



---

19. Mai 2017

# Parameterschätzung - Bayes'sche Schätzung



# Parameterschätzung

---

19. Mai 2017

Was haben wir bis jetzt behandelt?

- Klassifizierung
- Repräsentation unseres Wissens (durch Wahrscheinlichkeiten und Verteilungen)

Lernen:

- Gewinnen des Wissens aus Trainingsdaten



# Parameterschätzung

---

19. Mai 2017

Relativ unproblematisch ist normalerweise die Schätzung der a-Priori-Wahrscheinlichkeiten aus Beispieldaten

Schwieriger ist die Bestimmung der Likelihoodfunktionen  $p(\mathbf{x} | \omega_i)$ .

Einfacher wird das Problem, wenn nur Parameter einer parametrisierten Funktion bestimmt werden müssen, z.B. bei Normalverteilung  $\boldsymbol{\mu}_i$  und  $\mathbf{C}_i$ .



# Parameterschätzung

---

19. Mai 2017

Methoden für diese Parameterschätzung:

- Maximum likelihood
- Bayes'sche Schätzung

Was ist der Unterschied?

- Maximum likelihood sucht die Parameter, die die Trainingsdaten am besten erklären
- Bayes'sche Schätzung liefert nicht einen Wert, sondern Wahrscheinlichkeitsverteilung für Parameter.



# Parameterschätzung

19. Mai 2017

## Maximum likelihood

Wir bestimmen likelihood für die Trainingsdaten  $\mathcal{D}$ :

$$P(\mathcal{D} | \boldsymbol{\theta}) = \prod_{k=1}^n p(\mathbf{x}_k | \boldsymbol{\theta})$$

und suchen den Parametervektor  $\hat{\boldsymbol{\theta}}$ , der  $P(\mathcal{D} | \boldsymbol{\theta})$  maximiert.

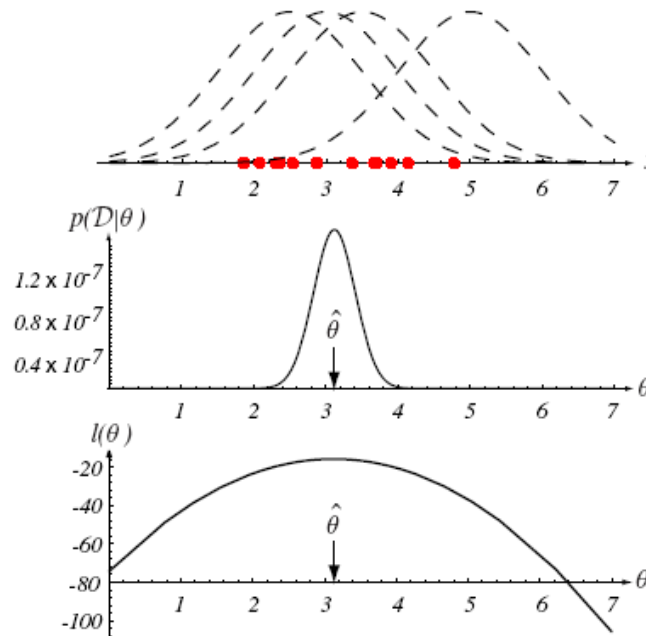
Einfacher ist es oft, die **log-likelihood** zu maximieren

$$l(\boldsymbol{\theta}) = \ln P(\mathcal{D} | \boldsymbol{\theta}) = \sum_{k=1}^n \ln p(\mathbf{x}_k | \boldsymbol{\theta})$$



