



Human Centered Multimedia
Institute of Computer Science

UNA Universität
Augsburg
University

Multimodale Analyse

Wiederholung von Klassifikation

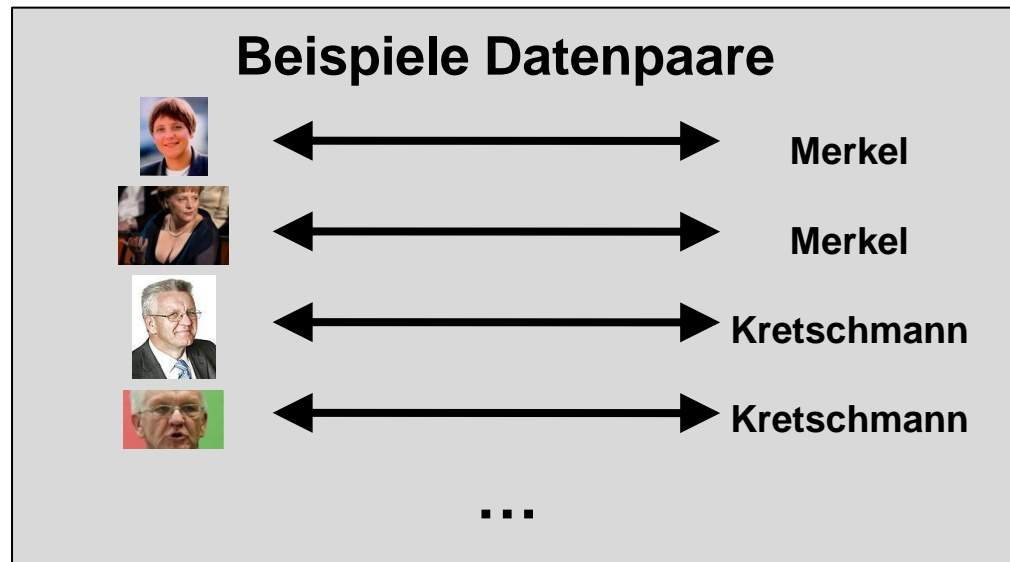
Elisabeth André
Stephan Hammer
Chi-Tai Dang



Human Centered Multimedia
Institute of Computer Science
Augsburg University
Universitätsstraße. 6a
86159 Augsburg, Germany

- Bisher haben wir erfahren, wie man eine Vielzahl von **Merkmalen** bei den unterschiedlichen Daten **berechnet**.
- Außerdem haben wir gelernt, wie man aus dieser Menge die geeigneten **Merkmale** zur Klassifikation **auswählt**.
- Wir werden nun Verfahren kennenlernen die genutzt werden um **Daten** in bestimmte **Klassen einzuordnen**
- Innerhalb dieser Klassen befinden sich dann diejenigen Daten, deren Merkmale sich mehr oder weniger ähneln.
→ Diesen Vorgang nennt man **Klassifikation**

- Anhand von Beispielen erlernt das System eine Zuordnung von Objekten in vorhandene **Klassen**



?

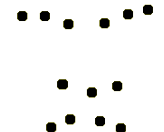
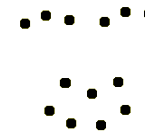
- **Was kann Klassifikation?**
 - Automatisches Sortieren von Äpfeln und Birnen in Güteklassen
 - Einteilung von Kunden in Kundengruppen (z.B. „*Top-Kunden*“, „*Normal-Kunden*“, „*Problemfälle*“ und „*potentielle Wechsler*“)
 - Unterscheidung von Museumsbesuchern anhand ihrer Navigationsmuster („*Schmetterling*“, „*Grashüpfer*“, etc.)
 - Automatische Handschrift- oder Zeichenerkennung: Varianten eines Buchstabens werden der gleichen Klasse zugeordnet
 - Klassifikation von Gesichtern anhand von biometrischen Merkmalen oder sogar Identifikation von Gesichtern („*FBI*“)
 - Erkennung von Emotionsklassen anhand von Sprache, Gesten, Körperhaltung, Gesichtsmimik und biosensorischen Signalen.

