



基於 Spark 環境上的電影評論語意分析 —
以 meta-critic 網站之評論為例

指導教授:劉敦仁老師

組 員:0453406 郭晉晏

0453433 陳楷仁

0453435 簡大鈞

目錄

一、動機	1
二、資料集與分析工具	1
(一) 資料集	1
(二) 分析工具	2
三、實作流程	3
(一) 資料擷取	4
(二) 分類方法實作	5
(三) 實驗結果	8
四、分析結果與結論	9

一、動機

在現在時代網路的發展越來越蓬勃，我們每天所製造以及所獲得的資訊量與日俱增，日益龐大的資訊量漸漸無法為人工所負荷，龐大的資料無法被人們有效率的紀錄並吸收，因此許多人開始嘗試以電腦去處理複雜的問題，而其中一個複雜的問題就是文字探勘的語意分析，電影評論的語意分析就是其中一個相當泛用的例子，現在有許多已經發展的相當有規模的電影評論網站像是爛番茄，IMDB，以及 meta-critic，這些網站上都有相當程度的網友在上面留下自己的影評，若想要歸納這些影片使用傳統的方法運用人力的方式去標籤每個評論，將會曠日費時，因此若能用機器學習來運用演算法產生出近似一般人加入感情所產生的結果，將能有效率的幫助我們歸納這些影評。這學期在資料探勘研究與實務中，在實作作業中我們有學到決策樹的實作，課程中也學習到貝式分類的概念，因此在這次的資料探勘期末專題我們選擇運用 meta-critic 的影評資料來實作電影評論的語意分析，透過 Spark 裡的 MLlib 中的決策樹和 Naïve Bayes 的函式庫來幫助我們建立一個能處理龐大影評資料的機器學習模型，並用此模型能夠正確分類出正反評論，作為是否觀看電影時之參考依據。

二、資料集與分析工具

(一) 資料集

在資料探勘的過程中通常會建立預測模型，之後以此模型針對目標進行預測。因此資料集在這方面的應用是個不可或缺的角色。

meta-critics 是一個外國相當熱門之影評網站，上面的影評基本分為兩種：user-critics、meta-critics。網站中 user-critics 是指一般使用者的評論，而 meta-critic 則是專業的影評或影評團隊留下的電影評論，而兩者評論也皆會有對該電影之評分。在此我們使用 python 撰寫爬蟲程式，協助我們抓取一般使用者的 user-critics 以及相對應評分。

(二) 分析工具

隨著大數據概念的興起，相關應用在硬體方面有著嚴苛的需求，因此近年來，分散式系統的相關應用也不斷地被發展來協助計算任務的完成，這裡我們使用了以下工具及語言：

Hadoop

是一款支援資料密集型分布式應用的開源軟體框架，隨著巨量資料應用帶來的龐大硬體資源需求，其也成為相當熱門之工具。Hadoop 提供了分布式檔案系統，用以存儲所有計算節點的資料，這為整個集群帶來了非常高的帶寬。MapReduce 和分布式檔案系統的設計，使得整個框架能夠自動處理節點故障。

Spark

是一個開源叢集運算框架，最初是由加州大學柏克萊分校 AMPLab 所開發。相對於 Hadoop 的 MapReduce 會在執行完工作後將中介資料存放到磁碟中，Spark 運用了記憶體內運算速度較快的概念，在資料尚未寫入硬碟時即在記憶體內分析運算。因此 Spark 在運算速度的表現上能做到比 Hadoop MapReduce 快上數倍，非常適合用於機器學習演算法。

Python

是一種物件導向、直譯式的電腦程式語言。它包含了一組功能完備的標準庫，能夠輕鬆完成很多常見的任務，另外也提供了豐富的擴充套件，可以幫助使用者透過 `pip install` 指令輕鬆安裝所需之功能，這邊我們用了以下的函式庫協助我們在實驗的進行：

