3. Bayes'sche Klassifikation

Definitionen:

D: Lernstichprobe

T: Teststichprobe

S: Gesamtstichprobe $S = D \cup T$ mit $D \cap T = \emptyset$

Probabilistische Beschreibung von Klassen und Merkmalen

Mathematisches Modell: Umwelt als Zufallsgenerator, der Paare (**m**, ω) erzeugt. Die Gesamtheit (das Ensemble) dieser Paare mit ihrer Variabilität und ihren Häufigkeiten wird durch eine Wahrscheinlichkeitsverteilung (WV) beschrieben.

Zufallsgenerator $P(\mathbf{m}, \omega)$ $\longrightarrow \omega$

Der Zufallsgenerator liefere dabei unabhängige Ausgaben in dem Sinne, dass die Verbundwahrscheinlichkeit von N Ausgaben (\mathbf{m}^n , ω^n) gleich dem Produkt $\prod_{n=1}^N P(\mathbf{m}^n, \omega^n)$ ist.

Probabilistische Beschreibung von Klassen und Merkmalen

$$P(\mathbf{m}, \omega) = p(\mathbf{m} \mid \omega)P(\omega) = P(\omega \mid \mathbf{m})p(\mathbf{m})$$

Verbund-WV

$$P(\omega) = \int_{\mathbb{IM}} P(\mathbf{m}, \omega) \, d\mathbf{m}$$

A Priori WV

$$p(\mathbf{m} \mid \omega)$$

Klassenspezifische Merkmals-WV, "Likelihood"

$$p(\mathbf{m}) = \sum_{\omega = \omega_1}^{\omega_c} P(\mathbf{m}, \omega) = \sum_{\omega = \omega_1}^{\omega_c} p(\mathbf{m} \mid \omega) P(\omega)$$

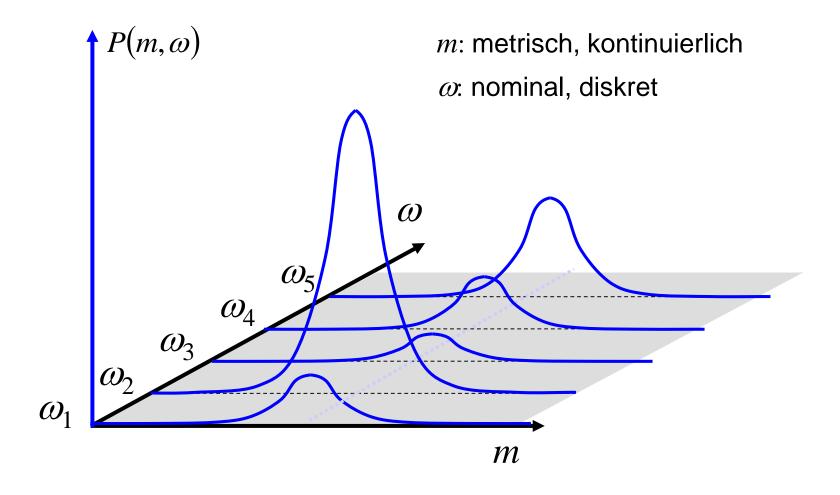
Gesamt-WV der Merkmale

Bayes'sche Formel:

$$P(\omega \mid \mathbf{m}) = \frac{p(\mathbf{m} \mid \omega)P(\omega)}{p(\mathbf{m})}$$

A Posteriori WV

Beispiel: Wahrscheinlichkeitsverteilung gemischter Größen



Entscheidungsraum (Klassifikationsraum): IK

$$\Omega/\sim = \{\omega_1, ..., \omega_c\}$$

$$\Omega/\sim \rightarrow \{\mathbf{\omega}_1, ..., \mathbf{\omega}_c\}$$

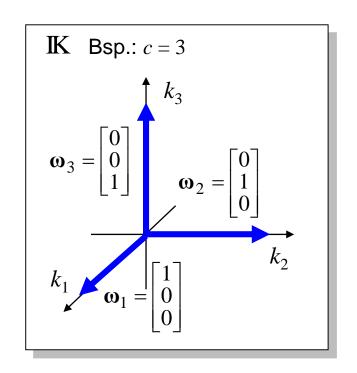
$$\omega_i \mapsto \mathbf{\omega}_i, \quad i = 1, ..., c$$

$$\omega_i \Leftrightarrow \mathbf{\omega}_i \coloneqq [0, ..., 1, ..., 0]^T \in \mathbb{R}^c$$

$$i\text{-te Komponente}$$

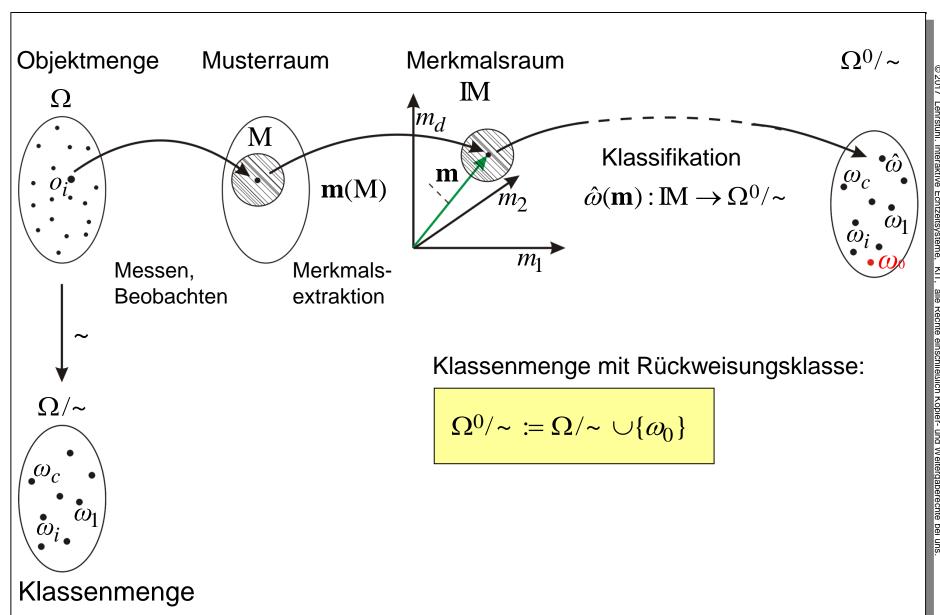
 ω_i : Zielvektor (target vector)

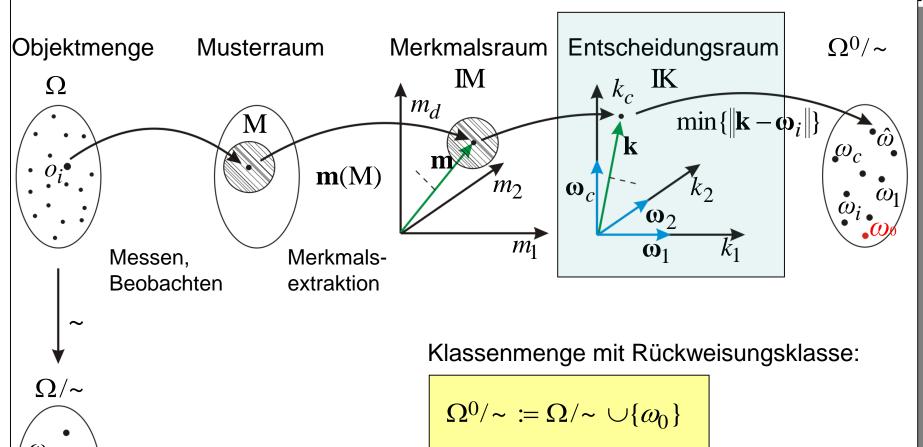
$$\|\mathbf{\omega}_i\| = 1$$
 $\mathbf{\omega}_i^{\mathrm{T}} \mathbf{\omega}_i = \delta_i^j$



Ordinalität der Nummerierung der ω_i suggeriert Nachbarschaften zwischen den Klassen, die mit Falle einer nominalen Klassenstruktur keine Bedeutung haben.

Die Vektorrepräsentation eliminiert diesen Anschein.





Minimale Distanz Klassifikation:

$$\hat{\omega}(\mathbf{m}) = \omega_j \text{ mit } j = \underset{i}{\operatorname{argmin}} \{ \| \mathbf{k}(\mathbf{m}) - \mathbf{\omega}_i \| \}$$

Klassenmenge

Entscheidungsvektor (Klassifikationsvektor):

$$\mathbf{k}(\mathbf{m}) : \mathbb{IM} \to \mathbb{IK} = \operatorname{span}(\boldsymbol{\omega}_1, ..., \boldsymbol{\omega}_c) = \mathbb{IR}^c$$

Komponenten: Entscheidungsfunktionen

$$k_i(\mathbf{m}), i = 1,...,c$$

Alle Klassifikatoren (und mithin alle Entscheidungsgrenzen im Merkmalsraum) lassen sich mit diesem Formalismus beschreiben.

Parametrische Entscheidungsfunktionen: $k(m;\theta)$

 $\mathbf{\theta} = [\theta_1, ..., \theta_k]^T \in \Theta$: Parametervektor

• : Parameterraum

θ anhand der Daten D bestimmen bedeutet: Lernen aus Beispielen.