19. Mai 2017

## maximum likelihood



**FIGURE 3.1.** The top graph shows several training points in one dimension, known or assumed to be drawn from a Gaussian of a particular variance, but unknown mean. Four of the infinite number of candidate source distributions are shown in dashed lines. The middle figure shows the likelihood  $p(D|\theta)$  as a function of the mean. If we had a very large number of training points, this likelihood would be very narrow. The value that maximizes the likelihood is marked  $\hat{\theta}$ ; it also maximizes the logarithm of the likelihood—that is, the log-likelihood  $l(\theta)$ , shown at the bottom. Note that even though they look similar, the likelihood  $p(D|\theta)$  is shown as a function of  $\theta$  whereas the conditional density  $p(x|\theta)$  is shown as a function of  $x$ . Furthermore, as a function of  $\theta$ , the likelihood  $p(D|\theta)$  is not a probability density function and its area has no significance. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



# Parameterschätzung

19. Mai 2017

Die Parameter  $\hat{\boldsymbol{\theta}}$ , die die maximum likelihood Bedingung erfüllen, sind also:

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} l(\boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} \sum_{k=1}^n \ln p(\mathbf{x}_k | \boldsymbol{\theta})$$

Sie müssen folgende Bedingung erfüllen:

$$\nabla l(\boldsymbol{\theta}) = \frac{\partial l(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \begin{pmatrix} \frac{\partial l(\boldsymbol{\theta})}{\partial \theta_1} \\ \vdots \\ \frac{\partial l(\boldsymbol{\theta})}{\partial \theta_p} \end{pmatrix} = \sum_{k=1}^n \nabla \ln p(\mathbf{x}_k | \boldsymbol{\theta}) = \mathbf{0}$$



# Maximum Likelihood

19. Mai 2017

Beispiel: Gaußverteilung  $p(\mathbf{x}_k | \boldsymbol{\theta}) = N(\boldsymbol{\mu}, \mathbf{C})(\mathbf{x}_k)$

Eindimensionale log-likelihood für einzelnen Datenpunkt:

$$\ln p(x_k | \boldsymbol{\theta}) = -\frac{1}{2} \ln 2\pi\sigma^2 - \frac{1}{2\sigma^2} (x_k - \mu)^2$$

Die Ableitungen nach  $\mu$  und  $\sigma^2$ :

$$\nabla \ln p(x_k | \boldsymbol{\theta}) = \begin{pmatrix} \frac{1}{\sigma^2} (x_k - \mu) \\ -\frac{1}{2\sigma^2} + \frac{(x_k - \mu)^2}{2\sigma^4} \end{pmatrix}$$





# Maximum Likelihood

19. Mai 2017

Für die gesamte log-likelihood (alle Daten) bekommen wir:

$$\sum_{k=1}^n \frac{1}{\sigma^2} (x_k - \mu) = 0 \quad \Rightarrow \quad \hat{\mu} = \frac{1}{n} \sum_{k=1}^n x_k$$

$$-\sum_{k=1}^n \frac{1}{\sigma^2} + \sum_{k=1}^n \frac{(x_k - \mu)^2}{\sigma^4} = 0 \quad \Rightarrow \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{k=1}^n (x_k - \mu)^2$$



# Maximum Likelihood

19. Mai 2017

Die Schätzung für die Varianz benötigt den korrekten Wert von  $\mu$ , setzten wir stattdessen  $\hat{\mu}$  ein, so erhalten wir ein zu kleines Ergebnis. Man kann zeigen, daß die korrekte Schätzung unter Verwendung von  $\hat{\mu}$  lautet:

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_{k=1}^n (x_k - \hat{\mu})^2$$



# Bayes'sche Schätzung

19. Mai 2017

Bei der Bayes'schen Schätzung betrachtet man statt der wahrscheinlichsten Parameterkombination  $\hat{\theta}$  eine Plausibilitätsverteilung der Parameter  $p(\theta | \mathcal{D})$

Ziel: Berechnung von  $P(\omega_i | \mathbf{x}, \mathcal{D})$

Bayes Formel:

$$P(\omega_i | \mathbf{x}, \mathcal{D}) = \frac{p(\mathbf{x} | \omega_i, \mathcal{D}) P(\omega_i | \mathcal{D})}{\sum_{j=1}^c p(\mathbf{x} | \omega_j, \mathcal{D}) P(\omega_j | \mathcal{D})}$$



# Bayes'sche Schätzung

---

19. Mai 2017

Von der gesuchten klassenabhängigen Dichte

$p(\mathbf{x} | \omega_i, \mathcal{D})$  nehmen wir an, dass sie nur von den Daten der jeweiligen Klasse abhängt und betrachten daher jeweils nur eine Klasse, können also  $\omega_i$  weglassen.

