

### ASSIGNMENT 4.3

(a)

tokens	unigram probability
A	0.018407
AND	0.017863
AT	0.004313
AS	0.003992
AN	0.002999
ARE	0.002990
ABOUT	0.001926
AFTER	0.001347
ALSO	0.001310
ALL	0.001182
A.	0.001026
ANY	0.000632
AMERICAN	0.000612
AGAINST	0.000596
ANOTHER	0.000428
AMONG	0.000374
AGO	0.000357
ACCORDING	0.000348
AIR	0.000311
ADMINISTRATION	0.000292
AGENCY	0.000280
AROUND	0.000277
AGREEMENT	0.000263
AVERAGE	0.000259
ASKED	0.000258
ALREADY	0.000249
AREA	0.000231
ANALYSTS	0.000226
ANNOUNCED	0.000227
ADDED	0.000221
ALTHOUGH	0.000214
AGREED	0.000212
APRIL	0.000207
AWAY	0.000202

(b)

Most likely words	bigram probability
<UNK>	0.615020
U.	0.013372
FIRST	0.011720
COMPANY	0.011659
NEW	0.009451

(c)

$L_u : -64.509440$

$L_b : -44.740469$

The bigram model has a higher log-likelihood.

(d)

$L_u : -41.643460$

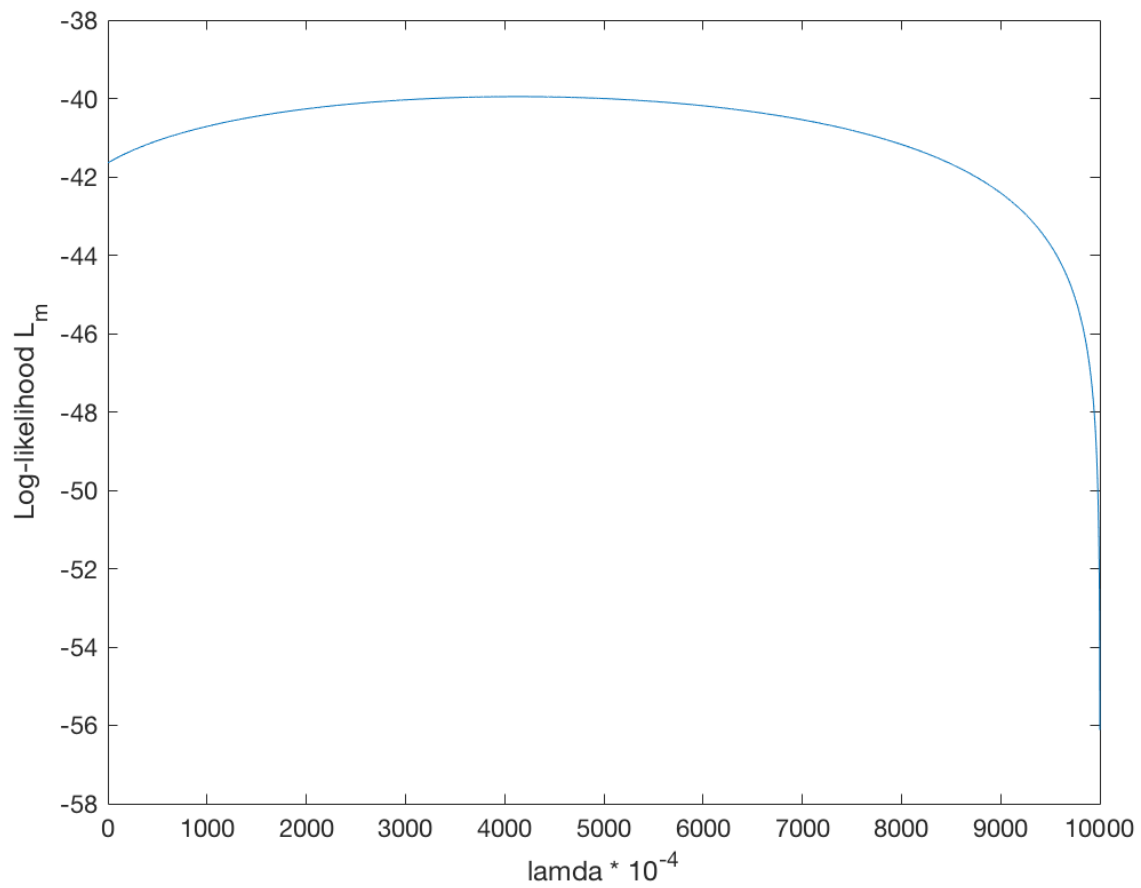
$L_b : -\infty$

"NINETEEN OFFICIALS" and "SOLD FIRE" are not observed in the training corpus.

