

Configuration Manual

MSc Research Project
Programme Name

Xiyao Xu
Student ID: x18107834

School of Computing
National College of Ireland

Supervisor: XXX

National College of Ireland
MSc Project Submission Sheet



School of Computing

Student Name:xiyao xu.....
Student ID:x18107834.....
Programme:Msc in Data Analysis..... **Year:** ...2018-2019
Module:Research paper.....
Lecturer:Bahman Honari.....
Submission Due Date:December 12th 2019.....

Project Title: Analysing the impact of social media on online shopping platform sales by using sentiment analysis with text mining.....

Word Count:150..... **Page Count:**17.....

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

ALL internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

Signature:xiyao xu.....
Date:December 12th.....

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST

Attach a completed copy of this sheet to each project (including multiple copies)	<input type="checkbox"/>
Attach a Moodle submission receipt of the online project submission, to each project (including multiple copies).	<input type="checkbox"/>
You must ensure that you retain a HARD COPY of the project, both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer.	<input type="checkbox"/>

Assignments that are submitted to the Programme Coordinator Office must be placed into the assignment box located outside the office.

Office Use Only

Signature:	
Date:	
Penalty Applied (if applicable):	

Configuration Manual

Xiyao Xu

Student ID: x18107834

1 Software Requirement

Data extraction:

The data is extract by data extraction software: bazhuayu. <https://www.bazhuayu.com/>



Register for the software. Then click on the following icon.



* 任务名: 微博网页-关键词搜索

* 任务组: 我的任务组

* 用户名: User name of weibo

* 密码: Password of weibo

* 翻页次数: Pages needed to be used

搜索(1~10000): keyword

2 Data pre-processing

I used python to do pre-processing. Software used: Pycharm

Python download: <https://www.python.org/downloads/release/python-2716/>

Pycharm download: <https://www.jetbrains.com/pycharm/>

Codes for sentiment analysis pre-processing

```

1  # -*- coding: utf-8 -*-
2  __author__ = 'Bai Chenjia'
3
4  import ...
5
6  print "user dictionary loading..."
7  import sys
8  reload(sys)
9  sys.setdefaultencoding("utf-8")
10 jieba.load_userdict('F://paper/word_dictionary/pos_dic.txt')
11 jieba.load_userdict('F://paper/word_dictionary/neg_dic.txt')
12
13 # segmentation, return list
14 def segmentation(sentence):
15     seg_list = jieba.cut(sentence)
16     seg_result = []
17     for w in seg_list:
18         seg_result.append(w)
19     #print seg_result[:]
20     return seg_result
21
22
23 def postagger(sentence):
24     pos_data = pseg.cut(sentence)
25     pos_list = []
26     for w in pos_data:
27         pos_list.append((w.word, w.flag))
28     #print pos_list[:]
29     return pos_list
30
31
32 # sentence segment
33 def cut_sentence(words):
34     words = words.decode('utf8')
35     start = 0
36     i = 0
37     token = 'meaningless'
38     sents = []
39     punc_list = ',.!?;~, . ! ? ; ~... '.decode('utf8')
40     #print "punc_list", punc_list
41     for word in words:
42         #print "word", word
43         if word not in punc_list: # if it is not Apostrophe
44             segmentation()
45             for w in seg_list

```

```

41     #print "word", word
42     if word not in punt_list:  # if it is not Apostrophe
43         #print "word1", word
44         i += 1
45         token = list(words[start:i+2]).pop()
46         #print "token:", token
47     elif word in punt_list and token in punt_list:  # process Ellipsis
48         #print "word2", word
49         i += 1
50         token = list(words[start:i+2]).pop()
51         #print "token:", token
52     else:
53         #print "word3", word
54         sents.append(words[start:i+1])  # stop the sentence
55         start = i + 1
56         i += 1
57     if start < len(words):
58         sents.append(words[start:])
59     return sents
60
61 def read_lines(filename):
62     fp = open(filename, 'r')
63     lines = []
64     for line in fp.readlines():
65         line = line.strip()
66         line = line.decode("utf-8")
67         lines.append(line)
68     fp.close()
69     return lines
70
71     # remove stop words
72 def del_stopwords(seg_sent):
73     stopwords = read_lines("F://paper/stop_words.txt")  # read stop word L
74     new_sent = []  # stop removed sentence