Python 函式庫	說明
NumPy	一個經常數學運算用途之函式庫，支援較高維度的陣列及矩陣運算。
Pandas	一個可用於資料格式處理之函式庫，在資料分析的前處理步驟中經常被使用到。

NLTK(Natural Language Toolkit)	一個經常被用於自然語言分析的函式庫，這裡我們取它在情感分析上的處理加以應用。
MLlib	一個建立於 Spark 之上的有關於機器學習函式庫，可被 JAVA、Scala、Python、R 等語言實作，內有許多分類、分群及回歸等機器學習相關之方法。
Splearn	一個 spark 上的擴充套件，可以讓 scikit-learn 的一些方法在 pySpark 上實現

我們將使用 Hadoop 環境和 spark 作為分析工具軟體，而使用的工具主要為 python 的 library 包括 NumPy、Pandas、NLTK 以及 Spark 上的 library MLlib，其中 NumPy 和 Pandas 主要是作為資料處理用，NLTK 有機器學習演算法功能也可以做資料處理，而最主要的部份我們將會使用到 pySpark 上的分散是機器學習 library MLlib 來運作 decision tree。

三、實作流程

此次分析主要利用兩種方式作分析

- (1) 利用 scikit-learn 將資料作 TF-IDF，以建立出 frequency matrix，然後利用 Naïve Bayes 將資料作分析。
- (2) 利用 nltk 中的”sentiment”套件先將資料作語意分析會分析出’compound’（代表總和情緒的高低，介於-1~1，負值表示為負情緒，正值表示為正情緒）、’neg’（負面情緒字詞比例）、’neu’（中間情緒字詞比例）、’pos’（正面情緒字詞比例）。將分析出來的 4 項數據使用 decision tree 作分析。

以上分析皆將資料的 75%設為 training data，剩下 25%為 test data；(1)的部分均在 Spark 上進行，(2)僅有 decision 的部分在 spark 上進行。

(一) 資料擷取

1. 至 metacritic(<http://www.metacritic.com/>) 抓取資料，抓取的部分為評論以及分數，總共抓取資料為 858 筆。

ROGUE ONE: A STAR WARS STORY

WALT DISNEY STUDIOS MOTION PICTURES | RELEASE DATE: DECEMBER 16, 2016

7.7

USER SCORE

Generally favorable reviews
based on 1728 Ratings

USER RATING DISTRIBUTION

POSITIVE: 1,352
MIXED: 213
NEGATIVE: 163

USER REVIEWS

USER SCORE BY DATE MOST HELPFUL

分數

10

ColourSurround

Jan 4, 2017

I came in with meager expectations and was genuinely blown away with how the film handled the story, paced the action, and played out as a whole. The acting was superb, the effects made you feel like you were watching a Star Wars film - and... Expand

評論

0 of 1 users found this helpful 0 1

All this user's reviews

10

Commonsense

Dec 19, 2016

Critics are pretty wrong on this one. Characters could have been developed a little more, but the tension, cinematography and, overall story didn't disappoint. No one in my theater seemed unhappy and many even clapped and cheered and this was... Expand

2 of 8 users found this helpful 2 6

All this user's reviews

圖 (一) metacritic 上我們資料擷取之目標，有評論內容和對電影之分數

2. 利用 python 中的 "BeautifulSoup" 套件 parser html 中的資料，在將 parse 出的資料轉成 csv 格式並輸出，其輸出檔案命名為 review.csv，如圖(二)。

	A	B	C	D	E
1	SCORE	REVIEW			
2	10	I came in with meager expectations and was genuinely blown away with how the film handled the story, paced the action, and played out as a whole. The acting was superb.			
3	10	Critics are pretty wrong on this one. Characters could have been developed a little more, but the tension, cinematography and, overall story didn't disappoint. No one in my			
4	10	was 12 when the last decent Star Wars movie came out.I've been burned ever since: Special Editions, Prequels, Force Awakens (sorry, it sucks).So the bar was pretty low g			
5	10	Rogue One was spectacular and everything I wanted it to be. I consider this movie an apology for the abysmal Force Awakens and for this I can only say one thing: Apology accepted!			
	評分	網友影評			

