



5. Mai 2017

Klassifizierung

Klassifizierung: Beispiel

5. Mai 2017

Automatische Unterscheidung
zw. Lachs und Seebarsch (aus
Duda et. al.: Pattern
classification)

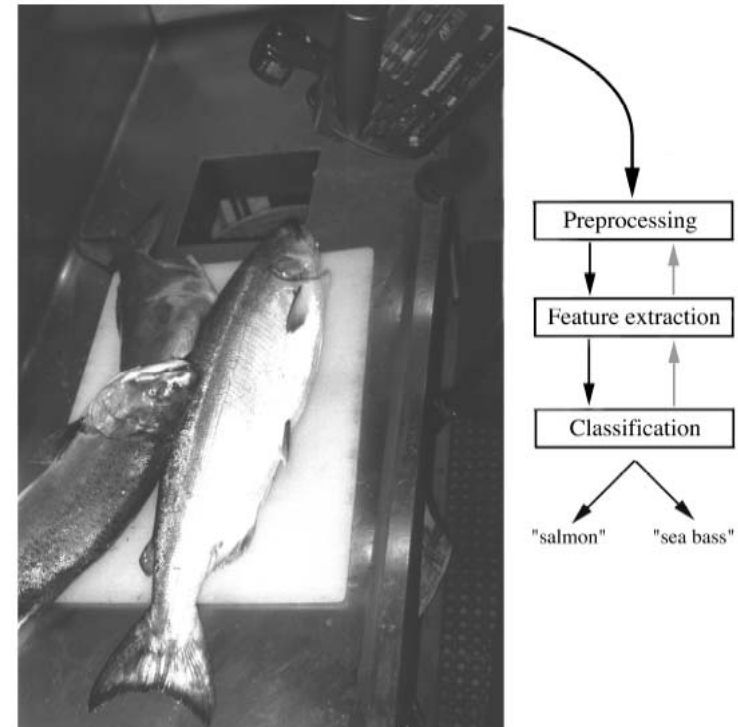


FIGURE 1.1. The objects to be classified are first sensed by a transducer (camera), whose signals are preprocessed. Next the features are extracted and finally the classification is emitted, here either "salmon" or "sea bass." Although the information flow is often chosen to be from the source to the classifier, some systems employ information flow in which earlier levels of processing can be altered based on the tentative or preliminary response in later levels (gray arrows). Yet others combine two or more stages into a unified step, such as simultaneous segmentation and feature extraction. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



5. Mai 2017

Länge:

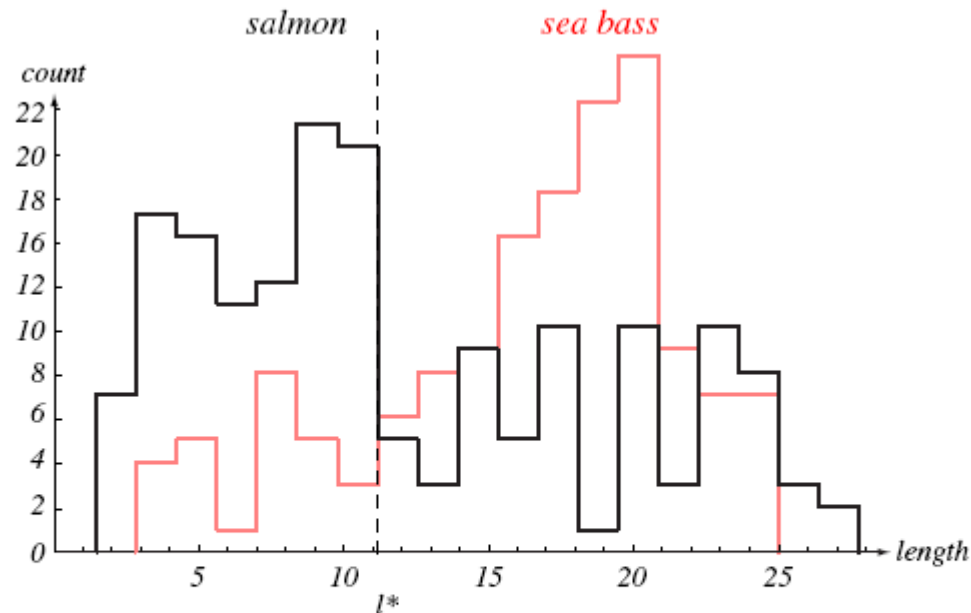


FIGURE 1.2. Histograms for the length feature for the two categories. No single threshold value of the length will serve to unambiguously discriminate between the two categories; using length alone, we will have some errors. The value marked l^* will lead to the smallest number of errors, on average. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



5. Mai 2017

Wahl der Schwelle:

- Berücksichtigung von Kosten für Falschentscheidung bzw. Gewinn f. richtige Entscheidung
- Wähle Schwelle, die minimale Kosten verursacht