Unser (Zwischen-)Ziel ist also die Bestimmung von  $p(\mathbf{x} | \mathcal{D})$ , das als Approximation für  $p(\mathbf{x})$  dienen soll.



# Bayes'sche Schätzung

19. Mai 2017

Wir setzen voraus, daß die wirkliche Verteilung parametrisch geschrieben werden kann, die Form der Funktion  $p(\mathbf{x} | \boldsymbol{\theta})$  ist also bekannt.

Bei der Bayes'schen Schätzung verwenden wir für die gesuchte Verteilung  $p(\mathbf{x} | \mathcal{D})$  die Formel:

$$\begin{aligned} p(\mathbf{x} | \mathcal{D}) &= \int p(\mathbf{x}, \boldsymbol{\theta} | \mathcal{D}) d\boldsymbol{\theta} \\ &= \int p(\mathbf{x} | \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathcal{D}) d\boldsymbol{\theta} \end{aligned}$$



# Bayes'sche Schätzung

19. Mai 2017

Diese Verteilung  $p(\mathbf{x} | \mathcal{D}) = \int p(\mathbf{x} | \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathcal{D}) d\boldsymbol{\theta}$   
hat also nicht unbedingt die parametrische  
Form  $p(\mathbf{x} | \boldsymbol{\theta})$ , von der wir eigentlich wissen, daß sie  
stimmt.

Im Gegensatz dazu verwendeten wir bei maximum  
likelihood:

$$p(\mathbf{x} | \hat{\boldsymbol{\theta}})$$

mit

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} P(\mathcal{D} | \boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} l(\boldsymbol{\theta})$$



# Bayes'sche Schätzung

---

19. Mai 2017

Wie rechnen wir nun

$$p(\mathbf{x} | \mathcal{D}) = \int p(\mathbf{x} | \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathcal{D}) d\boldsymbol{\theta}$$

aus?

Wir brauchen:  $p(\boldsymbol{\theta} | \mathcal{D})$

und verwenden dafür wieder die Bayes'sche Formel:

$$p(\boldsymbol{\theta} | \mathcal{D}) = \frac{p(\mathcal{D} | \boldsymbol{\theta}) p(\boldsymbol{\theta})}{\int p(\mathcal{D} | \boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta}}$$



# Bayes'sche Schätzung

19. Mai 2017

Beispiel: Normalverteilung  $p(x | \mu) = N(\mu, \sigma^2)$  ; dabei ist  $\sigma$  bekannt,  $\mu$  ist unbekannt.

$$p(\mu | \mathcal{D}) = \frac{p(\mathcal{D} | \mu) p(\mu)}{\int p(\mathcal{D} | \mu) p(\mu) d\mu}$$

Irgendwie müssen wir den Prior festlegen, und tun das z.B. so:  $p(\mu) = N(\mu_0, \sigma_0^2)$

Wir erhalten:  $p(\mu | \mathcal{D}) = \alpha \prod_{k=1}^n p(x_k | \mu) p(\mu)$





# Bayes'sche Schätzung

19. Mai 2017

Als Ergebnis kommt raus:

$$p(\mu | \mathcal{D}) = N(\mu_n, \sigma_n^2)$$

mit

$$\mu_n = \frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \hat{\mu}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \mu_0$$

