

CSE 250A. Assignment 4

Hao-en Sung (A53204772)
 wrangle1005@gmail.com

October 23, 2016

4.1 Maximum likelihood estimation of a multinomial distribution

(a) Log-likelihood

Sol.

$$L = \log P(\text{data}) = \log \left(\prod_{t=1}^T P(X_t) \right) \quad (1)$$

$$= \sum_{n=1}^{2N} \sum_{t=1}^{C_n} \log P(X_i = n) \quad (2)$$

$$= \sum_{n=1}^{2N} C_n \log p_n \quad (3)$$

□

(b) Maximum likelihood estimate

Sol. After apply the Lagrange framework, I have

$$L := C_n \log p_n - \lambda \left(\sum_{n=1}^{2N} p_n - 1 \right), \quad (4)$$

$$\frac{\partial L}{\partial p_n} = \frac{C_n}{p_n} - \lambda. \quad (5)$$

At optimal, it can be derived that

$$p_n = \frac{C_n}{\lambda}. \quad (6)$$

Based on the original constraint, it is said

$$\sum_{n=1}^{2N} p_n = \frac{\sum_{n=1}^{2N} C_n}{\lambda} = 1 \quad (7)$$

$$\lambda = \sum_{n=1}^{2N} C_n. \quad (8)$$

Replace λ with $\sum_{n=1}^{2N} C_n$ in Eq. 6, then it is clear that $p_n = \frac{C_n}{\sum_{n=1}^{2N} C_n}$. Since C_n is nonnegative, p_n is also nonnegative. □

(c) Even versus odd

Sol.

$$\begin{aligned} & \sum_{n=1}^{2N} (-1)^n p_n = 0 \\ \Rightarrow & - \sum_{n=1}^N p_{2n-1} + \sum_{n=1}^N p_{2n} = 0 \\ \Rightarrow & \sum_{n=1}^N p_{2n-1} = \sum_{n=1}^N p_{2n} \end{aligned}$$

□

(d) Maximum likelihood estimate

Sol. From (c) I know that $\sum_{n=1}^N p_{2n-1} = \sum_{n=1}^N p_{2n} = \frac{1}{2}$. Beside that, I can rewrite Eq. 3 in (a) to $\sum_{n=1}^N C_{2n-1} \log p_{2n-1} + \sum_{n=1}^N C_{2n} \log p_{2n}$. Since odd die and even die are independently constrained, I can downscale origin problem to two sub-problems. Thus, I need to solve

$$\begin{aligned} \max_{p_{2n-1}} \sum_{n=1}^N C_{2n-1}, \quad \text{s.t.} \quad \sum_{n=1}^N p_{2n-1} &= \frac{1}{2}, \\ \max_{p_{2n}} \sum_{n=1}^N C_{2n}, \quad \text{s.t.} \quad \sum_{n=1}^N p_{2n} &= \frac{1}{2}. \end{aligned}$$

Follow similar step from Eq. 4 to 8, I can get

$$\begin{aligned} p_{2n-1} &= \frac{C_{2n-1}}{\sum_{n=1}^N C_{2n-1}} \\ p_{2n} &= \frac{C_{2n}}{\sum_{n=1}^N C_{2n}}. \end{aligned}$$

□

4.2 Maximum likelihood estimation in belief networks

(a) Express the maximum likelihood estimates for the CPTs in G_1 in terms of these counts.

Sol. From figure for G_1 , I know CPT tables are

$$\begin{aligned} P_1(X_1 = x_1) &= \frac{P_1(X_1 = x_1)}{P_1(X_1)} \\ P_{i+1}(X_{i+1} = x_{i+1} | X_i = x_i) &= \frac{P_i(X_i = x_i, X_{i+1} = x_{i+1})}{P_i(X_i = x_i)}, \quad \forall i \in \{1, \dots, n-1\}. \end{aligned}$$

Take advantage of the conclusion from Section 4.1, I am aware that the probability p_n maximize log-likelihood when it equals to $\frac{C_n}{\sum_{n=1}^N C_n}$.

