

Prof. Dr. Andreas Hotho,
M.Sc. Janna Omeljanenko
Lecture Chair X for Data Science, Universität Würzburg

2. Exercise for “Sprachverarbeitung und Text Mining”

12.11.2021

1 Knowledge Questions

1. What is Part of Speech Tagging and what may Part of Speech Tagging be used for?

Part of Speech Tagging refers to the classification of words of a sentence into syntactic units (word types).

As input features for:

- Parsers
- Information Extraction
- Speech synthesis
- Machine Translation

2. What is the difference between open and closed classes of word types in POS-tagging?

Closed classes, in contrast to open classes, contain only a finite set of words, and new ones cannot be added.

Example: Pronouns, Auxiliaries, Prepositions

3. Which two typical approaches are used for POS-tagging?

- Rule-based approaches (z.B. ENGTWOL)
- Stochastic approaches (z.B. HMM, MEMM, CRF)

4. Which two types of probabilities are needed in the modeling approach of the Hidden Markov Model? What elements does a Hidden Markov Model generally consist of?

- Transition probabilities $P(T_i|T_{i-1})$
- Likelihood probabilities $P(W_i|T_i)$

Generally, an HMM consists of:

- A list **Q** with the possible states
- a matrix **A** with the transition probabilities
- a list **O** with the possible observations
- a matrix **B** with observation likelihoods
- start and end states

5. What does the Markov-assumption state when applied to POS-tagging?

The probability of a tag depends only on the previous tag.

6. Which POS-tagsets do you know for the German and English language?

- German: Besides some minor modifications, the STTS (Stuttgart-Tübingen Tagset) is actually always used in German.
- English: Penn TreeBank tagset

7. All POS-tagging methods presented in the lecture require a list (dictionary) that lists all possible POS-tags for a word. However, such lists are usually not complete. Give a solution for POS-tagging with unknown words.

- Assume that each tag can be assigned to the unknown word with equal probability.

- Use additional information about word properties and compute probability for each POS-tag given observed properties (multipliable under assumption of independence).
 - e.g. English word ending with 's' → noun in the plural with $x\%$ probability.
 - e.g. English word beginning with capital letter → proper name with $y\%$ probability.

8. Name three types of evaluation of POS-tagging methods.

- Overall error rate
- Error rates on particular tags
- Error rates on particular words
- Confusion-Matrix

with respect to sample solution (gold standard)

2 POS Tagging

Part-of-speech taggers model the words as observations (observed variables), and the associated word types (part-of-speech tags) as hidden variables (hidden states). We assume in the following a highly simplified set of generalized word types (tag-set):

| | | |
|-----|------------|--------------------------|
| DET | Determiner | the,a,... |
| N | Noun | year,home,costs,time,... |
| PRO | Pronoun | he,their,you |
| V | Verb | said,took,saw |

Using the table below, calculate the most likely sequence of tags for the following sentence:

