

# PTT推啦分析

## 機器學習導論期末專題

105061225 周宜星 105061254 林士平

(註: 此code需要在code資料夾下才能執行)

### 一、專題目標

針對不同的版，分析文章推啦數和特定關鍵字之間的關係，最終目標是針對不同版各訓練出一個模型，將文章丟入便可以推測它的推啦數區間。

### 二、專題流程



### 三、專題內容 – 網路爬蟲介紹

1. 利用request得到網頁資料，再利用beautiful soup從html格式擷取我們需要的資訊。
2. 我們找了ptt較熱門的11個版，而且沒有封鎖推啦功能。
3. 由於ptt上有些版是18+的，所以爬蟲程式中需要額外做自動過濾18+的確認動作。
4. 我們每版的data做200頁，非18+的需要40分鐘、18+的版需要90分鐘。
5. 得到了11個database，文章的推文數-噓文數為label，文章內文去除引用及推噓內文內容為feature。

下面為database基本資料與對應的code:

```
In [11]: from IPython.display import Image
Image(filename='image/專題資料簡介.jpg', width=1000)
```

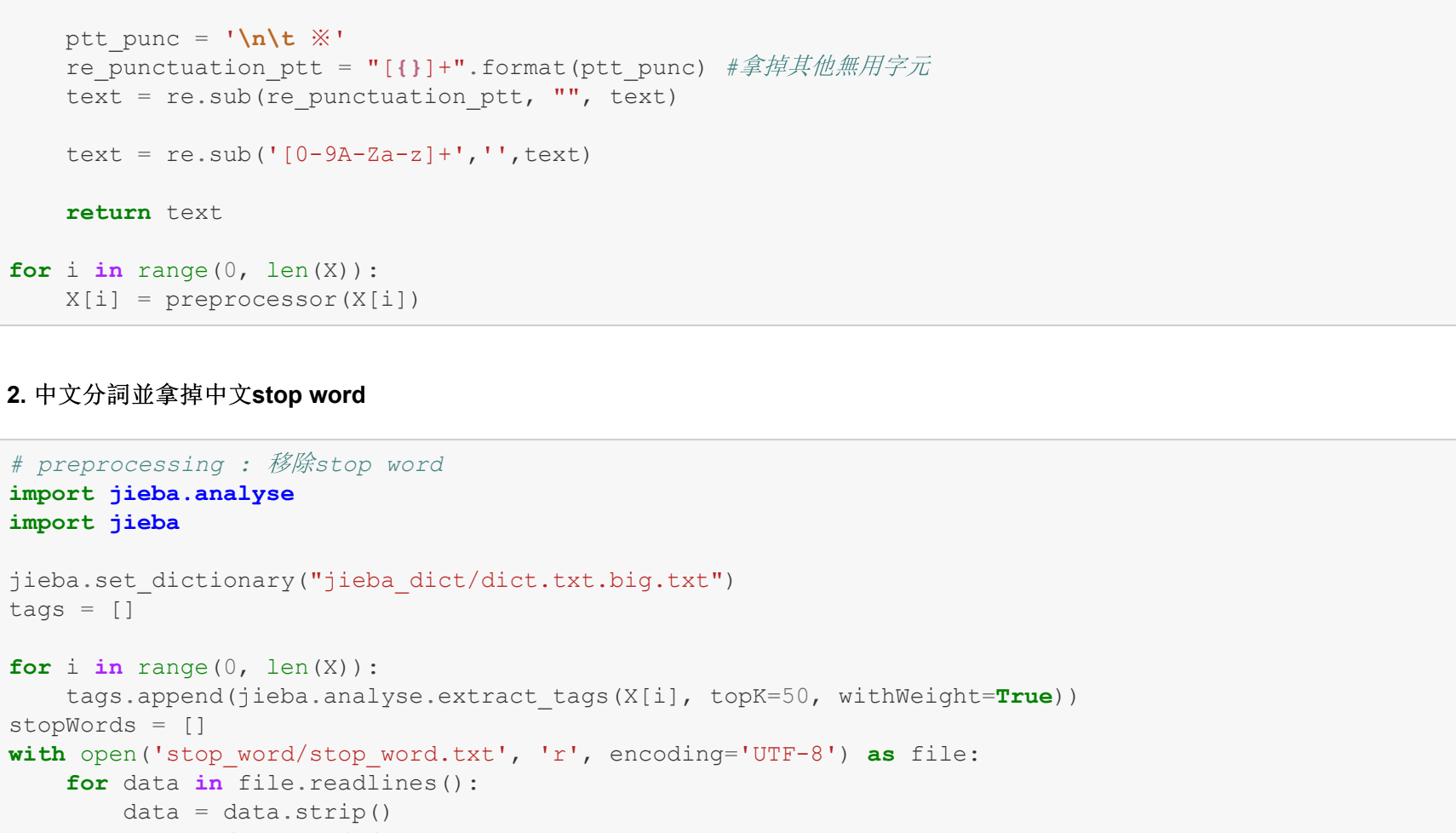
Out [11]:

	版	文章量(篇)	對應的.csv檔	對應的程式檔
	媽寶版(BabyMother board)	3958	babymother_max.csv	train_babymother_new.ipynb
	八卦版(Gossiping board)	3541	gossiping_max.csv	train_gossiping_new.ipynb
	政黑版(HatePolitics board)	3774	hatepolitics_max.csv	train_hatepolitics_new.ipynb
	韓星版(KoreaStar board)	3983	koreastar_max.csv	train_koreastar_new.ipynb
	省錢版(lifeismoney board)	3969	lifeismoney_max.csv	train_lifeismoney_new.ipynb
	婚姻版(marriage board)	3905	marriage_max.csv	train_marriage_new.ipynb
	行動通訊/手機(mobilecomm board)	3949	mobilecomm_max.csv	train_mobilecomm_new.ipynb
	電影版(movie board)	3970	movie_max.csv	train_movie_new.ipynb
	sex版(sex board)	3945	sex_max.csv	train_sex_new.ipynb
	股市版(stock board)	3962	stock_max.csv	train_stock_new.ipynb
	女版(womantalk board)	3967	womantalk_max.csv	train_womantalk_new.ipynb

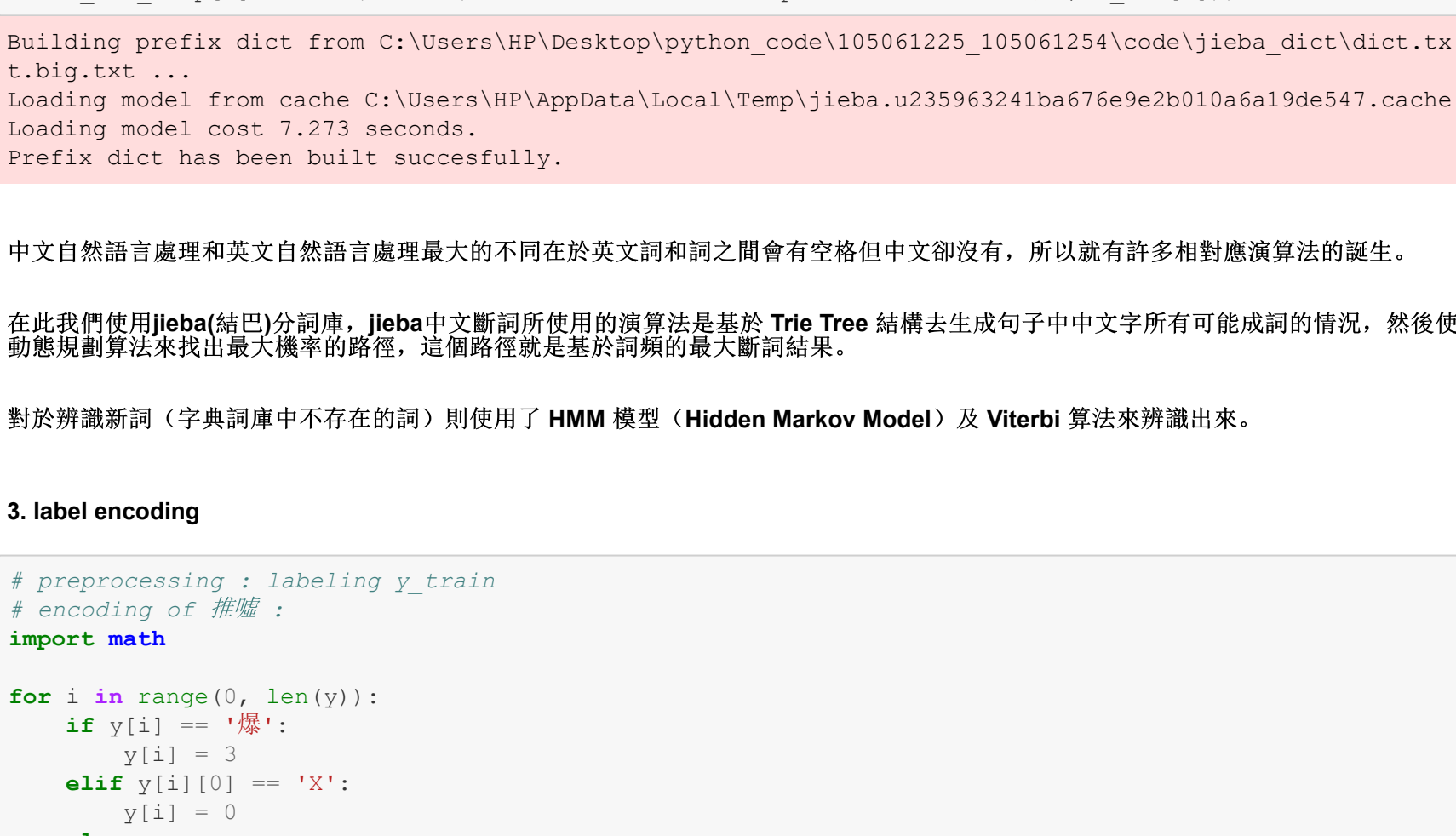
(註: 網路爬蟲程式檔名codeWeb Crawler.ipynb; database檔案均在code資料夾內)

### 四、專題內容 – 預處理、萃取feature vector與label encoding

(註: 此版將針對train\_gossiping\_new.ipynb為例子，其他的程式檔案見code資料夾)