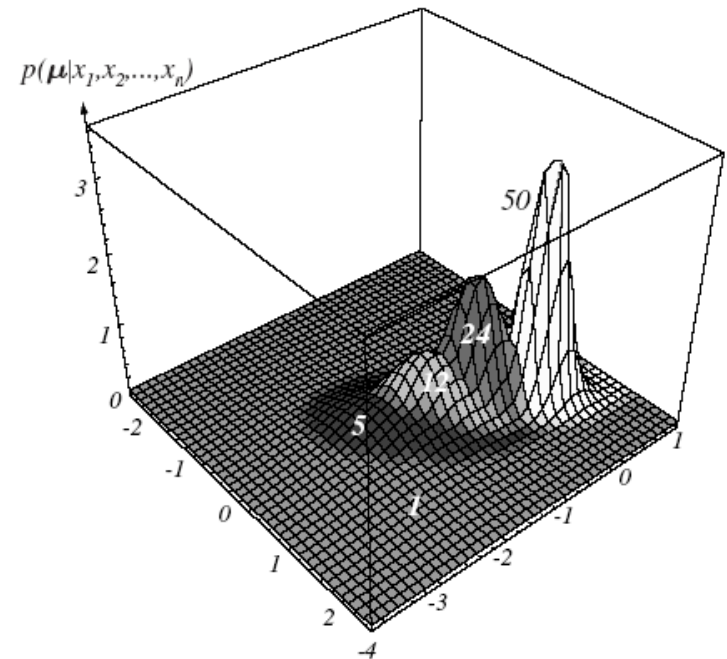
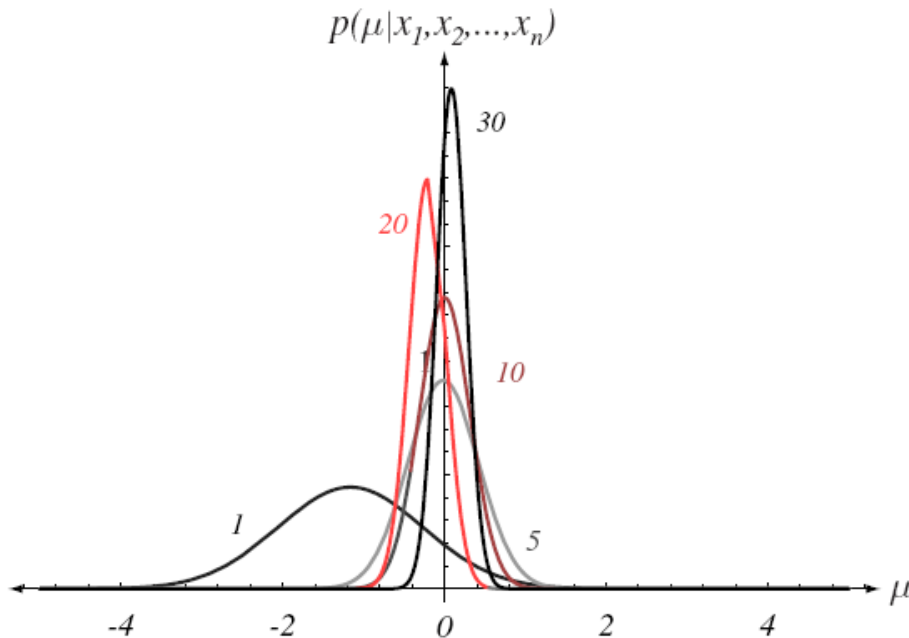
und

$$\sigma_n^2 = \frac{\sigma_0^2 \sigma^2}{n\sigma_0^2 + \sigma^2}$$



19. Mai 2017

Als Ergebnis kommt raus:  $p(\mu | \mathcal{D}) = N(\mu_n, \sigma_n^2)$



**FIGURE 3.2.** Bayesian learning of the mean of normal distributions in one and two dimensions. The posterior distribution estimates are labeled by the number of training samples used in the estimation. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



# Bayes'sches Lernen

---

19. Mai 2017

Um klassifizieren zu können müssen wir nur  
noch  $p(x | \mathcal{D})$  berechnen:

$$p(x | \mathcal{D}) = \int p(x | \mu) p(\mu | \mathcal{D}) d\mu$$

und erhalten:  $p(x | \mathcal{D}) = N(\mu_n, \sigma^2 + \sigma_n^2)$



# Bayes'sches Lernen

19. Mai 2017

Allgemein: Wie konvergiert  $p(\mathbf{x} | \mathcal{D})$  gegen  $p(\mathbf{x})$ ?

