

# Final Presentation

<Team 4>

12146323 Jeong Sang-Heon

14146320 Lee Jun-Ha

14146326 Hong Jun-Ki

16102278 Sung Chang-Keu

# Final Presentation

<Team 4>

12146323 Jeong Sang-Heon

14146320 Lee Jun-Ha

14146326 Hong Jun-Ki

16102278 Sung Chang-Keu

## Contents

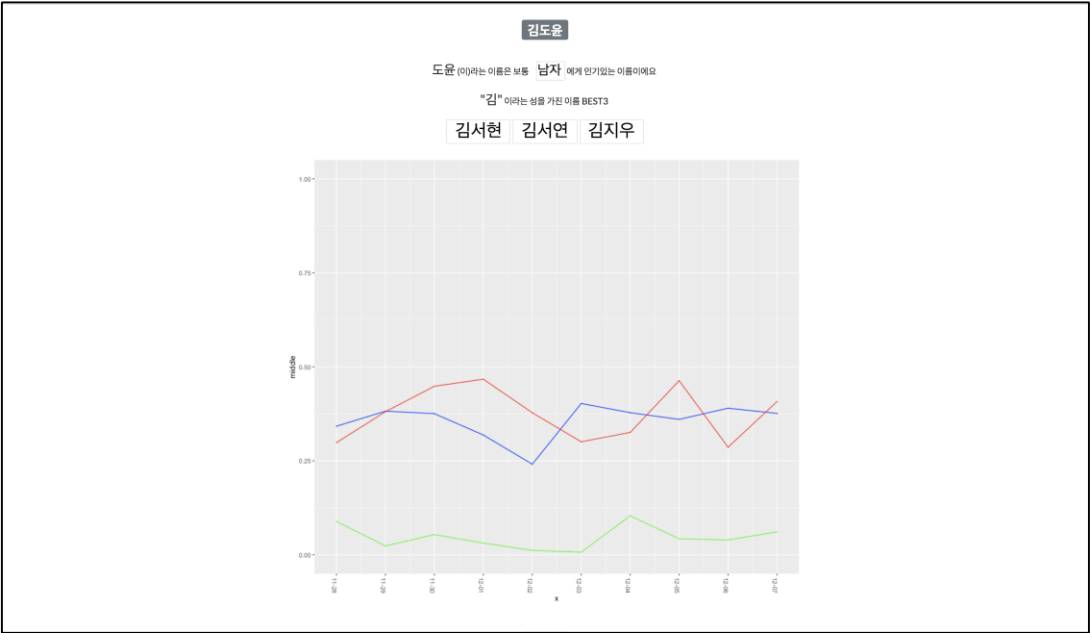
1. The result and methods of implementation
2. The used techniques
3. The overall steps of testing
4. The steps of analyzing the problems and debugging

# 1. The result and methods of implementation

## 1.1 Result of implementation

### 1.1.1 User Interface

- Tweets & Statistic survey of name
- Based on web service
- Search specific name → sentimental score graph will be shown
- Show the names of the top three people who have same last name.



### 1.1.2 Compare with Requirement specification

| Functional Requirement     | System | Sentiment Database     | Select the appropriate information from the DB and print it to the user        |
|----------------------------|--------|------------------------|--|
|                            | User   | Web Page               | Users search through web pages to get the information they want                |
| Non-functional Requirement | System | Language               | System language type   |
|                            |        | Accessibility          | The users can use service without difficulty and don't need guide to implement |
|                            |        | Update Cycle           | Update system's modification   |
|                            |        | Data-backup Cycle      | Backup the system data   |
|                            |        | Information Accuracy   | Whether the system can provide users with validated data                       |
|                            |        | Cope with Errors       | How to deal with errors  |
|                            |        | Data Access            | Correct access for the data  |
|                            | User   | Search the Preferences | Users can check the preferences for smartphone                                 |

<Requirement specification – SRS>

# 1. The result and methods of implementation

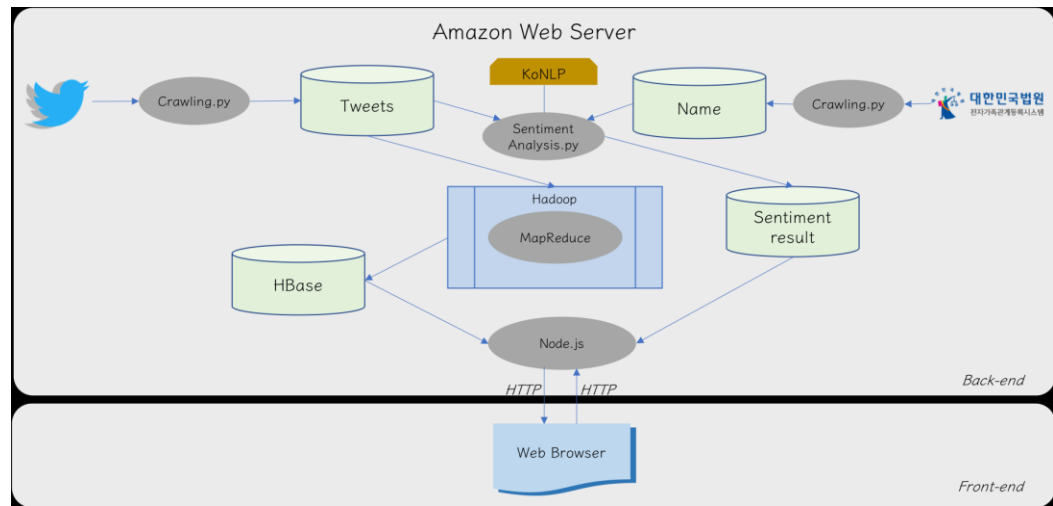
## 1.1.2 Compare with Requirement specification

| Requirements           | Score | Detailed evaluation   |
|------------------------|-------|---|
| Sentiment database     | 4     | Database deliver accurate information to certain input  |
| Web page               | 5     | User can view all the information provided on the web page  |
| Language               | 1     | Only Korean is supported  |
| Accessibility          | 5     | Web pages are easy to use because they are simple   |
| Update Cycle           | 3     | Since the maximum range of data crawl using Twitter is 10 days, the data must be updated once every 10 days. - Increased cost |
| Data backup Cycle      | 3     |   |
| Information accuracy   | 2     | Insufficient validation of result   |
| Cope with errors       | 4     | Can handle the error which is caused by user input  |
| Data access            | 5     | Provide results using only information about the name user enter  |
| Search the preferences | 4     | Shows the user's preference for the name they enter in scores and graphs  |