```

Package pip, jieba, numpy are needed.

```
print
75
76     if word in stopwords:
77         continue
78     else:
79         new_sent.append(word)
80 return new_sent
81
82 # get the value of 6 lists, return the list according to the requirement
83 def read_value(request):
84     degree_dict = []
85     if request == "one":
86         degree_dict = read_lines("C://Users/dell/Desktop/Sentiment_dict/Sentiment_dict/1")
87     elif request == "two":
88         degree_dict = read_lines("C://Users/dell/Desktop/Sentiment_dict/Sentiment_dict/2")
89     elif request == "three":
90         degree_dict = read_lines("C://Users/dell/Desktop/Sentiment_dict/Sentiment_dict/3")
91     elif request == "four":
92         degree_dict = read_lines("C://Users/dell/Desktop/Sentiment_dict/Sentiment_dict/4")
93     elif request == "five":
94         degree_dict = read_lines("C://Users/dell/Desktop/Sentiment_dict/Sentiment_dict/5")
95     elif request == "six":
96         degree_dict = read_lines("C://Users/dell/Desktop/Sentiment_dict/Sentiment_dict/6")
97     else:
98         pass
99     return degree_dict
100
101
102
103 if __name__ == '__main__':
104     test_sentence1 = "这款手机大小合适"
105     test_sentence2 = "这款手机大小合适，配置也还可以，很好用，只是屏幕有点小。。。总之，戴妃+是"
106     test_sentence3 = "这手机的画面挺好，操作也比较流畅。不过拍照真的太烂了！系统也不好。"
107     """
108     seg_result = segmentation(test_sentence3) # word segmentation, sentence input, list
109     for w in seg_result:
```

3 Data analysis

3.1 Sentiment analysis

Sentiment analysis used python and pycharm, which the same with section2.

Code:


```

1      # -*- coding: utf-8 -*-
2      __author__ = 'Bai Chenjia'
3
4      import text_process as tp
5      import numpy as np
6
7
8
9      print "Begin"
10     posdict = tp.read_lines("C://Users/dell/Desktop/Sentiment_dict/Sentiment_dict/emotion_di
11     negdict = tp.read_lines("C://Users/dell/Desktop/Sentiment_dict/Sentiment_dict/emotion_di
12
13     most_degree = tp.read_lines('C://Users/dell/Desktop/Sentiment_dict/Sentiment_dict/degree
14     very_degree = tp.read_lines('C://Users/dell/Desktop/Sentiment_dict/Sentiment_dict/degree
15     more_degree = tp.read_lines('C://Users/dell/Desktop/Sentiment_dict/Sentiment_dict/degree
16     ish_degree = tp.read_lines('C://Users/dell/Desktop/Sentiment_dict/Sentiment_dict/degree_
17     insufficient_degree = tp.read_lines('C://Users/dell/Desktop/Sentiment_dict/Sentiment_dic
18     inverse_degree = tp.read_lines('C://Users/dell/Desktop/Sentiment_dict/Sentiment_dict/degre
19     print "The degree files are loaded"
20
21     emotion_level1 = "sad"
22     emotion_level2 = "angry"
23     emotion_level3 = "calm"
24     emotion_level4 = "peaceful"
25     emotion_level5 = "happy"
26
27     emotion_level6 = "Emotion fluctuation is small"
28     emotion_level7 = "Emotion fluctuation is big"
29
30
31     # 2.Adverb of degree processing, according to the type of degree adverb, multiplied by d
32     def match(word, sentiment_score):
33         if word in most_degree:
34             sentiment_score *= 2.0

