

Sprachverarbeitung - 4. Woche

Andreas Wendemuth



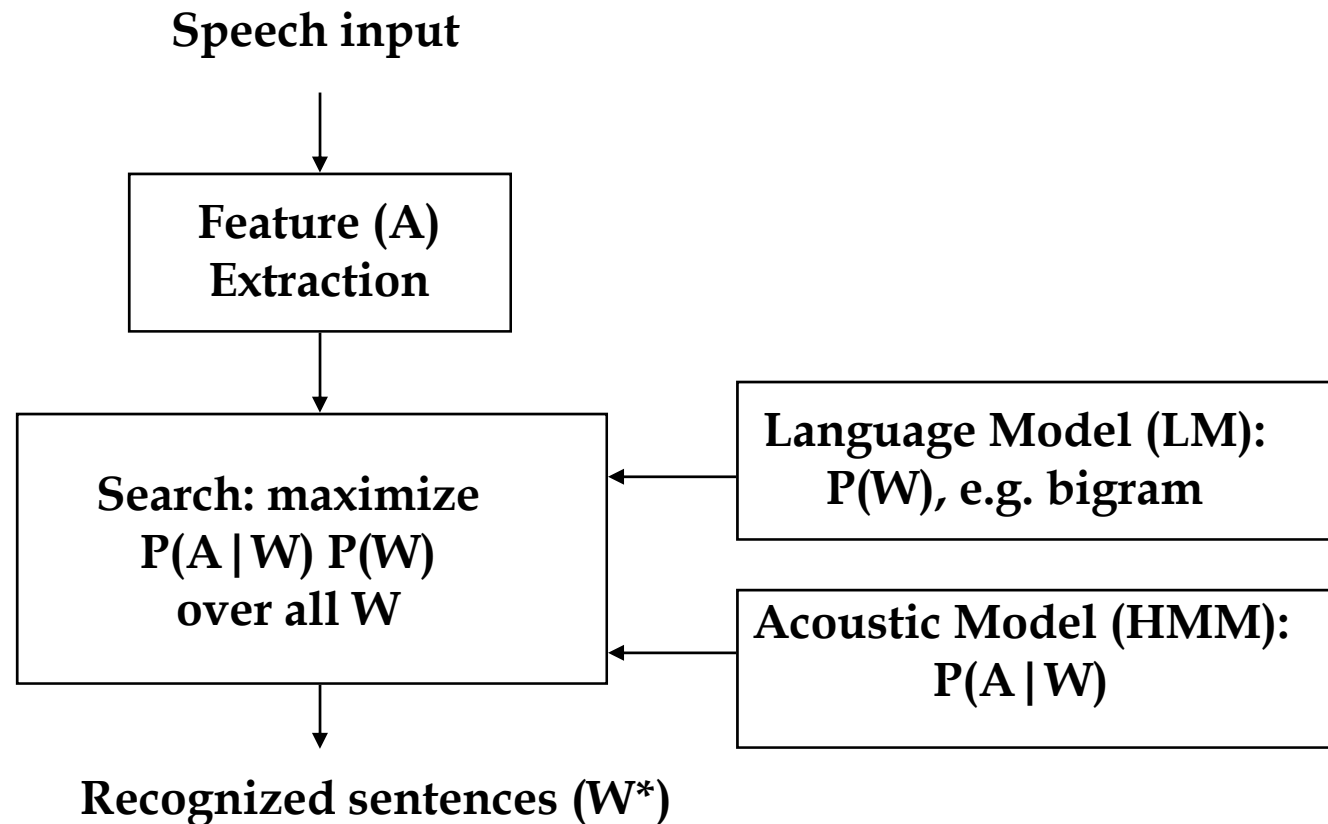
Letzte Woche:

1. Overview Speech Recognition Systems & Architectures
2. Acoustic modeling & feature extraction (1)
3. Feature extraction (2)
4. Klassifikation in HMM Modellen
5. Wortmodellierung (trigramme, tying)
6. search/decoding, lattices,
wordgraphs, confidence measures
7. acoustic adaptation
8. language models and grammars, Language model
adaptation, lexica, phonology
9. speech understanding, dialogue control
10. Design of computer speech recognition systems

