

Klassifikation

Agenda

I Begriff Klassifikation

- Abgrenzung
- Anforderungen
- Anwendungsgebiete

III Dimensionsreduktion

Umsetzung in Software

Vergleich der Algorithmen

Quellenangaben

Fragen

II Klassifikatoren

- Entscheidungsbaum
- naive Bayes – Algorithmus
- Künstlich neuronale Netze
- k-nearest neighbour – Methode
- Support Vector Machine

- eine Data-Mining-Methode, die der Zusammenfassung von Datenobjekten zu disjunkten Klassen dient
- Klasse = Menge von zusammengehörigen Datenobjekten
- im Gegensatz zum Clustering müssen die Klassen bekannt sein
- zwei Phasen:
 1. Lernphase
 2. Klassifizierung von Datenobjekten
- Attribute (nominal, numerisch) der Datenobjekte

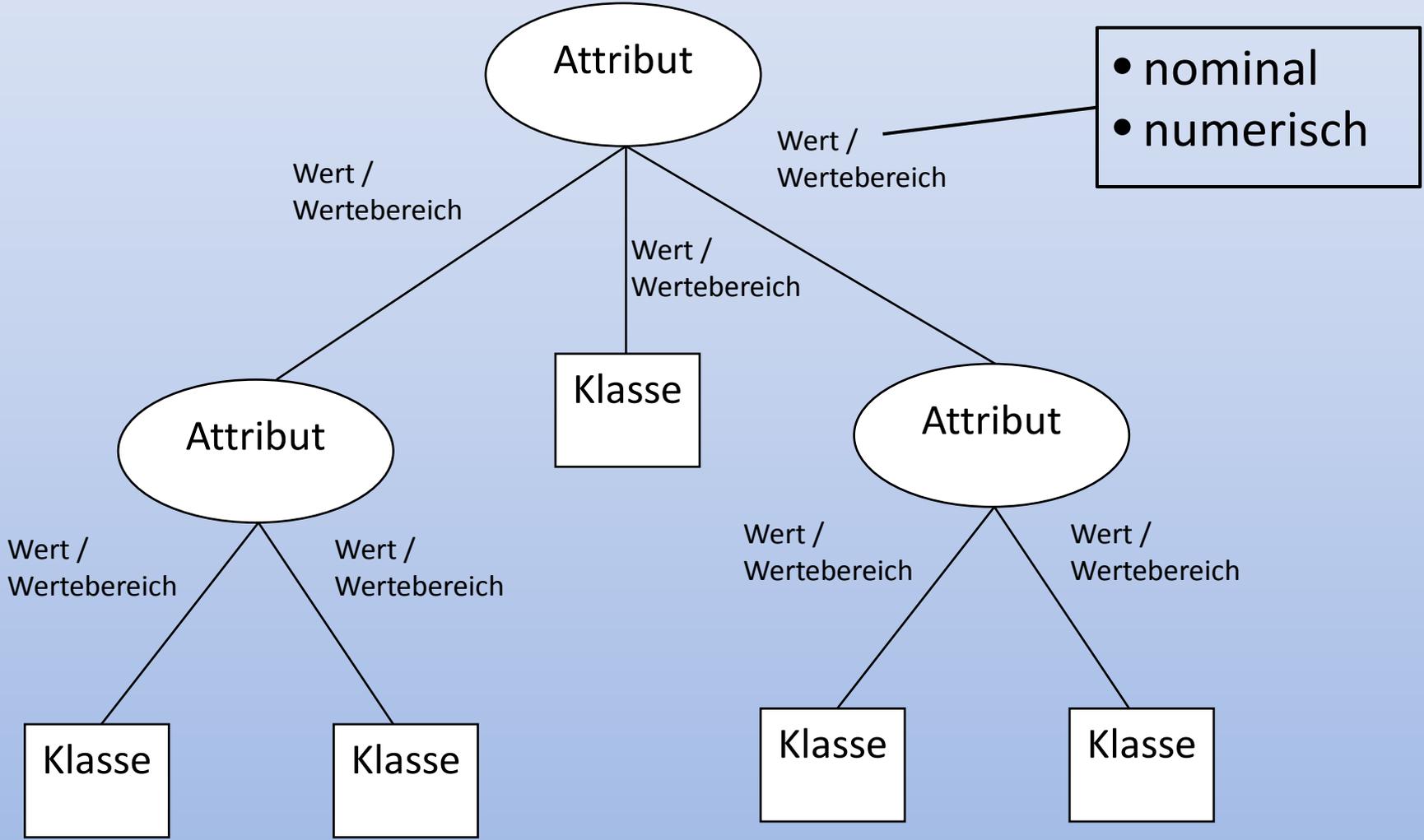
Anforderungen an einen Klassifikator und Daten

- möglichst korrekte Klassifizierung der Daten
- eine akzeptable Performance und Skalierbarkeit
- Ausreißern, fehlende Attribute
- Overfitting vermeiden
- Trainingsdaten sind notwendig

