

CMPT 413 - Spring 2011 - Midterm #2

Please write down “Midterm #2” on the top of the answer booklet.

When you have finished, return your answer booklet along with this question booklet.

Mar 15, 2011

- (1) Consider a language model over character sequences that computes the probability of a word based on the characters in that word, so if word $w = c_0, c_1, \dots, c_n$ then $P(w) = P(c_0, \dots, c_n)$. Let us assume that the language model is defined as a bigram character model $P(c_i | c_{i-1})$ where

$$P(c_0, \dots, c_n) = \prod_{i=1,2,\dots,n} P(c_i | c_{i-1}) \quad (1)$$

For convenience we assume that we have explicit word boundaries: $c_0 = \text{bos}$ and $c_n = \text{eos}$ where *bos* stands for *begin sentence marker* and *eos* stands for *end of sentence marker*.

Based on this model, for the English word *booking* the probability would be computed as:

$$P(\text{booking}) = P(b | \text{bos}) \times P(o | b) \times P(o | o) \times P(k | o) \times P(i | k) \times P(n | i) \times P(g | n) \times P(\text{eos} | g)$$

The inflection *ing* is a suffix and is generated after the stem *book* with probability

$$P(\text{ing}) = P(i | k) \times P(n | i) \times P(g | n) \times P(\text{eos} | g)$$

In Semitic languages, like Arabic and Hebrew, the process of inflection works a bit differently. In Arabic, for a word like *kitab* the stem would be *k-t-b* where the place-holders ‘-’ for inflection characters have been added for convenience. We will assume that each word is made up of a sequence of consonant-vowel sequences CVCVCV... and the vowels always form the inflection.

- a. Provide the definition of an n -gram model that will compute the probability for the word *kitab* and *k-t-b* as follows:

$$P(\text{kitab}) = P(k | \text{bos}) \times P(t | k) \times P(b | t) \times P(i | b) \times P(a | i) \times P(\text{eos} | a)$$

$$P(\text{k-t-b}) = P(k | \text{bos}) \times P(t | k) \times P(b | t) \times P(- | b) \times P(- | -) \times P(\text{eos} | -)$$

Write down the equation for this n -gram model in the same mathematical notation as equation (1).

Answer:

$$P(c_0, \dots, c_n) = \begin{cases} \prod_{i=1}^n P(c_i | c_{i-1}) & \text{if } n \leq 3 \\ \left(P(c_1 | c_0) \times \prod_{i=3,5,\dots}^{\ell} P(c_i | c_{i-2}) \right) \times & \text{if } n > 3 \\ \left(P(c_2 | c_{\ell_o}) \times \prod_{i=4,6,\dots}^{\ell} P(c_i | c_{i-2}) \times P(c_n | c_{\ell_e}) \right) & \end{cases}$$

Define $\ell = n - (n \bmod 2)$ and ℓ_o is the last odd number less than ℓ and ℓ_e is the last even number less than ℓ . As long as the boundary cases are right for the bigrams, we don’t penalize off by one in the length, and we don’t penalize for $n \leq 3$.

- b. Using your n -gram model show how $P(kitab) = P(ktb) \times P(ia)$.

Answer:

$$\begin{aligned} P(kitab) &= P(c_0 = \text{bos}, c_1 = k, c_2 = i, c_3 = t, c_4 = a, c_5 = b, c_6 = \text{eos}) \\ &= P(ktb) \times P(ia, \text{eos}) \end{aligned}$$

$$P(ktb) = P(c_1 = k \mid c_0 = \text{bos}) \times P(c_3 = t \mid c_1 = k) \times P(c_5 = b \mid c_3 = t)$$

this term corresponds to the first bracket in the eqn above

$$P(ia) = P(c_2 = i \mid c_{\ell_o} = c_5 = b) \times P(c_4 = a \mid c_2 = i) \times P(c_n = c_6 = \text{eos} \mid c_{\ell_e} = c_4 = a)$$

corresponds to the second bracket in the eqn above

- (2) The probability model $P(t_i \mid t_{i-2}, t_{i-1})$ is provided below where each t_i is a part of speech tag, e.g. $P(D \mid N, V) = \frac{1}{3}$. Also provided is $P(w_i \mid t_i)$ that a word w_i has a part of speech tag t_i , e.g. $P(\text{flies} \mid V) = \frac{1}{2}$.

