Aufgabenstellung der Parameterschätzung

Gegeben:

- Mathematisches Modell einer WV $p(\mathbf{m}|\boldsymbol{\theta})$ für eine beobachtbare Größe \mathbf{m}
- Stichprobe $D = \{\mathbf{m}_1, ..., \mathbf{m}_N\}$ für festen, aber unbekannten Parametervektor $\boldsymbol{\theta}$.

Gesucht:

"Wahrer" Wert des Parametervektors θ.

Ansatz:

Schätzfunktion θ̂(D) (Schätzer), welche die Stichprobe D statistisch auswertet und in gewissem Sinne dem "wahren" Wert θ möglichst nahe kommt.

Methodiken der Parameterschätzung

Likelihoodmethodik (Klassische Statistik)

- wird als unbekannte, konstante, nicht-stochastische Größe angesehen.
- Maximum-Likelihood-Schätzung: Man wählt $\hat{\theta}(D)$ so, dass die gegebenen Beobachtungen $D = \{\mathbf{m}_1, ..., \mathbf{m}_N\}$ maximal wahrscheinlich werden.

Bayes'sche Methodik (Bayes'sche Statistik)

- Wird als Zufallsvariable angesehen und über eine WV beschrieben.
- Z.B.: Maximum-A Posteriori-Schätzung: Man wählt $\hat{\theta}(D)$ als den Wert θ mit der höchsten A Posteriori Wahrscheinlichkeit nach Beobachtung von D.

4. Parameterschätzung – Einschub: Bedeutung von Wahrscheinlichkeit

Axiome nach Kolmogorov:

A1: Nichtnegativität $Pr(A) \ge 0$

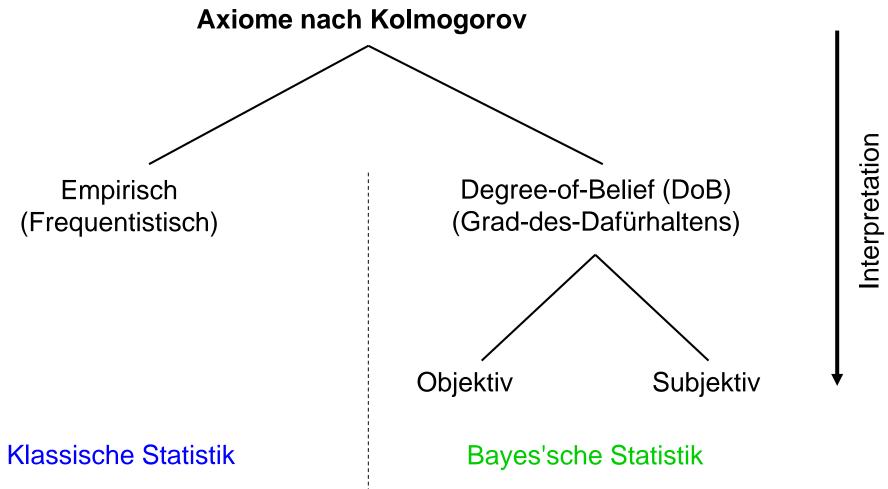
A2: Normierung Pr(Sicheres Ereignis) = 1

A3: Additivität $A \cap B = \emptyset \implies \Pr(A \cup B) = \Pr(A) + \Pr(B)$

Bedingte Wahrscheinlichkeit: $Pr(A|B) := \frac{Pr(A \cap B)}{Pr(B)}$

Die Kolmogorovschen Axiome legen fest, wie mit Wahrscheinlichkeiten gerechnet wird (Syntax), machen aber keine Aussage darüber, was Wahrscheinlichkeiten bedeuten (Semantik).

Mit den Axiomen verträgliche Interpretationen bezüglich der Bedeutung von Wahrscheinlichkeit:



Parameterschätzung im Kontext der Mustererkennung

■ Schätzung der Parameter der klassenspezifischen WVen $p(\mathbf{m}|\mathbf{\theta}_i, \omega_i)$ i = 1,...,c aus den Daten D.

Annahmen:

- Struktur der klassenspezifischen WVen ist gegeben. Nur die Parameter θ_i sind unbekannt.
- Überwachtes Lernen: die Klassenzugehörigkeiten der einzelnen Stichproben in D sind bekannt. → Partitionierung von D in $D_1,...,D_c$ mit den zur jeweiligen Klasse gehörenden Stichproben.
- Die Stichproben in D_i haben keinen Einfluss auf $p(\mathbf{m}|\mathbf{\theta}_i, \omega_i)$ $i \neq j$.

Konsequenz:

Schätzung der Parameter der klassenspezifischen WVen zerfällt in c separate Parameterschätzaufgaben, die unabhängig voneinander behandelt werden können. → Klassenunterscheidende Kennzeichnung kann in diesem Kapitel weitgehend unterdrückt werden.

Eigenschaften von Schätzern: Erwartungtreue

Klassische Statistik

Erwartungstreue (unbiased estimator):

$$E\{\hat{\boldsymbol{\theta}}\} = \boldsymbol{\theta}$$

$$E\{\hat{\boldsymbol{\theta}}\} = \int_{\Theta} \hat{\boldsymbol{\theta}} p(\hat{\boldsymbol{\theta}}) d\hat{\boldsymbol{\theta}} = \int_{\mathbb{M}} \hat{\boldsymbol{\theta}}(\mathbf{m}) p(\mathbf{m} \mid \boldsymbol{\theta}) d\mathbf{m}$$

Bias (systematischer Fehler):

$$b(\mathbf{\theta}) = \mathbf{E}\{\hat{\mathbf{\theta}}\} - \mathbf{\theta}$$

Bemerkung: Simple Korrektur durch Subtraktion nur möglich, falls b nicht vom Parameter θ abhängt.

Bayes'sche Statistik

Erwartungstreue (unbiased estimator):

$$E\{\hat{\boldsymbol{\theta}}\} = E\{\boldsymbol{\theta}\}$$

$$E\{\hat{\boldsymbol{\theta}}\} = \int_{\Theta} \hat{\boldsymbol{\theta}} p(\hat{\boldsymbol{\theta}}) d\hat{\boldsymbol{\theta}} = \int_{\Theta} \int_{\mathbb{M}} \hat{\boldsymbol{\theta}}(\mathbf{m}) p(\mathbf{m}, \boldsymbol{\theta}) d\mathbf{m} d\boldsymbol{\theta}$$

Bias:

$$b = \mathbf{E}\{\hat{\mathbf{\theta}}\} - \mathbf{E}\{\mathbf{\theta}\}$$

Eigenschaften von Schätzern: Varianz

Varianz (stochastischer Fehler):
$$Var\{\hat{\theta}\} = E\{(\hat{\theta} - E\{\hat{\theta}\})^2\} = E\{\hat{\theta}^2\} - E\{\hat{\theta}\}^2$$

Schätzverfahren mit minimaler Varianz werden als effiziente (wirksame) Schätzverfahren bezeichnet.

Cramer-Rao-Schranke (CRB) für erwartungstreue Schätzer (klassische Statistik) für einen skalaren Parameter:

 $D = \{m_1, ..., m_N\}$, Beobachtungen m_i seien stochastisch unabhängig.

$$|\operatorname{Var}\{\hat{\theta}(\mathbf{D})\}| \ge \frac{1}{N \operatorname{E}\left\{\left(\frac{\partial \ln p(\mathbf{m}|\theta)}{\partial \theta}\right)^{2}\right\}} = \frac{1}{N \int_{\mathbb{IM}} \left(\frac{\partial \ln p(\mathbf{m}|\theta)}{\partial \theta}\right)^{2} p(\mathbf{m}|\theta) d\mathbf{m}}$$

Details und Verallgemeinerungen siehe z.B.: K. Kroschel, "Statistische Informationstechnik", Springer 2004

- CRB ist eine untere Schranke für die Varianz aller erwartungstreuen Schätzer.
- Effiziente Schätzverfahren erreichen die CRB.