Anwendungsgebiete

- Medizin
- Marketing & Marktanalysen
- Versicherung
- Mustererkennung

Entscheidungsbaum



Konstruktion eines Entscheidungsbaumes

Informationsmenge:

- Maß für den Informationsgehalt, der nach dem Wissen über den Attributwert vorhanden ist
- auch als Entropie bezeichnet
- Einheit: bits
- $\text{info}([p_1, p_2, \dots, p_n]) = -p_1 \log_2 p_1 - p_2 \log_2 p_2 \dots - p_n \log_2 p_n$
n – Anzahl der Klassen
- $p_i = \frac{\text{\# des Attributwerts}}{\text{\# aller Werte dieses Attributs}}$
- $n = 1 \rightarrow \text{info}(p_1) = 0$

Beispiel - Wetterdaten

Informationsgewinn Wurzel: $info([9,5]) = -\frac{5}{14} \log_2 \frac{5}{14} - \frac{9}{14} \log_2 \frac{9}{14} = 0.940$

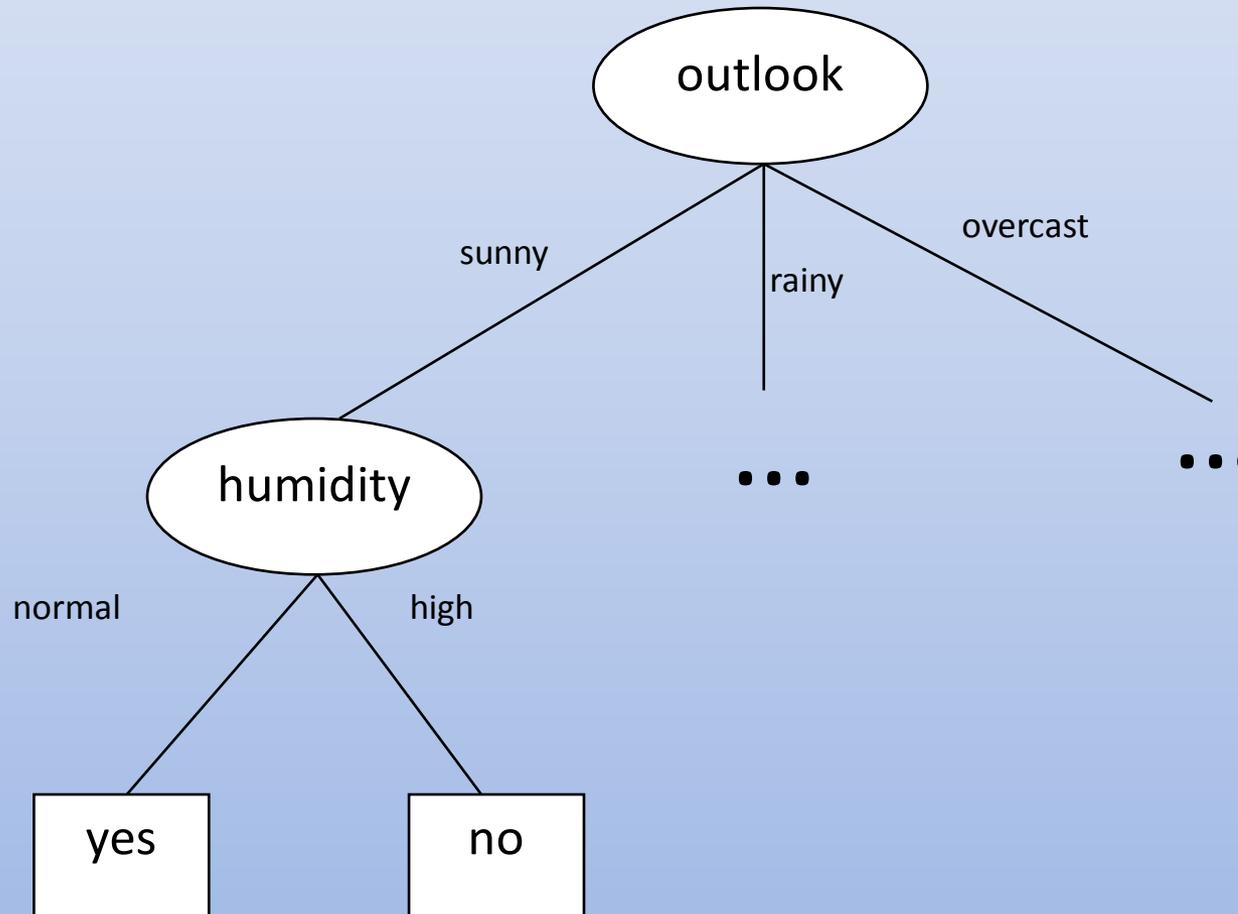
Informationsgewinn aller Attribute:

