (-1, 1)$

Gewichtsmatrix  $W$ :

$$W = \begin{pmatrix} -1 & 1 \\ 1 & -1 \\ -1 & 1 \end{pmatrix}$$

Berechnung:

$$\text{sgn}(xW) = \text{sgn} \left( (1, -1, 1) \begin{pmatrix} -1 & 1 \\ 1 & -1 \\ -1 & 1 \end{pmatrix} \right) = \text{sgn}(-3, 3) \mapsto (-1, 1) = y$$

$$\text{sgn}(Wy^T) = \text{sgn} \left( \begin{pmatrix} -1 & 1 \\ 1 & -1 \\ -1 & 1 \end{pmatrix} \begin{pmatrix} -1 \\ 1 \end{pmatrix} \right) = \text{sgn} \begin{pmatrix} 2 \\ -2 \\ 2 \end{pmatrix} \mapsto \begin{pmatrix} 1 \\ -1 \\ 1 \end{pmatrix} = x^T$$

## BAM: Training

(wie heteroassoziativer Speicher)

- ▶ für ein zu speicherndes Musterpaar  $x \in \{-1, 1\}^m, y \in \{-1, 1\}^n$ :  $W = x^T y$   
(Korrelationsmatrix von  $x$  und  $y$ )
- ▶ für mehrere zu speichernde Musterpaare  $(x^{(1)}, y^{(1)}), \dots, (x^{(k)}, y^{(k)})$ :

$$W = \sum_{i=1}^k W^{(i)} = \sum_{i=1}^k \left(x^{(i)}\right)^T y^{(i)}$$

Für alle  $k$  Trainingsmuster gleichzeitig:

- ▶ Eingabemuster  $X \in \{-1, 1\}^{k \times m}$
- ▶ Ausgabemuster  $Y \in \{-1, 1\}^{k \times n}$
- ▶  $XX^T \in \{-1, 1\}^{k \times k}$  (alle Einträge positiv)
- ▶  $\text{sgn}(Y) = \text{sgn}(XX^T Y) \in \{-1, 1\}^{k \times n}$
- ▶  $\text{sgn}(Y) = \text{sgn}(XW) \in \{-1, 1\}^{k \times n}$  mit  $W = X^T Y$

# Hopfield-Netz

(autoassoziativer Musterspeicher)

Idee: autoassoziativer BAM (gleiche Anzahl  $n$  von Ein- und Ausgaben)  
jedes Paar von Ein- und Ausgabeknoten identifiziert

**Topologie:**  $K_n$  mit symmetrischer Gewichtsmatrix  $W \in \mathbb{R}^n$   
Jedes Neuron ist zugleich Ein- und Ausgang  $x \in \{-1, 1\}^n$   
keine Selbstrückkopplung, also  $\forall i \in \{1, \dots, n\} : w_{ii} = 0$

**Zustand:** Aktivierung aller Neuronen

**Eingabe:** Startzustand (initiale Zustände aller Neuronen)

**Berechnung:** Folge von Schritten (Zustandsübergängen):  
Aktualisierung: Berechnung der Aktivierung aller Neuronen  
wiederholen, bis stabiler Zustand erreicht oder Abbruch

**Konvergenz** : Erreichen eines stabilen Zustandes (ändert sich bei  
Aktualisierung beliebiger Neuronen nicht)

**Ausgabe** : Zustand des stabilen Netzes

Aktualisierung in jedem Schritt:

**synchron:** gleichzeitige Zustandsänderung für alle Neuronen

**asynchron:** Zustandsänderung eines zufällige gewählten Neurons  
(faire Auswahl)

## Hopfield-Netz – Beispiele

$$W = \begin{pmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix} \quad x = (1, -1, 1)$$

$$W = \begin{pmatrix} 0 & -1 & 1 & 1 \\ -1 & 0 & -1 & -1 \\ 1 & -1 & 0 & 1 \\ 1 & -1 & 1 & 0 \end{pmatrix} \quad x = (-1, -1, 1, 1)$$

$$W = \begin{pmatrix} -1 & -1 & -1 \\ -1 & -1 & -1 \\ -1 & -1 & -1 \end{pmatrix} \quad x = (1, 1, 1)$$

$w_{ii} \neq 0 \rightarrow$  Oszillation

# Hopfield-Netz – Training

direkte Berechnung der Gewichte möglich,  
kein Training notwendig

- ▶ für ein zu speicherndes Muster  $x \in \{-1, 1\}^m$ :

$$W = x^T x$$

mit Modifikation: alle Diagonalelemente 0

- ▶ für mehrere zu speichernde Muster  
 $x^{(1)} \in \{-1, 1\}^m, \dots, x^{(k)} \in \{-1, 1\}^m$ :

$$W = \sum_{i=1}^n \left( x^{(i)} \right)^T x^{(i)}$$

mit Modifikation: alle Diagonalelemente 0

Beispiel (Tafel):

- ▶ ein Muster:  $x = (1, -1, 1, 1)$
- ▶ mehrere Muster:  $x^{(1)} = (-1, 1, -1)$  und  $x^{(2)} = (1, -1, 1)$

# Unüberwachtes Lernen

bei unbekannter Zielfunktion

zur Analyse von Datenmengen

Ziel: Gruppierung ähnlicher Daten (Clustering)

z. B. erste Schicht in RBF-Netzen

Generalisierung durch Einordnung neuer Daten in vorhandene Gruppen (Cluster)

Training durch Menge von Trainingmustern

$$T = \{x^i \mid i \in \{1, \dots, k\} \wedge x^i \in \mathbb{R}^m\}$$

Methode: Wettbewerbslernen (competitive learning)

nur „Gewinner“-Neuron feuert

# Bündeln von Mustern (Clustering)

Eingabe: Menge von Trainingsmustern  $\{x^{(i)}, \dots, x^{(m)}\}$

Ziel: Gruppierung **ähnlicher** Muster

Anordnung von Mustern in Bündeln:

- ▶ Ähnlichkeit aller Muster eines Clusters
- ▶ Trennung von Mustern mit wenig Gemeinsamkeiten

## Ähnlichkeit im $\mathbb{R}^n$

Zwei Punkte  $x, y \in \mathbb{R}^n$  sind **einander ähnlich**, falls sie einen **geringen Abstand** voneinander haben

Beispiele für Abstandsfunktionen: siehe RBF-Netze

statt Euklidischem Abstand

$$d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

quadrierter Euklidischer Abstand (einfachere Berechnung)

$$d(x, y) = \sum_{i=1}^n (x_i - y_i)^2$$

# Zuordnung von Mustern zu Bündeln

Eingaben:

- ▶ Menge  $\{1, \dots, k\}$  von Bündeln,
- ▶ Muster  $x$

$x$  wird dem Bündel  $j \in \{1, \dots, k\}$  zugeordnet, von welchem es den **geringsten Abstand** hat

Idee: Jedes Bündel  $j \in \{1, \dots, k\}$  hat ein **Zentrum**  $p_j$   
(durchschnittliche Position aller Punkte im Bündel, Prototyp)

neues Muster  $x$  wird dem Bündel  $j$  genau dann zugeordnet, wenn der Abstand von  $x$  zum Zentrum  $p_j$  des Bündels  $j$  minimal ist, d.h.

$$\text{gdw. } \forall i \in \{1, \dots, k\} : d(x, p_j) \leq d(x, p_i)$$

geometrisch: Zerlegung des  $R^n$  in Gebiete  
z.B. im  $\mathbb{R}^2$ : Voronoi-Diagramme

# Selbstorganisierende Karte (SOM)

(Teuvo Kohonen, ca. 1980)

Netz-Topologie: Ein-Schicht-Feed-Forward-Netz mit

- ▶ Eingabe  $(x_1, \dots, x_m) \in \mathbb{R}^m$ , interpretiert als Punkt im  $m$ -dimensionalen Eingaberaum  $\mathbb{R}^m$
- ▶ Ausgabeneuronen  $(p_1, \dots, p_n)$  (Bündelneuronen) interpretiert als  $n$  Punkte  $p_i = (w_{1i}, \dots, w_{mi})$  im  $\mathbb{R}^m$ , Zentren der Bündel
- ▶ vollständig verbunden (alle Vorwärtskanten), Eingabegewichte zum Neuron  $j$ :  $(w_{1j}, \dots, w_{mj})$  interpretiert als Koordinaten des Zentrums des Bündels  $j$

# Selbstorganisierende Karte – Funktionen

Funktionen des Ausgabeneurons  $j$ :

- ▶ Eingabefunktion: Distanzfunktion (Metrik)  
z. B. quadrierter Euklidischer Abstand:

$$l_j(x_1, \dots, x_m) = \sum_{i=1}^m (w_{ij} - x_i)^2$$

- ▶ Aktivierungsfunktion  
(Auswahl des Neurons mit dem kleinsten Abstand von  $x$ )

$$A_j(l_j(x)) = \begin{cases} 1 & \text{falls } \forall l \in \{1, \dots, n\} : l_j(x) \leq l_l(x) \\ 0 & \text{sonst} \end{cases}$$

Wettbewerbslernen: genau **ein** Ausgabeneuron feuert

- ▶ Ausgabefunktion: Identität

# Selbstorganisierende Karte – Lernen

Beginn: zufällige Eingangsgewichte der Ausgabeneuronen  
(Anordnung der Bündelzentren)

Trainieren der Gewichte

für jede Eingabe  $x = (x_1, \dots, x_m)$ :

- ▶ Eingabe  $x$
- ▶ Anpassung der Eingangsgewichte des aktivierten Bündelneurons  $j$   
(Verschiebung des Bündelzentrums  $p_j$  in Richtung des aktuellen Eingabevektors  $x$ )

$$w'_{ij} = w_{ij} + \eta(x_i - w_{ij})$$

# Topologie-erhaltende SOM

Idee: Topologie-erhaltende Abbildung in andere (meist geringere) Dimension

**Anordnung** der Bündelneuronen (z.B. linear, eben, räumlich)

Blockparty:

Berücksichtigung der **Nachbarschaft** (räumlichen Beziehungen) zwischen den Bündelneuronen (in der definierten Anordnung) bei der Aktualisierung der Positionen

z.B. Graph, Gitter, Abstandsfunktionen

# Training Topologie-erhaltender SOM

- ▶ Start mit zufälligen Gewichten (Bündelzentren)
- ▶ Anpassung der Eingangsgewichte des aktivierten Bündelneurons  $p_j$  (in Richtung des Eingabevektors)
- ▶ Anpassung der Eingangsgewichte aller Bündelneuronen  $p_k$  in einer (zuvor festgelegten) **Nachbarschaft** des aktivierten Bündelneurons  $p_j$

$$w'_{ik} = w_{ik} + n(p_k, p_j)\eta(x_i - w_{ik})$$

mit Nachbarschaftsfunktion  $n(p_k, p_j)$   
(Einfluss fällt mit wachsendem Abstand auf 0,  
z.B. RBF-Funktion)

