# Information Extraction Project 3

Tianyi Chen*

May 11, 2017

## 1 Code compile and usage

A **C++** program has been implemented for handling this project. To compile the **C++** program, in project directory, run the below command in terminal.

```
1 make
```

Successful compiling can generate an executable program called **test_tree**. **test_tree** can deliver all required experiments with suitable setting. The usage of running **test_tree** is as the below:

```
1 ./test_tree -m mode
```

where option **-m** can control the which part of experiment to run, feasible values are in $\{1, 2, 3\}$

- **1**: Part 1: Construct agglomerative clustering tree.

- **2**: Part 2: Construct Bit-encoding based decision tree.

- **3**: Part 3: Construct Chou's decision tree.

## 2 Agglomerative clustering tree

To complete the required experiment, run the below command in terminal

```
1 ./test_tree -m 1
```

We get the below output:

```
1  Initialize agglomerative clustering tree...
2  Agglomerative clustering tree is growing...
3  Print out growed agglomerative clustering tree...
4  level is 9,  letters: o, a, u, i, e,  , z, v, k, m, f, q, j, b, p, w, c, t, h, x, r, l, n, y, s, g, d,
5  level is 4,  letters: o, a, u, i, e,  ,
6  level is 8,  letters: z, v, k, m, f, q, j, b, p, w, c, t, h, x, r, l, n, y, s, g, d,
7  level is 3,  letters: o, a, u, i, e,
8  level is 0,  letters:  ,
9  level is 7,  letters: z, v, k, m, f, q, j, b, p, w, c, t, h,
10 level is 4,  letters: x, r, l, n, y, s, g, d,
11 level is 2,  letters: o, a, u, i,
12 level is 0,  letters: e,
13 level is 6,  letters: z, v, k, m, f, q, j, b, p, w, c, t,
14 level is 0,  letters: h,
15 level is 3,  letters: x, r, l, n,
16 level is 2,  letters: y, s, g, d,
17 level is 1,  letters: o, a,
18 level is 1,  letters: u, i,
19 level is 5,  letters: z, v, k, m, f, q, j, b, p,
20 level is 2,  letters: w, c, t,
21 level is 2,  letters: x, r, l,
```

---

*Email: tchen59@jhu.edu.

```
22  level is 0,   letters: n,
23  level is 1,   letters: y, s,
24  level is 1,   letters: g, d,
25  level is 0,   letters: o,
26  level is 0,   letters: a,
27  level is 0,   letters: u,
28  level is 0,   letters: i,
29  level is 4,   letters: z, v, k, m, f,
30  level is 3,   letters: q, j, b, p,
31  level is 1,   letters: w, c,
32  level is 0,   letters: t,
33  level is 1,   letters: x, r,
34  level is 0,   letters: l,
35  level is 0,   letters: y,
36  level is 0,   letters: s,
37  level is 0,   letters: g,
38  level is 0,   letters: d,
39  level is 3,   letters: z, v, k, m,
40  level is 0,   letters: f,
41  level is 2,   letters: q, j, b,
42  level is 0,   letters: p,
43  level is 0,   letters: w,
44  level is 0,   letters: c,
45  level is 0,   letters: x,
46  level is 0,   letters: r,
47  level is 2,   letters: z, v, k,
48  level is 0,   letters: m,
49  level is 1,   letters: q, j,
50  level is 0,   letters: b,
51  level is 1,   letters: z, v,
52  level is 0,   letters: k,
53  level is 0,   letters: q,
54  level is 0,   letters: j,
55  level is 0,   letters: z,
56  level is 0,   letters: v,
```

For the questions shown in the description of this project. The below are the answers.