5. Mai 2017

Helligkeit:

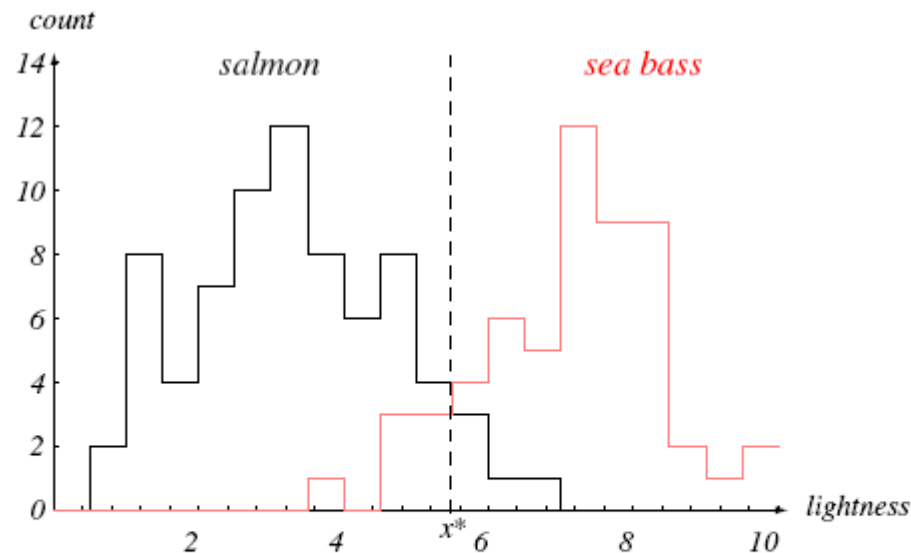


FIGURE 1.3. Histograms for the lightness feature for the two categories. No single threshold value x^* (decision boundary) will serve to unambiguously discriminate between the two categories; using lightness alone, we will have some errors. The value x^* marked will lead to the smallest number of errors, on average. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



5. Mai 2017

Kombination mehrerer Merkmale:

- Finden von Entscheidungsgrenzen in mehrdimensionalen Merkmalsräumen

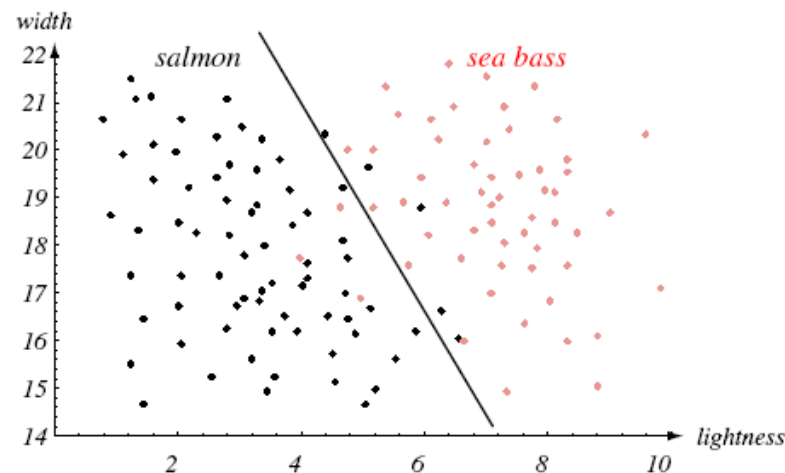


FIGURE 1.4. The two features of lightness and width for sea bass and salmon. The dark line could serve as a decision boundary of our classifier. Overall classification error on the data shown is lower than if we use only one feature as in Fig. 1.3, but there will still be some errors. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



5. Mai 2017

Komplizierte Grenze:

- Perfekt für Trainingsdaten:

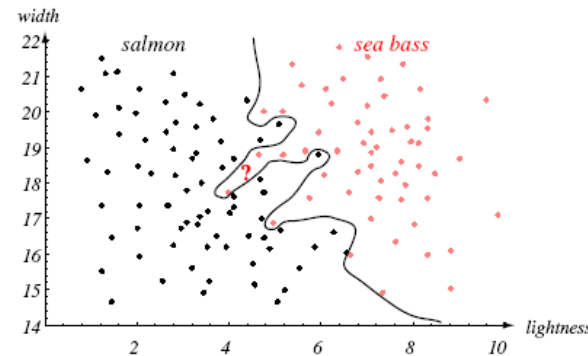


FIGURE 1.5. Overly complex models for the fish will lead to decision boundaries that are complicated. While such a decision may lead to perfect classification of our training samples, it would lead to poor performance on future patterns. The novel test point marked ? is evidently most likely a salmon, whereas the complex decision boundary shown leads it to be classified as a sea bass. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

- Schlecht für neue Daten.