Beweisskizze CRB: $\mathbf{m} \in \mathbb{R}^d$, $\theta \in \mathbb{R}$, N = 1

Voraussetzungen:
$$\frac{\partial}{\partial \theta} \int p(\mathbf{m} \mid \theta) d\mathbf{m} = \int \frac{\partial p(\mathbf{m} \mid \theta)}{\partial \theta} d\mathbf{m}$$

 $\hat{\theta}(\mathbf{m})$ ist erwartungstreu: $\mathbf{E} \{ \hat{\theta}(\mathbf{m}) \} = \theta$

$$\frac{\partial}{\partial \theta} \int (\hat{\theta}(\mathbf{m}) - \theta) p(\mathbf{m} \mid \theta) d\mathbf{m} = \int \frac{\partial}{\partial \theta} (\hat{\theta}(\mathbf{m}) - \theta) p(\mathbf{m} \mid \theta) d\mathbf{m} = 0$$

mit
$$\frac{\partial \hat{\theta}(\mathbf{m})}{\partial \theta} = 0$$
 folgt:

$$-\int p(\mathbf{m} \mid \theta) d\mathbf{m} + \int (\hat{\theta}(\mathbf{m}) - \theta) \frac{\partial p(\mathbf{m} \mid \theta)}{\partial \theta} d\mathbf{m} = 0$$

mit
$$\frac{\partial p(\mathbf{m}|\theta)}{\partial \theta} = \frac{\partial \ln p(\mathbf{m}|\theta)}{\partial \theta} p(\mathbf{m}|\theta)$$
 folgt:

$$\int \left(\hat{\theta}(\mathbf{m}) - \theta\right) p(\mathbf{m} \mid \theta) \frac{\partial \ln p(\mathbf{m} \mid \theta)}{\partial \theta} d\mathbf{m} = 1$$

$$\int \left(\hat{\theta}(\mathbf{m}) - \theta\right) p(\mathbf{m} \mid \theta) \frac{\partial \ln p(\mathbf{m} \mid \theta)}{\partial \theta} d\mathbf{m} = 1$$

Mit der Schwarzschen Ungleichung: $\int x^2(\mathbf{m}) d\mathbf{m} \int y^2(\mathbf{m}) d\mathbf{m} \ge \left(\int x(\mathbf{m}) y(\mathbf{m}) d\mathbf{m}\right)^2$ und $\left(\int (\hat{\theta}(\mathbf{m}) - \theta) \sqrt{p(\mathbf{m} \mid \theta)} \sqrt{p(\mathbf{m} \mid \theta)} \frac{\partial \ln p(\mathbf{m} \mid \theta)}{\partial \theta} d\mathbf{m}\right)^2 = 1$ folgt:

$$\underbrace{\int \left(\hat{\theta}(\mathbf{m}) - \theta\right)^{2} p(\mathbf{m} \mid \theta) d\mathbf{m} \int \left(\underbrace{\frac{\partial \ln p(\mathbf{m} \mid \theta)}{\partial \theta}\right)^{2} p(\mathbf{m} \mid \theta) d\mathbf{m}}_{\mathbf{Var} \left\{\hat{\theta}(\mathbf{m})\right\}} = \mathbf{E} \left\{ \underbrace{\left(\underbrace{\frac{\partial \ln p(\mathbf{m} \mid \theta)}{\partial \theta}\right)^{2}}_{\mathbf{E} \left\{\hat{\theta}(\mathbf{m})\right\}} : \text{ Fisher-Information} \right\}$$

→ Cramer-Rao-Schranke

Plausibilität: Wie stark ändert sich \mathbf{m} mit θ ?

$$\mathbf{m} \sim p(\mathbf{m} \mid \theta) \xrightarrow{\frac{\partial}{\partial \theta}} \frac{\frac{\partial p(\mathbf{m} \mid \theta)}{\partial \theta}}{\frac{\partial}{\partial \theta}} \xrightarrow{\frac{\partial}{\partial \theta}} \frac{1}{\frac{\partial}{\partial \theta}} = \frac{\partial \ln p(\mathbf{m} \mid \theta)}{\partial \theta}$$

$$\xrightarrow{()^{2}} \left(\frac{\partial \ln p(\mathbf{m} \mid \theta)}{\partial \theta}\right)^{2} \xrightarrow{E\{\}} E\{\left(\frac{\partial \ln p(\mathbf{m} \mid \theta)}{\partial \theta}\right)^{2}\} = J(\theta)$$

Quantifizierung der Empfindlichkeit von \mathbf{m} bezüglich Veränderungen von θ .

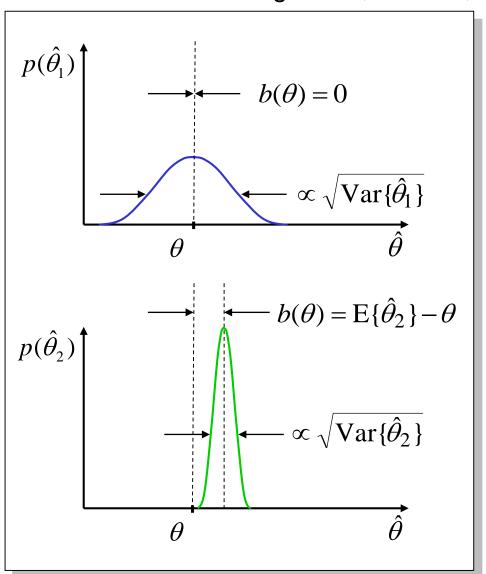
Eigenschaften von Schätzern: Konsistenz

Ein Schätzverfahren heißt konsistent, falls mit wachsender Zahl N der zur Schätzung herangezogenen Beobachtungen die Wahrscheinlichkeit, mit der die Schätzung vom wahren Wert um ein beliebiges ε abweicht, gegen Null konvergiert:

$$\lim_{N \to \infty} P\{ \| \hat{\mathbf{\theta}}(\mathbf{D}) - \mathbf{\theta} \| \ge \varepsilon \} = 0, \ \forall \varepsilon > 0$$

$$D = \{\mathbf{m}_1, ..., \mathbf{m}_N\}$$

Diskussion: Erwartungstreue, Varianz, mittlerer quadratischer Fehler



 θ : sei wahrer Wert

Schätzer $\hat{\theta}_1$: erwartungstreu

Schätzer $\hat{\theta}_2$: nicht erwartungstreu, aber:

$$\sqrt{\operatorname{Var}\{\hat{\theta}_2\}} < \sqrt{\operatorname{Var}\{\hat{\theta}_1\}}$$

Var{.}: stochastischer Fehler

Bias b: systematischer Fehler

$$E\{(\hat{\theta}-\theta)^2\} = b^2(\theta) + Var\{\hat{\theta}\}:$$

mittlerer quadratischer Fehler

Welcher Schätzer ist besser?

Beispiel: N stochastisch unabhängig gewonnene Stichproben einer Zufallsvariablen m, für die gilt: $m \sim p(m \mid \mu)$ (unabhängige identisch verteilte Stichprobe: i.i.d.). Es gelte:

$$\mu = \mathbf{E}\{m\}$$

$$\sigma^2 = \text{Var}\{m\} = \text{E}\{(m-\mu)^2\}$$

Erwartungswert μ soll durch arithmetischen Mittelwert geschätzt werden.

$$\hat{\mu} = \frac{1}{N} \sum_{k=1}^{N} m_k$$

$$E\{\hat{\mu}\} = E\left\{\frac{1}{N}\sum_{k=1}^{N}m_k\right\} = \frac{1}{N}\sum_{k=1}^{N}E\{m_k\} = \frac{1}{N}\sum_{k=1}^{N}\mu = \mu \quad \Rightarrow \text{Erwartungstreue}$$

Beispiel: Fortsetzung

Varianz des Schätzers:

$$\operatorname{Var}\{\hat{\mu}\} = \operatorname{E}\left\{(\hat{\mu} - \mu)^{2}\right\} = \operatorname{E}\left\{\left(\frac{1}{N}\sum_{k=1}^{N}m_{k} - \mu\right)^{2}\right\} = \frac{1}{N^{2}}\sum_{k=1}^{N}\sum_{l=1}^{N}\operatorname{E}\left\{(m_{k} - \mu)(m_{l} - \mu)\right\}$$

$$Var{\{\hat{\mu}\}} = \frac{1}{N^2} \sum_{k=1}^{N} E\{(m_k - \mu)^2\} = \frac{1}{N^2} \sum_{k=1}^{N} \sigma^2 = \frac{1}{N} \sigma^2$$

- Varianz verschwindet mit wachsendem $N \to Konsistenz$ aufgrund der Ungleichung von Tschebyscheff: $P\{|\hat{\theta} E\{\hat{\theta}\}| \ge \varepsilon\} \le Var\{\hat{\theta}\} / \varepsilon^2, \forall \varepsilon > 0$
- Abnahme der Standardabweichung des Schätzers $\sqrt{\operatorname{Var}\{\hat{\mu}\}}$ mit $1/\sqrt{N}$ ist typisch für die meisten praktisch relevanten Aufgabenstellungen.

