

一、編譯環境：

gcc version 7.2.0 executed under Mac OSX terminal

二、主要檔案內容：

為了完成最終測得 testing data accuracy 的目的，我共寫了三個 cpp File (train.cpp, test.cpp, acc.cpp)，以下分別簡介各 cpp 檔的內容與作用。

- train.cpp

此檔案是用來接收 hmm initial mode 並且預期透過多次 BW algorithm training 來優化 $\lambda (\pi, A, B)$ 。一開始先計算各項求得 $\pi$ 、A、B 三指標所會需要的資料：forward algorithm  $\alpha[t][i]$ 、backward algorithm  $\beta[t][i]$ 、 $\gamma[t][i]$  以及 $\epsilon[t][i][j]$ ，再把他們 run 過每筆 sample 的數值累積起來，最後再透過講義公式求得新的 $\lambda (\pi, A, B)$ ，每次 iteration 再從這個新的 $\lambda (\pi, A, B)$ 開始計算。在 training 的時候可以發現 $\pi$ 、A、B 三個指標的資料會隨著 iteration 增加而逐漸收斂為一大約的數值。

- test.cpp

此檔案負責接收 testing\_data.txt，然後實作 viterbi algorithm 透過 $\delta[t][i]$ 與 $\delta[t+1][i]$ 的前後比較，在每筆 sample 中求出一條機率最大的路徑，此路徑就是最有可能的 state sequence (即為檔案裡的 maxRouteProb)，再去分別計算 model\_01.txt~model\_05.txt 的 maxRouteProb，當中最大的即為此筆 sample 所對應最有可能的 model 檔案，最後再把 testing 出的結果 output 到 result1.txt 中。

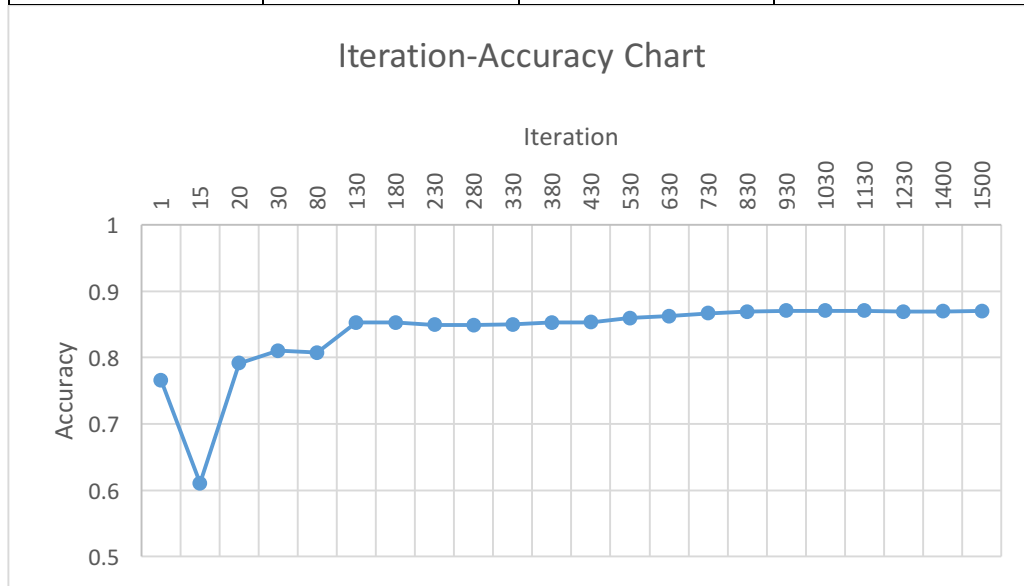
- acc.cpp

此檔案會接收剛剛 testing 出的結果 (result1.txt)，然後與 testing\_answer.txt 逐一比對，做正確率的計算，然後將結果 output 到 acc.txt 中。

三、數據分析與討論：

Iteration	Accuracy	Iteration	Accuracy
1	0.7660	430	0.8532
15	0.6104	530	0.8592
20	0.7912	630	0.8624

30	0.8104	730	0.8668
80	0.8072	830	0.8692
130	0.8528	930	0.8704
180	0.8528	1030	0.8704
230	0.8492	1130	0.8704
280	0.8488	1230	0.8692
330	0.8496	1400	0.8696
380	0.8524	1500	0.8700



從以上數據表格跟圖表可以發現，accuracy 一開始隨 iteration 增加而上升，但當 iteration 增至 15 時有明顯下降，之後繼續便不斷緩慢爬升，大概 iterate 到 900~1100 次的時候 accuracy 到達飽和峰值，若 iteration 再增加的話 accuracy 便開始在峰值 0.87 上下微小增減。

在實作 hmm training 後我發現，隨著 iteration 增加，準確率大致會不斷上升然後逐漸收斂於一高峰值，因為多次 BW-Algorithm 不斷地將 hmm 參數做優化，五個 model 的  $P(O|\lambda)$  也不斷地被 maximize，他們每筆 sample 的識別力增強，所以分析出來的準確率便越來越高。而 accuracy 一開始在 iteration=15 處的明顯下降，推測是因為在 iteration 總數較小時，進行每次 iteration 時 hmm 參數數值的變動皆較大，而 hmm 的參數剛好在 training 到那時候偏離收斂值太多，造成 accuracy 較低。