

GWV – Grundlagen der Wissensverarbeitung

Tutorial 9: Hidden Markov Models

A Part-of-Speech(PoS) tagger is a program that assigns word classes (i.e. PoS) to words of a text. For example, for the input

The dog barked

the corresponding output should be the following¹:

The\DET dog\NN barked\VBD

As always, you can find the relevant data at our wiki. The data we use will be one word per line, with each line consisting of word [tab] PoS-tag.

Let

$Word_i$ be the variable denoting the i th Word

Tag_i be the variable denoting the PoS-tag for the i th word

Exercise 1.1: (A simple PoS-Tagger)

Write a PoS tagger that uses a model where the probability distribution of $Word_i$ only depends on the state of Tag_i and Tag_i only depends on Tag_{i-1} . The tagger model (i.e. the transition and emission probabilities) should be trained from the file provided in the wiki. (4 Pt.)

Create a function that takes a list of words (possibly from the command line) and uses filtering to produce a corresponding list of PoS tags. (2 Pt.)

Make sure that your tagger can cope with input that includes words that are not in the training data. (2 Pt.)

Implement either the forward-backward algorithm to output several PoS tags per word or the viterbi algorithm to output the most probable tag sequence. (4 Pt.)

Hint: You can assess how good your tagger is by tagging the test file and comparing the output produced by the tagger with the original file, e.g. with

```
diff -u original_file generated_file | grep '^+' | wc -l
```

Version: December 8, 2014
Achievable score on this sheet: 12

of
12

¹The tags for this example do not correspond to the ones in the data file (because the data we work on is German). The German PoS tags are explained here: <http://www.ims.uni-stuttgart.de/forschung/ressourcen/lexika/TagSets/stts-table.html>