Kriterium : Einfachheit der Grenze

Einfachere Grenze:

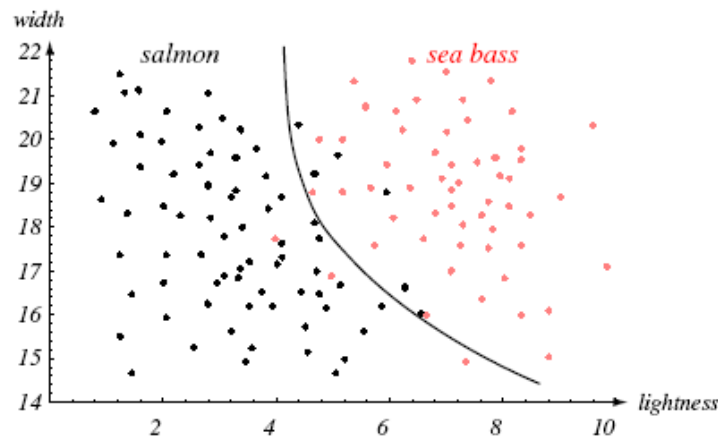


FIGURE 1.6. The decision boundary shown might represent the optimal tradeoff between performance on the training set and simplicity of classifier, thereby giving the highest accuracy on new patterns. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



5. Mai 2017

Bayes'sche Entscheidungstheorie



5. Mai 2017

Parameterschätzung

Grundlagen nach „Data Analysis – A Bayesian Tutorial“ von D. S.
Sivia, Oxford University Press 1996



Parameterschätzung

5. Mai 2017

Allgemeine Auffassung oft:

- Wahrscheinlichkeitsrechnung = Mathematik
- Statistik = Rezepte aus obskuren Kochbüchern
- Zitat aus Numerical Recipes:

In other words, we identify the probability of the data given the parameters (which is a mathematically computable number), as the *likelihood* of the parameters given the data. This identification is entirely based on intuition. It has no formal mathematical basis in and of itself; as we already remarked, statistics is *not* a branch of mathematics!



5. Mai 2017

- Lösung: Bayes'sche Theorie



Reverend Thomas Bayes (1702-1761)



Grundlagen

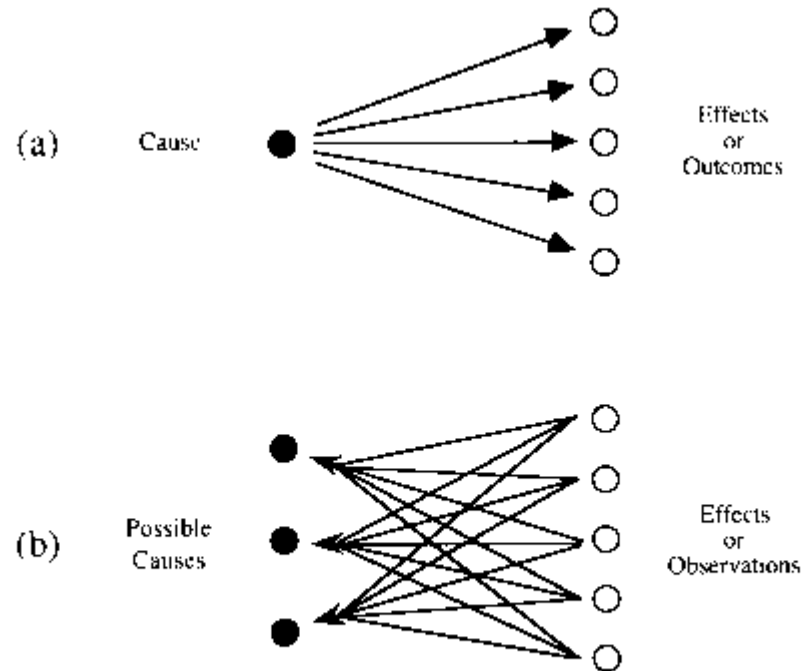


Fig. 1.1 A schematic representation of (a) deductive logic, or pure mathematics, and (b) plausible reasoning, or inductive logic.

[aus Sivia: Data Analysis, Oxford 1996]



Grundlagen

5. Mai 2017

Grundlagen

- Logisch konsistenter Rahmen für quantitative Festlegung der „Glaubwürdigkeit“ (Cox 1946 auf der Grundlage von Bayes (1763) und Laplace (1812))
- Reelle Zahl als Maß für das Vertrauen, daß eine Aussage richtig ist
- Zwei einfache Axiome:
 - Festlegung wie glaubhaft A ist, legt fest wie glaubhaft (nicht A) ist.
 - Festlegung, wie glaubhaft Y ist und wie glaubhaft X ist, wenn Y wahr ist legt fest, wie glaubhaft (X und Y) ist.



5. Mai 2017

Ergebnis:

- Zahlen für „Glaubhaftigkeit“ müssen sich abbilden lassen auf positive Zahlen, die der Summenregel und der Produktregel der Wahrscheinlichkeitsrechnung gehorchen.