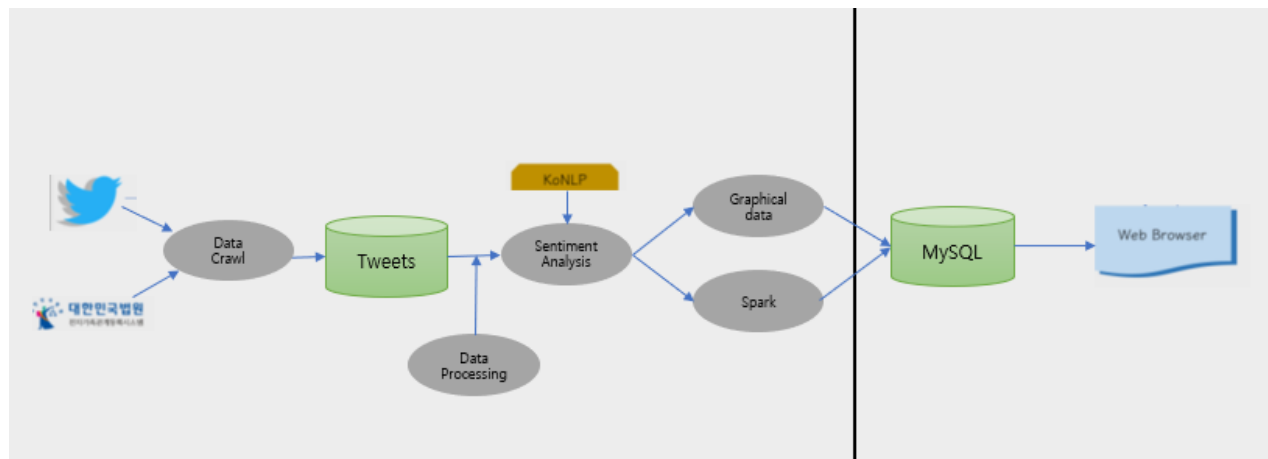
# 1. The result and methods of implementation

## 1.2 Method of implementation

### 1.2.1 Initial service flow



### 1.2.2 Final service flow



| Methods            | Language    | Implementation                                   | Input         | Output        |
|--------------------|-------------|--|---------------|---------------|
| Crawl              | Python      | Crawl tweets from web site                       | (tweets) txt  | json          |
| Data Preprocessing | Python      | JSON type transformation                         | json          | json*         |
| Sentiment Analysis | Java Script | Utilize sentiment analyzer library               | json          | (score) json  |
| Graph generating   | R           | Generate time-series graph                       | json          | jpg           |
| Spark processing   | pyspark     | Transform DataFrame to calculate key-value pairs | DataFrame     | txt           |
| Database           | MySQL       | ETL graph and score data                         | txt / jpg dir | txt / jpg dir |

## 2. The used techniques

### 2.1 The techniques for crawling the data – Tweepy package



```
def get_tweets(keyword, num_limit):
    i=0

    os.mkdir("\\Users\\hong\\Desktop\\크롤링데이터\\"+keyword)

    for tweet in tweepy.Cursor(api.search, q=keyword, since='2018-12-08', until='2018-12-06').items(num_limit):
        StatusObject = tweet.json
        dict1 = {
            'id': StatusObject['id_str'],
            'permalink': "",
            'username': StatusObject['user']['name'],
            'text': StatusObject['text'],
            'date': StatusObject['created_at'],
            'retweets': StatusObject['retweet_count'],
            'favorites': StatusObject['favorite_count'],
            'mentions': StatusObject['entities']['user_mentions'],
            'hashtags': StatusObject['entities']['hashtags'],
            'geo': StatusObject['geo']
        }

        print(str(i)+'-'+str(dict1))

        ##### 날짜 형태 변환하기 #####
        unformatted = StatusObject['created_at']
        # Use re to get rid of the milliseconds.
        remove_ms = lambda x: re.sub("\\d+", "", x)
        # Make the string into a datetime object.
        mk_dt = lambda x: datetime.strptime(remove_ms(x), "%a %b %d %H:%M:%S %Y")
        # Format your datetime object.
        my_fmt = lambda x: "%Y-%m-%d".format(mk_dt(x))
```

- 2018-11-29.json
- 2018-11-30.json
- 2018-12-01.json
- 2018-12-02.json
- 2018-12-03.json
- 2018-12-04.json
- 2018-12-05.json
- 2018-12-06.json

Tweets mentioning the name



Codes for crawling the tweets



Crawled data saved as JSON files

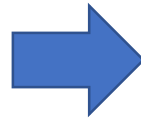
## 2. The used techniques

### 2.2 The techniques for processing the data – Tweepy package

```
[{"id": "1068567487413813248",
 "permalink": "",
 "username": "상요",
 "text": "RT @eungalawyer: 1방남자들이(기바람 선동으로) 여학생 대상으로 인기투표해서 신경안쓰는척하지만",
 "date": "Fri Nov 30 18:08:36 +0000 2018",
 "retweets": 29,
 "favorites": 0,
 "mentions": [
  {
    "screen_name": "eungalawyer",
    "name": "응가변호사",
    "id": "1025015594771349504",
    "id_str": "1025015594771349504",
    "indices": [
      3,
      15
    ]
  }
 ],
 "hashtags": [],
 "geo": null
 },
 {"id": "1068551330048892928",
 "permalink": "",
 "username": "강현우",
 "text": "(그걸 또 들은건지 해맑게 웃으며 재현이 끌어안아) 당당히 강현우는 유재현꺼쨌~ 아, 귀여워!",
 "date": "Fri Nov 30 17:04:24 +0000 2018",
 "retweets": 0,
 "favorites": 0,
 "mentions": [],
 "hashtags": [],
 "geo": null
 },
 {"id": "1068490776345174017",
 "permalink": "",
 "username": "유재현",
 "text": "흔해도 예쁜 이름이 있잖아요. 난 마음에 드는데 강현우",
 "date": "Fri Nov 30 13:03:47 +0000 2018",
 "retweets": 0
 }
 ]

path = '/Users/hong/Desktop/크롤링데이터/'+names
extension = 'json'
os.chdir(path)
result = [i for i in glob.glob('*.*').format(extension))]

for file in result:
    open_file = open(file, 'r', encoding = 'UTF8')
    read_file = open_file.read()
    new_content = "[" + read_file + "]"
    new_content2 = new_content.replace(",]",",]")
    write_file = open(file, 'w')
    write_file.write(new_content2)
    write_file.close()
```



