
Einführung in Web- und Data-Science

Prof. Dr. Ralf Möller

Universität zu Lübeck

Institut für Informationssysteme

Tanya Braun (Übungen)



Teilnehmerkreis und Voraussetzungen

Studiengänge

- Bachelor **Informatik**
- Bachelor **IT-Sicherheit**
- Bachelor **Mathematik in Medizin und Lebenswissenschaften**
- Bachelor **Medizinische Informatik**
- Bachelor **Medieninformatik**
- Bachelor **Robotik und Autonome Systeme**

Voraussetzungen

- Keine



Organisatorisches: Übungen

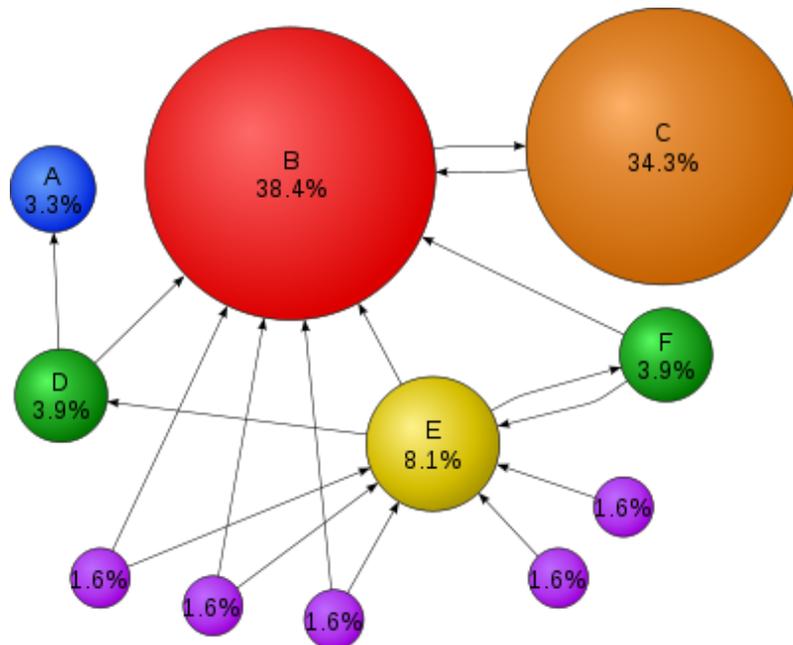
- **Vorlesung:** Donnerstags, ab dem 19. Oktober 2017
- **Übungen:** Mittwochs 13-14 Uhr, ab 25.10.
Anmeldung über Moodle nach dieser Veranstaltung
- **Übungsaufgaben** stehen jeweils nach der Vorlesung ca. ab 18 Uhr über Moodle bereit (erstes Übungsblatt erscheint am 19.10.2017)
- **Abgabe der Lösungen** erfolgt bis **Montag 12 Uhr** in der IFIS-Teeküche (jeweils in der zweiten Woche nach der Ausgabe in den Kasten für die jeweilige Übungsgruppe)
- Aufgaben sollen in einer **2-er Gruppe** bearbeitet werden (also bitte Name(n), Matrikelnummer(n) und Übungsgruppennummer vermerken)
- In den Übungen am Mittwoch wird der Übungszettel besprochen, dessen Lösungen bis zum jeweils vorigen Montag abgegeben werden, und auch Fragen zum jeweils neuen Übungszettel geklärt.

Organisatorisches: Prüfung

- Die **Eintragung in den Kurs** und in eine Übungsgruppe ist **Voraussetzung**, um an dem Modul Non-Standard-Datenbanken teilnehmen zu können
- Am Ende des Semesters findet eine **Klausur** statt
- **Voraussetzung** zur Teilnahme an der Klausur sind mindestens **50% der gesamtöglichen Punkte aller Übungszettel**

Web und Data Science

- Web Science
 - Analyse von Strukturen im Web (Mensch und Computer)
 - Formalisierung durch große Graphstrukturen und entsprechende Entscheidungsprobleme über Graphen
 - Beispiel: Pagerank (Bewertung von Webseiten)



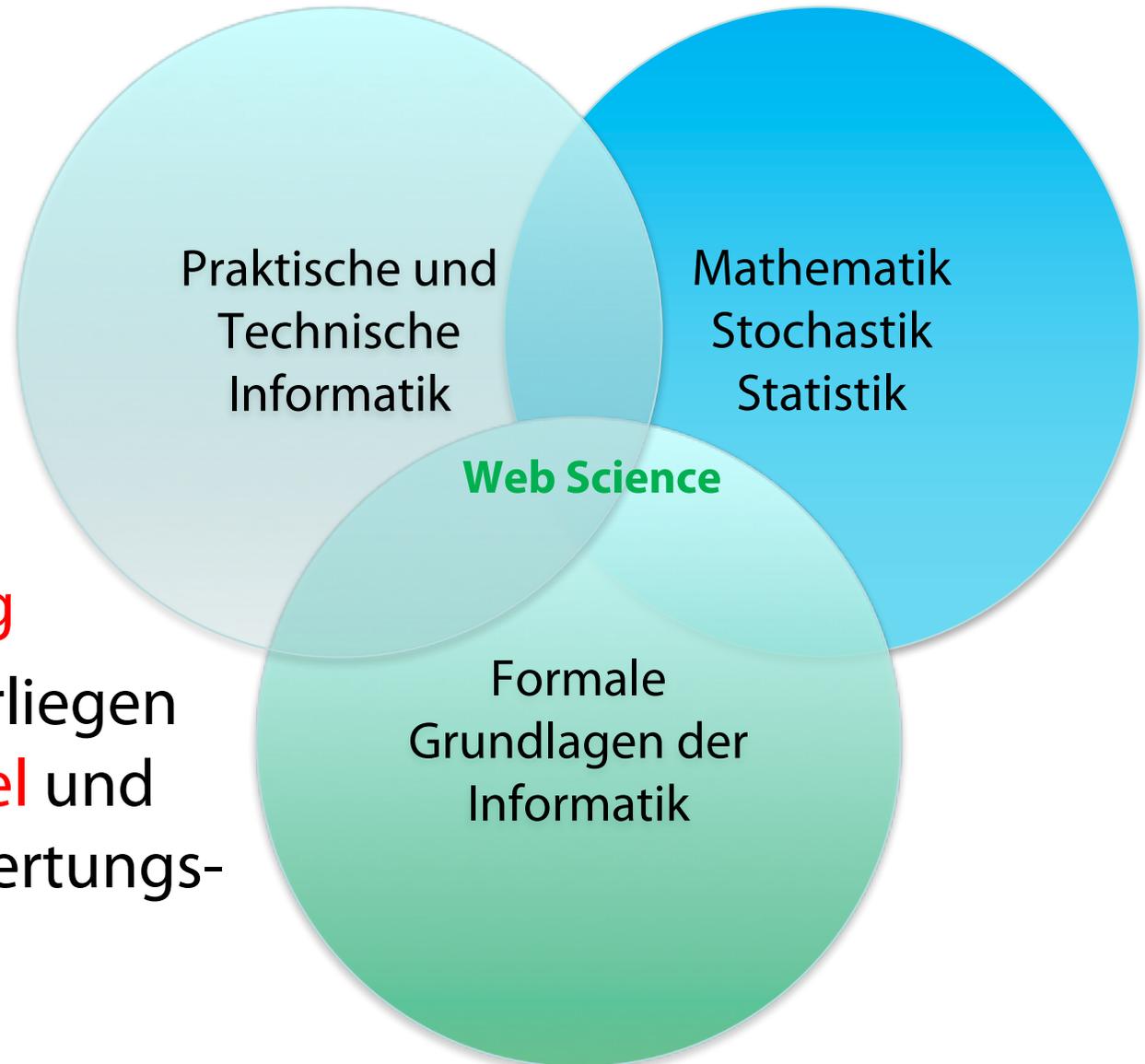
Zufallssurfer-Modell:

Größe der Kreise in etwa proportional der relativen Häufigkeit, mit der sich ein Surfer auf einer Seite befindet

[Wikipedia]

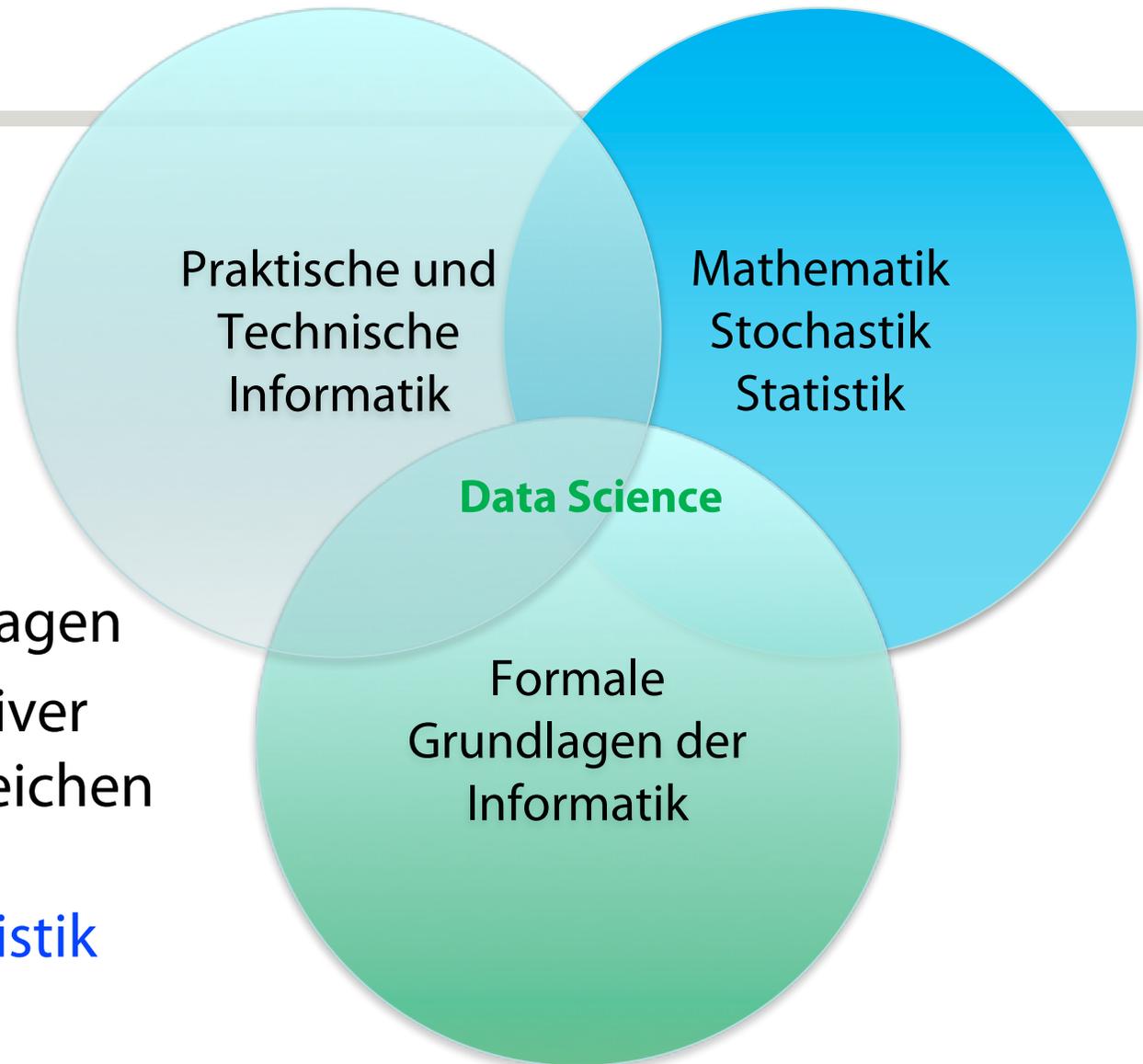
Herausforderungen für die Informatik

- Graphstrukturen extrem **groß**
- Verfahren zur Lösung von Entscheidungsproblemen extrem **aufwendig**
- Graphdaten unterliegen ständigem **Wandel** und so auch die Auswertungsergebnisse



Data Science

- Extraktion von Wissen aus Daten (u.a. Graphdaten)
- Begriff schon vor 50 Jahren für Informatik vorgeschlagen
- Entwicklung innovativer Konzepte in den Bereichen **Logik, Datenbanken und Stochastik / Statistik** (Datenanalyse und Wissensentdeckung)
- Verwendung von **LADS und Analysis**



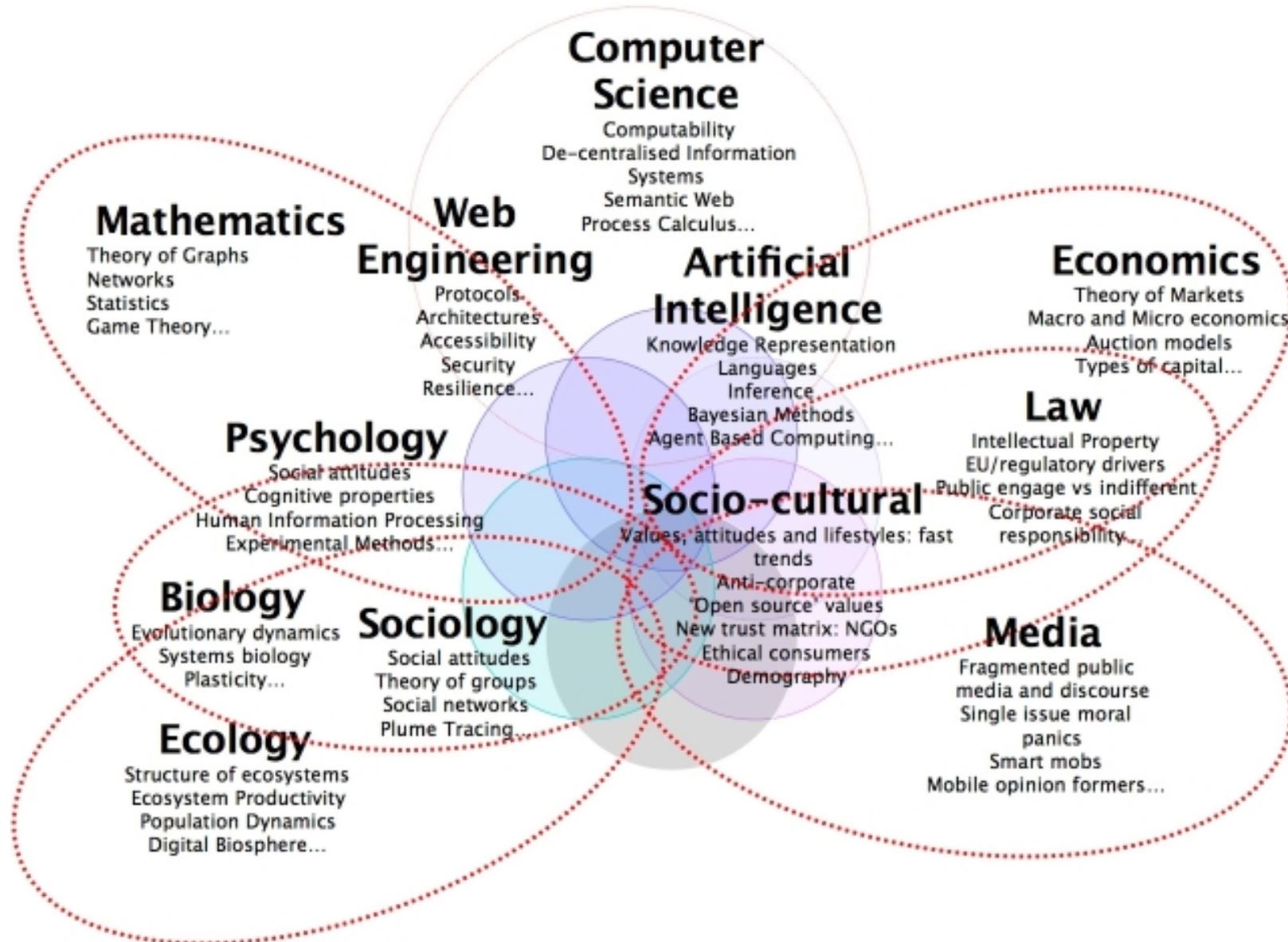
Herausforderungen für die Informatik

- Große Datenbestände
 - Speicher und Zugriffstechnologie
- Starker Zuwachs an Daten, hohe Dynamik
 - Hohe Datenraten und Echtzeitanforderungen
- Heterogene Datenbestände
 - Verteiltes Datenmanagement
 - Datenintegration
- Robuste Modelle der menschlichen Interpretationsvorgänge (Kognition) notwendig für Auswertung und Präsentation von Ausgaben
 - Verarbeitung von Text- und Graphikdaten nicht nur syntaktisch
 - Sprach und Videodatenverarbeitung in den Kinderschuhen

• ...

