

Appendix zu:

D: Grundlegende Lernverfahren

- **Motivation: Fallbeispiel Sägemühle**
- **Kalkül und Regel von Bayes: Zweiklassenproblem**
- **Generalisierung: Minimale Fehler und Risiken der Klassifikation nach Bayes**
- **Bayes-Kalkül und Diskriminantenfunktionen**
- **Fall der Normalverteilung (NV) als klassenbedingte Wahrscheinlichkeitsdichtefunktion**
- **NV und Diskriminantenfunktionen (linear und quadratisch)**

Bayessche Entscheidungstheorie

(Quelle: R. Duda, P. Hart (1972): Pattern Classification and Scene Analysis. Wiley & Sons)

- Historie
 - Thomas Bayes (1702-1763), brit. Pastor und Mathematiker
 - posthume Veröffentlichung des Kalküls
- Problemstellung: Sägemühle
 - 2-Klassen-Problem: Esche und Birke
 - Problem (Fig. 1.2): Merkmal Helligkeit unzureichend für „sichere“ Entscheidung

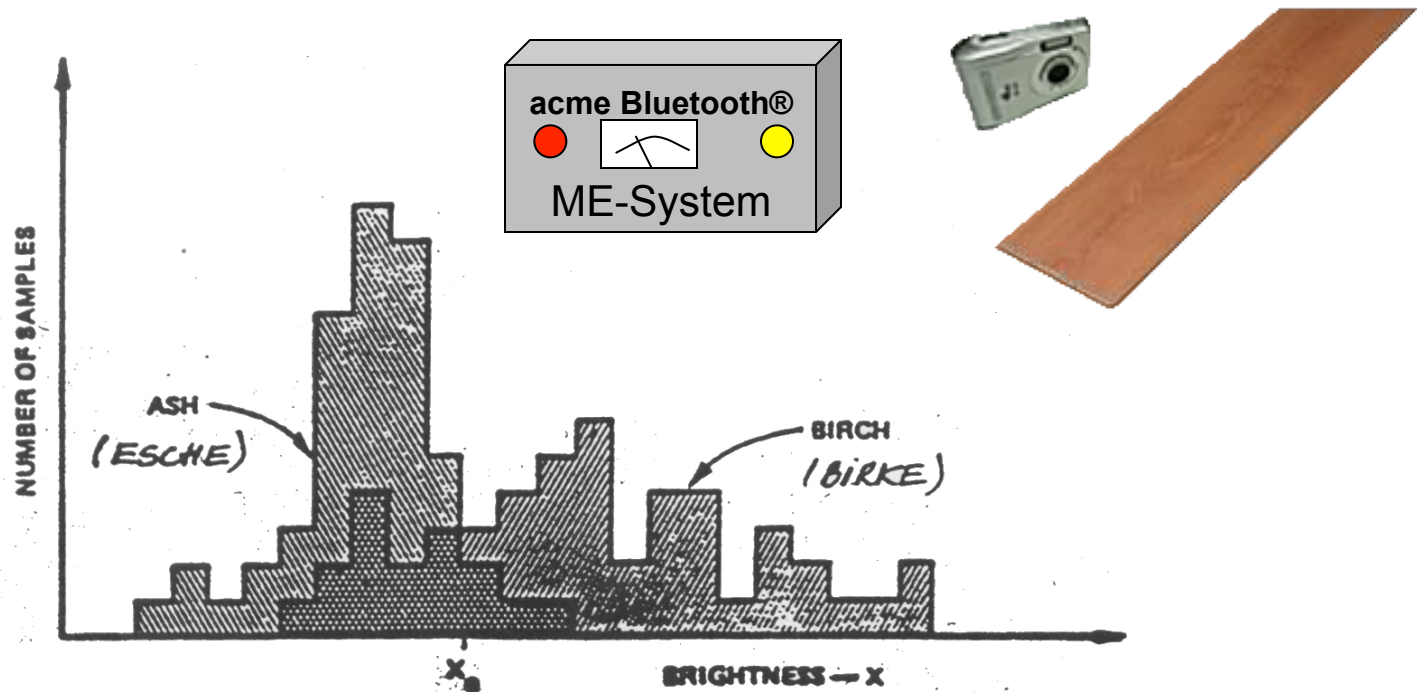


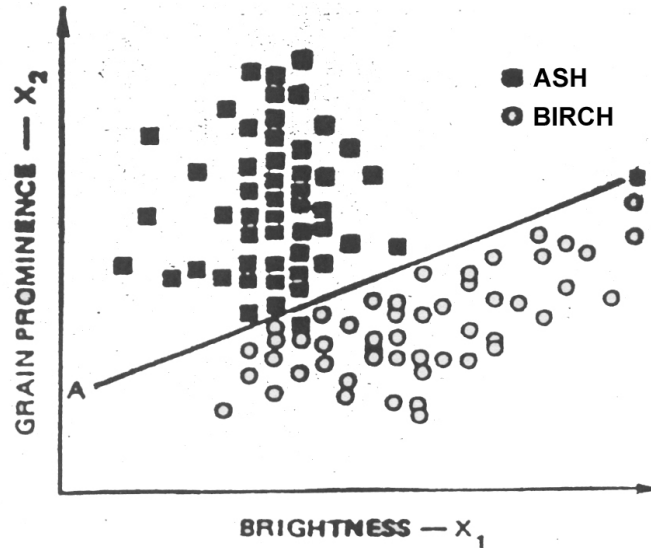
FIGURE 1.2. Histograms for the brightness feature.

(diskrete Häufigkeitsverteilung  Histogramm der „beobachteten“ Helligkeit) 2

Bayessche Entscheidungstheorie

- Lösung: Erhöhung der Dimension des Merkmalraumes:

Helligkeit und Maserung (x_1 und x_2), damit Merkmalvektor $\vec{x} = (x_1, x_2)^T$



(Fig. 1.2. und 1.3. aus Duda & Hart, 1972)

FIGURE 1.3. Scatter diagram for the feature vectors

(2-d Merkmalraum: Helligkeit vs. Maserung)

jedoch: weiterhin lokaler Fehler beobachtbar!

Lösung: Erhöhung der Dimension

Wahl anderer (d.h. besserer) Merkmale

Verwendung nichtlinearer Trennfunktion

aber: „Kosten“ der Dimensionserhöhung etc., d.h. mathematischer, algorithmischer und technischer Aufwand

Bayessche Entscheidungstheorie

- zentrale wiss.-methodische Fragen
 - Berücksichtigung der statistischen Eigenschaften der Daten (hier: Elemente der Stichprobe als Vektoren im Merkmalraum)?
 - Klassifikation mit Minimierung der Fehlerwahrscheinlichkeit?
 - Bezug zu Diskriminantenfunktionen?
 - Perzeptron als statistischer Klassifikator?
 - theoretische Grenzen?
- mathematisches Modell von Bayes für Entscheiden/Klassifizieren bei statistischer Unsicherheit mit minimalem Fehler
- Bayes 'sches Kalkül
 - Annahmen (allgemein)
 - 1) wahrscheinlichkeitstheoretische Formulierung des Problems
 - 2) probabilistische Parameterwerte vollständig bekannt (!)

Bayessche Entscheidungstheorie

„*Apriori versus Aposteriori*“

Es ist, wie man weiß, IMMANUEL KANT gewesen, der in seinen kritischen Schriften die Grenzen der Vernunft und der Urteilskraft untersuchte²⁶; und er schälte aus ihnen jene Voraussetzungen heraus, die nicht aus der Erfahrung heraus stammen können, weil sie selbst für den elementarsten Erfahrungs-Gewinn die Voraussetzung sind. Dies sind die *Apriori* der Vernunft und der Urteilskraft. Ermutigend in KANTS Ergebnis ist nur die Präzision. Wir erfahren nun sehr genau, worin das Dilemma der Vernunft besteht. Das Problem der Vernunft aber wird anerkanntermaßen nicht gelöst, es wird präzisiert. Denn das *Apriori*, was wörtlich: »vom früheren her« bedeutet, läßt sich nicht hinterfragen. Die Kette der Voraussetzungen von den Voraussetzungen kann kein Ende finden. Und damit erweist sich gerade das, worauf sich all unsere Vernunft gründen muß, durch eben diese Vernunft als durchaus unbegründbar.