Pre-process JSON format data to implement  
Morpheme Analyze Server & Twitter Sentiment Analysis  
without errors securely

```
pip3 install konlpy
ksa-server localhost 7000
```

(Morpheme analyzer server)

```
pip3 install konlpy
pip3 install JPype1
ksa-client localhost 7000 2017-01-01 2017-02-01 ./sourceDir ./targetDir
```

(Twitter sentiment analysis)

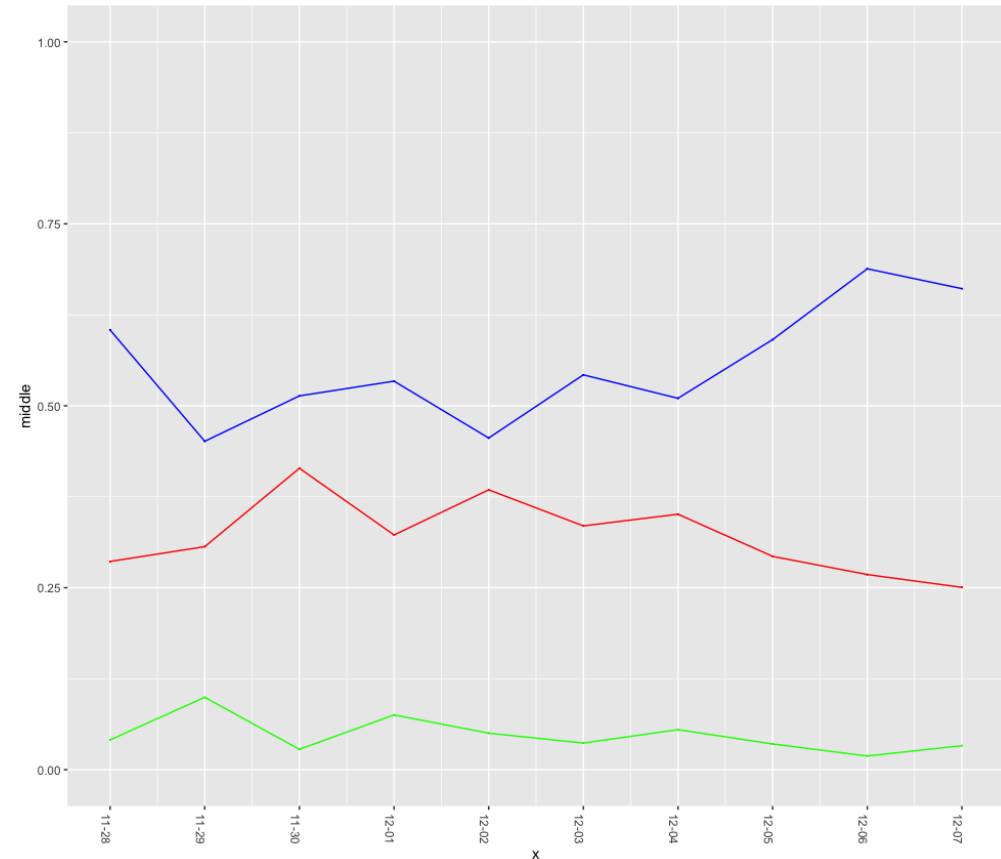
```
[{"id": "1067740711427104768",
 "permalink": "",
 "username": "라오",
 "text": "도윤이 생각난김에 비설 올라자 에버노트에 올린건 아닌데 이거 생각하고 편",
 "date": "Wed Nov 28 11:23:17 +0000 2018",
 "retweets": 0,
 "favorites": 4,
 "mentions": [],
 "hashtags": [],
 "geo": null,
 "raw": "도윤이 생각난김에 비설 올라자 에버노트에 올린건 아닌데 이거 생각하고 편",
 "polarity": {"com": 0.15002740612142465,
 "pos": 0.4907180830727546,
 "neg": 0.25983683296288046,
 "neut": 0.03922011399439713,
 "none": 0.060197563848543204},
 "intensity": {
  "high": 0.19258227442857145,
  "low": 0.05744768228571429,
  "meium": 0.6888359331428572,
  "none": 0.06113411014285715},
 "expressive": {
  "dir-action": 0.0015509855716501408,
  "dir-explicit": 0.49228916149889845,
  "dir-speech": 0.265855048752265,
  "indirect": 0.18314091745473443,
  "writing-device": 0.057163886722451984}},
 {"id": "1067591229850845184",
 "permalink": "",
 "username": "초키해보고지구부순 사보",
 "text": "김도윤 심장에 스며든다,,,진짜 미쳐ㅏㅏㅏ 미세먼지 오지는날에도 힘",
 "date": "Wed Nov 28 01:29:18 +0000 2018",
 "retweets": 0,
 "favorites": 0,
 "mentions": [],
 "hashtags": [],
 "geo": null,
 "raw": "김도윤 심장에 스며든다,,,진짜 미쳐ㅏㅏㅏ 미세먼지 오지는날에도 힘",
 "polarity": {"com": 0.15002740612142465,
 "pos": 0.4907180830727546,
 "neg": 0.25983683296288046,
 "neut": 0.03922011399439713,
 "none": 0.060197563848543204},
 "intensity": {
  "high": 0.19258227442857145,
  "low": 0.05744768228571429,
  "meium": 0.6888359331428572,
  "none": 0.06113411014285715},
 "expressive": {
  "dir-action": 0.0015509855716501408,
  "dir-explicit": 0.49228916149889845,
  "dir-speech": 0.265855048752265,
  "indirect": 0.18314091745473443,
  "writing-device": 0.057163886722451984}}
```

(The results of sentiment analysis  
which are shown as  
positive/negative/neutral value,  
will be used for making the graph by  
R, and displaying the most preferred  
names)