Summenregel:

$$\text{prob}(X | I) + \text{prob}(\bar{X} | I) = 1$$

Produktregel:

$$\text{prob}(X, Y | I) = \text{prob}(X | Y, I) \times \text{prob}(Y | I)$$

Daraus lassen sich alle anderen Regeln unter Verwendung der Boole'schen Algebra herleiten.



Boole'sche Algebra

Für Aussagen gelten die Regeln der Boole'schen Algebra. Diese können wir mit einfachen Axiomen definieren: Eine Boole'schen Algebra ist eine Menge mit zwei Verknüpfungen (Produkt und Summe) auf dieser Menge so dass für Elemente A und B aus dieser Menge gilt:

Es gibt zwei Elemente 0 and 1, die nicht gleich sind.

[A1]

$$AB = BA$$

$$A+B = B+A$$

[A2]

$$A(B+C) = (AB)+(AC)$$

$$A+(BC) = (A+B)(A+C)$$

[A3]

$$1A = A$$

$$0+A = A$$

[A4]

$$A\bar{A} = 0$$

$$A+\bar{A} = 1$$

[A5]

In einer Zeile stehen jeweils zwei duale Axiome, die sich durch Vertauschung von Produkt und Summe, sowie von 0 und 1 ergeben. Im folgenden werden die linken Axiome mit a, die rechten mit b gekennzeichnet([A2a] bedeutet also $AB = BA$).

Aus Harri Valpola: <http://users.ics.aalto.fi/harri/thesis/bayesformulas.html>

Aus den Axiomen werden folgende Lemmata abgeleitet:

$$\neg\neg A = A \quad [L1]$$

$$AA = A \quad A+A = A \quad [L2]$$

$$\neg 1 = 0 \quad \neg 0 = 1 \quad [L3]$$

$$AB = 0 \ \& \ A+B = 1 \Rightarrow B = \neg A \quad [L4]$$

$$0A = 0 \quad 1+A = 1 \quad [L5]$$

$$A(A+B) = A \quad A+AB = A \quad [L6]$$

$$A(BC) = (AB)C \quad A+(B+C) = (A+B)+C \quad [L7]$$

$$\neg A(AB) = 0 \quad \neg A+(A+B) = 1 \quad [L8]$$

$$\neg(AB) = \neg A + \neg B \quad \neg(A+B) = \neg A \neg B \quad [L9]$$

$$AB = 1 \Rightarrow A = 1 \quad A+B = 0 \Rightarrow A = 0 \quad [L10]$$

Wenn nur die Elemente 0 und 1 verwendet werden, erhält man die Boole'sche Logik. 0 wird dabei als "falsch", 1 als "wahr" interpretiert. Produkt steht für das "und", die Summe für das "oder".

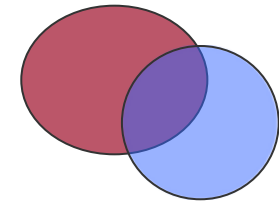
Aus Harri Valpola: <http://users.ics.aalto.fi/harri/thesis/bayesformulas.html>



5. Mai 2017

Beispiel: Herleitung der erweiterten Summenregel:

$P(A+B \mid C) =$	[L1]
$P(\neg\neg(A+B) \mid C) =$	[L7b]
$P(\neg(\neg A \neg B) \mid C) =$	[Sum Rule]
$1 - P(\neg A \neg B \mid C) =$	[Product Rule]
$1 - P(\neg A \mid C) P(\neg B \mid \neg AC) =$	[Sum Rule]
$1 - P(\neg A \mid C) [1 - P(B \mid \neg AC)] =$	
$1 - P(\neg A \mid C) + P(\neg A \mid C) P(B \mid \neg AC) =$	[Sum Rule]
$P(A \mid C) + P(\neg A \mid C) P(B \mid \neg AC) =$	[Product Rule]
$P(A \mid C) + P(\neg AB \mid C) =$	[A2a]
$P(A \mid C) + P(B \neg A \mid C) =$	[Product Rule]
$P(A \mid C) + P(B \mid C) P(\neg A \mid BC) =$	[Sum Rule]
$P(A \mid C) + P(B \mid C) [1 - P(A \mid BC)] =$	
$P(A \mid C) + P(B \mid C) - P(B \mid C) P(A \mid BC) =$	[Product Rule]
$P(A \mid C) + P(B \mid C) - P(BA \mid C) =$	[A2a]
$P(A \mid C) + P(B \mid C) - P(AB \mid C)$	



Aus Harri Valpola: <http://users.ics.aalto.fi/harri/thesis/bayesformulas.html>



Aus der Produktregel folgt u.a. der Satz von Bayes:

$$\text{prob}(X | Y) = \frac{\text{prob}(Y | X) \times \text{prob}(X)}{\text{prob}(Y)}$$