Was begründet unser Vertrauen auf Wahrscheinlichkeit, die uns eine, wenn auch nur ungefähre, Voraussicht vorspiegelt, über eine Voraussicht, die wir nicht besitzen können? Wie viele weiße Schwäne müssen wir sehen, um in unserem Schluß: alle Schwäne werden weiß sein, dennoch zu irren? Wer macht uns vertrauen, daß Mehrheit Wahrheit bedeuten könnte? »Der Verrückte, der sich für ein Rührei hält«, sagt der verzweifelte BERTRAND RUSSELL, wäre dann »nur deshalb abzulehnen, weil er sich in der Minderheit befindet.«²⁷ Hierin liegt schon das Problem der *a-priori*-Wahrscheinlichkeit, sowie das HUME-KANT-POPPERSche Induktionsproblem, auf dem alles Schließen vom Speziellen auf das Allgemeine ruht. Ein Schluß also, mit dessen Berechtigung alle induktive Wissenschaft, das ist die ganze Naturwissenschaft, stehen oder fallen muß. Keine Gewißheit, kein zureichender Grund ist, wie RUDOLF CARNAP und WOLFGANG STEGMÜLLER zeigen²⁸, selbst in der modernen induktiven Logik für diesen Schluß zu finden. Und KARL POPPER erklärt ihn als einen Widerspruch in sich selbst.²⁹“

(aus: R. Riedl (1981): *Biologie der Erkenntnis*, Verlag P. Parey)

Bayessche Entscheidungstheorie

- Bayessches Kalkül
 - Annahmen (allgemein)
 - 1) wahrscheinlichkeitstheoretische Formulierung des Problems
 - 2) probabilistische Parameterwerte vollständig bekannt (!)
 - 2-Klassen-Problem (z.B. Sägemühle)
 - Klassenvektor \mathbb{Y} Zustandsvektor $\vec{\omega} = (\omega_1, \omega_2)^T$
 - a priori-Wahrscheinlichkeit eines Zustandes $P(\omega_i)$, mit
$$P(\omega_1), P(\omega_2) \geq 0$$
$$P(\omega_1) + P(\omega_2) = 1$$
d.h. Wahrscheinlichkeit, dass sich das System im Zustand ω_i befindet
 - Merkmalsvektor \vec{x} als kontinuierliche Zufallsvariable
 - klassenbedingte Wahrscheinlichkeitsdichtefunktion $p(\vec{x} | \omega_j)$, d.h. die Wahrscheinlichkeit, \vec{x} zu beobachten, wenn ω_j vorliegt (oder: wenn das System im Zustand ω_j ist)
 - spezielle Annahmen
 - $P(\omega_i)$ und $p(\vec{x} | \omega_j)$ seien bekannt
 - \vec{x} wurde gemessen bzw. beobachtet

Bayessche Entscheidungstheorie

- Bayessche Regel

$$P(\omega_j | \vec{x}) = \frac{p(\vec{x} | \omega_j) \cdot P(\omega_j)}{p(\vec{x})} \quad \text{mit } p(\vec{x}) = \sum_{j=1}^2 p(\vec{x} | \omega_j) \cdot P(\omega_j) \quad ^1$$

a posteriori-Wahrscheinlichkeit; d.h. Wahrscheinlichkeit, dass Klasse ω_j vorliegt (bzw. Wahrscheinlichkeit, dass sich das System im Zustand ω_j befindet), nachdem (!) \vec{x} beobachtet/gemessen wurde.

$$\Rightarrow \begin{cases} P(\omega_1 | \vec{x}) > P(\omega_2 | \vec{x}), & \text{dann } \vec{x} \sim \omega_1 \\ P(\omega_1 | \vec{x}) < P(\omega_2 | \vec{x}), & \text{dann } \vec{x} \sim \omega_2 \end{cases} \quad ^2$$

mit minimaler mittlerer Fehlerwahrscheinlichkeit, d.h. über alle möglichen Entscheidungen gemittelte Wahrscheinlichkeit für eine falsche Entscheidung wird minimal!

- Fehler

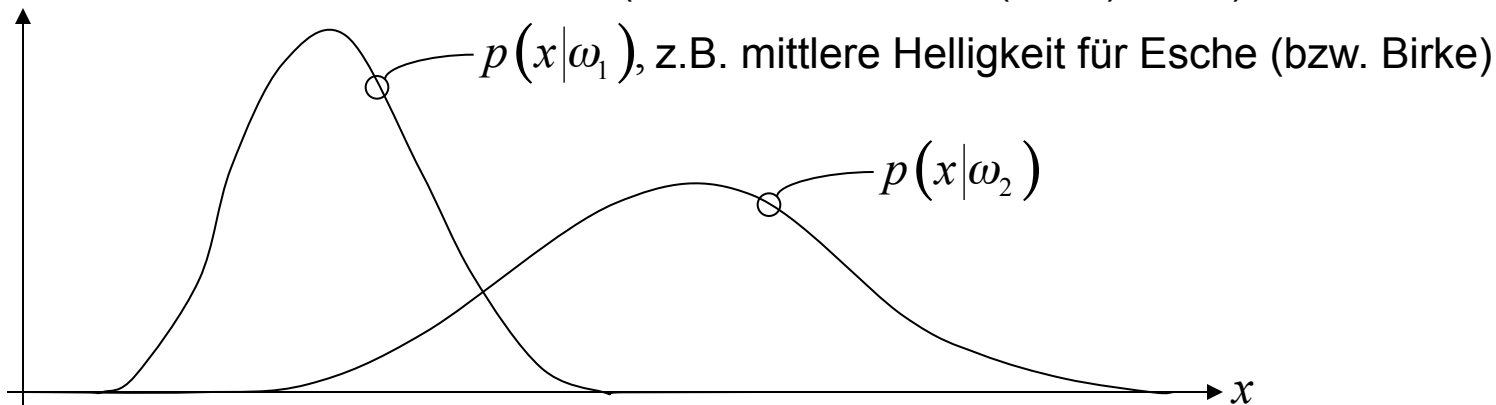
$$P(\text{Fehler} | \vec{x}) = \begin{cases} P(\omega_1 | \vec{x}), & \text{für } \vec{x} \sim \omega_2 \\ P(\omega_2 | \vec{x}), & \text{für } \vec{x} \sim \omega_1 \end{cases}$$

¹: nach Duda & Hart (1972) "... just a scale factor", damit $P(\omega_1 | \vec{x}) + P(\omega_2 | \vec{x}) = 1$

²: lies $\vec{x} \sim \omega_1$: \vec{x} wird der Klasse ω_1 zugeordnet

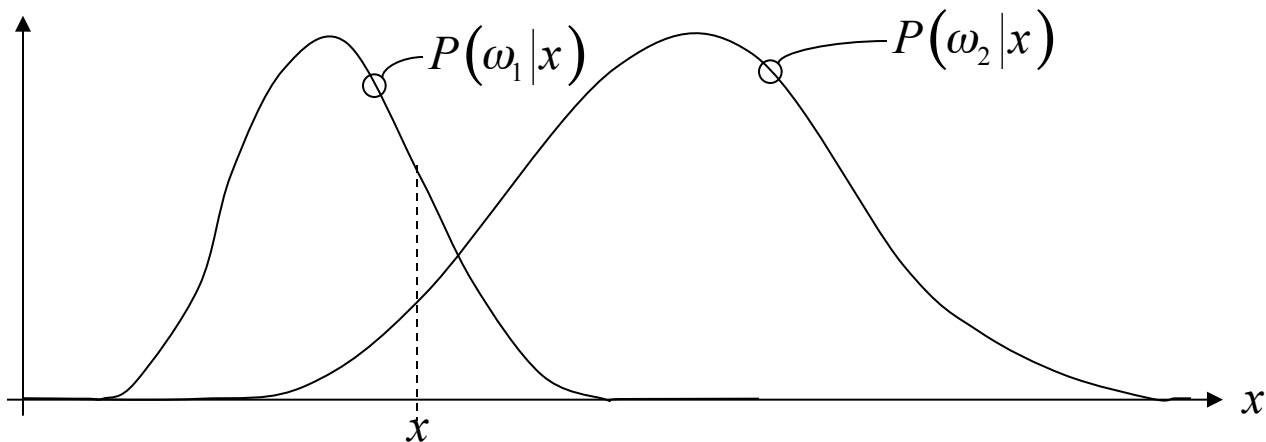
Bayessche Entscheidungstheorie

- Wahrscheinlichkeits-Funktionen (nach: Duda & Hart (1972), S.12)