## 2. The used techniques

### 2.2 The techniques for processing the data – R script

```
70 ggplot.sentiment <- function (dataset) {  
71   # 감성그래프  
72   names = as.Date(unlist(list.map(dataset, return(data))))  
73   pos = unlist(list.map(dataset, return(polarity$pos)))  
74   neg = unlist(list.map(dataset, return(polarity$neg)))  
75   neut = unlist(list.map(dataset, return(polarity$neut)))  
76   df = data.frame(  
77     x = names,  
78     middle = pos*0+0.5,  
79     pos = pos,  
80     pos.label = 'Positive',  
81     neg = neg,  
82     neg.label = 'Negative',  
83     neut = neut,  
84     neut.label = 'Neutral'  
85   )  
86   color.pos = "blue"  
87   color.neg = "red"  
88   color.neut = "green"  
89   size.dot = 0  
90   gp = ggplot(df, aes(x, middle)) ## + geom_line()  
91   gp + ylim(0, 0.5) +  
92     geom_point(aes(x, pos), color = color.pos, size = size.dot) + geom_line(aes(x, pos), color = color.pos) +  
93     geom_point(aes(x, neg), color = color.neg, size = size.dot) + geom_line(aes(x, neg), color = color.neg) +  
94     geom_point(aes(x, neut), color = color.neut, size = size.dot) + geom_line(aes(x, neut), color = color.neut) +  
95     scale_x_date(date_labels = "%m-%d", date_breaks = "1 day") +  
96     theme(axis.text.x = element_text(angle = 270, hjust = 1))  
97     ## theme(axis.text.x = element_blank(),  
98     ##       axis.ticks = element_blank())  
99 }
```



Visualize the result of sentiment analysis to graph by R



## 2. The used techniques

### 2.2 The techniques for processing the data – Spark

<JSON object>

/ sentiment / 2018-12-06.json

```
[{"id":"1070710713209896960","permalink":"","username":"Today 김 다원[도윤] / KR / 27Y / 185cm ","date":"Thu Dec 06 16:05:01 +0000 2018","retweets":0,"favorites":0,"mentions":[],"hashtags":[],"geo":null,"raw":"김 다원[도윤] / KR / 27Y / 185cm https://t.co/iSJoYSdiuI","polarity":{"com":0.007142857,"pos":0.214285714,"neg":0.328571429,"neut":0.014285714,"none":0.435714286},"intensity":{"high":0.142857143,"low":0.114285714,"meium":0.307142857,"none":0.435714286},"expressive":{"dir-action":0,"dir-explicit":0.1857142858142857,"dir-speech":0.2999999969999996,"indirect":0.0857142859142857,"writing-device":0.42857142857142855}},{"id":"1070708466321260545","permalink":"","username":"n인 트친소 맞팔 시운","text":"김도윤 오빠만 ","date":"Thu Dec 06 15:56:05 +0000 2018","retweets":0,"favorites":1,"mentions":[],"hashtags":[],"geo":null,"raw":"김도윤 오빠만 https://t.co/vGzNPi41G3","polarity":{"com":0.0035714285,"pos":0.607142857,"neg":0.1642857145,"neut":0.007142857,"none":0.217857143},"intensity":{"high":0.0714285715,"low":0.057142857,"meium":0.6535714285,"none":0.217857143},"expressive":{"dir-action":0,"dir-explicit":0.09285714295357143,"dir-speech":0.64999999675,"indirect":0.042857142978571426,"writing-device":0.21428571439285712}},{"id":"1070683636238700545","permalink":"","username":"해나","text":"@N00Y0D 제가 원피스가 참 안 어울리더라고요 그 김에 다이어트 해서 저거 입을게용. 입구 도윤님이랑 데이트 가야징 ~~~♥","date":"Thu Dec 06 14:17:25 +0000 2018","retweets":0,"favorites":0,"mentions":[{"screen_name":"N00Y0D","name":"도윤","id":1032808823856160800,"id_str":"1032808823856160768","indices":[0,7]}],"hashtags":[],"geo":null,"raw":"@N00Y0D 제가 원피스가 참 안 어울리더라고요 그 김에 다이어트 해서 저거 입을게용. 입구 도윤님이랑 데이트 가야징 ~~~♥ https://t.co/tpPg aPDugv","polarity":{"com":0.20587027919999998,"pos":0.3965106731999999,"neg":0.28581691299999995,"neut":0.054864532,"none":0.05693760259999999},"intensity":{"high":0.31763136293647376,"low":0.046715927790656815,"meium":0.5746305416850739,"none":0.061027216758779556}
```

<Data parse in type of data frame>

```
cloudera@quickstart:~  
File Edit View Search Terminal Help  
>>>  
>>> df = sqlContext.read.json ("/sentiment/2018-12-06.json")  
>>> df2 = df.select("polarity")  
>>> df2.collect()  
[Row(polarity=Row(com=0.007142856999999999, neg=0.328571429000000003, neut=0.014285714, none=0.0035714284999999999, neg=0.164285714500000001, neut=0.0071428569999999999, none=0.2178571999999, neg=0.28581691299999995, neut=0.054864532000000001, none=0.056937602599999991, pos=0.3071428571429000000003, neut=0.014285714, none=0.435714286000000001, pos=0.21428571399999999)),  
87, neut=0.023765480285714288, none=0.21051374728571429, pos=0.54828300821428566)), Row(polarity=Row(com=0.0035714285, pos=0.607142857, neg=0.1642857145, neut=0.007142857, none=0.217857143, pos=0.60714285700000004)),  
Row(polarity=Row(com=0.20587027919999998, neg=0.28581691299999995, neut=0.054864532000000001, none=0.056937602599999991, pos=0.39651067319999999)),  
Row(polarity=Row(com=0.0071428569999999999, neg=0.328571429000000003, neut=0.014285714, none=0.435714286000000001, pos=0.21428571399999999)),  
Row(polarity=Row(com=0.0082070268571428576, neg=0.20923073735714287, neut=0.023765480285714288, none=0.21051374728571429, pos=0.54828300821428566)),  
Row(polarity=Row(com=0.025030213644645017, neg=0.40104256917150305, neut=0.12338104079452722, none=0.090807898077914867, pos=0.35973827831140986))]  
>>>  
>>>  
>>>  
>>>  
>>>  
>>>  
>>>  
>>>  
>>>  
>>>  
>>> df = sqlContext.read.json ("/sentiment/2018-12-06.json")  
>>> df2 = df.select("polarity")  
>>> df2.collect()  
[Row(polarity=Row(com=0.0071428569999999999, neg=0.328571429000000003, neut=0.014285714, none=0.435714286000000001, pos=0.21428571399999999)), Row(polarity=Row(com=0.0035714284999999999, neg=0.164285714500000001, neut=0.0071428569999999999, none=0.217857143, pos=0.60714285700000004)),  
Row(polarity=Row(com=0.20587027919999998, neg=0.28581691299999995, neut=0.054864532000000001, none=0.056937602599999991, pos=0.39651067319999999)),  
Row(polarity=Row(com=0.0071428569999999999, neg=0.328571429000000003, neut=0.014285714, none=0.435714286000000001, pos=0.21428571399999999)),  
Row(polarity=Row(com=0.0082070268571428576, neg=0.20923073735714287, neut=0.023765480285714288, none=0.21051374728571429, pos=0.54828300821428566)),  
Row(polarity=Row(com=0.025030213644645017, neg=0.40104256917150305, neut=0.12338104079452722, none=0.090807898077914867, pos=0.35973827831140986))]  
>>>
```

