

# Computational Linguistics

LIN 467/567 – Spring 2018

Homework 3

Due date: end of March 23rd

## Instructions

Note: **Answers not submitted as specified below will not be graded.**

Upload to UBLearn a single zip file containing your answers to the two questions below. The zip file must be named 'hw3UBitNAME.zip', where UBitNAME is your UB email's name. Don't forget to click the button after uploading the file, or it will not be submitted.

Students can work with others but must ultimately create their own individual answers. Any homework assignments that are sufficiently similar may be considered instances of plagiarism, and will be penalized and reported. See syllabus for more information about UB's official policy on ethical behavior.

## Exercises

1. There is a total of twelve logically possible taggings for the sentence *Kids like sprouts* ('Kids' = NN / VB, 'like' = VB / PREP / ADJ, 'sprouts' = NN / VB). Construct the (maxed) trellis chart in a spreadsheet (Excel, OpenOffice, or GoogleDocs, for example) and determine what is the most likely tag sequence for that string given the HMM below. Your answer must be submitted as a PDF file called UBitNAME.pdf.

[3 points]

$\sigma(\text{NN}, \text{VB}) = 0.6$	$\tau(\text{kids}, \text{NN}) = 0.4$
$\sigma(@, \text{VB}) = 0.4$	$\tau(\text{kids}, \text{VB}) = 0.1$
$\sigma(\text{NN}, \text{PREP}) = 0.3$	$\tau(\text{like}, \text{VB}) = 0.7$
$\sigma(@, \text{NN}) = 0.5$	$\tau(\text{like}, \text{ADJ}) = 0.005$
$\sigma(\text{VB}, \text{VB}) = 0.15$	$\tau(\text{sprouts}, \text{NN}) = 0.4$
$\sigma(\text{VB}, \text{PREP}) = 0.25$	$\tau(\text{sprouts}, \text{VB}) = 0.2$
$\sigma(\text{VB}, \text{ADJ}) = 0.3$	$\tau(\text{like}, \text{PREP}) = 0.02$
$\sigma(\text{VB}, \text{NN}) = 0.6$	
$\sigma(\text{NN}, \text{ADJ}) = 0.01$	
$\sigma(\text{VB}, @) = 0.4$	
$\sigma(\text{PREP}, \text{VB}) = 0.5$	
$\sigma(\text{PREP}, \text{NN}) = 0.3$	
$\sigma(\text{ADJ}, \text{VB}) = 0.008$	
$\sigma(\text{ADJ}, \text{NN}) = 0.4$	
$\sigma(\text{NN}, @) = 0.6$	

2. The file `viterbi.pl` is a Prolog implementation of a modified Viterbi algorithm which lists all possible taggings and their probabilities instead of simply computing the best one. This is useful to determine how ambiguous the input is, how uncertain the tagger is, or if there is a need to backtrack to second most probable tagging, and so on. However, `viterbi.pl` is incomplete: it does not have the HMM model itself. To complete this exercise, you must:

- (a) Tag the words in the `DA.txt` mini corpus and create a tagged version called `TDA.txt`. Use the StanfordNLP suite MaxEnt tagger, by issuing the following (single) command:<sup>1</sup>

```
java -cp "*" edu.stanford.nlp.tagger.maxent.MaxentTagger -model
models/english-left3words-distsim.tagger -textFile DA.txt -tokenize true
-outputFormat tsv -outputFile TDA.txt
```

Before you can run the above, you must download [the MaxEnt tagger](#) and [JKD](#) installed (you may already have it, depending on your operative system).

- (b) Use bash commands to extract a Prolog-readable HMM from your `TDA.txt` and add it to `viterbi.pl`. Do not smooth the model.

Your answer should therefore consist of (i) the `TDA.txt` tagged corpus created from a normalized version of `DA.txt`, (ii) a text file `bashUBitNAME.txt` containing bash commands used to extract a Prolog-readable HMM from the tagged corpus `TDA.txt`, and (iii) an augmented `viterbi.pl` text file named `hmmUBitNAME.pl` that now contains the HMM you derived from `TDA.txt`. The `viterbi.pl` file must successfully run the Viterbi algorithm over word sequences (numbers and tags are fictitious):<sup>2</sup>

---

```
?- viterbi_all([he,can,can,a,can],X).
X = tag(2.5515e-12, [pron, aux, v, det, aux]) ;
X = tag(1.27575e-12, [pron, aux, v, det, v]) ;
X = tag(1.96466e-11, [pron, aux, v, det, n]);
false
```

---

[3 points]

---

<sup>1</sup>To read more go to <https://nlp.stanford.edu/software/tagger.shtml>. This tagger used the [Penn Treebank Tagset](#) (see J&M'08 for more discussion).

<sup>2</sup>The number  $2.5515e-12$  means  $2.5515 \times 10^{-12} = 0.0000000000025515$ . This is the underflow problem discussed in class, caused by multiplying probabilities. If you wish, modify `viterbi.pl` so that it works in log-space instead.