(hypothetische) klassenbedingte Wahrscheinlichkeits-Dichtefunktion

(mit: x entspricht Helligkeit im Sägemühlen-Bsp.), d.h. „Streuung“ des/der beobachteten Merkmale(s) in der Stichprobe in Abhängigkeit der bekannten Klasse



(hypothetische) *a posteriori*-Wahrscheinlichkeiten für $P(\omega_1) = \frac{2}{3}$, $P(\omega_2) = \frac{1}{3}$, d.h. $x \sim \omega_1$,
da $P(\omega_1|x) > P(\omega_2|x)$ mit Fehler $P(\text{Fehler}|x) = P(\omega_2|x)$

Bayessche Entscheidungstheorie

- Sei $P(\text{Fehler})$ die mittlere Fehlerwahrscheinlichkeit mit

$$P(\text{Fehler}) = \int_{-\infty}^{\infty} P(\text{Fehler}|x) \cdot p(x) dx \quad (\text{für ein Merkmal!})$$

somit: für $x \sim \omega_2$

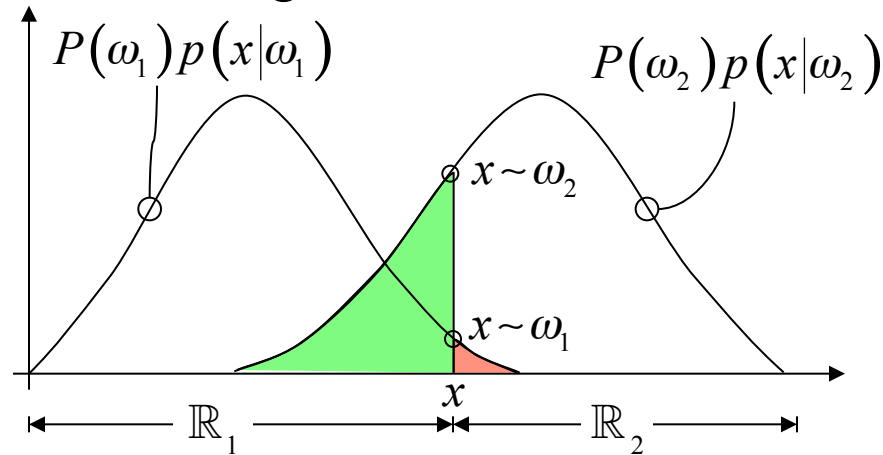
$$P(\text{Fehler}) = \int_{-\infty}^{\infty} \underbrace{P(\omega_1|x)}_{\frac{P(\omega_1)p(x|\omega_1)}{p(x)} \hat{=} \text{a posteriori-Wahrscheinlichkeit}} \cdot p(x) dx$$

$$P(\text{Fehler}) = \int_{-\infty}^{\infty} P(\omega_1) \cdot p(x|\omega_1) dx$$

und für $x \sim \omega_1$

$$P(\text{Fehler}) = \int_{-\infty}^{\infty} P(\omega_2) \cdot p(x|\omega_2) dx$$

Bayessche Entscheidungstheorie



d.h. Fehler bei

$$\begin{cases} x \sim \omega_1 : \int_{\mathbb{R}_1} P(\omega_2) p(x|\omega_2) dx \\ x \sim \omega_2 : \int_{\mathbb{R}_2} P(\omega_1) p(x|\omega_1) dx \end{cases}$$

da jedoch gemäß Bayesscher Regel gelten muss $x \sim \omega_2$

$$\Rightarrow \underbrace{\int_{\mathbb{R}_1} P(\omega_2) p(x|\omega_2) dx}_{\text{in Abb. grün dargestellt}} > \underbrace{\int_{\mathbb{R}_2} P(\omega_1) p(x|\omega_1) dx}_{\text{in Abb. rot dargestellt}}$$

d.h. Bayessche Regel minimiert mittlere Fehlerwahrscheinlichkeit

Hinweis: \mathbb{W}_1 und \mathbb{W}_2 in der obigen Abbildung sind willkürlich festgelegte Entscheidungs-Regionen zur Illustration!

Bayessche Entscheidungstheorie

Grenzfälle:

- i) $p(x|\omega_1) = p(x|\omega_2)$: Entscheidung hängt nur von *a priori*-Wahrscheinlichkeiten ab!
- ii) $P(\omega_1) = P(\omega_2)$: Entscheidung hängt nur von klassenbedingter Wahrscheinlichkeits-Dichtefunktion ab!
- iii) $P(\omega_1|x) = P(\omega_2|x)$: „zufällige“ Entscheidung mit falls gleiche Fehler-Wahrscheinlichkeit!

Bayessche Entscheidungstheorie

- Generalisierung (für den kontinuierlichen Fall)

- mehr als ein Merkmal: Merkmalsvektor \vec{x}
- mehr als 2 Klassen/Zustände: n -Klassen-Problem
- nicht nur Klassifikation, sondern auch andere Formen von „actions“, z.B.

Zurückweisungen

- nicht nur Betrachtung der Fehler-Wahrscheinlichkeit, sondern Berücksichtigung von Kosten-(Verlust-)Funktionen ('loss functions')

- $\Omega = \{\omega_1, \dots, \omega_n\}$: endliche Menge von Klassen

$A = \{\alpha_1, \dots, \alpha_n\}$: endliche Menge von „actions“ (zu verstehen als z.B. Klassenzuweisung, Systemabbruch etc.)

$\lambda(\alpha_i|\omega_j)$: Kosten (Verlust) bei der Entscheidung für α_i , falls „das System im Zustand ω_j “

\vec{x} : vektorielle Zufallsvariable mit d Komponenten, $(x_1, x_2, \dots, x_d)^T$

$p(\vec{x}|\omega_j)$: klassenbedingte Wahrscheinlichkeits-Dichtefunktion

$P(\omega_j)$: a priori-Wahrscheinlichkeit

a posteriori-Wahrscheinlichkeit (bei s Klassen):

$$P(\omega_j|\vec{x}) = \frac{p(\vec{x}|\omega_j) \cdot P(\omega_j)}{p(\vec{x})} \quad \text{mit} \quad p(\vec{x}) = \sum_{j=1}^s p(\vec{x}|\omega_j) \cdot P(\omega_j)$$

Bayessche Entscheidungstheorie

Annahme: \vec{x} wird beobachtet;

Entscheidung für 'action' α_i wird in Erwägung gezogen, ω_j ist jedoch wahre(r) Klasse (bzw. Zustand),

dann ist $\lambda(\alpha_i|\omega_j)$ 'loss' \mathbb{W} Kosten/Verlust

- bedingtes (mittleres) Risiko über alle s Klassen

$$R(\alpha_i|\vec{x}) = \sum_{j=1}^s \lambda(\alpha_i|\omega_j) \cdot P(\omega_j|\vec{x})$$

nach der Beobachtung/Messung von \vec{x} (wegen $P(\omega_j|\vec{x})$)

- Gesamtrisiko

für jedes \vec{x} mit Entscheidungs-Regel $\alpha(\vec{x}): \vec{x} \sim \{\alpha_1, \dots, \alpha_a\}$:

$$\mathbf{R} = \int_{-\infty}^{\infty} R(\alpha(\vec{x})|\vec{x}) p(\vec{x}) d\vec{x} \rightarrow \min,$$

d.h. Wahl von $\alpha(\vec{x})$ so, dass $R(\alpha(\vec{x})|\vec{x})$ minimal wird für jedes \vec{x} .

- optimale Bayessche Entscheidungsregel (Bayes-Klassifikator):

„Um das Gesamtrisiko zu minimieren, berechne das bedingte Risiko $R(\alpha_i(\vec{x})|\vec{x})$,

$i = 1, \dots, a$, und wähle α_i so, daß $R(\alpha_i(\vec{x})|\vec{x})$ ein Minimum ist.“

! minimales Gesamtrisiko \mathbb{W} Bayessches Risiko !