# SOM – Trainingsverlauf

Idee:

Anpassung (Verringerung) von Lernrate und Radius während der Lernphase

Heuristik (sinnvoller Trainingsverlauf):

- ▶ Start mit
  - ▶ großem Radius (nahe halbem Kartenradius)
  - ▶ großer Lernrate (nahe 1)unverändert über ca. 1000 Iterationen  
**Ordnungsphase**
- ▶ Nachbarschafts-Radius (Einflussbereich jedes Bündelneurons) und Lernrate werden mit der Zeit verringert (über ca. 10 000 Iterationen)  
**Feinabstimmung**

# SOM – Beispiele

Topologie-erhaltende Abbildung zwischen Räumen (evtl. verschiedener Dimensionen)

Beispiele:

- ▶ Projektion Kugel  $\rightarrow$  Ebene  
(Robinson-Projektion: Erde  $\rightarrow$  Weltkarte)
- ▶ Sensorische Karten:  
Abbildung von Reizen auf benachbarten Bereichen der Haut  
auf benachbarte Bereiche im Gehirn
- ▶ Abbildung von akustischen Reizen benachbarter Frequenzen  
auf benachbarte Bereiche im Gehirn

# Was bisher geschah

- ▶ Daten, Information, Wissen, Lernen, Intelligenz
- ▶ explizites und implizites Wissen

Ziel der KI:

Unterstützen / Treffen sinnvoller Entscheidungen

Methode:

Bewertung / Klassifikation der Optionen (z.B. Aktionen, Objekte)

- ▶ symbolische KI:  
heuristische Suche, Spielbäume / -graphen
- ▶ statistische KI:  
maschinelles Lernen  
künstliche neuronale Netze

# Zusammenfassung Maschinelles Lernen mit KNN

Ziel: Bewertung / Klassifikation

Methode: **Approximation von Funktionen** (gegeben als Trainingsmenge)  
schon bekannt aus Mathematik, z.B.

- ▶ **lineare** Approximation

Daten: Menge  $\{(\vec{x}, \vec{y})\}$  (Trainingsmenge)

**Ansatz:** lineare Funktion  $a(\vec{x}) = M\vec{x} + \vec{b}$

Approximation durch direkte Minimierung des (z.B. quadratischen) Fehlers

- ▶ Taylorreihen (in kleinen Umgebungen), Ansatz: Polynome

- ▶ Fourierreihen (für periodische Funktionen)

Ansatz: Linearkombinationen von Winkelfunktionen

in KNN (DNN, CNN, ...)

- ▶ Funktion an Stützstellen gegeben (Trainingsmenge)

- ▶ **Ansatz:** durch Topologie des Netzes festgelegte Funktion

- ▶ Approximation durch Training (Minimierung des Fehlers)

Nachteil: implizites Wissen, keine Erklärung für Bewertung möglich

# WH: Regelbasierte Systeme

Idee:

- ▶ Wissensrepräsentation:  
Darstellung des Kontextes als Menge von **Regeln**
- ▶ Wissensverarbeitung:  
Ableiten von Entscheidungen durch logisches Schließen
- ▶ Vorteil: explizites Wissen,  
nachvollziehbar, liefert Erklärung,  
praktisch häufig ausreichend,  
Transformation in alternative (auch für Nicht-Informatiker)  
intuitive Darstellungen (Entscheidungstabellen, -bäume,  
-diagramme) möglich
- ▶ Nachteil: bei unvollständigem / unsicherem Kontextwissen  
schwierig bis unmöglich

# WH: Regelbasierte Systeme

Regel : Implikation  $r = (\varphi \rightarrow \psi)$

normale Regel :  $r = ((b_1 \wedge \dots \wedge b_n \wedge \neg c_1 \wedge \dots \wedge \neg c_m) \rightarrow h)$   
mit (aussagen- oder prädikatenlogischen) Atomen  
 $b_1, \dots, b_n, c_1, \dots, c_m, h$

definite Regel (Hornklausel, ohne negative Voraussetzungen):  
Implikation  $r = (\varphi \rightarrow \psi)$  mit  
 $\varphi = (b_1 \wedge \dots \wedge b_n)$  und  $\psi = h$   
mit (aussagen- oder prädikatenlogischen) Atomen  
 $b_1, \dots, b_n, h$

Bestandteile der Regel  $r = (\varphi \rightarrow \psi)$ :

Kopf  $\psi$  (Folgerung)

Rumpf  $\varphi$  (Voraussetzung)

oft  $\varphi = (b_1 \wedge \dots \wedge b_n \wedge \neg c_1 \wedge \dots \wedge \neg c_m)$  mit  
positiven Voraussetzungen  $b_1, \dots, b_n$   
und negativen Voraussetzungen  $c_1, \dots, c_m$

# Regelbasiertes Problemlösen

gegeben:

- ▶ Wissensbasis  $P = F \cup R$  (z.B. logisches Programm)  
Jede Regel mit Variablen repräsentiert die Mengen aller Instanzen.
- ▶ Anfrage (Hypothese)  $h$   
evtl. mit Variablen

gesucht:

- ▶ Nachweis, dass die Hypothese  $h$  aus der Wissensbasis folgt  
Lösung: erfüllende Belegung der Variablen  
(falls in Anfrage enthalten)

übliche Verfahren:

**Rückwärtsverkettung:** Suche nach Antworten auf eine gezielte Anfrage (zielorientiert), z.B. Resolution, Prolog

**Vorwärtsverkettung:** Ableitung neuen Wissens als Folgerungen aus Faktenmengen (datenorientiert), Fixpunkt-Operator, z.B. Datalog

# Beispiel

Aus der Wissensbasis

F1 Tom ist ein Baby.  $b(t)$

F2 Tom ist männlich.  $m(t)$

F3 Anna ist weiblich.  $w(a)$

R1 Babies sind Kinder.  $\forall x : b(x) \rightarrow k(x)$

R2 Männliche Kinder sind Jungen.  $\forall x : (m(x) \wedge k(x)) \rightarrow j(x)$

R3 Weibliche Kinder sind Mädchen.  $\forall x : (w(x) \wedge k(x)) \rightarrow g(x)$

- ▶ Antwort auf Frage „Wer ist ein Junge“ (Für welche  $x$  gilt  $j(x)$ ) durch Rückwärtsverkettung:  $x = t$
- ▶ Durch Vorwärtsverkettung (ohne gezielte Anfrage) folgt:  
Tom ist ein Kind. Tom ist ein Junge.

# Interpretation von Regeln

(normale) Regel

$$r = ((b_1 \wedge \dots \wedge b_n \wedge \neg c_1 \wedge \dots \wedge \neg c_m) \rightarrow h)$$

mögliche Interpretation des Regelkopfes  $h$  als

**Aussage** , z.B. Kälte  $\wedge$  Niederschlag  $\rightarrow$  Schnee

**Deduktionsregeln**

**Aktion** , z.B. heiße Stirn  $\rightarrow$  Fieber messen

**Aktionsregeln**, Produktionsregeln

# Konflikte bei Aktionsregeln

gegeben: Faktenmenge  $F$ , Regelmenge  $R$  (z.B. definit)

Regel  $b_1 \wedge \dots \wedge b_n \rightarrow h$  in  $F$  genau dann **anwendbar**, wenn ihr Rumpf  $b_1 \wedge \dots \wedge b_n$  in  $F$  erfüllt ist (also wenn  $\{b_1, \dots, b_i\} \in F$ )

Problem:

- ▶ Faktenmenge (Zustand)  $F'$  mit mehreren anwendbaren Regeln aus  $R$ , die gegenteilige Aktionen auslösen
- ▶ Aktionsteil welcher Regel soll ausgeführt werden?

**Konfliktmenge** zu einer Faktenmenge  $F'$ :

Menge aller in  $F'$  anwendbaren Regeln aus  $R$

# Beispiel

Wissensbasis:

- R1 Wenn Besuch zum Essen kommt,  
gibt es Weiß- oder Rotwein.
- R2 Zum Fisch gibt es Weißwein.
- R3 Wenn Bob da ist, gibt es Rotwein.
- F Bob kommt zum Fischessen.

Frage: Welches Getränk soll serviert werden?

# Konfliktlösestrategien

Konfliktlösung durch **Auswahl** einer Teilmenge aller anwendbaren Regeln, durch deren Anwendung kein Konflikt entsteht.

oft durch (schrittweises) Entfernen ausgewählter Regeln

(einige) Konfliktlösestrategien:

**Refraktionsprinzip:** Keine Regelinstanz darf zweimal nacheinander feuern

**Präferenzen:** Bei der Wissensrepräsentation definierte Ordnung auf Regeln  
Regeln niederer Präferenzen feuern nicht, wenn eine Konflikt-Regel höherer Präferenz angewendet wurde.

