

Report

R10922132

吳峻銘

1. environment

software: 使用提供的 docker 跑，並且用 vs code 來連接 docker 做編輯

hardware : macbook pro, 16GB ram

2. implement

mapping.py

使用 python 基本的讀檔操作就好，但 encoding 的部分要記得改成 big5 hkscs

mydisambig.cpp

使用 srilm 內建的讀檔操作會省掉很多步驟接著就實作 viterbi 就好

```
int main(int argc, char *argv[])
{
    VocabMap map(voc_key, voc_value);
    Ngram lm(voc, 2);

    File mapfile(argv[2], "r");
    map.read(mapfile);
    mapfile.close();

    File lmfile(argv[3], "r");
    lm.read(lmfile);
    lmfile.close();

    File inputfile(argv[1], "r");
    FILE *fp = fopen(argv[4], "w");

    char *inputline = NULL;
    while(inputline = inputfile.getline())
    {
        VocabString word_of_line[MAXLINE] = {};
        char decode[MAXLINE][3] = {};
        int len = Vocab::parseWords(inputline, word_of_line, MAXLINE);

        viterbi(word_of_line, decode, lm, map, len);

        char output[MAXLINE_LEN] = {};
        for(int i=0; i<len; ++i)
        {
            strcat(output, decode[i]);
            strcat(output, " ");
        }
        fprintf(fp, "<s> %s</s>\n", output);
    }
    inputfile.close();
    fclose(fp);
}
```

Viterbi 的部分只有實作 bigram 參考課程網頁 FAQ 上提供的公式

Bigram Derivation of Viterbi

To derive bigram part of Viterbi algorithm, define:

$$\delta_t(q_i) = \max_{W_{1:t-1}} P(W_1, \dots, W_{t-1}, W_t = q_i)$$

where:

$$\delta_t(q_i) = \max_{W_{1:t-1}} P(q_i | W_{t-1}) P(W_{1:t-1}) = \max_{q_j} P(q_i | q_j) \delta_{t-1}(q_j)$$

For $t = 1$, initialize first timestep like:

$$\delta_1(q_i) = P(W_1 = q_i)$$

3. observe

```
1 <s> 忽視新聞開場迎喜李四端金素梅明搭檔雙主播 </s>
2 <s> 華社新聞將在明天年第一天推出約旦雙主播 </s>
3 <s> 由王牌主播李四端與剛出爐的新科立委高金素梅也同播報新聞 </s>
4 <s> 嘗試啼聲的播報內容 </s>
5 <s> 就是高金素梅自己明天上午應總統府之邀 </s>
6 <s> 參加元旦升旗典禮唱歌歌的新聞 </s>
7 <s> 華社今天召開董事會 </s>
8 <s> 總經理周蓉參與經營團隊展現上任一年多達成績單 </s>
9 <s> 在節目經營組織重鎮市都有新氣象 </s>
10 <s> 在會長董事郭力昕提案討論與台商策略聯盟的可行性 </s>
11 <s> 董事會決定以客源節流 </s>
12 <s> 對外積極合作 </s>
13 <s> 尋求電視台中型 </s>
14 <s> 忽視在明年首先提出誠意十字的新聞 </s>
15 <s> 讓全國觀眾觀眾耳目以新 </s>
16 <s> 李四端特別在新的年開始邀請高金素梅搭檔雙主播 </s>
17 <s> 其心準備當華社主播初體驗的她還將美國九十年代最後一天 </s>
18 <s> 獻給了李四端 </s>
19 <s> 兩人在試播搭檔後十分有信心 </s>
20 <s> 相信一定會在明天發揮完美表現 </s>
```

1.txt

基本上大致的 **decode** 都表現不錯，但有些詞組成句字會怪怪的，像是第一句應該要是「華視新聞」而不是「忽視新聞」，推測是「忽視」的 **bigram** 分數比「華視」高太多導致最後 **decode** 出來的結果是「忽視」，這種問題可以透過增 **corpus** 或是改變 **bigram** 背後的模型來解決。

4. challenge

用 vs code 當 IDE 時會因為編碼用 UTF-8 導致輸入中文字去 vocab 在編譯時 **decode** 不出來，後來改變 vs code 的編碼就解決了。

一開始實作 **viterbi** 演算法時陣列開太大導致 **segment fault**，所以可以了解實作 **viterbi** 不能用在太長的語音序列，不然 **DP** 策略會沒辦法執行。