圖 (二) 利用爬蟲程式所儲存之格式 CSV 格式

(二) 分類方法實作

Mllib 是 spark 官方所提供的有關於機器學習的套件，其中有相當多的機器學習方法可以用於分析像是分類的分析方法如 Naïve Bayes、Decision Tree、SVM，其他也有分群、回歸方法可供利用，另外亦提供範例程式(scala、JAVA、python) 等供大家參考，這裡的方法我們選擇實作學期初所學到的決策樹分類法，我們除了使用 Mllib 外也使用 splearn 來作貝式分類器(Naïve Bayes)。

資料前處理

將資料作處理，這次主要的目的為判斷語意與分數高低之關係，因此須先將分數區分出為高或低，故將 0~5 分設為” 0” 表示為低分，6~10 分設為” 1” 表示為高分，程式部分如下。

```
df = pd.read_csv("review.csv",header=None,encoding='latin1')
df[0] =df[0].apply(lambda death: 0 if death <= 5 else 1)
df=df.dropna()
```

(a)決策樹

我們運用 nltk 的 SentimentIntensityAnalyzer 將 User Review 產生一個-1~+1 之間的值，-1 代表負面情緒的最極端，+1 代表正面情緒的最極端。將 review 轉換成 SVM 檔案格式的 odm 檔，轉換格式是為了符合 mllib decision tree 的 input 格式，而 data.odm 格式如下：

```
1 1:0.9157 2:0.094 3:0.73 4:0.176
1 1:0.2671 2:0.177 3:0.628 4:0.195
1 1:0.8243 2:0.061 3:0.767 4:0.152
1 1:0.4199 3:0.640 4:0.151
1 1:0.9804 2:0.044 3:0.763 4:0.194
1 1:0.8988 3:0.766 4:0.294
1 1:0.89 2:0.12 3:0.765 4:0.115
1 1:-0.9485 2:0.276 3:0.724
1 1:-0.7081 2:0.263 3:0.627 4:0.11
1 1:0.861 2:0.125 3:0.635 4:0.238
1 1:0.765 2:0.128 3:0.631 4:0.241
1 1:0.9484 2:0.09 3:0.713 4:0.196
1 1:0.3182 3:0.9379999999999999 4:0.062
1 1:-0.8391 2:0.17 3:0.776 4:0.054
1 1:0.9166 2:0.063 3:0.644 4:0.272
1 1:0.6001 2:0.155 3:0.601 4:0.244
1 1:-0.3084 2:0.235 3:0.569 4:0.197
1 1:0.8336 3:0.849 4:0.151
1 1:0.8673999999999999 2:0.121 3:0.595 4:0.284
1 1:0.8173 2:0.132 3:0.654 4:0.215
1 1:-0.9121 2:0.211 3:0.694 4:0.095
1 1:0.4753 2:0.1 3:0.72 4:0.18
1 1:-0.6998 2:0.173 3:0.698 4:0.13
1 1:0.8622 2:0.093 3:0.5699999999999999 4:0.34
1 1:0.653 2:0.12 3:0.749 4:0.131
1 1:0.9101 2:0.141 3:0.589 4:0.27
1 1:0.9223 3:0.616 4:0.384
1 1:-0.7008 2:0.219 3:0.588 4:0.193
1 1:-0.6062999999999999 2:0.248 3:0.5580000000000001 4:0.194
1 1:0.9276 2:0.08500000000000001 3:0.76 4:0.154
1 1:0.4391 2:0.162 3:0.581 4:0.257
1 1:0.7089 3:0.763 4:0.237
1 1:0.9409999999999999 2:0.092 3:0.6919999999999999 4:0.216
1 1:0.9852 2:0.05 3:0.632 4:0.318
1 1:0.9723999999999999 2:0.05 3:0.76 4:0.19
1 1:0.9227 2:0.096 3:0.759 4:0.146
1 1:0.4976 2:0.111 3:0.741 4:0.148