Bayessche Entscheidungstheorie

Beispiel: 2-Klassen-Problem

α_1 : Entscheidung für Klasse ω_1

α_2 : Entscheidung für Klasse ω_2

$\lambda_{ij} = \lambda(\alpha_i|\omega_j)$: Kosten für Entscheidung α_i (d.h. für ω_i), obwohl ω_j vorliegt

$$\mathbf{R}(\alpha(\vec{x})|\vec{x}) = \mathbf{L}\mathbf{P}(\vec{\omega}|\vec{x}) \Rightarrow \begin{aligned} R(\alpha_1|\vec{x}) &= \lambda_{11}P(\omega_1|\vec{x}) + \lambda_{12}P(\omega_2|\vec{x}) \\ R(\alpha_2|\vec{x}) &= \lambda_{21}P(\omega_1|\vec{x}) + \lambda_{22}P(\omega_2|\vec{x}) \end{aligned}$$

mit Verlust-/Kosten („loss“-) -Matrix $\mathbf{L} = \begin{pmatrix} \lambda_{11} & \lambda_{12} \\ \lambda_{21} & \lambda_{22} \end{pmatrix}^1$

$\vec{x} \sim \omega_1$ iff $R(\alpha_1|\vec{x}) < R(\alpha_2|\vec{x})$

bzw. iff $(\lambda_{21} - \lambda_{11}) \underbrace{P(\omega_1|\vec{x})}_{>} > (\lambda_{12} - \lambda_{22}) \underbrace{P(\omega_2|\vec{x})}_{>}$

bzw. iff $(\lambda_{21} - \lambda_{11})p(\vec{x}|\omega_1)P(\omega_1) > (\lambda_{12} - \lambda_{22})p(\vec{x}|\omega_2)P(\omega_2)$

bzw. iff $\underbrace{\frac{p(\vec{x}|\omega_1)}{p(\vec{x}|\omega_2)}}_{>} > \frac{(\lambda_{12} - \lambda_{22})P(\omega_2)}{(\lambda_{21} - \lambda_{11})P(\omega_1)} \stackrel{!}{=} \text{konstanter Schwellwert}$

'likelihood ratio': Bayessche Regel in Abhängigkeit eines von der Beobachtung \vec{x} unabhängigen Schwellwertes!

¹ Beachte: Kosten für einen Fehler sind normalerweise höher als für eine richtige Entscheidung, daher $\lambda_{21} > \lambda_{11}$ und $\lambda_{12} > \lambda_{22}$ im allgemeinen Fall

Bayessche Entscheidungstheorie

- Klassifikation mit minimaler Fehlerrate
 - symmetrische (oder 'zero-one') Kostenfunktion

$$\ddot{e}(\alpha_i | \omega_j) = \begin{cases} 0, & i = j \\ 1, & i \neq j \end{cases} \quad i, j = 1, \dots, c \text{ (Klassen)}$$

$$\Rightarrow \mathbf{L} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \quad \begin{array}{l} 0: \text{ keine Kosten} \\ 1: \text{ Einheitskosten bei Fehler} \end{array}$$

- d.h. Entscheidung für α_i ist gefallen, ω_j liegt tatsächlich vor, dann ist Entscheidung korrekt für $i = j$ (wegen Nullkosten)
- für das 2-Klassen-Problem gilt nun

$$\lambda_{11} = \lambda_{22} = 0 \quad \text{und} \quad \lambda_{12} = \lambda_{21} = 1$$

und

$$R(\alpha_1 | \vec{x}) = \lambda_{12} P(\omega_2 | \vec{x}) = P(\omega_2 | \vec{x})$$

$$R(\alpha_2 | \vec{x}) = \lambda_{21} P(\omega_1 | \vec{x}) = P(\omega_1 | \vec{x})$$

- da nach Bayes-Regel gilt:

$$\vec{x} \sim \omega_1 \quad \text{iff} \quad P(\omega_1 | \vec{x}) > P(\omega_2 | \vec{x}) \quad \text{mit Fehler} \quad P(\text{Fehler} | \vec{x}) = P(\omega_2 | \vec{x})$$

muss auch gelten: $R(\alpha_1 | \vec{x}) < R(\alpha_2 | \vec{x})$

und somit ist die Forderung nach der Minimierung des Risikos erfüllt!

Bayessche Entscheidungstheorie

- d.h. das mit dieser Kosten-/Verlust-Funktion verbundene Risiko entspricht der mittleren Fehler-Wahrscheinlichkeit, da

$$\begin{aligned}
 R(\alpha_i | \vec{x}) &= \sum_{j=1}^c \lambda(\alpha_i | \omega_j) \cdot P(\omega_j | \vec{x}) \\
 &= \sum_{j \neq i} P(\omega_j | \vec{x}) \\
 &= 1 - \underbrace{P(\omega_i | \vec{x})}
 \end{aligned}$$

bedingte Wahrscheinlichkeit, dass α_i korrekt ist!

- ergo: „Um die mittlere Fehler-Wahrscheinlichkeit zu minimieren, muss i so gewählt werden, dass die a posteriori- Wahrscheinlichkeit $P(\omega_i | \vec{x})$ maximiert wird“, d.h. minimale Fehlerrate nur für $\vec{x} \sim \omega_i$ wenn $P(\omega_i | \vec{x}) > P(\omega_j | \vec{x})$, $\forall i \neq j$
- Klassifikatoren, Diskriminantenfunktionen und Entscheidungs-(Trenn-)Hyperflächen (siehe: Teil C „Diskriminantenfunktionen...“ zur Einführung)
 - Diskriminantenfunktionen $g_i(\vec{x})$
 - Entscheidungsregel $\vec{x} \sim \omega_i$ wenn $g_i(\vec{x}) > g_j(\vec{x})$, $\forall i \neq j$, $i = 1, 2, \dots, c$ (Klassen)
- d.h. Aufteilung des Merkmalraumes durch Trenn-Hyperflächen in c „Entscheidungsgebiete“ $\mathbb{X}_1, \dots, \mathbb{X}_c$, so dass \vec{x} innerhalb von \mathbb{X}_i und $g_i(\vec{x}) = g_j(\vec{x})$ ist Gleichung der Trennfläche(n)

Bayessche Entscheidungstheorie

- Beziehung zu Bayes-Klassifikator

- allgemeiner Fall:
$$g_i(\vec{x}) = \underbrace{-R(\omega_i|\vec{x})}_{\text{minimales bedingtes Risiko}}$$

maximale Diskriminantenfunktion

- Fall der minimalen Fehlerrate:
$$g_i(\vec{x}) = \underbrace{P(\omega_i|\vec{x})}_{\text{maximale a posteriori-Wahrscheinlichkeit}}$$