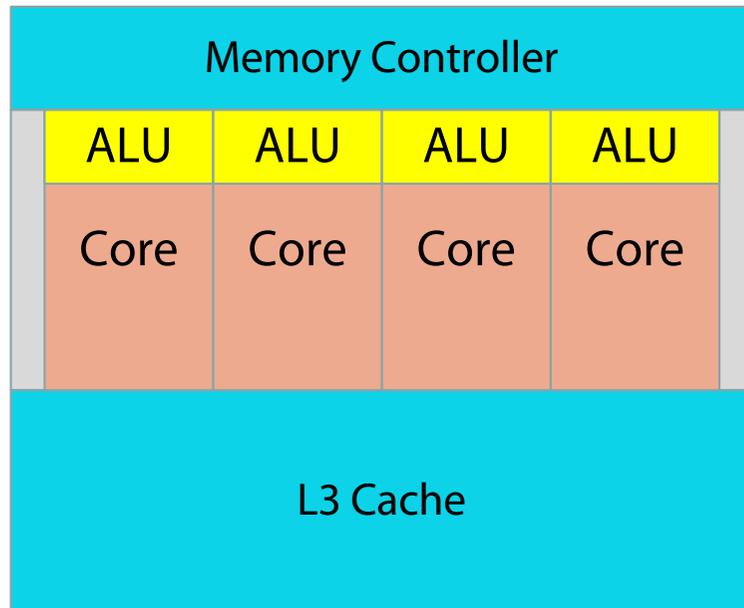


Web und Data Science: Sozio-technische Sicht

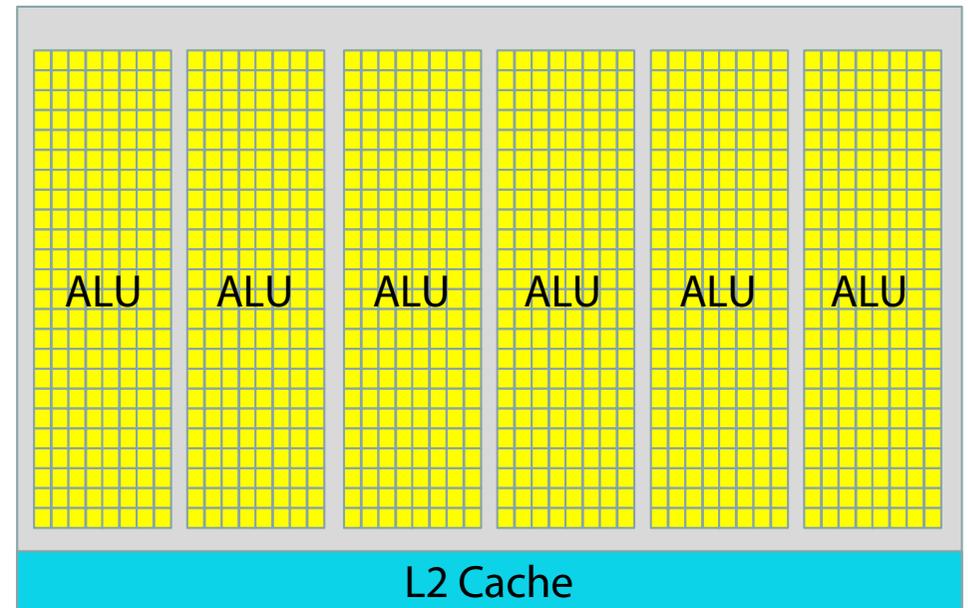


Hardware-Sicht: CPU vs GPU

Intel® CPU



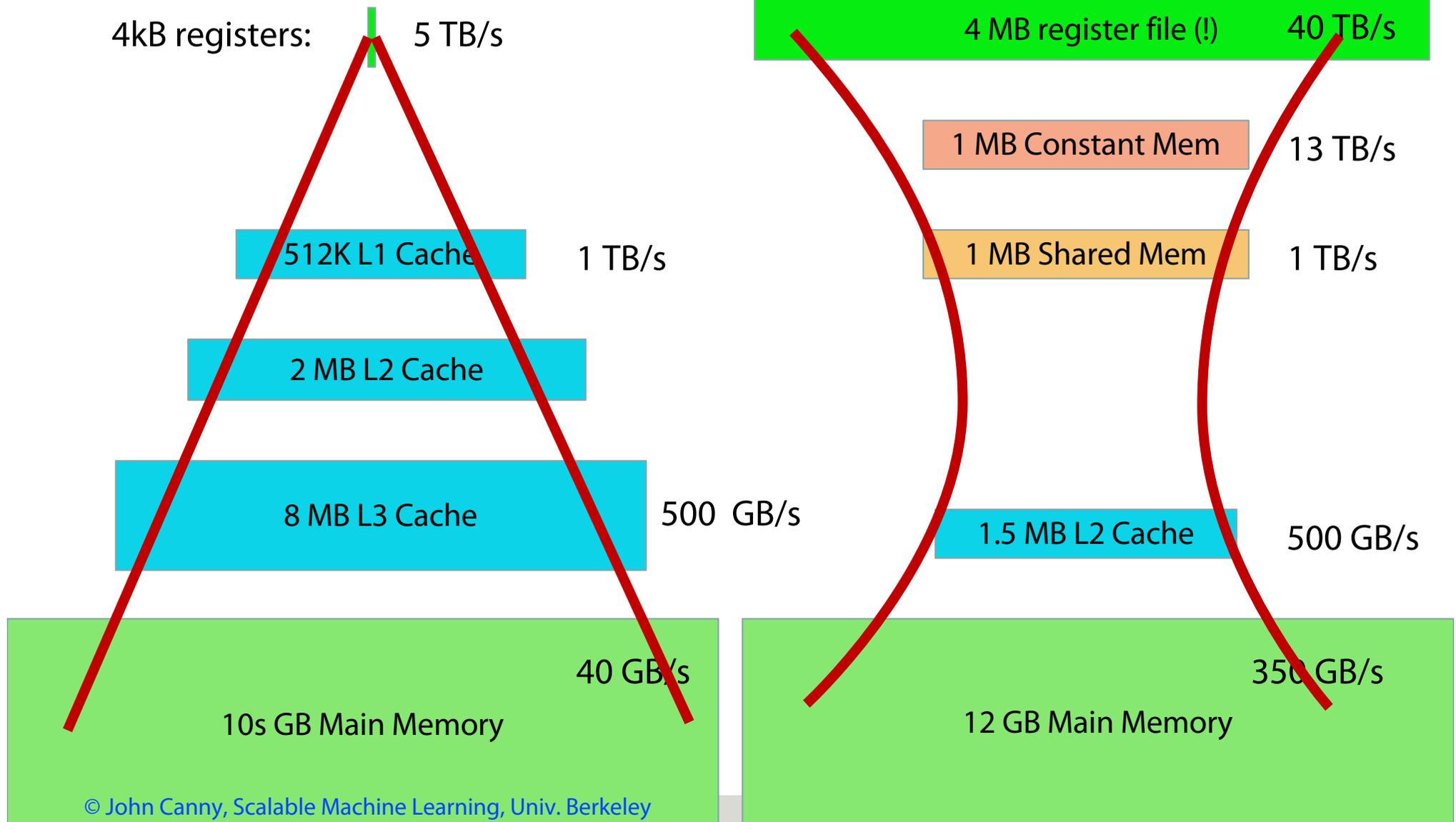
NVIDIA® GPU



Hardware-Sicht: CPU vs GPU

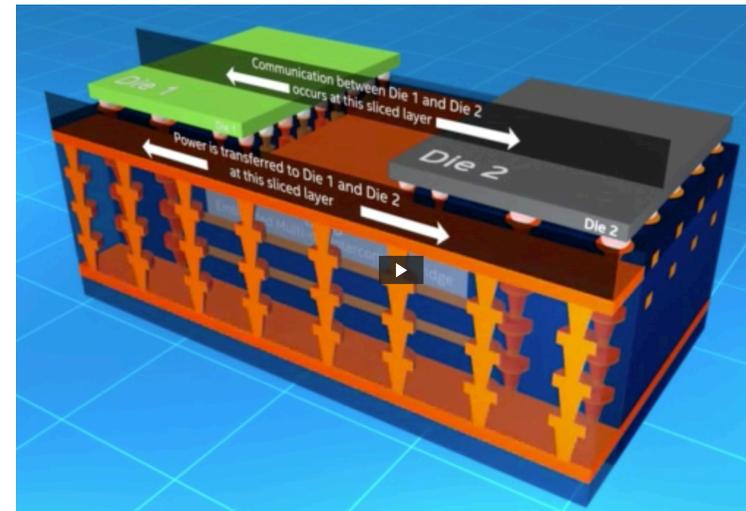
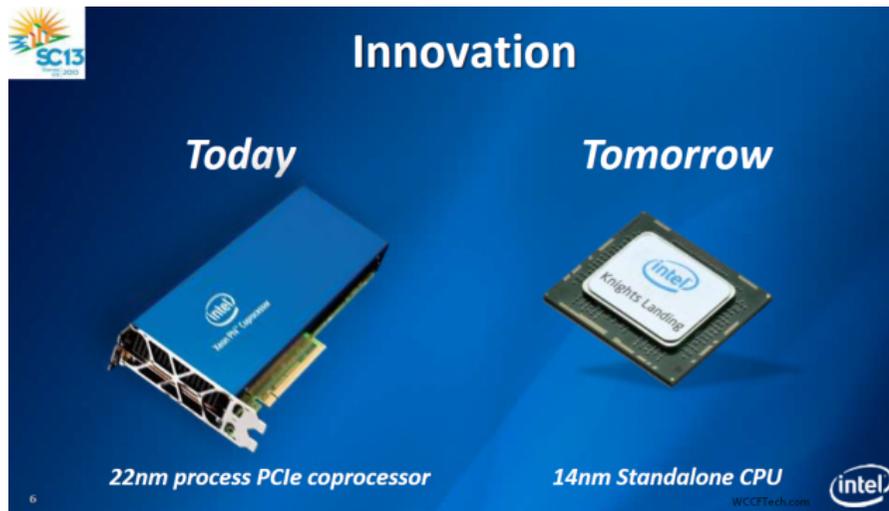
Intel® 8 core Sandy Bridge CPU

NVIDIA® GK110 GPU



Neue Hardwarestrukturen

- Heute: CPU (mit n Kernen) und externem Co-Prozessor (PCI-Bus)

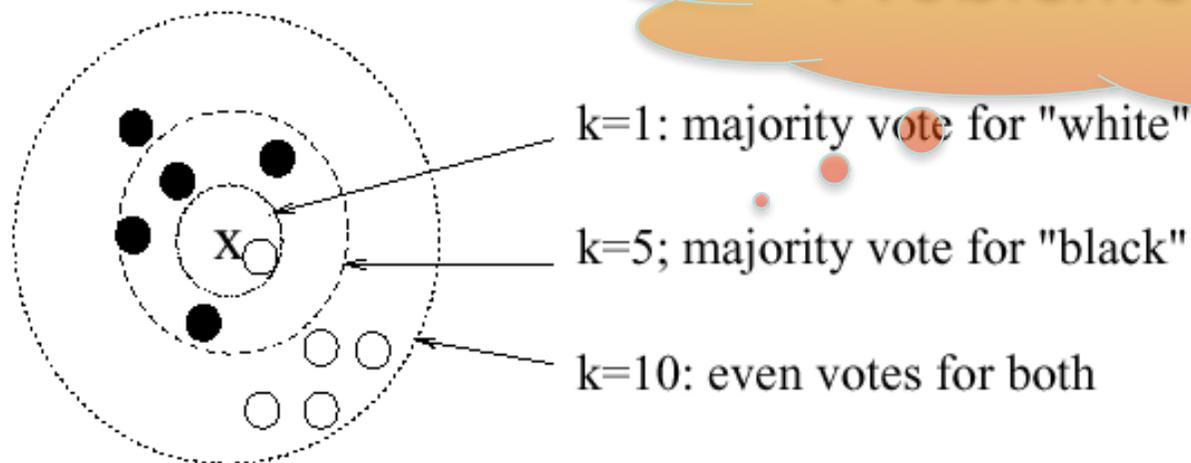


- Morgen: CPU und GPU integriert
 - z.B. Intel "Knight's Landing" Prozessor
- Morgen: CPU und FPGA eng verzahnt
 - Embedded Multi Die Interconnect Bridge
 - Dynamisch rekonfigurierbare Datenfluss-Hardware für spezielle Aufgaben

Speicheranforderungen für Anwendungen

- Annahme: Gegeben viele Datenpunkte
 - Beispielmerkmale: (x, y, Farbe) , $\text{Farbe} \in \{ \text{white, black} \}$
- Anfrage: Datenpunkt ohne Wert für bestimmtes Merkmal
 - Beispiel: Merkmal Farbe ohne Wert
- Anfragebeantwortung (Klassifikation des Anfragepunkts):
Mehrheitsvotum der k-nächsten Nachbarn (kNN-Verfahren)

Probleme erkennbar?