<p>outlook</p> <p>$info_{overcast}([4,0]) = info_{overcast}([4]) = 0$</p> <p>$info_{rainy}([2,3]) = 0.971$</p> <p>$info_{sunny}([3,2]) = 0.971$</p> <p>Mittelwert: $\frac{4}{14} \cdot 0 + \frac{5}{14} \cdot 0.971 + \frac{5}{14} \cdot 0.971 = 0.693$</p> <p>Informationsgewinn = $0.940 - 0.693 = 0.247$</p>	<p>temperature</p> <p>$info_{cool}([3,1]) = 0.811$</p> <p>$info_{mild}([4,2]) = 0.919$</p> <p>$info_{warm}([2,2]) = 1$</p> <p>Mittelwert = 0.911</p> <p>Informationsgewinn = 0.029</p>	<p>humidity</p> <p>$info_{high}([3,4]) = 0.985$</p> <p>$info_{normal}([6,1]) = 0.592$</p> <p>Mittelwert = 0.789</p> <p>Informationsgewinn = 0.152</p>	<p>windy</p> <p>$info_{true}([2,6]) = 0.811$</p> <p>$info_{false}([3,3]) = 1$</p> <p>Mittelwert = 0.892</p> <p>Informationsgewinn = 0.048</p>
--	---	---	---

outlook = sunny

<p>temperature</p> <p>$info_{cool}([1]) = 0$</p> <p>$info_{mild}([1,1]) = 1$</p> <p>$info_{hot}([2]) = 0$</p> <p>Mittelwert = 0.400</p> <p>Informationsgewinn = $0.971 - 0.400 = 0.571$</p>	<p>humidity</p> <p>$info_{high}([2]) = 0$</p> <p>$info_{normal}([3]) = 0$</p> <p>Mittelwert = 0</p> <p>Informationsgewinn = 0.971</p>	<p>windy</p> <p>$info_{true}([1,1]) = 1$</p> <p>$info_{false}([1,2]) = 0.918$</p> <p>Mittelwert = 0.951</p> <p>Informationsgewinn = 0.020</p>
---	---	---

Beispiel - Wetterdaten



fehlende Attributwerte

- fehlende Attributwerte stellen einen eigenen Zweig dar
- die in der Trainingsmenge gebräuchlichste Verzweigung wird gewählt
- jede Verzweigung des fehlenden Attributwertes wird durchlaufen, die Blätter werden anhand der Anzahl der Elemente der Trainingsinstanz numerisch gewichtet (sehr komplex)

Attribute mit vielen Verzweigungen

- Mehrfachverzweigung mit sehr vielen untergeordneten Knoten
- z.B. numerische Attribute: Wertebereiche für numerische Attribute
- IDs: höchster Informationsgewinn -> Verzweigung nach diesem Attribut ?
- statt Informationsgewinn wird das Gewinnverhältnis verwendet (siehe Literatur)

naive Bayes - Verfahren

- stochastisches Klassifizierungsmodell

Datenobjekt

z.B.: meine Prüfung in BWL

- gelernt: mittel
- Spicker: nein
- Wohlbefinden: gut
- Vorlesung besucht: regelmäßig

Klasse

Note 1

Note 2

Note 3

Note 4

0.1

0.4

0.3

0.2

formal:

n – #Klassen

m – #Merkmale/Attribute

$$\textit{Klassenzuordnung} = \underset{1 \leq i \leq n}{\textit{index}} \{ \max(P(C_i | X)) \}$$

C_i – Klasse i

X - Datenobjekt

$P(C_i)$ – a – priori – Wahrscheinlichkeit

$P(C_i | X)$ – a – posteriori – Wahrscheinlichkeit

$$P(C_i | X) = \frac{P(X | C_i) \cdot P(C_i)}{P(X)}$$

$$P(X | C_i) = \prod_{j=1}^m P(x_j | C_i)$$

- für kontinuierliche Merkmale ist eine Verteilungsfunktion nötig

Beispiel Wetter

Beispiel Wetter

Grundlage der Klassifikation

	Anzahl		relative Häufigkeit	
	P(yes)	P(no)	P(x _j yes)	P(x _j no)
	yes	no	yes	no
play	9	5	9/14	5/14
outlook				
sunny	2	3	2/9	3/5
overcast	4	0	4/9	0/5
rainy	3	2	3/9	2/5
temperature				
hot	2	2	2/9	2/5
mild	4	2	4/9	2/5
cool	3	1	3/9	1/5
humidity				
high	3	4	3/9	4/5
normal	6	1	6/9	1/5
windy				
false	6	2	6/9	2/5
true	3	3	3/9	3/5