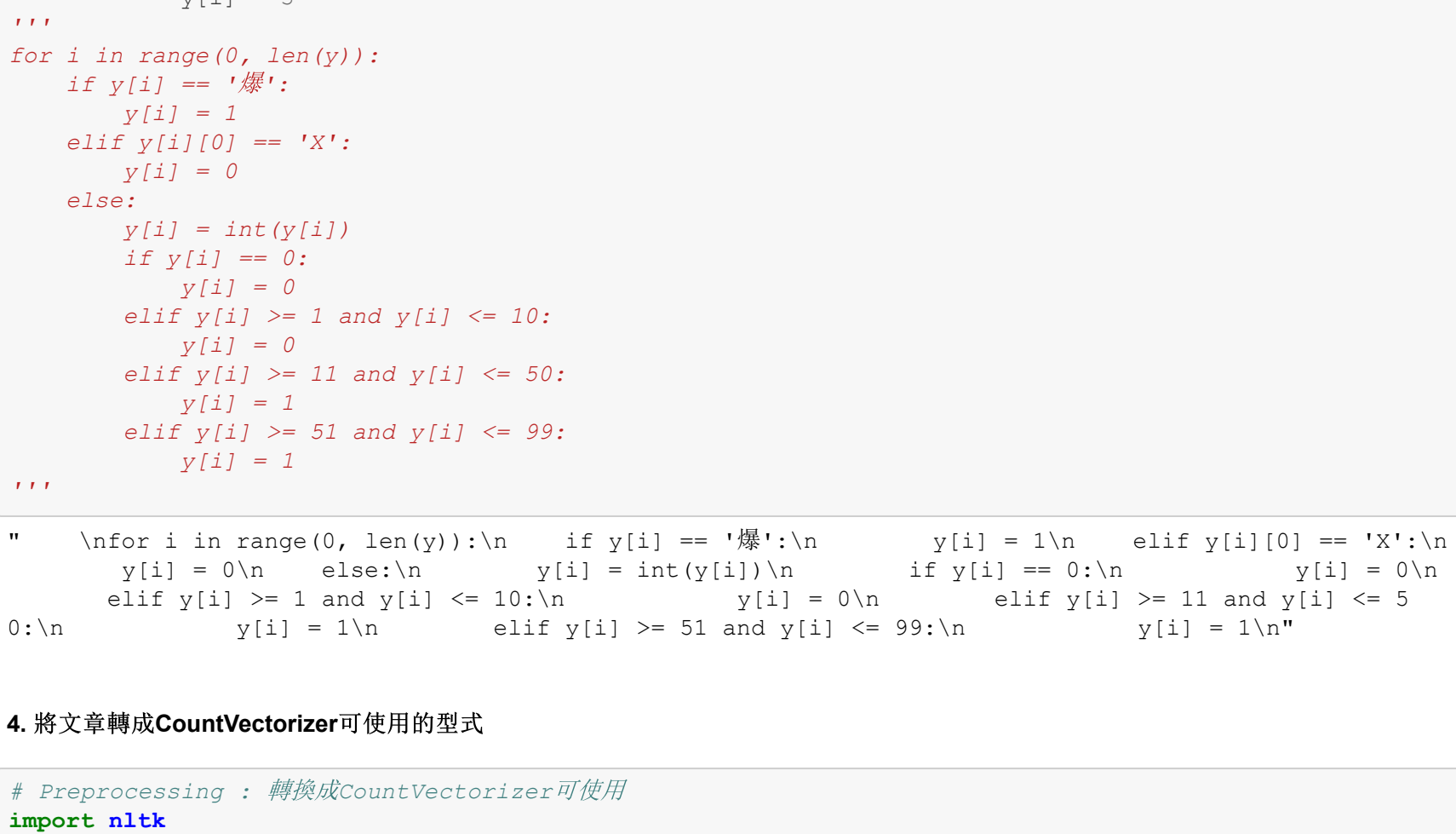


label encoding的方式如下圖:



接下來為預處理、萃取feature vector與label encoding的程式說明:

#### 1. 拿掉中英文標點以及其他無用字元



#### 2. 中文分詞並拿掉中文stop word

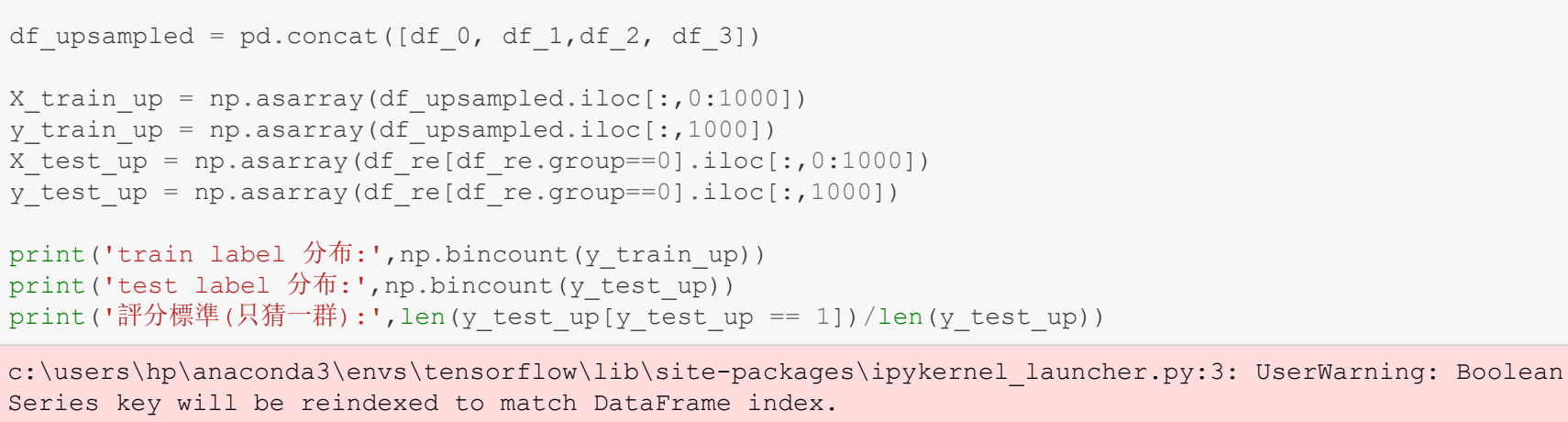


中文自然語言處理和英文自然語言處理最大的不同在於英文詞和詞之間有空格但中文卻沒有，所以就有許多相對應演算法的誕生。

在我們使用jieba進行分詞前，jieba中文斷詞所使用的演算法是基於 Trie Tree 結構去生成詞字中文字字串有可能成詞的情況，然後使用動態規畫算法找出最大機率的路徑，這個路徑就是從根到葉的最大詞結果。