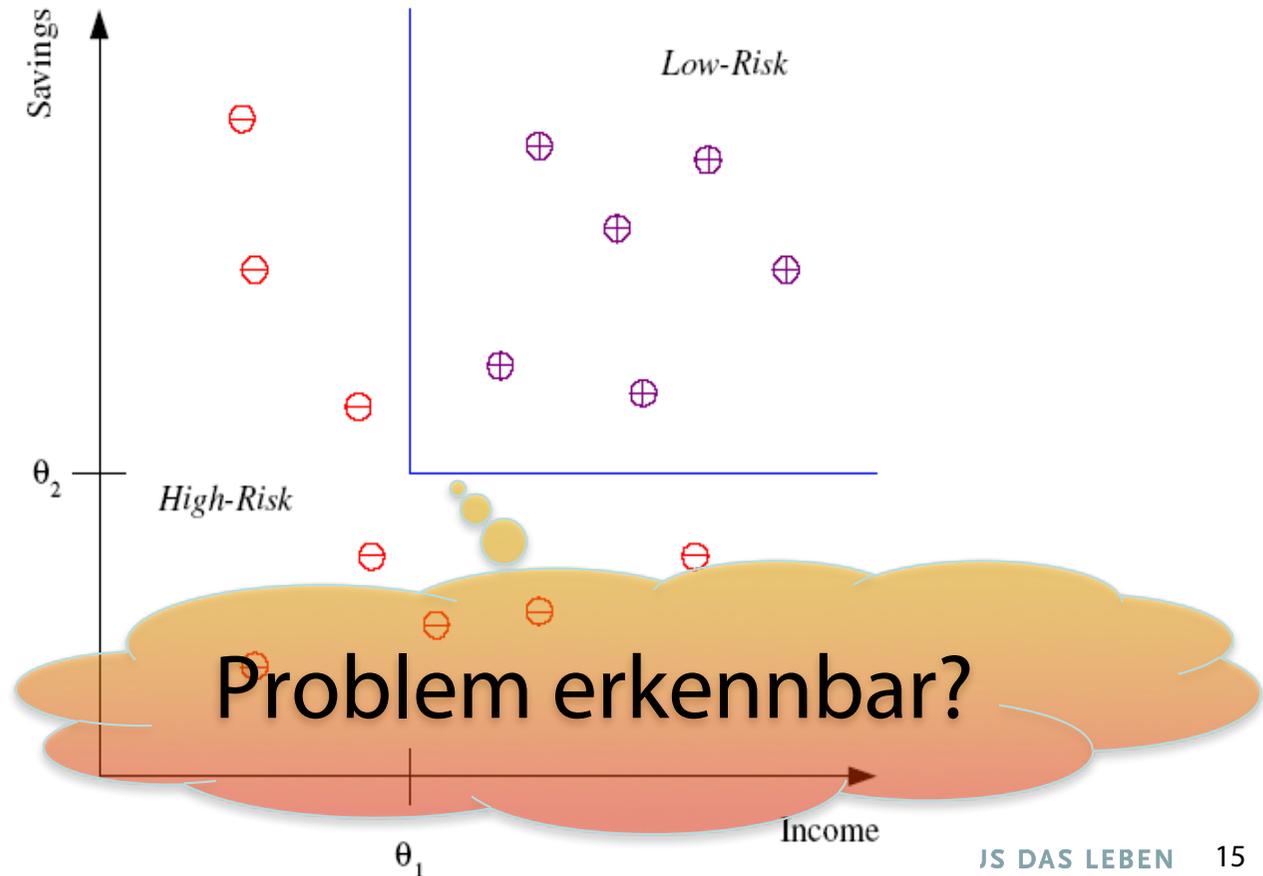
Fix, E., Hodges, J.L. Discriminatory analysis, nonparametric discrimination:
Consistency properties. Technical Report 4, USAF School of Aviation
Medicine, Randolph Field, Texas, 1951

Probleme mit kNN

- Klassifikationsergebnis stark von k abhängig
- Hoher Speicherbedarf
- Effizienter Zugriff auf "Nachbarn" erfordert weitere Maßnahmen (noch mehr Speicherbedarf)

Generalisierung von Daten möglich?

- Repräsentation der Daten durch Parameter
 - Wenn $(\text{Einkommen} > \theta_1 \wedge \text{Ersparnisse} > \theta_2)$, dann kreditwürdig (\oplus), sonst nicht (\ominus)
- Nur 2 Parameter nötig: (θ_1, θ_2)
- Modell fordert geringen Speicher

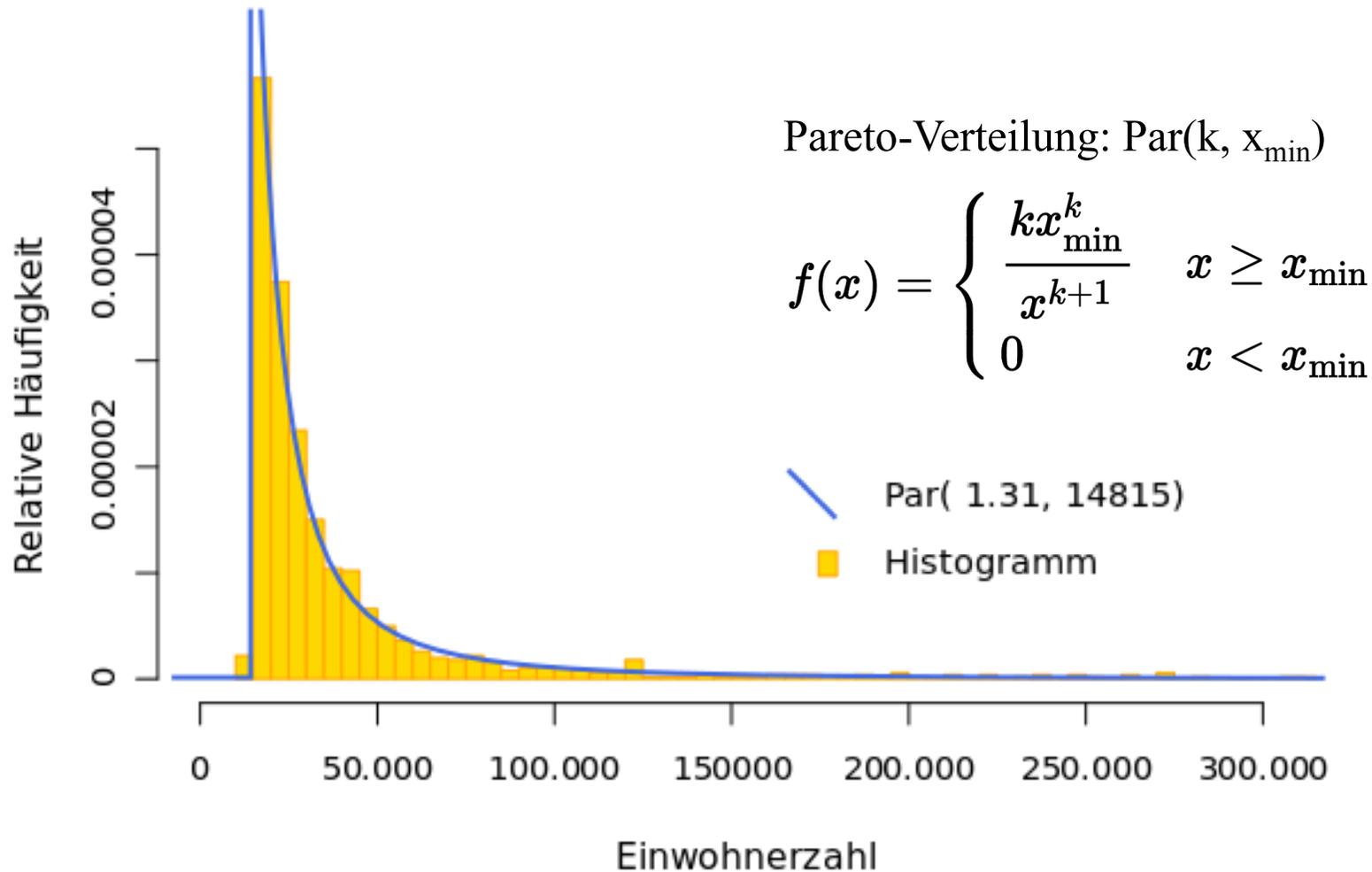


Beispiel

- Anzahl von Städten mit bestimmten Einwohnerzahlen schätzen
- Daten: Liste von Einwohnerzahlen (einige kommen mehrfach vor)
- Explizites Modell: Zählerfeld aufbauen
 - Zählerfeld ggf. **sehr groß**
- Implizites Modell: Potenzgesetz $y = ax^b$
 - a und b bestimmen
 - **Aufwendiges** Optimierungsproblem lösen

Begriff der "Verteilung"

Einwohnerzahlen Deutscher Städte



Rechnen auf realistischen Daten

Unterstützung von **großen Modellen** durch Cluster-basiertes maschinelles Lernen

Daten Parallelismus: Daten auf Clusterknoten aufgeteilt, jeder Knoten hat lokale Kopie des geteilten Modells

Modell-Parallelismus: Modell auf verschiedene Knoten verteilt. Keiner verarbeitet gesamte Daten,

In beiden Fällen gilt:
Knoten müssen kommunizieren,
um Modell zu bestimmen



Data Models vs. Algorithmic Models

Data Modeling

vs.

Algorithmic Modeling

$Y \leftarrow F(X, \text{random noise, parameters})$

$Y \leftarrow$ **Black Box** $\leftarrow X$
Random Forests

We understand the world

How well 'my data model' works
Statisticians, Data Analysts, Data Miners
Linear Regression
Logistic Regression
Known Distributions
Confidence Intervals
Predictor Variables & Goodness of Fit

We don't understand the world

The world produces data in a black-box
Data Scientists
Machine Learning, AI & Neural Nets
Random Forests, SVM, GBT
Unknown Multivariate Distributions
Iterative
Predictive Accuracy

"Statistical Modeling: The Two Cultures" Leo Breiman, 2001

Über den Kurs

- Repräsentationssprachen
 - Modellierungsansatz
 - Entscheidungs- und Berechnungsprobleme
 - Verfahren zu deren Lösung
- Gewinnung von Modellen aus Daten als Optimierungsproblem (Algorithm. Datenanalyse, Data Mining, Maschinelles Lernen)
 - Überwachtes Lernen,
 - Deep Learning, Transduktives Lernen, Reinforcement-Lernen
 - Unüberwachtes Lernen
 - Transfer-Lernen
- Verwendung von Modellen in Anwendungen
 - Prädikation, Forensik, Verstehen der Realität (Science-Aspekt!)
 - Autonome Aktivitäten/ Agenten

Fach-Sem.	KP	Kernbereich Informatik	101	13	Mathematik	32	5	F	Web and Data Sc.	36	7
1	30	CS1000 Einführung in die Programmierung	10	1	MA1000 Lin. Algebra u. Disk. Strukturen 1	8	1		CS1800 Efg in Web und Data Science	4	1
					MA2000 Analysis 1	8	1				
2	30	CS1001 Algorithmen und Datenstrukturen	8	1	MA1500 Lin. Algebra u. Disk. Strukturen 2	8	1				
		CS1201 Technische Grundlagen 1	6	1	MA2500 Analysis 2	4	1				
3	30	CS1202 Technische Grundlagen 2	6	1				CS	CS3050 Codierung und Sicherheit	4	1
		CS1002 Einführung in die Logik	4	1							
		CS2000 Theoretische Informatik	8	1							
		CS2300 SW-Engineering	6	1							
4	30	CS2301 SW-Engineering Praktikum	6	0	MA2510 Stochastik 1	4	1		MA3110 Numerik I	4	1
		CS2100 Rechnerarchitektur	4	1							
		CS2150 Betriebssysteme und Netze	8	1							
		CS2700 Datenbanken	4	1							
5	33	CS3000 Algorithmen-Design	4	1				CS			
		CS3010 Mensch-Computer-Interaktion	4	1				CS			
6	27	Wahlpflicht Informatik	8	1-2					CS3051 Parallelverarbeitung	4	1
		CS3990 Bachelorarbeit Informatik Kolloquium	12 3								

Bachelor



Fach-Sem.	KP	Kernbereich Informatik	101	13	Mathematik	32	5	Fachüberggr. Bereich	11	0	Web and Data Sc.	36	7
1	30	CS1000 Einführung in die Programmierung	10	1	MA1000 Lin. Algebra u. Disk. Strukturen 1	8	1				CS1800 Efg in Web und Data Science	4	1
					MA2000 Analysis 1	8	1						
2	30	CS1001 Algorithmen und Datenstrukturen	8	1	MA1500 Lin. Algebra u. Disk. Strukturen 2	8	1				CS3050 Codierung und Sicherheit	4	1
		CS1201 Technische Grundlagen 1	6	1	MA2500 Analysis 2	4	1						
3	30	CS1202 Technische Grundlagen 2	6	1				CS2450 Werkzeuge für wiss. Arbeiten	2	0	MA3110 Numerik I	4	1
		CS1002 Einführung in die Logik	4	1									
		CS2000 Theoretische Informatik	8	1									
		CS2300 SW-Engineering	6	1									
4	30	CS2301 SW-Engineering Praktikum	6	0	MA2510 Stochastik 1	4	1				CS3130 Algorithmische Datenanalyse	8	1
		CS2100 Rechnerarchitektur	4	1							CS3100 Signalverarbeitung	8	1
		CS2150 Betriebssysteme und Netze	8	1									
		CS2700 Datenbanken	4	1									
5	33	CS3000 Algorithmen-Design	4	1				CS					
		CS3010 Mensch-Computer-Interaktion	4	1				CS			CS3204 Künst. Intel. 1 oder CSxxx Logik für Inform.	4	1
6	27	Wahlpflicht Informatik	8	1-2									
		CS3990 Bachelorarbeit Informatik Kolloquium	12 3										