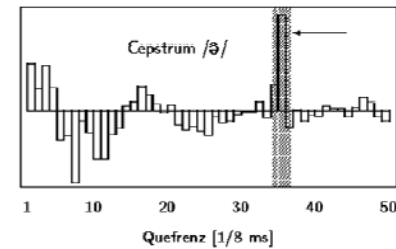
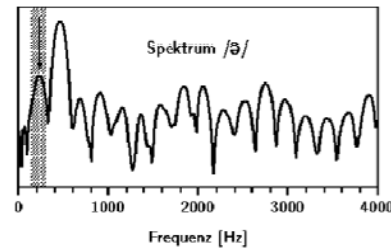
SR Architecture revisited

W = a word sequence (e.g. word/ sentence/ whole dictation)

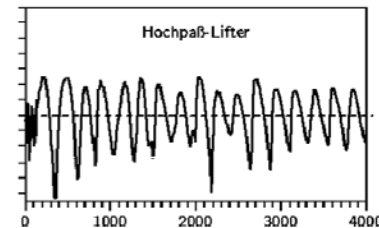
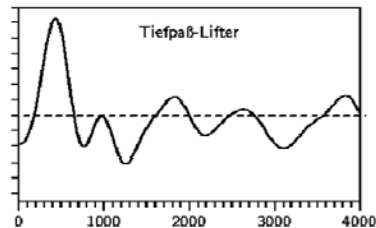
A = an acoustic feature vector sequence (the input for the recognizer)



Homomorphe Analyse (Cepstrum)



- Cepstrum \Leftrightarrow Spektrum des logarith. Betragsspektrums
- **Quefrenz** = Periodendauer in Sekunden (Einheit ist $1/f_A$)
- **Cepstralgipfel** markiert die Dauer der Grundperiode
- **Formantstruktur** = niedrige Cepstralkoeffizienten



- **Lifterung** — Rücktransformation in den Spektralbereich

$$\{\hat{C}_\nu^{(m)}\} = \text{DFT}\{\hat{c}_q^{(m)}\}$$

- Grobstruktur: $\hat{c}_q^{(m)} = 0$ für $q > 20$ (Tiefpaß)
- Feinstruktur: $\hat{c}_q^{(m)} = 0$ für $q < 20$ (Hochpaß)

Lineare Vorhersage

Lineare Vorhersageformel, Ordnung p

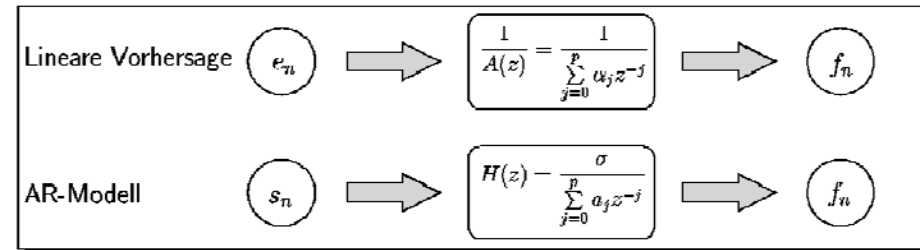
$$\hat{f}_n = - \sum_{j=1}^p \alpha_j f_{n-j}$$

Vorhersagefehler mit den Prädiktionskoeffizienten α_j

$$e_n = f_n - \hat{f}_n = \sum_{j=0}^p \alpha_j f_{n-j} \quad (\alpha_0 = 1)$$

Nach z -Transformation ergeben sich die Systemgleichungen

$$E(z) = F(z) \cdot A(z) \quad \text{und} \quad F(z) = E(z) \cdot \frac{1}{A(z)},$$



Wenn die Vorhersagekoeffizienten α_j mit den Modellparametern a_j übereinstimmen:

$$e_n = \sigma s_n \quad (\text{Fehler} = \text{Verstärkung} \times \text{Anregung})$$