```

```
text_process.py × dict_main.py ×
Q+
35     elif word in very_degree:
36         sentiment_score *= 1.75
37     elif word in more_degree:
38         sentiment_score *= 1.5
39     elif word in ish_degree:
40         sentiment_score *= 1.2
41     elif word in insufficient_degree:
42         sentiment_score *= 0.5
43     elif word in inverse_degree:
44         #print "inversedict", word
45         sentiment_score *= -1
46     return sentiment_score
47
48     # Example: [5, -2] → [7, 0]; [-4, 8] → [0, 12]
49     def change_to_positive_number(positive_count, negative_count):
50         pos_count = 0
51         neg_count = 0
52         if positive_count < 0 and negative_count <= 0:
53             neg_count += negative_count - positive_count
54             pos_count = 0
55         elif negative_count < 0 and positive_count < 0:
56             pos_count = positive_count - negative_count
57             neg_count = 0
58         elif positive_count < 0 and negative_count >= 0:
59             neg_count = -positive_count
60             pos_count = -negative_count
61         else:
62             pos_count = positive_count
63             neg_count = negative_count
64         return (pos_count, neg_count)
65
66
67     # Find the total sentiment score of a single Weibo sentence
68     def single_sentence_sentiment_score(weibo_sent):
```

```

68 def single_sentence_sentiment_score(weibo_sent):
69     single_review_senti_score = []
70     cuted_review = tp.cut_sentence(weibo_sent) #Sentence segmentation, analyzing each s
71
72     for sent in cuted_review:
73         seg_sent = tp.segmentation(sent) # segmentation
74         seg_sent = tp.del_stopwords(seg_sent)[:]  

75         #for w in seg_sent:
76         # print w,
77         i = 0 # Record the position of the scanned word
78         s = 0 # Record the position of the sentiment word
79         poscount = 0 # Record the positive emotion score in the segmented sentence
80         negcount = 0 # Record the negative emotion score in the segmented sentence
81
82         for word in seg_sent: # analysis sentence by sentence
83             #print word
84             if word in posdict:
85                 #print "posword:", word
86                 poscount += 1
87                 for w in seg_sent[s:i]:
88                     poscount = match(w, poscount)
89                 #print "poscount:", poscount
90                 s = i + 1
91
92             elif word in negdict:
93                 #print "negword:", word
94                 negcount += 1
95                 for w in seg_sent[s:i]:
96                     negcount = match(w, negcount)
97                 #print "negcount:", negcount
98                 s = i + -1
99
100
101     elif word == "! ".decode("utf-8") or word == "! ".decode('utf-8'):

```

```
text_process.py x dict_main.py x
Q
102         for w2 in seg_sent[::-1]:
103             if w2 in posdict:
104                 poscount += 2
105                 break
106             elif w2 in negdict:
107                 negcount += 2
108                 break
109         i += 1
110
111         single_review_senti_score.append(change_to_positive_number(poscount, negcount))
112     pos_result, neg_result = 0, 0 # Record the positive sentiment score and negative s
113     for res1, res2 in single_review_senti_score:
114         pos_result += res1
115         neg_result += res2
116     #print pos_result, neg_result
117     result = pos_result - neg_result # final score of the sentence
118     result = round(result, 1)
119     return result
120
121     """
122     # 测试
123     weibo_sent = "这手机的画面挺好，操作也比较流畅。不过拍照真的太烂了！系统也不好。"
124     score = single_review_sentiment_score(weibo_sent)
125     print score
126     """
127
128
129     def run_score():
130         fp_test = open('F://paper/code/test_data.txt', 'r')
131         contents = []
132         for content in fp_test.readlines():
133             content = content.strip()
134             content = content.decode("utf-8")
135             contents.append(content)
```



```
text_process.py x dict_main.py x
Q-
170
171     pos_number = len(pos_list)
172     neg_number = len(neg_list)
173     mid_number = len(middle_list)
174     total_number = pos_number + neg_number + mid_number
175     number_ratio = pos_number/neg_number
176     pos_number_ratio = round(float(pos_number)/float(total_number), 2)
177     neg_number_ratio = round(float(neg_number)/float(total_number), 2)
178     mid_number_ratio = round(float(mid_number)/float(total_number), 2)
179     text_pos_number = "积极微博条数为 " + str(pos_number) + " 条, 占全部微博比例的 %" + str(
180     text_neg_number = "消极微博条数为 " + str(neg_number) + " 条, 占全部微博比例的 %" + str(
181     text_mid_number = "中性情感微博条数为 " + str(mid_number) + " 条, 占全部微博比例的 %" +
182
183     pos_array = np.array(pos_list)
184     neg_array = np.array(neg_list)
185     total_array = np.array(total_list)
186     pos_mean = pos_array.mean()
187     neg_mean = neg_array.mean()
188     total_mean = total_array.mean()
189     mean_ratio = pos_mean/neg_mean
190     if pos_mean <= 6:
191         text_pos_mean = emotion_level4
192     else:
193         text_pos_mean = emotion_level5
194     if neg_mean >= -6:
195         text_neg_mean = emotion_level2
196     else:
197         text_neg_mean = emotion_level1
198     if total_mean <= 6 and total_mean >= -6:
199         text_total_mean = emotion_level3
200     elif total_mean > 6:
201         text_total_mean = emotion_level4
202     else:
203         text_total_mean = emotion_level2
```



```

203 |     text_total_mean = emotion_level12
204
205     pos_variance = pos_array.var(axis=0)
206     neg_variance = neg_array.var(axis=0)
207     total_variance = total_array.var(axis=0)
208     var_ratio = pos_variance/neg_variance
209     #print "pos_variance:", pos_variance, "neg_variance:", neg_variance, "var_ration:",
210     if total_variance > 10:
211         text_total_var = emotion_level7
212     else:
213         text_total_var = emotion_level6
214     #####
215     result_dict = {}
216     result_dict['pos_number'] = pos_number
217     result_dict['neg_number'] = neg_number
218     result_dict['mid_number'] = mid_number
219     result_dict['number_ratio'] = round(number_ratio, 1)
220     result_dict['pos_mean'] = round(pos_mean, 1)
221     result_dict['neg_mean'] = round(neg_mean, 1)
222     result_dict['total_mean'] = round(total_mean, 1)
223     result_dict['mean_ratio'] = abs(round(mean_ratio, 1))
224     result_dict['pos_variance'] = round(pos_variance, 1)
225     result_dict['neg_variance'] = round(neg_variance, 1)
226     result_dict['total_variance'] = round(total_variance, 1)
227     result_dict['var_ratio'] = round(var_ratio, 1)
228
229     result_dict['text_pos_number'] = text_pos_number
230     result_dict['text_neg_number'] = text_neg_number
231     result_dict['text_mid_number'] = text_mid_number
232     result_dict['text_pos_mean'] = text_pos_mean
233     result_dict['text_neg_mean'] = text_neg_mean
234     result_dict['text_total_mean'] = text_total_mean
235     result_dict['text_total_var'] = text_total_var
236     """