**Spezifität:** spezifischere Regelinstanzen werden bevorzugt

**Aktualität:** erstmals feuernde Regeln werden bevorzugt

**zufällige Auswahl** einer Regel

# Regelbasierte Systeme mit Unsicherheit

Problem bei Bestimmung von Merkmalswerten  
(Antworten auf Fragen), falls Wert

- ▶ unbekannt
- ▶ ungenau
- ▶ unsicher, unzuverlässig
- ▶ genauere Untersuchung unmöglich, zeitaufwendig, teuer

Abhilfe z.B. durch Hinzufügen von Heuristiken (Erfahrungswerte, Schätzungen, Wahrscheinlichkeiten) zu Regeln

# Unsichere Regelsysteme

Idee: Zuordnung eines Gewichtes zu jeder Regel  
repräsentiert Heuristik, Schätzungen, Wahrscheinlichkeit

Beispiel:

- ▶ Fakten (Merkmale):  $M = \{p_1, p_2, p_3, q_1, q_2, r\}$
- ▶ definierte Atome (Klassen) :  $K = \{h\}$
- ▶ Gewichte:  $G = \mathbb{R}$
- ▶ Regeln:

$$p_1 \wedge p_2 \wedge p_3 \rightarrow h \quad (0.7)$$

$$q_1 \vee q_2 \rightarrow h \quad (0.4)$$

$$r \rightarrow h \quad (-0.8)$$

auch dafür Transformation in Entscheidungstabellen, -bäume,  
-diagramme möglich (Scores, z.B. Wahl-O-Mat)

# Unsicheres Schließen

Ziele:

- ▶ Bewertung, Vergleich verschiedener Lösungen
- ▶ Auswahl der **besten** Lösung

Idee:

Schrittweise Berechnung der Werte der Lösungen  
über Zwischenlösungen  
durch Vorwärts- oder Rückwärtsverkettung

Wert der Lösung im Regelkopf wird berechnet aus:

- ▶ Werte der Voraussetzungen,
- ▶ Wert der Regel

## Beispiel

- ▶ Fakten (Merkmale):  $\{p_1, p_2, p_3, q_1, q_2, r\}$
- ▶ definierte Atome (Klassen):  $\{h\}$
- ▶ Gewichte:  $\in \mathbb{R}$
- ▶ Regeln:

$$p_1 \wedge p_2 \wedge p_3 \rightarrow h \quad (0.7)$$

$$q_1 \vee q_2 \rightarrow h, \quad (0.4)$$

$$r \rightarrow h, \quad (-0.8)$$

Werte der Fakten:

$$w(p_1) = w(p_2) = w(q_2) = 1, \quad w(p_3) = 0.6, \quad w(q_1) = w(r) = 0.5$$

Berechnung der Werte der Regelrümpfe mit  $\wedge \mapsto \min$ ,  $\vee \mapsto \max$

Anwendung der Regeln: Multiplikation mit Regelgewicht

Berechnung des Wertes von  $h$  (oft Maximum oder Summe)

## Beispiel MYCIN (1972)

- ▶ Regelbasiertes System  
einfache Regeln der Form: Merkmalswerte  $\rightarrow$  Klasse  
(vermeidet zyklische Auswertung)
- ▶ Regeln  $r$  (und Fakten) mit Gewissheitsfaktor  $\gamma_r \in [-1, 1]$
- ▶ logischen Junktoren zugeordnete Berechnung der Gewissheitsfaktoren:

$$\gamma_{\varphi \vee \psi} = \max(\gamma_{\varphi}, \gamma_{\psi}) \quad \gamma_{\varphi \wedge \psi} = \min(\gamma_{\varphi}, \gamma_{\psi}) \quad \gamma_{(h|B)} = \gamma_{(B \rightarrow h)} \max(0, \gamma_B)$$

- ▶ Regel für Kombination der Werte bei mehreren anwendbaren Regeln  
 $B_1 \rightarrow h, \dots, B_n \rightarrow h$ :

$$\gamma_{(h|B_1, \dots, B_n)} = \gamma_{(h|B_1, \dots, B_{n-1})} \oplus \gamma_{(h|B_n)} \quad \text{mit}$$

$$\text{falls } \gamma_1 \geq 0, \gamma_2 \geq 0 \text{ dann } \gamma_1 \oplus \gamma_2 = \gamma_1 + \gamma_2 - \gamma_1 \gamma_2$$

$$\text{falls } \gamma_1 < 0, \gamma_2 < 0 \text{ dann } \gamma_1 \oplus \gamma_2 = \gamma_1 + \gamma_2 + \gamma_1 \gamma_2$$

$$\text{falls } \gamma_1 \gamma_2 < 0 \text{ dann } \gamma_1 \oplus \gamma_2 = \frac{\gamma_1 + \gamma_2}{1 - \min(|\gamma_1|, |\gamma_2|)}$$

(Nachbildung wahrscheinlichkeitstheoretischer Regeln)

# Regelbasierte Systeme: Beobachtung

Regeln werden verwendet zum

- ▶ deduktiven Schließen  
z.B. Krankheit → Symptome
- ▶ Schließen von Merkmalen / Klassen auf Aktionen  
z.B. Symptome → Test, Krankheit → Therapie
- ▶ Schließen von Merkmalen / Klassen auf Klassen  
z.B. Symptome → Diagnosen, Verbrechen → Strafe

Richtung der Regeln entspricht also

- ▶ Richtung von (einfachen) Formulierungen, menschlichem Erklären
- ▶ aber nicht notwendig der Richtung **Ursache → Wirkung**

menschliches Denken wechselt intuitiv zwischen beiden Richtungen  
maschinelle Verfahren verwenden festgelegte Richtung

## WH: Bayes-Netze

Bayes-Netz (alternativ: believe network)  $N = (G, \Theta)$  mit

- ▶ DAG  $G$  mit Knoten: Zufallsvariablen
- ▶  $\Theta$ : für jeden Knoten  $X$  mit Eltern  $X_1, \dots, X_k$ :  
Wahrscheinlichkeiten  $P(X = e \mid X_1 = e_1 \cap \dots \cap X_k = e_k)$   
für alle  $[e, e_1, \dots, e_k] \in W^{k+1}$
- ▶ Semantik (induktiv):  
 $N$  beschreibt Wahrscheinlichkeits-Raum durch  
$$P(X = e) = P(X = e \mid \dots X_k = e_k \dots) \cdot \prod_k P(X_k = e_k)$$
  
IA: Quellen des DAG (ohne Vorgänger, Bedingungen, d.h.  
 $\prod \emptyset = 1$ )

Beispiel: Alarm

# Inferenz mit BN

- ▶ die Diagnose-Aufgabe: gegeben ein BN, gesucht sind bedingte Wahrscheinlichkeiten der Ursache(n), unter der Bedingung von Beobachtungen
- ▶ Bsp:  $P(\text{Einbruch} = 1 \mid \text{John} = 1 \cap \text{Mary} = 1)$
- ▶ kann exakt bestimmt werden, dauert jedoch  $2^{|\mathcal{N}|}$   
kann nicht besser gehen, weil aussagenlogische Erfüllbarkeit auf dieses Inferenzproblem reduziert werden kann
- ▶ die Alternative sind schnellere (Simulations)Verfahren, die einen Näherungswert liefern

## Bayes-Netze: Beobachtung

- ▶ Richtungen der Kanten repräsentieren bedingte Wahrscheinlichkeiten
- ▶ Kanten und deren Richtung entsprechen nicht notwendig den tatsächlichen kausalen Beziehungen.
- ▶ Modellierung von kausalen Beziehungen von / zu unbeobachtbaren (Stör-)Größen nicht möglich
- ▶ Ob die Richtung der kausalen Beziehung entspricht, kann nicht im Modell festgestellt werden.

Also: Modellierung gesucht, die

- ▶ unbeobachtbare Größen enthalten kann
- ▶ Bestätigung / Feststellung von kausalen Beziehungen ermöglicht

# Kausalitätsstufen

1. Assoziation (Wahrnehmen)  
Was ist?
2. Intervention (Beeinflussen)  
Was geschieht, wenn Aktion X ausgeführt wird?
3. Counterfactuals (Vorstellen, Reflektieren)  
**Warum** geschah ... ?  
Lag es an Aktion X?  
Was wäre, wenn Aktion X nicht stattgefunden hätte?

Methoden:

- ▶ Kausale Modelle (Diagramme)
- ▶ do-Kalkül
- ▶ Counterfactuals

# Kausalität

Was ist Ursache? Was Wirkung?

- ▶ Hahn kräht, Sonne geht auf
- ▶ brennende Kerze am Fenster, brennende Gardine

Gibt es andere Ursachen, die beides beeinflussen?

# Kausales Wissen

Wie erwerben Kinder Wissen über Zusammenhänge zwischen Ursache und Wirkungen?

- ▶ (passives) Beobachten
- ▶ aktive Experimente
- ▶ Begleitung und Unterstützung durch Bezugspersonen (Erklärungen, Hinweise auf Unterschiede)

Annahme:

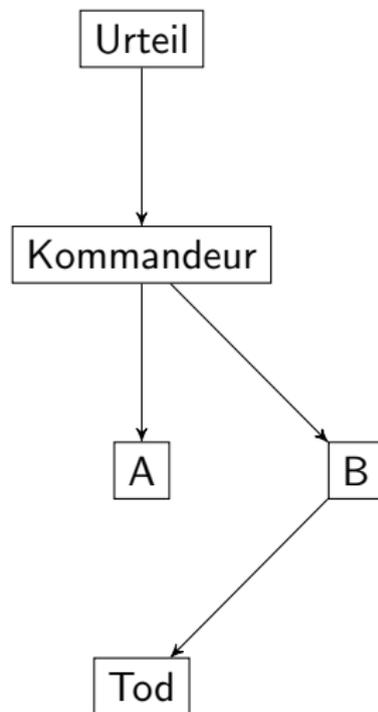
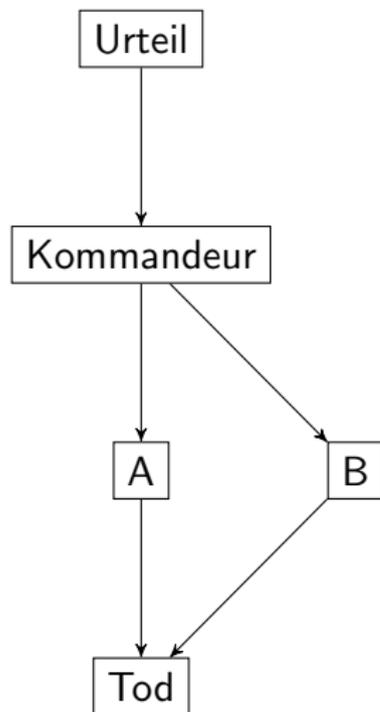
Menschen erwerben / haben **mentale Modelle**  
kompakte Darstellung der Informationen, die zum kausalen Schließen benötigt werden

Mini-Turing-Test:

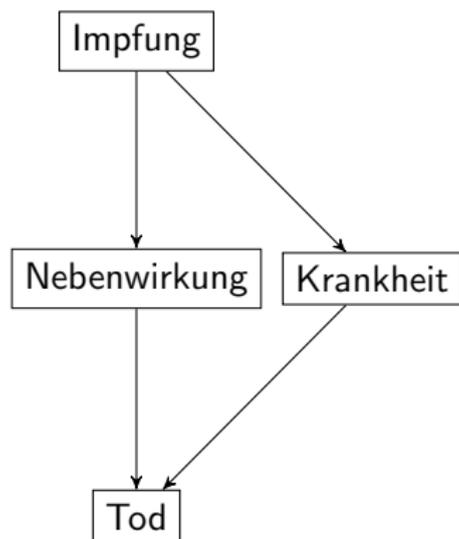
- ▶ Wie lässt sich **kausales** Wissen (analog) maschinell erwerben?
- ▶ Was sind geeignete (maschinelle) Repräsentationen des zum kausalen Schließen notwendigen Wissens?

