- 執行環境：
  ruby 1.9.3（Windows 7）

- 執行方式：
  在 Windows console 端執行指令 *ruby pa3.rb*（須安裝 ruby 1.9.3）。

- 執行時間（筆記型電腦，**i5 處理器、4GB 記憶體**）：

|  | 第一次 | 第二次 | 第三次 |
|---|---|---|---|
| **Feature Selection** | 8'30" | 11'20" | 10'00" |
| **Training & Testing** | 5'30" | 4'30" | 4'22" |
| 總時間 | 14'00" | 15'50" | 14'22" |

  平均執行時間：14'44"

- 程式執行步驟：

  **1. "training.txt"**

  首先讀取老師提供的"training.txt"檔案，把每個 class 和檔案的關係存成
  {"1" => [11, 19, 29,...], "2" => [1, 2, 3,...], ...} 的 hash 形式，方便之後讀取。

```ruby
 7    # Read "training.txt" into trainingFiles{}
 8
 9    array = open('training.txt', 'rb').readlines
10    trainingFiles = Hash.new
11    array.each do |arr|
12        a = arr.split
13        trainingFiles[a[0]] = a.drop(1)
14    end
```

  **2. Extract vocabulary**

  把 training document 的內容讀取出來，並且用之前寫的 extract 函式（寫
  在"extraction.rb"檔案中）做關鍵字的萃取，刪除多餘的 stopword 以及做
  stemming，把初步的關鍵詞彙存在 *vocab* 這個變數中。

```ruby
16    # Extract vocabulary from training document set (require extraction.rb)
17
18    print "Extracting vocabulary from training documents..."
19    vocab = Array.new
20    trainingFileTerms = Hash.new
21    trainingFiles.each { |key, docNames|
22        docNames.each do |dn|
23            f = open( INPUT_DIR + '/' + dn + '.txt', 'rb').read
24            extracted_array = extract(f)
25            trainingFileTerms[dn] = extracted_array
26            vocab.concat(extracted_array)
27        end
28    }
29    vocab = vocab.uniq
```

### 3. Feature selection

呼叫 selectFeatures 函式（寫在"featureSelection.rb"檔案中）。這邊選擇使用投影片中的「$\chi^2$ feature selection」方法，針對每一個 class，對前面 *vocab* 變數中的每一個 term 計算出現於該 class 的次數、出現於非該 class 中的次數，再計算 $\chi^2$ 統計值，並且在每個 class 中選出前 38 或 39 個 $\chi^2$ 統計值最高的 term（比較重要、比較有判別性的 term），放入最後真正會使用的關鍵字集合中。最後會湊滿 500 個關鍵字，存在 *vocabFS* 這個變數中。

```
31   # Feature selection: chi-square (require featureSelection.rb)
32
33   vocabFS = Array.new
34   for i in 1..CLASS_NUM
35       res = selectFeatures(trainingFiles, trainingFileTerms, vocab, CLASS_NUM, i.to_s)
36       vocabFS.concat(res)
37   end
```

```
1   def selectFeatures(trainingDocs, trainingDocTerms, vocab, class_num, c)
2
3   print "\nFeature Selection for class " + c + "..."
4   termUtility = Hash.new
5
6   #for each term in vocabulary
7   vocab.each do |t|
8       form = Hash.new
9
10      x = 0.0
11      countDocs = 0.0
12      #for each document in class c, check if t exists
13      trainingDocs[c].each do |dn|
14          if trainingDocTerms[dn].include?(t)
15              x += 1
16          end
17          countDocs += 1
18      end
19      form["on&present"] = x
20      form["on&absent"] = countDocs - x
21      on = countDocs
22
23      y = 0.0
```

### 4. Training

確定關鍵字集合之後就進入 training 階段。針對每一個 class，計算每一個 term 在該 class 中出現的條件機率，也就是說會有 C * M = 13 * 500 = 6500 個條件機率值，存在 condprob 這個變數中，以{"apple" => [0.01, 0.02,...], "tree" => [0.1, 0.032,...], ...} 這樣的 hash 形式儲存。

```
53    for i in 1..CLASS_NUM
54        print "\nOn class " + i.to_s + "..."
55        prior[i] = nClass[i] / nTotal
56
57        #concatenate all docs from class c
58        text = Array.new
59        trainingFiles[i.to_s].each do |dn|
60            text.concat(trainingFileTerms[dn])
61        end
62
63        termCount = Hash.new
64        termCountTotal = 0.0
65
66        vocabFS.each do |t|
67            termCount[t] = text.count(t).to_f
68            termCountTotal = termCountTotal + termCount[t] + 1
69        end
70
71        vocabFS.each do |t|
72            if condprob.has_key?(t) == false
73                a = Array.new
74                condprob[t] = a
75            end
76            condprob[t][i] = (termCount[t] + 1) / termCountTotal
77        end
78    end
```

### 5. Testing

　　Training 結束後進入 testing phase。將 training documents 以外的其他文件一個一個打開，利用前面計算出來的條件機率，對每份文件計算其在 13 個 class 的分數，選出分數最高的 class，作為這份文件所屬的 class。

```
91    #for each document
92    Dir.foreach( INPUT_DIR + '/') do |doc|
93        next if doc == '.' or doc == '..'
94        #if the document is not a training document
95        if trainingDocs.include?(doc.chomp(".txt")) == false
96            print "\nOn document " + doc + "..."
97
98            f = open( INPUT_DIR + '/' + doc, 'rb').read
99            terms = extract(f)
100           score = Hash.new
101
102           #compute score for each class
103           for i in 1..CLASS_NUM
104               score[i.to_s] = Math.log(prior[i])
105               terms.each do |t|
106                   if vocabFS.include?(t)
107                       score[i.to_s] += Math.log(condprob[t][i])
108                   end
109               end
110           end
111
112           scoreSorted = score.sort_by{|key, value| value}.reverse
113           classifyResult[doc.chomp(".txt")] = scoreSorted[0][0]
114       end
115   end
116
117   output = classifyResult.sort_by{|key, value| key.to_i}
```