Similarly, $P(G_1)$ is maximized as $P^*(G_1)$ when

$$\begin{aligned} P_1(X_1 = x_1) &= \frac{\text{COUNT}_1(X_1 = x_1)}{\text{COUNT}_1(X_1)} = \frac{\text{COUNT}_1(X_1 = x_1)}{T} \\ P_{i+1}(X_{i+1} = x_{i+1} | X_i = x_i) &= \frac{\text{COUNT}_i(X_i = x_i, X_{i+1} = x_{i+1})}{\text{COUNT}_i(X_i = x_i)}, \quad \forall i \in \{1, \dots, n-1\}. \end{aligned}$$

□

(b) Express the maximum likelihood estimates for the CPTs in G_2 in terms of these counts.

Sol. Similar to (a), I know CPT tables are

$$\begin{aligned} P_n(X_n = x_n) &= \frac{P_n(X_n = x_n)}{P_n(X_n)} \\ P_i(X_i = x_i | X_{i+1} = x_{i+1}) &= \frac{P_i(X_i = x_i, X_{i+1} = x_{i+1})}{P_{i+1}(X_{i+1} = x_{i+1})}, \quad \forall i \in \{1, \dots, n-1\}. \end{aligned}$$

When $P(G_2)$ is maximized as $P^*(G_2)$,

$$\begin{aligned} P_n(X_n) &= \frac{\text{COUNT}_n(X_n = x_n)}{\text{COUNT}_n(X_n)} = \frac{\text{COUNT}_n(X_n = x_n)}{T} \\ P_i(X_i = x_i | X_{i+1} = x_{i+1}) &= \frac{\text{COUNT}_i(X_i = x_i, X_{i+1} = x_{i+1})}{\text{COUNT}_{i+1}(X_{i+1} = x_{i+1})}, \quad \forall i \in \{1, \dots, n-1\}. \end{aligned}$$

□

(c) Using your answers from parts (a) and (b), show that the maximum likelihood CPTs for G_1 and G_2 from this data set give rise to the same joint distribution over the nodes $\{X_1, X_2, \dots, X_n\}$.

Sol. It can be derived that

$$\begin{aligned}
P(G_1 = g_1) &= P_1(X_1 = x_1) \cdot \prod_{i=1}^{n-1} P_{i+1}(X_{i+1} = x_{i+1} | X_i = x_i) \\
&= \frac{\text{COUNT}_1(X_1 = x_1)}{T} \cdot \prod_{i=1}^{n-1} \frac{\text{COUNT}_i(X_i = x_i, X_{i+1} = x_{i+1})}{\text{COUNT}_i(X_i = x_i)} \\
&= \frac{\prod_{i=1}^{n-1} \text{COUNT}_i(X_i = x_i, X_{i+1} = x_{i+1})}{T \cdot \prod_{i=2}^{n-1} \text{COUNT}_i(X_i = x_i)} \\
&= \frac{\prod_{i=1}^{n-1} \text{COUNT}_i(X_i = x_i, X_{i+1} = x_{i+1})}{T \cdot \prod_{i=1}^{n-2} \text{COUNT}_{i+1}(X_{i+1} = x_{i+1}) \cdot \text{COUNT}_n(X_n = x_n)} \cdot \text{COUNT}_n(X_n = x_n) \\
&= \prod_{i=1}^{n-1} \frac{\text{COUNT}_i(X_i = x_i, X_{i+1} = x_{i+1})}{\text{COUNT}_{i+1}(X_{i+1} = x_{i+1})} \cdot \frac{\text{COUNT}_n(X_n)}{T} \\
&= \prod_{i=1}^{n-1} \frac{P_i(X_i = x_i, X_{i+1} = x_{i+1})}{P_{i+1}(X_{i+1} = x_{i+1})} \cdot P_n(X_n = x_n) = P(G_2 = g_2).
\end{aligned}$$

□

(d) Suppose that some but not all of the edges in these DAGs were reversed, as in the graph G_3 shown below. Would the maximum likelihood CPTs for G_3 also give rise to the same joint distribution? (Hint: does G_3 imply all the same statements of conditional independence as G_1 and G_2 ?)

