

# Implementierung eines POS-Taggers mithilfe von Haskell

Andy Püschel  
Till Köhler

Institut für Theoretische Informatik  
Fakultät Informatik  
Technische Universität Dresden

19.01.2018

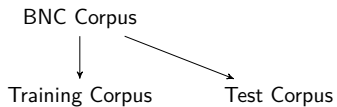
# Outline

- 1 Aufgabenstellung
- 2 Modellierung von Hidden-Markov-Modellen
- 3 Supervised Training mit Smoothing
- 4 Annotation mit Handling of Unknown Words
- 5 Evaluation des Modells
- 6 Ergebnisse

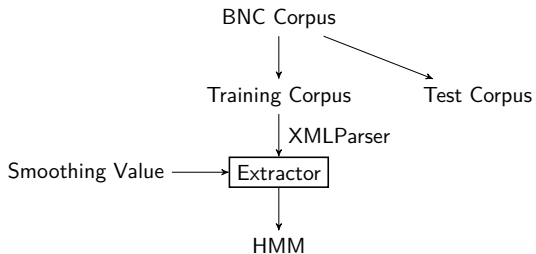
# Aufgabenstellung

- Spezifikation einer Haskell-Datenstruktur zur Darstellung von Hidden-Markov-Modellen
- Supervised Training eines Hidden-Markov-Modells aus einem mit POS-Tags annotierten Textkorpus
- Serialisierung von Hidden-Markov-Modellen aus der Haskell-Datenstruktur in menschenlesbares Format
- Deserialisierung zurück in die Haskell-Datenstruktur
- Annotation eines beliebigen Textkorpus mit POS-Tags
- Evaluation des gewonnenen Hidden-Markov-Modells