Lineare Vorhersage

Das Modellspektrum

Verknüpfen von Vorhersage- und Produktionsmodell durch die idealisierte Annahme

$$H(z) = \sigma/A(z)$$

Querschnittflächen A_j aus partieller Korrelation (PARCOR) k_n :

$$A_{j+1} = \frac{1 - k_j}{1 + k_j} \cdot A_j$$

Schätzung des **Vokaltraktfrequenzgangs**

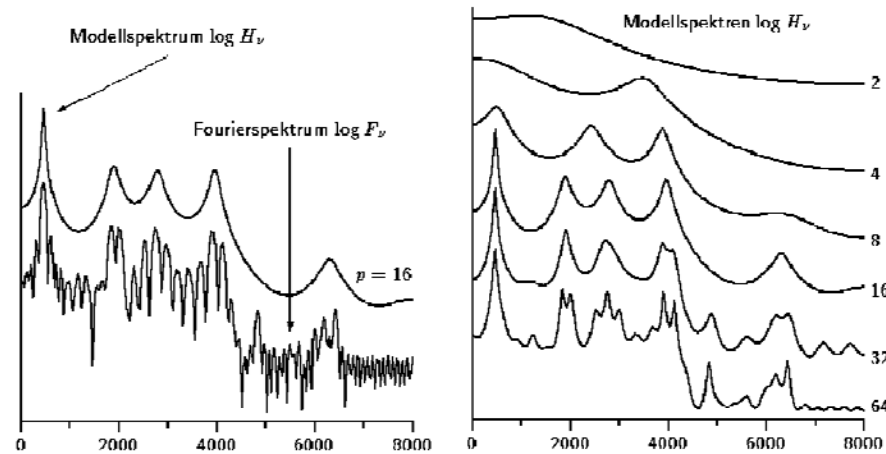
$$H_\nu = H(e^{2\pi i \nu / N}) = \frac{\sigma}{A(e^{2\pi i \nu / N})} = \frac{\sigma}{A_\nu}$$

Funktionswerte A_ν von $A(z)$ auf dem Einheitskreis durch

$$\text{DFT}\{1, \alpha_1, \dots, \alpha_p, \underbrace{0, \dots, 0}_{N-p-1 \text{ Nullen}}\}$$

Verstärkungsfaktor aus Autokorrelationsfunktion

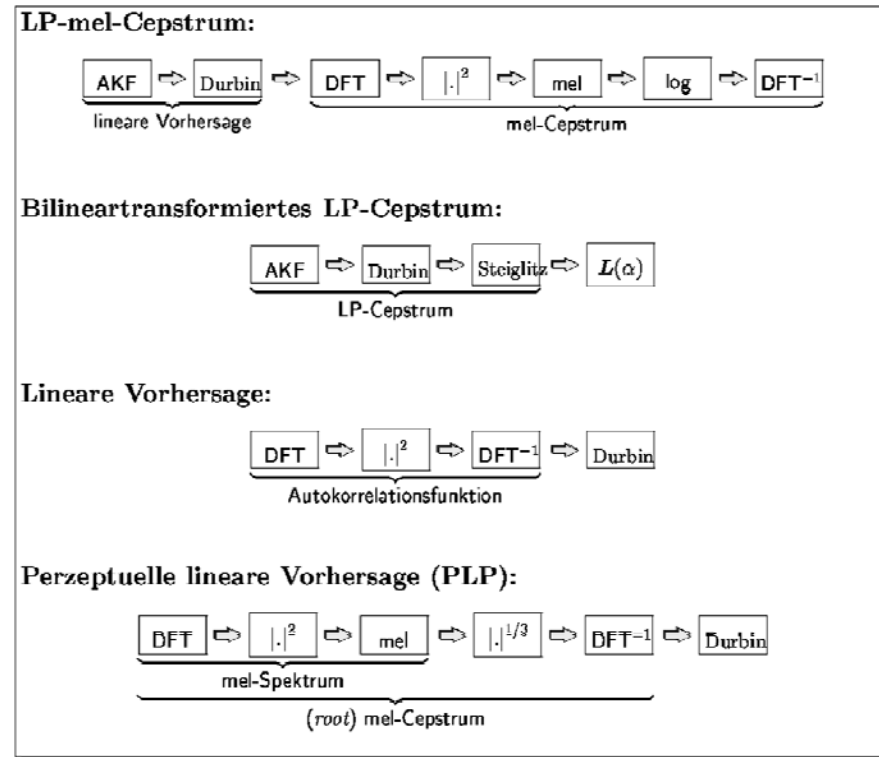
$$\sigma^2 = \sum_{j=0}^p \alpha_j r_j$$



Verzerrung der Frequenzachse

Gehörrichtige Frequenzverzerrung

Integration der Mel-Skalierung in die Berechnung der Vorhersagemerkmale



Bemerkung 3.4

- $L(\alpha)$ ist eine geeignete **Bilineartransformation** zur Frequenzverzerrung im Cepstralbereich
- PLP nach dem Satz von Wiener & Khintchine:

$$\text{FT}\{r_n\} = |\text{FT}\{f_n\}|^2 \quad \text{bzw.} \quad \{r_n\} = \text{FT}^{-1}\{|\text{FT}\{f_n\}|^2\}$$