Herleitung:

$$\begin{aligned} \text{prob}(X, Y) &= \text{prob}(X | Y) \times \text{prob}(Y) \\ &= \text{prob}(Y, X) = \text{prob}(Y | X) \times \text{prob}(X) \end{aligned}$$



Wichtig für Datenanalyse und Parameterschätzung:

- X : Hypothese
- Y : Daten

$$\begin{aligned} & \text{prob}(Hypothese \mid Daten) \\ = & \frac{\text{prob}(Daten \mid Hypothese) \times \text{prob}(Hypothese)}{\text{prob}(Daten)} \end{aligned}$$



Namen:

- A-priori-Wahrscheinlichkeit für Hypothese

$$\text{prob}(Hypothese)$$

- Likelihood-Funktion

$$\text{prob}(Daten \mid Hypothese)$$

- A-posteriori-Wahrscheinlichkeit

$$\text{prob}(Hypothese \mid Daten)$$



5. Mai 2017

- A-priori-Wahrscheinlichkeit für Daten („evidence“)

$\text{prob}(\textit{Daten})$



5. Mai 2017

Da $\text{prob}(Daten)$ unabhängig von der Hypothese ist,
erhalten wir:

$$\begin{aligned} & \text{prob}(Hypothese \mid Daten) \\ & \propto \text{prob}(Daten \mid Hypothese) \times \text{prob}(Hypothese) \end{aligned}$$



„Marginalisation“, d.h. Elimination von nicht benötigten Variablen:

$$\text{prob}(X) = \sum_{k=1}^M \text{prob}(X, Y_k)$$

das gilt, wenn die Y_k sich zu eins addieren und gegenseitig ausschließen.



Kontinuierliche Form der Gleichung, z.B. wenn Y der Wert eines kontinuierlichen Parameters ist:

$$\text{prob}(X) = \int_{-\infty}^{\infty} \text{prob}(X, Y) dY$$

In diesem Fall bezeichnet $\text{prob}(X, Y)$ eine Wahrscheinlichkeitsdichtefunktion:

$$\text{prob}(X, Y = y) = \lim_{\delta y \rightarrow 0} \frac{\text{prob}(X, y \leq Y < y + \delta y)}{\delta y}$$



Klassische Wahrscheinlichkeit vs. Plausibilität im Bayes'schen Sinn:

- Klassische Wahrscheinlichkeit als relative Häufigkeit bei unendlich vielen Experimenten
- Wie aber anwenden auf Schätzung z.B. der Masse von Saturn aus wenigen Beobachtungen (Laplace)?
Ensemble von Welträumen?



Laplace berechnete die A-posteriori-Wahrscheinlichkeitsdichte für die Masse unter der Voraussetzung der verfügbaren Daten. Er schloß daraus, daß

‘... it is a bet of 11,000 to 1 that the error of this result is not 1/100th of its value’.

und behielt bis heute recht (0.63% daneben).

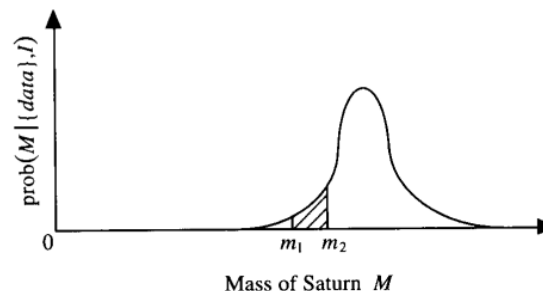


Fig. 1.2 A schematic illustration of the result of Laplace's probability analysis of the mass of Saturn.



Wahrscheinlichkeit in der Bayes'schen Auffassung repräsentiert den Wissensstand über eine physikalische Wirklichkeit, nicht diese selbst.

Beispiel:

- Man zieht verdeckt einen von 5 roten und 7 grünen Bällen: Wahrscheinlichk. f. rot: $5/12$
- Beim zweiten Zug (ohne Zurücklegen) hängt Wahrscheinlichkeit vom ersten Zug ab.



5. Mai 2017

Wie ist Wahrscheinlichkeit beim ersten Zug (dessen Ergebnis wir noch nicht kennen), wenn das Ergebnis des zweiten schon bekannt ist?

Immer noch $5/12$?

Nein, offensichtlich können wir mehr sagen (deutlich, wenn ursprünglich nur eine rote und eine grüne Kugel da waren).

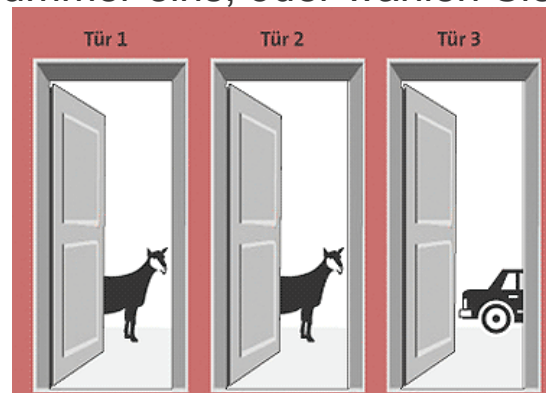
Die bedingten Wahrscheinlichkeiten stellen also logische, keine kausalen Zusammenhänge dar.

