

Sprachsynthese: Graphem-Phonem-Konvertierung

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

8. Januar 2021

Inhalt

- Einflussfaktoren
- Alinierung
- Konvertierung
 - Table Lookup with Defaults (van den Bosch et al., 1993)
 - Maschinelles Lernen
 - Entscheidungsbäume
- Silbifizierung
 - in Graphemfolge
 - in Phonemfolge
- Wortbetonung

Einflussfaktoren

Notation: Grapheme: in spitzen Klammern; Phoneme' in Schrägstriche nach SAM-PA;
Graphem/Phonem/Wortvariablen: kursive Kleinbuchstaben; Graphemklassen (Vokal): kursive Großbuchstaben (V);

__#_σ] = 'vor Silbengrenze'

- Ist die Abbildung $\langle \text{Graphem} \rangle \rightarrow / \text{Phonem} /$ eindeutig?

- Nein.** Beispiel $\langle s \rangle$:

$\langle s \rangle \rightarrow /s/$ (*was*)

$\langle s \rangle \rightarrow /z/$ (*Vase*)

$\langle s \rangle \rightarrow /S/$ (*stehen*)

$\langle s \rangle \rightarrow / _ /$ (*Wasser*)

- weiteres Beispiel: $\langle u \rangle$ in $\langle \text{Bund} \rangle$ und $\langle \text{Quelle} \rangle$

Einflussfaktoren

Graphemkontext

$\langle s \rangle \rightarrow /z/ \mid V_V$; **aber:** *losen vs. Loserwerb*

$\langle s \rangle \rightarrow /S/ \mid _ \langle t \rangle$; **aber:** *Stabilität vs. Rost*

...reicht nicht aus

Silbenstruktur

- Auslautverhärtung, keine Beeinflussung durch Graphemumgebung über Silbengrenzen hinweg
- $\langle g \rangle \rightarrow /g/ \mid k/$: *Weg*e, *Weg*
- $\langle s \rangle \rightarrow /z/ \mid s/ \mid S/$: *Vase*, *Häuschen*

Einflussfaktoren

Morphologie

- morphologischer Einfluß direkt und über die Silbenstruktur manifestiert
- **direkt:** Phonem-Identität abhängig von Morphemklasse
- **Beispiele:**
 - <er> in *Erlöser* (/ʔE6l2:z6/); morph. $er_{prefix} + lös_{verb} + er_{suffix}$.
 <er> → /ʔE6/ | im Präfix
 <er> → /6/ | im Suffix
 - <e> in *geben* (/ge:b@n/); morph. $geb_{verb} + en_{infl}$.
 <e> → /e:/ | im einsilbigen Verbstamm
 <e> → /@/ | in Flektionsendung

Einflussfaktoren

- **indirekt:** morphologische Struktur bestimmt Silbenstruktur und damit Phonem-Identität
- **Beispiele:**
 - $\langle ng \rangle$ in *Angel* ($/\text{?}aN@l/\text{)}$ vs. *Angelegenheit* ($/\text{?}ang@le:g@nhalt/\text{)}$:
 morph. $angel_{noun}$ vs. $an_{prefix} + ge_{prefix} + leg_{verb} + en_{suffix} + heit_{suffix}$.
 $\langle ng \rangle \rightarrow /N/$ | keine Silbengrenze
 $\langle ng \rangle \rightarrow /n.g/$ | (hier: morphologisch bedingte) Silbengrenze
 - $\langle se \rangle$ in *losen* ($/lo:z@n/\text{)}$ vs. *Losentscheid* ($/lo:s\text{?}EntSalt/\text{)}$:
 morph. $los_{verb} + en_{infl}$ vs. $los_{noun} + ent_{prefix} + scheid_{verb}$.
 $\langle se \rangle \rightarrow /z@/$ | keine Silbengrenze
 $\langle se \rangle \rightarrow /s.\text{?}E/$ | (hier: morphologisch bedingte) Silbengrenze

Alinierung: Grundlagen

- Zur Gewinnung des G2P-Trainingsmaterials
- Zuordnung zusammengehöriger Abschnitte in Graphem- und Phonemsequenzen

O	b	e	r	s	c	h	u	l	z	e	u	g	n	i	s
?+o:	b	6	__	S	__	__	u:	l	t+s	OY	__	k	n	l	s

- Alinierungsproblem formuliert als **Minimierung der Distanz zwischen Graphem- und Phonemsequenz**