Bachelor

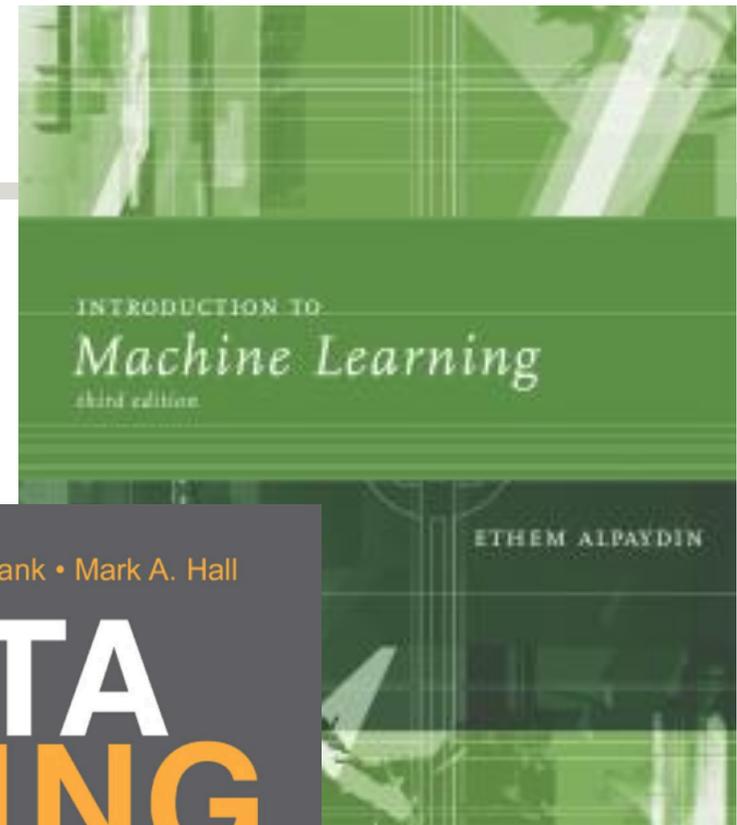
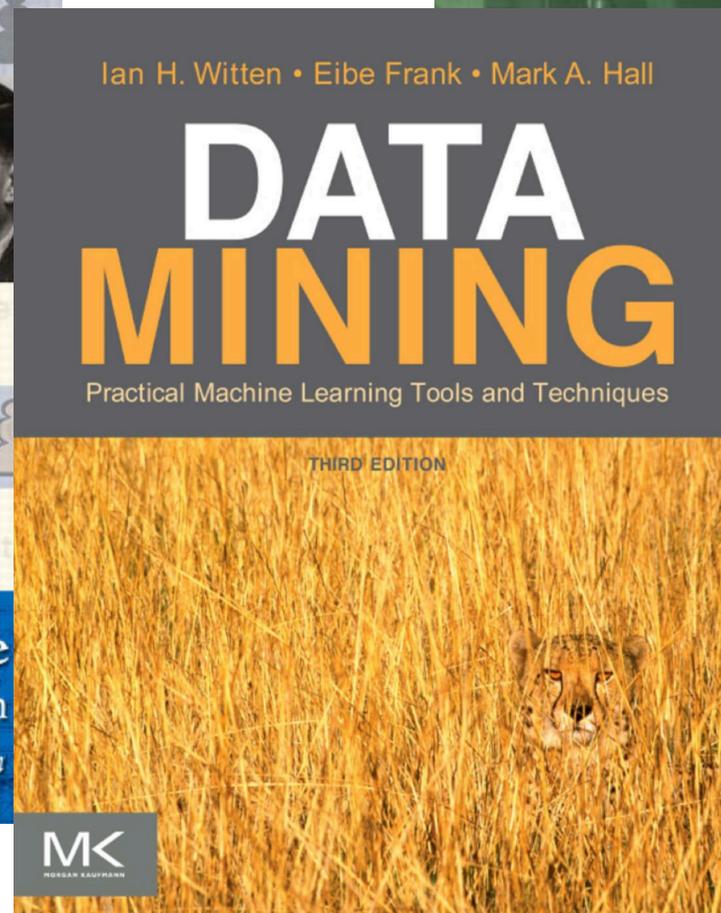
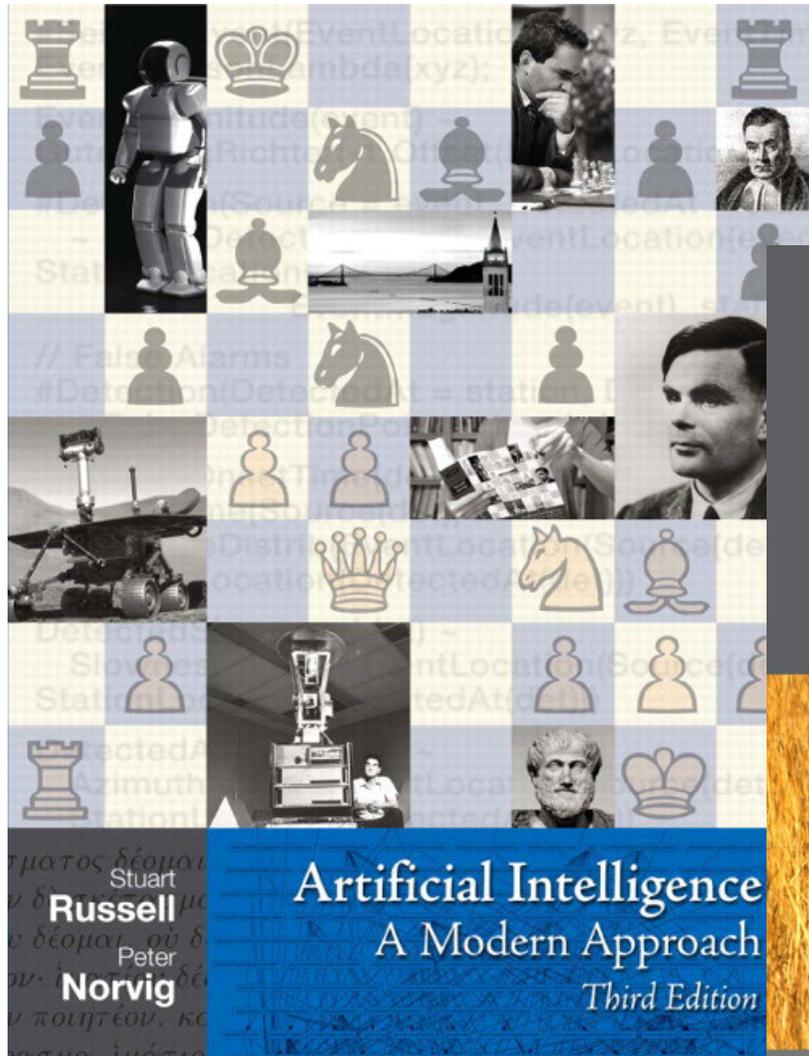


Fach-Sem.	Kernbereich Informatik	KP Prüf	Fachübe	Web and Data Science	KP Prüf
1	Basismodul Theoretische Informatik CS4000 Algorithmik (WS) oder CS4020 Spezifikation und Modellierung (SS)	6 1	mindest		
	Basismodul Praktische Informatik CS4130 Webbasierte Informationssysteme (SS) oder CS4150 Verteilte Systeme (WS)	6 1		CS5131 Web and DB Mining Agents (WS)	8 1
	Basismodul Technische Informatik CS4160 Echtzeitsysteme (WS) oder CS4170 Parallelrechnersysteme (SS)	6 1			
2 + 3	3 Vertiefungsmodule a 12 KP aus folgender Liste	36 3	CS5840 er (W	MA4030	Vert. Th. Web and Data Science: CS4250 Computer Vision MA4610 Stochastische Prozesse CS5275 Ausgewählte Methoden der Signalanalyse und Verbesserung MA4940 Test- und Schätztheorie MA4341 Zeitreihenanalyse MA4040 Numerik 2
	CS4501 Algorithmik, Logik und Komplexität CS4502 Parallele und verteilte Systeme CS4503 Ambient Computing und Anwendungen CS4504 Cyber Physical Systems CS4505 Systemarchitektur CS4506 Datensicherheit CS4507 Softwareverifikation CS4508 Datenmanagement CS4509 Internet-Technologien CS4510 Signalanalyse CS4511 Lernende Systeme CS4512 Bildgeb. Systeme und inverse Probleme CS4520 Fallstudie zur prof. Produktentwicklung		EC4001 Al (W EC4010 W (W PS5810 W Lehrtätigk PS5830 St. Business		
				CSxxx Projektprakt. Web and Data Science	4 0

Master

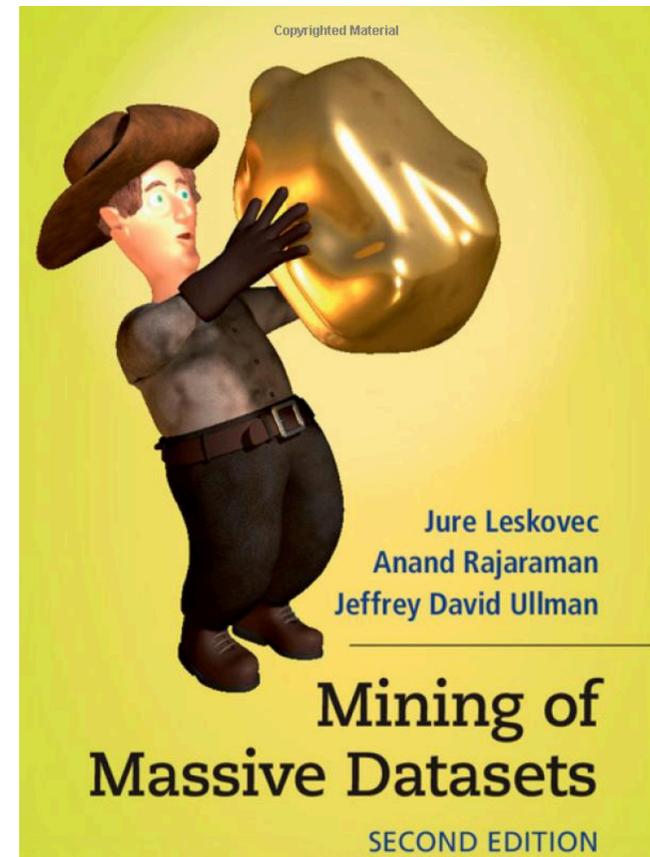


Literatur



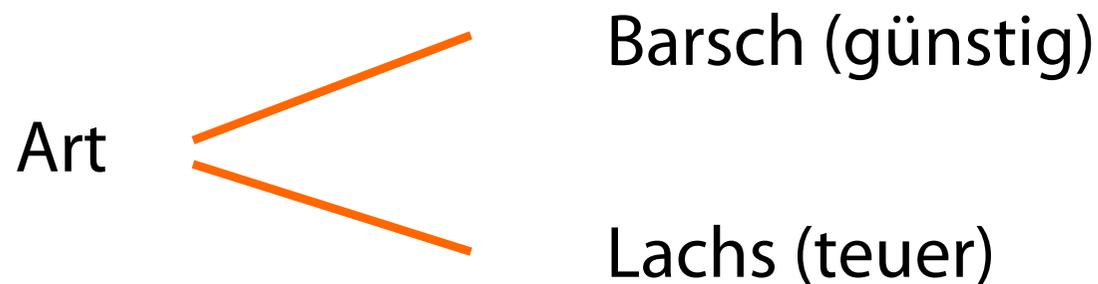
Literatur

- Stuart Russell, Peter Norvig, **Artificial Intelligence – A Modern Approach**, Pearson, 2009 (oder 2003er Ed.)
- Ian H. Witten, Eibe Frank, Mark A. Hall, Data Mining: **Practical Machine Learning Tools and Techniques**, Morgan Kaufmann, 2011
- Ethem Alpaydin, **Introduction to Machine Learning**, 3rd Ed., MIT Press, 2014
- Jure Leskovec, Anand Rajaraman, Jeffrey D. Ullman, **Mining of Massive Datasets**, 2nd Ed., Cambridge University Press, 2014
- Viele zusätzliche Bücher, Präsentationen, und Videos im Web



Ein erweitertes Beispiel

“Sortierung von Fischen auf einem Förderband nach Arten durch Bildverarbeitung”



Problemanalyse

Verwende Kamera und nehme Bilder auf,
um Merkmale zu bestimmen:

- Länge
- Helligkeit
- Breite
- Anzahl und Form der Flossen
- Position des Mundes usw.

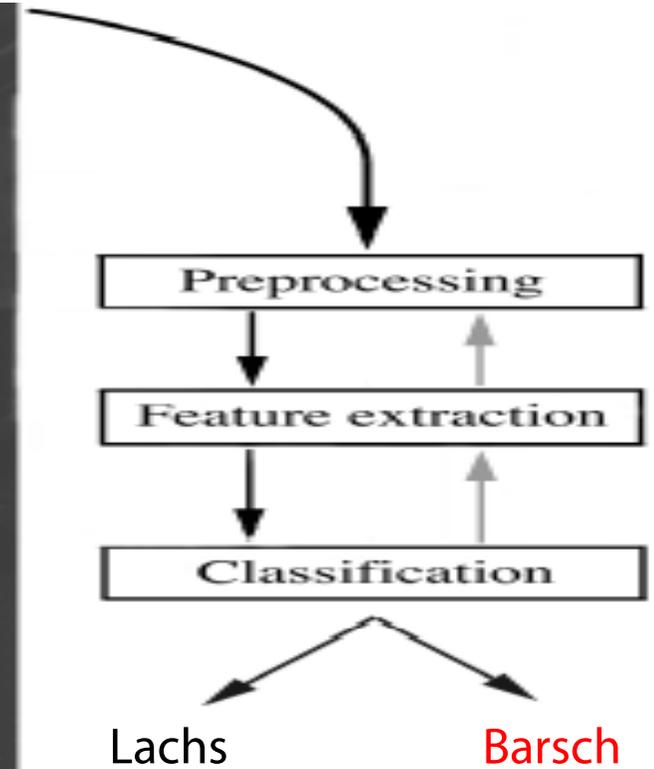
Menge aller möglichen Merkmale

Ziel: Wähle die relevanten aus

Vorverarbeitung

- Verwende Segmentierungsoperator, um Fisch und Hintergrund zu unterscheiden
- Fischregion wird verwendet, um Merkmalswerte zu extrahieren (geeignete Merkmale vorher festgelegt)
- Merkmalswerte weiterverarbeitet durch Klassifikator (Barsch oder Lachs)