# Kausaldiagramm – Beispiel

(Beispiel nach J. Pearl: Book of WHY)



## Kausaldiagramm – Beispiel



$$p(I) = 0.99, \quad p(\neg I) = 0.01$$

$$p(N|I) = 0.01, \quad p(N|\neg I) = 0$$

$$p(K|I) = 0, \quad p(K|\neg I) = 0.02$$

$$p(T|N) = 0.01$$

$$p(T|K) = 0.2$$

bei Impfrate 99%:

40 Tote durch Krankheit + 99 Tote durch Impfung

Counterfactual:

Was wäre ohne Impfung ( $p(I) = 0$ )?

4000 Tote durch Krankheit

# Kausal-Diagramme

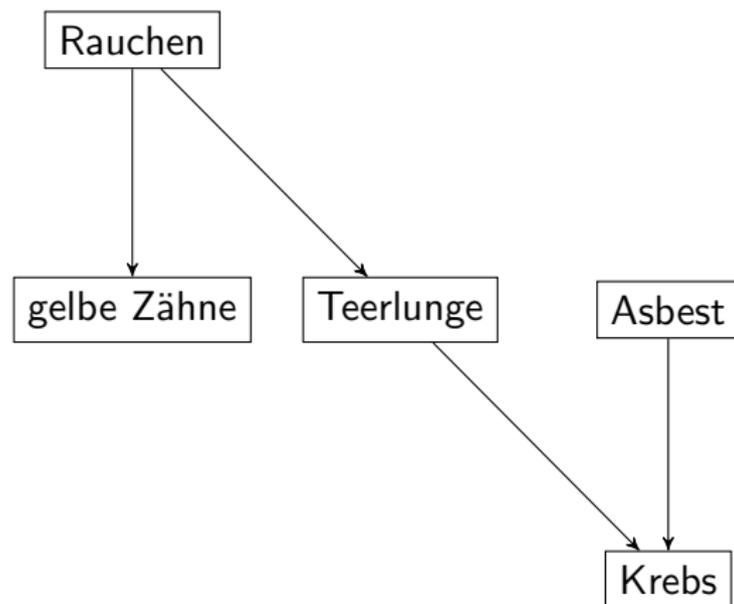
Graphen wie in Bayes-Netze:

- ▶ DAG
- ▶ Knoten: Zufallsvariablen
- ▶ Eltern-Knoten repräsentieren direkte Ursachen
- ▶ Kind-Knoten repräsentieren direkte Wirkungen
- ▶ Vorfahren repräsentieren indirekte Ursachen
- ▶ Nachkommen repräsentieren indirekte Wirkungen

# Arten von Knoten

- ▶ beobachtet (seeing)  
Wahrscheinlichkeiten bekannt  
(experimentell, Studien)
- ▶ beeinflussbar (doing)  
z.B. festen Wert (nicht Wahrscheinlichkeit!) zuordnen  
unbeobachtet,  
Wahrscheinlichkeiten unbekannt

## Kausaldiagramm – Beispiel



Kausalmodell repräsentiert

$$p(Z, R, A, T, K) = p(R)p(A)p(T|R)p(Z|R)p(K|A, T)$$

# Interventionen

Modifikation des Graphen

Beispiel:

- ▶ Bildung  $\rightarrow$  Einkommen  $\rightarrow$  Steuern
- ▶ Intervention: Erbschaft

Interventionen

- weich:** neue Knoten mit Kanten in das bestehende Modell  
im Bsp: (kleine) Erbschaft  $\rightarrow$  Einkommen,
- hart:** neue Knoten mit Kanten in das bestehende Modell  
und Löschen bestehender Kanten  
im Bsp: (große) Erbschaft  $\rightarrow$  Einkommen,  
Eltern von Einkommen (Bildung) verlieren Einfluss  
(Kante wird gelöscht)

# Prominente Teilgraphen

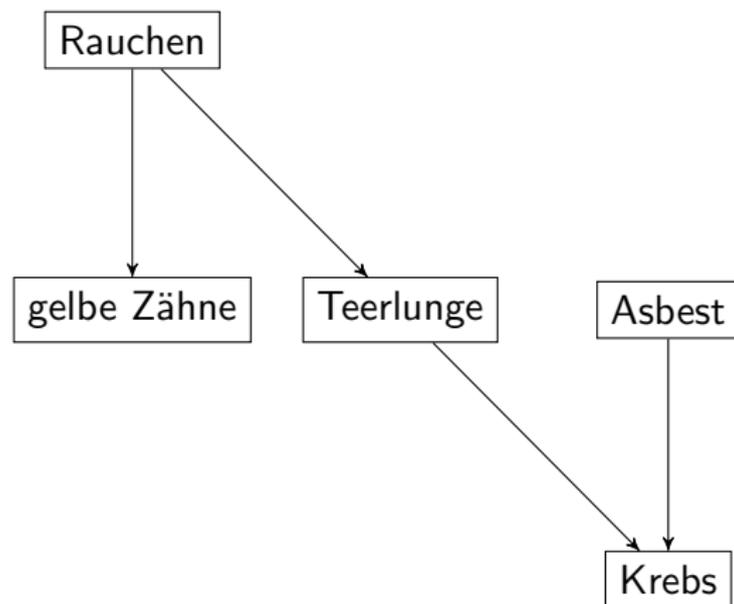
(gerichtete) Pfade aus drei Knoten:

Kette  $\boxed{u} \rightarrow \boxed{v} \rightarrow \boxed{w}$   
 $v$  heißt **Mediator**

Verzweigung  $\boxed{u} \leftarrow \boxed{v} \rightarrow \boxed{w}$   
aus gemeinsamem Eltern-Knoten  $v$   
 $v$  ist **direkte gemeinsame Ursache**

$v$ -Struktur  $\boxed{u} \rightarrow \boxed{v} \leftarrow \boxed{w}$   
 $v$  ist **direkte gemeinsame Wirkung** (collider)

## Kausaldiagramm – Beispiel



- ▶ R ist Mitte der Verzweigung  $Z \leftarrow R \rightarrow T$
- ▶ T ist Mitte der Kette  $R \rightarrow T \rightarrow K$  (Mediator)
- ▶ K ist Mitte der v-Struktur  $T \rightarrow K \leftarrow A$  (Collider)

# Was bisher geschah

- ▶ WH: Bayes-Netze
- ▶ Kausalitäts-Stufen (ladder of causality)
- ▶ Kausaldiagramme
- ▶ typische Teilgraphen (Pfade mit drei Knoten):
  - ▶ Kette
  - ▶ Verzweigung
  - ▶ v-Struktur (collider)
- ▶ Interventionen: hart, weich

# Blockierte Pfade in Kausaldiagrammen

gegeben:

- ▶ Kausaldiagramm  $G = (V, E)$ ,
- ▶ Menge  $E \subseteq V$  von (beobachteten) Knoten (Evidenz)

$$N(p) = \{q \in V \mid (p, q) \in E^*\}$$

Menge aller (gerichteten) Nachfahren von  $p$  (incl.  $p$ )

$$V(p) = \{q \in V \mid (q, p) \in E^*\}$$

Menge aller (gerichteten) Vorfahren von  $p$  (incl.  $p$ )

für Knotenmengen  $U \subseteq V$ :  $N(U) = \bigcup_{u \in U} N(u)$ ,  $V(U) = \bigcup_{u \in U} V(u)$

Idee: (ungerichtete) Pfade repräsentieren indirekte Einflüsse

(ungerichteter) Pfad  $P$  in  $G$  ist **blockiert durch  $E$**  gdw.

1.  $P$  enthält Kette  $\boxed{x} \rightarrow \boxed{p} \rightarrow \boxed{y}$  mit  $p \in E$  oder
2.  $P$  enthält Verzweigung  $\boxed{x} \leftarrow \boxed{p} \rightarrow \boxed{y}$  mit  $p \in E$  oder
3.  $P$  enthält v-Struktur  $\boxed{x} \rightarrow \boxed{p} \leftarrow \boxed{y}$  mit  $N(p) \cap E = \emptyset$

# Blockade in prominenten Teilgraphen

Kette  $\boxed{u} \rightarrow \boxed{v} \rightarrow \boxed{w}$

z.B. Feuer  $\rightarrow$  Rauch  $\rightarrow$  Alarm

- ▶  $E = \emptyset$ : nicht blockiert
- ▶  $E = \{u\}$ : nicht blockiert
- ▶  $E = \{v\}$ : blockiert
- ▶  $E = \{u, v\}$ : blockiert

Verzweigung  $\boxed{u} \leftarrow \boxed{v} \rightarrow \boxed{w}$

z.B. Schuhgröße  $\leftarrow$  Alter  $\rightarrow$  Lesefähigkeit

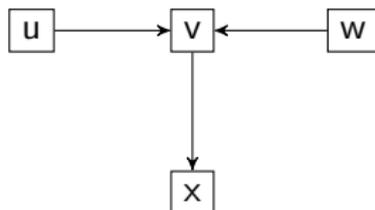
- ▶  $E = \emptyset$ : nicht blockiert
- ▶  $E = \{u\}$ : nicht blockiert
- ▶  $E = \{v\}$ : blockiert
- ▶  $E = \{u, v\}$ : blockiert

v-Struktur  $\boxed{u} \rightarrow \boxed{v} \leftarrow \boxed{w}$

z.B. Talent  $\rightarrow$  Berühmtheit  $\leftarrow$  Schönheit

- ▶  $E = \emptyset$ : blockiert
- ▶  $E = \{u\}$ : blockiert
- ▶  $E = \{v\}$ : nicht blockiert
- ▶  $E = \{u, v\}$ : nicht blockiert

# Beispiel



Menge  $E = \emptyset$ :

- ▶ Pfad  $u \rightarrow v \leftarrow w$  blockiert
- ▶ Pfad  $u \rightarrow v \rightarrow x$  nicht blockiert
- ▶ Pfad  $x \leftarrow v \leftarrow w$  nicht blockiert