maximale Diskriminantenfunktion

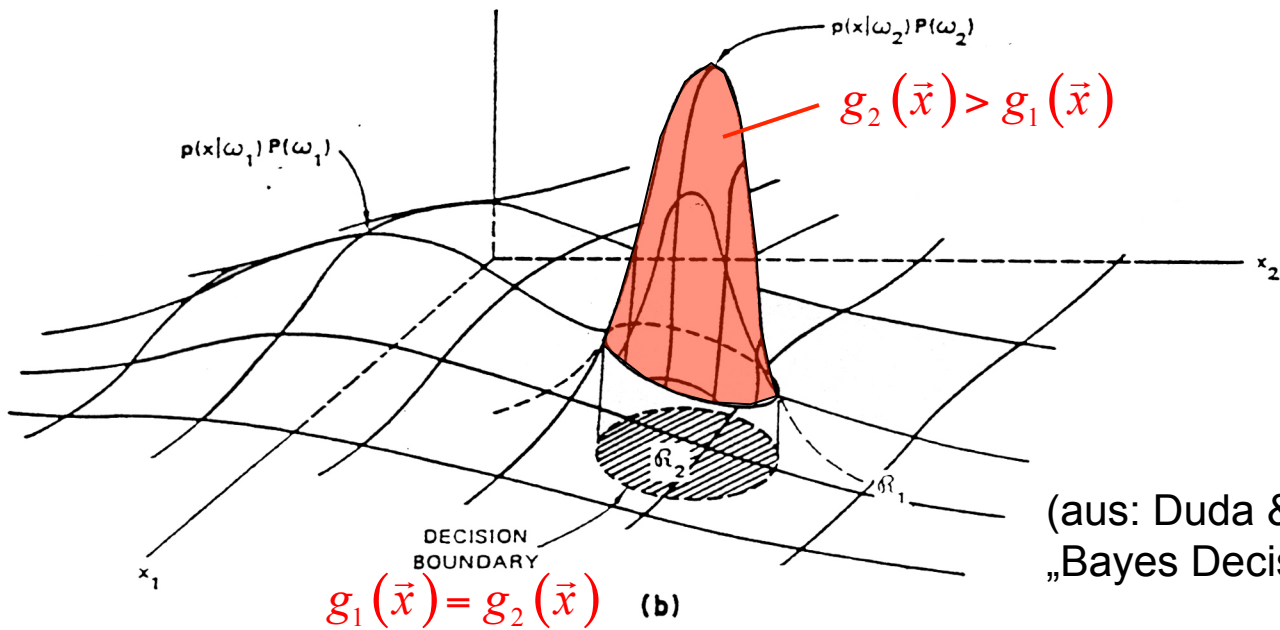
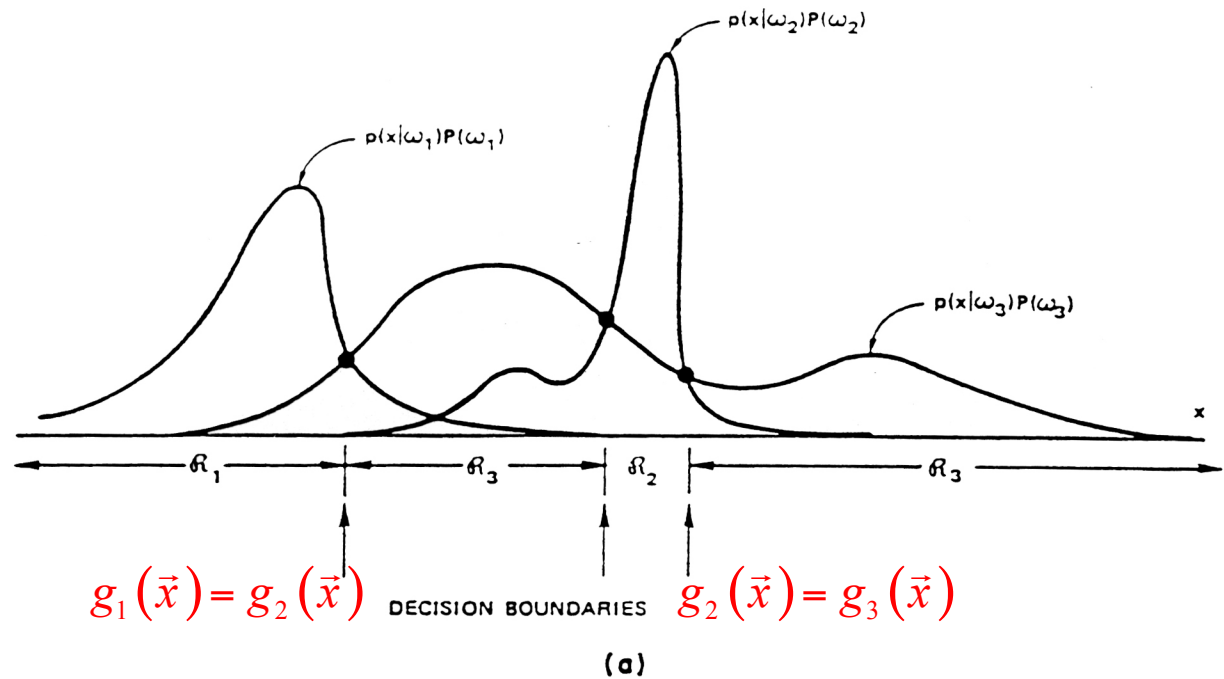
! Multiplikation von $g_i(\vec{x})$ mit einer positiven Konstanten oder Addition einer Konstanten beeinflusst nicht die Entscheidung, d.h. z.B. Ersetzung von $g_i(\vec{x})$ durch $f(g_i(\vec{x}))$, wobei f eine streng monoton ansteigende Funktion sei, hat keinen Einfluss auf Klassifikationsresultat.

⇒ z.B. verschiedene Ansätze zur 'minimum error rate'-Klassifikation mit verschiedenem **Aufwand** aber identischen Klassifikationsergebnissen:

$$g_i(\vec{x}) = \underbrace{P(\omega_i|\vec{x})}_{\text{maximale a posteriori-Wahrscheinlichkeit}}$$
$$g_i(\vec{x}) = \frac{p(\vec{x}|\omega_i) \cdot P(\omega_i)}{\sum_{j=1}^c p(\vec{x}|\omega_j) \cdot P(\omega_j)}$$
$$g_i(\vec{x}) = p(\vec{x}|\omega_i) \cdot P(\omega_i)$$
$$g_i(\vec{x}) = \log p(\vec{x}|\omega_i) + \log P(\omega_i)$$

Aufwandsreduktion!
(Klassengrenzen verschieben sich nicht!)

Beispiel:



(aus: Duda & Hart, 1973, S. 19, Kap. 2 „Bayes Decision Theory“)

FIGURE 2.4. Examples of decision boundaries and decision regions.

Bayessche Entscheidungstheorie

Beispiel: 2-Klassen-Fall

anstatt: $\vec{x} \sim \omega_1$ wenn $g_1(\vec{x}) > g_2(\vec{x})$

üblicher: $g(\vec{x}) = g_1(\vec{x}) - g_2(\vec{x})$ so dass

$\vec{x} \sim \omega_1$ wenn $g(\vec{x}) > 0$



technische Aufwandsreduktion, da nur
Vorzeichenprüfung (anstatt Maximumberechnung)!

mit 'minimum error rate'-Diskriminantenfunktion $g(\vec{x})^*$:

$$g(\vec{x}) = P(\omega_1|\vec{x}) - P(\omega_2|\vec{x})$$

$$\text{oder } g(\vec{x}) = \log \frac{p(\vec{x}|\omega_1)}{p(\vec{x}|\omega_2)} + \log \frac{P(\omega_1)}{P(\omega_2)}$$

(* ... „following two are particularly convenient...“, Duda & Hart, 1973)

Bayessche Entscheidungstheorie

- Fehler-Wahrscheinlichkeiten und Fehler-Integrale

- Entscheidungs-„Gebiete“ \mathbb{W}_1 und \mathbb{W}_2 im 2-Klassen-Fall

- 2 Fehlertypen:

- \bar{x} fällt in \mathbb{W}_2 , aber wahrer Zustand ist ω_1
- \bar{x} fällt in \mathbb{W}_1 , aber wahrer Zustand ist ω_2

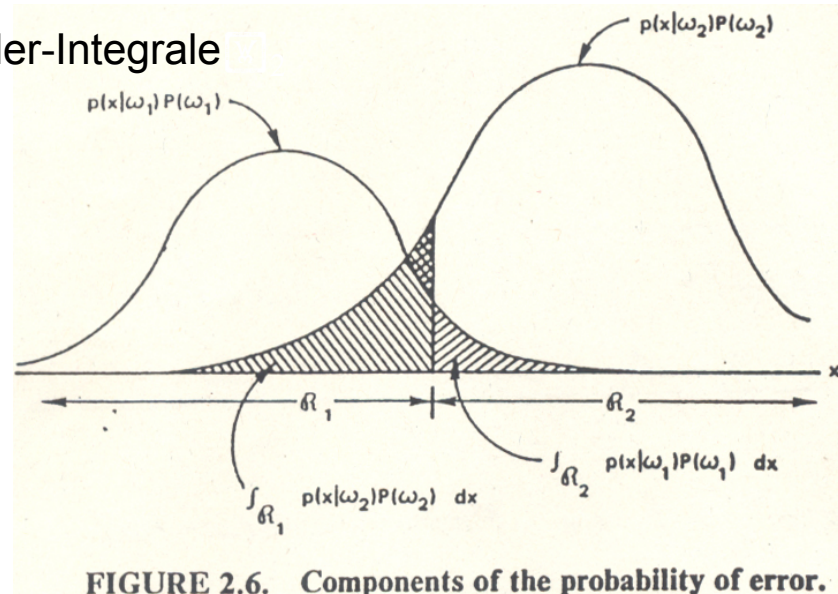


FIGURE 2.6. Components of the probability of error.