Klassifikation

X₁

outlook = sunny
 temperature = cool
 humidity = high
 windy = true

$$index\{\max(P(C_i | X))\}_{1 \leq i \leq n}$$

$$P(C_i | X) = \frac{P(X | C_i) \cdot P(C_i)}{P(X)}$$

Klassifikation

$$P(\text{yes} | X_1) = \frac{\frac{2}{9} \cdot \frac{3}{9} \cdot \frac{3}{9} \cdot \frac{3}{9} \cdot \frac{9}{14}}{P(X)} = \frac{0.0053}{0.0053 + 0.0206} = 0.205$$

$$P(\text{no} | X_1) = \frac{\frac{3}{5} \cdot \frac{1}{5} \cdot \frac{4}{5} \cdot \frac{3}{5} \cdot \frac{5}{14}}{P(X)} = \frac{0.0206}{0.0053 + 0.0206} = 0.795$$

-> Zuordnung zu Klasse ,no‘

Problem

Attributwert in den Trainingsdaten gibt es nicht in Kombination mit einer Klasse (siehe outlook=overcast) -> Laplace-Schätzung

	Anzahl		relative Häufigkeit		Laplace Schätzung	
	P(yes)	P(no)	P(x _j yes)	P(x _j no)	yes	no
	yes	no	yes	no		
outlook						
sunny	2	3	2/9	3/5	3/12	4/8
overcast	4	0	4/9	0/5	5/12	1/8
rainy	3	2	3/9	2/5	4/12	3/8

fehlende Werte

fehlende Wert stellen keine Probleme dar

-> Trainingsdaten: Instanz wird bei der Häufigkeitszählung für die Attributwert-Klassenkombination nicht berücksichtigt

-> Klassifikation: Attribut wird ausgelassen

numerische Attribute

nominale Werte

temperature	P(yes)	P(no)	P(x _i yes)	P(x _i no)
hot	2	2	2/9	2/5
mild	4	2	4/9	2/5
cool	3	1	3/9	1/5

numerische
Werte

temperature	yes	no
numerische Attributwerte	83	85
	70	80
	68	65
	64	72
	69	71
	75	
	75	
	72	
	81	
	Mittelwert μ	73
Standardabweichung σ	6,2	7,9

$$\sigma = \frac{1}{n-1} \sum_{i=1}^n (x_i - \mu)^2$$

numerische Attribute

temperature	yes	no
numerische Attributwerte	83	85
	70	80
	68	65
	64	72
	69	71
	75	
	75	
	72	
	81	
Mittelwert μ	73	74,6
Standardabweichung σ	6,2	7,9

$$f(x) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

rein nominale Attribute

$P(\text{temperature} = \text{mild} \mid \text{yes})$

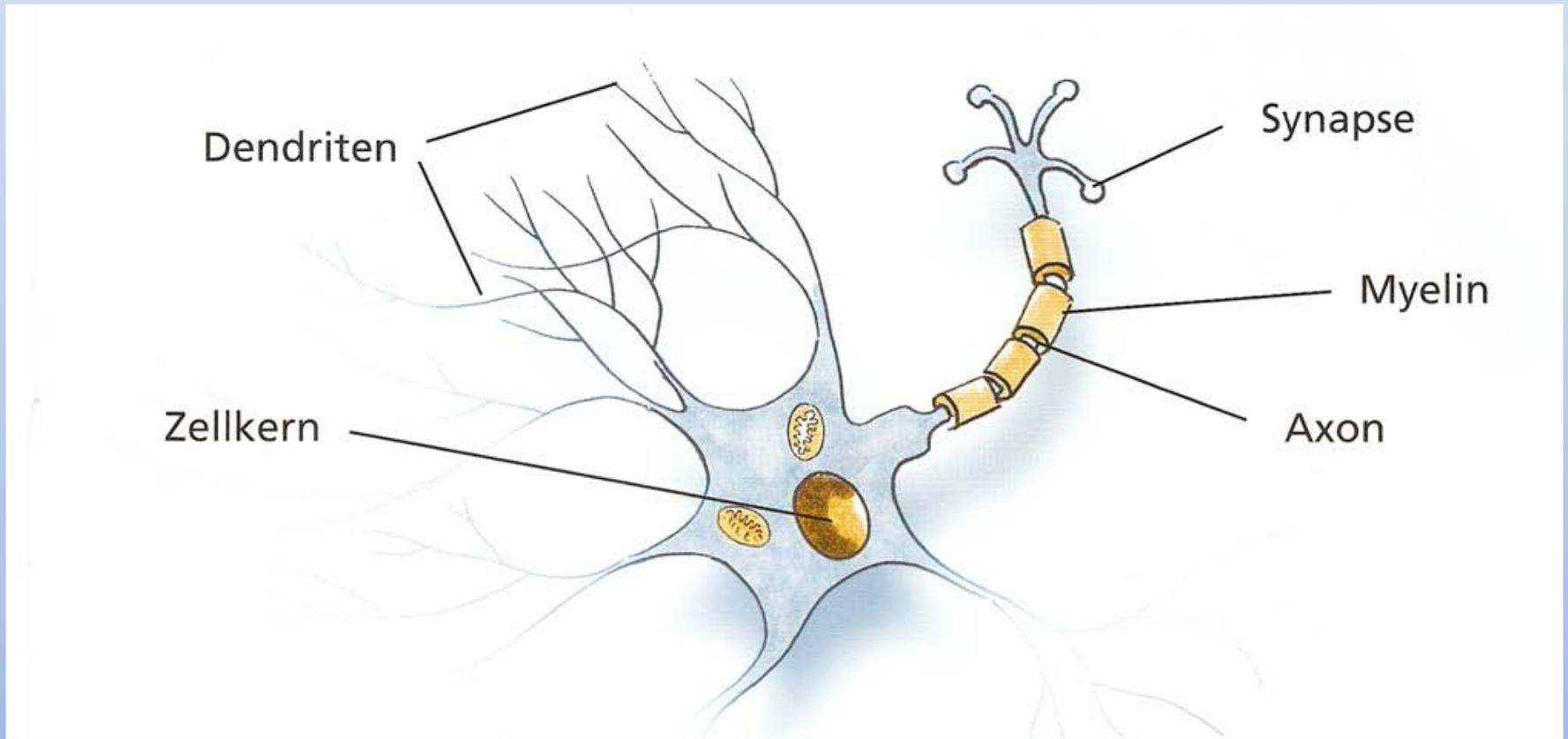
$$P(\text{yes} \mid X_1) = \frac{\frac{2}{9} \cdot \frac{3}{9} \cdot \frac{3}{9} \cdot \frac{3}{9} \cdot \frac{9}{14}}{P(X)}$$