Gegeben:

N stochastisch unabhängig gewonnene Stichproben eines Zufallsvektors \mathbf{m} , für den gilt: $\mathbf{m} \sim p(\mathbf{m} \mid \mathbf{\theta})$ (unabhängige identisch verteilte Stichprobe: i.i.d.).

Dann gilt:

$$p(\mathbf{D} \mid \mathbf{\theta}) = \prod_{k=1}^{N} p(\mathbf{m}_k \mid \mathbf{\theta}) \qquad \mathbf{D} = \{\mathbf{m}_1, ..., \mathbf{m}_N\}$$

 $p(D|\theta)$ als Funktion von θ wird als Likelihoodfunktion bezeichnet.

Maximum Likelihood Ansatz: Man schätze θ derart, dass die Wahrscheinlichkeit (Wahrscheinlichkeitsdichte) $p(D|\theta)$ der vorliegenden Beobachtungen D maximal wird.

$$\hat{\boldsymbol{\theta}}_{\mathrm{ML}} = \arg\max_{\boldsymbol{\theta}} \{ p(\mathbf{D} \mid \boldsymbol{\theta}) \}$$

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

4.1. Maximum Likelihood Schätzung

Wegen der Monotonie der Logarithmusfunktion, kann man äquivalent die logarithmierte Likelihoodfunktion maximieren; das bringt oft analytische Vorteile.

$$l(\mathbf{\theta}) := \ln(p(D \mid \mathbf{\theta}))$$

Loglikelihoodfunktion

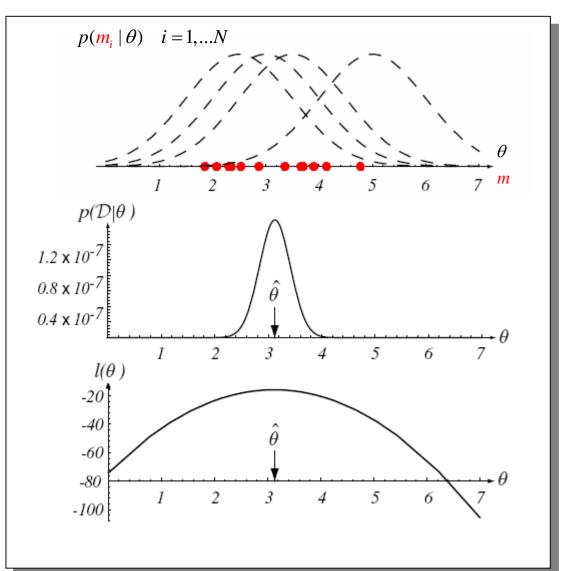
$$l(\mathbf{\theta}) = \sum_{k=1}^{N} \ln(p(\mathbf{m}_k \mid \mathbf{\theta}))$$

$$\hat{\boldsymbol{\theta}}_{\mathrm{ML}} = \arg\max_{\boldsymbol{\theta}} \{l(\boldsymbol{\theta})\}$$

Notwendig für Maximum:

$$\nabla_{\boldsymbol{\theta}} l(\boldsymbol{\theta}) = \sum_{k=1}^{N} \nabla_{\boldsymbol{\theta}} \ln(p(\mathbf{m}_{k} \mid \boldsymbol{\theta})) \stackrel{!}{=} \mathbf{0} \quad \text{mit} \quad \nabla_{\boldsymbol{\theta}} \equiv \left[\frac{\partial}{\partial \theta_{1}}, \dots, \frac{\partial}{\partial \theta_{c}} \right]^{T}$$

Beispiel: rote Punkte sind Stichproben einer normalverteilten Zufallsvariablen m mit bekannter Standardabweichung σ und unbekanntem, zu schätzendem Erwartungswert θ .



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

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

$$\ln p(\mathbf{m}_k \mid \boldsymbol{\mu}) = -\frac{1}{2} (\mathbf{m}_k - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{m}_k - \boldsymbol{\mu}) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\boldsymbol{\Sigma}|$$

$$\nabla_{\mathbf{\mu}} \ln p(\mathbf{m}_k \mid \mathbf{\mu}) = \mathbf{\Sigma}^{-1}(\mathbf{m}_k - \mathbf{\mu})$$

$$\sum_{k=1}^{N} \nabla_{\boldsymbol{\mu}} \ln p(\mathbf{m}_{k} | \boldsymbol{\mu}) = \sum_{k=1}^{N} \boldsymbol{\Sigma}^{-1} (\mathbf{m}_{k} - \hat{\boldsymbol{\mu}}) = 0$$

$$\sum_{k=1}^{N} (\mathbf{m}_k - \hat{\mathbf{\mu}}) = 0$$

$$\Rightarrow \qquad \hat{\mathbf{\mu}}_{\mathrm{ML}} = \frac{1}{N} \sum_{k=1}^{N} \mathbf{m}_{k}$$

Beispiel: $m \sim N(m; \mu, \sigma^2)$, μ, σ^2 unbekannt

$$\boldsymbol{\theta} = (\mu, \sigma^2)^{\mathrm{T}} = (\theta_1, \theta_2)^{\mathrm{T}}$$

$$\ln p(m_k \mid \mathbf{\theta}) = -\frac{1}{2\theta_2} (m_k - \theta_1)^2 - \frac{1}{2} \ln 2\pi \theta_2$$

$$\sum_{k=1}^{N} \nabla_{\mathbf{\theta}} \ln p(m_k \mid \mathbf{\theta}) = \sum_{k=1}^{N} \left[\frac{1}{\theta_2} (m_k - \theta_1), -\frac{1}{2\theta_2} + \frac{1}{2\theta_2^2} (m_k - \theta_1)^2 \right]^{\mathrm{T}} = \mathbf{0} \iff$$

$$\sum_{k=1}^{N} \frac{1}{\hat{\theta}_{2}} (m_{k} - \hat{\theta}_{1}) = 0$$

$$\sum_{k=1}^{N} \left(-\frac{1}{2\hat{\theta}_{2}} + \frac{1}{2\hat{\theta}_{2}^{2}} (m_{k} - \hat{\theta}_{1})^{2} \right) = 0$$

$$\Leftrightarrow \qquad \hat{\theta}_{1} = \hat{\mu}_{ML} = \frac{1}{N} \sum_{k=1}^{N} m_{k}$$

$$\hat{\theta}_{2} = \hat{\sigma}_{ML}^{2} = \frac{1}{N} \sum_{k=1}^{N} (m_{k} - \hat{\mu})^{2}$$

$$\hat{\theta}_1 = \hat{\mu}_{\text{ML}} = \frac{1}{N} \sum_{k=1}^{N} m_k$$

$$\hat{\theta}_2 = \hat{\sigma}_{ML}^2 = \frac{1}{N} \sum_{k=1}^{N} (m_k - \hat{\mu})^2$$