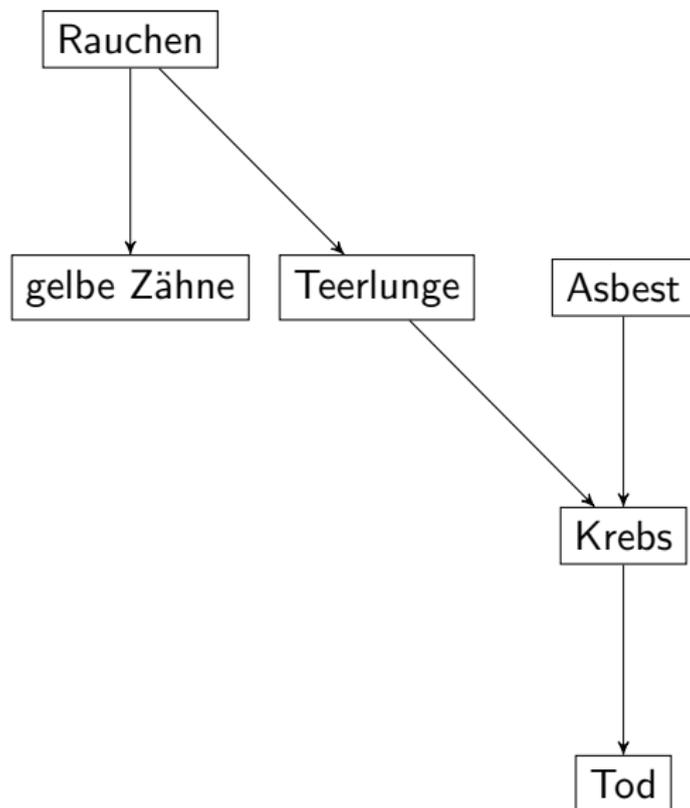
Menge  $E = \{v\}$ :

- ▶ Pfad  $u \rightarrow v \leftarrow w$  nicht blockiert
- ▶ Pfad  $u \rightarrow v \rightarrow x$  blockiert
- ▶ Pfad  $x \leftarrow v \leftarrow w$  blockiert

Menge  $E = \{x\}$ :

- ▶ Pfad  $u \rightarrow v \leftarrow w$  nicht blockiert
- ▶ Pfad  $u \rightarrow v \rightarrow x$  nicht blockiert
- ▶ Pfad  $x \leftarrow v \leftarrow w$  nicht blockiert

## Beispiel



## d-separierte Knoten

gegeben:

- ▶ Diagramm  $G = (V, E)$ ,
- ▶ Menge  $E \subseteq V$  von (beobachteten) Knoten (Evidenz)
- ▶ Knoten  $u, v \in V$

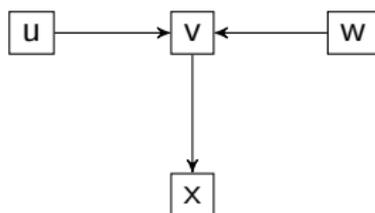
$u$  und  $v$  heißen genau dann in  $G$  **d-separiert durch  $E$**  ( $u \perp\!\!\!\perp v | E$ ), wenn

- ▶ in  $G$  kein (ungerichteter) Pfad zwischen  $u$  und  $v$  existiert oder
- ▶ jeder (ungerichtete) Pfad zwischen  $u$  und  $v$  blockiert ist.

d-Separation ist symmetrisch:  $u \perp\!\!\!\perp v | E$  gdw.  $v \perp\!\!\!\perp u | E$

Zusammenhang mit Statistik (Bestätigung des Kausaldiagramms):  
Sind  $u$  und  $v$  d-separiert durch  $E$ , dann  
sind die Zufallsgrößen  $u$  und  $v$  bedingt unabhängig unter  $E$

# Beispiel



Menge  $E = \emptyset$ :

- ▶  $u \perp\!\!\!\perp w|E$  (und  $w \perp\!\!\!\perp u|E$ ) über Pfad  $u \rightarrow v \leftarrow w$
- ▶ alle anderen Knotenpaare nicht d-separiert

Menge  $E = \{v\}$ :

- ▶  $u \perp\!\!\!\perp x|E$  (und  $x \perp\!\!\!\perp u|E$ ) über Pfad  $u \rightarrow v \rightarrow x$
- ▶  $x \perp\!\!\!\perp w|E$  (und  $w \perp\!\!\!\perp x|E$ ) über Pfad  $x \leftarrow v \leftarrow w$
- ▶ alle anderen Knotenpaare nicht d-separiert

Menge  $E = \{x\}$ :

- ▶  $u \perp\!\!\!\perp w|E$  (und  $w \perp\!\!\!\perp u|E$ ) über Pfad  $u \rightarrow v \leftarrow w$
- ▶ alle anderen Knotenpaare nicht d-separiert

# d-separierte Knotenmengen

gegeben:

- ▶ Diagramm  $G = (V, E)$ ,
- ▶ Menge  $E \subseteq V$  von (beobachteten) Knoten (Evidenz)
- ▶ Mengen  $X, Y \subseteq V$

Die Mengen  $X$  und  $Y$  heißen genau dann **d-separiert durch  $E$**  ( $X \perp\!\!\!\perp Y | E$ ), wenn jedes Knotenpaar  $(x, y) \in (X \times Y)$  d-separiert ist.

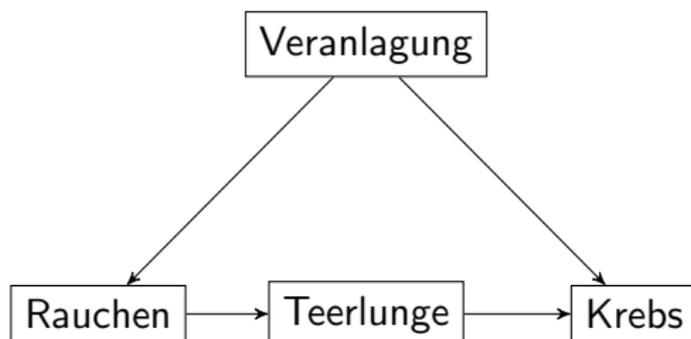
Beispiele: Tafel

# Modifikation von Kausaldiagrammen

zu Kausaldiagramm  $G = (V, E)$  und Knoten  $u \subseteq V$  sind

- ▶  $G_{\bar{u}} = (V, E')$  mit  $E' = E \setminus \{(v, u) \mid v \in V\}$
- ▶  $G_{\underline{u}} = (V, E')$  mit  $E' = E \setminus \{(u, v) \mid v \in V\}$

Beispiele (Tafel):



# do-Kalkül

Idee:

Simulation von Interventionen

(die unmöglich, unmoralisch, aufwendig, ... sind)

Beispiel:

Wenn jemand zu rauchen beginnt, wie stark erhöht sich (seine individuelle) Wahrscheinlichkeit, dass er Krebs bekommt?

Achtung:

Das ist nicht die Frage nach der bedingten Wahrscheinlichkeit  $p(\text{Krebs} = 1 \mid \text{Rauchen} = 1)$

Darstellung der Interventionen in Anfragen durch do-Ausdrücke  
z.B.  $p(\text{Krebs}=1 \mid \text{do}(\text{Rauchen}=1))$

Syntax der Anfragen (analog Statistik):

für Menge  $\{X_i \mid i \in \{1, \dots, n\}\}$  von Zufallsvariablen  
(Knoten im Kausaldiagramm)

$p(x_i \mid b_1, \dots, b_k)$  mit

$\forall j \in \{1, \dots, k\} : b_j \in \{x_k \mid k \in \{1, \dots, n\}\} \cup \{\text{do}(x_k) \mid k \in \{1, \dots, n\}\}$

( $x_i$  ist Wert der Zufallsvariable  $X_i$ )

## do-Kalkül – Umformregeln

- ▶ Daten bieten nur (bedingte) Wahrscheinlichkeiten.  
( Antworten auf do-freie Anfragen )
- ▶ Gesucht sind auch Auswirkungen von Interventionen  
( Antworten auf Anfragen mit do )  
ohne tatsächliche Ausführung der Interventionen.

Methode:

Transformation von Anfragen mit do zu do-freien Anfragen durch Anwendung bedingter Regeln.

Jede Regel besteht aus

- ▶ Modifikation des Kausaldiagramms und
- ▶ syntaktische Umformung der Anfrage

Eliminierung der do-Ausdrücke abhängig von

- ▶ Kausaldiagramm und
- ▶ Anfrage

Es gibt Kombinationen von Anfragen und Kausaldiagrammen, für welche die Eliminierung der do-Ausdrücke nicht möglich ist.

## do-Kalkül – Umformregeln

(bedingte) Regeln:

1. falls  $(Y \perp\!\!\!\perp Z|X, W)_{G_{\bar{X}}}$ :

$$p(y|\text{do}(x), z, w) = p(y|\text{do}(x), w)$$

( Beobachtung  $z$  kann entfernt / hinzugefügt werden )

2. falls  $(Y \perp\!\!\!\perp Z|X, W)_{G_{\bar{X}, \underline{Z}}}$ :

$$p(y|\text{do}(x), \text{do}(z), w) = p(y|\text{do}(x), z, w)$$

( Austausch von Aktion  $\text{do}(z)$  und Beobachtung  $z$  )

3. falls  $(Y \perp\!\!\!\perp Z|X, W)_{G_{\bar{X}, \overline{Z(W)}}$ :

$$p(y|\text{do}(x), \text{do}(z), w) = p(y|\text{do}(x), w)$$

wobei  $Z(W) = Z \setminus V(W)$

( Aktion  $\text{do}(z)$  kann entfernt / hinzugefügt werden )

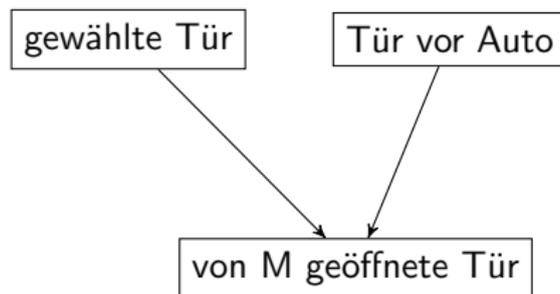
# Was bisher geschah

- ▶ WH: Bayes-Netze
- ▶ Kausalitäts-Stufen (ladder of causality)
- ▶ Kausaldiagramme
- ▶ typische Teilgraphen (Pfade mit drei Knoten):
  - ▶ Kette
  - ▶ Verzweigung
  - ▶ v-Struktur (collider)
- ▶ Interventionen: hart, weich
- ▶ blockierte (ungerichtete) Pfade in Kausaldiagrammen
- ▶ (bedingt) d-separierte Knoten
- ▶ Syntax Anfragen im do-Kalkül

# Paradox: Monty Hall

Gewinnspiel Let's Make a Deal:

- ▶ 3 Türen: dahinter ein Auto, zwei Ziegen
- ▶ Spieler gewinnt Objekte hinter geöffneten Türen.
- ▶ Spieler wählt eine Tür (bleibt zu).
- ▶ Moderator (Monty Hall) öffnet eine andere Tür (vor Ziege).
- ▶ Spieler öffnet Tür, darf aber zuvor seine Wahl ändern. Sollte er?
- ▶ Paradox: Ja, er sollte