Sensoreingabe



Mustererkennung

Feature Extraktion



Klassifikator

Vorhersage

Ärger ○ Freude
+



Human Centered Multimedia
Institute of Computer Science

UNA Universität
Augsburg
University

Multimodale Analyse

Grundlagen der Wahrscheinlichkeitstheorie

Elisabeth André
Stephan Hammer
Chi-Tai Dang



Human Centered Multimedia
Institute of Computer Science
Augsburg University
Universitätsstraße. 6a
86159 Augsburg, Germany

- Eine **Zufallsgröße** ist eine Größe, deren Wert wir nicht exakt kennen bzw. vorhersagen können. Wir können den möglichen Werten nur bestimmte **Wahrscheinlichkeiten** für ihr Auftreten zuordnen.
- Zufallsgrößen sind bei uns **Merkmale** und **Klassen**
- Der Fall, dass eine Zufallsgröße einen bestimmten Wert annimmt nennen wir **Ereignis** oder auch **Beobachtung**.
- Die **Wahrscheinlichkeit** eines Ereignisses A wird bezeichnet als $P(A)$ (P steht dabei für **Probability**).

■ Wahrscheinlichkeiten erfüllen die folgenden **Axiome**:

- Ist A ein Ereignis, dann gilt für die Wahrscheinlichkeit von A :

$$0 \leq P(A) \leq 1$$

- Ist A ein Ereignis, dann gilt für das Gegenereignis von A :

$$P(\bar{A}) = 1 - P(A)$$

- Sind A und B Ereignisse, die nicht gleichzeitig eintreten können:

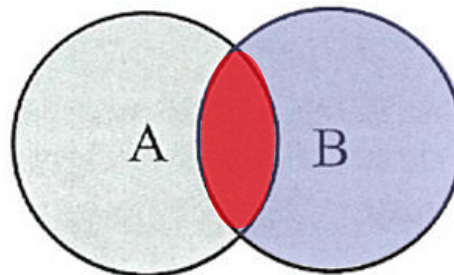
$$P(A \cup B) = P(A) + P(B)$$

- Sind A und B Ereignisse, die auch gleichzeitig eintreten können:

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

- Angenommen wir wollen die Wahrscheinlichkeit für ein Ereignis A bestimmen unter der Bedingung, dass wir schon wissen dass ein anderes Ereignis B eintritt.
- Die **bedingte Wahrscheinlichkeit** eines Ereignisses A gegeben, dass ein Ereignis B gilt, wird berechnet durch:

$$P(A | B) = \frac{P(A \cap B)}{P(B)}, \quad \text{wenn } P(B) > 0$$



- Durch Umformung der Definitionsformel der bedingten Wahrscheinlichkeit entsteht der **Multiplikationssatz**:

$$P(A \cap B) = P(A | B)P(B)$$

- Die Verallgemeinerung des **Multiplikationssatz** lautet:

$$P\left(\bigcap_{i=1}^n A_i\right) = \prod_{i=1}^n P\left(A_i \mid \bigcap_{j=1}^{i-1} A_j\right)$$

- Vgl. Kettenregel:

$$P\left(\bigcap_{i=1}^n A_i\right) = P(A_1)P(A_2 | A_1)P(A_3 | A_1 \cap A_2) \dots P\left(A_n \mid \bigcap_{i=1}^{n-1} A_i\right)$$

- Kennt man nur die bedingten Wahrscheinlichkeiten von einem Ereignis A , sowie die Wahrscheinlichkeiten der bedingenden Ereignisse so kann man die sogenannte **totale Wahrscheinlichkeit** berechnen:

$$P(A) = P(A | B)P(B) + P(A | \bar{B})P(\bar{B})$$

- Verallgemeinert lautet die **totale Wahrscheinlichkeit**:

$$P(A) = \sum_{i=1}^n P(A \cap B_i) = \sum_{i=1}^n P(A | B_i)P(B_i)$$

wenn die Ereignisse B_1, B_2, \dots, B_n eine Partition des Wahrscheinlichkeitsraums bilden und $P(B_i) > 0$ gilt.

- Der **Satz von Bayes** ergibt sich direkt aus der Definition der bedingten Wahrscheinlichkeit und der Definition des Multiplikationssatzes durch einfache Umformung:

$$P(A | B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B | A)P(A)}{P(B)}, \quad \text{wenn } P(A), P(B) > 0$$

- Der Nenner $P(B)$ dieser Formel kann dabei mit Hilfe des Satzes der totalen Wahrscheinlichkeit berechnet werden.