Alinierung: Levenshtein-Distanz

Levenshtein-Distanz

- **Minimal nötige Editierkosten** um Sequenz v (Grapheme) in Sequenz w (Phoneme) umzuwandeln
- **Standard-Editieroperationen:**
 - **Substitution** von v_i durch w_j : $\langle u \rangle \rightarrow /u:/$
 - **Löschung** von v_i : $\langle r \rangle \rightarrow _$
 - **Einfügung** von w_j : $_ \rightarrow /?/$

Alinierung: Kostenfunktion

Naiv

- **Substitution**

$$c(v_i, w_j) = \begin{cases} 0 & : v_i == w_j \\ 1 & : \text{else.} \end{cases} \quad (1)$$

- **Löschung, Einfügung**

$$c(v_i, _) = 1 \quad (2)$$

$$c(_, w_j) = 1 \quad (3)$$

- brauchbar für Wortvergleiche (z.B. automatische Rechtschreibkorrektur)

- unbrauchbar für Graphem-Phonem-Alignment: $\langle x \rangle \neq /x/$

Konvertierung: Table Lookup with Defaults

Table Lookup with Defaults (van den Bosch et al., 1993)

- **Training:** Speicherung des jeweils kürzesten Graphemkontexts für ein eindeutiges Graphem-Phonem-Mapping in Tabelle G

k a m m	/a/
u ß	/u:/

- **Sortierung** nach Länge der Graphemsequenz
- **2 Default-Tabellen:** Graphem-Fenster + am häufigsten damit ko-okkurrierendes Phonem

$\langle v_{i-1} \rangle \langle v_i \rangle \langle v_{i+1} \rangle$	/w _j /
$\langle v_i \rangle$	/w _j /

Konvertierung: Table Lookup with Defaults

- **Konvertierung:**
 - Suche nach passendem Graphem-Muster (von lang nach kurz) in Tabelle *G*
 - Falls nicht vorhanden, Rückgriff auf Default-Tabellen
 - Beispiel:
 - zu konvertieren: $\langle u \rangle$ in *Fuß*
 - in Tabelle *G* gefundenes Muster: *uβ*
- Ausgabe: */u:/*
- **Vorzüge**
 - rein datenbasierter Ansatz
- kein Expertenwissen nötig, sprachunabhängig

Maschinelles Lernen

- **Ziel:** Erlernen des Zusammenhangs zwischen Zielwerten (Kategorien oder kontinuierliche Werte) für Objekte und deren Attributen.
- bezogen auf Graphem-Phonem-Konvertierung
 - **Objekte:** Grapheme
 - **Attribute:** Graphem-Identität, umgebende Grapheme, Position des Graphems innerhalb der Silbe, ...
 - **Zielwerte:** Phoneme

Maschinelles Lernen

- Objekte als **Attributbündel** (*Feature vectors*) repräsentiert.
Beispiel:
 - **Attribute** für Graphem v_i :
[$\langle v_{i-1} \rangle, \langle v_i \rangle, \langle v_{i+1} \rangle, \text{Morphemtyp}, _ \#_\sigma$]
 - Mögliche Attributwerte: [a–z, a–z, a–z, *frei|gebunden*, 0|1]
(*gebundenes* Morphem kann allein kein Wort bilden)
 - Feature vector für erstes $\langle e \rangle$ in *geben*:
[*g, e, b, gebunden, 1*]
 - **Zielwert (hier kategorial)**: /e:/

Maschinelles Lernen

- **Attribute:** kategorial oder kontinuierlich
 - **kategorial:** Graphem-Identität, Position in Silbe, Phonemklasse, Wortbetonung, Morphstatus
 - **kontinuierlich:** relative Position des Graphems im Wort, Lautdauer, F0-Wert

Maschinelles Lernen

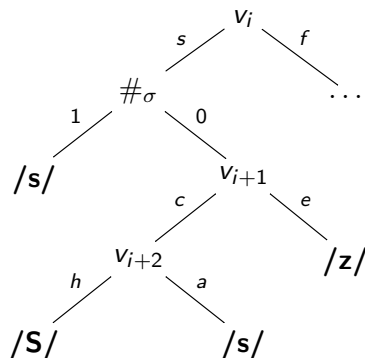
- **Überwachtes Lernen:** Zielwerte in Trainingsdaten bekannt; C4.5, CART, neuronale Netze (ANN)
- **Unüberwachtes Lernen:** Zielwerte nicht bekannt; Clustering, ANN
- **Versch. Methoden für Attribut- → Zielwert-Typen:**
 - C4.5: kategorial/kontinuierlich → kategorial; z.B. Akzent
 - CART: kategorial/kontinuierlich → kategorial/kontinuierlich; z.B. Lautdauer
 - ANN: kontinuierlich → kategorial/kontinuierlich