對於辨識新詞(字典詞庫中不存在的詞)則使用了 HMM 模型 (Hidden Markov Model) 及 Viterbi 算法來辨識出來。

#### 3. label encoding



經過jieba分詞後的文章必須經由以上的轉換才能使用CountVectorizer和TfidfTransformer做後續的處理

#### 5. train-test split(70%-30%)

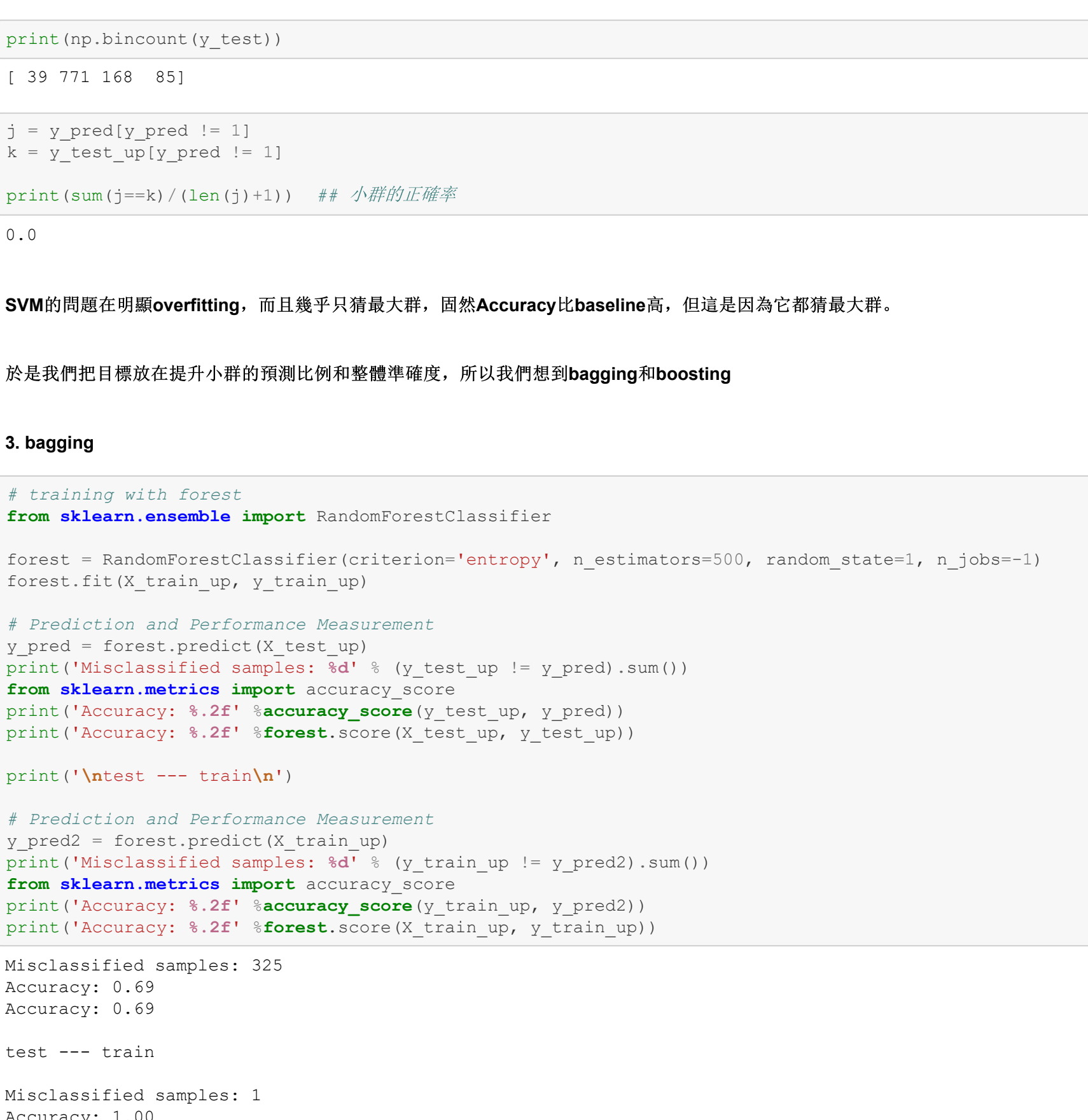


每一個label的sample數量不平衡，所以要做training set做upsample.

