

School of Computing and Information Systems
The University of Melbourne
COMP90042
NATURAL LANGUAGE PROCESSING (Semester 1, 2021)
Sample solutions: Week 4

Discussion

1. What is a **POS tag**?

- A POS tag, AKA word classes, is a label assigned to a word token in a sentence which indicates some grammatical (primarily syntactic) properties of its function in the sentence.

(a) POS tag (by hand) the following sentence: `Pierre Vinken, 61 years old, will join the board as a nonexecutive director Nov. 29.` according to the Penn Treebank tags. (Note that some of the tags are somewhat obscure.)

- According to the Penn Treebank:

```
NNP Pierre
NNP Vinken

    ' '
    CD 61
NNS years
JJ old

    ' '
MD will
VB join
DT the
NN board
IN as
DT a
JJ nonexecutive
NN director
NNP Nov.
CD 29

. .
```

(b) What are some common approaches to POS tagging? What aspects of the data might allow us to predict POS tags systematically?

- **Unigram:** Assign a POS tag to a token according to the most common observation in a tagged corpus; many words are unambiguous, or almost unambiguous.
- **N-gram:** Assign a POS tag to a token according to the most common tag in the same sequence (based on the sentence in which the token occurs) of n tokens (or tags) in the tagged corpus; context helps disambiguate.
- **Rule-based:** Write rules (relying on expertise of the writer) that disambiguate unigram tags.

- **Sequential:** Learn a Hidden Markov Model (or other model) based on the observed tag sequences in a tagged corpus.
 - **Classifier:** Treat as a supervised machine learning problem, with tags from a tagged corpus as training data.
2. What are the assumptions that go into a **Hidden Markov Model**? What is the time complexity of the **Viterbi algorithm**? Is this practical?
- **Markov assumption:** the likelihood of transitioning into a given state depends only on the current state, and not the previous state(s) (or output(s))
 - **Output independence assumption:** the likelihood of a state producing a certain word (as output) does not depend on the preceding (or following) state(s) (or output(s)).
 - The time complexity of the Viterbi algorithm, for an HMM with T possible states, and a sentence of length W , is $\mathcal{O}(T^2W)$. In POS tagging, there might typically be approximately 100 possible tags (states), and a typical sentence might have 10 or so tokens, so ... yes, it is practical (unless we need to tag really, really, quickly, e.g. tweets as they are posted).
- (a) How can an HMM be used for POS tagging a text? For the purposes of POS tagging:
- Tokens (possibly ignoring punctuation) correspond to outputs, and POS tags correspond to states of the HMM.
 - We will build a model based around sentences in a tagged corpus. We could instead use a document-based model; in which case, the model will probably discover that the initial distribution of states is very similar to the transition from end-of-sentence punctuation (. ? ! etc.), which probably represents a needless complication.
- How can the initial state probabilities π be estimated?
 - Record the distribution of tags for the first token of each sentence in a tagged corpus.
 - How can the transition probabilities A be estimated?
 - For each tag, record the distribution of tags of the immediately following token in the tagged corpus. (We might need to introduce an end-of-sentence dummy for the probabilities to add up correctly.)
 - How can the emission probabilities B be estimated?
 - For each tag, record the distribution of corresponding tokens in the tagged corpus.
- (b) Estimate π , A and B for POS tagging, based on the following tagged corpus:
- silver-JJ wheels-NNS turn-VBP
 - wheels-NNS turn-VBP right-JJ
 - right-JJ wheels-NNS turn-VBP
- For π , there are three sentences: two begin with JJ and one with NNS. Consequently:

$$\pi[JJ, NNS, VBP] = [\frac{2}{3}, \frac{1}{3}, 0]$$

- For A , we need to observe the distribution for each tag.
 - For JJ, two instances are immediately followed by NNS.
 - For NNS, all three instance are followed by VBP.
 - For VBP, it is followed by JJ once.

A	JJ	NNS	VBP
(from) JJ	0	1	0
NNS	0	0	1
VBP	1	0	0

- For B , we can simply read off the corresponding words for each tag in the corpus:
 - For JJ, there is one instance of *silver* and two of *right*.
 - For NNS, there are three instances of *wheels*.
 - For VBP, there are three instances of *turn*.
- Consequently:

B	right	silver	turn	wheels
(from) JJ	$\frac{2}{3}$	$\frac{1}{3}$	0	0
NNS	0	0	0	1
VBP	0	0	1	0

3. Consider using the following Hidden Markov Model to tag the sentence *silver wheels turn*:

$$\pi[\text{JJ}, \text{NNS}, \text{VBP}] = [0.3, 0.4, 0.3]$$

A	JJ	NNS	VBP	B	silver	wheels	turn
JJ	0.4	0.5	0.1	JJ	0.8	0.1	0.1
NNS	0.1	0.4	0.5	NNS	0.3	0.4	0.3
VBP	0.4	0.5	0.1	VBP	0.1	0.3	0.6

(a) Visualise the HMM as a graph.