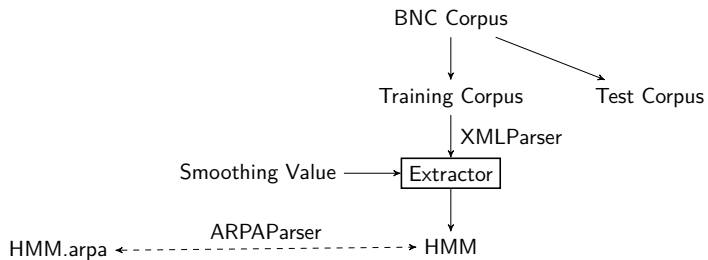
# Schematische Darstellung



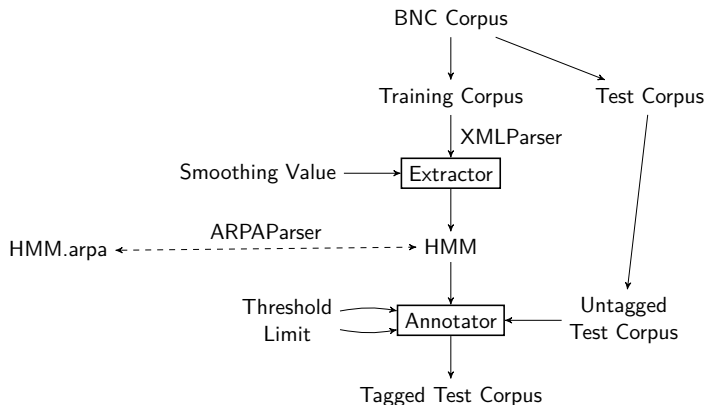
# Schematische Darstellung



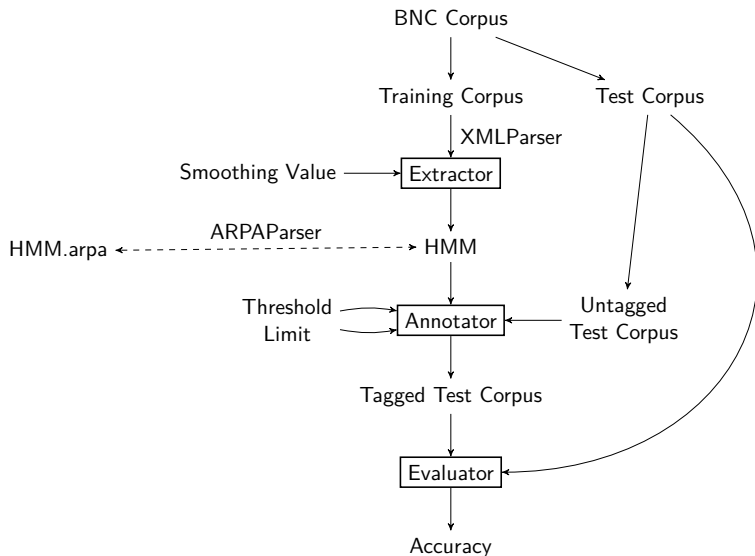
# Schematische Darstellung



# Schematische Darstellung



# Schematische Darstellung





# Outline

- 1 Aufgabenstellung
- 2 Modellierung von Hidden-Markov-Modellen**
- 3 Supervised Training mit Smoothing
- 4 Annotation mit Handling of Unknown Words
- 5 Evaluation des Modells
- 6 Ergebnisse

# Modellierung von Hidden-Markov-Modellen

HMM  $H = (Q, V, \#, \tau, \varepsilon)$  , wobei:

- $Q$  eine endliche Menge von Zuständen (POS-Tags),
- $V$  eine endliche Menge von Beobachtungen (Wörter),
- $\#$  der Startzustand,
- $\tau: Q \rightarrow Q \rightarrow [0, 1]$  die bed. Wsk.-Verteilung der Transitionen,
- $\varepsilon: Q \rightarrow V \rightarrow [0, 1]$  die bed. Wsk.-Verteilung der Emissionen.

# Modellierung von Hidden-Markov-Modellen

HMM  $H = (Q, V, \#, \tau, \varepsilon)$

```
1 data HMM =
2   HMM { -- | the initial state
3         initialState :: POSTag
4         -- | the list of states
5         , states      :: [POSTag]
6         -- | the list of observations
7         , observations :: [Word]
8         -- | the map of all transitions
9         , transitions  :: M.Map (POSTag, POSTag) Prob
10        -- | the map of all emissions
11        , emissions    :: M.Map (POSTag, Word) Prob
12        } deriving (Show)
```

# Modellierung von Hidden-Markov-Modellen

## ARPA-Format

```
1 \data\  
2  
3 \initial:  
4 START  
5  
6 \transitions:  
7 4.647932276069019e-2 AJ0 AJ0  
8 ...  
9  
10 \emissions:  
11 5.2954882440160984e-5 AJ0 Abdominal  
12 ...  
13  
14 \end\  

```

# Outline

- 1 Aufgabenstellung
- 2 Modellierung von Hidden-Markov-Modellen
- 3 Supervised Training mit Smoothing**
- 4 Annotation mit Handling of Unknown Words
- 5 Evaluation des Modells
- 6 Ergebnisse

# Supervised Training

- Trainingsgegenstand ist ein mit POS-Tags annotierter Textkorpus  $C$  (in unserem Fall: BNC im XML-Format)
- Supervised Training mithilfe relativer Häufigkeitsschätzung:

$$\text{Unigramme: } \tau_1(t_i) = \frac{f(t_i)}{|C|}$$

$$\text{Transitionen (Bigramme): } \tau_2(t_i | t_{i-1}) = \frac{f(t_i, t_{i-1})}{f(t_{i-1})}$$

$$\text{Emissionen: } \varepsilon(w_i | t_i) = \frac{f(w_i, t_i)}{f(t_i)}$$

mit  $w_j$  ... Wörter,  $t_j$  ... POS-Tags

## Lineare Interpolation (Jelinek und Mercer, 1980)

- Gefahr bei kleinen Trainingskorpora: zu geringes Vorkommen für jedes POS-Tag-Bigramm, um verlässliche Wahrscheinlichkeit zu bestimmen
- Bestimmung des Transitionswahrscheinlichkeiten durch lineare Interpolation der zugehörigen POS-Tag-Unigramme und -Bigramme:

$$\tau(t_i, t_{i-1}) = \lambda_1 \tau_1(t_i) + \lambda_2 \tau_2(t_i | t_{i-1})$$

mit  $\lambda_1 + \lambda_2 = 1$ .