Diese Woche:

1. Overview Speech Recognition Systems & Architectures
2. Acoustic modeling & feature extraction (1)
3. Feature extraction (2)
4. LDA, Grundlagen der Klassifikation
5. Klassifikation in HMM Modellen
6. Wortmodellierung (trigramme, tying)
7. search/decoding, lattices,
wordgraphs, confidence measures
8. acoustic adaptation
9. language models and grammars, Language model
adaptation, lexica, phonology
10. speech understanding, dialogue control
11. Design of computer speech recognition systems

Merkmals Transformationen zur Dimensionsreduktion

- Suchraum kann massiv eingeschränkt werden durch Dimensionsreduktion
- Rauschen kann in „subspaces“ angenommen werden, die durch Dimensionsreduktion verschwinden können
- Übersicht:
 1. Karhunen-Loeve-Transformation (Hauptachsen)
 2. (Quotient) Singular Value Decomposition (Q)SVD
„Subspace Identification for linear Systems“,
z.B. Buch (1996): von Overschee and De Moor
 3. Lineare Diskriminanten Analyse (LDA)

Hauptachsentransformation

- - bestimme Drehung des Raumes (Unitäre Transf.)
- - Funktion darstellbar als Produkt in neuen Koordin.
- - Verzerrungsfaktoren sind die Eigenwerte (EW)
- - z.B. schätze die Matrix einer Normalverteilung als Korrelationsmatrix der Beobachtungswerte, auf Hauptachsen stellt die Korr.matrix Ellipsen dar
- - behalte nur die Koordinaten mit den grössten EW
- Problem : - nicht o.k. bei Adidas-Verteilung(s.LDA)
 - berücksichtigt keine Klassenzugehörigkeit

Lineare Merkmalstransformationen

Karhunen-Loève-Transformation

Un/vollständige Entwicklung von x nach der Basis ϕ_1, \dots, ϕ_D :

$$x = \sum_{i=1}^D y_i \phi_i \quad \text{bzw.} \quad \hat{x} = \sum_{i=1}^d y_i \phi_i \quad \text{mit Koeffizienten } y_i = \phi_i^\top x$$

Mittlerer quadratischer Fehler (Dimension d)

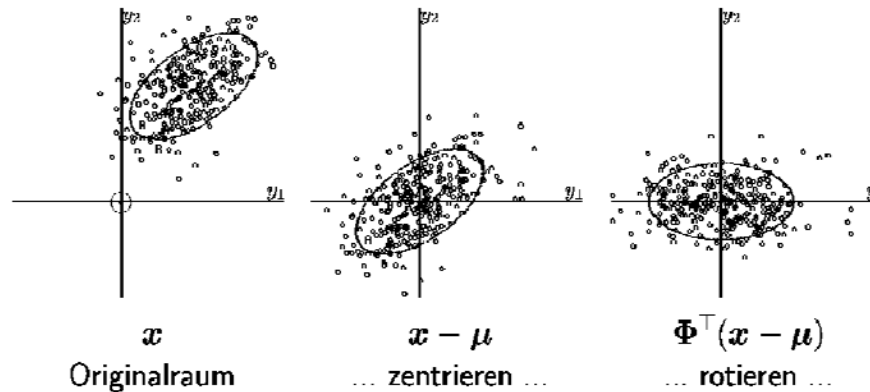
$$\varepsilon_d = \mathcal{E}[\|x - \hat{x}\|^2] = \mathcal{E}[\|\sum_{i=d+1}^D y_i \phi_i\|^2] = \sum_{i=d+1}^D \phi_i^\top \underbrace{\mathcal{E}[xx^\top]}_{:S} \phi_i$$

Schätzung der Momentenmatrix S (Lernstichprobe ω):

$$\hat{S} = \frac{1}{N} \sum_{j=1}^N x_j x_j^\top, \quad \omega = (x_1, \dots, x_N) \quad (S = \Sigma + \mu\mu^\top)$$

Fehlerminimierung \leadsto Eigenwertaufgabe $S\phi_i = \lambda_i \phi_i$

$$S = \Phi^\top \Lambda \Phi, \quad \Lambda = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_D \end{pmatrix}, \quad \lambda_i > 0$$