Das Ziegenproblem

5. Mai 2017

Die Ausgangssituation des Ziegenproblems lautet folgendermaßen: Sie sind Kandidat einer Fernsehshow und dürfen eine von drei verschlossenen Türen auswählen. Hinter einer der Türen wartet der Hauptgewinn, ein prachtvolles Auto, hinter den anderen beiden steht jeweils eine meckernde Ziege.

Frohgemut zeigen Sie auf eine der Türen, sagen wir Nummer eins. Doch der Showmaster, der weiß, hinter welcher Tür sich das Auto befindet, lässt sie nicht sofort öffnen, sondern sagt geheimnisvoll: »Ich zeige Ihnen mal was!« Er lässt eine andere Tür öffnen, sagen wir Nummer drei – und hinter dieser steht eine Ziege und glotzt erstaunt ins Publikum. Nun fragt der Showmaster lauernd: »Bleiben Sie bei Tür Nummer eins, oder wählen Sie doch lieber Nummer zwei?« Was sollten Sie tun?



aus: <http://www.zeit.de/2004/48/N-Ziegenproblem>, wo Sie auch die Lösung finden.



Parameterschätzung

5. Mai 2017

Beispiel: Ist die Münze in Ordnung?

- Wappen oder Zahl, Wahrscheinlichkeit Z für Zahl sollte $\frac{1}{2}$ sein.
- 4 von 11 waren Wappen, ist die Münze in Ordnung?
- Ansatz: Berechne $\text{prob}(Z/\{Daten\})$, also eine Wahrscheinlichkeitsdichtefunktion. Nach Satz von Bayes gilt:

$$\text{prob}(Z | Daten) \propto \text{prob}(Daten | Z) \times \text{prob}(Z)$$



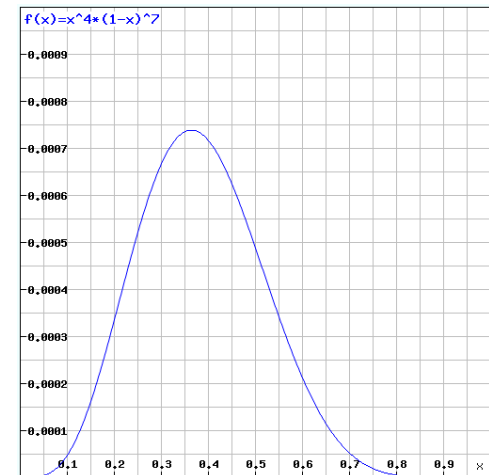
5. Mai 2017

Wir setzen an (A-priori-Wahrscheinlichkeit):

$$\text{prob}(Z) = \begin{cases} 1 & 0 \leq Z \leq 1 \\ 0 & \text{sonst} \end{cases}$$

Und Wahrscheinl. für R von N mal Zahl (Likelihood-Funktion):

$$\text{prob}(\text{Daten} \mid Z) \propto Z^R (1-Z)^{N-R}$$



Was ist wahrscheinlicher bei fairer Münze?

Diese Ziffernfolge ist die ASCII-Codierung vom Anfang des Johannesevangeliums:

“Im Anfang war das Wort”

oder

[illegible]



Bayes'sche Entscheidungstheorie

5. Mai 2017

Was brauchen wir, um uns zwischen zwei Hypothesen entscheiden zu können, wenn wir keine Daten haben?

- A-priori-Wahrscheinlichkeit (**Prior**)
- Wenn nicht mehr bekannt, lautet die **Entscheidungsregel**:

$$\omega_1, \text{ wenn } P(\omega_1) > P(\omega_2),$$

$$\omega_2 \text{ sonst.}$$



5. Mai 2017

Wenn x bekannt ist, brauchen wir zusätzlich die Likelihood-Funktionen $p(x|\omega_1)$ und $p(x|\omega_2)$