The unobserved pairs have caused some likelihoods calculated in the bigram model turned to 0. Thus, the total log-likelihood of bigram model has been  $\log(0)$ , which is invalid. So we defined the result as  $-\infty$  in this case.

(e)

The optimal value:  $\lambda = 0.41$ , where  $L_m = -39.95$



# 1 Source Code

Listing 1: hw4.3.m

```
1 clear all
2
3 [words] = textread('vocab.txt', '%s');
4 [counts] = textread('unigram.txt', '%d');
5 [w1,w2,count] = textread('bigram.txt', '%d_%d_%d');
6
7 %% ——— (a) ———
8
9 P_unigram = zeros(size(counts));
10 total = sum(counts);
11 for i = 1 : length(counts)
12     P_unigram(i) = counts(i)/total;
13 end
14 % sorted by frequency in ascending order
15 i = 1;
16 fprintf('———_(a)_____\\n');
17 fprintf('tokens:\\t\\t_unigram_probability:\\t\\n');
18 while(i < 500)
19     if(words{i}(1) == 'A')
20         fprintf('%s\\t\\t_%f\\n', words{i}, P_unigram(i));
21     end
22     i = i+1;
23 end
24
25 %% ——— (b) ———
26
27 P_bigram = zeros(size(w1));
28 total = sum(count);
29
30 for i = 1 : length(w1)
31     P_bigram(i) = count(i)/counts(w1(i));
32 end
33 Bigram = [w1 w2 count P_bigram];
34 % extract tuples where w1 == THE
35 tmp_the = [];
36 for i = 1 : length(w1)
37     if(Bigram(i,1) == 4)
38         tmp_the = [tmp_the; Bigram(i,:)];
39     end
40 end
41 % sorting
42 [sorted_p, index] = sort(tmp_the(:,4), 'descend');
```



```

89 P_b_2 = zeros(size(sentence_2));
90 indexes_2(1) = 2;
91 for i = 2 : length(sentence_2)
92     for j = 1 : length(words)
93         if(strcmpi(words{j}, sentence_2(i)))
94             P_u_2(i) = P_unigram(j);
95             indexes_2(i) = j;
96         end
97     end
98 end
99 fprintf('———_(d)———\n');
100 for i = 2 : length(sentence_2)
101     for j = 1 : length(w1)
102         if(w1(j) == indexes_2(i-1) && w2(j) == indexes_2(i))
103             P_b_2(i) = P_bigram(j);
104         end
105     end
106 end
107 end
108 fprintf('two_pairs:\n');
109 for i = 2 : length(P_b_2)
110     if(P_b_2(i) == 0)
111         fprintf('w1_: %s\t w2_: %s\n', words{indexes_2(i-1)}, words{indexes_2(i)});
112     end
113 end
114 L_u_2 = 0;
115
116 for i = 2 : length(sentence_2)
117     L_u_2 = L_u_2 + log(P_u_2(i));
118 end
119
120 fprintf('L_u: %f\n', L_u_2);
121
122 %% ——— (e) ———
123
124 lamda = 0 : 0.0001 : 1;
125 L_m = zeros(size(lamda));
126 for i = 1 : length(lamda)
127     P_m = (1-lamda(i))*P_u_2 + lamda(i)*P_b_2;
128     for j = 2 : length(sentence_2)
129         L_m(i) = L_m(i) + log(P_m(j));
130     end
131 end
132 plot(L_m);
133 xlabel('lamda*10^{-4}');
134 ylabel('Log-likelihood L_m');

```

```

135 hold on
136 [ymax indmax]=max(L_m);
137 xmax=lamda(indmax);
138 fprintf('———(e)———\n');
139 fprintf('max_lamda=%4.2f , L_m=%4.2f\n',xmax,ymax);

```

*Submitted by Xiaowen Mao on Oct 25.*