- Die Allgemeine Form des **Satz von Bayes** lautet dann:

$$P(A_i | B) = \frac{P(B | A_i)P(A_i)}{P(B)} = \frac{P(B | A_i)P(A_i)}{\sum_{j=1}^n P(B | A_j)P(A_j)}$$



- Wenn zwei Ereignisse A und B **unabhängig** sind, dann soll die totale Wahrscheinlichkeit von A niemals davon beeinflusst werden, ob wir schon wissen, dass B eintritt:

$$P(A) = P(A | B)$$

- Daraus ergibt sich dann für die Wahrscheinlichkeit dass beide Ereignisse eintreten die folgende Umformung:

$$P(A) = P(A | B) = \frac{P(A \cap B)}{P(B)} \Leftrightarrow P(A \cap B) = P(A) \cdot P(B)$$



Human Centered Multimedia
Institute of Computer Science

UNA Universität
Augsburg
University

Multimodale Analyse

Bayessche Statistik und Naive Bayes Klassifikation

Elisabeth André
Stephan Hammer
Chi-Tai Dang



Human Centered Multimedia
Institute of Computer Science
Augsburg University
Universitätsstraße. 6a
86159 Augsburg, Germany

- In der **Bayesschen Statistik** werden die Regeln für die Klassifikation von Beobachtungen mit dem **Satz von Bayes** als bedingte Wahrscheinlichkeiten formuliert.
- Seien $h_1, h_2, \dots, h_n \in H$ Hypothesen die nicht gleichzeitig auftreten können, von denen aber immer eine eintritt, dann kann man folgende Fragestellungen beantworten

Was ist die **wahrscheinlichste Hypothese** $h_i \in H$, also die **wahrscheinlichste Klassifikation** einer neuen Dateninstanz bei bekannten Trainingsdaten?

- $P(h_i)$ sei die Anfangswahrscheinlichkeit dafür, dass die Hypothese h_i erfüllt ist (**A-Priori Probability** von h_i).
- $P(D)$ sei die Anfangswahrscheinlichkeit dafür, dass die Daten D beobachtet werden.
- $P(D|h_i)$ sei die Wahrscheinlichkeit dafür, dass man D beobachtet, sofern die Hypothese h_i erfüllt ist.
- $P(h_i|D)$ sei die Wahrscheinlichkeit dafür, dass h_i bei den Daten D gilt (**A-Posteriori Probability** von h_i).

- Die **A-Priori Wahrscheinlichkeit** für ein Ereignis A beschreibt die Wahrscheinlichkeit für A , wenn keine Informationen über andere Ereignisse vorliegen.
- In der Bayesschen Statistik modellieren A-Priori Wahrscheinlichkeiten unser **Vorwissen** über die Häufigkeit einer Klasse und das Auftreten von Merkmalen.
 - $P(\text{Klasse}=\text{apfel}) = 0.25$ (d.h. 25% der Objekte sind Äpfel)
 - $P(\text{Klasse}=\text{birne}) = 0.40$ (d.h. 40% der Objekte sind Birnen)
 - $P(\text{Klasse}=\text{orange}) = 0.35$ (d.h. 35% der Objekte sind Orangen)

- Die **A-Posteriori-Wahrscheinlichkeit** für ein Ereignis A beschreibt die Wahrscheinlichkeit für A , wenn wir auch zusätzlich wissen, dass ein bestimmtes Ereignis eintritt.
- In der Bayesschen Statistik modellieren damit dann die **Zusammenhänge** zwischen Klassen und Merkmalen:
 - $P(\text{Form=rund} \mid \text{Klasse=orange}) = 1.0$
(d.h. wenn ich weiß, dass ich eine Orange betrachte, dann ist die Wahrscheinlichkeit dafür, dass sie rund ist 100%)

- Man betrachtet eine Menge H von Hypothesen und ist daran interessiert, die wahrscheinlichste Hypothese h aus H für die gegebenen Daten D zu finden.
- Die **Maximum A-Posteriori Hypothese** ist

$$h_{MAP} = \arg \max_{h \in H} P(h | D) =$$

$$\arg \max_{h \in H} \frac{P(D | h)P(h)}{P(D)} = \quad (\text{Mit dem Satz von Bayes})$$

$$\frac{1}{P(D)} \arg \max_{h \in H} P(D | h)P(h) \quad (\text{Herausziehen der Konstante } 1/P(D), \text{ die von } h \text{ unabhängig ist})$$