Die zentrierte Hauptachsen- / KL-transformation

- ... besitzt den Rekonstruktionsfehler $\varepsilon_d = \sum_{i=d+1}^D \lambda_i$
- ... hat die Eingabe decorreliert (y besitzt Kovarianzmatrix Λ)

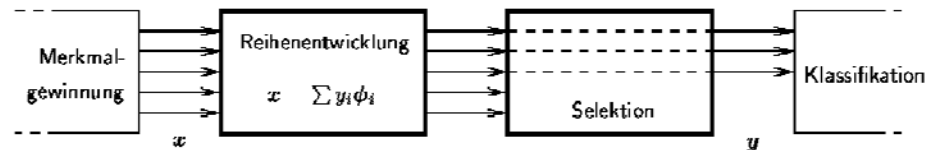
Merkmalraumtransformationen

- **Reduktion der Vektordimension**
 - ⇒ Geringerer Klassifikationsaufwand
 - ⇒ Robustere Schätzung statistischer Parameter
 - (*.. möglichst geringer Informationsverlust ...*)

Lineare Transformationsvorschrift

$$\varphi : \begin{cases} \mathbb{R}^D & \rightarrow \mathbb{R}^d \\ \mathbf{x} & \mapsto \mathbf{y} = \varphi(\mathbf{x}) = \Phi^\top \mathbf{x} \end{cases} \quad \text{mit } d \leq D$$

Typischerweise ist Φ eine *Rotationsmatrix*, d.h. $\Phi\Phi^\top = \mathbf{E}$



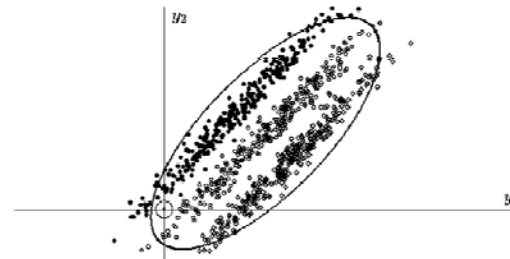
Gütekriterien für $\varphi(\cdot)$

- Varianzmaximierung \rightsquigarrow *Karhunen-Loève-Transformation*
- Klassenseparation \rightsquigarrow *Lineare Diskriminanzanalyse*
- MI-Kriterium \rightsquigarrow *Maximum-Likelihood-Rotationen*
- Fehlerrate \rightsquigarrow ???

Klassenbezogene Transformation

LINEARE DISKRIMINANZANALYSE

- die **Ballungsgebiete kompakt** halten
- ihre **Zentren** voneinander **entfernen**
- in eine der **Klassentrennung** förderliche Lage **drehen**



Zerlegung der Gesamtkovarianz $\Sigma = S_w + S_b$

$$S_w = \sum_{\kappa=1}^K p_{\kappa} \Sigma_{\kappa} \quad (\text{Intraklassen-Streuungsmatrix})$$

$$S_b = \sum_{\kappa=1}^K p_{\kappa} (\mu_{\kappa} - \mu)(\mu_{\kappa} - \mu)^{\top} \quad (\text{Interklassen-Streuungsmatrix})$$

Zentrierte LDA-Transformation

$$y = \Phi^{\top}(x - \mu) \quad \text{mit} \quad S_w^{-1} S_b \Phi = \Phi \Lambda$$

Translation — Drehung — Reskalierung

Gauß-Quantisierung des LDA-Cepstrums

Merkmalberechnung

Z.B. die Cepstralkoeffizienten & Ableitungen



Beschaffung einer Laut-Etikettierung

Per „bootstrap“ aus einem initialen Erkennungssystem



Berechnung der LDA-Transformationsmatrix Φ

Eigenwertanalyse der zerlegten Gesamtstreuung



Transformation der Lernstichprobe

Alle $\mathbf{x} \in \omega$ werden gemäß $\Phi : \mathbb{R}^{72} \rightarrow \mathbb{R}^{20}$ reduziert