Ansatz: Lernen anhand der Stichprobe $D = \{\mathbf{m}_1, ..., \mathbf{m}_N\}$ mit bekannter Klassenzugehörigkeit \Leftrightarrow Bestimmung der Parameter θ derart, dass $\mathbf{k}(\mathbf{m})$ die wahren $\mathbf{\omega}(\mathbf{m}_i) \in \{\mathbf{\omega}_1, ..., \mathbf{\omega}_c\}$ für alle $\mathbf{m}_i \in D$ in gewissem Sinne möglichst gut approximiert.

Zusammenhang:

Entscheidungsgebiete R_i im Merkmalsraum:

$$R_{i} \subset IM \Leftrightarrow \left\{ \mathbf{k} \middle| \|\mathbf{k} - \mathbf{\omega}_{i} \| < \|\mathbf{k} - \mathbf{\omega}_{j} \|, \forall j \neq i \right\} \subset IK$$

$$R_{i} = \left\{ \mathbf{m} \middle| \|\mathbf{k}(\mathbf{m}; \mathbf{\theta}) - \mathbf{\omega}_{i} \| < \|\mathbf{k}(\mathbf{m}; \mathbf{\theta}) - \mathbf{\omega}_{j} \|, \forall j \neq i \right\}$$

$$\partial R_{i} = \left\{ \mathbf{m} \middle| \|\mathbf{k}(\mathbf{m}; \mathbf{\theta}) - \mathbf{\omega}_{i} \| \leq \|\mathbf{k}(\mathbf{m}; \mathbf{\theta}) - \mathbf{\omega}_{j} \|, \forall j \neq i \right\} \setminus R_{i}$$

k ist parametrisiert durch $\theta \Rightarrow \partial R_i$ sind parametrisiert durch θ .

Struktur und Parameter von $\mathbf{k}(.,\theta)$ legen die Grenzen $\partial \mathbf{R}_i$ der Entscheidungsgebiete in IM fest.

Entscheidungsvektor mit kleinster mittlerer quadratischer Abweichung

$$f := \mathrm{E}\{\|\mathbf{k}(\mathbf{m}) - \boldsymbol{\omega}\|^2\} = \sum_{i=1}^{c} \int_{\mathbf{M}I} \|\mathbf{k}(\mathbf{m}) - \boldsymbol{\omega}_i\|^2 P(\mathbf{m}, \boldsymbol{\omega}_i) d\mathbf{m} \xrightarrow{!} \min_{\mathbf{k}(\mathbf{m})}$$

Problem der Variationsrechnung.

Sei k die optimale Lösung. Dann gilt:

$$f(\mathbf{k} + \delta \mathbf{k}) > f(\mathbf{k}), \ \forall \delta \mathbf{k} \neq \mathbf{0}$$

k und δ **k** sind Funktionen von **m**.

$$f(\mathbf{k} + \delta \mathbf{k}) = \mathbf{E} \left\{ \| \mathbf{k} + \delta \mathbf{k} - \mathbf{\omega} \|^{2} \right\}$$

$$= \mathbf{E} \left\{ \| \mathbf{k} - \mathbf{\omega} \|^{2} \right\} + 2\mathbf{E} \left\{ \delta \mathbf{k}^{\mathrm{T}} (\mathbf{k} - \mathbf{\omega}) \right\} + \mathbf{E} \left\{ \| \delta \mathbf{k} \|^{2} \right\}$$

$$f(\mathbf{k}) = \mathbf{E} \left\{ \| \mathbf{k} - \mathbf{\omega} \|^{2} \right\}$$

Entscheidungsvektor mit kleinster mittlerer quadratischer Abweichung

$$f(\mathbf{k} + \delta \mathbf{k}) > f(\mathbf{k}) \Leftrightarrow E\{\|\delta \mathbf{k}\|^2\} + 2E\{\delta \mathbf{k}^{\mathrm{T}}(\mathbf{k} - \mathbf{\omega})\} > 0, \forall \delta \mathbf{k} \neq \mathbf{0}$$

Die Ungleichung ist erfüllt, wenn gilt: $E\{\delta \mathbf{k}^{T}(\mathbf{k} - \mathbf{\omega})\} = 0$

$$E\left\{\delta\mathbf{k}^{\mathrm{T}}(\mathbf{k}-\boldsymbol{\omega})\right\}=0$$

$$E\{\delta \mathbf{k}^{\mathrm{T}}(\mathbf{k} - \mathbf{\omega})\} = \int_{\mathbf{IM}} \sum_{i=1}^{c} \delta \mathbf{k}^{\mathrm{T}}(\mathbf{k} - \mathbf{\omega}_{i}) P(\mathbf{m}, \mathbf{\omega}_{i}) d\mathbf{m} =$$

$$= \int_{\mathbf{IM}} \delta \mathbf{k}^{\mathrm{T}} \left[\sum_{i=1}^{c} (\mathbf{k} - \mathbf{\omega}_{i}) P(\mathbf{\omega}_{i} | \mathbf{m}) \right] p(\mathbf{m}) d\mathbf{m} = 0$$

Dieser Ausdruck wird für alle möglichen $\delta \mathbf{k}$ zu Null, falls der geklammerte Ausdruck Null ist.

Entscheidungsvektor mit kleinster mittlerer quadratischer Abweichung

$$\sum_{i=1}^{c} (\mathbf{k} - \mathbf{\omega}_i) P(\mathbf{\omega}_i \mid \mathbf{m}) = \mathbf{k} \sum_{i=1}^{c} P(\mathbf{\omega}_i \mid \mathbf{m}) - \sum_{i=1}^{c} \mathbf{\omega}_i P(\mathbf{\omega}_i \mid \mathbf{m}) = 0$$

Der optimale Entscheidungsvektor lautet somit:

$$\mathbf{k}(\mathbf{m}) = \sum_{i=1}^{c} \mathbf{\omega}_{i} P(\mathbf{\omega}_{i} \mid \mathbf{m}) = \mathbf{E}\{\mathbf{\omega} \mid \mathbf{m}\}$$

Setzt man nun die Definitionsgleichungen der Vektoren ω_i ein und beachtet $P(\omega_i | \mathbf{m}) = P(\omega_i | \mathbf{m})$, so folgt:

$$\mathbf{k}(\mathbf{m}) = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} P(\omega_1 \mid \mathbf{m}) + \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} P(\omega_2 \mid \mathbf{m}) + \dots + \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix} P(\omega_c \mid \mathbf{m}) = \begin{bmatrix} P(\omega_1 \mid \mathbf{m}) \\ P(\omega_2 \mid \mathbf{m}) \\ \vdots \\ P(\omega_c \mid \mathbf{m}) \end{bmatrix} =: \mathbf{p}$$

Entscheidungsvektor mit kleinster mittlerer quadratischer Abweichung

Der optimale Entscheidungsvektor ist der Vektor der A Posteriori Wahrscheinlichkeiten der Klassen ω bei gegebenem Merkmalsvektor \mathbf{m} .

Entscheidung für Klasse ω_i mit minimalem Abstand zwischen ω_i und $\mathbf{k}(\mathbf{m})$.

$$\|\mathbf{k}(\mathbf{m}) - \mathbf{\omega}_i\|^2 = \|\mathbf{p} - \mathbf{\omega}_i\|^2 \xrightarrow{!} \text{minimal}$$

$$\|\mathbf{p} - \mathbf{\omega}_i\|^2 = \| \begin{bmatrix} P_1 \\ \vdots \\ P_i - 1 \\ \vdots \\ P_c \end{bmatrix} \|^2 = \sum_{j \neq i} P_j^2 + (P_i - 1)^2 = \sum_{j=1}^c P_j^2 - 2P_i + 1 = \|\mathbf{p}\|^2 - 2P_i + 1$$

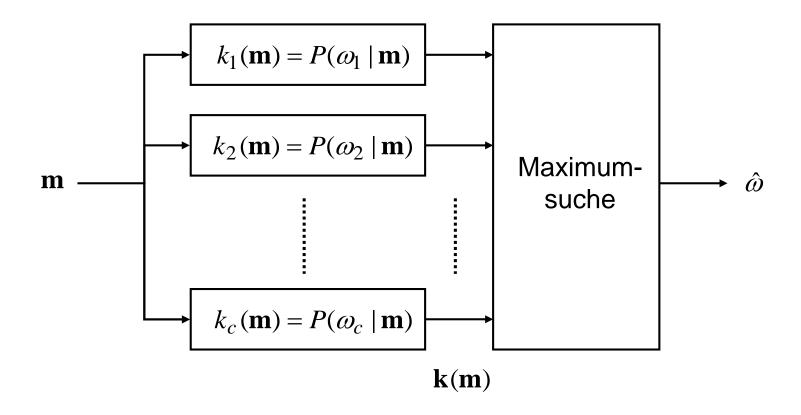
Ausdruck wird minimal für das größte $P_i := P(\omega_i \mid \mathbf{m})$.