## 2. The used techniques

### 2.2 The techniques for processing the data – Spark

<Data parse in type of Spark Context>

```
>
> df = sqlContext.read.json("/sentiment/*")
> df2 = df.select("polarity")
> df2.coalesce(1).rdd.saveAsTextFile("/sentiment/result.txt")
>
-
```

/ sentiment / result.txt / part-00000

```
Row(polarity=Row(com=0.15002740612142465, neg=0.25983683296288046, neut=0.03922011399439
7128, none=0.060197563848543204, pos=0.49071808307275461))
Row(polarity=Row(com=0.12971109175000001, neg=0.30588705012500006, neut=0.196922183375,
none=0.228433507625, pos=0.13904616712499995))
Row(polarity=Row(com=0.025864422012932212, neg=0.32989081016494542, neut=0.0294358505147
17924, none=0.21785714310892859, pos=0.39695177419847588))
Row(polarity=Row(com=0.15002740612142465, neg=0.25983683296288046, neut=0.03922011399439
7128, none=0.060197563848543204, pos=0.49071808307275461))
Row(polarity=Row(com=0.12971109175000001, neg=0.30588705012500006, neut=0.196922183375,
none=0.228433507625, pos=0.13904616712499995))
Row(polarity=Row(com=0.025864422012932212, neg=0.32989081016494542, neut=0.0294358505147
17924, none=0.21785714310892859, pos=0.39695177419847588))
Row(polarity=Row(com=0.14464285703616073, neg=0.42142857160535718, neut=0.02142857125535
7142, none=0.10892857152723215, pos=0.30357142857589281))
Row(polarity=Row(com=0.0066244239999999994, neg=0.29531490025000001, neut=0.024539170500
000002, none=0.19925115199999999, pos=0.47427035325))
Row(polarity=Row(com=0.0037274219999999998, neg=0.246174056, neut=0.016650246200000003,
none=0.22622331700000001, pos=0.50722495879999996))
Row(polarity=Row(com=0.008832565333333332, neg=0.56041986700000013, neut=0.0327188940000
00005, none=0.15455709166666667, pos=0.24347158199999999))
Row(polarity=Row(com=0.14464285703616073, neg=0.42142857160535718, neut=0.02142857125535
7142, none=0.10892857152723215, pos=0.30357142857589281))
```

```
///
>>> rdd = sc.textFile("/sentiment/result.txt/*")
>>> rdd.take(1)
[[u'Row(polarity=Row(com=0.15002740612142465, neg=0.25983683296288046, neut=0.039
220113994397128, none=0.060197563848543204, pos=0.49071808307275461))']]
>>>
>>> rdd1 = rdd.map(lambda line:line.split(",")[1])\
...             .map(lambda line:line.split("="))\
...             .map(lambda x: (x[0], float(x[1])))
>>> rdd1.take(1)
[[('neg', 0.25983683296288046)]]
>>> neg_sum = rdd1.reduceByKey(lambda v1,v2: v1+v2)
>>> print neg_sum
PythonRDD[87] at RDD at PythonRDD.scala:43
>>> neg_sum.collect()
[('neg', 20.807473277488825)]
>>> █
(negative_sum)
```

