

Name: Solutions

Email: hollink@gmail.com

CMSC/LING 723 - Computational Linguistics I  
Midterm Exam 18 Oct 2011

Question	Points
1	7
2	5
3a	10
3b	3
4	5
5	5
6	5
7a	10
7b	15
7c	5
7d	20
7e	10
Total	100

For Question 1, consider the following set of sentences as a corpus:

---

<s> I am Sam </s>  
<s> Sam I am </s>  
<s> I do not like green eggs and ham </s>  
<s> I do not like them , Sam I am </s>

---

**Question 1a. (1 point)** How many bigrams are there in this corpus (total)?

27 total; 18 unique.

**Question 1b. (1 point)** What is the most frequent bigram in this corpus?

“I am” and “<s> I” are tied, at 3 occurrences each.

**Question 1b. (2 points)** What is the probability of the bigram “*I am*” under an unsmoothed model (i.e.,  $P_{MLE}(am | I)$ )?

$$\frac{c(I \text{ am})}{c(I)} = \frac{3}{5} = 0.6$$

**Question 1b. (2 points)** What is the probability of the bigram “*I do*” under an unsmoothed model (i.e.,  $P_{MLE}(do | I)$ )?

$$\frac{c(I \text{ do})}{c(I)} = \frac{2}{5} = 0.4$$

**Question 1c. (1 point)** What is the most frequent co-occurrence bigram in this corpus?

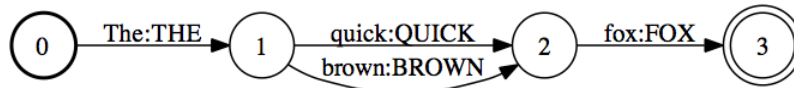
“<s> I” and “I </s>” are tied at 5 occurrences each.

**Question 2. (5 points)** What is the Soundex representation for the name “*Engkebrethson*”?

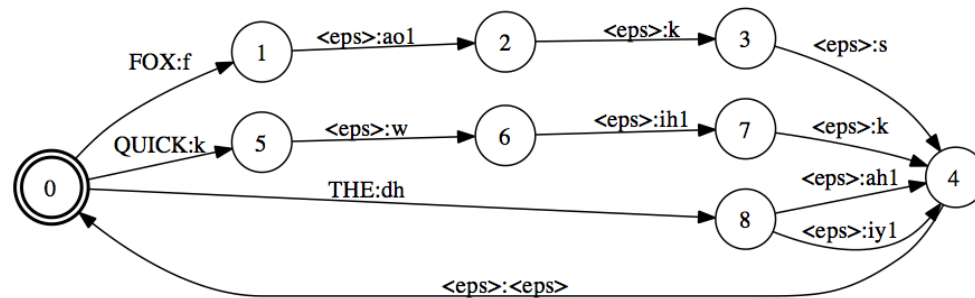
E522 (E526 is incorrect, and indicates improper treatment of the “gek” subsequence.)

For Question 3, consider the following two FSTs.

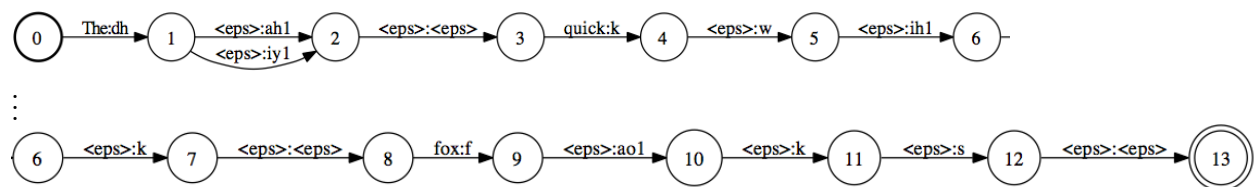
FST1:



FST2:



**Question 3a. (10 points)** Draw the composition of these two transducers (i.e.,  $\text{FST2} \circ \text{FST1}$ ).



**Question 3b. (3 points)** Is FST1 a deterministic transducer? Is FST2? Is your composed FST for 3a deterministic?

Yes, no, no.

**Question 4. (5 points)** Explain the Markov Assumption and why we use it. You may use equations in your answer if necessary, and/or give an example of how we might use the Markov Assumption.

The Markov Assumption is that events are conditionally independent of other events that occurred  $n + 1$  time steps prior (for a Markov order- $n$  model). We use the Markov assumption because it is computationally intractable to condition on every previous event in our history; the assumption can also resolve some sparsity issues.

**Question 5. (5 points)** Explain the difference between interpolation models and backoff models of smoothing.

Interpolation models use information from all of the models; backoff models use information from only one model, “backing off” to a lower-order model only if there is not enough evidence to use information from a higher-order model.