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

$$\boldsymbol{\theta} = (\boldsymbol{\mu}, \boldsymbol{\Sigma})^T$$

$$\ln p(\mathbf{m}_k \mid \mathbf{\theta}) = -\frac{1}{2} (\mathbf{m}_k - \mathbf{\mu})^{\mathrm{T}} \mathbf{\Sigma}^{-1} (\mathbf{m}_k - \mathbf{\mu}) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\mathbf{\Sigma}|$$

$$\sum_{k=1}^{N} \nabla_{\boldsymbol{\theta}} \ln p(\mathbf{m}_{k} \mid \boldsymbol{\theta}) =$$

$$\sum_{k=1}^{N} \left[\mathbf{\Sigma}^{-1} (\mathbf{m}_{k} - \boldsymbol{\mu}), -\frac{1}{2} \frac{\partial}{\partial \mathbf{\Sigma}} (\mathbf{m}_{k} - \boldsymbol{\mu})^{\mathrm{T}} \mathbf{\Sigma}^{-1} (\mathbf{m}_{k} - \boldsymbol{\mu}) - \frac{1}{2} \frac{\partial}{\partial \mathbf{\Sigma}} \ln |\mathbf{\Sigma}| \right]^{\mathrm{T}} \stackrel{!}{=} \mathbf{0}$$

$$\Leftrightarrow \begin{cases} \sum_{k=1}^{N} \hat{\mathbf{\Sigma}}^{-1}(\mathbf{m}_{k} - \hat{\mathbf{\mu}}) = \mathbf{0} \\ \sum_{k=1}^{N} \left[-\frac{1}{2} \frac{\partial}{\partial \hat{\mathbf{\Sigma}}} (\mathbf{m}_{k} - \hat{\mathbf{\mu}})^{\mathrm{T}} \hat{\mathbf{\Sigma}}^{-1} (\mathbf{m}_{k} - \hat{\mathbf{\mu}}) - \frac{1}{2} \frac{\partial}{\partial \hat{\mathbf{\Sigma}}} \ln |\hat{\mathbf{\Sigma}}| \right] = \mathbf{0} \end{cases}$$

Beispiel, Fortsetzung

$$\begin{split} \sum_{k=1}^{N} & \left[-\frac{1}{2} \frac{\partial}{\partial \hat{\Sigma}} (\mathbf{m}_{k} - \hat{\mathbf{\mu}})^{\mathrm{T}} \hat{\Sigma}^{-1} (\mathbf{m}_{k} - \hat{\mathbf{\mu}}) - \frac{1}{2} \frac{\partial}{\partial \hat{\Sigma}} \ln |\hat{\Sigma}| \right] = \mathbf{0} \\ & - \frac{\partial}{\partial \hat{\Sigma}} \sum_{k=1}^{N} \left[(\mathbf{m}_{k} - \hat{\mathbf{\mu}})^{\mathrm{T}} \hat{\Sigma}^{-1} (\mathbf{m}_{k} - \hat{\mathbf{\mu}}) \right] - N \frac{\partial}{\partial \hat{\Sigma}} \ln |\hat{\Sigma}| = \mathbf{0} \\ & \frac{\partial}{\partial \hat{\Sigma}} \sum_{k=1}^{N} \left[(\mathbf{m}_{k} - \overline{\mathbf{m}} + \overline{\mathbf{m}} - \hat{\mathbf{\mu}})^{\mathrm{T}} \hat{\Sigma}^{-1} (\mathbf{m}_{k} - \overline{\mathbf{m}} + \overline{\mathbf{m}} - \hat{\mathbf{\mu}}) \right] + N \frac{\partial}{\partial \hat{\Sigma}} \ln |\hat{\Sigma}| = \mathbf{0} \\ & \frac{\partial}{\partial \hat{\Sigma}} \sum_{k=1}^{N} \left[(\mathbf{m}_{k} - \overline{\mathbf{m}})^{\mathrm{T}} \hat{\Sigma}^{-1} (\mathbf{m}_{k} - \overline{\mathbf{m}}) \right] + N \frac{\partial}{\partial \hat{\Sigma}} (\overline{\mathbf{m}} - \hat{\mathbf{\mu}})^{\mathrm{T}} \hat{\Sigma}^{-1} (\overline{\mathbf{m}} - \hat{\mathbf{\mu}}) + N \frac{\partial}{\partial \hat{\Sigma}} \ln |\hat{\Sigma}| = \mathbf{0} \\ & \frac{\partial}{\partial \hat{\Sigma}} \operatorname{tr} \left[\hat{\Sigma}^{-1} \sum_{k=1}^{N} (\mathbf{m}_{k} - \overline{\mathbf{m}}) (\mathbf{m}_{k} - \overline{\mathbf{m}})^{\mathrm{T}} \right] + N \frac{\partial}{\partial \hat{\Sigma}} (\overline{\mathbf{m}} - \hat{\mathbf{\mu}})^{\mathrm{T}} \hat{\Sigma}^{-1} (\overline{\mathbf{m}} - \hat{\mathbf{\mu}}) + N \frac{\partial}{\partial \hat{\Sigma}} \ln |\hat{\Sigma}| = \mathbf{0} \\ & N \frac{\partial}{\partial \hat{\Sigma}} \operatorname{tr} \left(\hat{\Sigma}^{-1} \widetilde{\mathbf{S}} \right) + N \frac{\partial}{\partial \hat{\Sigma}} (\overline{\mathbf{m}} - \hat{\mathbf{\mu}})^{\mathrm{T}} \hat{\Sigma}^{-1} (\overline{\mathbf{m}} - \hat{\mathbf{\mu}}) + N \frac{\partial}{\partial \hat{\Sigma}} \ln |\hat{\Sigma}| = \mathbf{0} \\ & \text{mit } \widetilde{\mathbf{S}} := \frac{1}{N} \mathbf{S} = \frac{1}{N} \sum_{k=1}^{N} (\mathbf{m}_{k} - \overline{\mathbf{m}})^{\mathrm{T}} (\mathbf{m}_{k} - \overline{\mathbf{m}})^{\mathrm{T}} (\mathbf{m}_{k} - \overline{\mathbf{m}}) \end{aligned}$$

Beispiel, Fortsetzung

Beispiel, Fortsetzung
$$\frac{\partial}{\partial \hat{\Sigma}} \operatorname{tr}(\hat{\Sigma}^{-1} \tilde{\mathbf{S}}) + \frac{\partial}{\partial \hat{\Sigma}} (\overline{\mathbf{m}} - \hat{\boldsymbol{\mu}})^{\mathrm{T}} \hat{\Sigma}^{-1} (\overline{\mathbf{m}} - \hat{\boldsymbol{\mu}}) - \frac{\partial}{\partial \hat{\Sigma}} \ln |\hat{\Sigma}^{-1}| = \mathbf{0}$$

$$\frac{\partial}{\partial \hat{\Sigma}} \operatorname{tr}(\mathbf{V} \tilde{\mathbf{S}}) + \frac{\partial}{\partial \hat{\Sigma}} (\overline{\mathbf{m}} - \hat{\boldsymbol{\mu}})^{\mathrm{T}} \mathbf{V} (\overline{\mathbf{m}} - \hat{\boldsymbol{\mu}}) - \frac{\partial}{\partial \hat{\Sigma}} \ln |\mathbf{V}| = \mathbf{0}$$

$$\frac{\partial \ln p(\mathbf{m}_k | \mathbf{0})}{\partial \hat{\Sigma}} = \frac{\partial \ln p(\mathbf{m}_k | \mathbf{0})}{\partial \mathbf{V}} \frac{\partial \mathbf{V}}{\partial \hat{\Sigma}}$$

$$\frac{\partial \mathbf{V}}{\partial \hat{\Sigma}} = \hat{\Sigma}^{-2}$$

$$\frac{\partial}{\partial \mathbf{V}} \operatorname{tr}(\mathbf{V}\widetilde{\mathbf{S}}) + \frac{\partial}{\partial \mathbf{V}} (\overline{\mathbf{m}} - \hat{\mathbf{\mu}})^{\mathrm{T}} \mathbf{V} (\overline{\mathbf{m}} - \hat{\mathbf{\mu}}) - \frac{\partial}{\partial \mathbf{V}} \ln |\mathbf{V}| = \mathbf{0}$$

