



Grammatikmodelle

Vorlesungsunterlagen

Speech Communication 2, SS 2004

Franz Pernkopf/Erhard Rank

Institute of Signal Processing and Speech Communication

University of Technology Graz

Inffeldgasse 16c , 8010 Graz, Austria

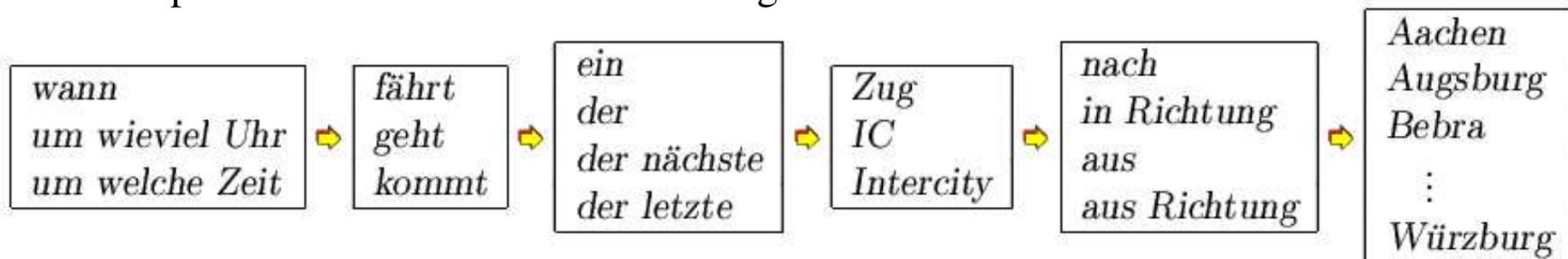
Tel: +43 316 873 4431

E-Mail: pernkopf/rank@inw.tugraz.at

Typische Fehlentscheidungen bei der automatischen Spracherkennung:

(1)	„den <u>nächsten</u> Zug“	↷	„den <u>nächste</u> Zug“	Kongruenzfehler
(2)	„um sechs <u>in</u> Bonn“	↷	„um sechs <u>den</u> Bonn“	Vertauschung, Auslassung oder Einfügung von Funktionswörtern
(3)	„von <u>Essen</u> nach Kiel“	↷	„von <u>ist in</u> nach Kiel“	Vertauschung von Inhaltswörtern durch ein oder mehr Funktionswörter
(4)	„nach <u>Köln</u> fahren“	↷	„nach <u>Ulm</u> fahren“	grammatisch unauffällige Vertauschungen

Satzbauplan für einfache Benutzeräußerungen am Bahnhof



Bayes-Regel

$$w^* = \operatorname{argmax}_{w \in W^*} P(w | \mathbf{X})$$

mit den *a posteriori* Wahrscheinlichkeiten

$$P(w | \mathbf{X}) = \frac{\overbrace{P(\mathbf{X} | w)}^{\text{ASM}} \cdot \overbrace{P(w)}^{\text{LSM}}}{P(\mathbf{X})}$$

GESUCHT ist ein Wahrscheinlichkeitsmodell

$$P(w_1 w_2 \dots w_T)$$

für Wortfolgen $w = w_1 \dots w_T$ der Länge T

Formel für bedingte Wahrscheinlichkeiten

$$P(w_1 w_2 \dots w_T) = P(w_T | w_1 w_2 \dots w_{T-1}) \cdot P(w_1 w_2 \dots w_{T-1})$$



Kettenregel

$$\begin{aligned} P(w_1 w_2 \dots w_T) &= P(w_1) \\ &\quad \cdot P(w_2 \mid w_1) \\ &\quad \cdot P(w_3 \mid w_1 w_2) \\ &\quad \cdot P(w_4 \mid w_1 w_2 w_3) \\ &\quad \cdot \dots \dots \dots \\ &\quad \cdot P(w_T \mid w_1 \dots w_{T-1}) \\ &= \prod_{t=1}^T P(w_t \mid w_1 \dots w_{t-1}) \end{aligned}$$