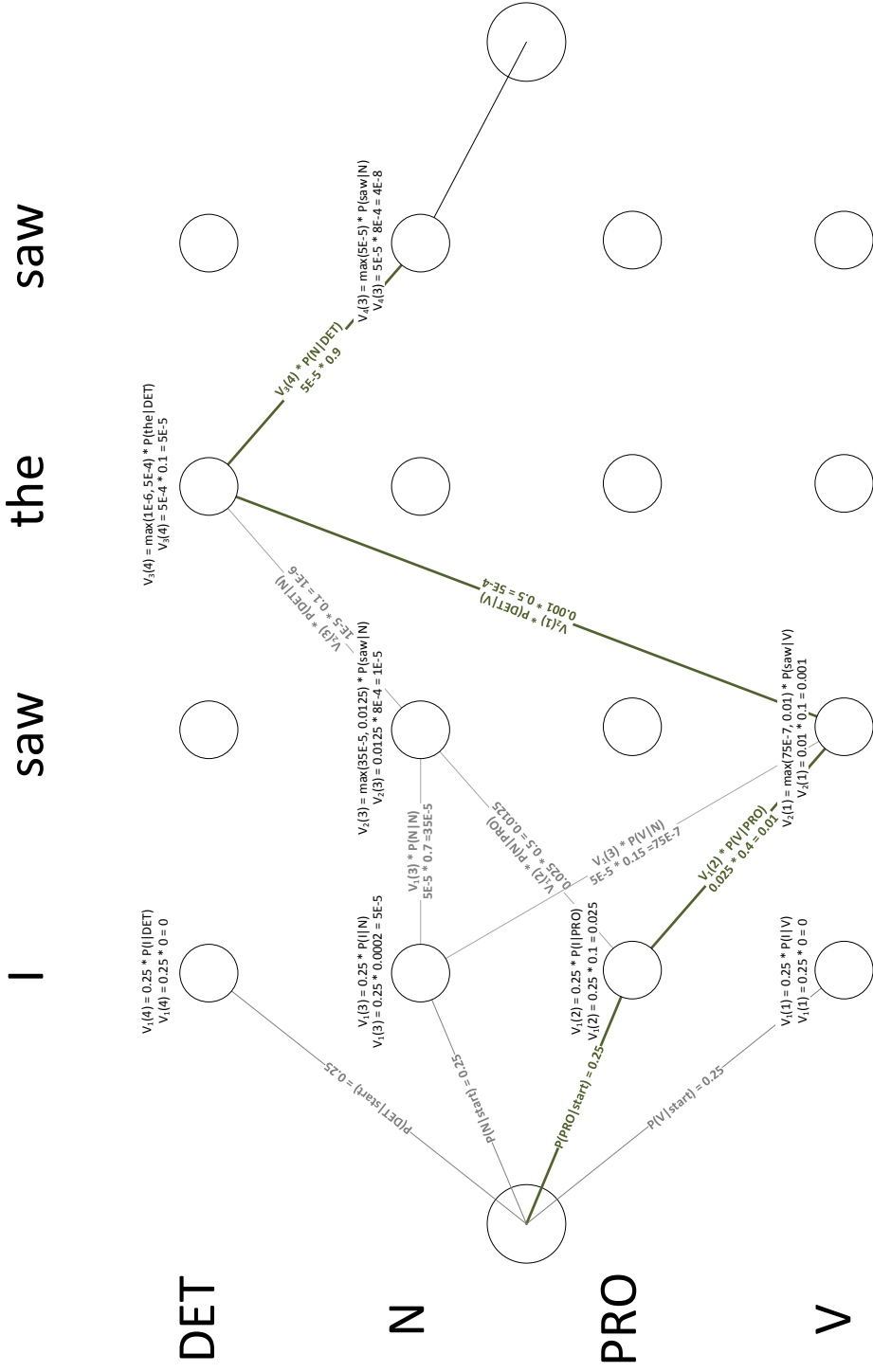
„I saw the saw“

1. Apply the Viterbi algorithm.
2. Which are the main factors influencing the complexity of the calculation?

Assume uniformly distributed start transition probabilities.

| from\to | DET | N | PRO | V |
|---------|------|-----|------|------|
| DET | 0.05 | 0.9 | 0.05 | 0 |
| N | 0.1 | 0.7 | 0.05 | 0.15 |
| PRO | 0.05 | 0.5 | 0.05 | 0.4 |
| V | 0.5 | 0.1 | 0.4 | 0 |

| w\t | DET | N | PRO | V |
|-----|-----|--------|-----|-----|
| I | 0 | 0.0002 | 0.1 | 0 |
| saw | 0 | 0.0008 | 0 | 0.1 |
| the | 0.1 | 0 | 0 | 0 |



Note: The unreachable states or states with probability 0 are not shown in the graph for readability reasons.

Main influence factors are the length of the sequence (linear) and the amount of available POS tags (quadratical).