The part of speech tag definitions are: bos (*begin sentence marker*), N (*noun*), V (*verb*), D (*determiner*), P (*preposition*), eos (*end of sentence marker*).

$P(t_i \mid t_{i-2}, t_{i-1})$	t_{i-2}	t_{i-1}	t_i
1	bos	bos	N
$\frac{1}{2}$	bos	N	N
$\frac{1}{2}$	bos	N	V
$\frac{1}{2}$	N	N	V
$\frac{1}{2}$	N	N	P
$\frac{1}{3}$	N	V	D
$\frac{1}{3}$	N	V	V
$\frac{1}{3}$	N	V	P
1	V	D	N
1	V	V	D
1	N	P	D
1	V	P	D
1	P	D	N
1	D	N	eos

$P(w_i \mid t_i)$	t_i	w_i
1	D	an
$\frac{2}{5}$	N	time
$\frac{2}{5}$	N	arrow
$\frac{1}{5}$	N	flies
1	P	like
$\frac{1}{2}$	V	like
$\frac{1}{2}$	V	flies
1	eos	eos
1	bos	bos

- a. Consider a Jelinek-Mercer style interpolation smoothing scheme for $P(w_i \mid t_i)$:

$$P_{jm}(w_i \mid t_i) = \Lambda[t_i] \cdot P(w_i \mid t_i) + (1 - \Lambda[t_i]) \cdot P(w_i)$$

Λ is an array with a value $\Lambda[t_i]$ for each part of speech tag t_i , such that $0 \leq \Lambda[t_i] \leq 1$. Provide a condition on Λ that must be satisfied to ensure that P_{jm} is a well-defined probability model.

Answer: Because of the following fact about $P(w_i \mid t_i)$:

$$\sum_{w_i} P(w_i \mid t_i) = 1$$

and in P_{jm} we are given t_i , so to interpolate with $P(w_i)$ the following condition has to hold:

$$\sum_{t_i} \Lambda[t_i] = 1$$

- b. Provide a Hidden Markov Model (*hmm*) that uses the trigram part of speech probability $P(t_i | t_{i-2}, t_{i-1})$ as the transition probability $P_{hmm}(s_j | s_k)$ and the probability $P(w_i | t_i)$ as the emission probability $P_{hmm}(w_j | s_j)$.

Important: Provide the *hmm* in the form of two tables as shown below. The first table contains transitions between states in the *hmm* and the transition probabilities and the second table contains the words emitted at each state and the emission probabilities. Do not provide entries with zero probability.

from-state s_k	to-state s_j	$P(s_j s_k)$	state s_j	emission w	$P(w s_j)$

Hint: In your *hmm* the state $\langle N, eos \rangle$ will have emission of word *eos* with probability 1 and will not have transitions to any other states.

Answer: Here are the two tables that define the HMM, the transition table on the left and the emission table on the right:

from-state s_k	to-state s_j	$P(s_j s_k)$
<i>bos, bos</i>	<i>bos, N</i>	$P(N bos, bos)$
<i>bos, N</i>	<i>N, N</i>	$P(N bos, N)$
<i>bos, N</i>	<i>N, V</i>	$P(V bos, N)$
<i>N, N</i>	<i>N, V</i>	$P(V N, N)$
<i>N, N</i>	<i>N, P</i>	$P(P N, N)$
<i>N, V</i>	<i>V, D</i>	$P(D N, V)$
<i>N, V</i>	<i>V, V</i>	$P(V N, V)$
<i>N, V</i>	<i>V, P</i>	$P(P N, V)$
<i>V, D</i>	<i>D, N</i>	$P(N V, D)$
<i>V, V</i>	<i>V, D</i>	$P(D V, V)$
<i>N, P</i>	<i>P, D</i>	$P(D N, P)$
<i>V, P</i>	<i>P, D</i>	$P(D V, P)$
<i>P, D</i>	<i>D, N</i>	$P(N P, D)$
<i>D, N</i>	<i>N, eos</i>	$P(eos D, N)$