$$\frac{\partial}{\partial \mathbf{X}} \ln |\mathbf{X}| = 2\mathbf{X}^{-1} - \operatorname{diag}(\mathbf{X}^{-1})$$

$$\frac{\partial}{\partial \mathbf{X}} \operatorname{tr}(\mathbf{A}\mathbf{X}) = \frac{\partial}{\partial \mathbf{X}} \operatorname{tr}(\mathbf{X}\mathbf{A}) = \mathbf{A}^{\mathrm{T}} + \mathbf{A} - \operatorname{diag}(\mathbf{A})$$

$$2\tilde{\mathbf{S}} - \operatorname{diag}(\tilde{\mathbf{S}}) + \frac{\partial}{\partial \mathbf{V}} (\overline{\mathbf{m}} - \hat{\boldsymbol{\mu}})^{\mathrm{T}} \mathbf{V} (\overline{\mathbf{m}} - \hat{\boldsymbol{\mu}}) - 2\hat{\boldsymbol{\Sigma}} + \operatorname{diag}(\hat{\boldsymbol{\Sigma}}) = \mathbf{0}$$

Beispiel, Fortsetzung

$$\begin{cases} \hat{\mathbf{\mu}} = \frac{1}{N} \sum_{k=1}^{N} \mathbf{m}_{k} \\ 2\tilde{\mathbf{S}} - \operatorname{diag}(\tilde{\mathbf{S}}) + \frac{\partial}{\partial \mathbf{V}} (\overline{\mathbf{m}} - \hat{\mathbf{\mu}})^{\mathrm{T}} \mathbf{V} (\overline{\mathbf{m}} - \hat{\mathbf{\mu}}) - 2\hat{\mathbf{\Sigma}} + \operatorname{diag}(\hat{\mathbf{\Sigma}}) = \mathbf{0} \end{cases}$$

$$\Leftrightarrow \begin{cases} \hat{\mathbf{\mu}} = \frac{1}{N} \sum_{k=1}^{N} \mathbf{m}_{k} \\ 2\tilde{\mathbf{S}} - 2\hat{\mathbf{\Sigma}} + \operatorname{diag}(\hat{\mathbf{\Sigma}} - \tilde{\mathbf{S}}) = \mathbf{0} \end{cases} \Leftrightarrow$$

ML-Schätzer:

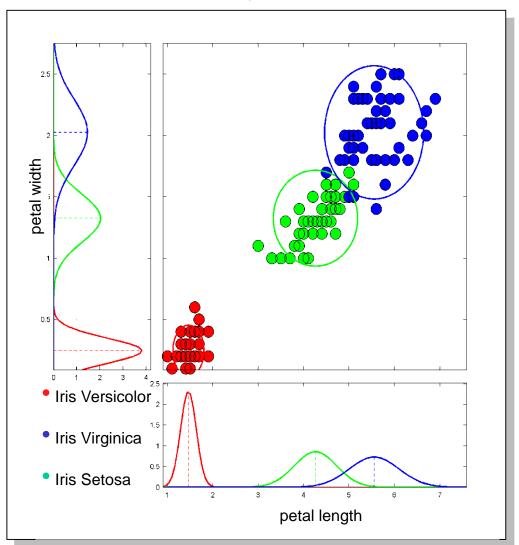
$$\hat{\boldsymbol{\mu}}_{\text{ML}} = \frac{1}{N} \sum_{k=1}^{N} \mathbf{m}_{k}$$

$$\hat{\boldsymbol{\Sigma}}_{\text{ML}} = \frac{1}{N} \sum_{k=1}^{N} (\mathbf{m}_{k} - \hat{\boldsymbol{\mu}})(\mathbf{m}_{k} - \hat{\boldsymbol{\mu}})^{\text{T}}$$

Differentiation nach Vektoren und Matrizen, siehe z.B.:

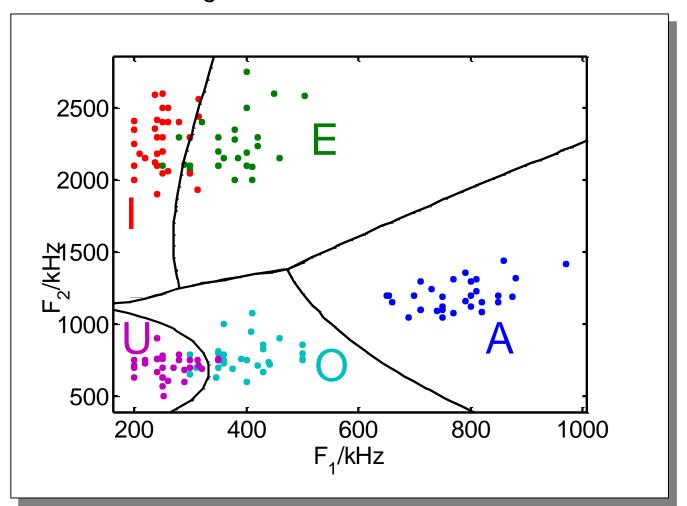
- T. K. Moon; W. C. Stirling: Mathematical Methods and Algorithms for Signal Processing, Prentice Hall, 2000.
- C. Voigt; J. Adamy: Formelsammlung der Matrizenrechnung, Oldenbourg Verlag, 2007.

Beispiel: Iris Datensatz nach R. A. Fisher, Modell: Gauß'sche WDFen mit unkorrelierten Merkmalen, Maximum Likelihood Schätzung der Parameter



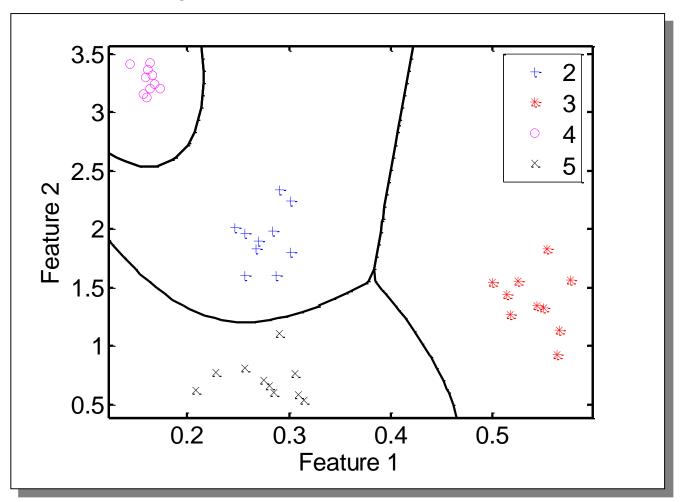
Ellipsen kennzeichnen die 2σ - Grenzen der Wahrscheinlichkeitsdichten.

Beispiel: Vokalerkennung



Bayes'sche Klassifikation, Modell: Normalverteilung für klassenbedingte Merkmals-WDFen, Parameterschätzung: Maximum Likelihood

Beispiel: Grifferkennung Hand



Bayes'sche Klassifikation, Modell: Normalverteilung für klassenbedingte Merkmals-WDFen, Parameterschätzung: Maximum Likelihood

Fundamentalgröße der Bayes'schen Klassifikation:

A Posteriori Wahrscheinlichkeitsverteilungen der Klassen

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

Problem: WVen unbekannt → Müssen aus den Daten D geschätzt werden.

Ansatz:

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

Die $p(\mathbf{m} | \omega_i, \mathbf{D})$ approximieren die $p(\mathbf{m} | \omega_i)$.