Klassifikation



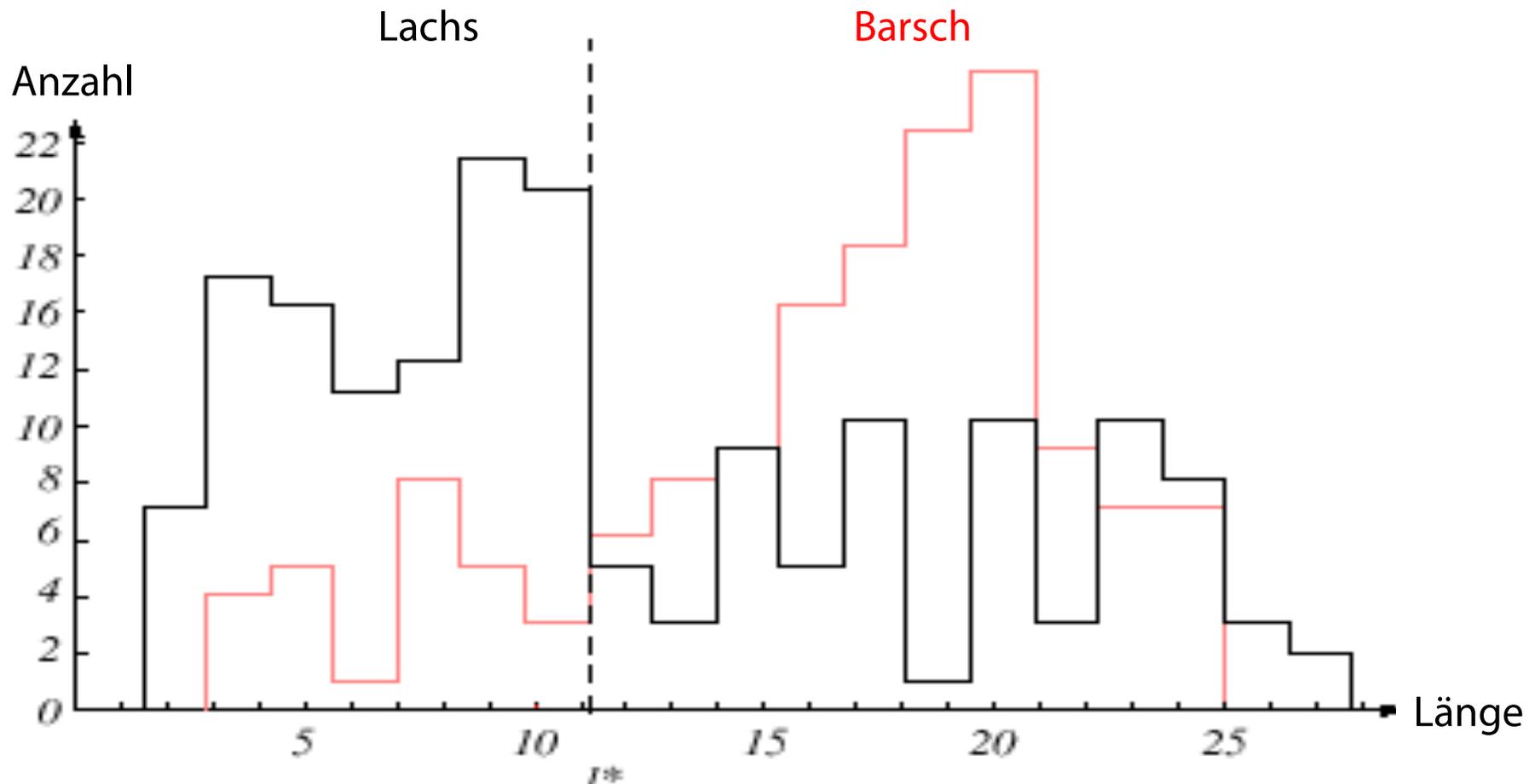
Bestimmung geeigneter Merkmale

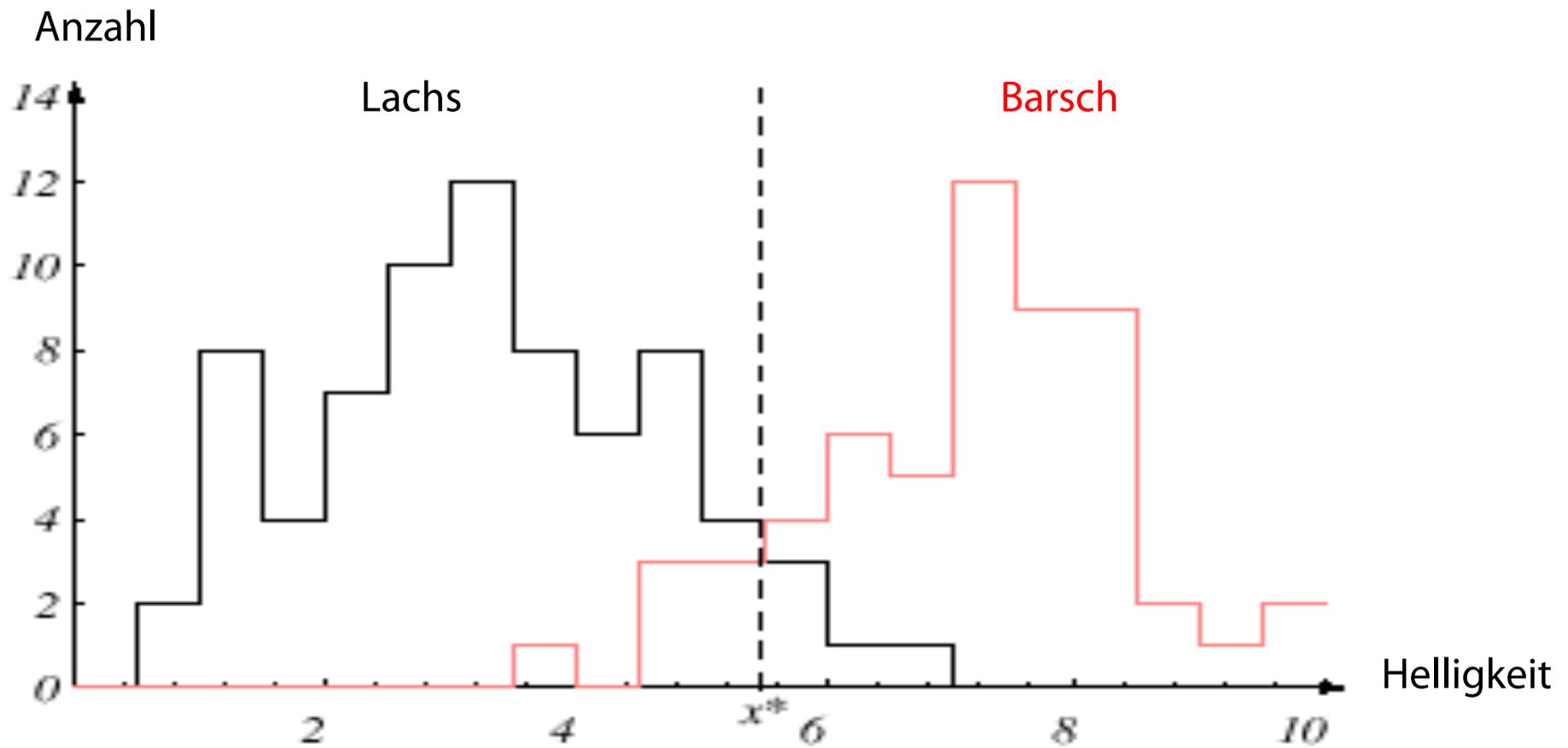
- Wir benötigen einen Experten, um die Merkmale festzulegen, mit denen man Barsche und Lachse richtig klassifizieren kann
- Wie wäre es mit Länge als Merkmal zur Unterscheidung?

Länge allein ist kein gutes Merkmal!

→ Hohe Kosten bei Fehlentscheidung

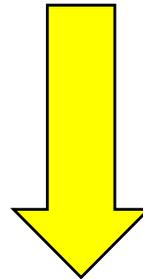
Wie wäre es mit Helligkeit?





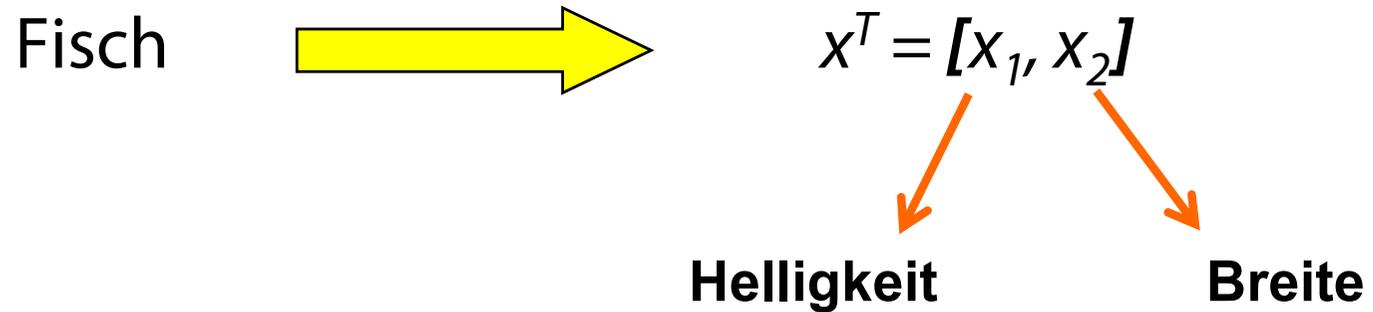
Schwellwert-Entscheidungsgrenze und induzierte Kosten

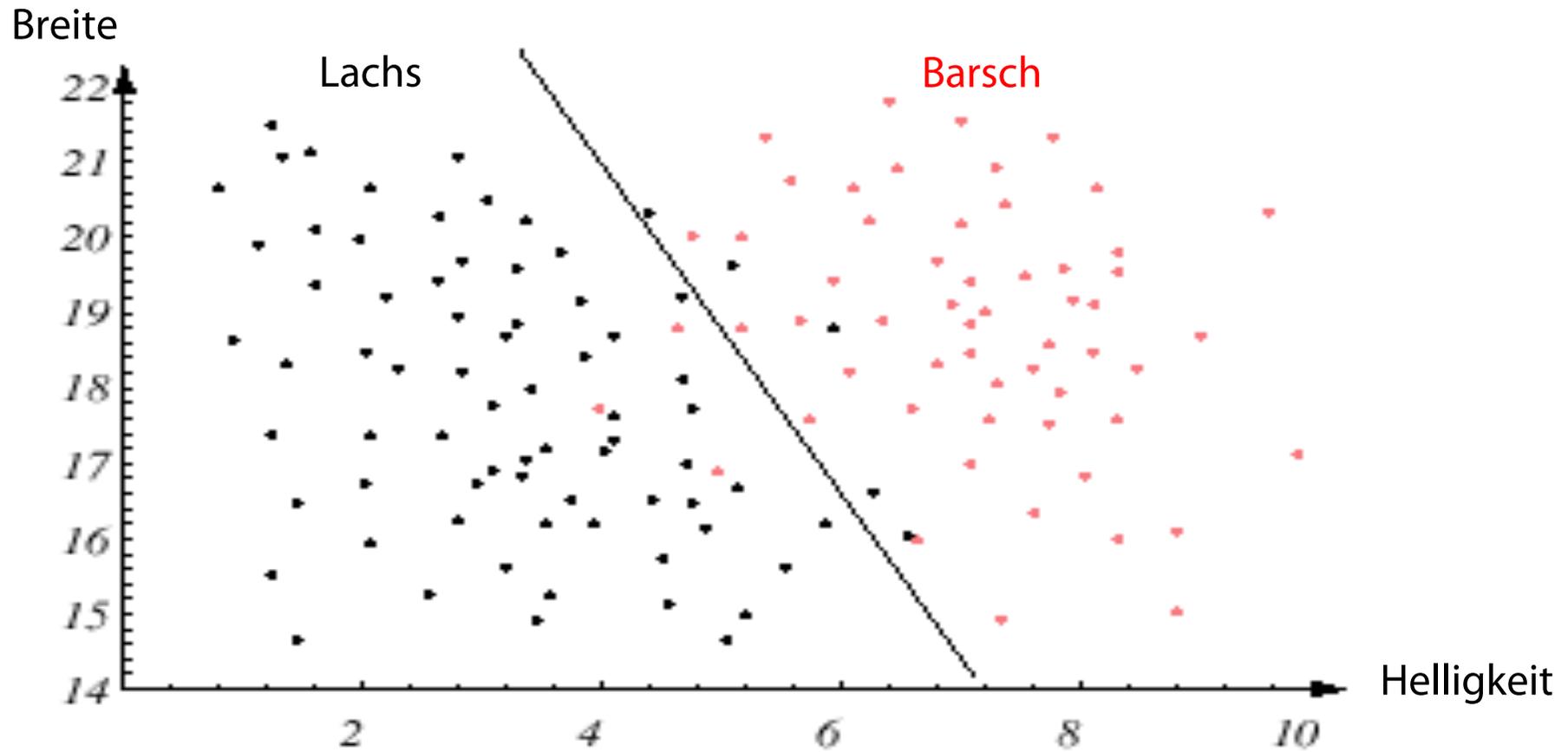
- Schwellwert-Entscheidungsgrenze in Richtung mittlerer Helligkeitswerte minimiert die Kosten (der Fehlklassifikation)



Untersuchung in der sog.
Entscheidungstheorie

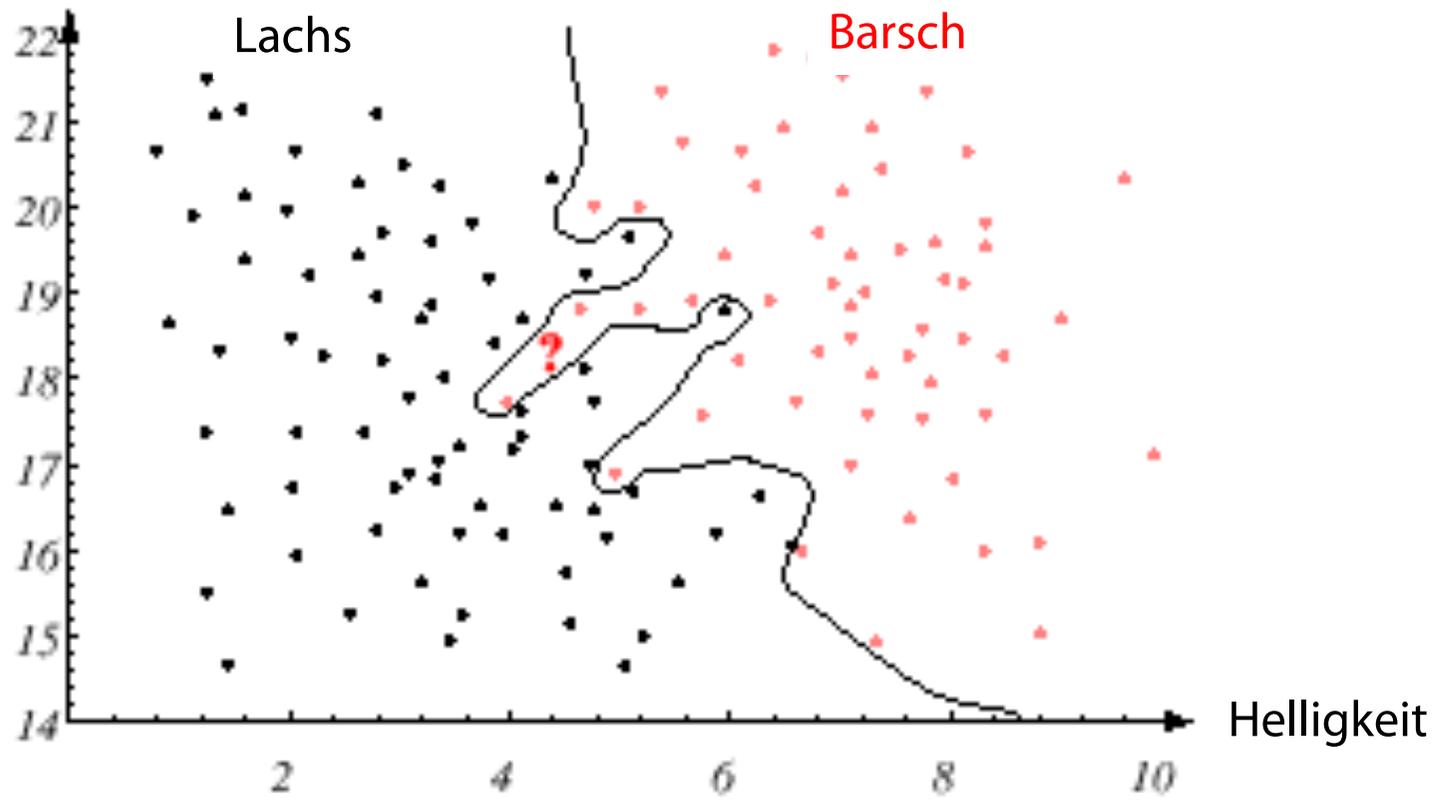
Passe Helligkeit an
und verwende zusätzlich die Breite des Fisches



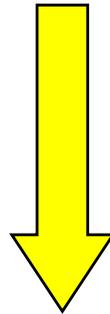


-
- Weitere Merkmale, die nicht direkt zu Helligkeit und Breite in Beziehung stehen, könnten hinzukommen
 - Vorsicht aber vor Reduktion durch "verrauschte Merkmale"
 - Wünschenswerterweise ergibt die **beste Entscheidungsgrenze** eine **optimale Performanz** (im Sinne einer Verlustminimierung durch Falschklassifikation)