#### 6. 使用bag of word和tfidf得到feature vector

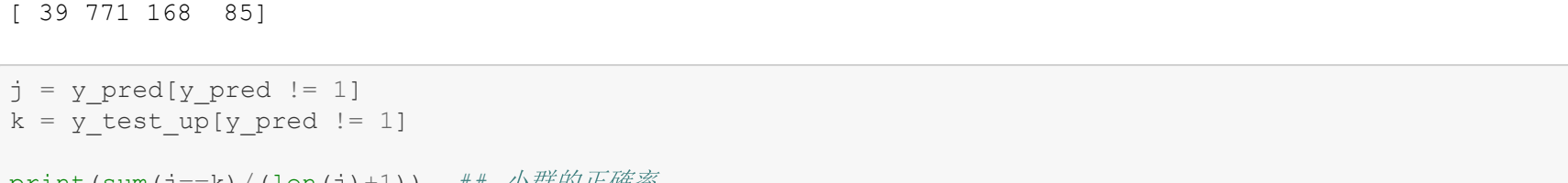


#### 7. resample使training set中每個label的sample數量相同

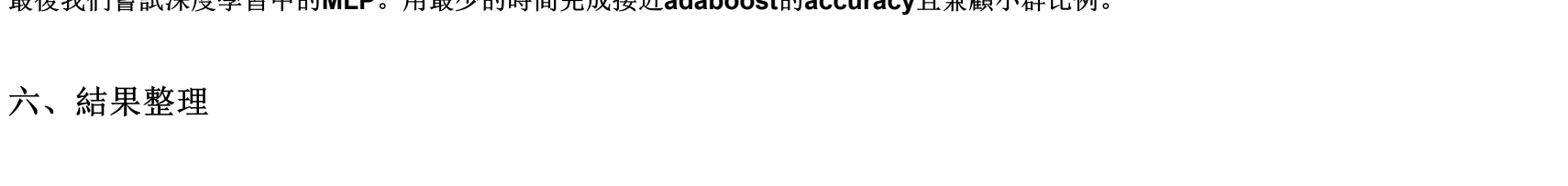
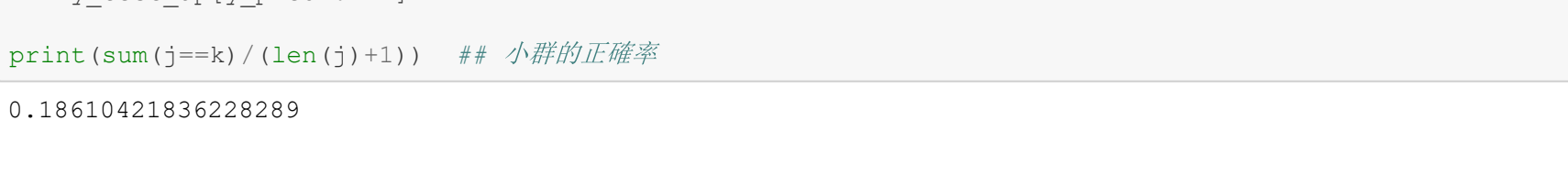


### 五、專題內容 – 使用一系列機器學習演算法

#### 1. 我們以logistic regression作為baseline



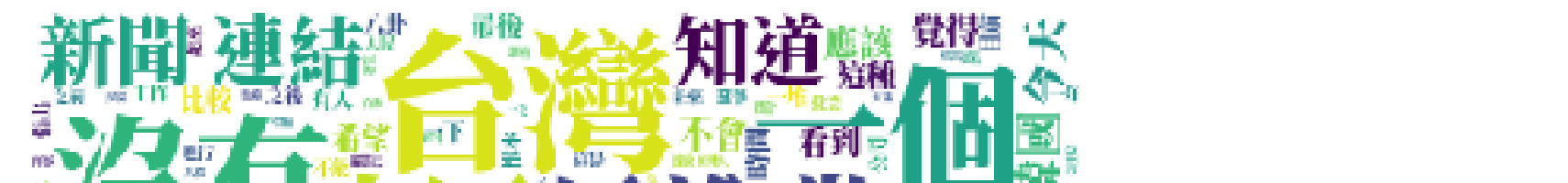
#### 2. 接著試試SVM



SVM的問題在明顯overfitting，而且幾乎只猜最大群，固然Accuracy比baseline高，但是是因為它都猜最大群。

於是我們把目標放在提升小群的預測比例和整體準確度，所以我們想到bagging或boosting

#### 3. bagging



bagging使用Randomforest，可以觀察到相對SVM大幅解決只猜最大群的問題，且accuracy比logistic regression高。

#### 4. boosting



boosting使用adaboost，與前者比較稍微犧牲準確度，但是針對小群的預測也更多，這是由於boosting的設計本來就是針對特殊情況，但是有些缺點，例如某些版收敛速度極慢，至少要10個小時可能才符合需求。

#### 5. Multilayer perceptron(MLP)



最後我們嘗試深度學習中的MLP，用最少的時間完成接近adaboost的accuracy且明顯小群比例。

### 六、結果整理

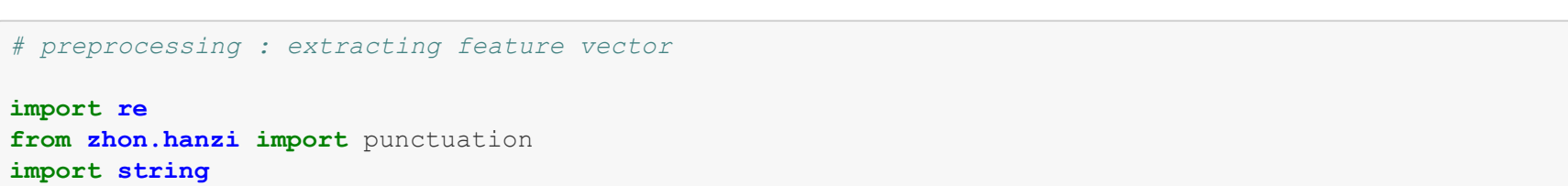
1. 首先是各個版的字雲分析結果，可以看到不同版的之間關鍵字確實有明顯不同，不過也可以觀察到某些版之間有重疊的關鍵字