- Wenn kein Vorwissen über die Wahrscheinlichkeiten der einzelnen Hypothesen besteht, macht man die spezielle Annahme, dass jede Hypothese aus H die gleiche Anfangswahrscheinlichkeit hat, d.h.

$$P(h_i) = P(h_j) \quad \forall h_i, h_j \in H$$

- Eine Hypothese, die $P(D|h)$ maximiert, wird dann die **Maximum Likelihood Hypothesis** h_{ML} genannt.

$$h_{ML} = \arg \max_{h \in H} P(D | h)$$

- Bei der **Naiven Bayes Klassifikation** fragen wir nach den Wahrscheinlichkeiten der Hypothesen, wenn wir wissen, dass ein bestimmtes Ereignis eingetroffen ist.
- **Gegeben:**
 - Eine endliche Menge von **Klassen** $V = \{v_1, v_2, \dots, v_m\}$
 - Eine Menge von klassifizierten **Trainingsbeispielen**
 - Eine **Instanz** als ein Tupel (a_1, a_2, \dots, a_n) von **Attributen**
- **Gesucht:**
 - Die wahrscheinlichste Klasse für die neue Instanz

- Die wahrscheinlichste Klasse v_{MAP} wird berechnet durch:

$$v_{MAP} = \arg \max_{v_j \in V} P(v_j | a_1, a_2 \dots a_n)$$

- Anwendung des Satz von Bayes ergibt diesen Ausdruck:

$$v_{MAP} = \arg \max_{v_j \in V} \frac{P(a_1, a_2 \dots a_n | v_j) P(v_j)}{P(a_1, a_2 \dots a_n)}$$

- Herausziehen der Normalisierungskonstante ergibt:

$$v_{MAP} = \arg \max_{v_j \in V} P(a_1, a_2 \dots a_n | v_j) P(v_j)$$

- Wir schätzen die Terme $P(a_1, a_2 \dots a_n | v_j)$ und $P(v_j)$ auf Basis der bekannten Trainingsdaten durch **Abzählen**.
- Abschätzung von $P(v_j)$ durch Bestimmung der relativen **Häufigkeit**, der Klasse v_j innerhalb der Trainingsdaten.
- Abschätzung von $P(a_1, a_2 \dots a_n | v_j)$ ist **problematisch**
 - Bei wenig Trainingsdaten ist es nicht ungewöhnlich, dass wir keine solche Kombinationen haben, d.h. wir würden für diesen Ausdruck dann lediglich 0 erhalten.
 - Gute Abschätzung nur bei **riesiger Trainingsmenge**.

- Zur Vereinfachung nimmt man an, dass alle $a_1, a_2 \dots a_n$ **bedingt unabhängig** voneinander sind (**naiver Ansatz**)

$$P(a_1, a_2 \dots a_n \mid v_j) = \prod_{i=1}^n P(a_i \mid v_j)$$

- Einsetzen ergibt dann den **Naiven Bayes Klassifikator**

$$v_{NB} = \arg \max_{v_j \in V} P(v_j) \prod_{i=1}^n P(a_i \mid v_j)$$

- Die **A-Posteriori-Wahrscheinlichkeit** für die Klasse v_{MAP} mit der höchsten Wahrscheinlichkeit wird gegeben durch

$$P(v_{MAP} | a_1, a_2 \dots a_n) = \frac{P(a_1, a_2 \dots a_n | v_{MAP})P(v_{MAP})}{\sum_{i=1}^m P(a_1, a_2 \dots a_n | v_i)P(v_i)}$$

- Die **wichtigsten Formeln** in der Bayeschen Statistik:
 - **Naive Bayes Klassifikator:**

$$v_{NB} = \arg \max_{v_j \in V} P(v_j) \prod_{i=1}^n P(a_i | v_j)$$

- **A-Posteriori- Wahrscheinlichkeit:**

$$P(v_{MAP} | a_1, a_2 \dots a_n) = \frac{P(a_1, a_2 \dots a_n | v_{MAP}) P(v_{MAP})}{\sum_{i=1}^m P(a_1, a_2 \dots a_n | v_i) P(v_i)}$$

- Das Schätzen der Wahrscheinlichkeit $P(v_j)$ geschieht durch Abzählen der relativen Häufigkeiten der Klassen innerhalb der vorhandenen Menge von Trainingsdaten.