Breite

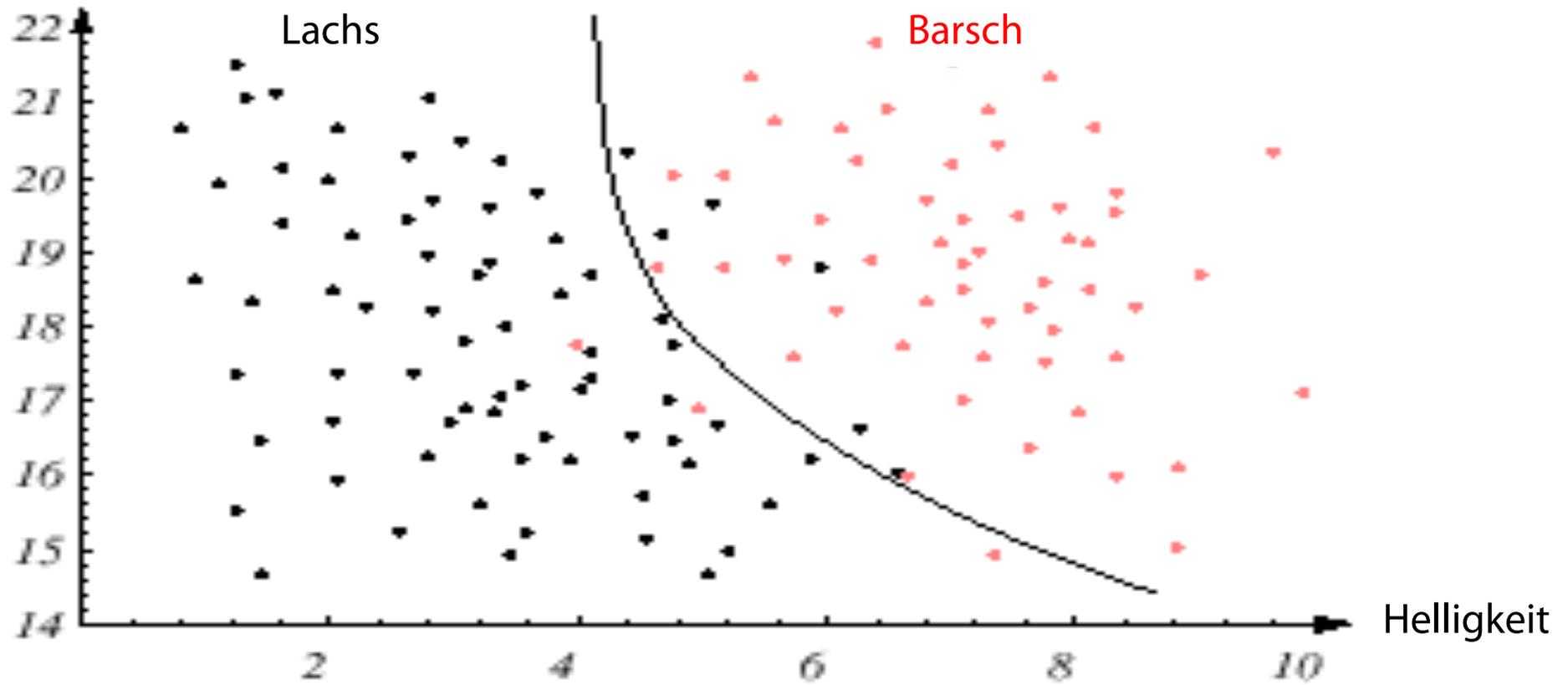


Vorfremde über Klassifikationsleistung auf Testdaten kann verfrüht sein. Wichtig ist die Leistung auf neuen Daten!



Generalisierungsfähigkeit zählt!

Breite



Klassifikatorentwicklung

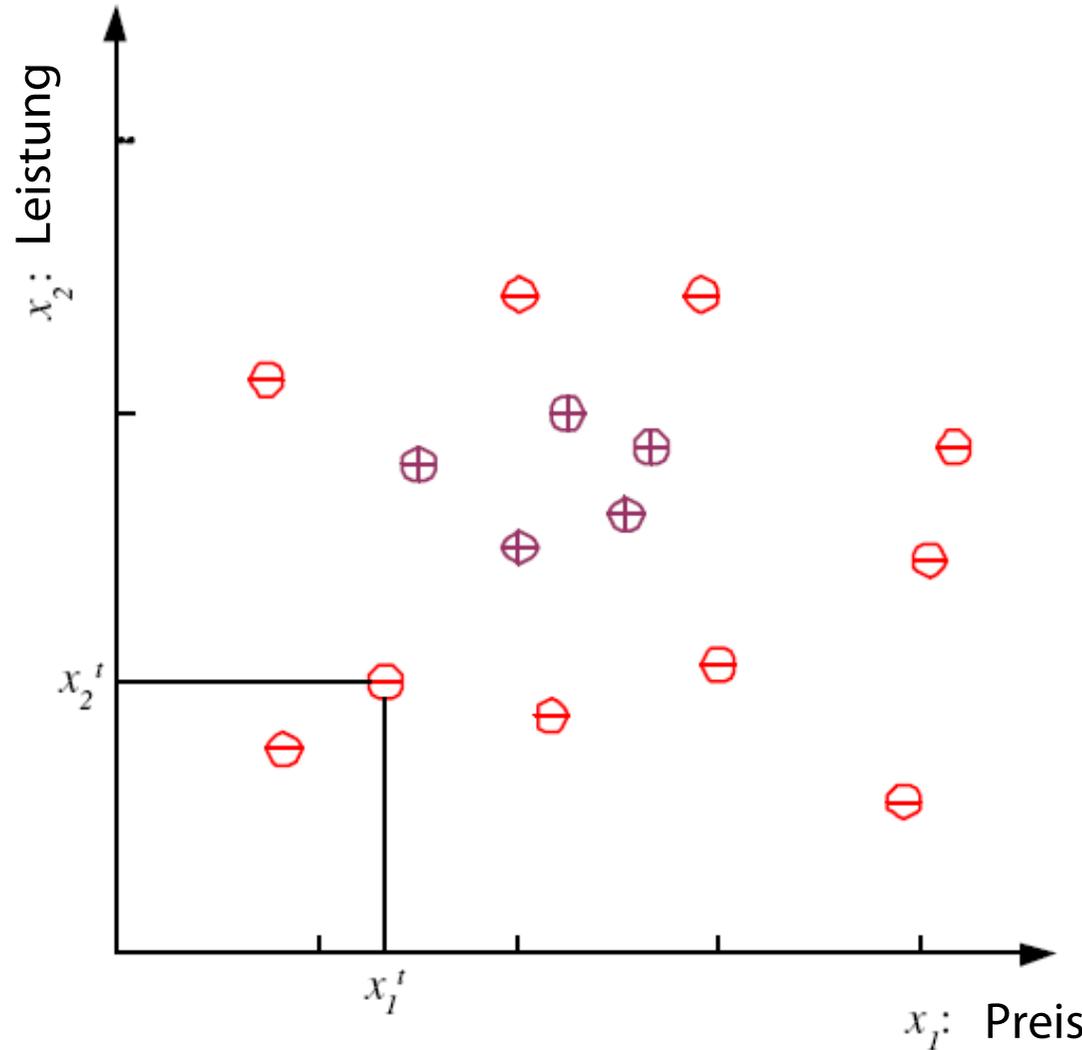
- Relevante Merkmale automatisch bestimmbar?
 - Deep Learning
- Dynamisch Anpassung des Klassifikators möglich?
 - Ohne spezielle Groundtruth Daten: Transduktives Lernen
 - Mit Belohnungssystem: Reinforcement-Lernen
- Übertragung eines Klassifikators auf neue Anwendung?
 - Transfer-Lernen

Bewertung eines Klassifikators

- Klasse C eines “Familienautos”
 - **Vorhersage:** Ist Auto x ein Familienauto?
 - **Wissensextraction:**
Was erwarten Menschen von einem Familienauto?
- Ausgabe:
 - Positive (+) und negative (–) Beispiele
- Repräsentation der Eingabe:
 - x_1 : Preis, x_2 : Leistung

Frage: Wie gut funktioniert ein bestimmter Klassifikator?

Trainingsmenge \mathcal{X}

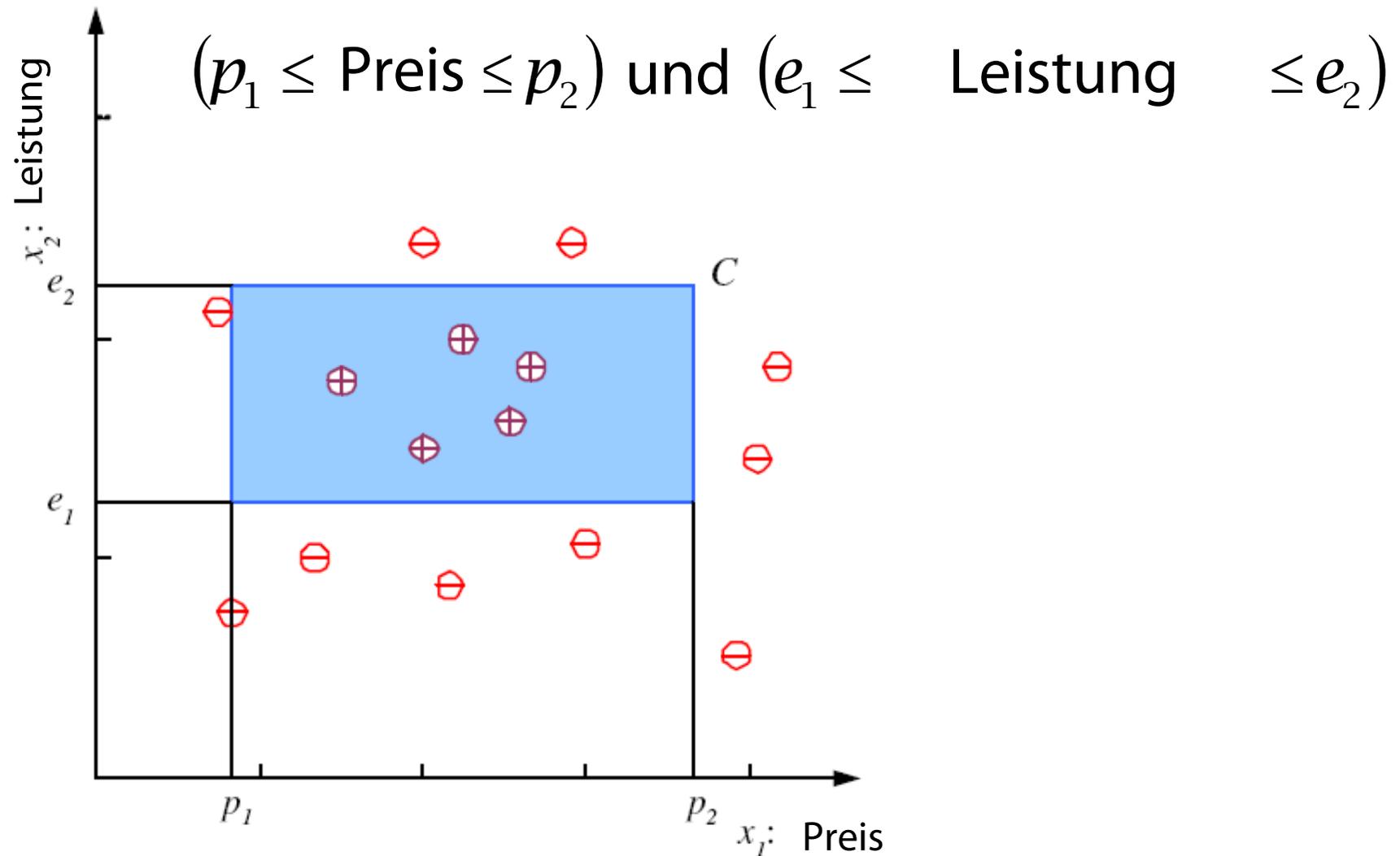


$$\mathcal{X} = \{\mathbf{x}^t, r^t\}_{t=1}^N$$

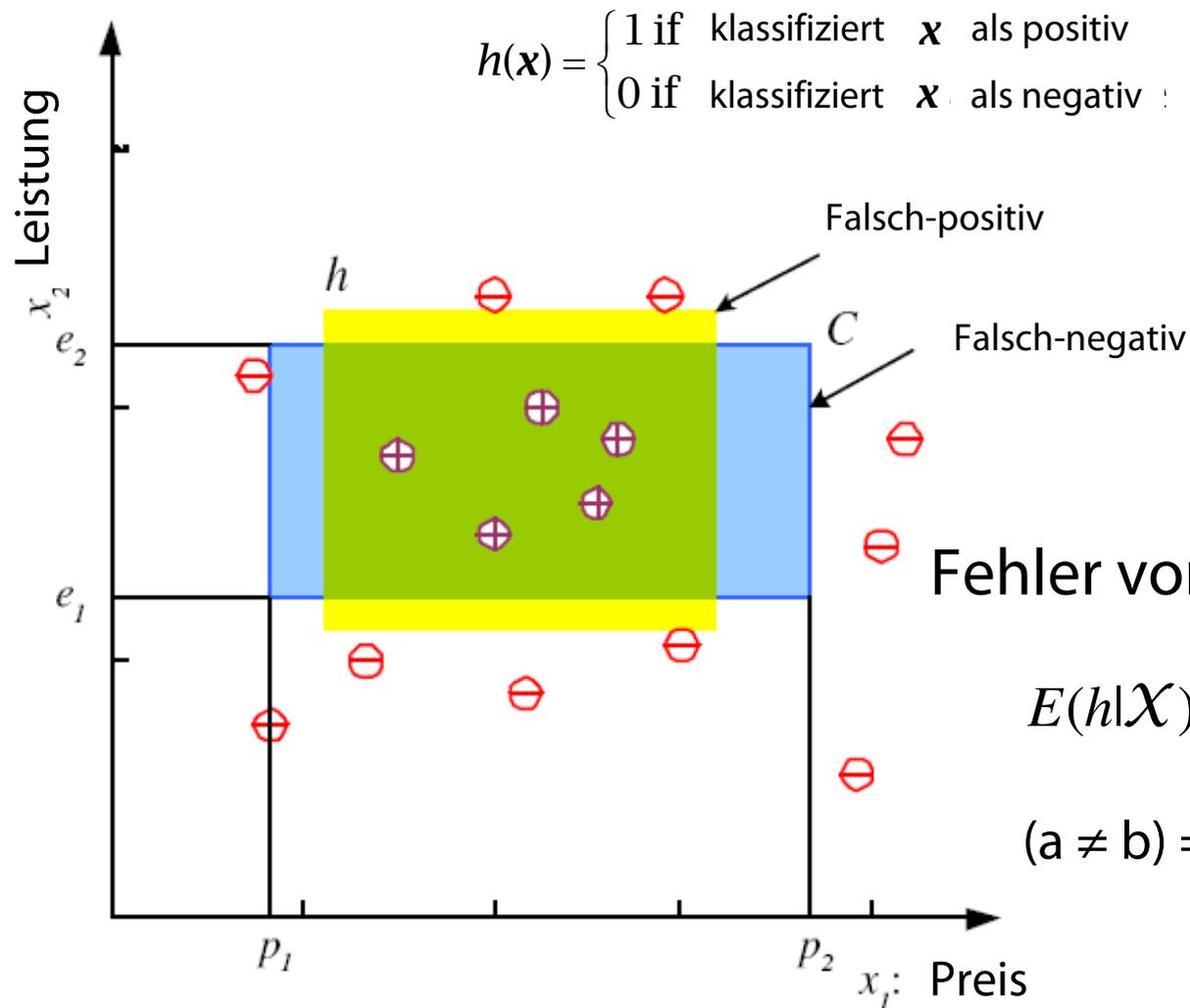
$$r = \begin{cases} 1 & \text{if } \mathbf{x} \text{ ist positiv} \\ 0 & \text{if } \mathbf{x} \text{ ist negativ} \end{cases}$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