- $P(w_t \mid w_1 \dots w_{t-1})$ ist die **bedingte** Wortwahrscheinlichkeit
- die Wortkette $w_1 \dots w_{t-1}$ heißt **Vergangenheit** oder **Kontext** des aktuellen Wortes w_t

Modelle mit beschränktem Kontext

$$P(w_t | w_1 \dots w_{t-1}) \approx P(w_t | \underbrace{w_{t-n+1} \dots w_{t-1}}_{(n-1)\text{-Gramm-Kontext}})$$

Unigramm-, Bigramm- und Trigramm-Modelle

$$P_{1G}(\mathbf{w}) = \prod_{t=1}^T P(w_t)$$

$$P_{2G}(\mathbf{w}) = P(w_1) \cdot \prod_{t=2}^T P(w_t | w_{t-1})$$

$$P_{3G}(\mathbf{w}) = P(w_1) \cdot P(w_2 | w_1) \cdot \prod_{t=3}^T P(w_t | w_{t-2}w_{t-1})$$

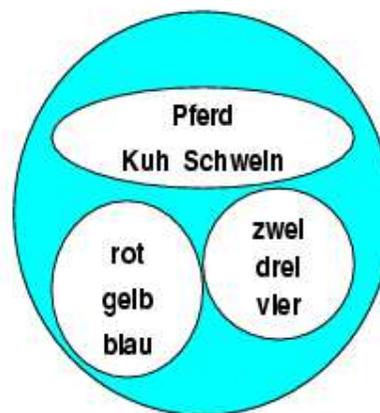
n -Gramm-Sprachmodelle

$$P_{nG}(\mathbf{w}) = \prod_{t=1}^T P(w_t | w_{t-n+1} \dots w_{t-1})$$

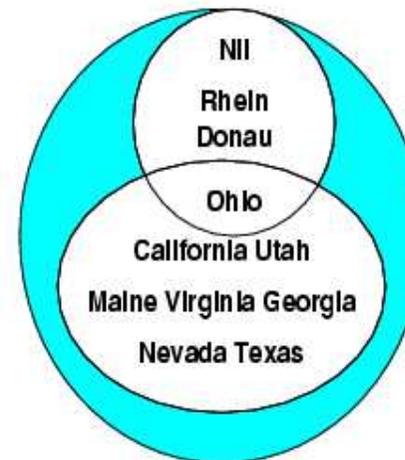
Gruppierung von Wörtern mit ähnlichem grammatischen u/o statistischen Verhalten

- **syntaktisch:**
Nomina, Verben, Adjektive von gleichem Genus, Kasus, Tempus, ...
- **semantisch:**
Ordinal- oder Kardinalzahlwörter;
Wörter, die Farbe, Größe, Temperatur, ... charakterisieren;
- **pragmatisch:**
Eigennamen für Personen, Städte, Straßen, ... ;
Nomina für Tiere, Pflanzen, Materialien, klinische Diagnosen, ...

Repräsentation von Wortassoziationen



disjunkte Kategorien



überlappende Kategorien

Maximum-Likelihood-Schätzung:

$$\hat{P}(v^i|\mathbf{v}) = \frac{\#(\mathbf{v}v^i)}{\#(\mathbf{v})} = \frac{\text{„Anzahl der } n\text{-Gramme } \mathbf{v}v^i\text{“}}{\text{„Anzahl der } (n-1)\text{-Gramme } \mathbf{v}\text{“}}$$

Problem beschränkten Stichprobenumfangs

- Wenn $\#(\mathbf{v}v^i) = 0$, so ist auch $\hat{P}(v^i|\mathbf{v}) = 0$
- Wenn $\#(\mathbf{v}) = 0$, so ist $\hat{P}(v^i|\mathbf{v})$ undefiniert!
- Insbesondere ist $\hat{P}(w) = 0$ für „neue“ oder „ungesehene“ Wörter



Glättung der Schätzwerte \hat{P}

Verfärbung der Zählfunktion

$$\hat{p}_k^{\text{ML}} = \frac{N_k}{T} \quad \Rightarrow \quad \hat{p}_k = \frac{N_k^*}{T^*} = \frac{N_k^*}{L + \sum_{j=1}^L N_j^*}$$

Sobald $N_k^* > 0$ für alle k , so gilt auch $\hat{p}_k > 0$.