Nehmen wir unabhängige Samples an, so erhalten wir

$$p(\mathcal{D}^n | \boldsymbol{\theta}) = p(\mathbf{x}^n | \boldsymbol{\theta}) p(\mathcal{D}^{n-1} | \boldsymbol{\theta})$$

Eingesetzt in Bayes Formel:

$$p(\boldsymbol{\theta} | \mathcal{D}^n) = \frac{p(\mathbf{x}^n | \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathcal{D}^{n-1})}{\int p(\mathbf{x}^n | \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathcal{D}^{n-1}) d\boldsymbol{\theta}}$$



# Bayes'sches Lernen

19. Mai 2017

Wir fangen also mit keinen Daten an und erhalten unseren Prior:

$$p(\boldsymbol{\theta}) = p(\boldsymbol{\theta} | \mathcal{D}^0)$$

und bekommen dann nacheinander:

$$p(\boldsymbol{\theta} | \mathcal{D}^1) = p(\boldsymbol{\theta} | \mathbf{x}_1)$$

$$p(\boldsymbol{\theta} | \mathcal{D}^2) = p(\boldsymbol{\theta} | \mathbf{x}_1, \mathbf{x}_2)$$

$$p(\boldsymbol{\theta} | \mathcal{D}^3) = p(\boldsymbol{\theta} | \mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3)$$

$\vdots$

Man nennt diese Art der Schätzung **rekursive Bayes-Schätzung**, ein Beispiel für ein **inkrementelles** oder **On-Line-Lernverfahren**



# Bayes'sches Lernen

19. Mai 2017

## Beispiel: Uniforme Verteilung

$$p(x | \theta) = U(0, \theta) = \begin{cases} 1/\theta & 0 \leq x \leq \theta \\ 0 & \text{sonst} \end{cases}$$

Als Prior nehmen wir an:  $p(\theta) = U(0, 10)$

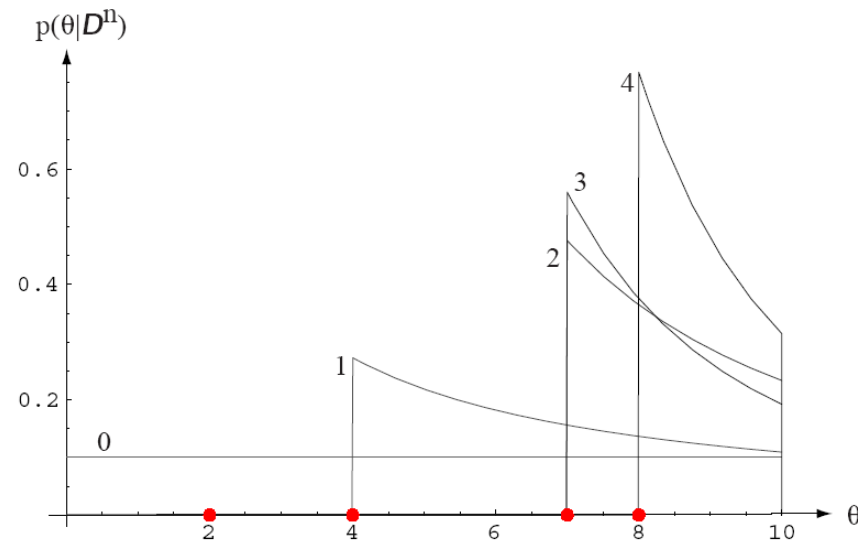
Bei einem Datenpunkt  $x_1 = 4$  erhalten wir:

$$p(\theta | \mathcal{D}^1) \propto p(x | \theta) p(\theta | \mathcal{D}^0) = \begin{cases} 1/\theta & 4 \leq \theta \leq 10 \\ 0 & \text{sonst} \end{cases}$$