$$f(\text{temperature} = 66 \mid \text{yes}) = 0.034$$

$$P(\text{yes} \mid X_1) = \frac{\frac{2}{9} \cdot 0.034 \cdot \frac{3}{9} \cdot \frac{3}{9} \cdot \frac{9}{14}}{P(X)}$$

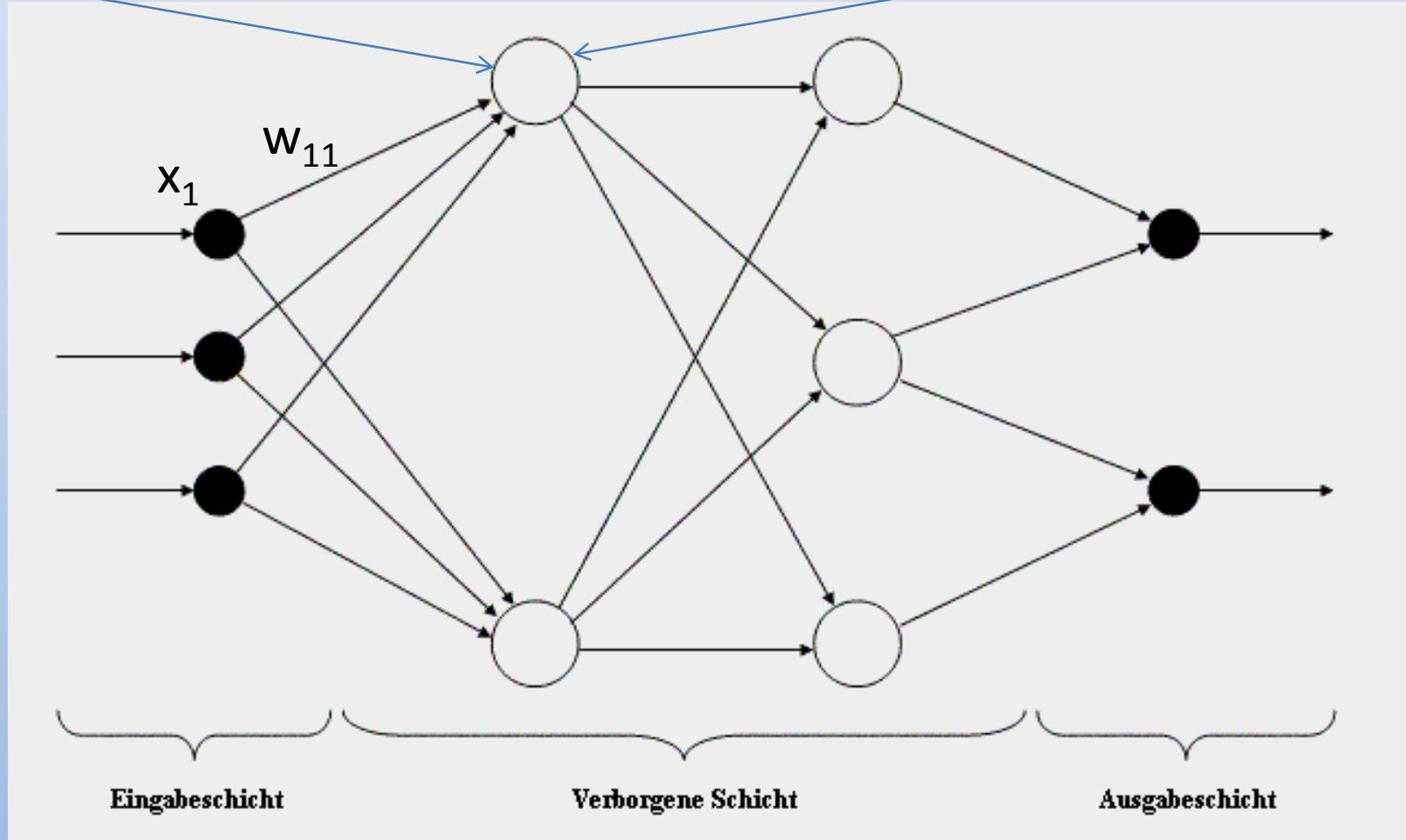
Künstliche Neuronale Netze

Neuron (Nervenzelle)



Aufbau eines KNN

$$f_{\text{Eingabe}} = \sum_{i=1}^n x_i w_{ij} \quad n\text{-\# Neuronen der Vorgängerschicht, } j\text{- Neuron der Schicht} \quad f_{\text{Aktivierung}}$$



Aktivierungsfunktion

- Schwellwert bzw. Sprungfunktion
- lineare Funktion
- Sigmoidfunktion

Backpropagation-Algorithmus

- Trainingsdaten werden in das Netz gespeist
- die Differenz zwischen Netzausgabe und gewünschten Ausgabe wird als Fehler berechnet
- liegt der Fehler über der Toleranzgrenze wird er zurück propagiert -> die Gewichtungen der Neuronenverknüpfungen werden abhängig vom Einfluss auf den Fehler angepasst
- Vorgang wiederholt sich bis der Fehler innerhalb der Toleranzgrenze liegt
- Backpropagation stellt ein Näherungsverfahren dar

k – nearest neighbour