Statistischer Vektorquantisierer-Entwurf

EM-Algorithmus \rightsquigarrow Kodebuch ($K = 256$) für ω^Φ



Quantisierung der Merkmalvektoren

Abbildung der $\mathbf{y} \in \mathbb{R}^{20}$ auf 8-bit-Klassenindizes

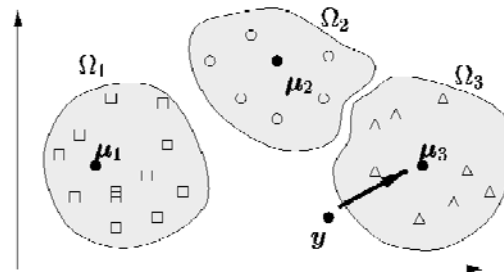
Klassifikation

Klassifikation D -dimensionaler Muster

- Merkmalvektoren $\boldsymbol{x} \in \mathbb{R}^D$
- Musterklassen $\Omega_1, \dots, \Omega_K$
- Etikettierte Lernstichprobe $\omega_1, \dots, \omega_K$
- Klassifikationsaufgabe

$$\boldsymbol{y} \in \mathbb{R}^D \rightsquigarrow \Omega_k \quad (k = ?)$$

⇒ ÜBERWACHTES LERNEN

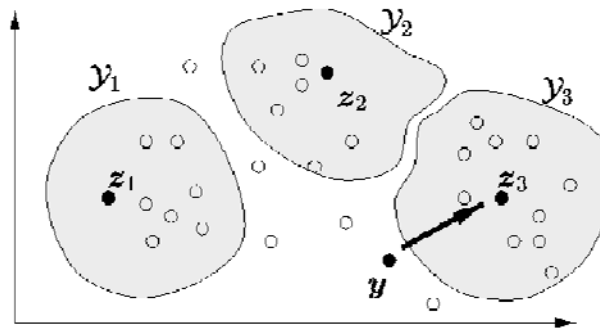


Vektorquantisierung

- Merkmalvektoren $\mathbf{x} \in \mathbb{R}^D$
- Repräsentantenvektoren $\mathbf{z}_1, \dots, \mathbf{z}_K$
- Nichtetikettierte Lernstichprobe $\omega \subseteq \mathbb{R}^D$
- Quantisierungsaufgabe

$$\mathbf{y} \in \mathbb{R}^D \rightsquigarrow \mathbf{z}_k \quad (k = ?)$$

⇒ UNÜBERWACHTES LERNEN



Numerische Klassifikatoren

GEGEBEN:

- D -dimensionale, reelle Merkmalvektoren

$$\mathbf{x} = (x_1, \dots, x_D) \in \mathbb{R}^D$$

- Inventar von K Musterklassen

$$\Omega_1, \Omega_2, \dots, \Omega_K$$

Statistisches Erzeugungsmodell:

- ⇒ Diskreter Zufallsprozeß
mit a priori Klassenwahrscheinlichkeiten

$$P(\Omega_\kappa), \quad \kappa = 1, \dots, K \quad \text{mit} \quad \sum_{\kappa=1}^K P(\Omega_\kappa) = 1$$

- ⇒ Multivariat-kontinuierlicher Prozeß
mit den bedingten Wahrscheinlichkeitsdichtefunktionen

$$P(\mathbf{x} | \Omega_\kappa) \quad \text{mit} \quad \int_{\mathbb{R}^D} P(\mathbf{y} | \Omega_\kappa) d\mathbf{y} = 1$$

Numerische Klassifikatoren

Optimale Entscheidungsregel

„Unschärfe“ Entscheidungsregel:

$$\delta(\Omega_\kappa | \mathbf{x}) \quad \text{mit} \quad \sum_{\kappa=1}^K \delta(\Omega_\kappa | \mathbf{x}) = 1 \quad \text{für alle } \mathbf{x} \in \mathbb{R}^D$$

Risikofunktion:

$$R(\delta) = \sum_{\kappa=1}^K P(\Omega_\kappa) \sum_{\lambda=1}^K r_{\kappa\lambda} \int_{\mathbb{R}^D} \delta(\Omega_\lambda | \mathbf{x}) P(\mathbf{x} | \Omega_\kappa) d\mathbf{x}$$