Der zweite Datenpunkt  $x_2 = 7$  liefert:

$$p(\theta | \mathcal{D}^2) \propto p(x | \theta) p(\theta | \mathcal{D}^1) = \begin{cases} 1/\theta^2 & 7 \leq \theta \leq 10 \\ 0 & \text{sonst} \end{cases}$$





Aus der Verteilung für  $\theta$  können wir  $p(x | \mathcal{D}^2)$  berechnen:

$$\begin{aligned} p(x | \mathcal{D}^2) &= \int_0^{10} p(x | \theta) p(\theta | \mathcal{D}^2) d\theta \\ &\propto \int_7^{10} U(0, \theta) \frac{1}{\theta^2} d\theta \end{aligned}$$





# Bayes'sches Lernen

---

19. Mai 2017

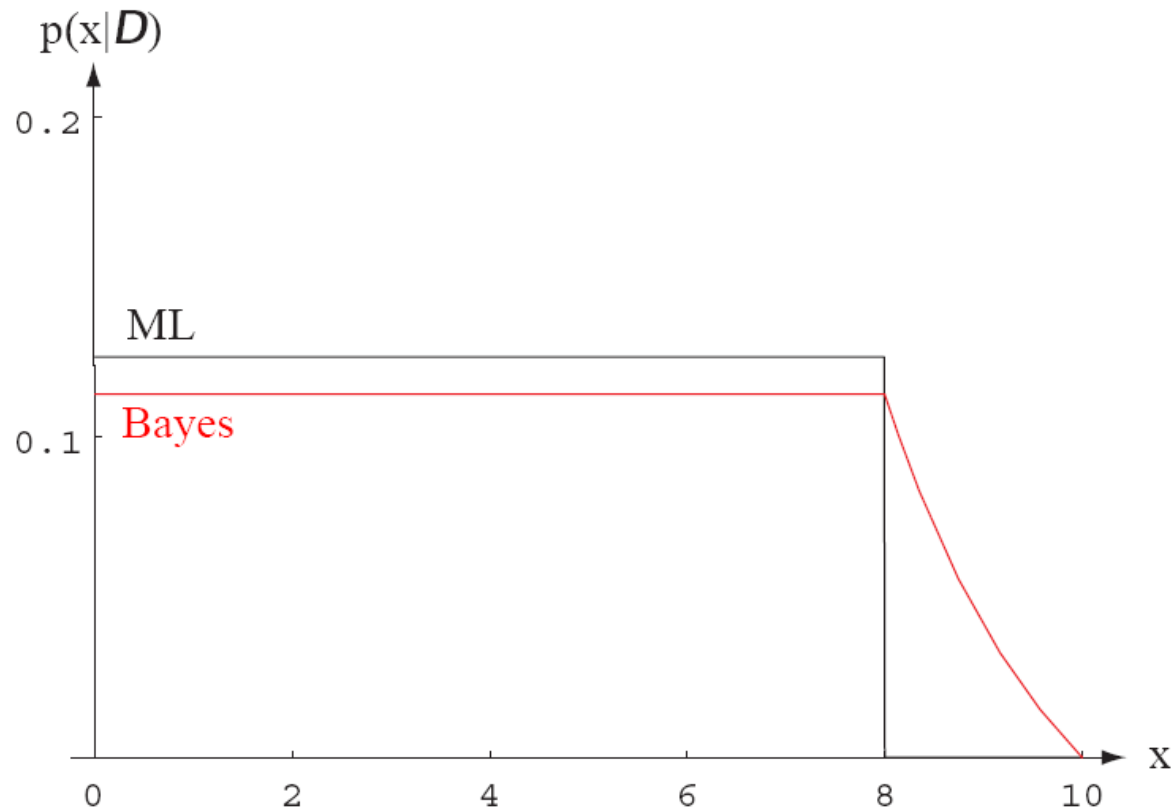
Frage: Was wäre Max Likelihood?