1 1:-0.8979 2:0.212 3:0.62 4:0.168
1 1:-0.8462 2:0.139 3:0.761
1 1:-0.05 2:0.122 3:0.759 4:0.12
1 1:-0.8683 2:0.175 3:0.757 4:0.068
1 1:-0.6239 2:0.177 3:0.6860000000000001 4:0.137
1 1:-0.6216 2:0.162 3:0.762 4:0.076
1 1:0.9378 2:0.063 3:0.823 4:0.114
1 1:0.6477999999999999 2:0.099 3:0.799 4:0.162
1 1:0.964 2:0.081 3:0.5679999999999999 4:0.352
1 1:0.9899 2:0.122 3:0.634 4:0.244
1 1:0.8617 2:0.056 3:0.726 4:0.219
1 1:0.3864 2:0.16 3:0.633 4:0.207
1 1:0.8216 2:0.127 3:0.578 4:0.295
1 1:-0.0555 2:0.207 3:0.793
1 1:0.9117 2:0.106 3:0.6870000000000001 4:0.207
1 1:-0.952 2:0.163 3:0.763 4:0.073
1 1:0.1827 2:0.214 3:0.553 4:0.233
1 1:0.743 2:0.123 3:0.63 4:0.247
1 1:0.3818 2:0.098 3:0.776 4:0.126
1 1:0.8993 2:0.164 3:0.553 4:0.284
1 1:0.7177 2:0.104 3:0.677 4:0.219
1 1:-0.9771 2:0.151 3:0.736 4:0.114
1 1:0.6486 2:0.128 3:0.649 4:0.224
1 1:0.9316 2:0.093 3:0.621 4:0.286
1 1:-0.8258 2:0.214 3:0.552 4:0.234
1 1:0.3402 2:0.074 3:0.352 4:0.374
1 1:0.8488 2:0.044 3:0.686 4:0.35
1 1:-0.6705 2:0.206 3:0.537 4:0.177
1 1:0.9538 2:0.05 3:0.648 4:0.302
1 1:0.9324 2:0.08699999999999999 3:0.486 4:0.427
1 1:0.8875 2:0.11 3:0.668 4:0.223
1 1:0.975 2:0.093 3:0.444 4:0.463
1 1:0.8777 2:0.099 3:0.704 4:0.197
1 1:0.8507 2:0.193 3:0.742 4:0.154
1 1:0.9951 2:0.065 3:0.791 4:0.144
1 1:0.6597 3:0.838 4:0.162
1 1:0.85 3:0.786 4:0.214
1 1:0.472 2:0.096 3:0.781 4:0.123
1 1:0.8431999999999999 2:0.044 3:0.802 4:0.154
1 1:0.1513 2:0.139 3:0.718 4:0.144
1 1:0.8666 2:0.12 3:0.757 4:0.122
1 1:0.4526 2:0.135 3:0.736 4:0.13
1 1:0.6057 2:0.133 3:0.6919999999999999 4:0.175
1 1:0.7783 2:0.221 3:0.467 4:0.311
1 1:-0.25 2:0.153 3:0.673 4:0.173
1 1:0.8081 2:0.232 3:0.471 4:0.298
1 1:0.901 2:0.043 3:0.727 4:0.23
1 1:-0.5941 2:0.168 3:0.712 4:0.119
1 1:0.9387 2:0.045 3:0.645 4:0.31
1 1:0.9336 2:0.057 3:0.783 4:0.159
1 1:0.6839 2:0.135 3:0.665 4:0.2
1 1:0.7238 2:0.103 3:0.769 4:0.128
1 1:0.9269999999999999 2:0.107 3:0.6069999999999999 4:0.202
1 1:0.3382 2:0.042 3:0.873 4:0.06500000000000001
1 1:-0.2714 2:0.164 3:0.718 4:0.118
1 1:0.9941 2:0.039 3:0.61 4:0.351
```