**Question 6. (5 points)** Consider the following scenario: we have a text corpus of 10,000 words, and are provided with a vocabulary list of 100 words for the corpus. There are 5 words in our vocabulary that were never observed in the corpus. If we have observed the word “*the*” 900 times in our corpus, what is the expected frequency estimate of this word according to Laplace’s Law?

$N=10,000$

$V=100$

$C(\text{the})=900$

$$C_{LAP} = P_{LAP} * N = \frac{C+1}{N+V} * N = \frac{901}{10,100} * 10,000 = 892.08$$

For Question 7, consider the following string: *bats can fly*

An HMM POS-tag model is described by the following two tables:

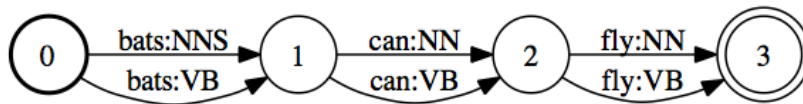
$$b_j(v_k) = \mathbf{P}(v_k|\tau_j)$$

	$k:$	1	2	3
$j$		bats	can	fly
1	NN	0	$\frac{1}{2}$	$\frac{1}{2}$
2	NNS	1	0	0
3	VB	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$

$$a_{ij} = \mathbf{P}(\tau_j|\tau_i)$$

	$j:$	0	1	2	3
$i$		</s>	NN	NNS	VB
0	<s>	0	0	$\frac{1}{2}$	$\frac{1}{2}$
1	NN	$\frac{1}{4}$	$\frac{1}{4}$	0	$\frac{1}{2}$
2	NNS	0	$\frac{1}{2}$	0	$\frac{1}{2}$
3	VB	$\frac{1}{2}$	$\frac{1}{4}$	0	$\frac{1}{4}$

**Question 7a. (10 points)** Based on the model given above, draw an (unweighted) FST of word:tag sequences for the string “*bats can fly*”, using as few states as possible. Show only the word:tag sequences which have probability greater than zero.



**Question 7b. (15 points)** Use the following table (trellis) to fill in the Viterbi forward probability  $\alpha_i^{\max}(t)$ , as well as the corresponding back pointers, given the model from the previous page.

(see 1st table on next page;  $v$ =Viterbi, **red**  $v^{\max}$ =Viterbi-best values,  $\psi$ =back pointers)

**Question 7c. (5 points)** What is the Viterbi-best tag sequence according to the above chart?

**NNS NN VB**

**Question 7d. (20 points).** What is the probability of the string “*bats can fly*” according to the aforementioned model? You may re-use the table (trellis) on the previous page, or re-create another below.

$$\frac{105+224}{13,824} = \frac{329}{13,824} = 0.023799$$

(see 2nd table on next page for work)