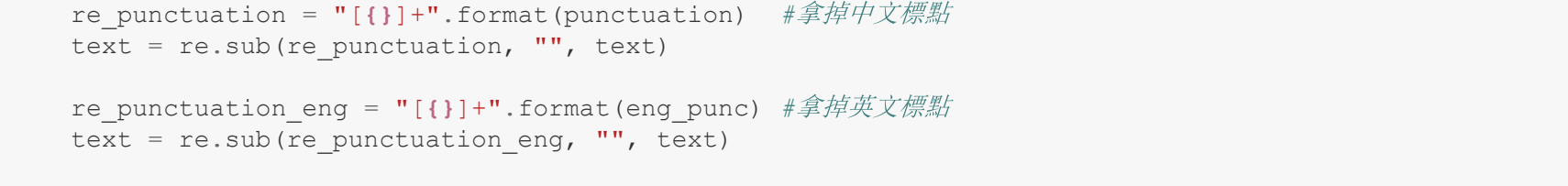
(1) 媽寶版(BabyMother board)



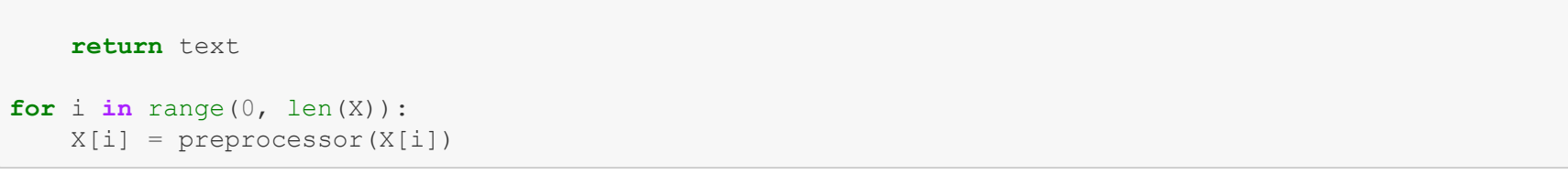
(2) 八卦版(Gossiping board)



(3) 政黑版(HatePolitics board)



(4) 韓星版(KoreaStar board)



(5) 省錢版(Lifeismoney board)



In [12]: Image(filename='image/lifeismoney\_wordcloud.png', width=600)



(6) 婚姻版(marriage board)

In [13]: Image(filename='image/marriage\_wordcloud.png', width=600)



(7) 行動通訊手機(mobilecomm board)

In [14]: Image(filename='image/mobilecomm\_wordcloud.png', width=600)



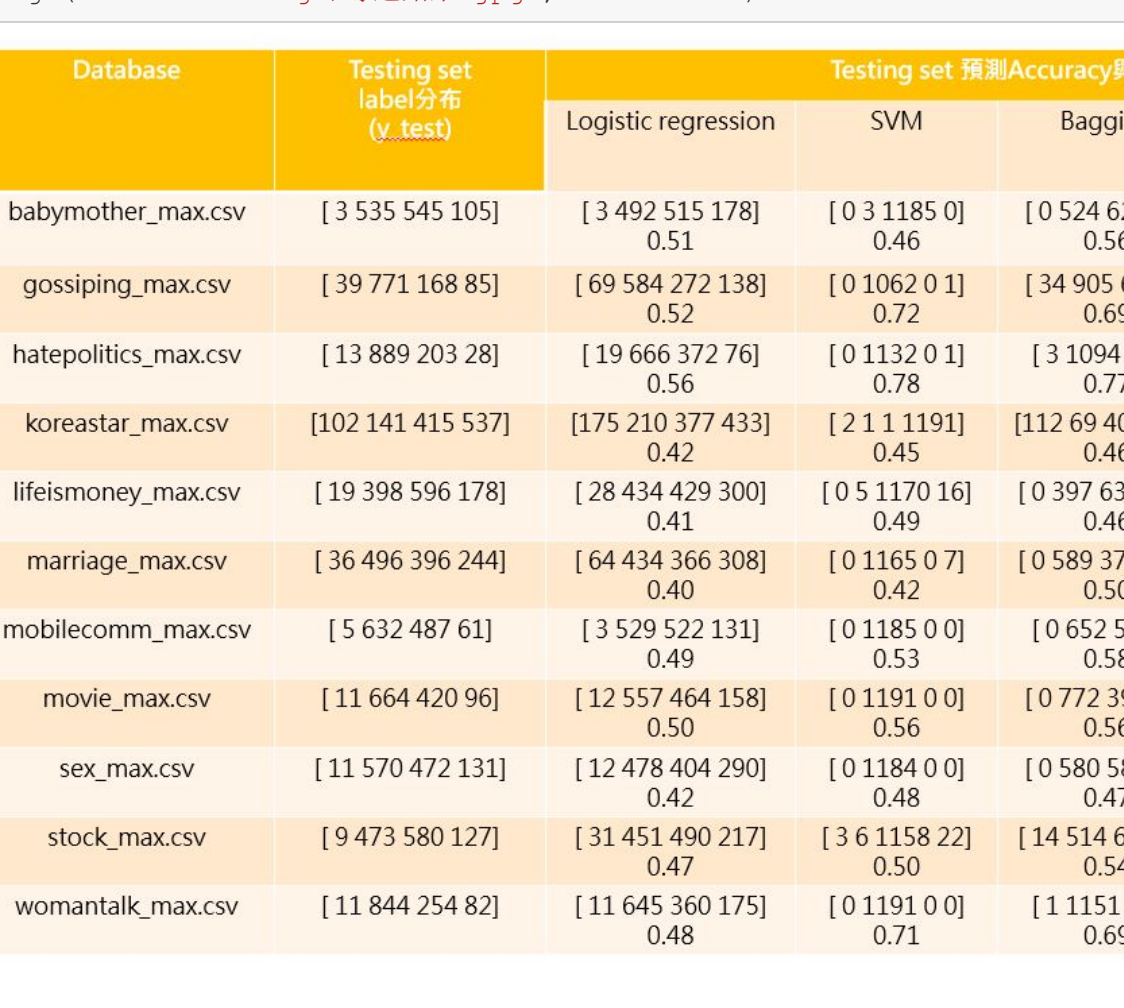
(8) 電影版(movie board)