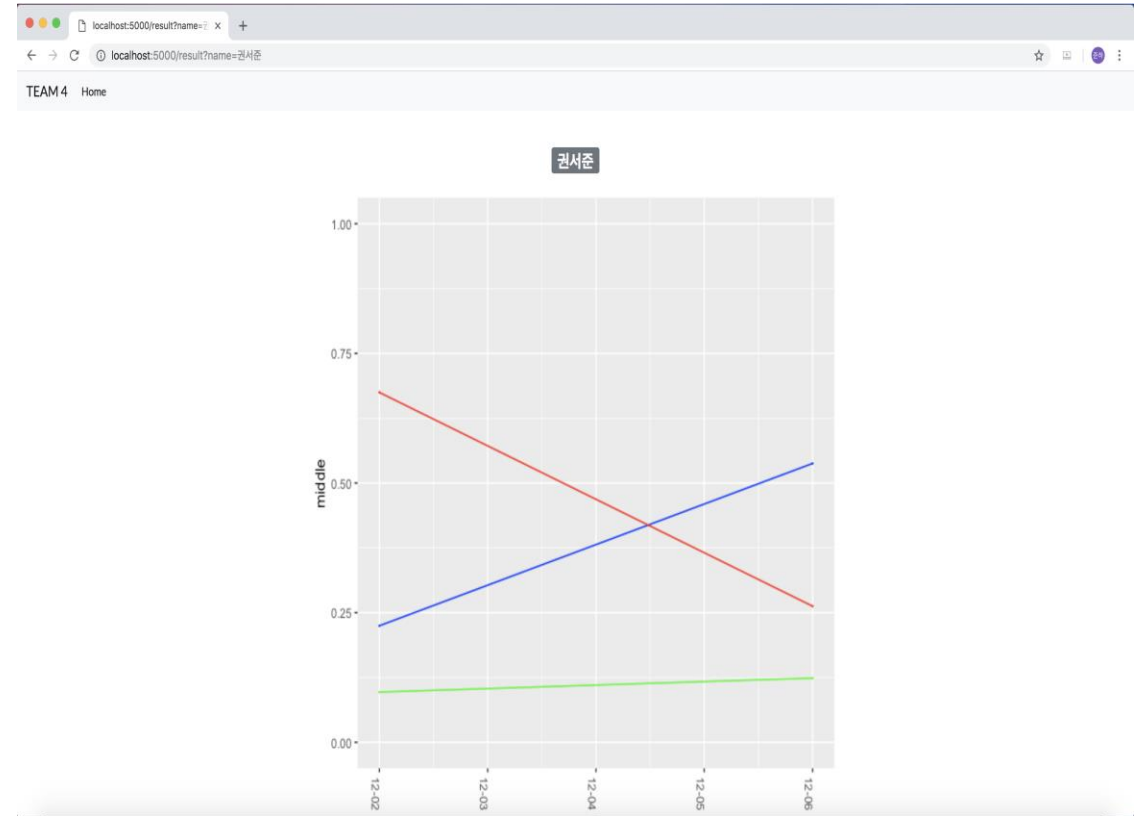
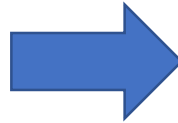
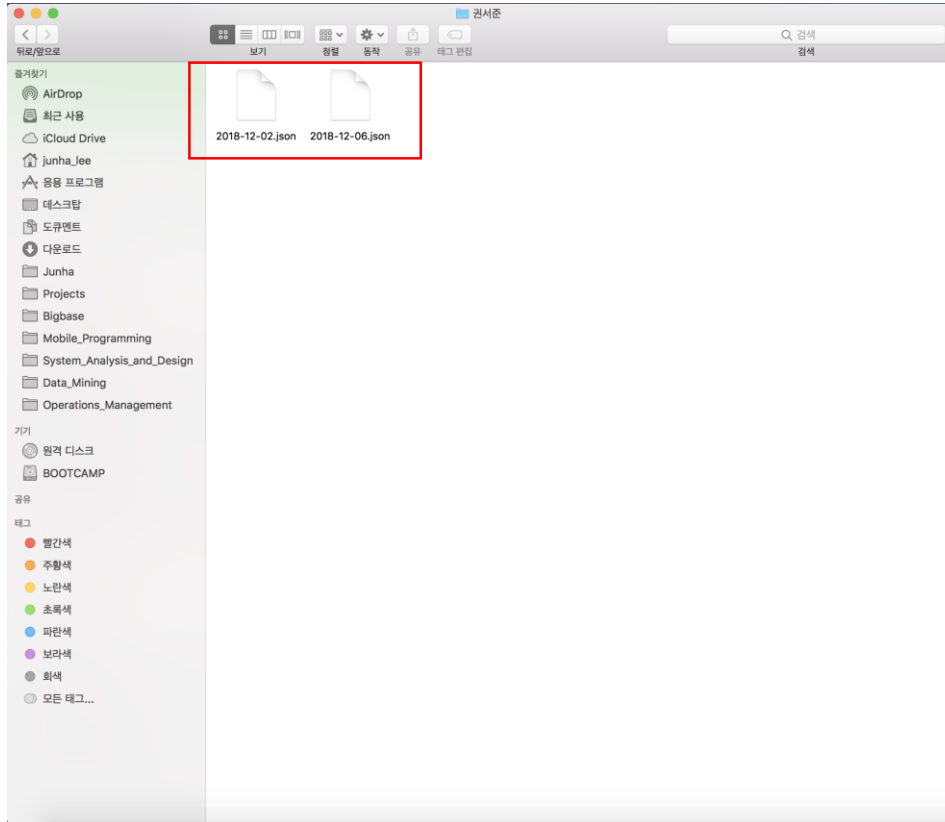
```
///
>>> rdd2 = rdd.map(lambda line:line.split(",")[2])\
...            .map(lambda line:line.split("="))\
...            .map(lambda x: (x[0], float(x[1])))
>>> rdd2.take(1)
[('neut', 0.039220113994397128)]
>>> neut_sum = rdd2.reduceByKey(lambda v1,v2: v1+v2)
>>> neut_sum.collect()
[('neut', 2.8746634063707814)]
>>> █
(neutral_sum)
```

```
>>> rdd3 = rdd.map(lambda line:line.split(",")[4])\
...            .map(lambda line:line.split("="))\
...            .map(lambda x: (x[0], x[1]))
>>> rdd4 = rdd3.map(lambda (v1,v2): (v1,float( v2.split(")") [0])) )
>>> pos_sum = rdd4.reduceByKey(lambda v1,v2: v1+v2)
>>> pos_sum.collect()
[('pos', 19.955320971382207)]
>>> █
(positive_sum)
```