Entscheidungsbäume

- Quinlan (1993); <http://www.cse.unsw.edu.au/~quinlan>
- **Modellierung:**
 - **Objekte:** Pfade durch den Baum
 - **Attribute:** Kanten zwischen Knoten
 - **Zielwerte:** terminale Knoten
- Vorteil: Transparenz → Wissensakquirierung möglich

Entscheidungsbäume

Beispielausschnitt:



Entscheidungsbäume

Rekursiver Aufbau des Baums

- **Fall 1:** Ordne alle Objekte, die noch nicht durch einen vollständigen Pfad im Baum repräsentiert sind und den **gleichen Zielwert** haben, einem neuen Blatt (terminaler Knoten) zu.
- **Fall 2:** Genauso, wenn die Objekte zwar **verschiedene Zielwerte** haben, sich **aber** anhand der gegebenen Attribute **nicht mehr weiter unterscheiden** lassen.
- **Fall 3:** Haben Objekte **verschiedene Zielwerte und unterscheiden** sie sich in einer oder mehreren Attributen, so wähle das **zur Partitionierung der Objektmenge am 'besten geeignete' Attribut** und erzeuge einen Knoten, an dem sich der Baum in mehrere durch Werte des betrachteten Attributs vorgegebene Kanten aufspaltet.

Entscheidungsäume

G2P-Anwendung von Entscheidungsäumen

- **Attribute:**
 - Graphemkontext (n Grapheme vorher/nachher)
 - Morphologie: Morphemklasse, +/- folgende Morphemgrenze
 - Silben: Aufbau der Silbe Onset/kein Onset (bedeckt/nackt), Coda/kein Coda (geschlossen/offen), Position in Silbe (Onset, Nukleus, Coda, Gelenk *alle*)
 - Phonem: Vorgeschichte (n Phoneme vorher)
- **Zielwert:** Phonem, incl. leeres Phonem ($/_/\$), oder Phonem-Cluster ($/\?+a:/$)

Morphologische Zerlegung

Voraussetzung: Morphem-Lexikon mit Morphemklassifizierung

- 1 Teile jedes Wort w *rekursiv* von links nach rechts in String-Präfixe und -Suffixe bis eine erlaubte Segmentierung möglich ist oder das Wortende erreicht wird.
- 2 Eine Grenze, die den aktuellen String in Präfix und Suffix unterteilt dann akzeptiert wenn (i) das Präfix im Lexikon zu finden ist, und (ii) eine erlaubte Segmentierung des Suffixes möglich ist, und (iii) die Kombination 'Präfix-Klasse + Klasse des ersten Suffix-Teils' nicht der Morphotaktik widerspricht, und (iv) die Klasse des letzten Suffix-Teils kompatibel ist mit dem POS von w .

Morphologische Zerlegung

Beispiel: *Fassade – nkletterer*

- (i) String-Präfix *Fassade* ist im Lexikon
 - (ii) **Rekursion:** erlaubte Segmentierung des String-Suffixes *nkletterer* möglich (*n-kletter-er*)
 - (iii) Morphemklassen *Fassade/NN – n/Fugenmorphem* kompatibel
 - (iv) Morphemklasse des letzten String-Suffixes *er/NN-Suffix* kompatibel mit POS NN des Words
- Segmentierung *Fassade – nkletterer* möglich

Silbifizierung

In Graphemfolge

- **3 vorherzusagende Klassen:**
Silbengrenze folgt, folgt nicht, Ambisyllabizität
- **Feature-Auswahl:**
 - für jedes Graphem $\langle v_i \rangle$
 - innerhalb eines auf $\langle v_i \rangle$ zentrierten symmetrischen Graphem-Fensters
 - Graphem, Konsonant/ Vokal
 - ggf. Morphemgrenze (+/- relevant für Silbengrenze)

Silbifizierung

- für Silbifizierung **relevante morphologische Grenzen**:
 - vor allen Morphemen außer Flexionsendungen, Suffixen, Komparationsmorphemen und Fugen
 - vor Flexionsendungen, Suffixen mit initialem Konsonanten und eigenem Silbenkern (*schaffte*)
 - vor Flexionsendungen, Suffixen mit initialem Vokal, wenn das vorangehende Morphem auf Vokal endet (*bauen*)

Erweiterungen

Folgende Folien sind *kein Prüfungstoff* bis Abschnitt
"Evaluierung".