(b) Use the **Viterbi algorithm** to find the most likely tag for this sequence.

- The most likely tag sequence can be read right-to-left, based upon the maximum probability we've observed: in this case, 0.0144 when *turn* is a VBP; this value is derived from the NNS \rightarrow VBP transition, so we can infer that *wheels* is an NNS; that in turn comes from the JJ \rightarrow NNS transition, so *silver* is a JJ.

α		1:silver	2:wheels	3:turn
JJ:	JJ	$\pi[\text{JJ}]B[\text{JJ}, \text{silver}]$ $0.3 \times 0.8 = 0.24$		
NNS:	NNS	$\pi[\text{NNS}]B[\text{NNS}, \text{silver}]$ $0.4 \times 0.3 = 0.12$		
VBP:	VBP	$\pi[\text{VBP}]B[\text{VBP}, \text{silver}]$ $0.3 \times 0.1 = 0.03$		

α	1:silver	2:wheels		3:turn
JJ:	0.24	JJ \rightarrow JJ	$A[\text{JJ},\text{JJ}]B[\text{JJ}, \text{wheels}]$	
		0.24	$\times 0.4 \times 0.1 = \mathbf{0.0096}$	
		NNS \rightarrow JJ	$A[\text{NNS},\text{JJ}]B[\text{JJ}, \text{wheels}]$	
		0.12	$\times 0.1 \times 0.1 = 0.0012$	
		VBP \rightarrow JJ	$A[\text{VBP},\text{JJ}]B[\text{JJ}, \text{wheels}]$	
NNS:	0.12	0.03	$\times 0.4 \times 0.1 = 0.0012$	
		JJ \rightarrow NNS	$A[\text{JJ},\text{NNS}]B[\text{NNS}, \text{wheels}]$	
		0.24	$\times 0.5 \times 0.4 = \mathbf{0.048}$	
		NNS \rightarrow NNS	$A[\text{NNS},\text{NNS}]B[\text{NNS}, \text{wheels}]$	
		0.12	$\times 0.4 \times 0.4 = 0.0192$	
VBP:	0.03	VBP \rightarrow NNS	$A[\text{VBP},\text{NNS}]B[\text{NNS}, \text{wheels}]$	
		0.03	$\times 0.5 \times 0.4 = 0.006$	
		JJ \rightarrow VBP	$A[\text{JJ},\text{VBP}]B[\text{VBP}, \text{wheels}]$	
		0.24	$\times 0.1 \times 0.3 = 0.0072$	
		NNS \rightarrow VBP	$A[\text{NNS},\text{VBP}]B[\text{VBP}, \text{wheels}]$	
		0.12	$\times 0.5 \times 0.3 = \mathbf{0.018}$	
		VBP \rightarrow VBP	$A[\text{VBP},\text{VBP}]B[\text{VBP}, \text{wheels}]$	
		0.03	$\times 0.1 \times 0.3 = 0.0009$	

α	1:silver	2:wheels	3:turn	
JJ:	0.24	0.0096	JJ \rightarrow JJ	$A[\text{JJ},\text{JJ}]B[\text{JJ}, \text{turn}]$
		JJ \rightarrow JJ	0.0096	$\times 0.4 \times 0.1 = 0.000384$
			NNS \rightarrow JJ	$A[\text{NNS},\text{JJ}]B[\text{JJ}, \text{turn}]$
			0.048	$\times 0.1 \times 0.1 = 0.00048$
			VBP \rightarrow JJ	$A[\text{VBP},\text{JJ}]B[\text{JJ}, \text{turn}]$
NNS:	0.12		0.018	$\times 0.4 \times 0.1 = \mathbf{0.00072}$
		0.048	JJ \rightarrow NNS	$A[\text{JJ},\text{NNS}]B[\text{NNS}, \text{turn}]$
		JJ \rightarrow NNS	0.0096	$\times 0.5 \times 0.3 = 0.00144$
			NNS \rightarrow NNS	$A[\text{NNS},\text{NNS}]B[\text{NNS}, \text{turn}]$
			0.048	$\times 0.4 \times 0.3 = \mathbf{0.00576}$
VBP:	0.03		VBP \rightarrow NNS	$A[\text{VBP},\text{NNS}]B[\text{NNS}, \text{turn}]$
			0.018	$\times 0.5 \times 0.3 = 0.0027$
		0.018	JJ \rightarrow VBP	$A[\text{JJ},\text{VBP}]B[\text{VBP}, \text{turn}]$
		NNS \rightarrow VBP	0.0096	$\times 0.1 \times 0.6 = 0.000576$
			NNS \rightarrow VBP	$A[\text{NNS},\text{VBP}]B[\text{VBP}, \text{turn}]$
			0.048	$\times 0.5 \times 0.6 = \mathbf{0.0144}$
			VBP \rightarrow VBP	$A[\text{VBP},\text{VBP}]B[\text{VBP}, \text{turn}]$
			0.018	$\times 0.1 \times 0.6 = 0.00108$