Sol. One can observe that in (c) there is a T in denominator, which is caused by the one root node — either X_1 in G_1 or X_n in G_2 . However, in G_3 there are two root nodes X_2 and X_{n-1} , which generates two T in denominator. It is clear that G_3 will not have the same joint probability as G_1 or G_2 . □

4.3 Statistical language modeling

(a) Compute the maximum likelihood estimate of the unigram distribution $P_u(w)$ over words w . Print out a table of all the tokens (i.e., words) that start with the letter "A", along with their numerical unigram probabilities (not counts). (You do not need to print out the unigram *probabilities* for all 500 tokens.)

Sol. Unigram probability for words starting with character "A" is shown as Table 1.

□

(b) Compute the maximum likelihood estimate of the bigram distribution $P_b(w'|w)$. Print out a table of the five most likely words to follow the word "THE", along with their numerical bigram probabilities.

Sol. Bigram probability for words starting with word "THE" is shown as Table 2.

□

(c) Consider the sentence "Last week the stock market fell by one hundred points." Ignoring punctuation, compute and compare the log-likelihoods of this sentence under the unigram and bigram models. In the equation for the bigram log-likelihood, the token $\langle s \rangle$ is used to mark the beginning of a sentence. Which model yields the highest log-likelihood?

Sol.

$$\begin{aligned}
L_u &= -64.509440 \\
L_b &= -44.740469
\end{aligned}$$

It is clear that the bigram model yields higher probability.

□

(d) Consider the sentence "Last week the stock market fell by one hundred points." Ignoring punctuation, compute and compare the log-likelihoods of this sentence under the unigram and bigram models. Which pairs of adjacent words in this sentence are not observed in the training corpus? What effect does this have on the log-likelihood from the bigram model?

Sol.

$$\begin{aligned}
L_u &= -41.643460 \\
L_b &= -\infty
\end{aligned}$$

The pairs (NINETEEN, OFFICIALS) and (SOLD, FIRE) do not show up in corpus, which makes $\log(0)$ and pushes the overall log probability to $-\infty$. □

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

Table 1: Unigram Probability for words start with character "A"

<UNK>	0.615020
U.	0.013372
FIRST	0.011720
COMPANY	0.011659
NEW	0.009451

Table 2: Bigram Probability for words start with word "THE"

(e) Consider the so-called mixture model that predicts words from a weighted interpolation of the unigram and bigram models. Compute and plot the value of this log-likelihood L_m as a function of the parameter $\lambda \in [0, 1]$. From your results, deduce the optimal value of λ to two significant digits.

Sol. The maximum probability of mixed unigram and bigram algorithm is -39.953680 , which appears at $\lambda = 0.41$. I also examined different λ with 0.01 precision, as shown in Fig. 1.

□

(f) Submit a hard copy of your source code for the previous parts of this problem.

I first use C++ to write down all main functions to solve (a) to (d). For (e), I use C++ generate probabilities for mixed unigram and bigram model with various λ , then use MATLAB to render the figure. Both C++ and MATLAB are included in Appendix as 1 and 2, respectively.