damit:

$$P(\text{Fehler}) = P(\bar{x} \in \mathbb{R}_2, \omega_1) + P(\bar{x} \in \mathbb{R}_1, \omega_2)$$

$$= P(\bar{x} \in \mathbb{R}_2 | \omega_1) \cdot P(\omega_1) + P(\bar{x} \in \mathbb{R}_1 | \omega_2) \cdot P(\omega_2)$$

$$\Rightarrow P(\text{Fehler}) = \int_{\mathbb{R}_2} p(\bar{x} | \omega_1) \cdot P(\omega_1) d\bar{x} + \int_{\mathbb{R}_1} p(\bar{x} | \omega_2) \cdot P(\omega_2) d\bar{x}$$

„This result is illustrated in the one-dimensional case in Figure 2.6. The two terms in the sum are merely the areas in the tails of the functions $p(\mathbf{x}|\omega_i)P(\omega_i)$. Because the regions \mathbb{R}_1 and \mathbb{R}_2 were chosen arbitrarily, the probability of error is not as small as it might be. By moving the decision boundary to the left, it is clear that we can eliminate the dark „triangular“ area and reduce the probability of error. In general, if $p(\mathbf{x}|\omega_1)P(\omega_1) > p(\mathbf{x}|\omega_2)P(\omega_2)$, it is advantageous to have \mathbf{x} be in \mathbb{R}_1 so that the smaller quantity will contribute to the integral; this is exactly what the Bayes decision rule achieves.“

Bayessche Entscheidungstheorie

- Alternative, z.B. im Mehr-Klassen-Fall

$$\begin{aligned} P(\text{Korrektheit}) &= \sum_{i=1}^c P(\vec{x} \in \mathbb{R}_i, \omega_i) \quad (\text{Summation über alle Klassen!}) \\ &= \sum_{i=1}^c P(\vec{x} \in \mathbb{R}_i | \omega_i) \cdot P(\omega_i) \end{aligned}$$

$$P(\text{Korrektheit}) = \sum_{i=1}^c \int_{\mathbb{R}_i} p(\vec{x} | \omega_i) P(\omega_i) d\vec{x}$$

! Bayes-Klassifikator maximiert $P(\text{Korrektheit})$, da er \mathbb{R}_i so festlegt, dass Integranden maximal (und damit $P(\text{Fehler})$ minimal) werden.

„In the multiclass case, there are more ways to be wrong than to be right, and it is simpler to compute the probability of being correct.“
(Duda & Hart, 1973, S. 21)

Bayessche Entscheidungstheorie

- Normal-Verteilung (NV) als klassenbedingte Wahrscheinlichkeits-Dichtefunktion
 - Motivation:
 - Bedeutung von $p(\vec{x}|\omega_i)$ für Bayes-Klassifikator
 - analytische „Handhabbarkeit“
 - \vec{x} einer Klasse ω_i streuen um Prototypen $\vec{\mu}_i$ (NV als „Modell“)
 - univariate normalverteilte Dichtefunktion

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

mit μ : Mittelwert,
 σ : Standardabweichung,
 σ^2 : Varianz

wobei
$$\mu = \int_{-\infty}^{\infty} x \cdot p(x) dx = E(x)$$

$$\sigma^2 = \int_{-\infty}^{\infty} (x - \mu)^2 \cdot p(x) dx = E\left((x - \mu)^2\right)$$

Notation¹: $p(x) \sim N(\mu, \sigma^2)$

¹ Zur Erinnerung: $\int_{\mu-2\sigma}^{\mu+2\sigma} p(x) dx \approx 0,95$

Bayessche Entscheidungstheorie

- multivariate normalverteilte Dichtefunktion in d Dimensionen

- $p(\vec{x}) \sim N(\vec{\mu}, \Sigma)$

- $\vec{x} = (x_1, x_2, \dots, x_d)^T$

- $\vec{\mu} = (\mu_1, \mu_2, \dots, \mu_d)^T$

- Σ : $d \times d$ – Kovarianzmatrix (Σ^{-1} : Inverse, $|\Sigma|$: Determinante)

$$p(\vec{x}) = \frac{1}{(2\pi)^{d/2} \cdot |\Sigma|^{1/2}} \cdot \exp\left(-\frac{1}{2}(\vec{x} - \vec{\mu})^T \Sigma^{-1}(\vec{x} - \vec{\mu})\right)$$

mit $\vec{\mu} = E(\vec{x}); \mu_i = E(x_i)$

und $\Sigma = E\left(\left(\vec{x} - \vec{\mu}\right)\left(\vec{x} - \vec{\mu}\right)^T\right)^1$

$$= \begin{pmatrix} \sigma_{11} & \cdots & \sigma_{1n} \\ \vdots & \ddots & \vdots \\ \sigma_{n1} & \cdots & \sigma_{nn} \end{pmatrix},$$

σ_{ii} : Varianz von x_i

$$\sigma_{ij} = \underbrace{E\left(\left(x_i - \mu_i\right)\left(x_j - \mu_j\right)\right)}_{\text{Kovarianz von } x_j \text{ und } x_j}$$

! $\sigma_{ij} = 0$: x_i und x_j sind statistisch unabhängig

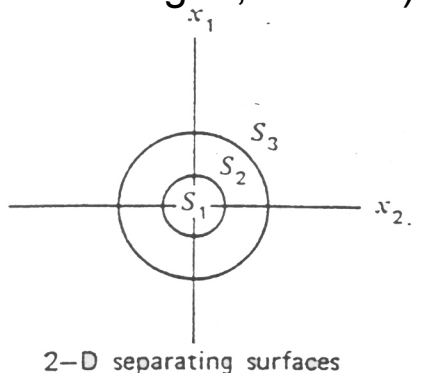
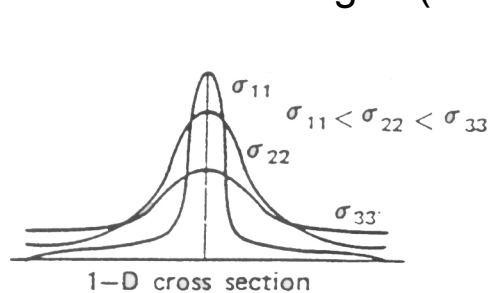
! beliebige Linearkombination von N : $\vec{y} = \mathbf{A}^T \vec{x} \Rightarrow p(\vec{y}) \sim N(\mathbf{A}^T \vec{\mu}, \mathbf{A}^T \Sigma \mathbf{A})$

¹ Σ ist immer symmetrisch und positiv-semidefinit, d.h. $\vec{x}^T \Sigma \vec{x} \geq 0, \forall \vec{x}$ (s. Lin. Alg.) 23

Bayessche Entscheidungstheorie

! Bedeutung der Kovarianzmatrix:

erlaubt Berechnung der Streuung ('dispersion') von Daten (Vektoren) im Merkmalraum bzw. von Datenhäufungen (auch: Ballungen, 'cluster') in z.B. jede beliebige Richtung



$$\begin{aligned} [\Phi_1] &= b_1 [I] \\ [\Phi_2] &= b_2 [I] \\ [\Phi_3] &= b_3 [I] \end{aligned}$$