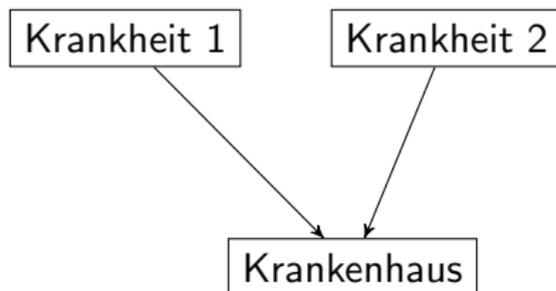


Information durch von M geöffnete Tür erzeugt (nicht-kausalen) Einfluss

## Berksons Paradox

Studie (Berkson, 1946) zu Zusammentreffen verschiedener (unabhängiger) Krankheiten bei Patienten im Krankenhaus  
Paradox:

- ▶ Krankheiten 1 und 2 sind (in der Gesamtbevölkerung) unabhängig voneinander, aber
- ▶ bei Patienten im Krankenhaus treten beide Krankheiten häufiger gleichzeitig auf



Information durch Krankenhausaufenthalt erzeugt (nicht-kausalen) Einfluss zwischen beiden Krankheiten

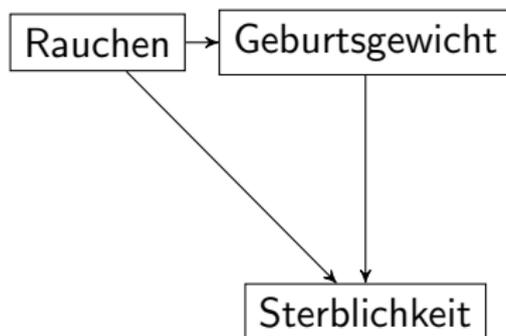
# Paradox: Geburtsgewicht

Zusammenhang zwischen

- ▶ Rauchen der Mutter und
- ▶ Sterblichkeit in ersten Wochen nach Geburt

bekannt

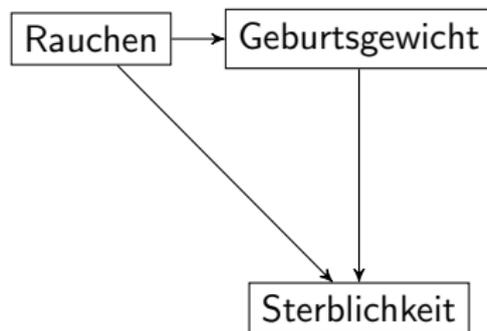
- ▶ Kinder von Raucherinnen haben geringeres Geburtsgewicht
- ▶ Sterblichkeit für Kindern mit gerigerem Geburtsgewicht höher



beobachtet: Geburtsgewicht, Sterblichkeit, Rauchen

## Paradox: Geburtsgewicht

Studie (Yerushalmy 1959) San Francisco Bay Area Gesundheitsdaten von 15000 Kindern vor / nach Geburt



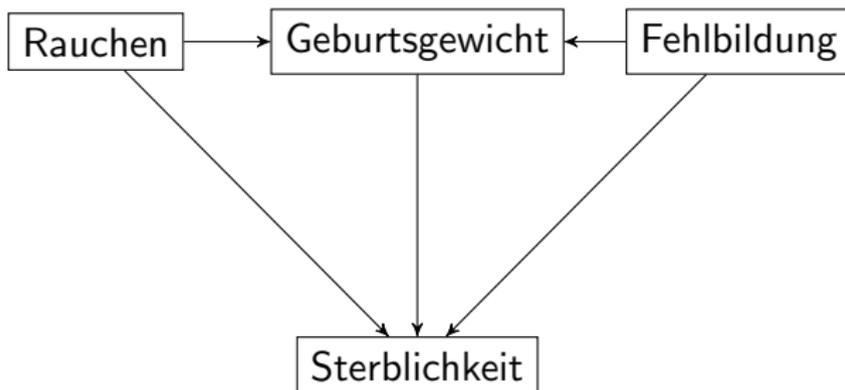
Paradox:

- ▶ Untergewichtig geborene Kinder haben höheres Sterberisiko, aber
- ▶ untergewichtige Kinder von Raucherinnen haben bei gleichem Geburtsgewicht bessere Überlebenschancen

Legt den Schluss nahe, dass (Beginn zu) Rauchen während der Schwangerschaft vorteilhaft für die Überlebenschance des Kindes ist.

## Paradox: Geburtsgewicht

zusätzliche Betrachtung versteckter (unbeobachteter) Merkmale, die Einfluss auf Geburtsgewicht und Sterblichkeit haben



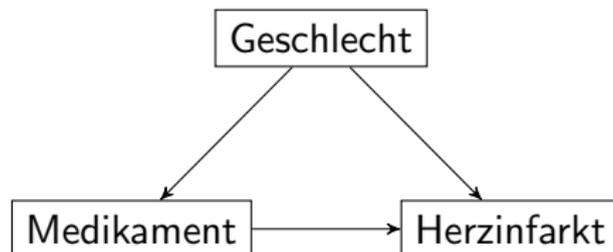
# Simpsons Paradox

Studie (Simpson, 1956) zur Wirkung eines neuen Medikamentes zur Verhinderung von Herzinfarkten, Teilnehmer wählten freiwillig, ob Sie das Medikament nehmen

Medikament → Herzinfarkt

Paradox:

- ▶ Behandlung senkt in Gesamtbevölkerung das Herzinfarkt-Risiko, aber
- ▶ Männer haben mit Behandlung höheres Herzinfarkt-Risiko und
- ▶ Frauen haben mit Behandlung höheres Herzinfarkt-Risiko



## WH: Interventionen

Beispiel: Untersuchung der Wirkung eines Behandlungsverfahrens  $X$  auf eine Krankheit  $Y$  statistisch /experimentell (mitunter) erreichbar durch **Zufallsstudien** (evtl. blind, doppelblind, ...):

- ▶ zufällige Zuordnung der Teilnehmer zu Gruppen  $X = 1$  (mit Behandlung) und  $X = 0$  (Kontrollgruppe ohne Behandlung zum Entfernen von Störgrößen, z.B. Plazebo-Effekt)

**Berechnung mit kausalen Modellen** (falls möglich):

(harte) Intervention auf  $X \subseteq V$  im Kausaldiagramm  $G = (V, E)$

- ▶ Intervention in  $G$ : Löschen aller Eingangskanten in  $X$
- ▶ unterbricht jeden direkten Einfluss anderer Knoten auf  $X$
- ▶ zerstört in  $G$  alle Pfade aus  $X$ , die mit einer Rückwärtskante beginnen, (Hintertür-Wege, back-door path) aus  $X$
- ▶ ersetzt die Kontrollgruppe durch Ausschalten von Störgrößen

## WH: do-Kalkül – Syntax

Idee: Simulation von Interventionen  
(die unmöglich, unmoralisch, aufwendig, ... sind)

Beispiel:  $p(K = 1 | \text{do}(R = 1))$

Wenn jemand zu rauchen beginnt, wie stark erhöht sich  
(seine individuelle) Wahrscheinlichkeit, dass er Krebs bekommt?

Achtung:

Das ist nicht die Frage nach der bedingten Wahrscheinlichkeit  
 $p(K = 1 | R = 1)$

Darstellung der Interventionen in Anfragen durch *do*-Ausdrücke

Syntax der Anfragen (analog Statistik):

für Menge  $\{X_i \mid i \in \{1, \dots, n\}\}$  von Zufallsvariablen  
(Knoten im Kausaldiagramm)

$p(x_i | b_1, \dots, b_k)$  mit

$\forall j \in \{1, \dots, k\} : b_j \in \{x_k \mid k \in \{1, \dots, n\}\} \cup \{\text{do}(x_k) \mid k \in \{1, \dots, n\}\}$

# do-Kalkül

Unterschied zwischen

- ▶ bedingte Wahrscheinlichkeit  $P(Y|X)$
- ▶ (harte) Intervention  $P(Y|\text{do}(X = x))$   
Ausschalten aller Ursachen für  $X$ :  
Löschen aller Eingangskanten zu  $X$  im Kausaldiagramm  
 $G \mapsto G_{\overline{X}}$

Unterschied entsteht durch Einfluss gemeinsamer Einflussgrößen (confounder) im Originaldiagramm  $G$

## WH: do-Kalkül – Auswertung von Anfragen

Ziel:

Auswertung der Anfrage durch Elimination der *do*-Ausdrücke durch Anwendung von bedingten Regeln.

Jede Regel besteht aus

- ▶ Modifikation des Kausaldiagramms und
- ▶ syntaktische Umformung der Anfrage

Eliminierung der *do*-Ausdrücke abhängig von

- ▶ Kausaldiagramm  $G$  und
- ▶ Anfrage

Es gibt Kombinationen von Anfragen und Kausaldiagrammen, für die die Eliminierung der *do*-Ausdrücke nicht möglich ist.

# WH: do-Kalkül – Variablen

gegeben:

- ▶ Anfrage  $p(Y|do(x), Z, W)$
- ▶ Kausaldiagramm  $G = (V, E)$
- ▶ disjunkte Mengen  $X, Y, Z, W \subseteq V$ 
  - ▶  $Z$  Beobachtung, Stelle der Intervention (vermutete Ursache)
  - ▶  $Y$  zu untersuchende Variable (vermutete Wirkung)
  - ▶ Kontextvariablen:  
 $X$  (Beobachtung / Intervention möglich),  
 $W$  (unbeobachtet)

zur Erinnerung:

$(Y \perp\!\!\!\perp Z|X, W)_G$  gdw.