### 3. The overall steps of testing

## Integration testing

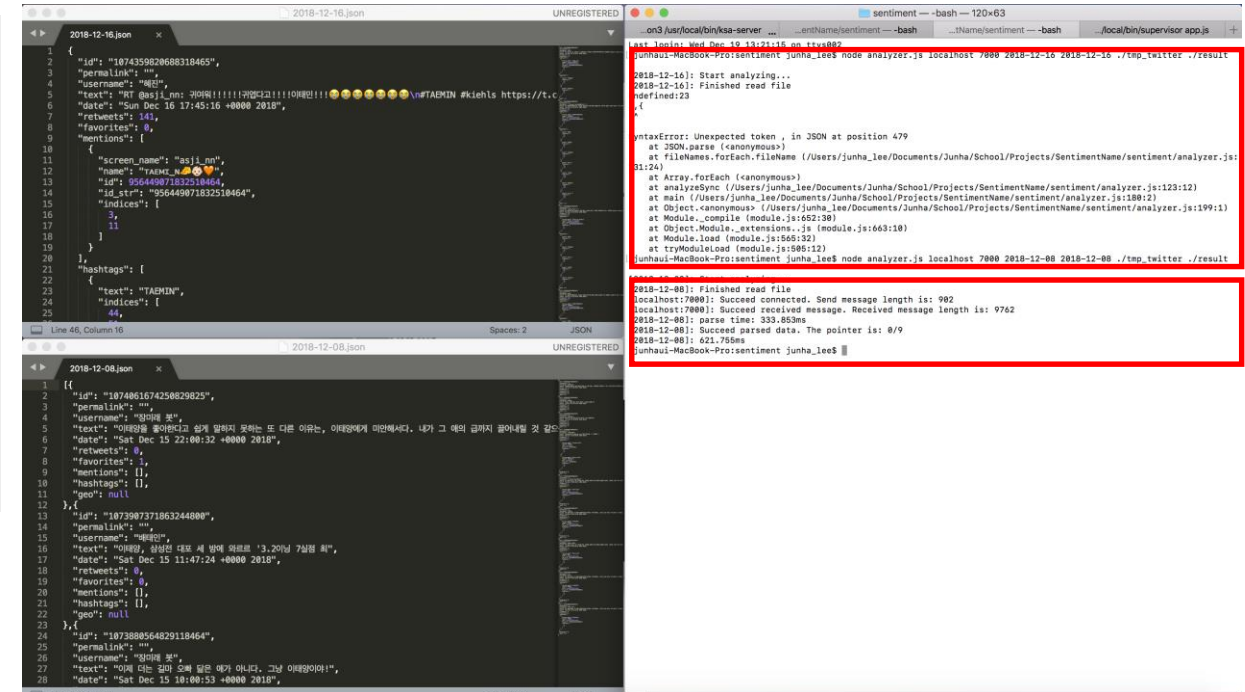
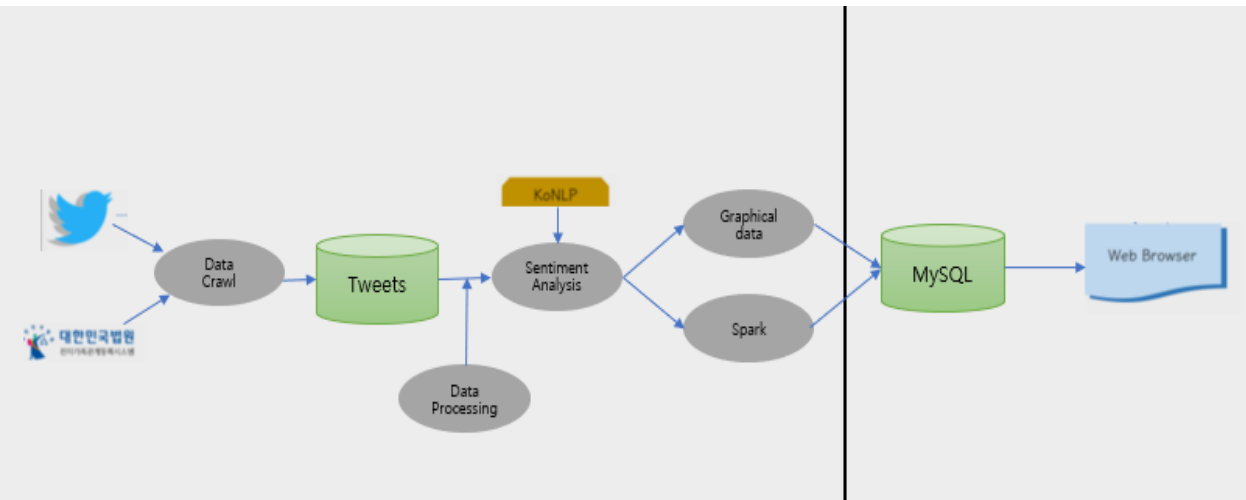
1) Test the case if the name which does not have enough data is inputted



→ Some of names don't have enough data to derive meaningful result

# Unit testing

2) Test whether the entire process is managed and implemented well as intended

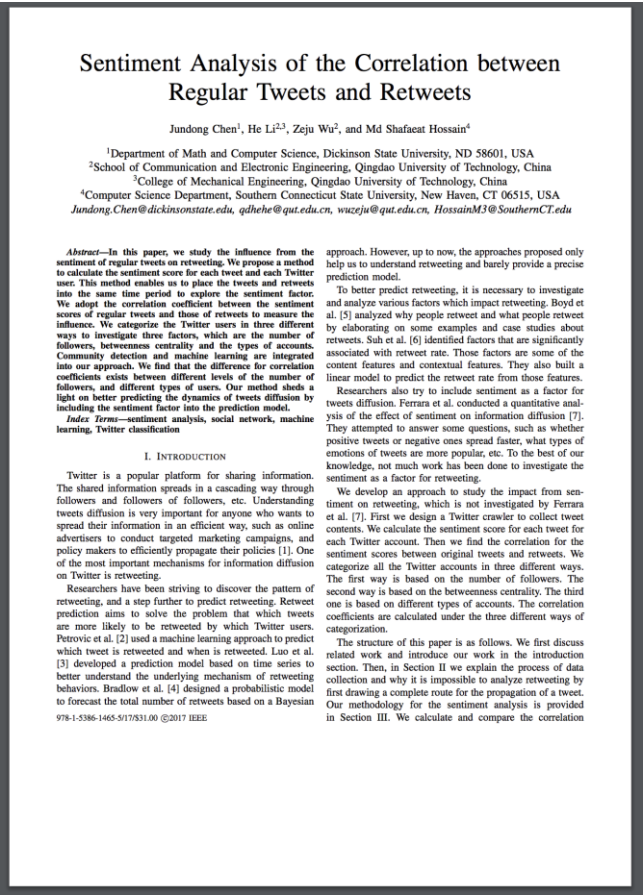


Although all the flow process is constructed, the entire functions should be checked several times to make sure that every steps are under the control and to avoid possible errors

3. The overall steps of testing

Unit testing

3) Check whether the retweets also have similar value with regular tweets in sentiment analysis



After checking the dataset, we find that only 8.4% of all the users have negative average sentiment scores for regular tweets, and 9.0% for retweets. That means most of the users have positive sentiment scores for their average one-day tweets.