## 2.1   Drawing clustering tree

The clustering tree after growing is shown in Figure 1 based on the information of output shown in Section 2.
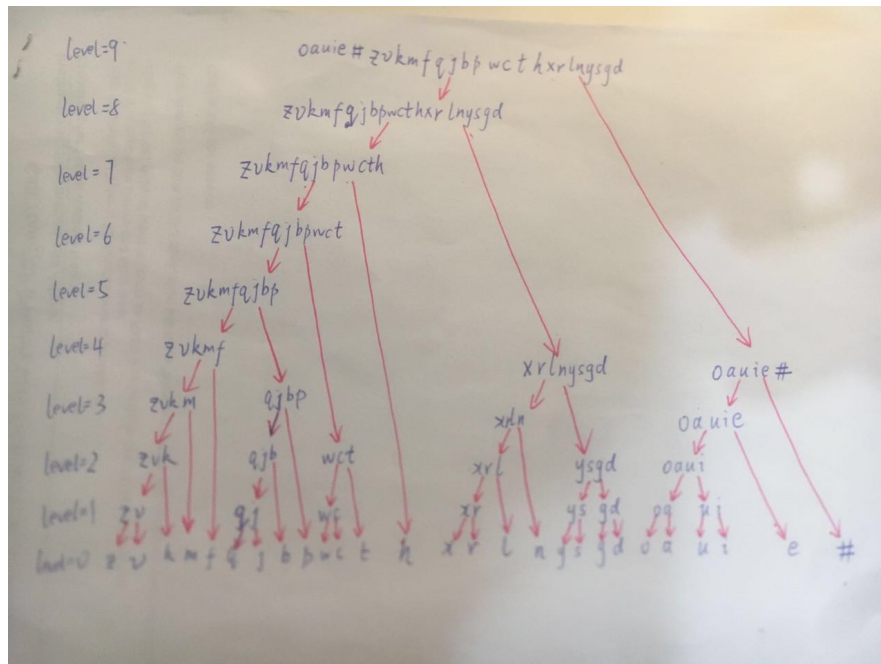


Figure 1: Clustering tree

## 2.2   The best 2-way clustering of letters

As shown in Figure 1, the best 2-way clustering of letters are $\{z, v, k, m, f, q, j, b, p, w, c, t, h, x, r, l, n, y, s, g, d\}$ and $\{o, a, u, i, e, \#\}$.

In Project 1, we found that in 2-state HMM, vowel and consonant have the same state with large probabilities. The best 2-way clustering is consistent with this observation in Project 1.

In Homework #9, we found that $a, e, i, o, u, \#$ have larger probabilities to belong to the same set by Chou's algorithm. The best 2-way clustering is also consistent with the observation in Homework #9.

## 2.3   The best 4-way clustering of letters

By the output shown in Section 2, we have the best 4-way clustering of letters is $\{o, a, u, i, e\}$, $\{\#\}$, $\{z, v, k, m, f, q, j, b, p, w, c, t, h\}$, and $\{x, r, l, n, y, s, g, d\}$.

It is also consistent with the observation in 4-state HMM in Project 1 that the letters in the same set have higher probabilities with the same state in 4-state HMM.