```

```

if __name__ == '__main__':
    results = run_score()
    write_results(results)
    result_dict = handel_result(results)

```

3.2 TF-IDF

It is processed on the google colab:
<https://colab.research.google.com/notebooks/welcome.ipynb>



```
#server connection
!apt-get install -y -qq software-properties-common python-software-properties modul
!add-apt-repository -y ppa:alessandro-strada/ppa 2>&1 > /dev/null
!apt-get update -qq 2>&1 > /dev/null
!apt-get -y install -qq google-drive-ocamlfuse fuse
from google.colab import auth
auth.authenticate_user()
from oauth2client.client import GoogleCredentials
creds = GoogleCredentials.get_application_default()
import getpass
!google-drive-ocamlfuse -headless -id={creds.client_id} -secret={creds.client_secret}
vcode = getpass.getpass()
!echo {vcode} | google-drive-ocamlfuse -headless -id={creds.client_id} -secret={cre

#direction
!ls

!mkdir -p drive
!google-drive-ocamlfuse drive
!ls

# go the the working directory
import os

os.chdir("drive/paper")
!ls

#keyword extraction

import jieba
import math
import jieba.analyse
```

Code:


```

import math
import jieba.analyse

class TF_IDF:

    def __init__(self, file, stop_file):
        self.file = file
        self.stop_file = stop_file
        self.stop_words = self.getStopWords()

    def getStopWords(self):
        swlist=list()
        for line in open(self.stop_file, "r", encoding="utf-8").readlines():
            swlist.append(line.strip())
        print("stop words loaded...")
        return swlist

    def loadData(self):
        dMap = dict()
        for line in open(self.file, "r", encoding="utf-8").readlines():
            id,title = line.strip().split("\t")
            dMap.setdefault(id, [])
            for word in list(jieba.cut(str(title).replace(" ", "")), cut_all=False):
                if word not in self.stop_words:
                    dMap[id].append(word)
        print("load products and its title, using jieba to segment and del")
        return dMap

    def getFreqWord(self, words):
        freqWord = dict()
        for word in words:
            freqWord.setdefault(word, 0)
            freqWord[word] += 1
        return freqWord