Optimalentscheidung: Klasse mit maximaler A Posteriori Wahrscheinlichkeit

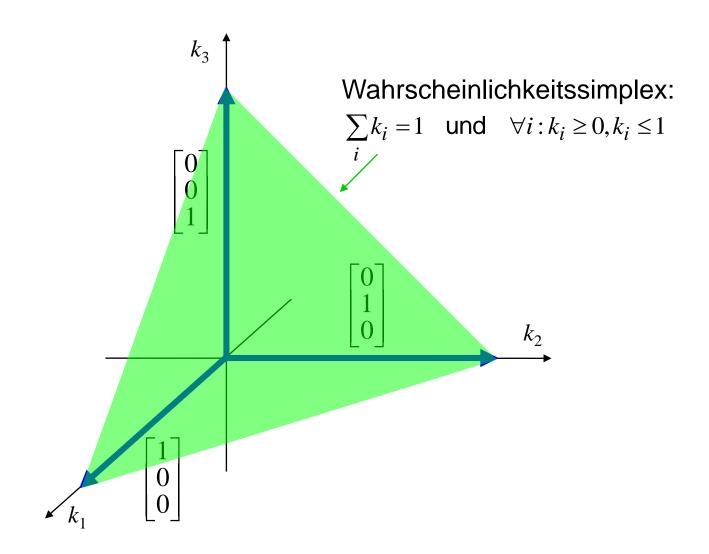
©2017 Lehrstuhl. Interaktive Echtzeitsysteme, KIT, alle Rechte einschließlich Kopier- und Weitergaberechte bei uns

3.1. Allgemeine Überlegungen zur Klassifikation

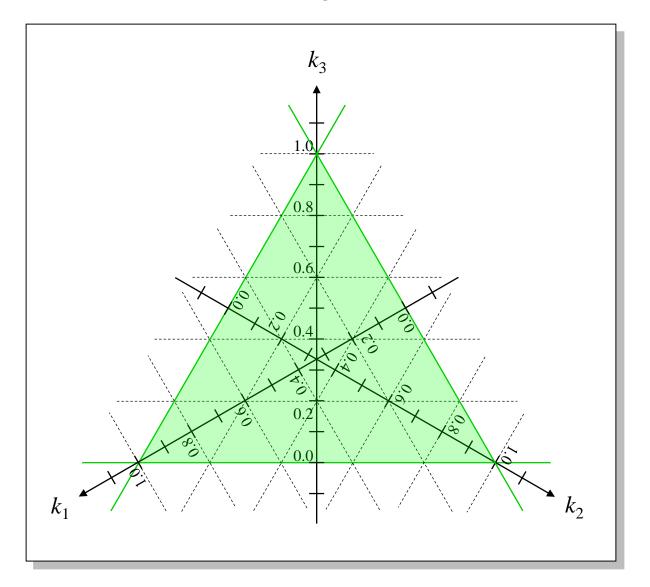
Maximum A Posteriori Entscheidung (MAP)



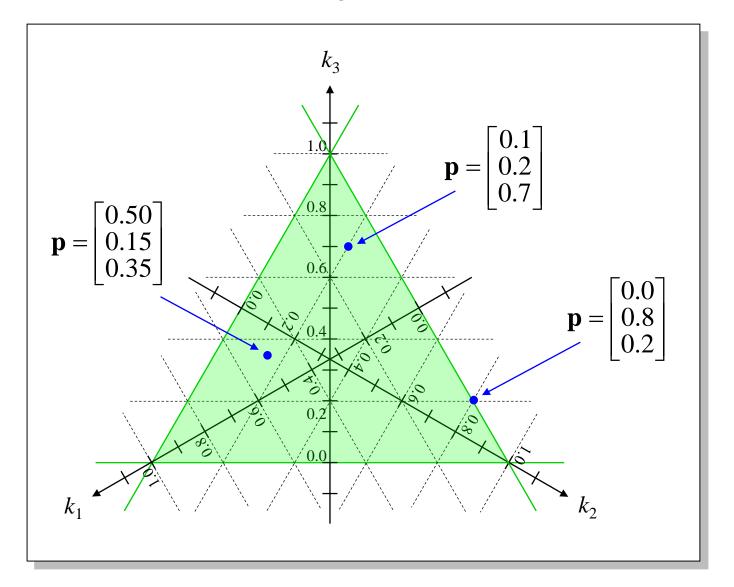
Wahrscheinlichkeiten im Entscheidungsraum



Wahrscheinlichkeiten im Entscheidungsraum



Wahrscheinlichkeiten im Entscheidungsraum



Allgemeiner Fall:

Kostenfunktion: $l(\hat{\omega}, \omega) : \Omega^0/\sim \times \Omega/\sim \to \mathbb{R}$

Beschreibt die Kosten für die Entscheidung für $\hat{\omega}$, wenn tatsächlich die Klasse ω zugrunde liegt.

Kostenmatrix:

$$\mathbf{L} \coloneqq \begin{bmatrix} l(\omega_0, \omega_1) & l(\omega_0, \omega_2) & \cdots & l(\omega_0, \omega_c) \\ \vdots & \vdots & \ddots & \vdots \\ l(\omega_l, \omega_1) & l(\omega_1, \omega_2) & \cdots & l(\omega_1, \omega_c) \\ \vdots & \vdots & \ddots & \vdots \\ l(\omega_c, \omega_1) & \cdots & \cdots & l(\omega_c, \omega_c) \end{bmatrix}$$

Ziel: Risiko
$$R = E\{l(\hat{\omega}, \omega)\} = \min_{\hat{\omega}(\mathbf{m})}^!$$

= minimale mittlere Kosten

Allgemeiner Fall:

$$R = E\{l(\hat{\omega}, \omega)\} = \int_{\mathbb{IM}} \sum_{\omega = \omega_{1}}^{\omega_{c}} l(\hat{\omega}(\mathbf{m}), \omega) P(\mathbf{m}, \omega) d\mathbf{m} =$$

$$= \int_{\mathbb{IM}} \left(\sum_{\omega = \omega_{1}}^{\omega_{c}} l(\hat{\omega}(\mathbf{m}), \omega) P(\omega | \mathbf{m}) \right) p(\mathbf{m}) d\mathbf{m} \qquad \stackrel{!}{=} \quad \underset{\hat{\omega}(\mathbf{m})}{\text{minimal}}$$

$$R(\hat{\omega} | \mathbf{m}) = \sum_{\omega = \omega_{1}}^{\omega_{c}} l(\hat{\omega}(\mathbf{m}), \omega) P(\omega | \mathbf{m}) \stackrel{!}{=} \underset{\hat{\omega}(\mathbf{m})}{\text{minimal}}$$

$$\hat{\omega}(\mathbf{m})$$

A Posteriori Risiko (bedingtes Risiko)

$$\mathbf{r} = [r_0, r_1, \dots, r_c]^{\mathrm{T}} := [R(\omega_0 \mid \mathbf{m}), R(\omega_1 \mid \mathbf{m}), \dots, R(\omega_c \mid \mathbf{m})]^{\mathrm{T}} = \mathbf{L}\mathbf{p}$$

A Posteriori Risiken:

$$r_i = R(\omega_i \mid \mathbf{m}) = \sum_{j=1}^{c} l(\omega_i, \omega_j) P(\omega_j \mid \mathbf{m})$$

Optimale Entscheidung:
$$r_i = \min\{r_0, r_1, \dots, r_c\} \iff \hat{\omega} = \omega_i$$

Diskussion:

Ingredienzen Bayes'sche Entscheidungstheorie:

- Klassenbedingte WV der Merkmale und A Priori WV der Klassen
 → A Posteriori WV der Klassen
- Kostenfunktion

Bayes'sche Entscheidungstheorie nutzt im Kontext des probabilistischen Ansatzes alles, was man über die Merkmale und die Klassen für die Gesamtheit aller Objekte der Domäne überhaupt wissen kann.

- → Leistung des Bayes'schen Optimalklassifikators stellt die obere Grenze für die Leistung anderer Klassifikatoren dar.
- → Bayes'sche Optimalklassifikation dient als Referenz für andere Klassifikatoren.

Problem: A Posteriori WV der Klassen $P(\omega|\mathbf{m})$ bzw. deren Bestimmungsstücke: die klassenspezifischen WVen der Merkmale $p(\mathbf{m}|\omega)$ und die A Priori WV $P(\omega)$ der Klassen sind i.d.R. nicht bekannt.

Praktische Vorgehensweise: Z.B. parametrische Modelle insbesondere für $p(\mathbf{m}; \theta | \omega)$ und Schätzung der Parameter θ anhand der Lernstichprobe D.

Beispiel für 2 Klassen:

Kosten: $l_{ij} = l(\hat{\omega} = \omega_i, \omega_j)$ i, j = 1, 2

Risiko: $R(\omega_1 \mid \mathbf{m}) = l_{11}P(\omega_1 \mid \mathbf{m}) + l_{12}P(\omega_2 \mid \mathbf{m})$

$$R(\omega_2 \mid \mathbf{m}) = l_{21}P(\omega_1 \mid \mathbf{m}) + l_{22}P(\omega_2 \mid \mathbf{m})$$

Entscheidung:

$$\hat{\omega} = \omega_1 \text{ wenn } R(\omega_1 \mid \mathbf{m}) < R(\omega_2 \mid \mathbf{m})$$

$$\Leftrightarrow (l_{21} - l_{11})P(\omega_1 \mid \mathbf{m}) > (l_{12} - l_{22})P(\omega_2 \mid \mathbf{m})$$