Richtige Klasse C



Hypothesenklasse \mathcal{H} (gelb)

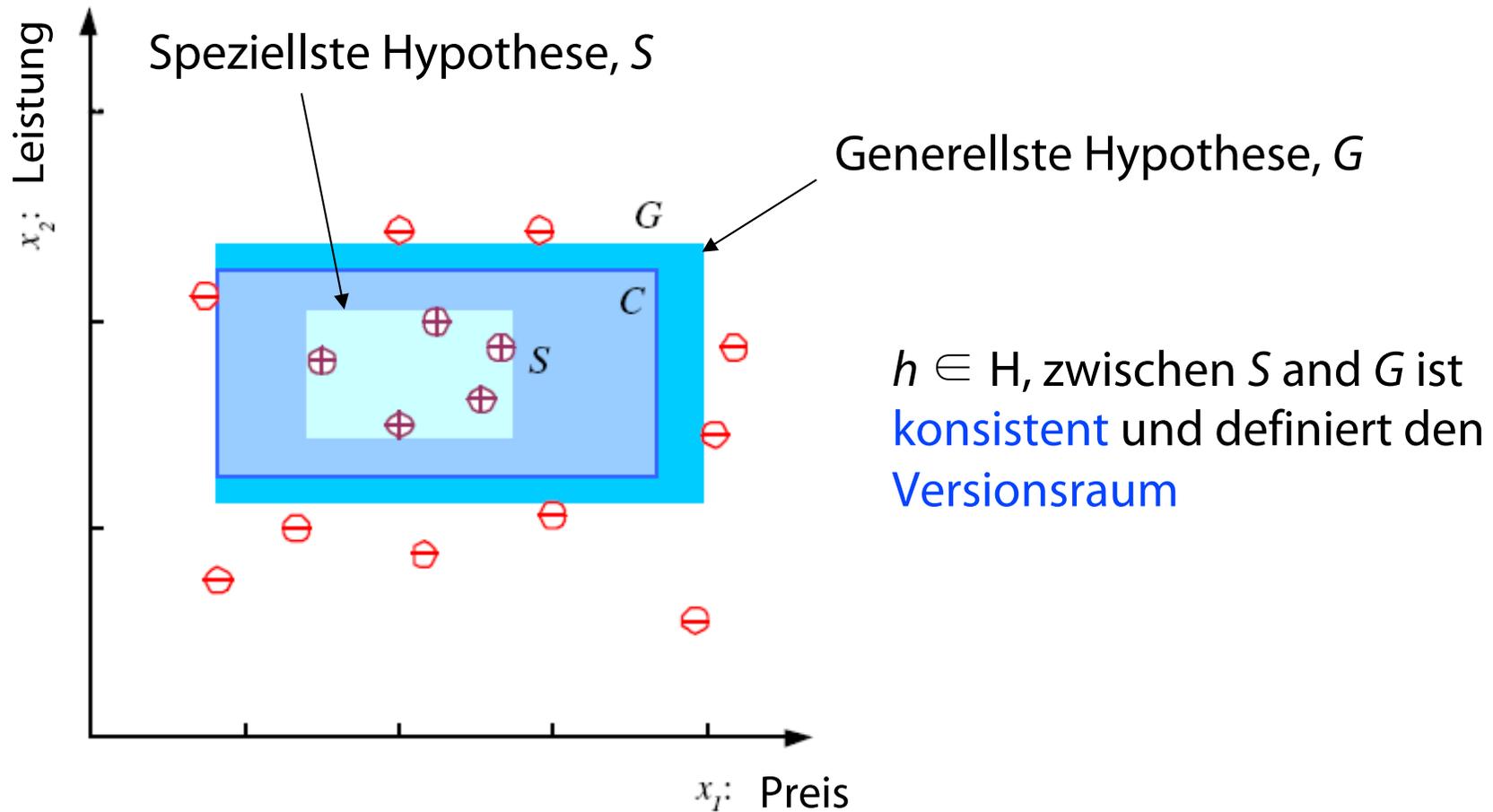


Fehler von h bzgl. \mathcal{H}

$$E(h|\mathcal{X}) = (1/N) \sum_{t=1}^N (h(\mathbf{x}^t) \neq r^t)$$

$$(a \neq b) = 1 \text{ if } \neq, 0 \text{ sonst}$$

S, G, and der Versionsraum (Version Space)

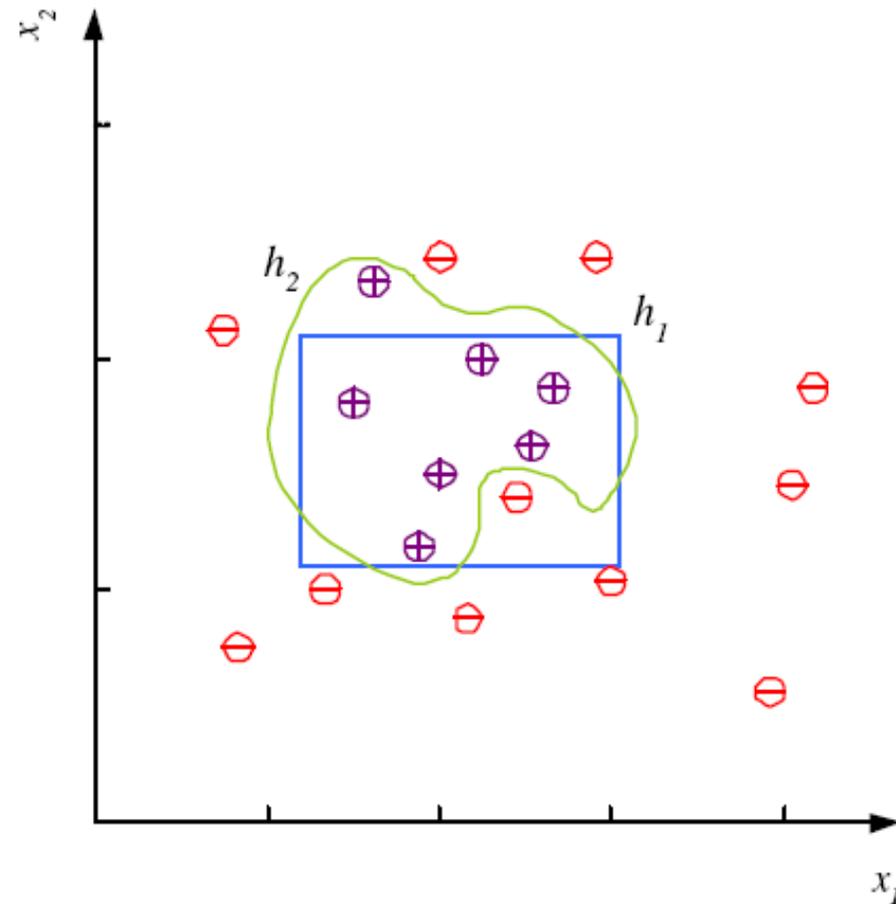


Rauschen und Modellkomplexität

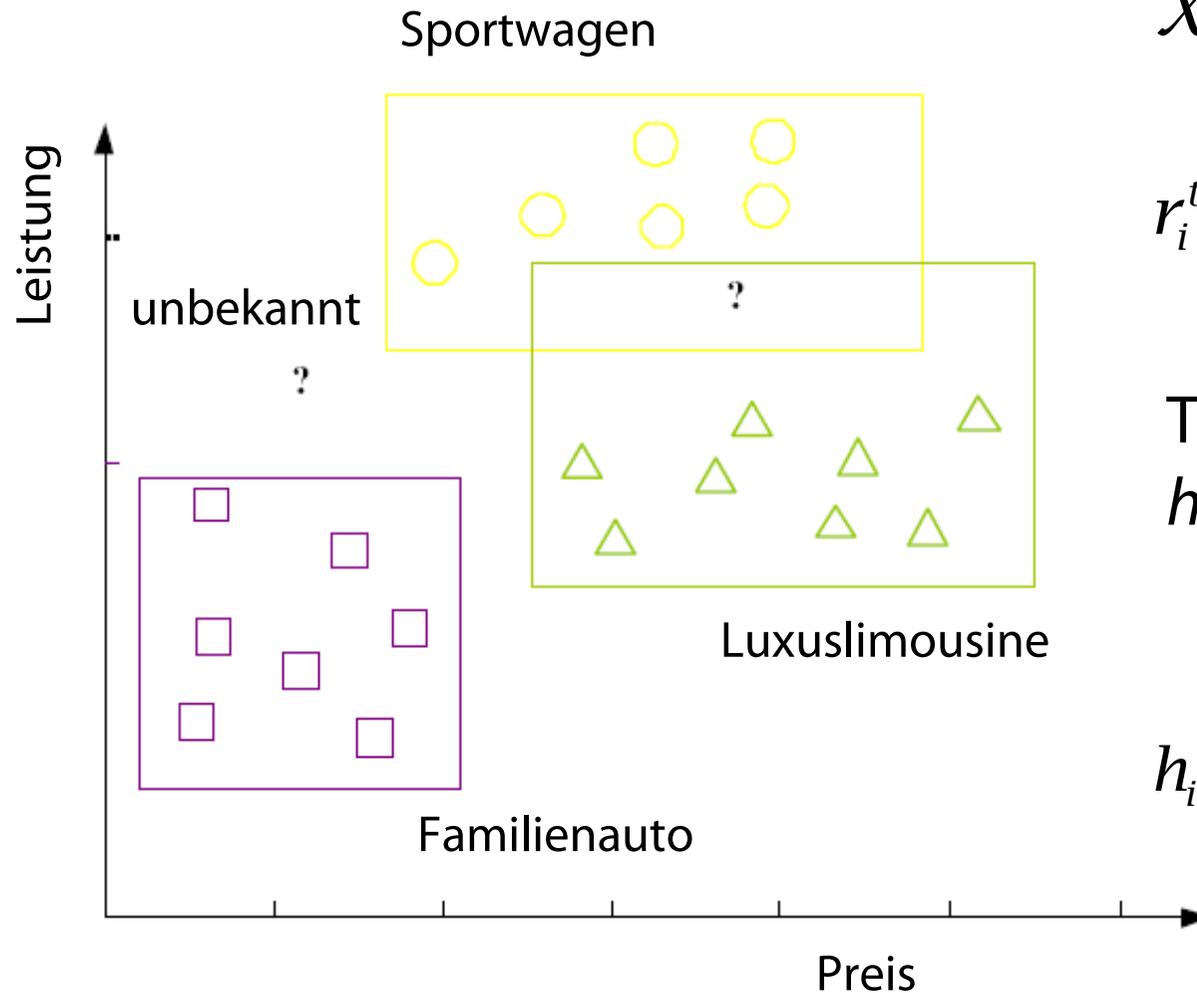
Verwende einfaches Modell:

- Einfacher zu verwenden
(weniger Berechnungsschritte)
- Leichter zu trainieren
(weniger Daten zu speichern)
- Leichter zu erklären
(besser interpretierbar)
- Bessere Generalisierung
(Occam's Razor)

Modellkomplexität:
"Größe" der Beschreibung



Clustering: Verschiedene Klassen, C_i $i=1, \dots, K$



$$\mathcal{X} = \{\mathbf{x}^t, r^t\}_{t=1}^N$$

$$r_i^t = \begin{cases} 1 & \text{if } \mathbf{x}^t \in C_i \\ 0 & \text{if } \mathbf{x}^t \in C_j, j \neq i \end{cases}$$

Trainiere Hypothesen

$h_i(\mathbf{x}), i=1, \dots, K:$

$$h_i(\mathbf{x}^t) = \begin{cases} 1 & \text{if } \mathbf{x}^t \in C_i \\ 0 & \text{if } \mathbf{x}^t \in C_j, j \neq i \end{cases}$$

Regression: Verschiedene Modellklassen

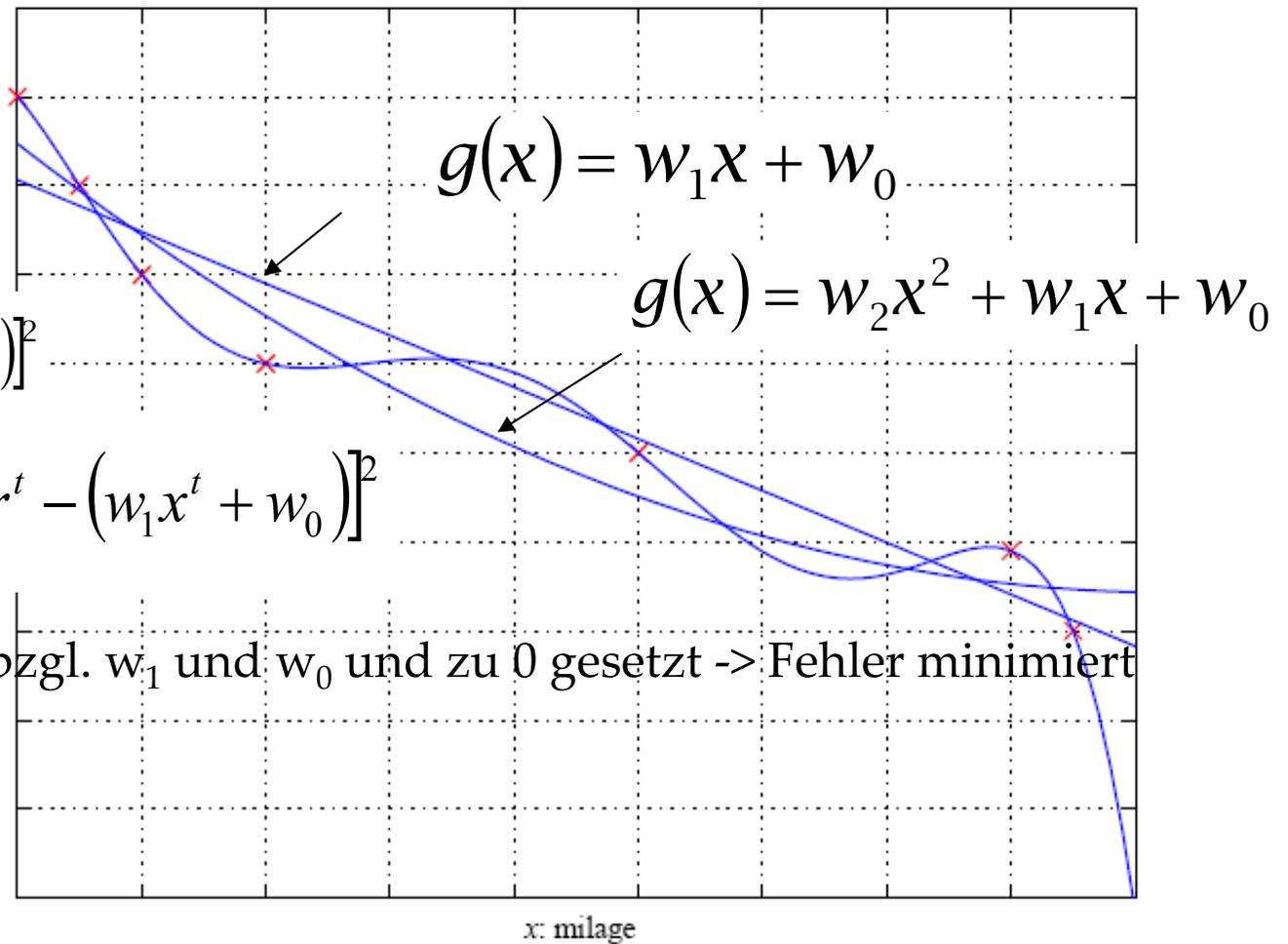
$$\mathcal{X} = \{x^t, r^t\}_{t=1}^N$$

$$r^t \in \mathbb{R}$$

$$r^t = f(x^t)$$

$$E(g | \mathcal{X}) = \frac{1}{N} \sum_{t=1}^N [r^t - g(x^t)]^2$$

$$E(w_1, w_0 | \mathcal{X}) = \frac{1}{N} \sum_{t=1}^N [r^t - (w_1 x^t + w_0)]^2$$