- **Beispiel:**

- $P(\text{Klasse}=A)=4/5$
- $P(\text{Klasse}=B)=1/5$

ID	Form	Farbe	Klasse
1	rund	blau	A
2	rund	grün	A
3	rund	gelb	A
4	eckig	grün	A
5	oval	weiß	B

- Die bedingten Wahrscheinlichkeiten $P(a_i | v_j)$ werden als relative Häufigkeiten durch folgende Formel abgeschätzt

$$P(a_i | v_j) = \frac{n_{i,j}}{n_j}$$

- Wobei n_j die Anzahl der Trainingsbeispiele für v_j ist und $n_{i,j}$ die Anzahl der Instanzen mit Attribut a_i und Klasse v_j .

- Beispiel:**

- $P(\text{Form}=\text{rund} | \text{Klasse}=\text{A})=3/4$
- $P(\text{Farbe}=\text{grün} | \text{Klasse}=\text{A})=2/4$
- $P(\text{Form}=\text{oval} | \text{Klasse}=\text{A})=0/4$

ID	Form	Farbe	Klasse
1	rund	blau	A
2	rund	grün	A
3	rund	gelb	A
4	eckig	grün	A
5	oval	weiß	B

- Bei kontinuierlichen metrischen Attributen:
 - Abschätzung durch eine **Diskrete Approximation** durch diskrete Intervalle:
 - $P(9.0 < \text{Durchmesser} \leq 9.5 \mid \text{Orange}) = 10\%$
 - $P(9.5 < \text{Durchmesser} \leq 10.0 \mid \text{Orange}) = 30\%$
 - $P(10.0 < \text{Durchmesser} \leq 10.5 \mid \text{Orange}) = 30\%$
 - $P(10.5 < \text{Durchmesser} \leq 11.0 \mid \text{Orange}) = 10\%$
 - $P(11.0 < \text{Durchmesser} \mid \text{Orange}) = 5\%$
 - Abschätzung durch **Wahrscheinlichkeitsdichtefunktionen** (meist Berechnung nach Normalverteilung):

$$P(x \mid C) = \frac{1}{\sqrt{2\pi}\sigma} e^{\frac{-(x-\mu)^2}{2\sigma^2}} \quad \text{wobei} \quad \mu = \frac{\sum_{x \in TR} x}{|TR|} \quad \text{und} \quad \sigma = \sqrt{\frac{\sum_{x \in TR} (x - \mu)^2}{|TR|}}$$

- Die Instanzen, welche betrachtet werden sind Tage.
- Jeder Tag besteht aus einem Tupel der Attribute
 - *Outlook*
 - *Temperature*
 - *Humidity*
 - *Wind*
- Ein neuer Tag soll dahingehend klassifiziert werden, ob Luigi an dem neuen Tag Tennis spielen wird oder nicht.
- Die Trainingsdatenmenge ist in einer Tabelle gegeben.



Naive Bayes Klassifikation

Beispiel des Tennisspielers

Day	Outlook	Temperature	Humidity	Wind	Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- Klassifiziert werden soll nun die folgende neue Instanz:
(*Outlook = sunny, Temperature = cool, Humidity = high, Wind = strong*)
- Wir wenden den Naive-Bayes Klassifikator an um die wahrscheinlichste Klasse (*yes* oder *no*) zu berechnen:

$$v_{NB} = \arg \max_{v_j \in V} P(v_j) \prod_i P(a_i | v_j)$$

$$= \arg \max_{v_j \in \{yes, no\}} P(v_j) \prod_i P(a_i | v_j)$$

$$= \arg \max_{v_j \in \{yes, no\}} P(v_j) \begin{pmatrix} P(Outlook = sunny | v_j) \\ P(Temperature = cool | v_j) \\ P(Humidity = high | v_j) \\ P(Wind = strong | v_j) \end{pmatrix}$$