**Question 7e. (10 points)** What is the perplexity of the aforementioned model on the string “*bats can fly*”?

$$0.023799^{-1/N}, N = 4 \\ = 2.5460$$

j	$t_i$	$\tau_k$	1	2	3	4
			bats	can	fly	$\langle s \rangle$
0						$v_{\langle s \rangle}(4) = v_{NN}(3) * a_{NN, \langle s \rangle} * b$ $= \frac{1}{64} * \frac{1}{4} * 1 = \frac{1}{256}$ $v_{\langle s \rangle}^{\max}(4) = v_{VB}(3) * a_{VB, \langle s \rangle} * b$ $= \frac{1}{48} * \frac{1}{2} * 1 = \frac{1}{96}$ $\psi_{\langle s \rangle}(4) = VB$
1	NN	0	$b_{NN, bats} = 0$ $v_{NN}^{\max}(2) = v_{NNS}(1) * a_{NNS, NN} * b_{NN, can}$ $= \frac{1}{2} * \frac{1}{2} * \frac{1}{2} = \frac{1}{8}$ $v_{NN}(2) = v_{VB}(1) * a_{VB, NN} * b_{NN, can}$ $= \frac{1}{6} * \frac{1}{4} * \frac{1}{2} = \frac{1}{48}$ $\psi_{NN}(2) = NNS$	$v_{NN}^{\max}(3) = v_{NN}(2) * a_{NN, NN} * b_{NN, fly}$ $= \frac{1}{8} * \frac{1}{4} * \frac{1}{2} = \frac{1}{64}$ $v_{NN}(3) = v_{VB}(2) * a_{VB, NN} * b_{NN, fly}$ $= \frac{1}{12} * \frac{1}{4} * \frac{1}{2} = \frac{1}{96}$ $\psi_{NN}(3) = NN$		
2	NNS	0	$v_{NNS}(1) = \alpha_{\langle s \rangle}(0) * a_{\langle s \rangle, NNS} * b_{NNS, bats}$ $= 1 * \frac{1}{2} * 1 = \frac{1}{2}$ $\psi_{NNS}(1) = \langle s \rangle$	$b_{NNS, can} = 0$	$b_{NNS, fly} = 0$	
3	VB	0	$v_{VB}(1) = \alpha_{\langle s \rangle}(0) * a_{\langle s \rangle, VB} * b_{VB, bats}$ $= 1 * \frac{1}{3} * \frac{1}{2} = \frac{1}{6}$ $\psi_{VB}(1) = \langle s \rangle$	$v_{VB}^{\max}(2) = v_{NNS}(1) * a_{NNS, VB} * b_{VB, can}$ $= \frac{1}{2} * \frac{1}{2} * \frac{1}{3} = \frac{1}{12}$ $v_{VB}(2) = v_{VB}(1) * a_{VB, VB} * b_{VB, can}$ $= \frac{1}{6} * \frac{1}{4} * \frac{1}{3} = \frac{1}{72}$ $\psi_{VB}(2) = NNS$	$v_{VB}^{\max}(3) = v_{NN}(2) * a_{NN, VB} * b_{VB, fly}$ $= \frac{1}{8} * \frac{1}{2} * \frac{1}{3} = \frac{1}{48}$ $v_{VB}(3) = v_{VB}(2) * a_{VB, VB} * b_{VB, fly}$ $= \frac{1}{12} * \frac{1}{4} * \frac{1}{3} = \frac{1}{144}$ $\psi_{VB}(3) = NN$	
j	$t_i$	$\tau_k$	1	2	3	4
			bats	can	fly	$\langle s \rangle$
0		1				$\alpha_{\langle s \rangle}(4) = \alpha_{NN}(3) * a_{NN, \langle s \rangle} * b$ $+ \alpha_{VB}(3) * a_{VB, \langle s \rangle} * b$ $= \frac{35}{1152} * \frac{1}{4} * 1 + \frac{216}{7} * \frac{1}{2} * 1$ $= \frac{1152}{4608} + \frac{7}{432} = \frac{105+224}{432} = \frac{13,824}{432}$ $= 0.023799$
1	NN	0	$b_{NN, bats} = 0$ $\alpha_{NN}(2) = \alpha_{NNS}(1) * a_{NNS, NN} * b_{NN, can}$ $+ \alpha_{VB}(1) * a_{VB, NN} * b_{NN, can}$ $= \frac{1}{2} * \frac{1}{2} * \frac{1}{6} + \frac{1}{6} * \frac{1}{4} * \frac{1}{2}$ $= \frac{1}{8} + \frac{1}{48} = \frac{7}{48}$	$\alpha_{NN}(3) = \alpha_{NN}(2) * a_{NN, NN} * b_{NN, fly}$ $+ \alpha_{VB}(2) * a_{VB, NN} * b_{NN, fly}$ $= \frac{7}{48} * \frac{1}{4} * \frac{1}{2} + \frac{7}{72} * \frac{1}{4} * \frac{1}{2}$ $= \frac{384}{384} + \frac{1152}{576} = \frac{1152}{576}$		
2	NNS	0	$\alpha_{NNS}(1) = \alpha_{\langle s \rangle}(0) * a_{\langle s \rangle, NNS} * b_{NNS, bats}$ $= 1 * \frac{1}{2} * 1 = \frac{1}{2}$	$b_{NNS, can} = 0$	$b_{NNS, fly} = 0$	
3	VB	0	$\alpha_{VB}(1) = \alpha_{\langle s \rangle}(0) * a_{\langle s \rangle, VB} * b_{VB, bats}$ $= \frac{1}{3} * \frac{1}{2} = \frac{1}{6}$	$\alpha_{VB}(2) = \alpha_{NNS}(1) * a_{NNS, VB} * b_{VB, can}$ $+ \alpha_{VB}(1) * a_{VB, VB} * b_{VB, can}$ $= \frac{1}{2} * \frac{1}{2} * \frac{1}{3} + \frac{1}{6} * \frac{1}{4} * \frac{1}{2}$ $= \frac{1}{12} + \frac{1}{72} = \frac{7}{72}$	$\alpha_{VB}(3) = \alpha_{NN}(2) * a_{NN, VB} * b_{VB, fly}$ $+ \alpha_{VB}(2) * a_{VB, VB} * b_{VB, fly}$ $= \frac{7}{48} * \frac{1}{2} * \frac{1}{4} + \frac{7}{72} * \frac{1}{4} * \frac{1}{3}$ $= \frac{288}{288} + \frac{864}{864} = \frac{216}{216}$	

**Extra Credit. (5 points)** Design and write up (including solutions) your own test question for a future Computational Linguistics I course.