In [15]: Image(filename='image/movie\_wordcloud.png', width=600)



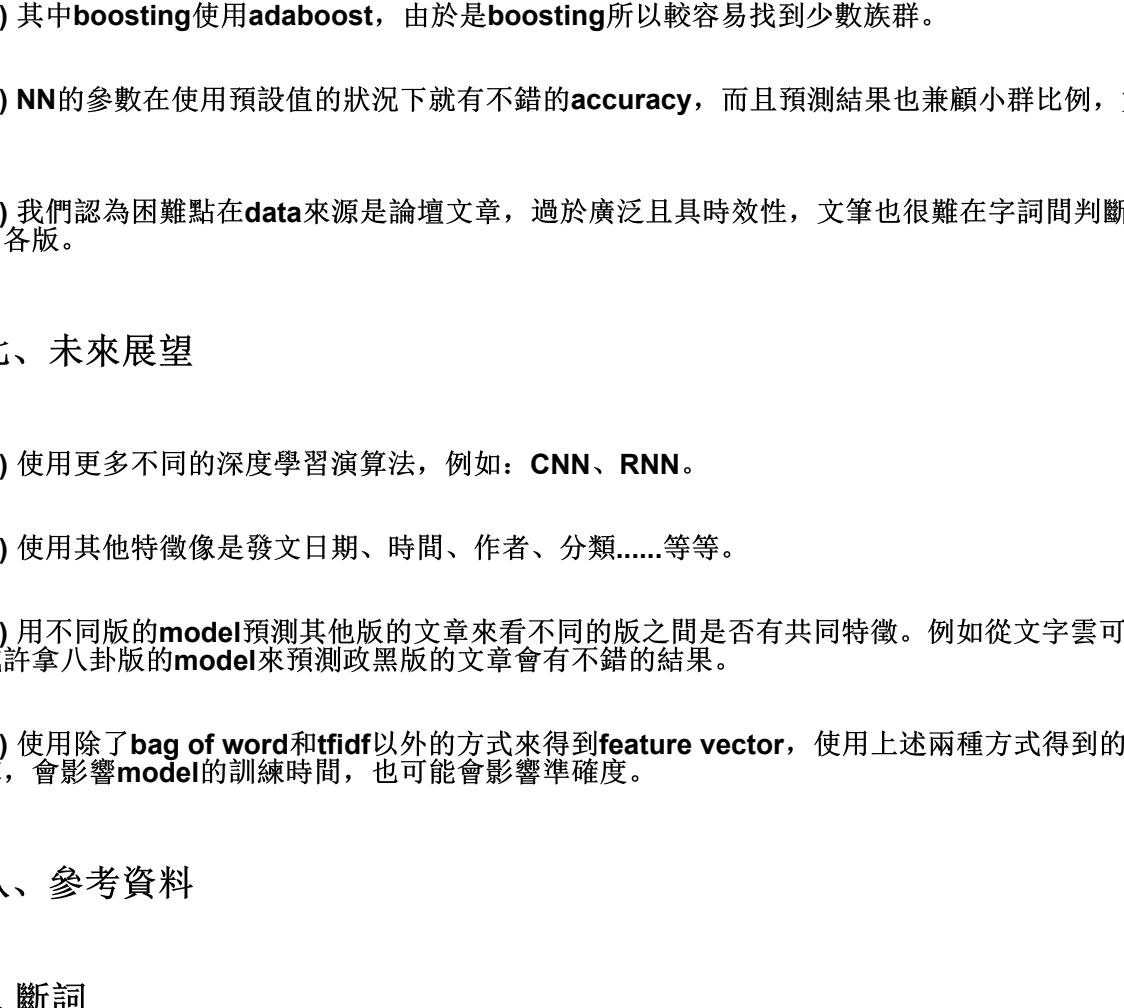
(9) 西版(sex board)

In [16]: Image(filename='image/sex\_wordcloud.png', width=600)



(10) 股版(stock board)

In [17]: Image(filename='image/stock\_wordcloud.png', width=600)



(11) 女版(womantalk board)

In [18]: Image(filename='image/womantalk\_wordcloud.png', width=600)



## 2. 各個版(training model)結果整理

In [3]: Image(filename='image/專題結果.jpg', width=1000)