```

3.3 Outlier process:

```

        count += 1
    for key in dMap.keys():
        if word in dMap[key]:
            count += 1
    return count

def getTFIDF(self, words, dMap):
    outDic = dict()
    freqWord = self.getFreqWord(words)
    for word in words:

        tf = freqWord[word]*1.0 / len(words)

        idf = math.log(len(dMap)/(self.getCountWordInFile(word, dMap)+1))
        tfidf = tf * idf
        outDic[word] = tfidf

    orderDic = sorted(outDic.items(), key=lambda x:x[1], reverse=True)
    return orderDic

def getTag(self, words):

    print(jieba.analyse.extract_tags(words, topK=200, withWeight=True))

if __name__ == "__main__":

    file = "./taobao_tfidf.txt"

    stop_file = "./stop_words.txt"

    tfidf=TF_IDF(file, stop_file)
    tfidf.getTag(open("./taobao_tfidf.txt", "r", encoding="utf-8").read(),)

```

Using R and R studio: link: <https://cran.r-project.org/bin/windows/base/>,
<https://rstudio.com/products/rstudio/download/>

Code:

Winsorizing:

```

1 data <- c(756153,534383,458519,24456,31138,2617734,75549,2905154,26520,3
2 ,374737,337104,1164074,261359,1898800,137491,486411,2178952,85561,56796,
3 ,1480950,388998,226727,1529637,48183,1510995,144798,115717,302831,442244
4 ,223083,509471,69630,64380,167594,301439,156588,766846,1258,213744,38219
5 ,450207,209583,84926,861646,888573,1805145,62196,231195,220568,929129,66
6
7 data
8
9 length(data)
10 summary(data)
11 benth <-783632 + 1.5*IQR(data)
12 benth
13
14 data[data > benth]
15 data[data > benth] <-benth
16 data
17 summary(data)
18 boxplot(data)
19
20

```

Robust regression:

```

1 install.packages("car")
2 install.packages("foreign")
3 install.packages("MASS")
4 library(foreign)
5 library(MASS)
6 require(MASS)
7 require(foreign)
8 library(car)
9 library(carData)
10
11 mydata<-read.csv("robost_regression.csv",header = T)
12 str(mydata)
13 attach(mydata)
14 summary(mydata)
15
16 plot(cooks.distance(lm(Sales ~ Popularity, data = mydata)))
17 qqnorm(Sales);qqline(Sales)
18
19 summary(ols <- lm(Sales ~ Popularity, data = mydata))
20
21 dist <- cooks.distance(ols)
22 dist<-data.frame(dist)
23 s <- stdres(ols)
24 a <- cbind(mydata, dist, s)
25
26 sabs <- abs(s)
27 a <- cbind(mydata, dist, s, sabs)
28 asorted <- a[order(-sabs), ]
29 asorted[1:10, ]
30
31 summary(rr.huber <- rlm(Sales ~ Popularity, data = mydata))
32
33 huber <- data.frame(usage = mydata$Sales, resid = rr.huber$resid, weight = rr.huber$w)
34 huber2 <- huber[order(rr.huber$w), ]
35 huber2[1:10, ]
36
37
38 rr.bisquare <- rlm(Sales ~ Popularity, data = mydata, psi = psi.bisquare)
39 summary(rr.bisquare)
40 bisqr <- data.frame(Sales = mydata$Sales, resid = rr.bisquare$resid, weight = rr.bisquare$w)
41 bisqr2 <- bisqr[order(rr.bisquare$w), ]
42 bisqr2[1:10, ]
43

```

References:

Csdn.net. (2019). [online] Available at:
https://blog.csdn.net/gamer_gyt/article/details/85690389 [Accessed 12 Dec. 2019].

Csdn.net. (2019). [online] Available at: <https://blog.csdn.net/chenpe32cp/article/details/77801600>
[Accessed 12 Dec. 2019].

Robust Regressions in R. (2018). YouTube. Available at:
<https://www.youtube.com/watch?v=IU25mxM4mhs> [Accessed 12 Dec. 2019].

Outlier Treatment in R - Part 2 - Winsorizing. (2016). YouTube. Available at:
<https://www.youtube.com/watch?v=JTgLDDyAuF0> [Accessed 12 Dec. 2019].