Vereinfachende Annahmen:

- A Priori Wahrscheinlichkeiten seien gegeben oder ausreichend genau bestimmt worden: $P(\omega_i) = P(\omega_i \mid D)$
- Die Stichproben in D_k haben keinen Einfluss auf $p(\mathbf{m}|\omega_i, D)$ $i \neq k$.

$$P(\omega_i \mid \mathbf{m}, \mathbf{D}) = \frac{p(\mathbf{m} \mid \omega_i, \mathbf{D}_i) P(\omega_i)}{\sum_{j=1}^{c} p(\mathbf{m} \mid \omega_j, \mathbf{D}_j) P(\omega_j)}$$

■ WDFen $p(\mathbf{m} \mid \omega_i)$ haben eine bekannte parametrisierte Form $p(\mathbf{m} \mid \mathbf{\theta}_i, \omega_i)$ mit $p(\mathbf{\theta}_i \mid \omega_i)$ bekannt.

$$p(\mathbf{m} \mid \omega_i, \mathbf{D}_i) = \int_{\Theta_i} p(\mathbf{m}, \mathbf{\theta}_i \mid \omega_i, \mathbf{D}_i) \, d\mathbf{\theta}_i = \int_{\Theta_i} p(\mathbf{m} \mid \mathbf{\theta}_i, \omega_i, \mathbf{p}_i) p(\mathbf{\theta}_i \mid \omega_i, \mathbf{D}_i) \, d\mathbf{\theta}_i$$

$$\Rightarrow p(\mathbf{m} \mid \omega_i, D_i) = \int_{\Theta_i} p(\mathbf{m} \mid \mathbf{\theta}_i, \omega_i) p(\mathbf{\theta}_i \mid \omega_i, D_i) d\mathbf{\theta}_i$$

Die Gesamtaufgabe zerfällt in c Teilaufgaben; eine pro Klasse. Für jede Klasse ω_i gilt es, die klassenspezifische WDF

$$p(\mathbf{m} \mid \omega_i, \mathbf{D}_i) = \int_{\Theta_i} p(\mathbf{m} \mid \mathbf{\theta}_i, \omega_i) p(\mathbf{\theta}_i \mid \omega_i, \mathbf{D}_i) d\mathbf{\theta}_i \qquad (*)$$

anhand der Daten D, zu bestimmen.

Zur einfacheren Schreibweise wird im Folgenden die Kennzeichnung der Klassen unterdrückt. Die Überlegungen gelten aber immer nur für eine Klasse ω_i .

→ Vereinfachte Schreibweise für (*):

$$p(\mathbf{m} \mid \mathbf{D}) = \int_{\Theta} p(\mathbf{m} \mid \mathbf{\theta}) p(\mathbf{\theta} \mid \mathbf{D}) d\mathbf{\theta}$$

Da $p(\mathbf{m} | \mathbf{\theta})$ als bekannt vorausgesetzt ist, gilt es im Folgenden, $p(\mathbf{\theta} | \mathbf{D})$ zu bestimmen.

Annahme: $\mu \sim N(\mu; \mu_0, \sigma_0^2)$

 (μ_0, σ_0^2) sind sogenannte Hyperparameter)

$$p(\mu \mid D) = \frac{p(D \mid \mu)p(\mu)}{\int p(D \mid \mu)p(\mu)d\mu} = \alpha(D) \prod_{k=1}^{N} p(m_k \mid \mu)p(\mu)$$

$$p(\mu \mid \mathbf{D}) = \alpha(\mathbf{D}) \prod_{k=1}^{N} \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2} \left(\frac{m_k - \mu}{\sigma}\right)^2\right] \frac{1}{\sqrt{2\pi}\sigma_0} \exp\left[-\frac{1}{2} \left(\frac{\mu - \mu_0}{\sigma_0}\right)^2\right]$$

$$p(m_k \mid \mu)$$

$$p(\mu)$$

$$= \alpha'(D) \exp \left[-\frac{1}{2} \left(\sum_{k=1}^{N} \left(\frac{\mu - m_k}{\sigma} \right)^2 + \left(\frac{\mu - \mu_0}{\sigma_0} \right)^2 \right) \right]$$

$$=\alpha''(\mathrm{D})\exp\left[-\frac{1}{2}\left[\left(\frac{N}{\sigma^2}+\frac{1}{\sigma_0^2}\right)\mu^2-2\left(\frac{1}{\sigma^2}\sum_{k=1}^N m_k+\frac{\mu_0}{\sigma_0^2}\right)\mu\right]\right] \quad \begin{array}{l} (\to \dots \text{ kann nur eigenfunction of the sein! Quadratisches the properties of the prop$$

(→ ... kann nur eine sein! Quadratische Ergänzung → ...)

Beispiel, Fortsetzung

Diese WDF lässt sich schreiben als: $p(\mu \mid D) = N(\mu; \mu_N, \sigma_N^2)$

$$p(\mu \mid \mathbf{D}) = \frac{1}{\sqrt{2\pi}\sigma_N} \exp \left[-\frac{1}{2} \left(\frac{\mu - \mu_N}{\sigma_N} \right)^2 \right]$$

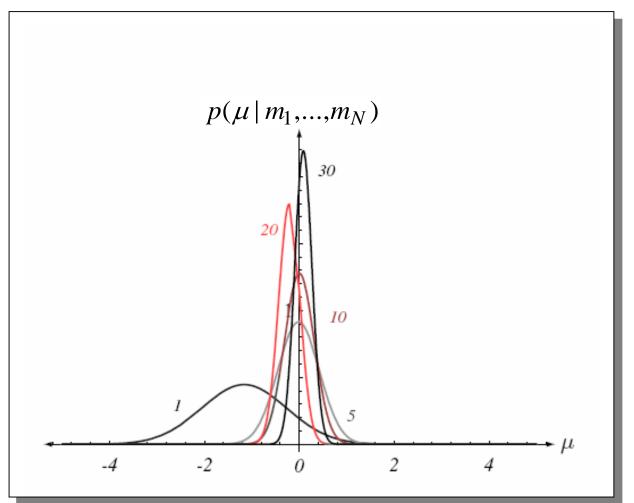
$$\mu_{N} := \left(\frac{N\sigma_{0}^{2}}{N\sigma_{0}^{2} + \sigma^{2}}\right) \hat{\mu}_{N} + \frac{\sigma^{2}}{N\sigma_{0}^{2} + \sigma^{2}} \mu_{0} \quad \sigma_{N}^{2} := \frac{\sigma_{0}^{2}\sigma^{2}}{N\sigma_{0}^{2} + \sigma^{2}} \qquad \hat{\mu}_{N} := \frac{1}{N} \sum_{k=1}^{N} m_{k}$$

$$\hat{\mu}_N \coloneqq \frac{1}{N} \sum_{k=1}^N m_k$$

Diskussion:

- A Priori + Empirische Information \rightarrow A Posteriori WDF $p(\mu | D)$
- die beste Schätzung (MAP) für μ nach N beobachteten Werten.
- beschreibt die Schätzunsicherheit (stochastischer Schätzfehler)
- μ_N liegt zwischen $\hat{\mu}_N$ und μ_0

Beispiel, Fortsetzung



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

A Posteriori Verteilungsdichte für wachsende Anzahl von Beobachtungen.