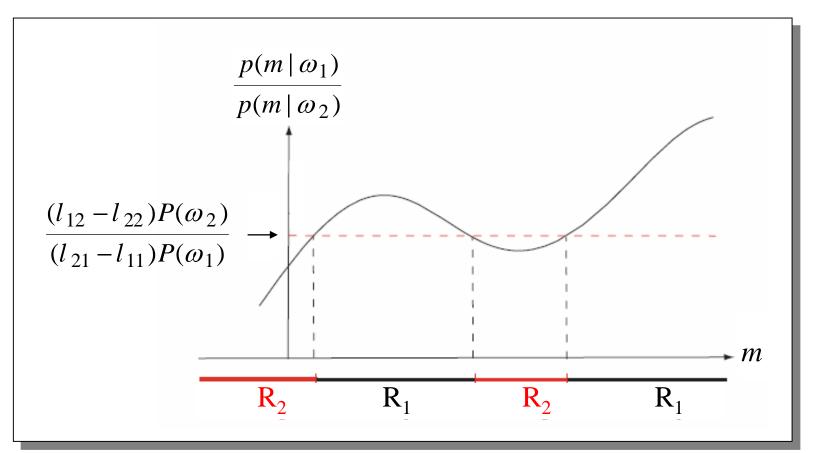
$$\Leftrightarrow (l_{21} - l_{11})p(\mathbf{m} \mid \omega_1)P(\omega_1) > (l_{12} - l_{22})p(\mathbf{m} \mid \omega_2)P(\omega_2)$$

$$\text{sonst } \hat{\omega} = \omega_2$$

$$\hat{\omega} = \omega_1 \quad \text{wenn} \quad \frac{p(\mathbf{m} \mid \omega_1)}{p(\mathbf{m} \mid \omega_2)} > \frac{(l_{12} - l_{22})P(\omega_2)}{(l_{21} - l_{11})P(\omega_1)} \quad \text{sonst} \quad \hat{\omega} = \omega_2$$

Likelihood-Verhältnis

Beispiel für 2 Klassen, 1 Merkmal:



Quelle: R. O. Duda, P. E. Hart, D. G. Stork: Pattern Classification

Minimum Fehler Klassifikation

$$\text{Kosten:} \quad l(\omega_i, \omega_j) \coloneqq 1 - \delta_i^{\ j} = \begin{cases} 0 & i = j \\ 1 & i \neq j \end{cases} \qquad i, j = 1, \dots, c$$

$$\Rightarrow R(\omega_i \mid \mathbf{m}) = \sum_{j=1}^{c} l(\omega_i, \omega_j) P(\omega_j \mid \mathbf{m}) = \sum_{j \neq i} P(\omega_j \mid \mathbf{m}) = 1 - P(\omega_i \mid \mathbf{m})$$

MAP-Entscheidung: "Maximum A Posteriori Entscheidung"

$$\hat{\omega} = \omega_i \text{ wenn } P(\omega_i \mid \mathbf{m}) > P(\omega_j \mid \mathbf{m}) \quad \forall j \neq i$$

Das Risiko ist dann gleich der Fehlerwahrscheinlichkeit.

Fehlerwahrscheinlichkeit

$$P(\text{Richtig}) = \sum_{i=1}^{c} P(\mathbf{m} \in \mathbf{R}_i, \omega_i) = \sum_{i=1}^{c} P(\mathbf{m} \in \mathbf{R}_i \mid \omega_i) P(\omega_i) = \sum_{i=1}^{c} \int_{\mathbf{R}_i} p(\mathbf{m} \mid \omega_i) P(\omega_i) d\mathbf{m}$$

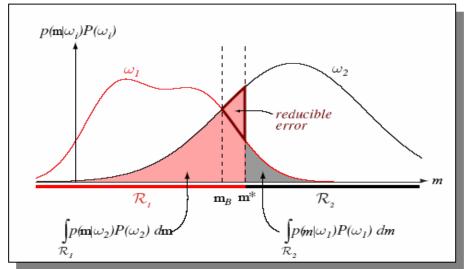
P(Fehler) = 1 - P(Richtig)

Für c = 2 Klassen:

$$P(\text{Fehler}) = P(\mathbf{m} \in \mathbf{R}_2, \omega_1) + P(\mathbf{m} \in \mathbf{R}_1, \omega_2)$$

$$= P(\mathbf{m} \in \mathbf{R}_2 \mid \omega_1) P(\omega_1) + P(\mathbf{m} \in \mathbf{R}_1 \mid \omega_2) P(\omega_2)$$

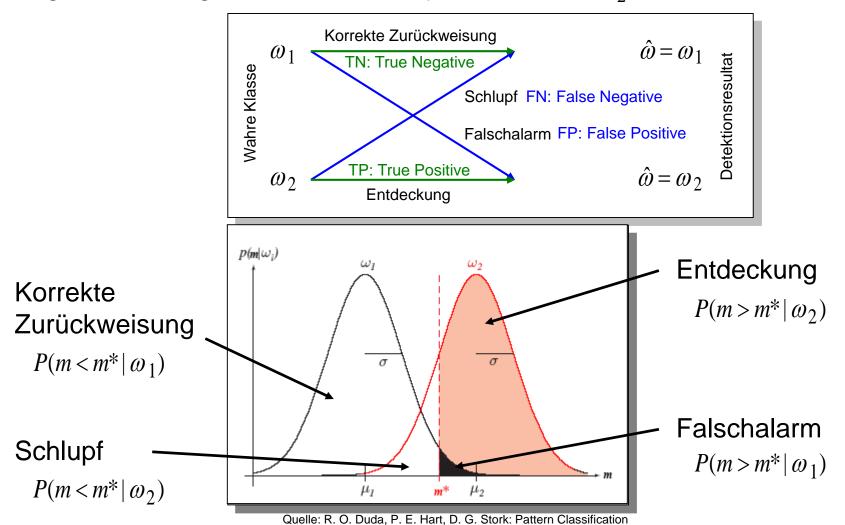
$$= \int_{\mathbf{R}_2} p(\mathbf{m} \mid \omega_1) P(\omega_1) d\mathbf{m} + \int_{\mathbf{R}_1} p(\mathbf{m} \mid \omega_2) P(\omega_2) d\mathbf{m}$$



Quelle: R. O. Duda, P. E. Hart, D. G. Stork: Pattern Classification

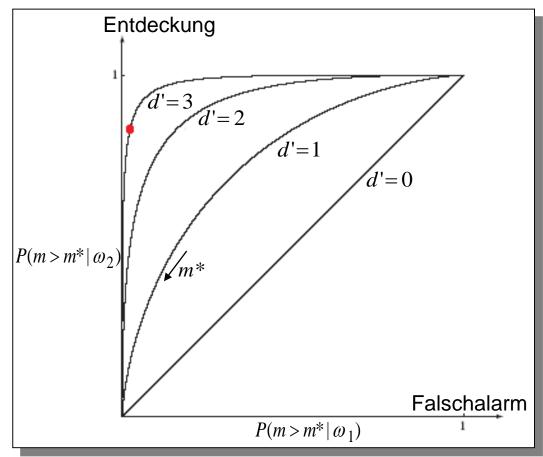
ROC (Receiver operating characteristic)

Aufgabenstellung Detektion: c=2; speziell *Klasse* ω_2 *soll* "entdeckt" werden. §



ehrstuhl. Interaktive Echtzeitsysteme, KIT, alle Rechte einschließlich Kopier- und Weitergaberechte bei uns-

ROC (Receiver operating characteristic)



Quelle: R. O. Duda, P. E. Hart, D. G. Stork: Pattern Classification

Die Kurven werden mit fallendem m^* durchlaufen.

$$d' := \frac{|\mu_2 - \mu_1|}{\sigma}$$

Referenzbeispiel: Optimale Entscheidungsgrenze nach Bayes

Zahl der Klassen: c=2

Klassenbedingte WDFen der Merkmale als Gauß'sche Mixtur:

$$p(\mathbf{m} \mid \omega_j) = \sum_{i=1}^7 \frac{1}{7} N \left(\mathbf{m}; \, \boldsymbol{\mu}_{ji}, \begin{pmatrix} \sigma_j^2 & 0 \\ 0 & \sigma_j^2 \end{pmatrix} \right) \quad j = 1, 2$$

$$\sigma_1^2 = \sigma_2^2 := 1$$

A priori Wahrscheinlichkeiten der Klassen: $P(\omega_j) := \frac{1}{2}$ j = 1, 2

Theoretischer Klassifikationsfehler = Fehlerwahrscheinlichkeit $\approx 10,85\%$

Mit diesem Modell wurden Stichproben erzeugt:

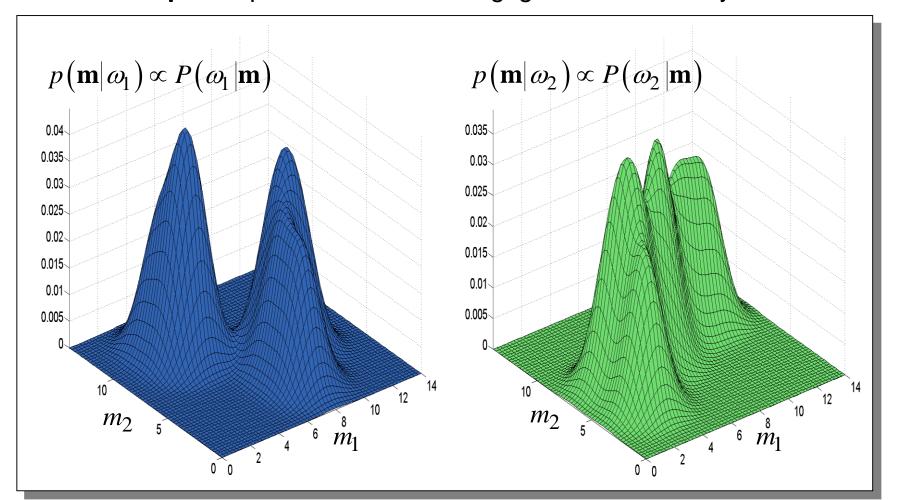
Lernstichprobe: D: N = 200, $N_1 = 100$, $N_2 = 100$,