- Abschätzen der Wahrscheinlichkeiten für *yes* und *no*

$$P(\text{Tennis} = \text{yes}) = 9/14$$

$$P(\text{Tennis} = \text{no}) = 5/14$$

- Bedingte Wahrscheinlichkeiten:

$$P(\text{Outlook} = \text{sunny} | \text{Tennis} = \text{yes}) = 2/9$$

$$P(\text{Temperature} = \text{cool} | \text{Tennis} = \text{yes}) = 3/9$$

$$P(\text{Humidity} = \text{high} | \text{Tennis} = \text{yes}) = 3/9$$

$$P(\text{Wind} = \text{strong} | \text{Tennis} = \text{yes}) = 3/9$$

$$P(\text{Outlook} = \text{sunny} | \text{Tennis} = \text{no}) = 3/5$$

$$P(\text{Temperature} = \text{cool} | \text{Tennis} = \text{no}) = 1/5$$

$$P(\text{Humidity} = \text{high} | \text{Tennis} = \text{no}) = 4/5$$

$$P(\text{Wind} = \text{strong} | \text{Tennis} = \text{no}) = 3/5$$

Day	Outlook	Tem	Hum	Wind	Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- Die Vorhersage des Zielwerts für die neue Instanz ergibt:

$$P(\text{yes}) \begin{pmatrix} P(\text{sunny} | \text{yes}) \\ P(\text{cool} | \text{yes}) \\ P(\text{high} | \text{yes}) \\ P(\text{strong} | \text{yes}) \end{pmatrix} \approx .0053, \quad P(\text{no}) \begin{pmatrix} P(\text{sunny} | \text{no}) \\ P(\text{cool} | \text{no}) \\ P(\text{high} | \text{no}) \\ P(\text{strong} | \text{no}) \end{pmatrix} \approx .0206$$

- Die Naive Bayes Klassifizierung wählt also *no* als das wahrscheinlichste Ergebnis für die neue Instanz.
- Die A-Posteriori Wahrscheinlichkeit für *no* kann man dann mittels der totalen Wahrscheinlichkeit berechnen.

$$P(\text{yes} | \dots) \approx \frac{0.0053}{0.0053 + 0.0206} \approx \frac{0.0053}{0.0259} \approx 0.205$$

$$P(\text{no} | \dots) \approx \frac{0.0206}{0.0053 + 0.0206} \approx \frac{0.0206}{0.0259} \approx 0.795$$

■ Problem:

- Die bisherige Methode zur Abschätzung führt zu einer Unterschätzung bei kleinem $n_{i,j}$ oder wenn $n_{i,j} \rightarrow 0$

■ Beispiel:

- Wenn wir keine Trainingsbeispiele haben, bei denen starker Wind vorlag und nicht Tennis gespielt wurde:

$$P(Wind = strong \mid PlayTennis = no) = 0$$

- Dieser Wert dominiert bei einer neuen Instanz mit *Wind=strong* alle anderen Beobachtungen und führt zur Unterschätzung der Gesamtwahrscheinlichkeit.

■ M-Estimate of Probability:

$$\frac{n_{i,j} + m \cdot p}{n_j + m}$$

- m Anzahl der neu hinzugefügten virtuellen Trainingsbeispiele
- n_j Anzahl der Trainingsbeispiele aus der Klasse v_j
- $n_{i,j}$ Anzahl der Beispiele mit Attribut a_i und Klasse v_j
- p A-Priori Wahrscheinlichkeit für Beispiele mit dem Attribut a_i und Klasse v_j . Um p zu wählen, nimmt man eine **Gleichverteilung** an, d.h falls das Attribut a_i k mögliche Werte hat, dann wählt man $1/k$

$$\frac{n_{i,j} + m \cdot p}{n_j + m}$$

■ Beispiel:

- Anzahl der Trainingsbeispiele mit *overcast* bei *no*:

$$n_{\text{overcast}, \text{no}} = 0$$

- Anzahl der Trainingsbeispiele für die Klasse *no*:

$$n_{\text{no}} = 5$$

- Annahme der Gleichverteilung für Attribut *outlook*:

$$p = 1/3$$

- M-Estimate für $P(\text{overcast} \mid \text{no})$ bei verschiedenen m :

$$m = 0 \rightarrow P(\text{overcast} \mid \text{no}) = (0 + 1/3 \cdot 0) / (5 + 0) = 0$$

$$m = 1 \rightarrow P(\text{overcast} \mid \text{no}) = (0 + 1/3 \cdot 1) / (5 + 1) = 0.055$$

$$m = 100 \rightarrow P(\text{overcast} \mid \text{no}) = (0 + 1/3 \cdot 100) / (5 + 100) = 0.318$$