Beispiel, Fortsetzung

$$p(m \mid D) = \int_{\Theta} p(m \mid \mu) p(\mu \mid D) d\mu$$

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

$$= \frac{1}{2\pi\sigma\sigma_N} \exp\left[-\frac{1}{2} \frac{(m-\mu_N)^2}{\sigma^2 + \sigma_N^2}\right] f(\sigma, \sigma_N)$$

$$\text{mit } f(\sigma, \sigma_N) \coloneqq \int \exp \left[-\frac{1}{2} \frac{\sigma^2 + \sigma_N^2}{\sigma^2 \sigma_N^2} \left(\mu - \frac{\sigma_N^2 m + \sigma^2 \mu_N}{\sigma_N^2 + \sigma^2} \right)^2 \right] \mathrm{d}\mu$$

$$p(m \mid D) = N(m; \mu_N, \sigma^2 + \sigma_N^2)$$

Beispiel: $p(\mathbf{m} | \boldsymbol{\mu}) = N(\mathbf{m}; \boldsymbol{\mu}, \boldsymbol{\Sigma}), \boldsymbol{\mu}$ unbekannt

Annahme: $p(\mu) = N(\mu; \mu_0, \Sigma_0)$

$$p(\mathbf{\mu} \mid \mathbf{D}) = \alpha \prod_{k=1}^{N} p(\mathbf{m}_{k} \mid \mathbf{\mu}) p(\mathbf{\mu})$$

$$= \alpha' \exp \left[-\frac{1}{2} \left(\mathbf{\mu}^{\mathrm{T}} (N \mathbf{\Sigma}^{-1} + \mathbf{\Sigma}_{0}^{-1}) \mathbf{\mu} - 2 \mathbf{\mu}^{\mathrm{T}} \left(\mathbf{\Sigma}^{-1} \sum_{k=1}^{N} \mathbf{m}_{k} + \mathbf{\Sigma}_{0}^{-1} \mathbf{\mu}_{0} \right) \right) \right]$$

$$= \alpha'' \exp \left[-\frac{1}{2} (\mathbf{\mu} - \mathbf{\mu}_{N})^{\mathrm{T}} \mathbf{\Sigma}_{N}^{-1} (\mathbf{\mu} - \mathbf{\mu}_{N}) \right]$$

$$p(\mathbf{\mu} \mid \mathbf{D}) = \mathbf{N}(\mathbf{\mu}; \mathbf{\mu}_N, \mathbf{\Sigma}_N)$$

$$\mathbf{\mu}_N = \mathbf{\Sigma}_0 \left(\mathbf{\Sigma}_0 + \frac{1}{N} \mathbf{\Sigma} \right)^{-1} \hat{\mathbf{\mu}}_N + \frac{1}{N} \mathbf{\Sigma} \left(\mathbf{\Sigma}_0 + \frac{1}{N} \mathbf{\Sigma} \right)^{-1} \mathbf{\mu}_0$$

$$\mathbf{\Sigma}_N = \mathbf{\Sigma}_0 (\mathbf{\Sigma}_0 + \frac{1}{N} \mathbf{\Sigma})^{-1} \frac{1}{N} \mathbf{\Sigma}$$

$$\hat{\mathbf{\mu}}_N = \frac{1}{N} \sum_{k=1}^{N} \mathbf{m}_k$$

Beispiel, Fortsetzung

$$p(\mathbf{m} \mid \mathbf{D}) = \int p(\mathbf{m} \mid \boldsymbol{\mu}) p(\boldsymbol{\mu} \mid D) d\boldsymbol{\mu}$$

$$p(\mathbf{\mu} \mid \mathbf{D}) = N(\mathbf{\mu}_N, \mathbf{\Sigma}_N)$$

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

$$p(\mathbf{m} \mid \mathbf{D}) = \mathbf{N}(\mathbf{m}; \boldsymbol{\mu}_N, \boldsymbol{\Sigma} + \boldsymbol{\Sigma}_N)$$

Überblick zum Verfahren; für jede Klasse i = 1,...,c separat anzuwenden.

- $p(\mathbf{m} | \mathbf{\theta}_i, \omega_i)$ strukturell als bekannt angenommen; Wert $\mathbf{\theta}_i$ unbekannt.
- $p(\theta_i \mid \omega_i)$ bekannt; enthält das A Priori Wissen über θ_i .
- Das weitere Wissen über θ_i ist in D_i enthalten; D_i besteht aus N_i unabhängigen Stichproben $\mathbf{m}_1,...,\mathbf{m}_{N_i}$ bezüglich der WDF $p(\mathbf{m}|\omega_i)$.

$$p(\mathbf{D}_{i} \mid \boldsymbol{\theta}_{i}, \omega_{i}) = \prod_{k=1}^{N} p(\mathbf{m}_{k} \mid \boldsymbol{\theta}_{i}, \omega_{i})$$

$$p(\boldsymbol{\theta}_{i} \mid \mathbf{D}_{i}, \omega_{i}) = \frac{p(\mathbf{D}_{i} \mid \boldsymbol{\theta}_{i}, \omega_{i}) p(\boldsymbol{\theta}_{i} \mid \omega_{i})}{\int p(\mathbf{D}_{i} \mid \boldsymbol{\theta}_{i}, \omega_{i}) p(\boldsymbol{\theta}_{i} \mid \omega_{i}) d\boldsymbol{\theta}_{i}}$$

$$p(\mathbf{m} \mid \mathbf{D}_{i}, \omega_{i}) = \int p(\mathbf{m} \mid \boldsymbol{\theta}_{i}, \omega_{i}) p(\boldsymbol{\theta}_{i} \mid \mathbf{D}_{i}, \omega_{i}) d\boldsymbol{\theta}_{i}$$

■ Wenn $p(\mathbf{\theta}_i | \mathbf{D}_i, \omega_i)$ sich stark um einen Wert $\mathbf{\theta}_i = \hat{\mathbf{\theta}}_i$ konzentriert, dort eine ausgeprägte Spitze ausbildet, dann kann $p(\mathbf{m} | \mathbf{D}_i, \omega_i)$ durch $p(\mathbf{m}/\hat{\mathbf{\theta}}_i, \omega_i)$ approximiert werden.

Vergleich: Bayes'sches Verfahren und Maximum Likelihood Schätzung.

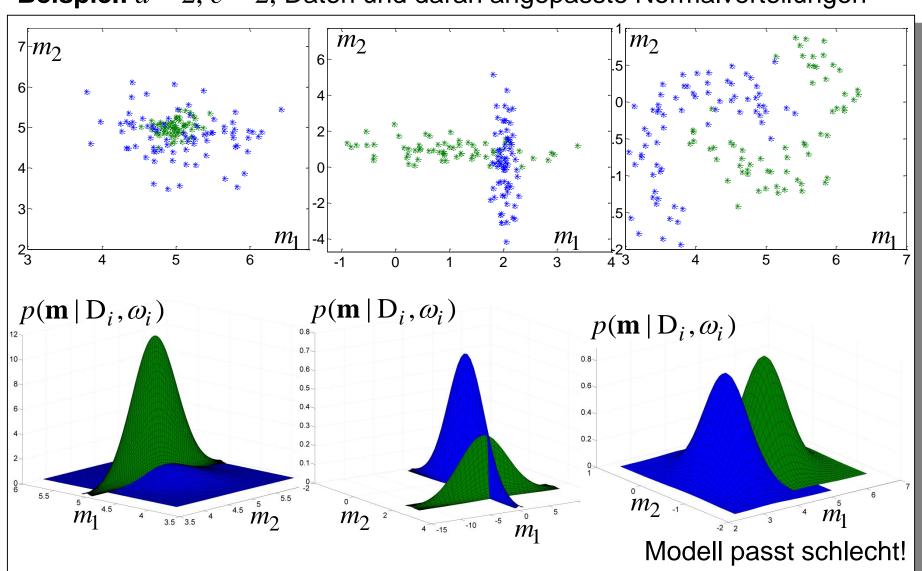
Bayes'sches Verfahren:

- + A Priori Wissen über θ kann eingebracht werden.
- Numerische Berechung mehrdimensionaler Integrale erforderlich.

Maximum Likelihood Verfahren:

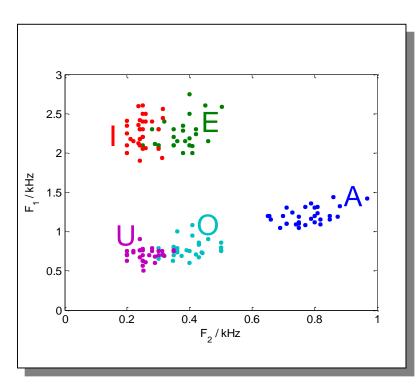
- Keine mehrdimensionalen Integrale, sondern nur Extremwertsuche.
- A Priori Wissen über θ kann nicht eingebracht werden.