Teststichprobe: T: N = 200, $N_1 = 100$, $N_2 = 100$,

© 2017 Lehrstuhl. Interaktive Echtzeitsysteme, KIT, alle Rechte einschließlich Kopier- und Weitergaberechte bei uns

3.2. Bayes'sche Entscheidungstheorie

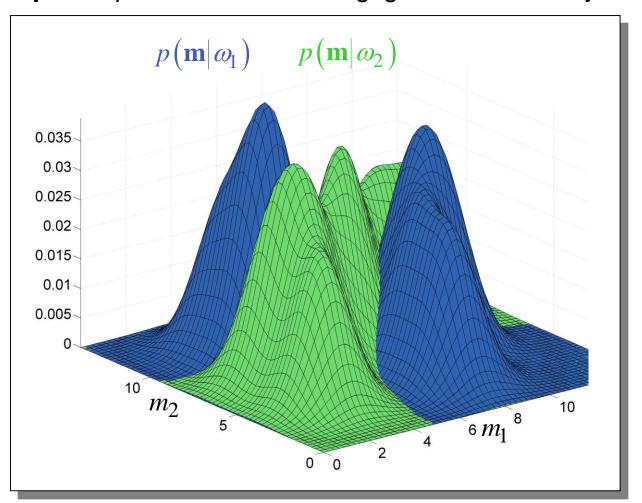
Referenzbeispiel: Optimale Entscheidungsgrenzen nach Bayes



Bedingte WDF der Merkmale im Falle der Klasse 1

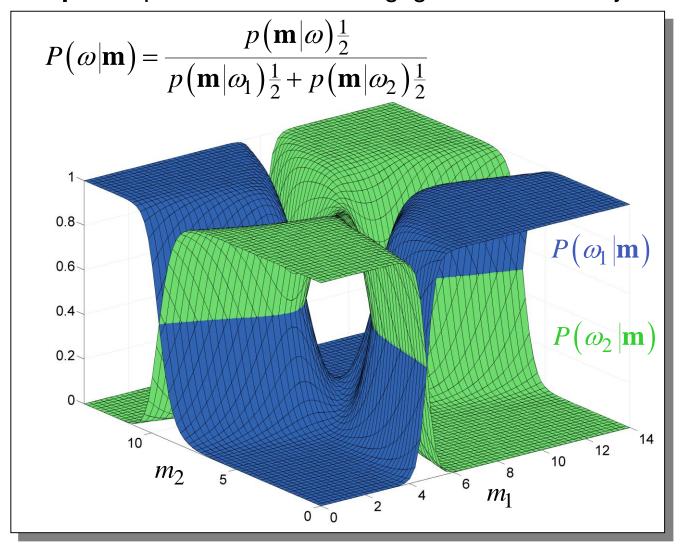
Bedingte WDF der Merkmale im Falle der Klasse 2

Referenzbeispiel: Optimale Entscheidungsgrenzen nach Bayes



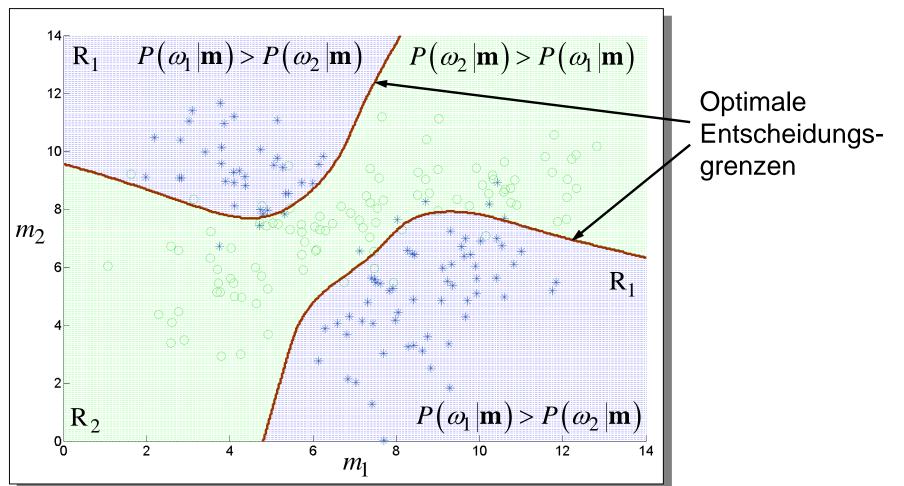
Bedingte WDFen der Merkmale für beide Klassen

Referenzbeispiel: Optimale Entscheidungsgrenzen nach Bayes



A posteriori Wahrscheinlichkeiten für beide Klassen

Referenzbeispiel: Optimale Entscheidungsgrenzen nach Bayes



Stichprobe: $N_1 = 100$, $N_2 = 100$, Empirischer Klassifikationsfehler: 9%

Minimax Entscheidungen

Ziel: Klassifikation, die *robust* ist bezüglich einer Variation der A Priori WV $P(\omega)$ der Klassen.

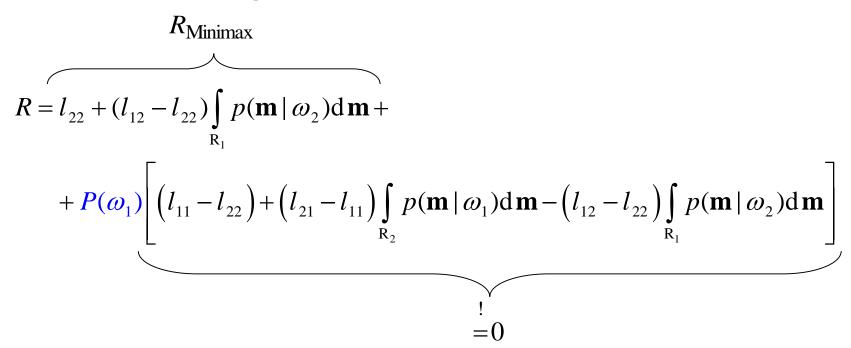
Beispiel für 2 Klassen:

Risiko:

$$R = \int R(\hat{\omega}(\mathbf{m}) \mid \mathbf{m}) p(\mathbf{m}) d\mathbf{m} = \int_{R_1} [l_{11}P(\omega_1)p(\mathbf{m} \mid \omega_1) + l_{12}P(\omega_2)p(\mathbf{m} \mid \omega_2)] d\mathbf{m} + \int_{R_2} [l_{21}P(\omega_1)p(\mathbf{m} \mid \omega_1) + l_{22}P(\omega_2)p(\mathbf{m} \mid \omega_2)] d\mathbf{m}$$

Mit
$$P(\omega_1) = 1 - P(\omega_2)$$
 und $\int_{R_i} p(\mathbf{m} \mid \omega_i) d\mathbf{m} = 1 - \int_{R_j} p(\mathbf{m} \mid \omega_i) d\mathbf{m}, \forall i \neq j$ folgt:

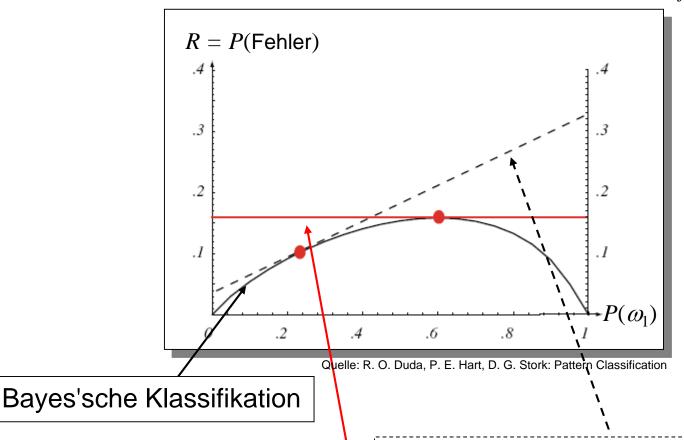
Minimax Entscheidungen, Beispiel für 2 Klassen



Beim Design des Minimax-Klassifikators werden die Entscheidungsgebiete R_i so festgelegt, dass die Abhängigkeit von der A Priori WV $P(\omega)$ der Klassen verschwindet.

$$R_{\text{Minimax}} = l_{22} + (l_{12} - l_{22}) \int_{R_1} p(\mathbf{m} \mid \omega_2) d\mathbf{m} = l_{11} + (l_{21} - l_{11}) \int_{R_2} p(\mathbf{m} \mid \omega_1) d\mathbf{m}$$

Minimax Entscheidungen, Beispiel für 2 Klassen, $l(\omega_i, \omega_j) = 1 - \delta_i^J$



Bayes'scher Klassifikator $P_{\text{Entwurf}}(\omega_1) = 0.25$

Minimax Entscheidungen

Diskussion:

Der Minimax-Klassifikator entspricht dem Bayes'schen Klassifikator für die *ungünstigste* A Priori WV der Klassen, d.h. für die A Priori WV mit dem maximalen Bayes'schen Risiko.

Der Minimax-Klassifikator minimiert das maximale Risiko.

Minimax-Entscheidungen sind vor allem in der Spieltheorie von Bedeutung, wo ein intelligenter Gegner versucht, maximal zu schaden. Bei Mustererkennungsaufgaben hingegen ist i.d.R. die Umwelt der "Gegner".