Human Centered Multimedia
Institute of Computer Science

UNA Universität
Augsburg
University

Multimodale Analyse

Naive Bayes - Bonusaufgabe

Elisabeth André
Stephan Hammer
Chi-Tai Dang



Human Centered Multimedia
Institute of Computer Science
Augsburg University
Universitätsstraße. 6a
86159 Augsburg, Germany

- EINZELABGABEN!!
- Sie nehmen jetzt an einem kurzen Experiment teil, das mit dem Naive-Bayes Klassifikator zu tun hat. Nächste Vorlesung erfahren Sie was es genau damit auf sich hat.
- Natürlich werden Ihre Daten anonym behandelt und es werden keine personenbezogenen Daten gespeichert.
- Bitte beantworten Sie zunächst die folgenden Fragen.
- Nummerieren Sie bitte die Fragen für die Auswertung.

1. Welches Bild gefällt Ihnen besser?

a)



b)



2. Welches Paar Schuhe halten sie für bequemer?

a)



b)



c)



3. Welche Art von Komödie gefällt Ihnen besser?

a)



b)



Naive Bayes Klassifikation

Bonusaufgabe

4. Wenn Sie ein Auto wären, welches Auto käme ihrer Vorstellung am nächsten?

a)



b)



5. Sehen Sie sich gerne im Spiegel an?

- a) Nein, ich mag das gar nicht.
- b) Es ist ok, mich im Spiegel zu sehen.
- c) Ja, ich mag meinen Anblick im Spiegel.

6. Was trifft eher zu?

- a) Nutzen Sie eher Android?
- b) Oder Nutzen Sie eher Apple?
- c) Oder Nutzen Sie eher Windows?



Human Centered Multimedia
Institute of Computer Science

UNA Universität
Augsburg
University

Multimodale Analyse

Naive Bayes - Textklassifikation

Elisabeth André
Stephan Hammer
Chi-Tai Dang



Human Centered Multimedia
Institute of Computer Science
Augsburg University
Universitätsstraße. 6a
86159 Augsburg, Germany

- NB Klassifikatoren in der **Textanalyse** weit verbreitet

Vorteile:

- Hohe Trainings- und Klassifizierungsgeschwindigkeit
- NB verbessert sich mit jeder neu klassifizierten Instanz
- Einsatz zur Kategorisierung von Dokumenten
 - Kategorien: z.B. Finanzen, Sport
 - Genres: z.B. Nachrichten, Kritiken
 - Meinungen: z.B. Positive, Negative
 - Labels: z.B. Spam und NoSpam

Die guten Klassifizierungseigenschaften von NB-Klassifikatoren machen sie insbesondere beim Einsatz in **Spam-Filtern** sehr beliebt

■ Gegeben:

- Trainingsmenge von Emails, von denen bekannt ist, ob sie in die Klasse *Spam* oder in die Klasse *No-Spam* klassifiziert wurden.
- Neue Email, die in *Spam* oder *No-Spam* klassifiziert werden soll

■ Gesucht:

- Klassifiziere die neue Email in die Klasse *Spam* oder *No-Spam*

■ Attribute:

- Häufigkeit bestimmter Worte aus einem Vokabular

- Vereinfachende Annahme:
 - Klassifikation ist unabhängig von der Aufeinanderfolge von Worten
 - Quasi eine Vernachlässigung der grammatikalischen Regeln
- Naive Bayes Klassifikation

positions: alle Textpositionen im zu klassifizierenden Dokument

w_i : Wort an der Position i

$$v_{NB} = \arg \max_{v_j \in V} P(v_j) \prod_{i \in \text{positions}} P(w_i | v_j)$$

- Abschätzen von $P(v_j)$ wie gewöhnlich durch Abzählen
- Abschätzen von $P(w_k/v_j)$ mit m-estimate

- Gleichverteilung der Wahrscheinlichkeiten aller Wörter w_1, \dots, w_l im Vokabular, da jedes Wort w_k einmal im Vokabular vorkommt
- Anzahl der neuen virtuellen Beispiele (m) ist Größe des Vokabulars

$$P(w_k | v_j) \approx \frac{n_k + 1}{n + |V|}$$

- n ist die Anzahl der Wortpositionen in allen Trainingsbeispielen der Klasse v_j
- n_k ist die Häufigkeit der Auftreten des Wortes w_k unter diesen n Wortpositionen.
- $|V|$ ist die Anzahl unterschiedlicher Worte in den Trainingsdaten