Idee

- Einordnung der Datenobjekt in gleiche Klasse wie ähnliche Trainingsobjekte
- Beschreibung einer Klasse durch ein Trainingsobjekt meist unzureichend
- Klassenzentrum bezeichnet ein Trainingsobjekt
- Verwendung einer Diskriminanzfunktion
- Berechnung des Abstands zu jedem Klassenzentrum zum Zeitpunkt der Klassifikation

Diskriminanzfunktion

$J(i)$ – Menge der Klassenzentren für Klasse i

$z_i(j)$ – j -te Zentrum der Klasse i

x – Datenobjekt

m – # Klassen

$$d_i(x) = \min_{1 \leq j \leq J(i)} \{ z_i^T(j) \cdot z_i(j) - 2 \cdot z_i^T(j) \cdot x \}$$

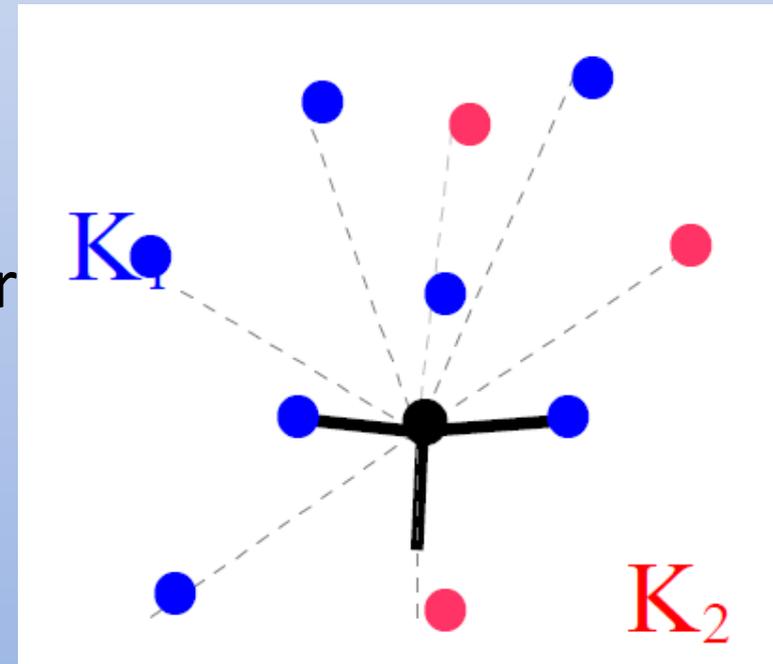
$$\text{Klassenindex} = \text{index} \left\{ \min_{1 \leq i \leq m} d_i(x) \right\}$$

NN -> kNN

- schwächere Berücksichtigung von ‚Ausreißern‘
wünschenswert -> Berechnung aller Werte

$$d_{ij}(x) = z_i^T(j) \cdot z_i(j) - 2 \cdot z_i^T(j) \cdot x$$

- Klassenzuordnung erfolgt
nach Bestimmung der
k-kleinsten $d_{ij}(x)$
- um so größer k desto schwächer
der Einfluss von Ausreißern
- Bsp: k=3