3.3. Normalverteilte Merkmale

Normalverteilung eindimensional

Eine kontinuierliche Zufallsvariable m heißt (univariat) normalverteilt mit Erwartungswert μ und Varianz σ^2 , wenn für ihre Wahrscheinlichkeitsdichtefunktion (WDF) gilt:

$$m \square N(m; \mu, \sigma^2)$$

$$p(m) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2}\left(\frac{m-\mu}{\sigma}\right)^2\right]$$

$$\mu := \mathrm{E}\{m\} = \int_{-\infty}^{\infty} mp(m) \,\mathrm{d}\,m \qquad \sigma^2 := \mathrm{E}\left\{ (m-\mu)^2 \right\} = \int_{-\infty}^{\infty} (m-\mu)^2 \, p(m) \,\mathrm{d}\,m$$

Zur Bedeutung der Normalverteilung (Gauß'sche Verteilung):

Zentraler Grenzwertsatz: Unter sehr allgemeinen Voraussetzungen für eine Folge stochastisch unabhängiger Zufallsvariablen ist die Summe dieser Folge asymptotisch normalverteilt;

Details findet man z.B. im Satz von Lindeberg-Feller, "Wahrscheinlichkeitsrechnung und mathematische Statistik", Fisz, VEB-Verlag, Berlin 1989.

Normalverteilung mehrdimensionnal

Eine kontinuierliche Zufallsvariable $\mathbf{m} \in \mathbb{R}^d$ heißt (multivariat) normalverteilt mit Erwartungswert μ und Kovarianzmatrix Σ , wenn für ihre WDF gilt:

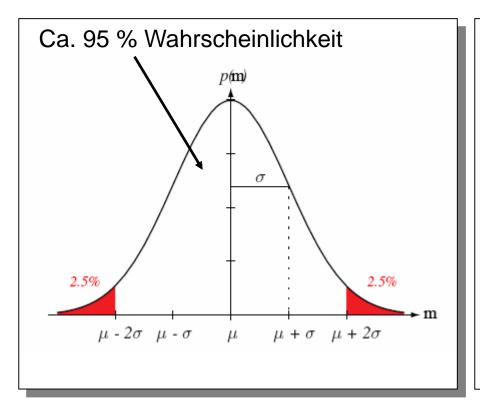
$$\mathbf{m} \sim N(\mathbf{m}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$p(\mathbf{m}) = \frac{1}{(2\pi)^{d/2} |\mathbf{\Sigma}|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{m} - \boldsymbol{\mu})^{\mathrm{T}} \mathbf{\Sigma}^{-1} (\mathbf{m} - \boldsymbol{\mu}) \right]$$

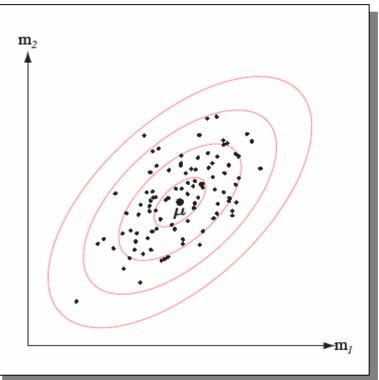
$$\boldsymbol{\mu} := \mathbf{E}\{\mathbf{m}\} = \int_{-\infty}^{\infty} ... \int_{-\infty}^{\infty} \mathbf{m} p(\mathbf{m}) \, d\mathbf{m}$$

$$\Sigma := E \left\{ (\mathbf{m} - \boldsymbol{\mu})(\mathbf{m} - \boldsymbol{\mu})^{\mathrm{T}} \right\} = \int_{-\infty}^{\infty} ... \int_{-\infty}^{\infty} (\mathbf{m} - \boldsymbol{\mu})(\mathbf{m} - \boldsymbol{\mu})^{\mathrm{T}} p(\mathbf{m}) d\mathbf{m}$$

Normalverteilung - Beispiele



Univariate Normalverteilung



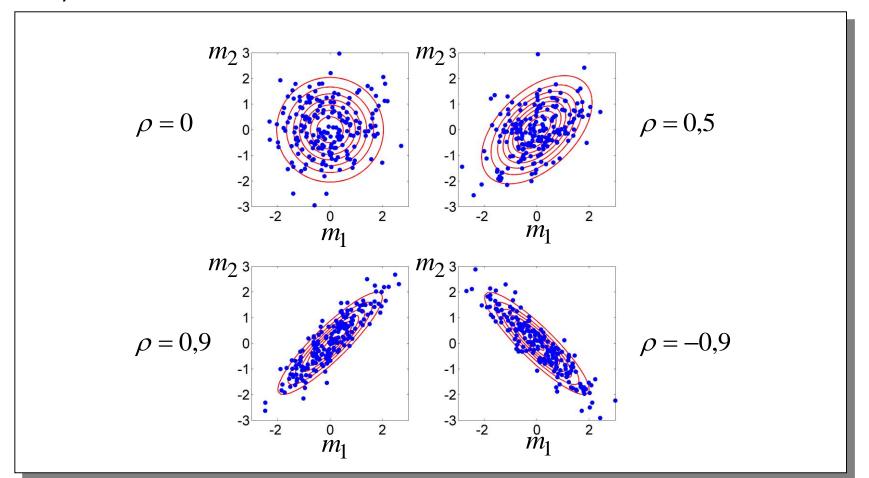
Quelle: R. O. Duda, P. E. Hart, D. G. Stork: Pattern Classification

Multivariate Normalverteilung

Ellipsen: Punkte konstanter Wahrscheinlichkeitssdichte.

Bsp.:
$$\mathbf{m} \sim N(\mathbf{m}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = N \left[\mathbf{m}; \boldsymbol{\mu}, \begin{bmatrix} \sigma_1^2 & \rho \sigma_1 \sigma_2 \\ \rho \sigma_1 \sigma_2 & \sigma_2^2 \end{bmatrix} \right] = N \left[\mathbf{m}; \boldsymbol{0}, \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix} \right]$$

 $\rho \in [-1, 1]$: Korrelationskoeffizient



MAP-Klassifikation bei normalverteilten Merkmalen:

$$\hat{\omega} = \omega_i \text{ wenn } P(\omega_i | \mathbf{m}) > P(\omega_j | \mathbf{m}) \quad \forall j \neq i \qquad P(\omega | \mathbf{m}) = \frac{p(\mathbf{m} | \omega)P(\omega)}{p(\mathbf{m})}$$

Normalverteilungsannahme für die klassenbedingte WV der Merkmale:

$$p(\mathbf{m} \mid \omega_i) = N(\mathbf{m}; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \quad i = 1, \dots, c$$

Normalverteilungsannahme macht nur bei mindestens *intervall-skalierten* Merkmalen Sinn.

MAP-Klassifikation bei normalverteilten Merkmalen:

$$\hat{\omega} = \omega_i \text{ wenn } P(\omega_i | \mathbf{m}) > P(\omega_i | \mathbf{m}) \quad \forall j \neq i$$

Wegen der strengen Monotonie der Logarithmusfunktion gilt:

$$P(\omega_i|\mathbf{m}) > P(\omega_i|\mathbf{m}) \Leftrightarrow \ln P(\omega_i|\mathbf{m}) > \ln P(\omega_i|\mathbf{m})$$

Bei normalverteilten Merkmalen rechnet es sich leichter mit den Logarithmen:

$$k_i(\mathbf{m}) := \ln p(\mathbf{m} \mid \omega_i) + \ln P(\omega_i)$$

$$p(\mathbf{m} \mid \omega_i) = N(\mathbf{m}; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$$

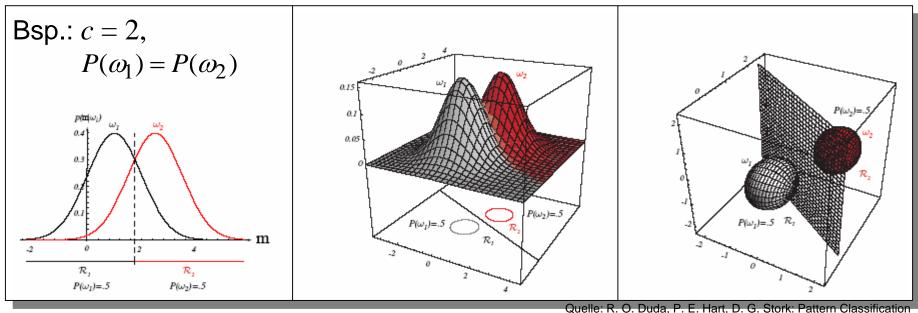
$$k_i(\mathbf{m}) = -\frac{1}{2}(\mathbf{m} - \boldsymbol{\mu}_i)^{\mathrm{T}} \boldsymbol{\Sigma}_i^{-1}(\mathbf{m} - \boldsymbol{\mu}_i) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\boldsymbol{\Sigma}_i| + \ln P(\omega_i)$$

Beispiel:
$$\Sigma_i = \sigma^2 \mathbf{I}$$
 unabhängig von i
$$k_i(\mathbf{m}) = -\frac{1}{2} (\mathbf{m} - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{m} - \boldsymbol{\mu}_i) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i)$$

$$(1/\sigma^2) \mathbf{I}$$

Entscheidungsrelevant:

$$k_i(\mathbf{m}) \coloneqq -\frac{1}{2\sigma^2} \| \mathbf{m} - \mathbf{\mu}_i \|^2 + \ln P(\omega_i)$$



Verteilungen sind *sphärisch* in d-Dimensionen, Grenzen sind d-1-dimensionale Ebenen (Hyperebenen) senkrecht zu den Verbindungsgeraden der Erwartungswerte.