Antwort:  $\hat{\theta} = 7$

und damit:  $p(x | \hat{\theta}) = (x | 7) = U(0, 7) = \begin{cases} 1/7 & 0 \leq x \leq 7 \\ 0 & \text{sonst} \end{cases}$



Vergleich Maximum Likelihood und Bayes nach der Berücksichtigung aller 4 Datenpunkte:





# Parameterschätzung

---

19. Mai 2017

Grundsätzliches:

- $p(\mathbf{x} | \boldsymbol{\theta})$  wird **identifizierbar** genannt, wenn  $\boldsymbol{\theta}$  eindeutig aus Verteilung bestimmt werden kann.
- Wenn mehrere  $\boldsymbol{\theta}$  in Frage kommen, erhält man asymptotisch keine  $\delta$ -Funktion, sondern mehrere Peaks.
- $p(\mathbf{x})$  ist aber (nach Integration über  $\boldsymbol{\theta}$ ) trotzdem eindeutig.



# Parameterschätzung

---

19. Mai 2017

Grundsätzliches:

Fehlerquellen bei der Klassifizierung

- Bayes Fehler (minimal möglicher Fehler)
- Modell-Fehler: Fehler durch falsche parametrische Form von  $p(\mathbf{x})$
- Schätzfehler: Falsches Modell durch zu wenige Trainingsdaten
  - Fehler = 0 für unendlich viele Trainingsdaten
  - ML und Bayes asymptotisch äquivalent



# Parameterschätzung

---

19. Mai 2017

Grundsätzliches:

Wahl der A-priori-Verteilung

- Nichtinformativer Prior  $p(\boldsymbol{\theta})$
- Kriterium: Invarianz



# Parameterschätzung

19. Mai 2017

Dimension:

Was bringt die Hinzunahme weiterer Merkmale?

Beispiel: Zwei Klassen mit gleicher Kovarianz

- Fehlerrate ist: 
$$P(e) = \int_{r/2}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{u^2}{2}} du$$

mit 
$$r^2 = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \mathbf{C}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$



# Parameterschätzung

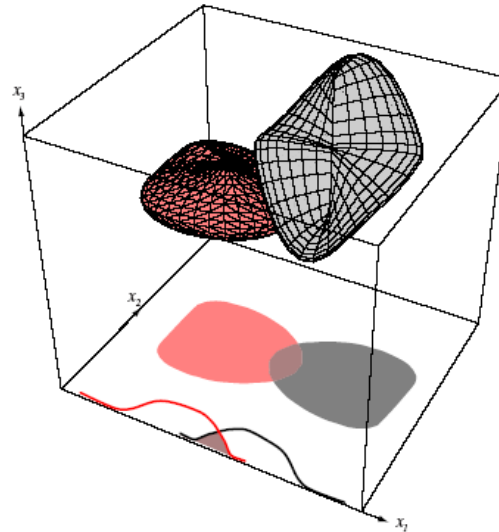
19. Mai 2017

Für Eigenwerte von  $\mathbf{C}$ , die nicht 0 sind, kann also der Fehler verbessert werden, wenn entsprechendes Merkmal dazugenommen werden. Im Grenzfall kann der Fehler sogar beliebig reduziert werden.

$$P(e) = \int_{r/2}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{u^2}{2}} du$$

Mahalanobis Distanz  $r^2 = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \mathbf{C}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$

Beispiel:



**FIGURE 3.3.** Two three-dimensional distributions have nonoverlapping densities, and thus in three dimensions the Bayes error vanishes. When projected to a subspace—here, the two-dimensional  $x_1 - x_2$  subspace or a one-dimensional  $x_1$  subspace—there can be greater overlap of the projected distributions, and hence greater Bayes error. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

In der Praxis leider manchmal Vergrößerung des Fehlers durch falsches Modell oder zu wenige Trainingsdaten