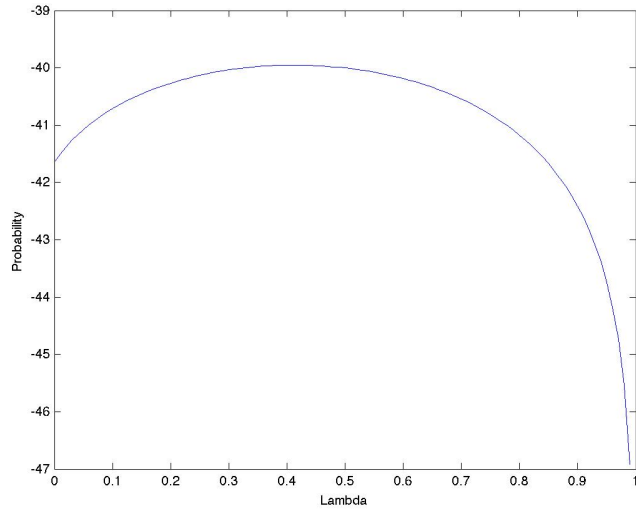


Figure 1: Mixed unigram and bigram model with different lambdas

Appendix

```

1 #include <stdio>
2 #include <stdlib>
3 #include <cmath>
4 #include <cstdio>
5 #include <vector>
6 #include <utility>
7 #include <string>
8 #include <unordered_map>
9 #include <algorithm>
10 #include <assert.h>
11
12 using namespace std;
13
14 const int MAXL = 10000;
15 const int NUMV = 500;
16 const int NUMS = 100;
17
18 char in[MAXL + 10];
19
20 unordered_map<string, int> umsi;
21 unordered_map<int, string> umis;
22
23 vector<string> vocab(NUMV, "");
24 vector<int> uni(NUMV, 0);
25 vector<vector<int>> bi(NUMV, vector<int>(NUMV, 0));
26
27 // helper function to calculate prob with different lambda
28 double callProb(vector<string>& lst, double lambda, bool warning) {
29     int size = lst.size();
30
31     int usum = 0;
32     for (int i = 0; i < NUMV; i++) {
33         usum += uni[i];
34     }
35
36     double lprob = 0;
37     for (int i = 0; i < size; i++) {
38         int tar = i == 0 ? umsi.at("<s>") : umsi.at(lst[i - 1]);
39         int bsum = 0;
40         for (int j = 0; j < NUMV; j++) {
41             bsum += bi[tar][j];
42         }
43
44         if (bi[tar][umsi.at(lst[i])] == 0 and lambda != 0 and warning) {
45             fprintf(stderr, "Warning: (%s, %s) does not appear in documents\n",
46                 umis.at(tar).c_str(), lst[i].c_str());
47         }
48
49         lprob += log((1 - lambda) * uni[umsi.at(lst[i])] / usum +
50             lambda * bi[tar][umsi.at(lst[i])] / bsum);
51     }
52
53     return lprob;
54 }
55
56 int main() {
57     { // read vocabulary
58         string FN = "../dat/vocab.txt";
59         FILE *pf = fopen(FN.c_str(), "r");
60         if (pf == NULL) {
61             fprintf(stderr, "Vocab file cannot read\n");