Database	Testing set label分布 (y_test)	Testing set 預測Accuracy與預測label分布(y_pred)					
		Logistic regression	SVM	Bagging	Boosting	MLP	
babymother_max.csv	[3 535 545 105]	[3 492 515 178] 0.51	[0 3 1185 0] 0.46	[0 524 621 43] 0.56	[0 708 431 49] 0.53	[0 528 527 133] 0.49	
gossiping_max.csv	[39 771 168 85]	[69 584 272 138] 0.52	[0 1062 0 1] 0.72	[34 905 63 61] 0.69	[25 788 85 165] 0.61	[42 661 243 117] 0.55	
hatepolitics_max.csv	[13 889 203 28]	[19 666 372 76] 0.56	[0 1132 0 1] 0.78	[3 1094 32 4] 0.77	[9 891 221 12] 0.67	[10 761 319 43] 0.61	
koreastar_max.csv	[102 141 415 537]	[175 210 377 433] 0.42	[2 1 1 1391] 0.45	[112 69 403 611] 0.46	[145 81 284 685] 0.45	[138 169 438 480] 0.41	
lifesimoney_max.csv	[19 398 596 178]	[28 434 429 300] 0.41	[0 5 1170 16] 0.40	[0 397 630 164] 0.46	[13 53 354 771] 0.51	[18 428 486 259] 0.40	
marriage_max.csv	[36 496 396 244]	[64 434 366 308] 0.40	[0 1165 0 7] 0.42	[0 589 378 205] 0.50	[22 597 381 202] 0.47	[25 442 392 313] 0.40	
mobilecomm_max.csv	[5 632 487 61]	[3 529 522 131] 0.49	[0 1185 0 0] 0.53	[0 652 530 3] 0.58	[2 545 606 32] 0.51	[0 550 570 65] 0.51	
movie_max.csv	[11 664 420 96]	[12 557 464 158] 0.50	[0 1191 0 0] 0.56	[0 772 396 23] 0.56	[2 339 710 140] 0.41	[5 588 460 138] 0.50	
sex_max.csv	[11 570 472 131]	[12 478 404 290] 0.42	[0 1184 0 0] 0.48	[0 580 585 19] 0.47	[1 275 606 302] 0.38	[3 507 449 225] 0.43	
stock_max.csv	[9 473 580 127]	[31 451 490 217] 0.47	[3 6 1158 22] 0.50	[14 514 603 58] 0.54	[3 675 460 51] 0.50	[13 473 521 182] 0.50	
womantalk_max.csv	[11 844 254 82]	[11 645 360 175] 0.48	[0 1191 0 0] 0.71	[1 1151 38 1] 0.69	[7 641 216 327] 0.45	[6 701 368 116] 0.51	

## 3. 結論

- (1) 可以看得出来SVM有明显overfitting的問題。
- (2) bagging和boosting可以解決svm只預選小群的問題，而且accuracy普遍也比baseline好。
- (3) 其中bagging使用random forest，效率高，但是有時會有overfitting的問題。
- (4) 其中boosting使用adaboost，由於是boosting所以較容易找到少數族群。
- (5) NN的參數在使用預設值的狀況下就有不錯的accuracy，而且預測結果也兼顧小群比例，如果未來再調整參數或許會有更好的結果。
- (6) 我們認為困難點在data來源是論壇文章，過於廣泛且具時效性，文章也很難在字詞間判斷，而且很難定義出分布均勻且有意義的label應用各版。

## 七、未來展望

- (1) 使用更多不同的深度學習演算法，例如：CNN、RNN。
- (2) 使用其他特徵像是發文日期、時間、作者、分類.....等等。
- (3) 用不同版的model預測其他版的文章來看不同的版之間是否有共同特徵。例如從文字雲可以看的出來八卦版和政黑版有類似的關鍵字，或許拿八卦版的model來預測政黑版的文章會有不錯的結果。
- (4) 使用除了bag of word/tfidf以外的方式來得到feature vector，使用上述兩種方式得到的feature vector最大的缺點是特徵矩陣為稀疏矩陣，會影響model的訓練時間，也可能會影響準確度。

## 八、參考資料

### 1. 斷詞

[NLTK 初學指南\(一\)- 簡單易上手的自然語言工具第一課-徐家宜](#)

[如何使用jieba 結巴中文分詞程式](#)

[以jieba 與gensim 探索文本主題- 五月天人生無限公司數說分析\(1\)](#)

### 2. jieba github

[original jieba github](#)

[Taiwan jieba github](#)

### 3. csv\_read for Chinese

[pandas\\_read\\_csv\(中文名義 中文路徑\)](#)

### 4. 拿掉標點符號

[jieba分詞，並去除了標點](#)

[zhon包](#)

[查找標點方法](#)

### 5. 文字雲

[python \(wordcloud\) 實例中文詞雲](#)

### 6. NLTK

[NLTK 初學指南\(二\)- 簡單易上手的自然語言工具第二課-徐家宜](#)

### 7. 移除中文停用自

[10-2 中文斷詞-移除停用詞](#)

### 8. 流程

[NLP入門- 文本預處理Pre-processing](#)

### 9. 關於NLP

[進入NLP 世界的最佳捷徑，喜給所有人的自然語言處理與深度學習入門指南](#)

[Practical Text Classification With Python and Keras](#)

### 10. 深度學習

[sklearn\\_neural\\_network.MLPClassifier](#)

[Convolutional Neural Network \(CNN\)](#)

[Deep Learning: An MIT Press book / Ian Goodfellow and Yoshua Bengio and Aaron Courville](#)