Nach abgeschlossener G2P Konvertierung Erweiterung durch:

- Phonetische Sibifizierung
- Wortbetonung (*lexical accent*)

Silbifizierung

In Phonemfolge

- 1 setze vor jedes Sonoritätsminimum eine Silbengrenze
- 2 Feinadjustierung gemäß Kohlers (1995) Silbenphonotaktik und silbengrenzrelevanter Morphemgrenzen

Beispiel:

/fE6hEltnls/ $\xrightarrow{1.}$ /fE6.hEl.tnls/ $\xrightarrow{2.}$ /fE6.hElt.nls/

Phonetische Silbifizierung

Silbenphonotaktik (Kohler, 1995)

$$\left(\left\{ \begin{array}{ccc} & K_{a,b,c} & \\ (K_a) & K_a & K_b \\ & K_a & K_c \\ (K_a) & K_a & K_a \end{array} \right\} \vee \left\{ \begin{array}{ccc} & K_{a,b} & \\ K_b & K_a & (K_a) \\ K_b & K_b & (K_a) \\ K_a & K_a & \end{array} \right\} \right) \left(\left\{ \begin{array}{c} K_a(+K_a) \\ +K_a(K_a) \end{array} \right\} \right)$$

- K_a : Plosive, Frikative
- K_b : Nasale, /l/, /r/
- K_c : /h/, /j/
- V : Vokale
- $+$: Morphemgrenze

Phonetische Silbifizierung

Silbenbeispiele:

- *Herbsts* /hE6psts/ $K_c VK_a K_a + K_a$
- *Psalm* /psalm/ $K_a K_a VK_b K_b$

Restriktionen gegen Übergeneralisierung

- **Beispiel:** für $(K_a)K_a K_b$ darf K_a nicht im Artikulationsort mit K_b übereinstimmen,
vgl. /fE6.hEl.tnls/ \rightarrow /fE6.hEl.t.nls/

Wortbetonung

Simplex-Wörter: Restriktionen

- **Drei-Silben-Fenster:** Wortbetonung kann nur auf eine der letzten drei Silben im Wort fallen (Ultima, Penultima, Antepenultima)
Ausnahme: *schwa* in Antepenultima (/ 'a:.b@n.tOY.6 /)
- **Closed penult:** geschlossene Penultima verhindern eine Betonung weiter links (/ hi.b'ls.kUs /)
- **Final schwa:** Betonung der Penultima, wenn der Silbenkern der Ultima aus einem Schwa besteht: / tsi.tr'o:.n@ /
Produktivität: / g'e:.nE.zls / → / gE.n'e:.z@ /

Wortbetonung

Allgemeine Tendenzen im Standard-Deutschen

- Wortbetonung eher hinten im Wort
- eher auf *schweren* Silben (Langvokal/ Diphthong, Coda)
- *schwa*-Silben nicht betonbar

Wortbetonung

Simplex+Affixe

- **betonte Morphemklassen:**
 - Verbpartikeln: **w'egfahren**
 - betonte Affixe: *Abstin'enz*, *pass'abel*
- **unbetonte Morphemklassen:**
 - Flektionsendungen
 - unbetonte Affixe: **entspr'echen**, *M'annschaft*, *'Ärgemis*

Komposita

- häufig Haupt- und Nebenbetonung
- **zweigliedrig:** Hauptbetonung auf erstem Glied (wenige Ausnahmen: *Lebew'ohl*)

Wortbetonung

- **mehrgliedrig:** Anwendung der **Compound-stress-Rule (CSR)** der metrischen Phonologie:

CSR: *im Kompositum AB ist B strong s, wenn es sich weiter verzweigt, ansonsten ist A strong und B weak w*



- Erklärung von **Ausnahmen:** Atomisierung lexikalisierter verzweigender Konstituenten: *K'unst#[hand#werk]*, daher: drei- → zweigliedrig → Hauptbetonung auf erstem Glied

Wortbetonung

weitere Schwierigkeiten

- **stress shift:** *D'oktor* > *Dokt'oren* (vgl. *Final schwa*-Restriktion); *'Ablauf* > *Progr'amm#abl"auf*
- **Homographen** unterschiedlicher Wortart, Valenz
 - *mod'ern/ADJ vs. m'odern/V*
 - *K'onstanz/NE vs. Konst'anz/NN*
 - *'allerhand/Indefpron,ADV vs. allerh'and/ADJD*
 - *d'arüber/PAV vs. dar'über/PTKVZ*
 - *d'urchlaufen/V(intrans) vs. durchl'aufen/V(trans)*