Support-Vector-Machine

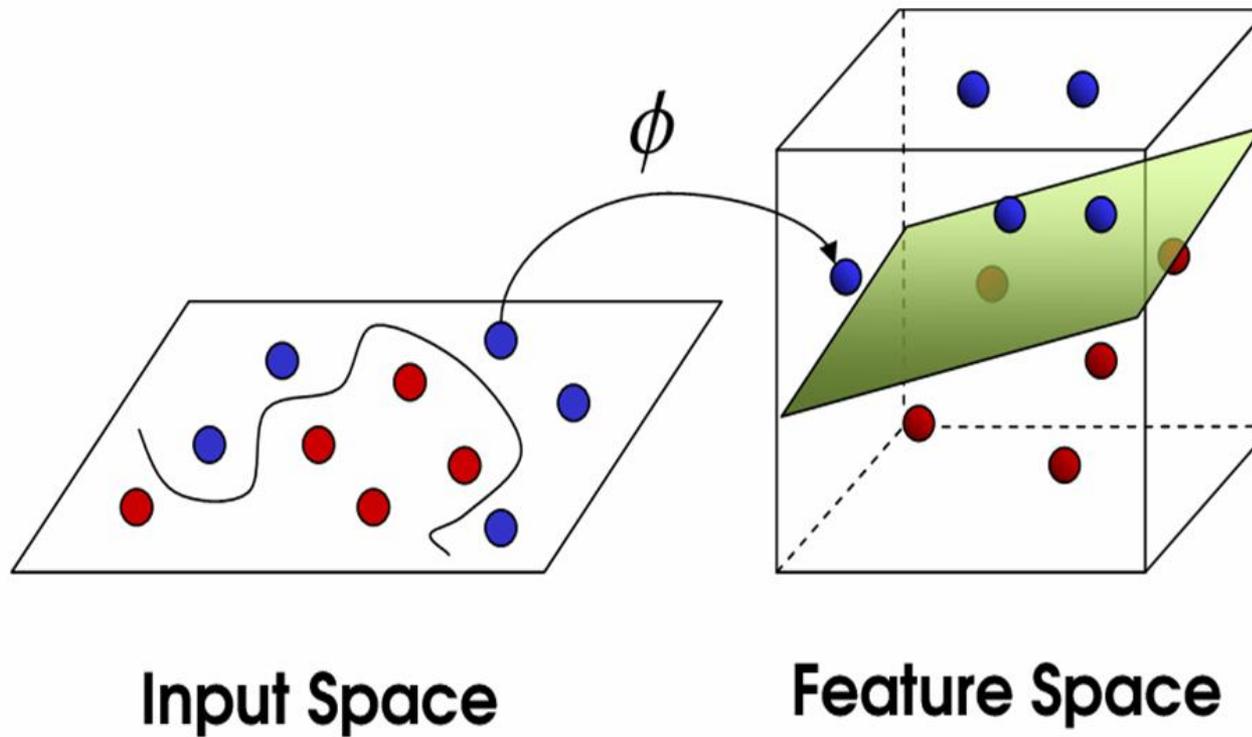
- Erweiterung der linearen Klassifikation
- benutzt Hyperebene als Trennfunktion
- es ist möglich auch nicht linear separierbare Datenobjekte durch eine Hyperebene von einander zu trennen -> Instanzraum wird in einen höher dimensional Raum abgebildet
- die Hyperebene wird maximal diskriminierende Hyperebene genannt

Maximal diskriminierende Hyperebene

- für die Objekte der beiden Klassen wird die konvexe Hülle gebildet
- die Ebene, die am weitesten von allen Objekten entfernt ist, heißt maximal diskriminierende Hyperebene
- liegt orthogonal zu den kürzesten Linien, die die beiden konvexen Hüllen verbindet
- die Objekte mit dem geringsten Abstand zu der maximal diskriminierenden Hyperebene werden support vectors / Stützvektoren genannt

Nichtlineare Klassengrenzen

Principle of Support Vector Machines (SVM)



Hauptkomponentenanalyse

Principal Component Analysis

Problem

- der Datensatz bezogen auf die Anzahl der Attribute ist meist zu groß
- wähle aus n Attributen die k ‚wesentlichen‘ Attribute aus
- Idee: Transformation in einen Sekundärraum mit k Attributen
- Methode: Karhunen-Loeve-Transformation

Karhunen-Loeve-Transformation

X – Datenobjekt mit n -Attributen

A – $n \times n$ Orthonormalmatrix mit a_1, \dots, a_n als Spaltenvektoren

Cov – Kovarianzmatrix zu X

$$X = A Y \rightarrow Y = A^{-1} X = A^T X$$

Beispiel

$$Cov = \begin{pmatrix} 4 & 1 \\ 1 & 4 \end{pmatrix}$$

charakteristische Polynom: $\begin{vmatrix} 4-\lambda & 1 \\ 1 & 4-\lambda \end{vmatrix} = \lambda^2 - 8\lambda + 15 \rightarrow \begin{matrix} \lambda_1 = 5 \\ \lambda_2 = 3 \end{matrix}$