state s_j	emission w	$P(w s_j)$
<i>bos, bos</i>	<i>bos</i>	1
<i>bos, N</i>	time	$\frac{1}{2}$
<i>bos, N</i>	arrow	$\frac{1}{3}$
<i>bos, N</i>	flies	$\frac{1}{3}$
<i>N, N</i>	time	$\frac{1}{2}$
<i>N, N</i>	arrow	$\frac{1}{2}$
<i>N, N</i>	flies	$\frac{1}{2}$
<i>N, V</i>	like	$\frac{1}{2}$
<i>N, V</i>	flies	$\frac{1}{2}$
<i>V, D</i>	an	1
<i>V, V</i>	like	$\frac{1}{2}$
<i>V, V</i>	flies	$\frac{1}{2}$
<i>N, P</i>	like	1
<i>V, P</i>	like	1
<i>P, D</i>	an	1
<i>D, N</i>	time	$\frac{1}{2}$
<i>D, N</i>	arrow	$\frac{1}{2}$
<i>D, N</i>	flies	$\frac{1}{2}$

- c. Based on your *hmm* constructed in 2b. what is the state sequence that would be provided by the Viterbi algorithm for the following input sentence:

bos bos time flies like an arrow eos

Answer:

Note that the only ambiguous words are *flies* (could be *N* or *V*) and *like* (could be *V* or *P*) and so all you need to do is compare the scores for the following sub-sequence. The bold-faced outcome wins for this sub-sequence which determines the best state sequence for the entire input.

flies	like	
(N, V)	(V, P)	$\frac{1}{2} \times \frac{1}{3} \times \mathbf{1}$
(N, V)	(V, V)	$\frac{1}{2} \times \frac{1}{3} \times \frac{1}{2}$
(N, N)	(N, P)	$\frac{1}{2} \times \frac{1}{2} \times 1$
(N, N)	(N, V)	$\frac{1}{2} \times \frac{1}{2} \times \frac{1}{2}$

Since the best state sequence is then (bos,N)–(N,V)–(V,P)–(P,D)–(D,N)–(N,eos) the output best state sequence will be *bos/bos, time/N, flies/V, like/P, an/D, arrow/N, eos/eos*.

Answer: The full table is given below but you do not need to compute the entire table to solve this

question.	bos	time	flies	like	an	arrow	eos	
	(bos,bos)	(bos,N)	(N,V)	(V,P)	(P,D)	(D,N)	(N,eos)	
	1	$\times 1 \times \frac{2}{5}$	$\times \frac{1}{2} \times \frac{1}{2}$	$\times \frac{1}{3} \times 1$	$\times 1 \times 1$	$\times 1 \times \frac{2}{5}$	$\times 1 \times 1$	$= \frac{1}{75}$
	(bos,bos)	(bos,N)	(N,V)	(V,V)	(V,D)	(D,N)	(N,eos)	
	1	$\times 1 \times \frac{2}{5}$	$\times \frac{1}{2} \times \frac{1}{2}$	$\times \frac{1}{3} \times \frac{1}{2}$	$\times 1 \times 1$	$\times 1 \times \frac{2}{5}$	$\times 1 \times 1$	$= \frac{1}{150}$
	(bos,bos)	(bos,N)	(N,N)	(N,P)	(P,D)	(D,N)	(N,eos)	
	1	$\times 1 \times \frac{2}{5}$	$\times \frac{1}{2} \times \frac{1}{5}$	$\times \frac{1}{2} \times 1$	$\times 1 \times 1$	$\times 1 \times \frac{2}{5}$	$\times 1 \times 1$	$= \frac{1}{125}$
	(bos,bos)	(bos,N)	(N,N)	(N,V)	(V,D)	(D,N)	(N,eos)	
	1	$\times 1 \times \frac{2}{5}$	$\times \frac{1}{2} \times \frac{1}{5}$	$\times \frac{1}{2} \times \frac{1}{2}$	$\times \frac{1}{3} \times 1$	$\times 1 \times \frac{2}{5}$	$\times 1 \times 1$	$= \frac{1}{750}$