啟用 Windows

圖 (三) 將資料轉為 SVM 格式

每筆的第一值為 Label，後面則為屬性，而每筆皆有四個屬性，每個值皆介於-1 到 1 之間：

屬性	說明
Compound	代表總和情緒的高低。
Pos	正面情緒字詞比例。
Neu	中間情緒字詞比例。
Neg	負情緒字詞比例。

使用決策樹進行分類

將 data.odm 檔以及我們的程式檔 decisiontree.py 上傳至 hadoop 後，我們便開啟 Spark 進程式。程式碼如圖（四）所示，我們將利用 predict 方法取得 TP、FP、PN、TN 等數值，再根據此四項數值計算出 Precision Rate、Recall Rate、Accuracy Rate、F1-Score 四個衡量指標。

```
from pyspark.mllib.tree import DecisionTree, DecisionTreeModel
from pyspark.mllib.util import MLUtils

# Load and parse the data file into an RDD of LabeledPoint.
data = MLUtils.loadLibSVMFile(sc, 'data.odm')
# Split the data into training and test sets (25% held out for testing)
(trainingData, testData) = data.randomSplit([0.75, 0.25])

# Train a DecisionTree model.
# Empty categoricalFeaturesInfo indicates all features are continuous.
model = DecisionTree.trainClassifier(trainingData, numClasses=2, categoricalFeaturesInfo={},
                                    impurity='gini', maxDepth=5, maxBins=32)

# predict test
predictions = model.predict(testData.map(lambda x: x.features))
#map true value and predict value
labelsAndPredictions = testData.map(lambda lp: lp.label).zip(predictions)

testErr = labelsAndPredictions.filter(lambda (v, p): v != p).count() / float(testData.count())
acc = labelsAndPredictions.filter(lambda (v, p): v == p).count() / float(testData.count())
#calculate TP FP FN TN
TP=labelsAndPredictions.filter(lambda (v, p): v ==1 and p == 1).count()
FP=labelsAndPredictions.filter(lambda (v, p): v ==0 and p == 1).count()
FN=labelsAndPredictions.filter(lambda (v, p): v ==1 and p == 0).count()
TN=labelsAndPredictions.filter(lambda (v, p): v ==0 and p == 0).count()
print(str(TP)+","+str(FP)+","+str(FN)+","+str(TN)+",")

precision=float(TP)/(TP+FP)
recall=float(TP)/(TP+FN)
print('Precision = ' + str(precision))
print('Recall= ' + str(recall))
print('Test Error = ' + str(testErr))
print('Accuracy = ' + str(acc))
print('Learned classification tree model:')
print(model.toDebugString())
```