(a) Equal means, uncorrelated covariances

$$[\Phi_1] = [\Phi_2] = [\Phi_3] \neq [I]$$

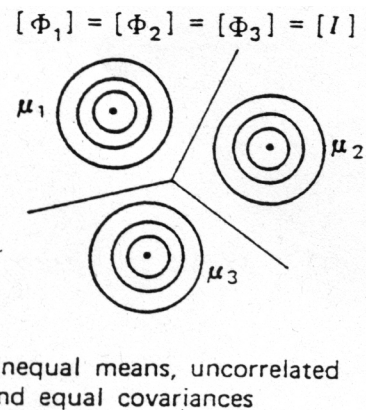
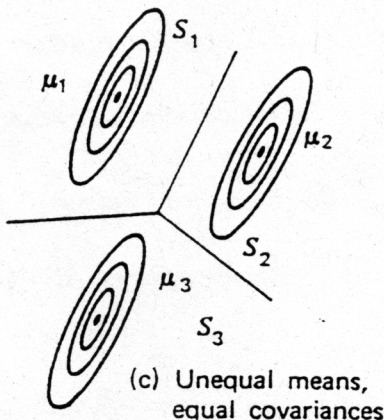
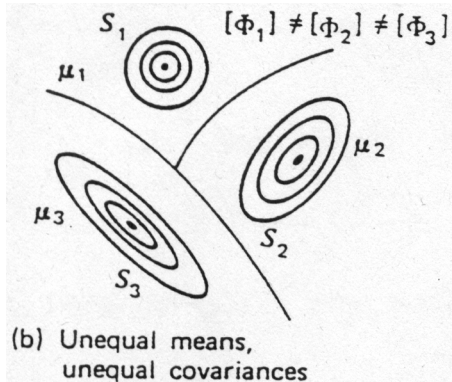
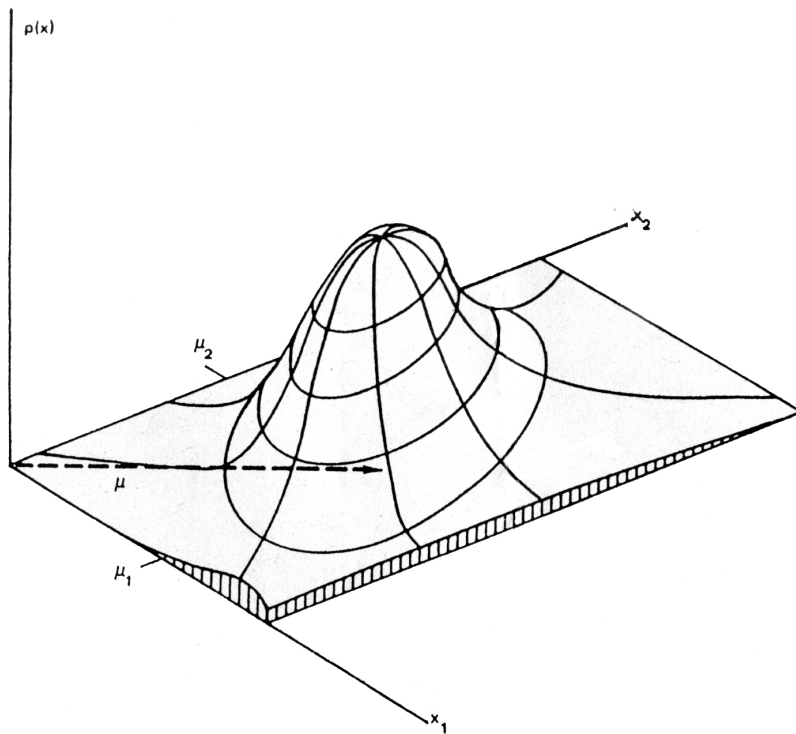


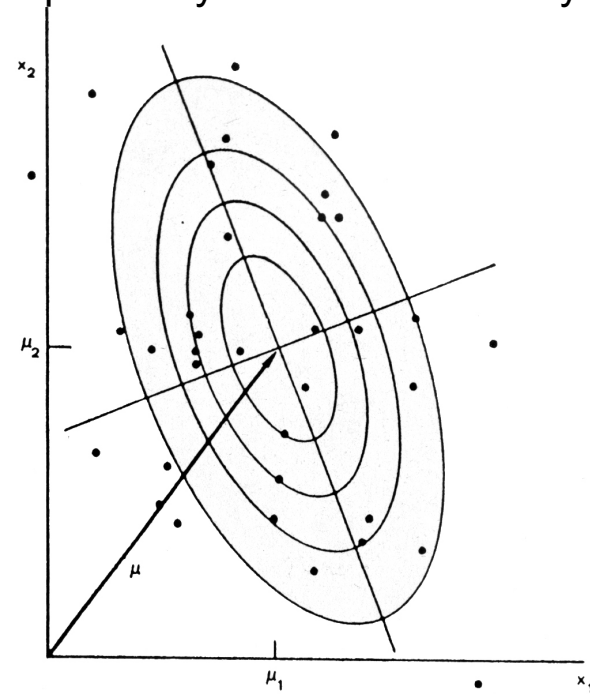
FIGURE 4.4. Some typical separating surfaces for $N(\mu_k, [\Phi_k])$, ($N = 2, k = 3$).

Bayessche Entscheidungstheorie

Beispiel Fig. 2.7 (aus: Duda & Hart, 1973, S. 25, Kap. 2 "Bayes Decision Theory")



(a) BIVARIATE NORMAL DENSITY



(b) SCATTER DIAGRAM

FIGURE 2.7. Two representations of a normal density.

- Interpretation

- Zentrum des 'cluster' bei $\vec{\mu} = (\mu_1, \mu_2)^T$
- „Form“ des 'cluster' gegeben durch Σ , d.h. Orte der Punkte gleicher Dichtefunktionswerte sind **Hyperellipsoide**, für die $(\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu}) = \text{const}$, wobei
 - 1) deren Hauptachsen definiert sind durch die Eigenvektoren von Σ
 - 2) Länge der Hauptachsen definiert durch Eigenwerte von Σ

Bayessche Entscheidungstheorie

- Mahalanobis-Distanz zwischen \vec{x} und $\vec{\mu}$

$$r^2 = (\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu})$$

d.h. r^2 ist konstant für Punkte auf Isodichte-Linien!

- Streuung ('scatter') von \vec{x} um $\vec{\mu}$ aus Volumina der Hyperellipsoiden V (für best. r)

$$V = V_d |\Sigma|^{1/2} r^d \quad \text{als 'scatter'-Maß}$$

$$\text{mit } V_d = \begin{cases} \frac{\pi^{d/2}}{(\frac{d}{2})!}, & \text{für } d \text{ gerade} \\ \frac{2^d \pi^{d-1/2} (\frac{d-1}{2})!}{d!}, & \text{für } d \text{ ungerade} \end{cases}$$

(V_d : Volumen einer d -dimensionalen Einheitskugel!)

Anmerkung: Für $d = 2$ gilt: $\Sigma = \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{pmatrix}$ und $V = \pi (\sigma_{11}\sigma_{22} - \sigma_{12}\sigma_{21}) r^2$,

d.h. für gegebenes d und beliebiges r hängt Streuung ('scatter') direkt von $|\Sigma|^{1/2}$ ab!!

Bayessche Entscheidungstheorie

- Diskriminantenfunktionen für die normalverteilte Dichtefunktion

- $g_i(\vec{x}) = \log p(\vec{x}|\omega_i) + \log P(\omega_i)$ für 'minimum-error-rate'-Klassifikation (s.o.)

- mit $p(\vec{x}|\omega_i) \sim N(\vec{\mu}_i, \Sigma_i)$:

$$g_i(\vec{x}) = -\frac{1}{2}(\vec{x} - \vec{\mu}_i)^T \Sigma_i^{-1}(\vec{x} - \vec{\mu}_i) - \frac{d}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_i| + \log P(\omega_i)$$

- Fall 1:
 - Merkmale statistisch unabhängig, $\sigma_{ij} = 0$
 - jedes Merkmal hat gleiche Varianz σ^2

$\Rightarrow \Sigma_i = \sigma^2 I$, d.h. gleichgroße Hyperkugel-"cluster" (pro Klasse) um $\vec{\mu}_i$

- $|\Sigma_i| = \sigma^{2d}$ und $\Sigma_i^{-1} = \left(\frac{1}{\sigma^2}\right) I$
 unabhängig von i : additive Konstanten für $g_i(\vec{x})$ werden vernachlässigt!

$$\Rightarrow g_i(\vec{x}) = -\frac{\|\vec{x} - \vec{\mu}_i\|^2}{2\sigma^2} + \log P(\omega_i)$$

mit euklidischer Norm $\|\cdot\|$, d.h. $\|\vec{x} - \vec{\mu}_i\|^2 = (\vec{x} - \vec{\mu}_i)^T (\vec{x} - \vec{\mu}_i)$

- falls $P(\omega_i)$ gleich für alle Klassen:

berechne $\|\vec{x} - \vec{\mu}_i\|$ ($\hat{=}$ euklidischer Abstand) von \vec{x} zu allen $\vec{\mu}_i$ bei c Klassen,

$\vec{x} \sim \omega_i$ wenn $\|\vec{x} - \vec{\mu}_i\| = \text{Minimum}$

\rightarrow „minimaler Abstand“-Klassifikator

Bayessche Entscheidungstheorie

• jedoch: $g_i(\vec{x}) = -\frac{(\vec{x} - \vec{\mu}_i)^T (\vec{x} - \vec{\mu}_i)}{2\sigma^2} + \log P(\omega_i)$

$$= -\frac{1}{2\sigma^2} (\underbrace{\vec{x}^T \vec{x}} - 2\vec{\mu}_i^T \vec{x} + \vec{\mu}_i^T \vec{\mu}_i) + \log P(\omega_i)$$

gleich für alle i (damit additive Konstante)!