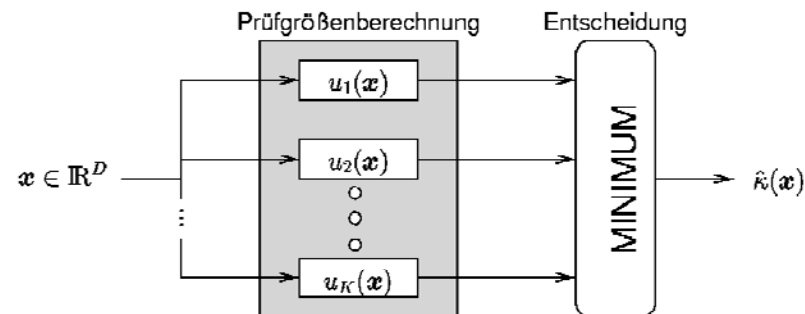
... mit den **Kosten** $r_{\kappa\lambda}$ (Verlust, Schaden) für die Verwechslung der Klassen $\Omega_\kappa, \Omega_\lambda$

Die **optimale Entscheidungsregel** ist deterministisch (!)

$$\delta^*(\Omega_\kappa | \mathbf{x}) = \begin{cases} 1 & \text{falls } u_\kappa(\mathbf{x}) = \min_\lambda u_\lambda(\mathbf{x}) \\ 0 & \text{sonst} \end{cases}$$

und verwendet die **Prüfgrößen** (Trennfunktionen)

$$u_\lambda(\mathbf{x}) = \sum_{\kappa=1}^K r_{\kappa\lambda} \cdot P(\Omega_\kappa) \cdot P(\mathbf{x} | \Omega_\kappa) \quad \text{für } \lambda = 1, \dots, K$$



Übernächste Woche:

1. Overview Speech Recognition Systems & Architectures
2. Acoustic modeling & feature extraction (1)
3. Feature extraction (2)
4. LDA, Grundlagen der Klassifikation
5. Klassifikation in HMM Modellen
6. Wortmodellierung (trigramme, tying)
7. search/decoding, lattices,
wordgraphs, confidence measures
8. acoustic adaptation
9. language models and grammars, Language model
adaptation, lexica, phonology
10. speech understanding, dialogue control
11. Design of computer speech recognition systems

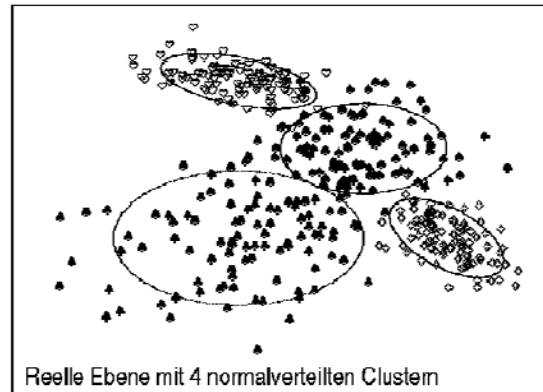
Statistische Clusteranalyse

ANNAIIME: die Vektoren \mathbf{x} sind **mischverteilt** gemäß

$$P(\mathbf{x} | \mathbf{p}, \boldsymbol{\theta}) = \sum_{\kappa=1}^K p_{\kappa} \cdot P(\mathbf{x} | \boldsymbol{\theta}_{\kappa})$$

Bemerkung 4.2 (*Identifizierbarkeit von Mischverteilungen*)

Für gemischte *Normalverteilungen* sind die Parameterwerte $\mathbf{p}, \boldsymbol{\theta}$ eindeutig bestimmbar, sofern der genaue Funktionsverlauf von $P(\mathbf{x} | \mathbf{p}, \boldsymbol{\theta})$ bekannt ist!



Maximum-Likelihood-Schätzung

$$\mathcal{L}_{\text{MIX}}(\mathbf{p}, \boldsymbol{\theta}) = \log \prod_{i=1}^N P(\mathbf{y}_i | \mathbf{p}, \boldsymbol{\theta}) = \sum_{i=1}^N \log \left(\sum_{\kappa=1}^K p_{\kappa} P(\mathbf{y}_i | \boldsymbol{\theta}_{\kappa}) \right)$$



Gekoppelt-transzendente Gleichungen für $\hat{p}_{\kappa}, \hat{\boldsymbol{\mu}}_{\kappa}, \hat{\boldsymbol{\Sigma}}_{\kappa}$

$$\hat{p}_{\kappa} = \frac{1}{N} \sum_{i=1}^N P(\Omega_{\kappa} | \mathbf{y}_i, \hat{\mathbf{p}}, \hat{\boldsymbol{\theta}})$$

Estimation - Maximization (EM-Algorithmus)