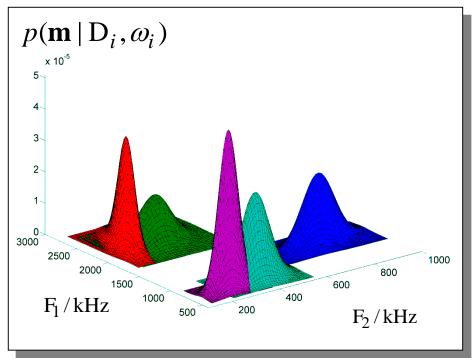
Beispiel: d = 2, c = 2, Daten und daran angepasste Normalverteilungen



Beispiel: Vokalerkennung, d = 2, c = 5,

Daten und daran angepasste Normalverteilungen





4.3. Bayes'sche Parameterschätzung

Wenn $p(\theta \mid D)$ sich stark um einen Wert $\theta = \hat{\theta}$ konzentriert, dort eine ausgeprägte Spitze ausbildet, dann kann $p(\mathbf{m} \mid D)$ durch $p(\mathbf{m}/\hat{\theta})$ approximiert werden.

Mit dieser Approximation wird der Rechenaufwand des Verfahrens aus Abschnitt 4.2. erheblich reduziert, da insbesondere das Integral

$$p(\mathbf{m} \mid \mathbf{D}) = \int p(\mathbf{m} \mid \mathbf{\theta}) p(\mathbf{\theta} \mid \mathbf{D}) d\mathbf{\theta}$$

nicht ausgewertet werden muss.

$$p(\mathbf{m} \mid \mathbf{D}) = \int p(\mathbf{m} \mid \mathbf{\theta}) p(\mathbf{\theta} \mid \mathbf{D}) d\mathbf{\theta} \approx \int p(\mathbf{m} \mid \mathbf{\theta}) \delta(\mathbf{\theta} - \hat{\mathbf{\theta}}) d\mathbf{\theta} = p(\mathbf{m} \mid \hat{\mathbf{\theta}})$$

Im Folgenden werden die wichtigsten Techniken zur Bayes'schen Parameterschätzung gezeigt.

4.3. Bayes'sche Parameterschätzung

θ wird als Zufallsvariable angesehen und über eine WV beschrieben.

Ansatz 1: Minimaler quadratischer Schätzfehler

$$\mathbf{E}\{\|\hat{\boldsymbol{\theta}}(\mathbf{D}) - \boldsymbol{\theta}\|^2\} = \int_{\mathbf{IM}^N \Theta} \int l(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta}) p(\boldsymbol{\theta}, \mathbf{D}) d\boldsymbol{\theta} d\mathbf{D} \overset{!}{\to} \text{minimal} d\mathbf{D} \coloneqq \prod_{i=1}^N d\mathbf{m}_i$$

Kostenfunktion:

$$l(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta}) := (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})^{\mathrm{T}} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})$$

$$E\{\|\hat{\boldsymbol{\theta}}(\mathbf{D}) - \boldsymbol{\theta}\|^2\} = \int_{\mathbf{IM}^N \Theta} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})^{\mathrm{T}} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathbf{D}) d\boldsymbol{\theta} \ p(\mathbf{D}) d\mathbf{D} \overset{!}{\to} \text{minimal}$$

$$\Leftrightarrow I(\mathbf{D}) \coloneqq \int_{\Theta} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})^{\mathrm{T}} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) p(\boldsymbol{\theta} \mid \mathbf{D}) d\boldsymbol{\theta} \text{ ist minimal.}$$

$$I(\mathbf{D}) = \int_{\Theta} \sum_{j=1}^{K} (\hat{\theta}_{j}(\mathbf{D}) - \theta_{j})^{2} p(\mathbf{\theta} \mid \mathbf{D}) d\mathbf{\theta}$$

$$\frac{\partial I(\mathbf{D})}{\partial \hat{\theta}_{i}(\mathbf{D})} \stackrel{!}{=} 0 \qquad \Leftrightarrow \qquad \int_{\Theta} 2(\hat{\theta}_{i}(\mathbf{D}) - \theta_{i}) p(\mathbf{\theta} \mid \mathbf{D}) d\mathbf{\theta} = 0$$

$$\Leftrightarrow \qquad \hat{\theta}_{i}(\mathbf{D}) = \int_{\Theta} \theta_{i} p(\mathbf{\theta} \mid \mathbf{D}) d\mathbf{\theta}$$

Hinreichende Bedingung für Minimum:

$$\frac{\partial^2 I(D)}{\partial (\hat{\theta}_i(D))^2} = 2 > 0$$

Schätzer mit minimalem mittleren quadratischen Fehler:

$$\hat{\mathbf{\theta}}(\mathbf{D}) = \int_{\Theta} \mathbf{\theta} p(\mathbf{\theta} \mid \mathbf{D}) d\mathbf{\theta} = \mathbf{E}\{\mathbf{\theta} \mid \mathbf{D}\}$$

mit
$$p(\mathbf{\theta} \mid \mathbf{D}) = \frac{p(\mathbf{D} \mid \mathbf{\theta})p(\mathbf{\theta})}{p(\mathbf{D})}$$

A Posteriori Erwartungswert der Zufallsvariablen θ

4.3. Bayes'sche Parameterschätzung

Ansatz 2: Konstante Gewichtung großer Fehler

Kostenfunktion:

$$l(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta}) := \begin{cases} 0 & \text{für} \\ 1 & \text{sonst} \end{cases} \|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\| < \Delta, \ \Delta > 0$$

$$E\{l(\hat{\boldsymbol{\theta}},\boldsymbol{\theta})\} = \int_{\mathbb{M}^N \Theta} \int l(\hat{\boldsymbol{\theta}},\boldsymbol{\theta}) p(\boldsymbol{\theta},D) d\boldsymbol{\theta} dD \xrightarrow{!} \text{minimal}$$

$$\Leftrightarrow$$

$$I(D) = \int_{\Theta} l(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta}) p(\boldsymbol{\theta} \mid D) d\boldsymbol{\theta}$$

ist minimal.

$$\Leftrightarrow$$

$$I(\mathbf{D}) = 1 - \int_{\{\boldsymbol{\theta} | \|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\| < \Delta\}} p(\boldsymbol{\theta} | \mathbf{D}) d\boldsymbol{\theta}$$

ist minimal.

Für $\Delta \rightarrow 0$ folgt:

$$\hat{\mathbf{\theta}}(\mathbf{D}) = \underset{\mathbf{\theta}}{\text{arg max}} \{ p(\mathbf{\theta} \mid \mathbf{D}) \}$$

MAP-Schätzer

4.3. Bayes'sche Parameterschätzung

Ansatz 2: Konstante Gewichtung großer Fehler

Wegen der Monotonie der Logarithmusfunktion, kann man äquivalent die logarithmierte A Posteriori WDF maximieren; das bringt oft analytische Vorteile.

$$\hat{\boldsymbol{\theta}}(\mathbf{D}) = \underset{\boldsymbol{\theta}}{\text{arg max}} \{ \ln p(\boldsymbol{\theta} \mid \mathbf{D}) \}$$

MAP-Schätzer

Falls $\ln p(\theta|\mathbf{m})$ differenzierbar ist, lautet die notwendige Bedingung für Maxima:

$$\nabla_{\boldsymbol{\theta}} \ln(p(\boldsymbol{\theta} \mid \mathbf{D})) = \sum_{k=1}^{N} \nabla_{\boldsymbol{\theta}} \ln p(\mathbf{m}_{k} \mid \boldsymbol{\theta}) + \nabla_{\boldsymbol{\theta}} \ln p(\boldsymbol{\theta}) \stackrel{!}{=} \mathbf{0}$$

In vielen praktisch relevanten Fällen liefern MAP-Schätzung und der A Posteriori Erwartungswert das gleiche Ergebnis.

Details siehe z.B. in: K. Kroschel, "Statistische Informationstechnik", Springer 2004

4.3. Bayes'sche Klassifikation – ergänzende Bemerkungen

Fehler bei der Bayes'schen Klassifikation

- Bayes'scher Fehler: ergibt sich durch die Überlappung der Verteilungsdichten; spiegelt die Diskriminanz der Merkmale wider.
- 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.