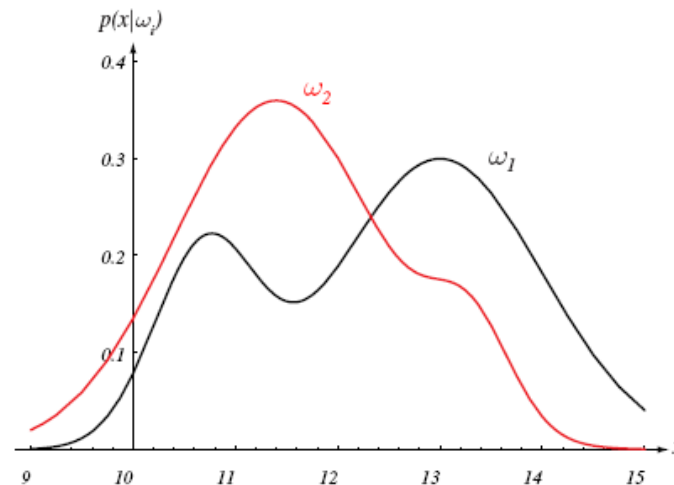


FIGURE 2.1. Hypothetical class-conditional probability density functions show the probability density of measuring a particular feature value x given the pattern is in category ω_i . If x represents the lightness of a fish, the two curves might describe the difference in lightness of populations of two types of fish. Density functions are normalized, and thus the area under each curve is 1.0. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



A-posteriori-Wahrscheinlichkeit (**Posterior**)

$$P(\omega_j | x) = \frac{p(x | \omega_j) \times P(\omega_j)}{p(x)}$$

Dabei ist **Evidence** (Wahrscheinlichkeit für Daten)
Skalierungsfaktor

$$p(x) = \sum_j p(x | \omega_j) \times P(\omega_j)$$



5. Mai 2017

A-posteriori Wahrscheinlichkeiten:

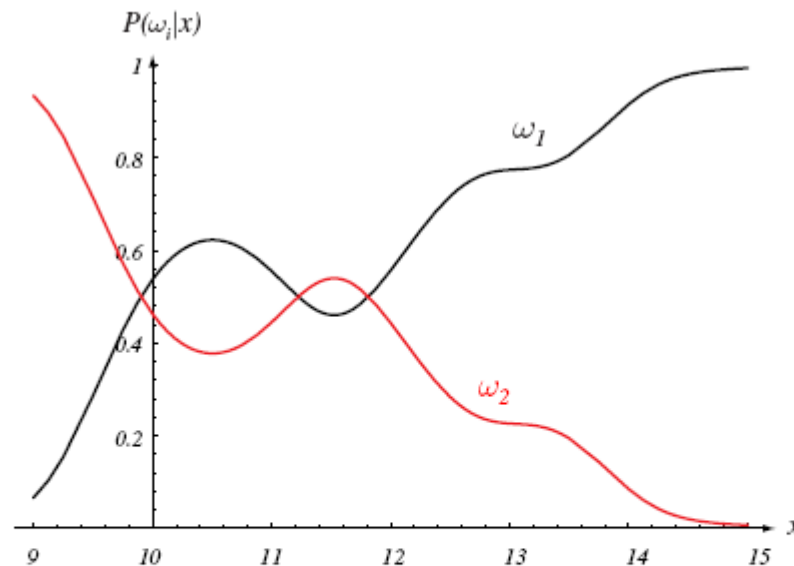


FIGURE 2.2. Posterior probabilities for the particular priors $P(\omega_1) = 2/3$ and $P(\omega_2) = 1/3$ for the class-conditional probability densities shown in Fig. 2.1. Thus in this case, given that a pattern is measured to have feature value $x = 14$, the probability it is in category ω_2 is roughly 0.08, and that it is in ω_1 is 0.92. At every x , the posteriors sum to 1.0. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



Eine mögliche Entscheidungsregel:

minimiere Wahrscheinlichkeit für Fehler:

$$P(\text{error}) = \int_{-\infty}^{\infty} p(\text{error}, x) dx = \int_{-\infty}^{\infty} P(\text{error} | x) p(x) dx$$

$$P(\text{error} | x) = \begin{cases} P(\omega_1 | x) & \text{wenn Entscheidung für } \omega_2 \\ P(\omega_2 | x) & \text{wenn Entscheidung für } \omega_1 \end{cases}$$

Als Entscheidungsregel erhalten wir somit:

$$\begin{array}{l} \omega_1, \text{ wenn } P(\omega_1 | x) > P(\omega_2 | x), \\ \omega_2 \text{ sonst.} \end{array}$$



5. Mai 2017

Für diese Regel ist die Evidence unerheblich, die Division durch $p(x)$ bewirkt nur, daß

$$P(\omega_1 | x) + P(\omega_2 | x) = 1$$



Erweiterung durch generalisierte **Loss** Funktion:

$$\lambda(\alpha_i | \omega_j)$$

sei der Verlust, der mit der Entscheidung α_i verbunden ist, wenn der wirkliche Zustand ω_j ist.

Der Erwartungswert des Verlustes der mit der Entscheidung α_i verbunden ist, wird Risiko (**Risk**) genannt und ist

$$R(\alpha_i | \mathbf{x}) = \sum_j \lambda(\alpha_i | \omega_j) P(\omega_j | \mathbf{x})$$



5. Mai 2017

Entscheidungsregel – Die **Entscheidungsfunktion**

$$\alpha(\mathbf{x})$$

gibt an, für welchen Zustand aus $\alpha_1 \dots \alpha_n$ wir uns bei Vorliegen von \mathbf{x} entscheiden.