Eigenvektoren: $Cov \cdot x = \lambda_i x \rightarrow \begin{pmatrix} 4 & 1 \\ 1 & 4 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \lambda_i \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$

$$\begin{aligned} \lambda_1 : x_1 = x_2 &\rightarrow e_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \xrightarrow{\text{normiert}} \begin{pmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{pmatrix} \\ \lambda_2 : x_1 = -x_2 &\rightarrow e_2 = \begin{pmatrix} 1 \\ -1 \end{pmatrix} \xrightarrow{\text{normiert}} \begin{pmatrix} \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \end{pmatrix} \end{aligned}$$

$$A = (e_1 \quad e_2) = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{pmatrix}$$

Primärvektor $x = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$

$y = A^T x = \begin{pmatrix} 0 \\ \sqrt{2} \end{pmatrix} \rightarrow$ Attribut im Sekundärraum $a = 1$

Umsetzung in Software

	Naive Bayes	Entscheidungsbaum	KNN	SVM	k-NN	Quelle
MS SQL Server 2008	x	x	x			msdn
Salford Systems - CART 6.0 ProEX		x	x	x		salford-systems.com
KNIME	x	x	x	x		www.knime.org
RapidMiner		x	x	x		rapid-i.com
Auton Lab - Auton Web					x	www.autonlab.org/autonweb

Vor-/Nachteile der Algorithmen

Entscheidungsbaum

- Regel leicht ersichtlich, aber der **Gesamtüberblick schwierig (Baum zu tief)** -> Pruning (siehe Literatur)
- **Attribute mit sehr viel möglichen Werten** ergeben Mehrfachverzweigungen
- zur Klassifikation relevante Attribute befinden sich in **Wurzelnähe**

naive Bayes

- arbeitet auf großen Datenmengen gut, mit hoher Genauigkeit
- nimmt Unabhängigkeit der Attribute an;
wenn fälschlicher Weise -> unzureichende Ergebnisse
- Geschwindigkeit ähnlich wie Entscheidungsbaum, KNN
- kontinuierliche Daten -> als Verteilungsfkt. wird meist fälschlicher Weise Normalverteilung vorausgesetzt

KNN

- tolerant gegenüber Ausreißern
- Attributkombinationen die nicht in der Trainingsmenge vorkommen, stellen kein Problem dar
- erlernte Gewichte sind kaum interpretierbar -> Klassifikationsergebnis nicht erklärbar
- lange Trainingsphase
- kategorielle Daten müssen in metrische Daten umgewandelt werden
- Netz ist individuell bezüglich Anzahl verborgener Schichten und Anzahl Neuronen zu erstellen

k - nearest neighbour

- für metrische und kategorielle Daten geeignet (durch Anpassung der Distanzfunktion)
- keine Lernphase, Daten werden zum Zeitpunkt der Klassifikation ausgewertet, **aber gesamte Trainingsmenge muss für die Klassifikation eines Objektes durchlaufen werden**
- **Aufwand steigt mit steigenden k**

SVM

- schnelle Klassifikation, weil nur Stützvektoren erforderlich
- effektiv auf Daten mit vielen Attributen
- finden der maximal diskriminierenden Hyperebene langwierig

Quellen

- [Los02] Loss, Dirk: Data Mining: Klassifikations- und Clusteringverfahren, 2002, http://www.dirk-loss.de/DM-class-cluster_dloss.pdf
- [Wit01] Witten, Ian H., Frank, Eibe: Data Mining – Praktische Werkzeuge und Techniken für das maschinelle Lernen, 2001, Hanser Verlag, ISBN: 3-446-21533-6
- [Fuh05] Prof. Dr.-Ing. Fuhr, Norbert: Informationssysteme, 2005, www.is.informatik.uni-duisburg.de/courses/dm_ws05, Grundlage: Witten & Frank
- [Wer00] Prof. Dr. Werner, Heinrich: der Forschungsgruppe Neuronale Netzwerke der Universität Kassel angehörig, Neuronale Netzwerke – Vorlesungsmaterial, 2000
- [Sch09] Prof. Dr. Schönherr, Siegfried: HTWK Leipzig, Mustererkennung – Vorlesungsmaterial, 2009

outlook	temperature	humidity	windy	play
overcast	cool	normal	true	yes
overcast	mild	high	true	yes
overcast	hot	high	false	yes
overcast	hot	normal	false	yes
rainy	cool	normal	false	yes
rainy	cool	normal	true	no
rainy	mild	high	false	yes
rainy	mild	normal	false	yes
rainy	mild	high	true	no
sunny	cool	normal	false	yes
sunny	mild	normal	true	yes
sunny	mild	high	false	no
sunny	hot	high	false	no
sunny	hot	high	true	no