

HMM-Synthese (Grundzüge)

Uwe Reichel
Institut für Phonetik und Sprachverarbeitung
Ludwig-Maximilians-Universität München
reichelu@phonetik.uni-muenchen.de

9. Januar 2017

Inhalt

- HMM-Grundlagen
- HMM und Phonemerkennung
- HMM-Synthese

HMM-Grundlagen

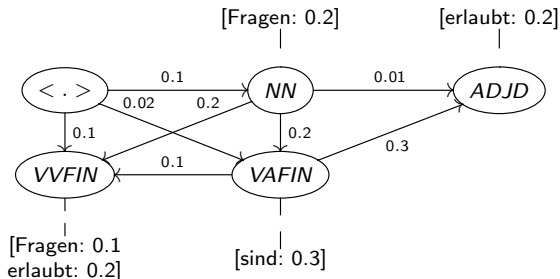
HMM: $\langle Q, K, S, A, B \rangle$

- $Q = \{q_i\}$: Menge von Zuständen
- K : Ausgabealphabet
- $S = \{s_i\}$: Startwahrscheinlichkeiten, dass man sich zu Beginn in Zustand q_i befindet
- $A = \{a_{ij}\}$ Transitionswahrscheinlichkeiten von Zustand q_i nach q_j
- $B = \{b_{jo}\}$ Emissionswahrscheinlichkeiten für Beobachtung o in Zustand q_j

HMM-Grundlagen

Symbolverarbeitung

- Beispiel Part-of-Speech-Tagging (vgl. POS-Folien)
- $B = \{b_{jo}\}$ sind **eindimensional**: $P(\text{Wort}|\text{POS})$
- Emission ist **kategorial** (hier: ein Wort)



HMM-Grundlagen

Signalverarbeitung

- Emissionswahrscheinlichkeiten $B = \{b_{jo}\}$ sind **mehrdimensional**: $P(e_1, e_2, \dots | q_j)$
- e_j : Signalcharakteristika (Dimensionen), z.B. Gesamtenergie, Mel-Cepstralkoeffizienten
- mögliche Modellierung als **mehrdimensionale Gaußglocken**
- Emissionen e_j sind **kontinuierlich**

HMM in der Phonemerkennung

HMM_x für Phonem x

- **Zustände** Q : Teilsegmente des Phonems x
- **Emissionswahrscheinlichkeiten** B : Wahrscheinlichkeiten für akustische Ausprägungen von x in den entsprechenden Teilsegmenten

HMM in der Phonemerkennung

Phonemerkennung

• Aufgabe:

- klassifiziere Sprachsignal s
- finde dasjenige HMM_x , durch das s mit der größten Wahrscheinlichkeit erzeugt wird

• Lösung:

- berechne für jedes HMM_x : $P(s|\text{HMM}_x)$ mittels **Viterbi** (vgl. *POS-Folien*)
- Zielphonem $\hat{x} = \arg \max_x P(s|\text{HMM}_x)$

HMM in der Sprachsynthese

Training

- **Full-Context-HMM**: je ein $HMM_{x,c}$ für ein Phonem x in einem bestimmten Kontext c
- c : *Phonemumgebung, prosodischer Kontext, Emotion, ...*
- Emissionen **kontinuierlich**, Emissionswahrscheinlichkeiten B **mehrdimensional**
 - Mel-Cepstral-Koeffizienten
 - F0
 - Segmentdauer
 - Intensität

HMM in der Sprachsynthese

- **Sparse-Data-Problem**
 - Trainingsmaterial reicht nicht aus, um für jedes Phonem x in jedem Kontext c die Parameter für $\text{HMM}_{x,c}$ verlässlich zu schätzen
 - Reduzierung der Anzahl unterschiedlicher Kontexte durch Clustering nach akustischer Ähnlichkeit
- Erstellung eines **Entscheidungsbaums**, der in Abhängigkeit von Phonem x und Kontext c das passende $\text{HMM}_{x,c}$ auswählt

HMM in der Sprachsynthese

Anwendung

- **Input:** segmentale, prosodische etc. **Zielspezifikationen** wie in Unit-Selection-Synthese
- Auswahl des passenden HMM mittels Entscheidungsbaum
- **kontextabhängige Phonemsequenz**
 $\langle x_1, c_1 \rangle \dots \langle x_n, c_n \rangle$:
Verkettung von $\text{HMM}_{x_1, c_1} \dots \text{HMM}_{x_n, c_n}$
- Generierung des Signals über den **wahrscheinlichsten Pfad** durch die HMM-Kette (mittels **Viterbi**)

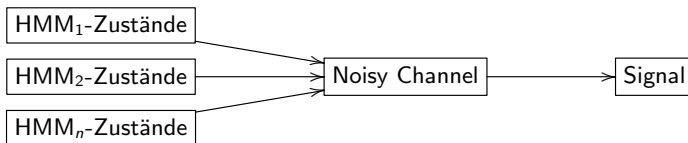
HMM in der Sprachsynthese

Gemeinsamkeiten zwischen Phonerkennung und Synthese

- Es werden **kontinuierliche Werte** emittiert.
- Emissionswahrscheinlichkeitsverteilungen sind **mehrdimensional**.

HMM in der Sprachsynthese

Unterschiede zwischen Phonemerkennung und Synthese



• Erkennung:

- ein HMM_x je **Phonem** x
- **bekannt: Kanalausgabe** (das akustische Signal)
- **gesucht: Kanalinput**, d.h. dasjenige HMM_x, das dem Signal am wahrscheinlichsten zugrundeliegt

HMM in der Sprachsynthese

- **Synthese:**
 - **Full-Context:** ein $\text{HMM}_{x,c}$ je **Phonem x und Kontext c**
 - **bekannt: Kanalinput**, d.h. das zu den segmentalen, prosodischen etc. **Zielspezifikationen** passende $\text{HMM}_{x,c}$
 - **gesucht: Kanalausgabe**, d.h. das durch $\text{HMM}_{x,c}$ am wahrscheinlichsten generierte akustische Signal

HMM in der Sprachsynthese

Modularisierung

- z.B. **Trennung von Quelle und Filter**
 - **Training** von $\text{HMM}_{q,c}$ und $\text{HMM}_{f,c}$
 - q Quelle: Parameter des Anregungssignals
 - f Filter: Filterparameter
 - **Anwendung:** getrennte Steuerung von Quell- und Filterparametern
 - Erhöhung der Flexibilität

HMM in der Sprachsynthese

Vorzüge

- **geringer Speicheraufwand:** HMM-Parameter statt Datenbank mit Sprachsignalen
- **hohe Flexibilität:**
 - kontinuierliche Steuerung der akustischen Parameter
 - getrennte Modellierung von Quelle und Filter
 - Generierung neuer Stimmen

Nachteile

- **Vollsynthese** → derzeit noch schlechtere Qualität als bei Unit-Selection