1. **LearnNaiveBayesText** (E , V)

- E Menge aller Beispieldokumente und V Menge aller Klassen
- Extrahiere ein Vokabular Voc aus der Trainingsbeispiel-Menge
- Berechne die Abschätzungen für $P(v_j)$ und $P(w_k|v_j)$ wie folgt

1. Für jede der Klassen v_j :

- $D_j \leftarrow$ Menge der Dokumente aus der Klasse v_j
- $P(v_j) \leftarrow |D_j| / |E|$
- $T_j \leftarrow$ Ein einziges Dokument bestehend aus allen D_j
- $n_j \leftarrow$ Anzahl der Wörter im Dokument T_j
- Für jedes Wort w_k aus dem Vokabular
 - $n_{k,j} \leftarrow$ Häufigkeit von Wort w_k im Dokument T_j
 - $P(w_k|v_j) \leftarrow (n_{k,j}+1) / (n_j + |Voc|)$

2. **ClassifyNaiveBayesText(D)**

- Sei D das zu klassifizierende Dokument und *positions* die Menge der Wortpositionen in D besetzt mit einem Wort aus dem Vokabular.
- Sei w_i das Wort an Position i im Dokument D dann berechne

$$v_{NB} = \arg \max_{v_j \in V} P(v_j) \prod_{i \in \text{positions}} P(w_i | v_j)$$

	Doc	Words	Class
Training	1	free all free	spam
	2	free lottery free	spam
	3	free trip	spam
	4	free erasmus place	no_spam
Test	5	free free free erasmus place	?

$$P(\text{spam}) = 3/4$$

$$P(\text{no_spam}) = 1/4$$

$$P(\text{free}|\text{spam}) = (5+1)/(8+6) = 3/7$$

$$P(\text{place}|\text{spam}) = (0+1)/(8+6) = 1/14$$

$$P(\text{erasmus}|\text{spam}) = (0+1)/(8+6) = 1/14$$

$$P(\text{free}|\text{no_spam}) = (1+1)/(3+6) = 2/9$$

$$P(\text{place}|\text{no_spam}) = (1+1)/(3+6) = 2/9$$

$$P(\text{erasmus}|\text{no_spam}) = (1+1)/(3+6) = 2/9$$

$$P(w_k|v_j) \leftarrow (n_{k,j}+1) / (n_j + |V|)$$

$n_{k,j}$ Häufigkeit von Wort w_k in Dokument D_j
 n_j Anzahl der Wörter in Dokument D_j

$$v_{NB} = \arg \max_{v_j \in V} P(v_j) \prod_{i \in \text{positions}} P(w_i | v_j)$$

$$P(\text{spam}|\text{doc5}) \sim 3/4 * (3/7)^3 * 1/14 * 1/14 = 0,0098$$

$$P(\text{no_spam}|\text{doc5}) \sim 1/4 * (2/9)^3 * 2/9 * 2/9 = 0,0023$$

- **Experiment von Joachims, 1996:**
 - Sammlung von je 1.000 Artikeln aus 20 Newsgroups
 - Davon 2/3 Trainingsbeispiele und 1/3 Testbeispiele
- **Vokabularbildung:**
 - Entferne die 100 häufigsten Worte („the“, „of“, ...)
 - Entferne Worte, die weniger als 3 mal auftauchen
- **Erkennungsrate:**
 - Zufälliges Raten: 5%
 - Mit Algorithmus: 89%

- Die vorgestellten Bayes-Methoden ermitteln A-Posteriori Wahrscheinlichkeiten für Hypothesen basierend auf den angenommenen oder irgendwie abgeschätzten A-Priori Wahrscheinlichkeiten und den neu beobachteten Daten.
- Mit jeder neu klassifizierten Instanz verbessert sich der Naive Bayes Klassifikator und bietet eine gute Performanz
- NB gibt sogar eine Wahrscheinlichkeitsverteilung an.
- Der Naive Bayes-Klassifikator basiert auf der Annahme, die Attribute seien bedingt unabhängig voneinander.
- Obwohl diese Annahme in der Praxis häufig verletzt wird liefert der Ansatz doch oft gute Klassifikationsresultate.

- T. Mitchell: Machine Learning
- Russel, Norvig: Artificial Intelligence: A Modern Approach.
- C.M. Bishop: Pattern Recognition and Machine Learning.