Partielle Ableitungen von E bzgl. w_1 und w_0 und zu 0 gesetzt -> Fehler minimiert

$$w_1 = \frac{\sum_t x^t r^t - \bar{x} \bar{r} N}{\sum_t (x^t)^2 - N \bar{x}^2}$$

$$w_0 = \bar{r} - w_1 \bar{x}$$

Überwachtes Lernen

- Beispiel: **Regression**
 - Gegeben Datenpunkte, bestimme Parameter, so dass für gegebene x-Werte, die y-Werte geschätzt werden können

- Beispiel: **Klassifikation**
 - Gegebene Datenpunkte jeweils mit Klassifikationswert, bestimme Klassifikationswert für Datenpunkte ohne diesen (binärer oder mehrwertiger Klassifikator)

Modellauswahl & Generalization

- Lernen kann als Optimierungsproblem angesehen werden:
 - Berechne Modell, so dass Fehlerfunktion minimiert
- Lernen ist ein **schlecht gestelltes Problem**;
In der Regel reichen die Daten nicht, eine eindeutige Lösung des Optimierungsproblems zu finden
- Es werden **Vorannahmen** benötigt (inductive bias),
Annahmen bzgl. H
- **Generalisierung**: Wie gut arbeitet das Modell auf neuen Daten?
- Überanpassung (Overfitting): H komplexer als C oder f
- Unteranpassung (Underfitting): H weniger komplex als C oder f

Drei-Wege-Austauschbeziehung

(Dietterich, 2003):

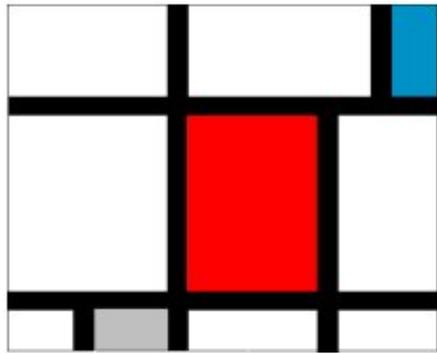
1. Komplexität von \mathcal{H} , $c(\mathcal{H})$,
 2. Trainingsmengengröße N ,
 3. Generalisierungsfehler, E , auf neuen Daten
- Wenn $N \uparrow$, $E \downarrow$
 - Wenn $c(\mathcal{H}) \uparrow$, gilt zuerst $E \downarrow$ und dann $E \uparrow$

Kreuzvalidierung

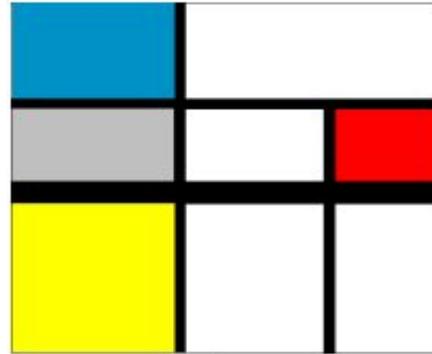
- Um den Generalisierungsfehler abzuschätzen, brauchen wir Daten, mit denen nicht trainiert wurde
- Aufteilung der Daten:
 - Trainingsmenge (50%)
 - Validierungsmenge (25%)
 - Testmenge (z.B. für Publikation) (25%)
- Neuabtastung, wenn wenige Daten vorhanden

Betrachtungsebenen überwachtes Lernen

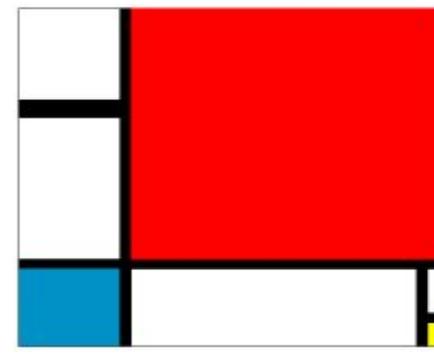
1. Modell: $g(\mathbf{x} | \theta)$
2. Fehlerfunktion
(Verlustfunktion): $E(\theta | \mathcal{X}) = \sum_t L(r^t, g(\mathbf{x}^t | \theta))$
3. Optimierungsverfahren: $\theta^* = \arg \min_{\theta} E(\theta | \mathcal{X})$



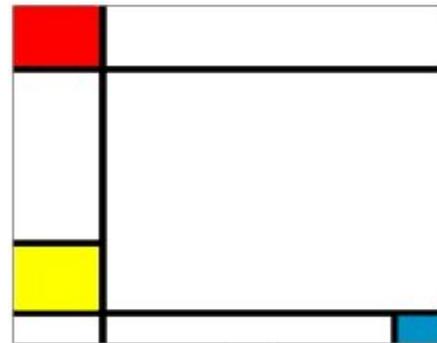
1



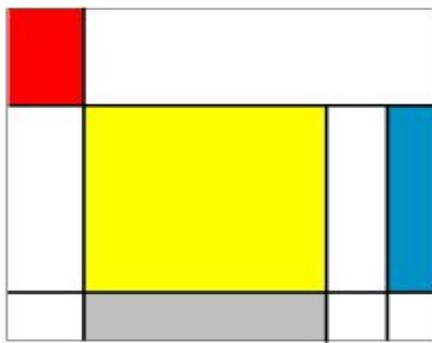
2



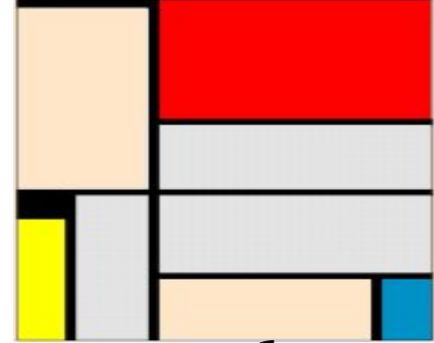
3



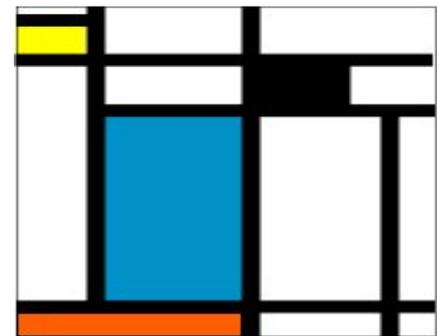
4



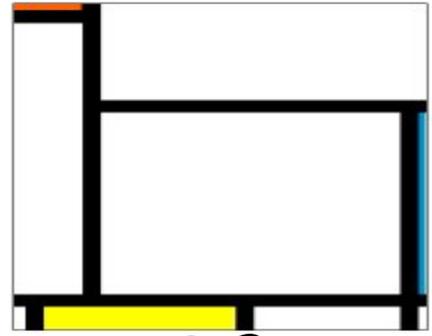
5



6



7



8 ?



Daten in Tabellarischer Form

Number	Lines	Line types	Rectangles	Colours	Mondrian?
1	6	1	10	4	No
2	4	2	8	5	No
3	5	2	7	4	Yes
4	5	1	8	4	Yes
5	5	1	10	5	No
6	6	1	8	6	Yes
7	7	1	14	5	No

Anfrage

Number	Lines	Line types	Rectangles	Colours	Mondrian?
8	7	2	9	4	

Analyse von Daten

- Betrachtung einer Spalte x mit n Werten
- Bestimmung des **Mittelwerts**: $\bar{x} = 1/n \cdot \sum_{i=1}^n x_i$
- Große und kleine Werte können sich aufheben
- Mittlere Abweichung vom Mittelwert betrachten (**Varianz**)
- Bestimmung der Varianz: $\text{var} = 1/n \cdot \sum_{i=1}^n (\bar{x} - x_i)^2$
- Meist betrachtet wird die sog. **Standardabweichung**: $\sigma = \sqrt{\text{var}}$

Halte Daten in normalisierter Form vor

Eine Möglichkeit zur Normalisierung:

$$x_t' \equiv \frac{x_t - \bar{x}_t}{\sigma_t}$$

Gemittelte Abweichung vom Mittel

Normalisierte Trainingsdaten

Number	Lines	Line types	Rectangles	Colours	Mondrian?
1	0.632	-0.632	0.327	-1.021	No
2	-1.581	1.581	-0.588	0.408	No
3	-0.474	1.581	-1.046	-1.021	Yes
4	-0.474	-0.632	-0.588	-1.021	Yes
5	-0.474	-0.632	0.327	0.408	No
6	0.632	-0.632	-0.588	1.837	Yes
7	1.739	-0.632	2.157	0.408	No

$$d(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{\sum_{t=1}^T [x_{it} - x_{jt}]^2}$$

Number	Lines	Line types	Rectangles	Colours	Mondrian?
8	1.739	1.581	-0.131	-1.021	

Normalisierte Trainingsdaten

Number	Lines	Line types	Rectangles	Colours	Mondrian?
1	0.632	-0.632	0.327	-1.021	No
2	-1.581	1.581	-0.588	0.408	No
3	-0.474	1.581	-1.046	-1.021	Yes
4	-0.474	-0.632	-0.588	-1.021	Yes
5	-0.474	-0.632	0.327	0.408	No
6	0.632	-0.632	-0.588	1.837	Yes
7	1.739	-0.632	2.157	0.408	No

$$\sqrt{(0 + 4,89 + 5,23 + 2,04)} = 3,489$$

Number	Lines	Line types	Rectangles	Colours	Mondrian?
8	1.739	1.581	-0.131	-1.021	

Distanz der Testinstanz von den Trainingsdaten

Beispiel	Distanz zum Test	Mondrian?
1	2.517	No
2	3.644	No
3	2.395	Yes
4	3.164	Yes
5	3.472	No
6	3.808	Yes
7	3.490	No

Klassifikation

1-NN	Yes
3-NN	Yes
5-NN	No
7-NN	No

Was verwenden wir bei reellwertiger Zielfunktion als Ausgabe?

- Mittel der k -nächsten Nachbarn

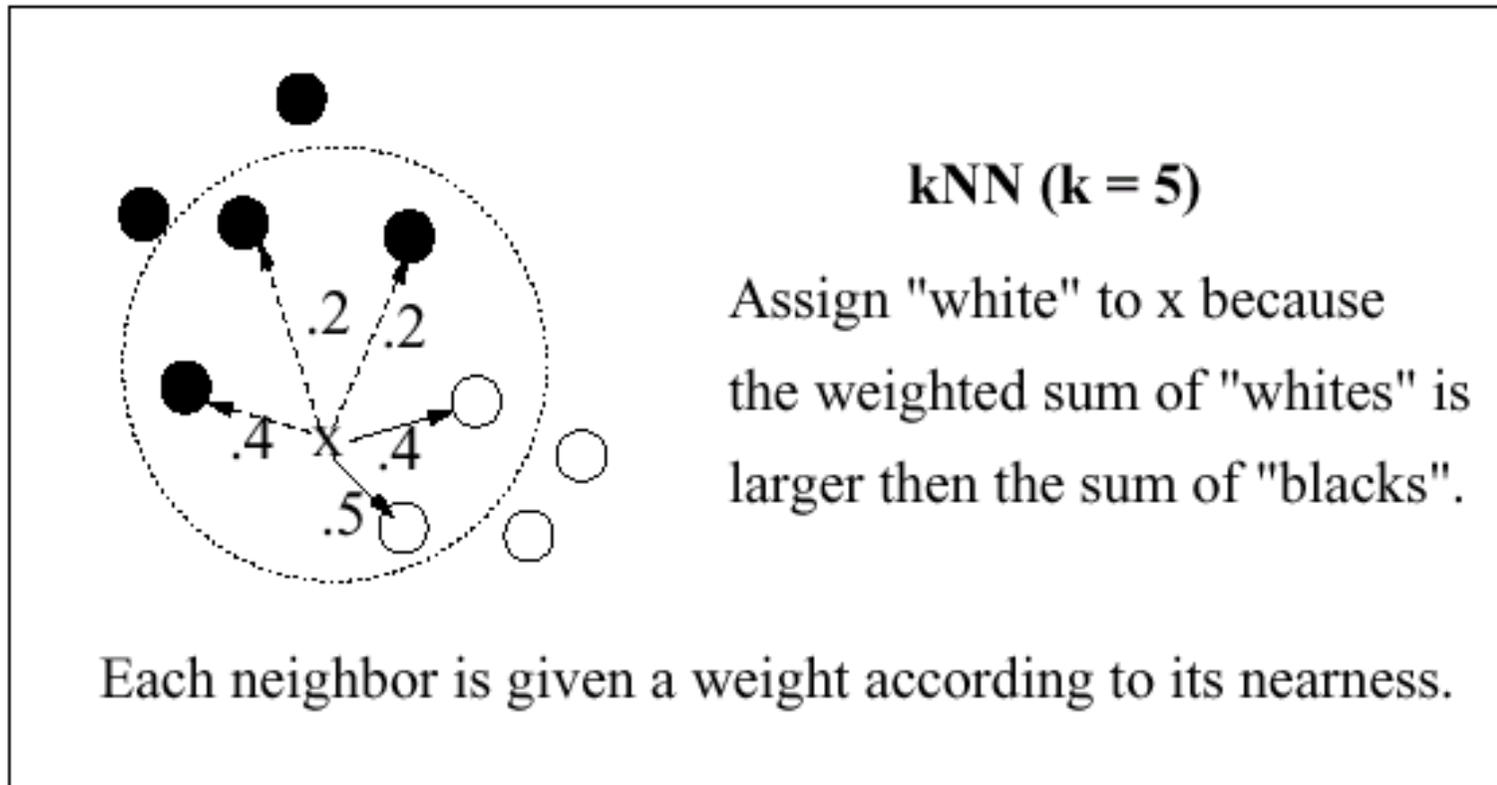
Variante von kNN: Distanzgewichtetes kNN

- Nähere Nachbarn haben mehr Einfluss

$$f(\mathbf{x}_q) := \frac{\sum_{i=1}^k w_i f(\mathbf{x}_i)}{\sum_{i=1}^k w_i} \quad \text{where } w_i = \frac{1}{d(\mathbf{x}_q, \mathbf{x}_i)^2}$$

Variante von kNN: Distanzgewichtetes kNN

k-NN using a weighted-sum voting scheme



Dann könnten wir statt nur **k** gleich **alle** Trainingsinstanzen (= Beispiele) nehmen

kNN: Zusammenfassung

- Sehr einfacher Ansatz
- Verhält sich auch noch gutartig, wenn Daten nicht einfach separiert werden können
- Rang 7 der 10 wichtigsten Data-Mining-Verfahren