017 Lehrstuhl. Interaktive Echtze

e, KIT, alle Rechte einschließlich Kopier- und Weitergabere

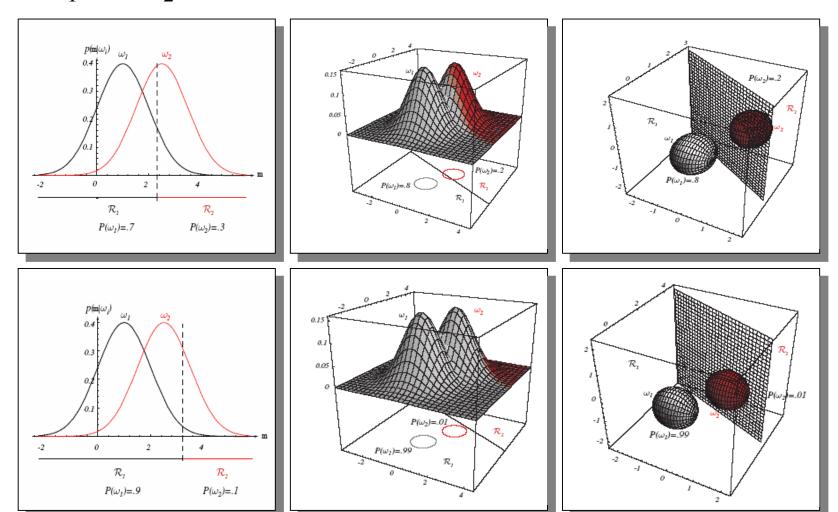
Entscheidungsgebietsgrenze zwischen R_i und R_j :

$$k_i(\mathbf{m}) - k_j(\mathbf{m}) = \frac{1}{\sigma^2} (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)^T \mathbf{m} - \frac{1}{2\sigma^2} (\|\boldsymbol{\mu}_i\|^2 - \|\boldsymbol{\mu}_j\|^2) + \ln \frac{P(\omega_i)}{P(\omega_j)} = 0$$

Hyperebene im Merkmalsraum

Beispiel: $\Sigma_i = \sigma^2 \mathbf{I}$

Für $P(\omega_1) \neq P(\omega_2)$: Grenze nicht mehr mittig zwischen den Erwartungswerten.



Quelle: R. O. Duda, P. E. Hart, D. G. Stork: Pattern Classification

Beispiel: $\Sigma_i = \Sigma$

$$k_{i}(\mathbf{m}) = -\frac{1}{2}(\mathbf{m} - \boldsymbol{\mu}_{i})^{\mathrm{T}} \boldsymbol{\Sigma}^{-1}(\mathbf{m} - \boldsymbol{\mu}_{i}) - \underbrace{\frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\boldsymbol{\Sigma}| + \ln P(\boldsymbol{\omega}_{i})}_{\text{unabhängig von } i}$$

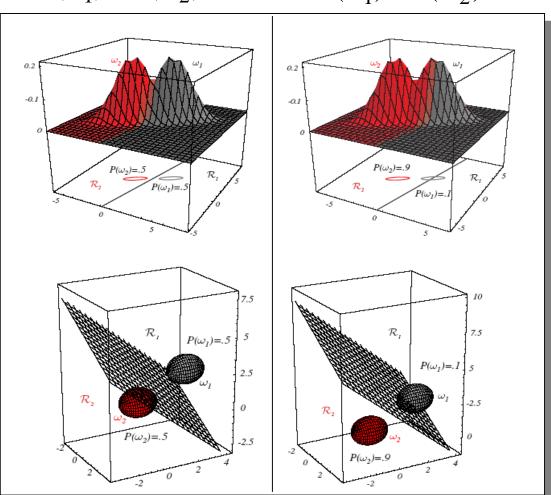
Entscheidungsrelevant:

$$k_i(\mathbf{m}) := -\frac{1}{2} (\mathbf{m} - \boldsymbol{\mu}_i)^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{m} - \boldsymbol{\mu}_i) + \ln P(\omega_i)$$

Beispiel: $\Sigma_i = \Sigma$

$$P(\omega_1) = P(\omega_2)$$

$$P(\omega_1) \neq P(\omega_2)$$



Quelle: R. O. Duda, P. E. Hart, D. G. Stork: Pattern Classification

WV sind "*elliptisch*". Grenzen sind d-1 dim. Ebenen (Hyperebenen); i.Allg. nicht senkrecht zu den Verbindungsgeraden der Erwartungswerte.

Allgemeiner Fall:

$$k_{i}(\mathbf{m}) = -\frac{1}{2}(\mathbf{m} - \boldsymbol{\mu}_{i})^{\mathrm{T}} \boldsymbol{\Sigma}_{i}^{-1}(\mathbf{m} - \boldsymbol{\mu}_{i}) - \underbrace{\frac{d}{2}\ln 2\pi - \frac{1}{2}\ln |\boldsymbol{\Sigma}_{i}| + \ln P(\boldsymbol{\omega}_{i})}_{\text{unabhängig von } i}$$

Entscheidungsrelevant:

$$k_{i}(\mathbf{m}) := \mathbf{m}^{\mathrm{T}} \left(-\frac{1}{2} \right) \mathbf{\Sigma}_{i}^{-1} \mathbf{m} + \frac{1}{2} \mathbf{\mu}_{i}^{\mathrm{T}} \mathbf{\Sigma}_{i}^{-1} \mathbf{m} + \frac{1}{2} \mathbf{m}^{\mathrm{T}} \mathbf{\Sigma}_{i}^{-1} \mathbf{\mu}_{i} - \frac{1}{2} \mathbf{\mu}_{i}^{\mathrm{T}} \mathbf{\Sigma}_{i}^{-1} \mathbf{\mu}_{i} - \frac{1}{2} \ln |\mathbf{\Sigma}_{i}| + \ln P(\omega_{i})$$

$$\mathbf{W}_{i}$$

$$\mathbf{W}_{i}$$

$$\mathbf{W}_{i}$$

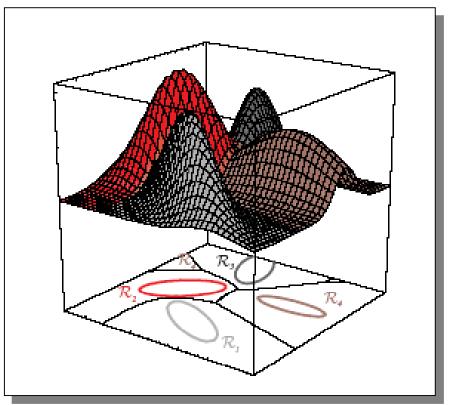
$$k_i(\mathbf{m}) = \mathbf{m}^{\mathrm{T}} \mathbf{W}_i \mathbf{m} + \mathbf{w}_i^{\mathrm{T}} \mathbf{m} + w_{i0}$$

Entscheidungsgebietsgrenze zwischen R_i und R_j :

$$k_i(\mathbf{m}) - k_i(\mathbf{m}) = \mathbf{m}^{\mathrm{T}}(\mathbf{W}_i - \mathbf{W}_i)\mathbf{m} + (\mathbf{w}_i - \mathbf{w}_i)^{\mathrm{T}}\mathbf{m} + w_{i0} - w_{i0} = 0$$

Die Entscheidungsgebietsgrenzen im Merkmalsraum sind Hyperquadriken.