# 3   Bit-encoding based decision tree

To complete the required experiment, run the below command in terminal

```
./test_tree -m 2
```

```
Constructing agglomerative clustering tree...
Construct bit encoding decision tree...
Encoding 11, 12, 13, 14...
Build development, and held-out data sets...
Bit encoding decision tree is growing...
Process: ............................................................................................................
No frontier node...
Bit tree growing complete...
Print out growed bit encoding decision tree...
path: , questioned bit 18
path: 0, questioned bit 19
path: 1, questioned bit 19
path: 00, questioned bit 20
path: 01, questioned bit 9
path: 10, questioned bit 9
path: 11, questioned bit 9
path: 000, questioned bit 21
path: 001, questioned bit 9
path: 010, questioned bit 0
path: 011, questioned bit 10
path: 100, questioned bit 10
path: 101, questioned bit 10
path: 110, questioned bit 10
path: 111, questioned bit 20
path: 0000, questioned bit 9
path: 0001, questioned bit 9
path: 0010, questioned bit 10
path: 0011, questioned bit 10
path: 0100, questioned bit 1
path: 0101, questioned bit 1
path: 0110, questioned bit 0
path: 0111, questioned bit 0
path: 1000, questioned bit 0
path: 1001, questioned bit 20
path: 1010, questioned bit 20
path: 1011, questioned bit 20
path: 1100, questioned bit 20
path: 1101, questioned bit 20
path: 1110, questioned bit 10
path: 1111, questioned bit 21
path: 00000, questioned bit 22
path: 00001, questioned bit 22
path: 00010, questioned bit 10
path: 00011, questioned bit 22
```

```
 45  path: 00100, questioned bit 11
 46  path: 00101, questioned bit 0
 47  path: 00110, questioned bit 11
 48  path: 00111, questioned bit 11
 49  path: 01100, questioned bit 1
 50  path: 01101, questioned bit 11
 51  path: 10000, questioned bit 1
 52  path: 10001, questioned bit 11
 53  path: 10010, questioned bit 21
 54  path: 10011, questioned bit 0
 55  path: 10100, questioned bit 0
 56  path: 10101, questioned bit 0
 57  path: 11000, questioned bit 21
 58  path: 11001, questioned bit 21
 59  path: 11110, questioned bit 22
 60  path: 11111, questioned bit 22
 61  path: 000000, questioned bit 10
 62  path: 000001, questioned bit 10
 63  path: 000010, questioned bit 10
 64  path: 000011, questioned bit 0
 65  path: 000100, questioned bit 22
 66  path: 000101, questioned bit 22
 67  path: 000110, questioned bit 10
 68  path: 000111, questioned bit 10
 69  path: 001100, questioned bit 0
 70  path: 001101, questioned bit 0
 71  path: 001110, questioned bit 12
 72  path: 001111, questioned bit 0
 73  path: 100100, questioned bit 22
 74  path: 100101, questioned bit 22
 75  path: 110000, questioned bit 22
 76  path: 110001, questioned bit 0
 77  path: 110010, questioned bit 22
 78  path: 110011, questioned bit 22
 79  path: 0000000, questioned bit 11
 80  path: 0000001, questioned bit 0
 81  path: 0000010, questioned bit 11
 82  path: 0000011, questioned bit 0
 83  path: 0000100, questioned bit 0
 84  path: 0000101, questioned bit 11
 85  path: 0000110, questioned bit 1
 86  path: 0000111, questioned bit 1
 87  path: 0001110, questioned bit 11
 88  path: 0001111, questioned bit 0
 89  path: 0011100, questioned bit 13
 90  path: 0011101, questioned bit 0
 91  path: 1001000, questioned bit 23
 92  path: 1001001, questioned bit 23
 93  path: 1001010, questioned bit 23
 94  path: 1001011, questioned bit 0
 95  path: 1100000, questioned bit 23
 96  path: 1100001, questioned bit 0
 97  path: 1100010, questioned bit 1
 98  path: 1100011, questioned bit 11
 99  path: 1100100, questioned bit 11
100  path: 1100101, questioned bit 0
101  path: 1100110, questioned bit 11
102  path: 1100111, questioned bit 11
103  path: 00000100, questioned bit 12
104  path: 00000101, questioned bit 0
105  path: 00001100, questioned bit 10
106  path: 00001101, questioned bit 10
107  path: 00111000, questioned bit 0
108  path: 00111001, questioned bit 0
109  path: 10010010, questioned bit 24
110  path: 10010011, questioned bit 0
111  path: 11000000, questioned bit 0
112  path: 11000001, questioned bit 11
113  path: 11000010, questioned bit 1
114  path: 11000011, questioned bit 1
115  path: 11000100, questioned bit 2
116  path: 11000101, questioned bit 11
117  path: 11000110, questioned bit 12
118  path: 11000111, questioned bit 1
119  path: 11001110, questioned bit 0
120  path: 11001111, questioned bit 0
121  path: 000011010, questioned bit 11
122  path: 000011011, questioned bit 11
123  path: 001110000, questioned bit 1
124  path: 001110001, questioned bit 1
125  path: 100100100, questioned bit 0
```

```
126  path: 100100101, questioned bit 0
127  path: 110000010, questioned bit 12
128  path: 110000011, questioned bit 0
129  path: 110001100, questioned bit 13
130  path: 110001101, questioned bit 13
131  path: 110011100, questioned bit 1
132  path: 110011101, questioned bit 1
133  path: 0000110100, questioned bit 12
134  path: 0000110101, questioned bit 2
135  path: 1100000100, questioned bit 13
136  path: 1100000101, questioned bit 0
137  path: 00001101000, questioned bit 13
138  path: 00001101001, questioned bit 13
139  path: 11000001000, questioned bit 0
140  path: 11000001001, questioned bit 0
141  path: 000011010010, questioned bit 14
142  path: 000011010011, questioned bit 2
143  Calculate perplexity on testing data set...
144  perplexity is: 8.17174
```