Gesucht ist die Funktion, die das Gesamtrisiko

$$R = \int R(\alpha(\mathbf{x}) | \mathbf{x}) p(\mathbf{x}) d\mathbf{x}$$

minimiert.



Die Lösung lautet:

$$\alpha(\mathbf{x}) = \arg \min \left(R(\alpha | \mathbf{x}) \right) = \arg \min \left(\sum_j \lambda(\alpha | \omega_j) P(\omega_j | \mathbf{x}) \right)$$

Das resultierende Gesamtrisiko R^* wird **Bayes risk** genannt und ist das kleinste Gesamtrisiko, das erreicht werden kann (optimale Entscheidungsregel)



Bayes'sche Entscheidungstheorie

5. Mai 2017

Spezialfall: Zwei Zustände ω_1 und ω_2 .

Wir kürzen ab: $\lambda_{ij} = \lambda(\alpha_i | \omega_j)$

Damit erhalten wir:

$$R(\alpha_1 | \mathbf{x}) = \lambda_{11}P(\omega_1 | \mathbf{x}) + \lambda_{12}P(\omega_2 | \mathbf{x})$$

$$R(\alpha_2 | \mathbf{x}) = \lambda_{21}P(\omega_1 | \mathbf{x}) + \lambda_{22}P(\omega_2 | \mathbf{x})$$

Die Entscheidungsregel lautet damit:

$$\alpha(\mathbf{x}) = \begin{cases} \alpha_1, & \text{wenn } (\lambda_{21} - \lambda_{11})P(\omega_1 | \mathbf{x}) > (\lambda_{12} - \lambda_{22})P(\omega_2 | \mathbf{x}) \\ \alpha_2 & \text{sonst} \end{cases}$$



Bayes'sche Entscheidungstheorie

5. Mai 2017

Die Ungleichung kann durch die Bayes Regel (und Multiplikation mit Evidence) umgeformt werden:

$$\begin{aligned} (\lambda_{21} - \lambda_{11}) P(\omega_1 | \mathbf{x}) &> (\lambda_{12} - \lambda_{22}) P(\omega_2 | \mathbf{x}) \\ (\lambda_{21} - \lambda_{11}) p(\mathbf{x} | \omega_1) P(\omega_1) &> (\lambda_{12} - \lambda_{22}) p(\mathbf{x} | \omega_2) P(\omega_2) \end{aligned}$$

Wenn $\lambda_{21} > \lambda_{11}$ können wir stattdessen auch schreiben:

$$\alpha(\mathbf{x}) = \alpha_1, \text{ wenn } \frac{p(\mathbf{x} | \omega_1)}{p(\mathbf{x} | \omega_2)} > \frac{(\lambda_{12} - \lambda_{22}) P(\omega_2)}{(\lambda_{21} - \lambda_{11}) P(\omega_1)}$$

Bayes'sche Entscheidungstheorie

5. Mai 2017

Das Verhältnis
$$\frac{p(\mathbf{x} | \omega_1)}{p(\mathbf{x} | \omega_2)}$$

wird **likelihood ratio** genannt und hängt nur von \mathbf{x} ab.

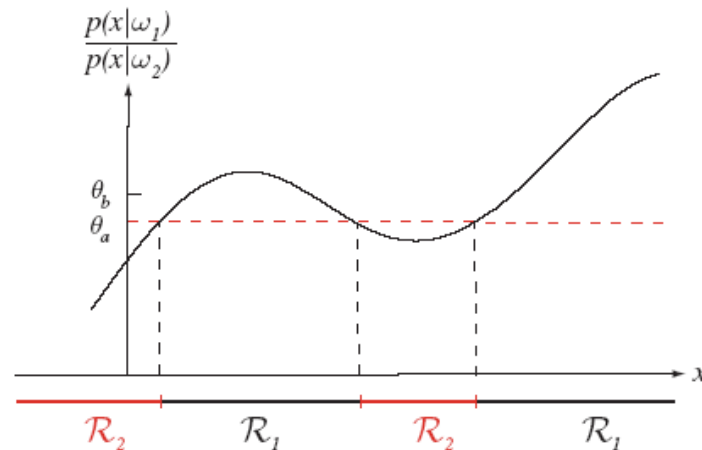


FIGURE 2.3. The likelihood ratio $p(x|\omega_1)/p(x|\omega_2)$ for the distributions shown in Fig. 2.1. If we employ a zero-one or classification loss, our decision boundaries are determined by the threshold θ_a . If our loss function penalizes miscategorizing ω_2 as ω_1 patterns more than the converse, we get the larger threshold θ_b , and hence \mathcal{R}_1 becomes smaller. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.