somit ist $g_i(\vec{x})$ linearisiert (s.a. Teil C: „Diskriminantenfunktionen...“):

$$\boxed{g_i(\vec{x}) = \vec{w}_i^T \vec{x} + w_{i_0}} \quad \boxtimes \text{ „lineare Maschine“}$$

mit $\vec{w}_i = \frac{1}{\sigma^2} \vec{\mu}_i$

und $w_{i_0} = -\frac{1}{2\sigma^2} \vec{\mu}_i^T \vec{\mu}_i + \log P(\omega_i)$

⇒ Trennfläche $g_i(\vec{x}) = g_j(\vec{x})$ (allgemein)

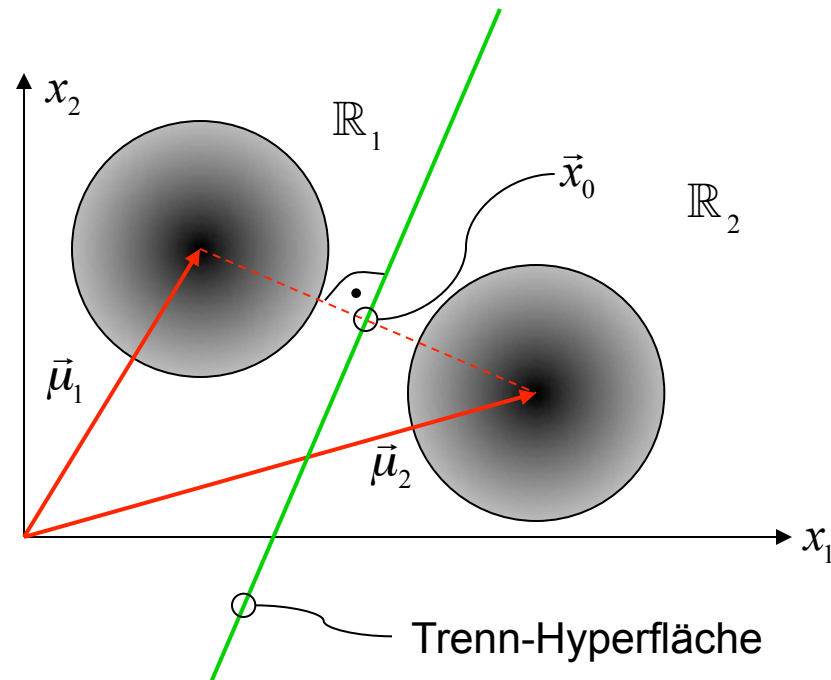
$$\vec{w}^T (\vec{x} - \vec{x}_0) = 0 \quad (\text{Fall 1})$$

mit $\vec{w} = \vec{\mu}_i - \vec{\mu}_j$

und $\vec{x}_0 = \frac{1}{2}(\vec{\mu}_i + \vec{\mu}_j) - \frac{\sigma^2}{\|\vec{\mu}_i - \vec{\mu}_j\|^2} \log \frac{P(\omega_i)}{P(\omega_j)} \cdot (\vec{\mu}_i - \vec{\mu}_j)$

Bayessche Entscheidungstheorie

d.h. Hyperebene durch \vec{x}_0 und orthogonal zu \vec{w} ;
da $\vec{w} = \vec{\mu}_i - \vec{\mu}_j$ ist Hyperebene orthogonal zu Verbindungslinie
zwischen $\vec{\mu}_i$ und $\vec{\mu}_j$ (für $P(\omega_i) = P(\omega_j)$ halbiert \vec{x}_0 diese Linie)



Bayessche Entscheidungstheorie

- Fall 2:
 - Kovarianz-Matrizen für alle Klassen gleich
 - d.h. $\Sigma_i = \Sigma$: gleiche Hyperellipsoid-"cluster" (pro Klasse um $\vec{\mu}_i$)
 - $|\Sigma_i|$ und $(d/2) \cdot \log 2\pi$ ebenso unabhängig von i

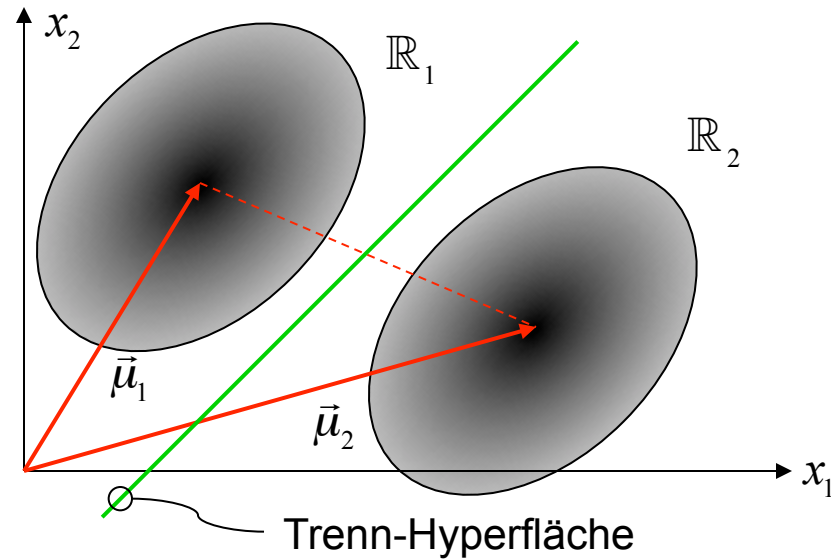
$$\Rightarrow \boxed{g_i(\vec{x}) = -\frac{1}{2} \underbrace{(\vec{x} - \vec{\mu}_i)^T \Sigma^{-1} (\vec{x} - \vec{\mu}_i)}_{r^2} + \log P(\omega_i)}$$

- falls $P(\omega_i)$ gleich für alle Klassen:
berechne Mahalanobis-Distanzen r_i^2 von \vec{x} zu allen $\vec{\mu}_i$
 $\vec{x} \sim \omega_i$ wenn $r_i^2 = \text{Minimum}$
- jedoch: $\vec{x}^T \Sigma^{-1} \vec{x}$ unabhängig von i und kann vernachlässigt werden,
somit $g_i(\vec{x}) = \vec{w}_i^T \vec{x} + w_{i_0}$
mit $\vec{w}_i = \Sigma^{-1} \vec{\mu}_i$
und $w_{i_0} = -\frac{1}{2} \vec{\mu}_i^T \Sigma^{-1} \vec{\mu}_i + \log P(\omega_i)$
 $\Rightarrow g_i(\vec{x})$ linearisiert!

Bayessche Entscheidungstheorie

- Trennfläche $\vec{w}^T (\vec{x} - \vec{x}_0) = 0$ mit $\vec{w} = \Sigma^{-1} (\vec{\mu}_i - \vec{\mu}_j)$ (!)

d.h. Hyperebene zwischen \mathbb{R}_1 und \mathbb{R}_2 i.a. nicht orthogonal, wie im Fall 1!



Bayessche Entscheidungstheorie

- Fall 3:
 - allgemeiner multivariater normalverteilter Fall, d.h. klassenspezifische Kovarianz-Matrizen
 - nur $(d/2) \cdot \log 2\pi$ kann vernachlässigt werden, d.h. quadratische Diskriminanten-Funktion

$$g_i(\vec{x}) = \vec{x}^T \mathbf{W}_i \vec{x} + \vec{w}_i^T \vec{x} + w_{i_0}$$

$$\text{mit } \mathbf{W}_i = -\frac{1}{2} \boldsymbol{\Sigma}_i^{-1}$$

$$\vec{w}_i = \boldsymbol{\Sigma}_i^{-1} \vec{\mu}_i$$

$$w_{i_0} = -\frac{1}{2} \vec{\mu}_i^T \boldsymbol{\Sigma}_i^{-1} \vec{\mu}_i - \frac{1}{2} \log |\boldsymbol{\Sigma}_i| + \log P(\omega_i)$$

⇒ 'hyperquadrics' als Trennflächen:

Paare von Hyperkugeln

oder Hyperellipsoiden

oder Hyperparaboloiden

oder Hyperboloiden und Hyperebenen

für $d = 2$ (also \mathbb{R}_1 und \mathbb{R}_2)

! Grenzen des Perzeptrons bei normalverteilten Daten im Merkmalraum