Wortbetonung

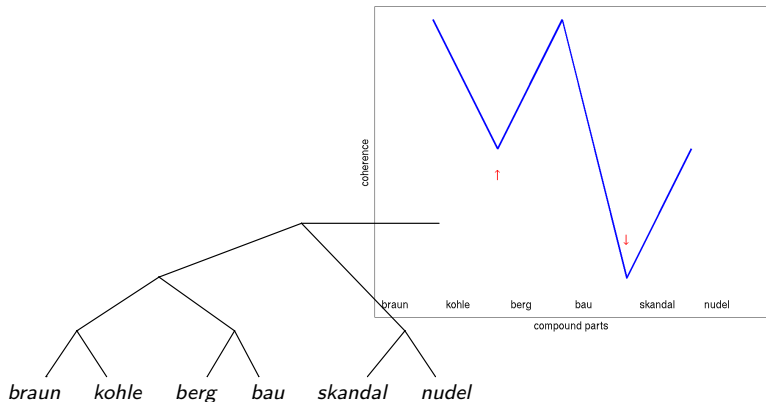
- **Kontrastakzent:**
 - *'Arbeitgeber vs. Arbeitg'eber und Arbeitn'ehmer*
 - *R'egel vs. ich habe Reg'el gesagt, nicht Regal*
- **unklare Fälle:** *m'erkwürdigerweise vs. merkwüdig'erw'eise*

Wortbetonung

Vorhersage

- 1 Kompositumzerlegung (mittels morphologischer Analyse, s.o.)
- 2 Bestimmung des betonten Kompositumglieds (mittels metrischer Bäume)
- 3 Lokalisierung der betonten Silbe im betonten Kompositumglied (mittels morphologischer Analyse zur Identifizierung betonter Affixe und maschinellem Lernen)

Betontes Kompositumglied: Induktion metrischer Bäume



Kohärenz: z.B. Bigramm-Wahrscheinlichkeiten: $P(\text{skandal}|\text{bau}) < P(\text{bau}|\text{berg})$

Wortbetonung: Lokalisierung der betonten Silbe

Mittels maschinellem Lernen

- **Objekte:** Silben oder Wörter (Kompositumglieder)
- **Zielwerte:**
 - für Silben: +/- betont
 - für Wörter: Index der betonten Silbe

Wortbetonung: Lokalisierung der betonten Silbe

Instanzbasiertes Lernen (Mustervergleich; Daelemans et al., 1994)

- **Training:** Abspeichern von Wörtern in Form von Merkmalsvektoren zusammen mit Index der betonten Silbe
- **Merkmale:** z.B. Silbengewicht “schwer, leicht” der letzten beiden Silben; **Zielwerte:** ultima, penultima, antepenultima

Wort	Merkmalsvektor	Betonung
Wanne	s l	p
Sonne	s l	p
Hibiskus	s s	p
genau	l s	u

Wortbetonung: Lokalisierung der betonten Silbe

- **Anwendung** – für Wort w : Übernahme des häufigsten Betonungsmusters unter den w -ähnlichsten Wörtern
- **Hamming-Distanz d** zweier Merkmalsvektoren: Anzahl der unterschiedlichen Werte

$w = \text{Wonne}$ (Merkmalsvektor: $[s \]]$)

$d(\text{Wonne}, \text{Wanne})=0$, $d(\text{Wonne}, \text{Sonne})=0$,

$d(\text{Wonne}, \text{Hibiskus})=1$, $d(\text{Wonne}, \text{genau})=2$

→ Übernahme der Betonung der ähnlichsten Wörter
Wanne, Sonne, also penultima

Wortbetonung: Lokalisierung der betonten Silbe

Entscheidungsbaum

- Objekte: Silben
- Features: Silbengewicht, Morphemklasse, POS, Silbenindex, Position in Kompositum
- Zielwert: +/- betont

Evaluierung

- anhand eines Testcorpus mit manueller Transkription (*gold standard*)
- **Word error rate**: Anteil der Wörter, in denen (mindestens) eine Abweichung vom *gold standard* auftritt
- **Phone error rate**: ermittelt über die **Levenshtein-Distanz** zwischen gewünschter (*gold standard*) und berechneter Transkription
- sehr stark abhängig von Sprache / Testkorpus!