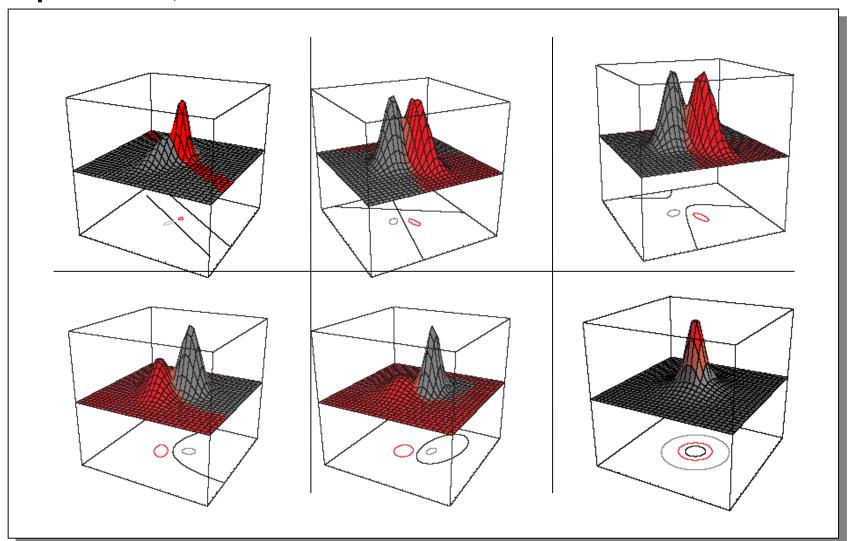
Bsp: c = 4, d = 2



Quelle: R. O. Duda, P. E. Hart, D. G. Stork: Pattern Classification

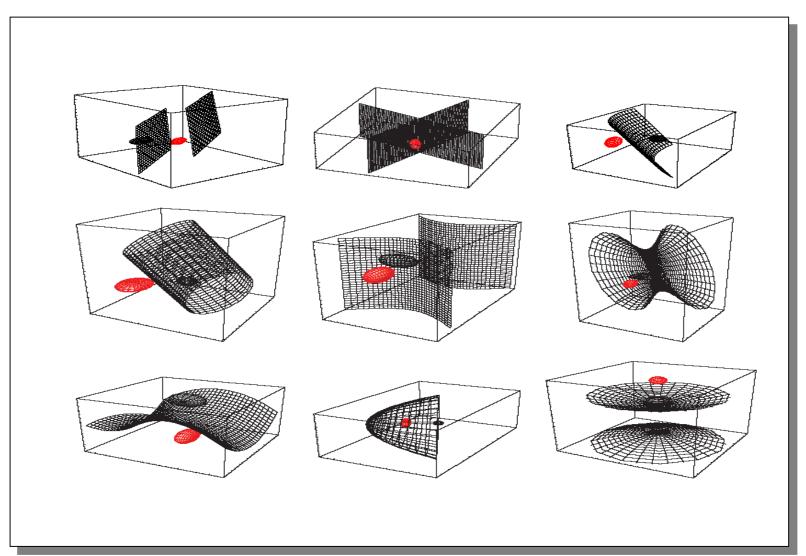
$$\partial \mathbf{R}_{i} = \bigcup_{j \neq i} \left\{ \mathbf{m} \middle| k_{i}(\mathbf{m}) = k_{j}(\mathbf{m}) \land k_{j}(\mathbf{m}) \ge k_{l}(\mathbf{m}) \forall l \neq i, j \right\}$$

Beispiele: c = 2, d = 2



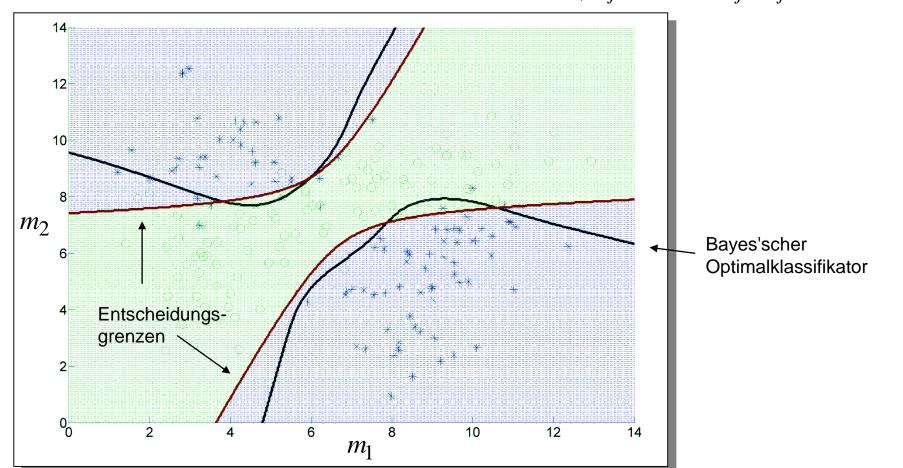
Quelle: R. O. Duda, P. E. Hart, D. G. Stork: Pattern Classification

Beispiele: c = 2, d = 3



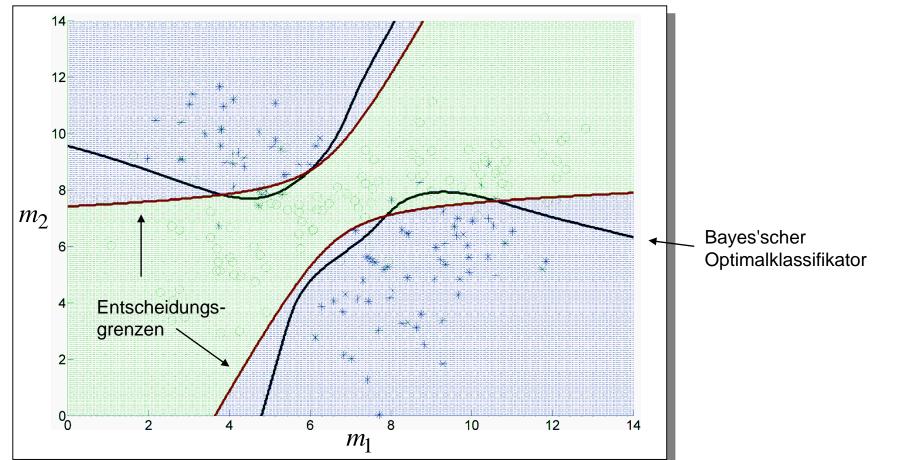
Quelle: R. O. Duda, P. E. Hart, D. G. Stork: Pattern Classification

Beispiel: Klassifikator für das Referenzbeispiel auf der Basis einer Gaußdichte für jede Klasse: $p(\mathbf{m} | \omega_j) = N(\mathbf{m}; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$ j = 1, 2



Parameter aus der Lernstichprobe geschätzt. → Trainingsfehler = 9,5%

 $p(\mathbf{m} | \omega_j) = N(\mathbf{m}; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) \quad j = 1, 2$



Parameter aus Lernstichprobe geschätzt.

Teststichprobe \rightarrow Testfehler = 11% Asymptotischer Testfehler $\approx 12,37\%$

3.4. Merkmale mit beliebiger Verteilung

Beliebige Verteilungen können mittels Gauß'schen Mixturen approximiert werden. Definition einer Gauß'schen Mixtur:

$$\mathbf{m} \square p(\mathbf{m}) = \sum_{k=1}^{K} \pi_k N(\mathbf{m}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \qquad \pi_k \ge 0, \ \sum_{k=1}^{K} \pi_k = 1$$

Gauß'sche Mixturen sind universelle Approximatoren. Mit entsprechend vielen Komponenten lassen sich beliebige WDFen beliebig genau nähern.

Jede klassenbedingte WDF der Merkmale $p(\mathbf{m}|\omega_j)$, j=1,...,c kann mittels einer klassenindividuellen Gauß'schen Mixtur approximiert werden. Komponentenanzahlen K_j für jede Klasse individuell definierbar.

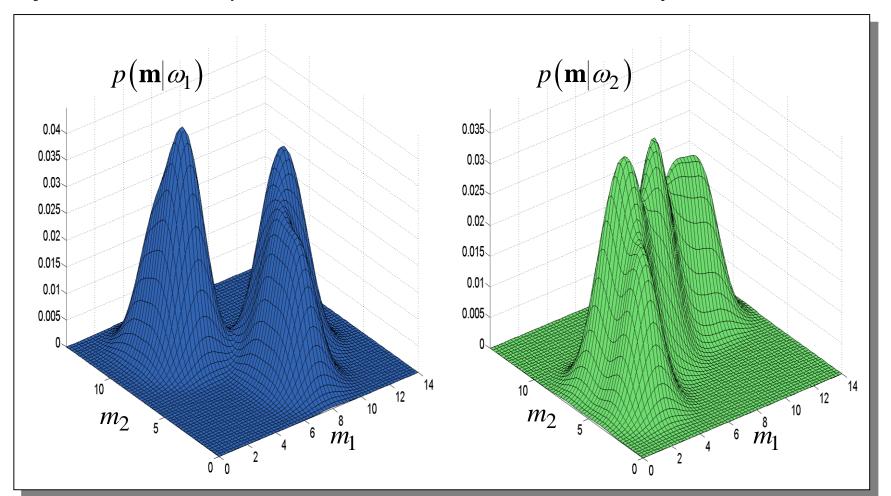
$$\mathbf{m} | \omega_j \square p(\mathbf{m} | \omega_j) = \sum_{k=1}^{K_j} \pi_{jk} N(\mathbf{m}; \boldsymbol{\mu}_{jk}, \boldsymbol{\Sigma}_{jk}) \qquad \pi_{jk} \ge 0, \sum_{k=1}^{K_j} \pi_{jk} = 1$$

Betrachtung für $K_j = K$: Pro Klasse sind auf Basis der Lerndaten $1/2K(d^2+3d+2)-1$ Parameter zu schätzen.

$$\underbrace{\pi_{j1},...,\pi_{jK}}_{K-1},\underbrace{\mu_{j1},...,\mu_{jK}}_{Kd},\underbrace{\Sigma_{j1},...,\Sigma_{jK}}_{\frac{1}{2}Kd(d+1)} \qquad j=1,...,c$$

3.4. Merkmale mit beliebiger Verteilung

Bsp.: Referenzbeispiel benutzt Gauß'sche Mixturen mit jeweils K = 7.



$$p(\mathbf{m} | \omega_j) = \sum_{k=1}^{7} \frac{1}{7} N(\mathbf{m}; \boldsymbol{\mu}_{jk}, \mathbf{I}) \quad j = 1, 2$$

4.3. Bayes'sche Klassifikation – ergänzende Bemerkungen

Fehler bei der Bayes'schen Klassifikation

- Bayes'scher Fehler: ergibt sich durch die Überlappung der der klassenbedingten Verteilungsdichten der Merkmale.
- Modellfehler: ergibt sich durch ein nicht passendes Modell.
 Auswahl eines passenden Modells mit Hilfe eines Anpassungstests für WVen, z.B. Chi-Quadrat-Test; Details in statistischer Grundlagenliteratur.
- Schätzfehler: ergibt sich aufgrund endlicher Datensätze.