圖（四）decisiontree.py 程式碼


```

module into which all the refactored classes and functions are moved. This module
will be removed in 0.20.
DeprecationWarning)
Traceback (most recent call last):
  File "/home/hduser/tfidf_bayes_1.py", line 28, in <module>
    conf = SparkConf().setAppName(appName)
NameError: name 'SparkConf' is not defined
hduser@master:~$ spark-submit --tfidf_bayes_1.py
/usr/local/lib/python2.7/dist-packages/sklearn/cross_validation.py:44: Deprecati
onWarning: This module was deprecated in version 0.18 in favor of the model_sele
ction module into which all the refactored classes and functions are moved. Also
note that the interface of the new CV iterators are different from that of this
module. This module will be removed in 0.20.
"This module will be removed in 0.20.", DeprecationWarning)
/usr/local/lib/python2.7/dist-packages/sklearn/grid_search.py:43: DeprecationWar
ning: This module was deprecated in version 0.18 in favor of the model_selection
module into which all the refactored classes and functions are moved. This modu
le will be removed in 0.20.
DeprecationWarning)
Traceback (most recent call last):
  File "/home/hduser/tfidf_bayes_1.py", line 28, in <module>
    conf = SparkConf().setAppName("tfidf_bayes_1")
NameError: name 'SparkConf' is not defined
hduser@master:~$ spark-submit --tfidf_bayes_1.py
/usr/local/lib/python2.7/dist-packages/sklearn/cross_validation.py:44: Deprecati
onWarning: This module was deprecated in version 0.18 in favor of the model_sele
ction module into which all the refactored classes and functions are moved. Also
note that the interface of the new CV iterators are different from that of this
module. This module will be removed in 0.20.
"This module will be removed in 0.20.", DeprecationWarning)
/usr/local/lib/python2.7/dist-packages/sklearn/grid_search.py:43: DeprecationWar
ning: This module was deprecated in version 0.18 in favor of the model_selection
module into which all the refactored classes and functions are moved. This modu
le will be removed in 0.20.
DeprecationWarning)
17/01/13 14:23:09 WARN NativeCodeLoader: Unable to load native-hadoop library fo
r your platform... using builtin-java classes where applicable
17/01/13 14:23:12 WARN Utils: Service 'SparkUI' could not bind on port 4040. Att
empting port 4041.
17/01/13 14:23:16 WARN TaskSetManager: Stage 0 contains a task of very large siz
e (513 KB). The maximum recommended task size is 100 KB.
17/01/13 14:23:20 WARN TaskSetManager: Stage 1 contains a task of very large siz
e (513 KB). The maximum recommended task size is 100 KB.
17/01/13 14:23:20 WARN TaskSetManager: Stage 2 contains a task of very large siz
e (513 KB). The maximum recommended task size is 100 KB.
17/01/13 14:23:22 WARN TaskSetManager: Stage 3 contains a task of very large siz
e (513 KB). The maximum recommended task size is 100 KB.
17/01/13 14:23:51 WARN TaskSetManager: Stage 4 contains a task of very large siz
e (513 KB). The maximum recommended task size is 100 KB.
{"f1": 0.8736842185263157, "recall": 1.0, "precision": 0.7757089345794392, "accu
racy": 0.7767441868465116}
hduser@master:~$

```

預測結果

圖（七）setiment_analyze.py 實驗結果

（三）實驗結果

實驗中，我們的正評筆數共有 166 筆，負評共有 49 筆，兩分類結果如下表：

分類方法	Precision rate	Recall rate	Accuracy
決策樹	78.4%	91.4%	74.6%
Naïve Bayes	77%	100%	77.67%

Decision tree		Actual	
		Pos	Neg
Predicted	Pos	160	44
	neg	15	14

Naïve Bayes		Actual	
		Pos	Neg
Predicted	Pos	166	48
	neg	0	1

四、分析結果與結論

從實驗結果發現其實兩種分類法在準確率上相差不遠，大約皆在 76-77%，而再進一步比較 Recall Rate、Precision Rate 發現，雖然 Precision Rate 兩者相差不遠，但在 Recall Rate 方面卻相差非常多，來到了 100%，導致 f-score 方面也有很大的差距，我們認為這種現象的產生有可能是兩個原因所造成的，第一是在資料前處理上，擷取屬性上所造成的差距，決策樹的 NLTK 僅產生四個屬性之值，然而 Naïve Bayes，由於是做 TFIDF，因此會產生非常多的屬性，導致屬性的矩陣相對稀疏，再來有可能是因為我們在正評和負評的筆數相差過多，也間接導致這樣的結果產生。

根據此問題，我們這裡提出了幾個想法或許可以改善問題：

- 可以找尋更多的訓練資料集來解決 TFIDF 屬性矩陣過於稀疏的問題。
- 在測試資料的正負評筆數上，可以試著更為均衡（正、負評測試資料比至少 $> \frac{1}{2}$ ）。
- 去除明顯正、負面之資料，以較為中立的影評用於訓練以及測試資料。

在未來的相關實作中，我們將會特別注意這方面問題，以期在各項數據能夠避免異常情況，使研究成果更具參考性。