- **Laplace- oder uniforme Bayes-Schätzung**

$$N_k^* := N_k + 1 \quad \Rightarrow \quad \hat{p}_k = \frac{N_k + 1}{T + L}$$

- **Jeffrey-Glättung**

$$N_k^* := N_k + 1/2 \quad \Rightarrow \quad \hat{p}_k = \frac{N_k + 1/2}{T + L/2}$$

- **Quadratmittelmethode** (minimiert Fehlschätzungsrisiko)

$$N_k^* := N_k + 1/2\sqrt{T}$$

Beispiel 7.1 ($\mathbf{p} = (0.5, 0.4, 0.1)$ und $\mathbf{N} = (4, 3, 0)$)

- ML-Schätzung — $\hat{\mathbf{p}} = (0.57, 0.43, 0)$
- Laplace-Schätzung — $\hat{\mathbf{p}} = (0.5, 0.4, 0.1)$

Informationen in Sprachmodellen *niedrigerer Ordnung*:

$$\hat{P}(w_3 | w_1 w_2) \Rightarrow \hat{P}(w_3 | w_2) \Rightarrow \hat{P}(w_3) \Rightarrow 1/L$$

Allgemeine Rückfallstrategie (Rekursionsgleichung)

$$\tilde{P}(w|\mathbf{v}) = \begin{cases} \hat{P}(w|\mathbf{v}) & \text{falls } \#(\mathbf{v}w) > 0 \\ \beta(\mathbf{v}) \cdot \tilde{P}(w|\mathbf{v}') & \text{falls } \#(\mathbf{v}w) = 0 \end{cases}$$

- $\hat{P}(w|\mathbf{v})$ ist ein Schätzer für $P(w|\mathbf{v})$ mit nichtverschwindenden Wahrscheinlichkeiten.
- \mathbf{v}' ist der abgemagerte Bedingungsteil des größeren Modells.

Interpolierte n -Gramme

$$\begin{aligned}\tilde{P}^{(n)}(w_t \mid w_1 \dots w_{t-1}) &= \lambda_0^{(n)} \cdot \frac{1}{L} + \lambda_1^{(n)} \cdot \hat{P}(w_t) \\ &+ \lambda_2^{(n)} \cdot \hat{P}(w_t \mid w_{t-1}) \\ &+ \sum_{i=3}^n \lambda_i^{(n)} \cdot \hat{P}(w_t \mid w_{t-i+1} \dots w_{t-1})\end{aligned}$$

Interpolierte n -Gramme

$$\begin{aligned} \tilde{P}^{(n)}(w_t \mid w_1 \dots w_{t-1}) &= \lambda_0^{(n)} \cdot \frac{1}{L} + \lambda_1^{(n)} \cdot \hat{P}(w_t) \\ &+ \lambda_2^{(n)} \cdot \hat{P}(w_t \mid w_{t-1}) \\ &+ \sum_{i=3}^n \lambda_i^{(n)} \cdot \hat{P}(w_t \mid w_{t-i+1} \dots w_{t-1}) \end{aligned}$$

- Falls $t < n$, fülle w_{t-n+1}, \dots, w_0 mit \$-Marken auf:

$$\tilde{P}^{(n)}(\mathbf{w}) = \prod_{t=1}^T \tilde{P}^{(n)}(w_t \mid w_1 \dots w_{t-1})$$

- Falls $t < n$, verwende $\tilde{P}^{(t)}$ oder $\tilde{P}^{(t+1)}$ statt $\tilde{P}^{(n)}$:

$$\tilde{P}^{(n)}(\mathbf{w}) = \prod_{t=1}^n \tilde{P}^{(t)}(w_t \mid w_1 \dots w_{t-1}) \cdot \prod_{t=n+1}^T \tilde{P}^{(n)}(w_t \mid w_1 \dots w_{t-1})$$

- Im Falle $\#(\mathbf{w}_{t-i+1}^{t-1}) = 0$ wird $\hat{P}(w_t \mid \mathbf{w}_{t-i+1}^{t-1}) = 1/L$ substituiert.