# Supervised Training mit Smoothing

```
1 extract content smoothingValue =
2   ...
3   where hmm = flip compute smoothingValue
4             . flip count emptyCounts
5             $ XP.parseXML content
```

```
1 smooth :: Prob    -- ^ the bigram probability
2         -> Prob    -- ^ the unigram probability
3         -> Double  -- ^ the smoothing value
4         -> Prob    -- ^ the resulting probability
5 smooth biprob uniprob sv
6   = (1 - sv) * biprob + sv * uniprob
```



# Outline

- 1 Aufgabenstellung
- 2 Modellierung von Hidden-Markov-Modellen
- 3 Supervised Training mit Smoothing
- 4 Annotation mit Handling of Unknown Words**
- 5 Evaluation des Modells
- 6 Ergebnisse

# Viterbi-Algorithmus

**Input:** HMM  $H = (Q, V, \#, \tau, \varepsilon)$  und Satz  $w = w_1 \dots w_n$ .

**Output:** Tag-Sequenz  $\hat{t} = t_1 \dots t_n$ , sodass:

$$\hat{t} = \operatorname{argmax}_{t_1 \dots t_n \in T^n} P(\mathbb{W} = w_1 \dots w_n, \mathbb{T} = t_1 \dots t_n)$$

mit

$$P(\mathbb{W} = w_1 \dots w_n, \mathbb{T} = t_1 \dots t_n) = \varepsilon(w_1 \mid t_1) \cdot \tau(t_1 \mid \#) \cdot \prod_{i=2}^n \varepsilon(w_i \mid t_i) \cdot \tau(t_i \mid t_{i-1}) \cdot \tau(\# \mid t_n)$$

# Beam Search (Rabiner, 1989)

- Grenzwert  $\theta$  ... gibt den Anteil des Maximalwertes an, der überschritten werden muss
- Limit  $l$  ... gibt die maximale Anzahl der Elemente (Beam-Breite) an
- Score  $\delta$  ... gibt den derzeitigen Score eines Elements an

$$\delta_{new} \geq \frac{\delta_{max}}{\theta}$$

## Annotation mit Beam Search

```
1  annotate :: String      -- ^ the text corpus
2         -> HMM          -- ^ the HMM
3         -> Maybe Double -- ^ the threshold
4         -> Maybe Int    -- ^ the limit
5         -> String       -- ^ the annotated corpus
6  annotate [] _ _ _      = []
7  annotate tc hmm th l =
8      let (s : ss) = lines tc
9          pq       = vInit th l
10         rawTags  = vIterate (words s) hmm pq
11         fillUp   = map (const "-") ([1 ..] :: [Int])
12         tags     = L.drop 1 rawTags ++ fillUp
13     in unwords (zipWith (\w t -> w ++ "/" ++ t)
14                   (words s) tags)
15         ++ "\n" ++ go (unlines ss) hmm th l
```

# Unknown Words (Samuelsson, 1993)

**Problem:** Trainingskorpus enthält unbekannte Wörter

**Lösung:** Suffix-Analyse  $\rightarrow$  Schätzung der Emissionen für das unbekannte Wort:

$$\varepsilon'(t_i | u_{l,m}) = \begin{cases} \tau_1(t) & , \text{ wenn } |u_{l,m}| = 0 \\ \frac{\varepsilon(t_i|u_{l,m}) + \sigma\varepsilon'(t_i|u_{l+1,m})}{1+\sigma} & , \text{ sonst} \end{cases}$$

mit  $\sigma$  ... Standardabweichung von  $\tau_1$ ,  $t_i$  ... POS-Tag und  $w_j = u_1 \dots u_m$  ... unbekanntes Wort.

# Outline

- 1 Aufgabenstellung
- 2 Modellierung von Hidden-Markov-Modellen
- 3 Supervised Training mit Smoothing
- 4 Annotation mit Handling of Unknown Words
- 5 Evaluation des Modells**
- 6 Ergebnisse

# Evaluation des Modells

- Evaluation des Hidden-Markov-Modells anhand eines Testkorpus mit entfernten POS-Tags
- Vergleich des getaggten Korpus und des originalen Korpus (vor dem Entfernen der POS-Tags):

$$accuracy = \frac{\text{Anzahl übereinstimmender POS-Tags}}{\text{Anzahl aller getaggten Wörter}}$$

# Outline

- 1 Aufgabenstellung
- 2 Modellierung von Hidden-Markov-Modellen
- 3 Supervised Training mit Smoothing
- 4 Annotation mit Handling of Unknown Words
- 5 Evaluation des Modells
- 6 Ergebnisse**



# Literaturverzeichnis