$X \cup W \subseteq V$  blockiert in  $G$  alle Pfade zwischen  $Y$  und  $Z$

# WH: do-Kalkül – Umformregeln

(bedingte) Regeln:

1. falls  $(Y \perp\!\!\!\perp Z | X, W)_{G_{\bar{X}}}$ :

$$p(Y | \text{do}(x), Z, W) = p(Y | \text{do}(x), W)$$

(Hinzufügen / Entfernen von Beobachtungen)

2. falls  $(Y \perp\!\!\!\perp Z | X, W)_{G_{\bar{X}, \underline{Z}}}$ :

$$p(Y | \text{do}(x), \text{do}(z), W) = p(Y | \text{do}(x), Z, W)$$

(Austausch von Aktion  $\text{do}(x)$  und Beobachtung  $x$ )

3. falls  $(Y \perp\!\!\!\perp Z | X, W)_{G_{\bar{X}, \overline{Z(W)}}$ :

$$p(Y | \text{do}(x), \text{do}(z), W) = p(Y | \text{do}(x), W)$$

wobei  $Z(W) = Z \setminus V_{G_{\bar{X}}}(W)$

(Menge aller Knoten in  $Z$ , die in  $G_{\bar{X}}$  keine gerichteten Vorfahren von  $W$  sind)

(Hinzufügen / Entfernen von Interventionen)

## Regel 1 – Analyse

Hinzufügen / Entfernen der Beobachtung  $Z$  möglich,  
falls  $Z$  keinen Einfluss auf  $Y$  hat

Bedingung:

$$(Y \perp\!\!\!\perp Z | X, W)_{G_{\bar{X}}}$$

graphisch im Kausaldiagramm:

nach Löschen aller Eingangskanten zu  $X$   
sind alle Pfade von  $Z$  nach  $Y$  blockiert

statistisch:

nach Intervention auf  $X$  ist  $Y$  (bedingt) unabhängig von  $Z$

Umformung:

$$p(Y | \text{do}(x), Z, W) = p(Y | \text{do}(x), Z)$$

Beobachtung  $Z$  kann (ohne Änderung der Bedeutung) entfernt  
oder hinzugefügt werden.

## Regel 1 – Spezialfälle

Falls  $(Y \perp\!\!\!\perp Z|X, W)_{G_{\bar{X}}}$ :  $p(y|do(x), z, w) = p(y|do(x), w)$

Spezialfälle:

- ▶  $X = W = \emptyset$ :  $G_{\bar{X}} = G$

Falls  $(Y \perp\!\!\!\perp Z)_G$  ( $Y$  unabhängig von  $Z$ ):  $p(y|z) = p(y)$

bekannt aus Statistik:  $Y$  unabhängig von  $Z$

- ▶  $X = \emptyset, W \neq \emptyset$  (passiv beobachtet):  $G_{\bar{X}} = G$

Falls  $(Y \perp\!\!\!\perp Z|W)_G$  ( $Y$  unter  $W$  bedingt unabhängig von  $Z$ ):

$$p(y|z, w) = p(y|w)$$

bekannt aus Statistik:  $Y$  bedingt (unter  $W$ ) unabhängig von  $Z$

- ▶  $W = \emptyset$  (unbekannt),  $X \neq \emptyset$ :

Falls  $(Y \perp\!\!\!\perp Z|X)_{G_{\bar{X}}}$  ( $Y$  bedingt unabhängig von  $Z$ ):

$$p(y|do(x), z) = p(y|do(x))$$

Statistik: nach starker Intervention  $X = x$  ist  $Y$  bedingt (unter  $W$ ) unabhängig von  $Z$

ÜA:

Analyse und Betrachtung der Spezialfälle für Regeln 2 und 3

## Regel 1 – Beispiel

$$G = \boxed{\text{Feuer}} \rightarrow \boxed{\text{Rauch}} \rightarrow \boxed{\text{Alarm}}$$

Bedingung für  $Y = \{A\}$ ,  $Z = \{F\}$ ,  $X = \{R\}$ :

$$(A \perp\!\!\!\perp F | R)_{G_{\bar{R}}}$$

ist erfüllt, denn in

$$G_{\bar{R}} = \boxed{\text{Feuer}} \not\rightarrow \boxed{\text{Rauch}} \rightarrow \boxed{\text{Alarm}}$$

existiert kein (nicht blockierter) Pfad von  $F$  zu  $A$

Statistik: Nach Intervention auf  $R$  ist  $A$  unabhängig von  $F$ .

Also ist Regel 1 anwendbar

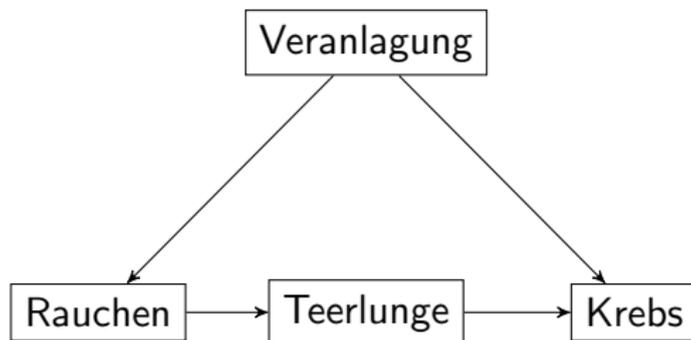
$$p(A|do(R), F) = p(A|do(R))$$

(Beobachtung von Feuer unter  $do(R)$  irrelevant)

## do-Kalkül – Beispiel

Tafel:

Transformation der Anfrage  $p(K|do(r))$   
zum Kausaldiagramm



# Was bisher geschah

- ▶ 3 Kausalitäts-Stufen (ladder of causality)
- ▶ Kausaldiagramme
- ▶ typische Teilgraphen (Pfade mit drei Knoten), je mit Bezeichnung des mittleren Knotens:
  - ▶ Kette, Mediator
  - ▶ Verzweigung, Confounder
  - ▶ v-Struktur, Collider
- ▶ blockierte (ungerichtete) Pfade in Kausaldiagrammen
- ▶ (bedingt) d-separierte Knoten
- ▶ Interventionen: hart, weich
- ▶ do-Kalkül: 3 Regeln
- ▶ Transformation von Anfragen in do-freie Anfragen (Antwort direkt aus statistischen Daten möglich)

# Counterfactuals

Counterfactual: Betrachtung der Form

Was wäre unter anderen Voraussetzungen (als den tatsächlichen) geschehen?

Beispiele:

- ▶ Hätte ich A auch bei doppeltem Preis gekauft?
- ▶ Hätte ich die Prüfung bestanden, wenn ich 2 Tage mehr in die Vorbereitung investiert hätte?
- ▶ Wäre A auch gestorben, wenn er vor fünf Jahren aufgehört hätte zu rauchen?
- ▶ Wäre A auch gestorben, wenn er eine Impfung bekommen hätte?
- ▶ Wäre A auch gestorben, wenn er keine Impfung bekommen hätte?

# Beispiel: Trolley-Probleme

nach ersten Folien in

[https:](https://imweb.imn.htwk-leipzig.de/~schwarz/forsch/talk-prw18.pdf)

[//imweb.imn.htwk-leipzig.de/~schwarz/forsch/talk-prw18.pdf](https://imweb.imn.htwk-leipzig.de/~schwarz/forsch/talk-prw18.pdf)

(einige) ethische Prinzipien:

**Utilitarismus:** Um Schaden zu vermeiden, darf  
(gesamtgesellschaftlich) geringerer Schaden für  
Unbeteiligte in Kauf genommen werden.  
(erfüllt in Bystander, Loop, Footbridge)

**double effekt:** Um Schaden zu vermeiden, darf Schaden für  
Unbeteiligte **nicht als Mittel** eingesetzt werden.  
(erfüllt in Bystander)

**triple effekt:** Um Schaden zu vermeiden, darf Schaden für  
Unbeteiligte **als Nebenwirkung** in Kauf genommen  
werden.  
(erfüllt in Bystander, Loop)

# Unterschied zwischen double und triple effect

Counterfactual zur Unterscheidung zwischen „Mittel“ und „Nebenwirkung“:

Hätte die Aktion den ursprünglichen Schaden auch vermieden, wenn die dadurch ausgelöste schädliche Wirkung (aufgrund einer anderen Ausgangsposition) nicht eingetreten wäre?

in den Trolley-Beispielen:

- ▶ Bystander: Hätte die Aktion (Weiche) beide auf dem Hauptgleis gerettet, wenn niemand auf dem Abstellgleis steht?
- ▶ Loop: Hätte die Aktion (Weiche) beide auf dem Hauptgleis gerettet, wenn niemand auf dem Umweg steht?
- ▶ Footbridge: Hätte die Aktion (Stoß) beide auf dem Hauptgleis gerettet, wenn niemand von der Brücke stürzt?

## Darstellung von Counterfactuals – Versuch

1. Versuch: Darstellung durch do-Ausdruck

kein Problem bei Interventionen auf  $S$ , also Anfragen der Form:  
 $p(D|\text{do}(S = 1))$ ,  $p(D|\text{do}(S = 0))$

aber Problem unter Betrachtung der tatsächlichen (beobachteten) Variable, z.B. Anfragen der Form:  $p(D|\text{do}(S = 0), D = 1)$

Unterscheidung nötig zwischen

tatsächlicher (beobachteter) Variable ( $D = 1$ )

hypothetischer (zu untersuchender) Variable ( $D' = 2$ )

Unterscheidung verschiedener „Welten“  $D$  und  $D'$  nötig

$p(D'|\text{do}(S = 0), D = 1)$

unter Bewahrung des Zusammenhanges zwischen den Welten  $D$  und  $D'$

(gemeinsame Abhängigkeit von anderen Variablen, Kontext  $U$ )

Darstellung durch do-Ausdrücke nicht möglich, **Versuch gescheitert**

# Darstellung von Counterfactuals

tatsächliche (beobachtete) Variable:  $Y$

hypothetische (zu untersuchende) Variable:  $Y_{X=x}$  (kurz  $Y_x$ )

im Beispiel:  $D' = D_{S=0}$

$$p(D_{S=0} | S = 1, D = 1)$$

( $S = 0$  gilt in der hypothetischen,  $S = 1$  in der tatsächlichen Welt)

## Counterfactuals als Anfragen

Anfrage der Form:  $Y_{X=x}(u) = y$

repräsentiert die Behauptung:

Würde in der Situation  $U = u$  (tatsächlicher Kontext)  $X = x$  gelten (statt tatsächlich  $X = x'$ ), dann wäre dort  $Y = y$ .