## 3.1   Draw Bit-encode tree

Based on the output shown in Section 3, we draw the Bit-encode tree in Figure 2. Only the first five levels are drawn. The following levels can be drew similarly by the path information. The number in node circle is the questioned bit at current Tree node.
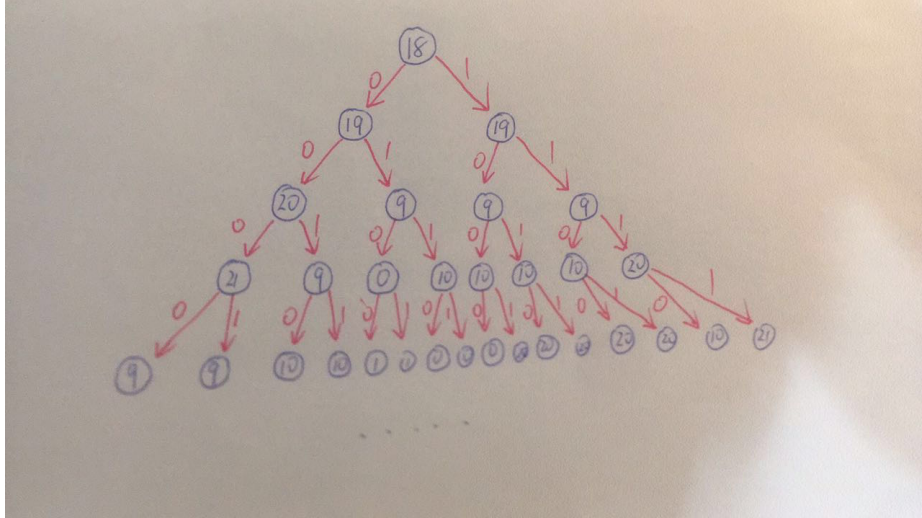


Figure 2: The first five levels Bit-encoded Tree

Based on the tree shown in 2, we can observe that $b_{11}$ is often asked after some $b_{3,j}$ has been asked. Therefore, the immediate preceding letter $l_3$ is more informative than the coarsest information about the letter $l_1$.

## 3.2   Compute Perplexity

The perplexity on test dataset is computed as well, which is shown in the output in Section 3. The perplexity is 8.17174 with default threshold 0.005.

# 4   Chou's decision tree

To complete the required experiment, run the below command in terminal

```
1  ./test_tree -m 3
```

The implementation is a little slow, typically requires several minutes to one hour to complete. The threshold for accepting candidate left and right nodes is still 0.005. Repeat the experiments four times, we obtain the below outputs.

```
1  Constructing agglomerative clustering tree...
2  Initialize Chou's decision tree...
3  Build development, and held-out data sets...
4  Chou's decision tree is growing...(A little slow...)
5  Process: .................
6  No frontier node...
7  Bit tree growing complete...
8  Calculate perplexity on testing data set...
9  perplexity is: 4.93075
```

```
1   Constructing agglomerative clustering tree...
2   Initialize Chou's decision tree...
3   Build development, and held-out data sets...
4   Chou's decision tree is growing...(A little slow...)
5   Process: .
6   No frontier node...
7   Bit tree growing complete...
8   Calculate perplexity on testing data set...
9   perplexity: 17.0487
10  perplexity is: 17.0487
```

```
1  Constructing agglomerative clustering tree...
2  Initialize Chou's decision tree...
3  Build development, and held-out data sets...
4  Chou's decision tree is growing...(A little slow, allow several minutes to one hour to complete...)
5  Process: .......
6  No frontier node...
7  Bit tree growing complete...
8  Calculate perplexity on testing data set...
9  perplexity is: 9.01359
```

```
1  Constructing agglomerative clustering tree...
2  Initialize Chou's decision tree...
3  Build development, and held-out data sets...
4  Chou's decision tree is growing...(A little slow, allow several minutes to one hour to complete...)
5  Process: ...........
6  No frontier node...
7  Bit tree growing complete...
8  Calculate perplexity on testing data set...
9  perplexity is: 15.5708
```

## 4.1 Tie breaking procedure

For the histories that neither $l_i \notin \mathcal{A}_i$ nor $l_i \notin \bar{\mathcal{A}}_i$, I simply randomly assign them into $\mathcal{A}_i$ or $\bar{\mathcal{A}}_i$ without no preference, that is 50 % into $\mathcal{A}_i$ or $\bar{\mathcal{A}}_i$. Such tie breaking procedure also works on test data set as well.

## 4.2 Compute Perplexity

I repeat the experiments 4 times. As shown in the output in Section 4, we can see the perplexity are 4.93075, 17.0487, 9.01359, 15.5708. We can see that Chou's Algorithm highly depends on its initialization. Sometimes, Chou's algorithm can outperform Bit-encoding tree, sometimes, it performs worse or competitively with Bit-encoding tree.