```

```

62         exit(EXIT_FAILURE);
63     }
64
65     int cnt = 0;
66     while (fscanf(pf, "%s", in) != EOF) {
67         umsi[in] = cnt;
68         umis[cnt] = in;
69         vocab[cnt++] = in;
70     }
71
72     fclose(pf);
73     assert(vocab.size() == NUMV);
74 }
75
76 { // read unigram
77     string FN = "../dat/unigram.txt";
78     FILE *pf = fopen(FN.c_str(), "r");
79     if (pf == NULL) {
80         fprintf(stderr, "Unigram file cannot read\n");
81         exit(EXIT_FAILURE);
82     }
83
84     int d;
85     int cnt = 0;
86     while (fscanf(pf, "%d", &d) != EOF) {
87         uni[cnt++] = d;
88     }
89
90     fclose(pf);
91     assert(cnt == NUMV);
92 }
93
94 { // read bigram
95     string FN = "../dat/bigram.txt";
96     FILE *pf = fopen(FN.c_str(), "r");
97     if (pf == NULL) {
98         fprintf(stderr, "Bigram file cannot read\n");
99         exit(EXIT_FAILURE);
100     }
101
102     int u, v, d;
103     while (fscanf(pf, "%d\\t%d\\t%d", &u, &v, &d) != EOF) {
104         u -= 1;
105         v -= 1;
106         assert(u >= 0 and u < NUMV and v >= 0 and v < NUMV);
107         bi[u][v] = d;
108     }
109
110     fclose(pf);
111 }
112
113 { // 4.3 (a)
114     int usum = 0;
115     for (int i = 0; i < NUMV; i++) {
116         usum += uni[i];
117     }
118     printf("4.3 (a)\\n");
119     for (int i = 0; i < NUMV; i++) {
120         if (vocab[i][0] == 'A') {
121             printf("%s %f\\n", vocab[i].c_str(), 1.0 * uni[i]/usum);
122         }
123     }
124     printf("\\n\\n");
125 }
126
127 { // 4.3 (b)
128     int tar = umsi.at("THE");
129     int sum = 0;
130     vector<pair<int, int>> vct;
131     for (int j = 0; j < NUMV; j++) {
132         sum += bi[tar][j];
133         vct.emplace_back(bi[tar][j], j);
134     }
135     sort(vct.rbegin(), vct.rend());
136     printf("4.3 (b)\\n");
137     for (int i = 0; i < 5; i++) {
138         printf("%s %f\\n", vocab[vct[i].second].c_str(), 1.0 * vct[i].first/sum);
139     }
140     printf("\\n\\n");
141 }
142
143 { // 4.3 (c)
144     vector<string> lst{"LAST", "WEEK", "THE", "STOCK", "MARKET", "FELL", "BY", "ONE", "HUNDRED", "POINTS"};
145     printf("4.3 (c)\\n");
146     { // I. unigram
147         double lprob = callLProb(lst, 0, true);
148         printf("unigram: %f\\n", lprob);
149     }
150     { // II. bigram
151         double lprob = callLProb(lst, 1, true);

```

```

152         printf("bigram: %f\n", lprob);
153     }
154     printf("\n\n");
155 }
156
157 { // 4.3 (d)
158     vector<string> lst{"THE", "NINETEEN", "OFFICIALS", "SOLD", "FIRE", "INSURANCE"};
159     printf("4.3 (d)\n");
160     { // I. unigram
161         double lprob = callLProb(lst, 0, true);
162         printf("unigram: %f\n", lprob);
163     }
164     { // II. bigram
165         double lprob = callLProb(lst, 1, true);
166         printf("bigram: %f\n", lprob);
167     }
168     printf("\n\n");
169 }
170
171 { // 4.3 (e)
172     string FN = "../res/lambda_data.csv";
173     FILE* pf = fopen(FN.c_str(), "w");
174     if (pf == NULL) {
175         fprintf(stderr, "Lambda file cannot write\n");
176         exit(EXIT_FAILURE);
177     }
178
179     vector<string> lst{"THE", "NINETEEN", "OFFICIALS", "SOLD", "FIRE", "INSURANCE"};
180     double mmax = -DBL_MAX;
181     double pmax = -1;
182     for (int i = 0; i < NUMS; i++) {
183         double pos = 1.0 * i / NUMS;
184         double lprob = callLProb(lst, pos, false);
185         if (lprob > mmax) {
186             mmax = lprob;
187             pmax = pos;
188         }
189         fprintf(pf, "%f, %f\n", pos, lprob);
190     }
191
192     printf("4.3 (d)\n");
193     printf("max prob: %f (at %.2f)", mmax, pmax);
194     fclose(pf);
195 }
196 }

```

Listing 1: Main Code

```

1 % read in data
2 M = csvread('../res/lambda_data.csv');
3
4 % plot figure
5 res = figure('visible', 'off');
6 plot(M(:,1), M(:,2));
7 saveas(res, '../res/lambda.jpg');

```

Listing 2: Draw Lambda Plot