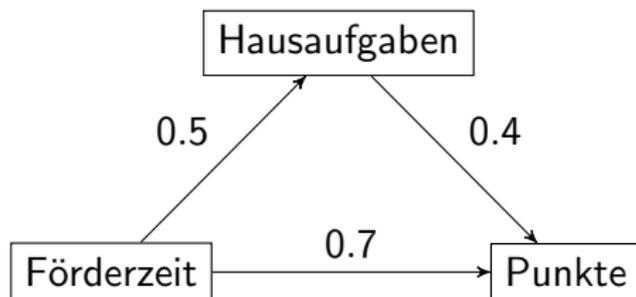
Idee:

- ▶  $u$  beschreibt die Merkmale einer speziellen Situation (Kontext, Individuum  $u$ )
- ▶  $X = x$  wird interpretiert als Anweisung zu minimaler Änderung des Modells  $M$  der aktuellen Situation (mit Eigenschaften  $u$ ) zu  $M_{X=x}$ , in dem  $X = x$  gilt (Intervention)

$$Y_{X=x}(u) = Y_{M_{X=x}}(u)$$

zusätzlich zum Kausaldiagramm notwendig: Information  $u$

## Strukturelles Gleichungsmodell – Beispiel



repräsentiert die Gleichungen:

$$F = U_F$$

$$H = aF + U_H$$

$$P = bF + cH + U_P$$

mit (Kantenbeschriftungen)  $a, b, c$ :

Parameter des (hier linear angenommenen) Zusammenhangs und fallspezifischem Kontext  $U_F, U_H, U_P$

Spezialfall (Individuum)  $e$ :  $F = 0.5, H = 1, P = 1.5$

# Auswertung von Counterfactuals

Anfrage  $Y_{X=x}(u) = y$

Auswertung in drei Schritten:

1. Abduktion:

Berechnung von  $U = u$  (unbeobachtet) aus Beobachtung des speziellen Falles  $E = e$

(Schritt in die Vergangenheit in der tatsächlichen Welt)

2. Modifikation des Modells zu  $M_x$ :

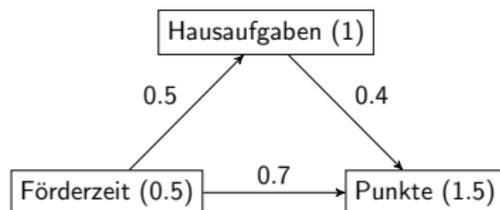
Ersetzen der Strukturgleichung für  $X$  durch  $X = x$  (Intervention)  
(Abbiegen in die hypothetische Welt)

3. Voraussage:

Berechnung von  $Y$  aus  $M_x$  und  $U = u$

(Schritt in die Gegenwart in der hypothetischen Welt)

# Beispiel



Wieviele Punkte hätte Joe (das hier untersuchte Individuum mit den beobachteten Merkmalen  $E = e$ ), wenn er doppelt so viel Zeit in die Hausaufgaben investiert hätte?

$$P_{H=2}(U) = ?$$

1. Abduktion von  $U$  aus  $e$ :  $F = 0.5, H = 1, P = 1.5$ :

$$U_F = F = 0.5, U_H = H - 0.5F = 1 - 0.25 = 0.75,$$

$$U_P = P - 0.7F - 0.4H = 1.5 - 0.7 \cdot 0.5 - 0.4 \cdot 1 = 1.5 - 0.35 - 0.4 = 0.75$$

2. modifiziertes Modell (Gleichungen):

$$F = 0.5$$

$$H = 2$$

$$P = 0.7F + 0.4H + 0.75$$

3. Auswertung der Anfrage (Voraussage):

$$P_{H=2}(U_F = 0.5, U_H = 0.75, U_P = 0.75) = 0.7 \cdot 0.5 + 0.4 \cdot 2 + 0.75 = 1.9$$

Die Studierenden sind in der Lage, Wissensrepräsentationen zur Modellierung zu benutzen, die über klassische Prädikatenlogik hinausgeht.

Insbesondere können sie dem Problem angemessene Wissensverarbeitungstechniken zur Simulation intelligenten Verhaltens auswählen.

Sie verstehen aktuelle Fachbeiträge und können die dort vorgestellten Ansätze verständlich aufbereiten und präsentieren.

# Prüfung

- ▶ Klausur 90 min  
am Dienstag, dem 27.07.2019 um 9:00-10:30 in N001-H
- ▶ Inhalt:
  - ▶ Vorlesungsinhalt
  - ▶ Aufgabentypen wie Übungsaufgaben
- ▶ Prüfungsvorleistung Beleg (PVB):  
≥ 2 MLV-Punkte
- ▶ (ausschließlich) zulässiges Hilfsmittel:  
A4-Blatt (beidseitig) handbeschrieben

# Hygienevorschriften zur Prüfung

aktuelle Informationen zum Hygienekonzept der HTWK

<https://www.htwk-leipzig.de/hochschule/aktuelles/coronavirus/hygienekonzept-der-htwk-leipzig>  
speziell zu Prüfungen:

<https://www.htwk-leipzig.de/studieren/im-studium/pruefungen/>

dort insbesondere: Mindeststandards für die Durchführung schriftlicher  
Präsenzprüfungen für Studierende  
(derzeitige) Highlights:

- ▶ ab 20 Minuten vor Prüfungsbeginn unter Einhaltung des (euklidischen) Mindestabstands von 1,50 m in Gebäude und Prüfungsraum gehen
- ▶ Vor Betreten des Prüfungsraumes QR-Code scannen
- ▶ Maske während der gesamten Prüfung auf
- ▶ Nachweis eines negativen qualifizierten Corona-Tests vorzulegen
- ▶ Formular gern ausgefüllt mitbringen, Unterschrift im Prüfungsraum:

[https://www.htwk-leipzig.de/fileadmin/portal/htwk/studieren/pruefungsplan/Formular\\_Corona\\_Praesenzpruefung.pdf](https://www.htwk-leipzig.de/fileadmin/portal/htwk/studieren/pruefungsplan/Formular_Corona_Praesenzpruefung.pdf)

- ▶ notfalls: Formular pandemiebedingter Prüfungsrücktritt

[https://www.htwk-leipzig.de/fileadmin/portal/htwk/studieren/download/Antrag\\_Corona\\_](https://www.htwk-leipzig.de/fileadmin/portal/htwk/studieren/download/Antrag_Corona_)

# Daten, Wissen, Intelligenz

Umwelt		Reize, Eindrücke
Agent	Wahrnehmen, Beobachten	Daten
	Erkennen, Verstehen	Information
	Anwenden, Können	Wissen
	Lernen	Wissenserwerb (Intelligenz?)
	Verstehen, Reflektieren, Begründen, Erklären, Erkennen der Grenzen	Intelligenz

# Daten, Information, Wissen, Intelligenz

**Daten** Darstellungsform (Syntax)  
Zeichenketten, Bilder, Ton, ...

**Information** Bedeutung der Daten (Semantik)  
in einem bestimmten Kontext

**Wissen** Information mit einem Nutzen,  
trägt zur Lösung eines Problemes bei,  
Nutzen abhängig von vorhandenem Kontextwissen

**Wissenserwerb** Erweiterung des Wissens in verschiedenen Kontexten  
durch logisches Schließen,  
Hinzufügen neuen Wissens durch Training

**Intelligenz** Fähigkeit zur Begründung, Erkennen eigener Grenzen,  
(vorherige) Einschätzen der Wirkung von Aktionen /  
Interventionen  
Reflexion, Retrospektive  
Schließen mit Counterfactuals:  
Betrachtung alternativer „Welten“  
(Was wäre eingetreten, wenn im gegebenen Kontext eine  
Eigenschaft nicht oder anders erfüllt gewesen wäre?)

# Lehrinhalte im Sommersemester 2021

- ▶ Einteilung symbolische / statistische KI
- ▶ Zustandsübergangssysteme (symb + evtl. stat)
  - ▶ Heuristische Suche
  - ▶ Spielbaum-Suche
- ▶ Künstliche neuronale Netze (stat)
  - ▶ mathematische Neuronenmodelle
  - ▶ Feed-Forward-Netze
  - ▶ rekurrente Netze
  - ▶ Assoziativspeicher
- ▶ kausales Schließen (symb + stat)

# Zustandsübergangssysteme

## Heuristische Suche

- ▶ Eigenschaften von Heuristiken
- ▶ Greedy-Suche
- ▶ Bestensuche
- ▶ A\*-Suche

## Spielbäume

- ▶ Expansion von Knoten in Spielbäumen
- ▶ MiniMax-Werte
- ▶  $\alpha$ - $\beta$ -Pruning

# Zustandsübergangssysteme – Wissensrepräsentation

Formalisierung des **Wissens**:

Problembeschreibung (Modellierung) durch  
Zustandsgraphen:

- ▶ Knoten: Zustände
- ▶ Kanten: zulässige Übergänge zwischen Zuständen
- ▶ ausgezeichnete Startzustände
- ▶ Eigenschaften der Zielzustände

mögliche gewünschte **Lösung**:

- ▶ ein Zielzustand (bei definierten Zielbedingungen)
- ▶ alle Zielzustände
- ▶ Zielzustände mit guter Bewertung (Zusatzinformation)
- ▶ Pfade zu Zielzuständen (Pläne, Strategien)

Anwendungen, z.B.

- ▶ Verifikation
- ▶ kombinatorische Suchprobleme
- ▶ Planen
- ▶ Spielsituationen

# Zustandsübergangssysteme: Wissensverarbeitung

Suchverfahren in Zustandsgraphen:

**Tiefensuche** (uninformiert)

Verwaltung entdeckter, aber noch nicht abgearbeiteter Knoten in Stack

**Breitensuche** (uninformiert)

Verwaltung entdeckter, aber noch nicht abgearbeiteter Knoten in Queue

**heuristische Suche**

von Lösungen mit zusätzlichen (Optimalitäts-) Eigenschaften oder zur Beschleunigung der Suche  
Verwaltung entdeckter, aber noch nicht abgearbeiteter Knoten in Priority-Queue mit geeignet gewählten Prioritäten

- ▶ Bestensuche (nur bisherige Kosten)
- ▶ Greedy (nur Heuristik)
- ▶  $A^*$  (bisherige Kosten + Heuristik)

**Spielbaum-Suche** (Zwei-Personen-Spiele)

Minimax-Suche,  $\alpha$ - $\beta$ -Pruning

# Maschinelles Lernen

- ▶ überwacht
  - ▶ korrigierend
  - ▶ bestärkend (reinforcement)
- ▶ unüberwacht

# Künstliche Neuronen

- ▶ biologisches Vorbild
- ▶ mathematisches Modell
- ▶ Eingangs-, Aktivierungs-, Ausgangsfunktion
- ▶ McCullochs-Pitts-Neuron
- ▶ Schwellwert-Neuron
- ▶  $\Delta$ -Lernregel
- ▶ Faltungs-Neuron
- ▶ ...
- ▶ geometrische Interpretationen

# Künstliche Neuronale Netze

- ▶ Schichten-Struktur
- ▶ Ein-, Mehr-Schicht-FFN
- ▶ Cognitron, CNN
- ▶ rekurrente Netze
- ▶ Assoziativspeicher: BAM, Hopfield-Netz
- ▶ Lernverfahren / Training
- ▶ Anwendungen

# Kausales Schließen

Einführung, angelehnt an Judea Pearl: Book of Why

- ▶ WH: Bayes-Netze
- ▶ Kausaldiagramme
- ▶ d-Separation in Kausaldiagrammen
- ▶ do-Kalkül
- ▶ Schließen mit Counterfactuals