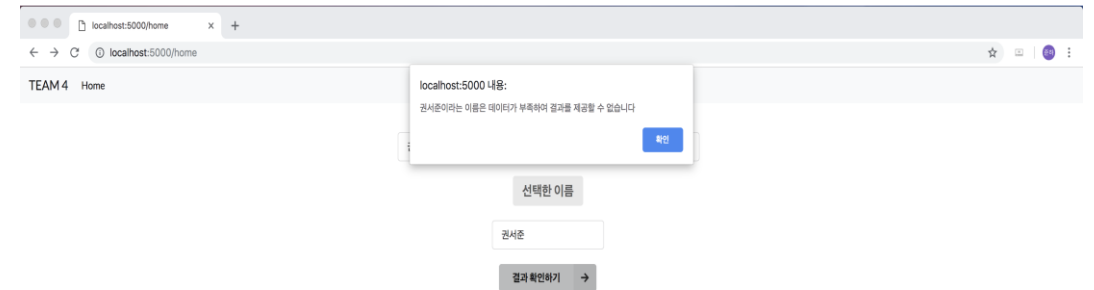
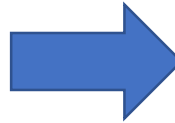
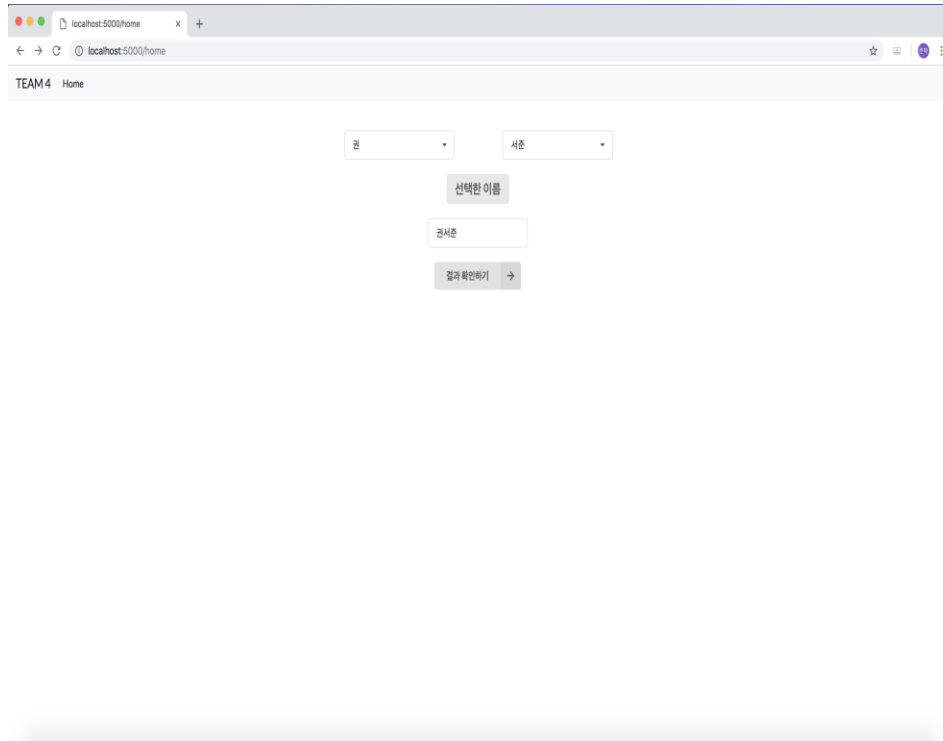
In this research, it is said that retweets(9.0%) show similar average sentiment scores with regular tweets(8.4%)

(Thesis about correlation between tweets-retweets)

## 4. The steps of analyzing the problems and debugging

### 4.1 Case 1 – The amount of saved data is not enough to show meaningful results

## Integration testing



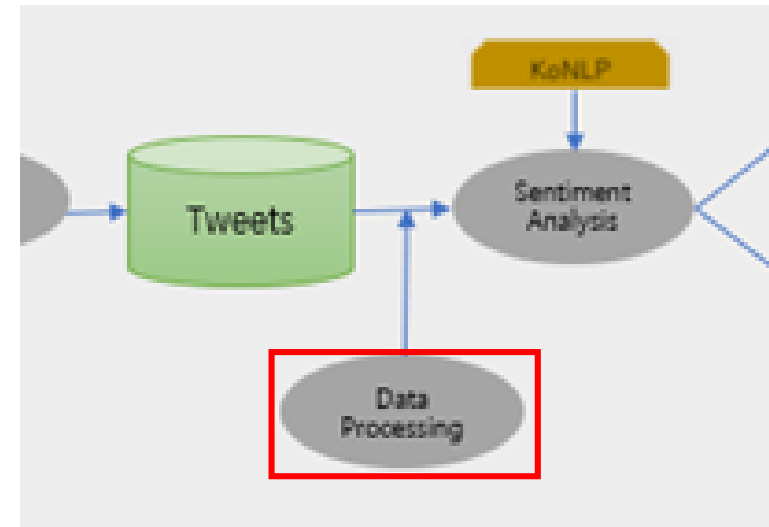
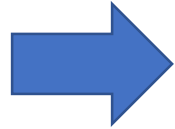
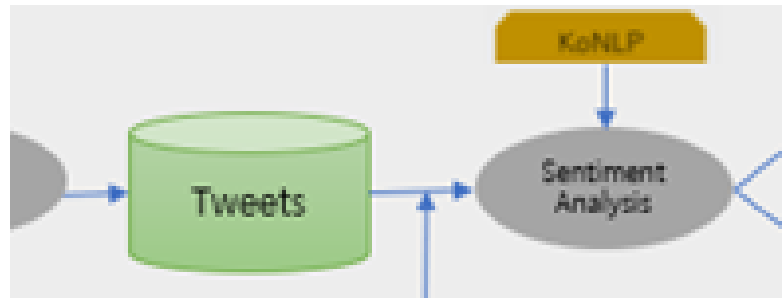
If the name with lack of data is inputted,  
web would show pop-up message to notify that the search is unavailable



## 4. The steps of analyzing the problems and debugging

## Unit testing

4.2 Case 2 – Cannot read the previous JSON format while implementing sentiment analysis



Couldn't read the JSON format to implement sentiment analysis

Add the steps for processing data before implementing sentiment analysis

## 4. The steps of analyzing the problems and debugging

## Unit testing

### 4.3 Case 3 – whether to accept retweets while crawling the tweets

```
30
31 #트윗 크롤링을 날짜 기준으로 저장합니다.
32 def get_tweets(keyword, num_limit):
33     i=0
34     for tweet in tweepy.Cursor(api.search, q=keyword, since='2018-11-01', until='2018-11-14', lang="ko").items(num_limit):
35         if (not tweet.retweeted) and ('RT @' not in tweet.text):
36             StatusObject = tweet._json
37             dict1 = {
38                 'id': StatusObject['id_str'],
```

While processing the result of analysis,  
we found that the retweets are also have enough  
value for representing the preference of keywords

```
30
31 #트윗 크롤링을 날짜 기준으로 저장합니다.
32 def get_tweets(keyword, num_limit):
33     i=0
34     for tweet in tweepy.Cursor(api.search, q=keyword, since='2018-11-01', until='2018-11-14', lang="ko").items(num_limit):
35         if (not tweet.retweeted) and ('RT @' not in tweet.text):
36             StatusObject = tweet._json
37             dict1 = {
38                 'id': StatusObject['id_str'],
```

Therefore, we deleted the line  
which exclude the retweets from crawled data





Thank you