# T-SA:
## Twitter keyword Search API based Tweet Analysis
## (트위터 키워드 검색 API기반 트윗 분석)

과      목    산학캡스톤디자인1(2019-1학기)

담 당 교 수   정현숙 교수님

팀      명    브이아이 (VI)

발   표   자   이 석 준

발 표 일 자   2019.04.11.

# T-SA: Contents

Twitter Keyword Search API based Tweet Analysis

# T-SA: Team Introduction

Twitter Keyword Search API based Tweet Analysis

| | |
|---|---|
| Name | Lee SeokJune |
| Student ID | 20165072 |
| Cell Phone | 010-4020-5717 |
| E-mail | op2se1@gmail.com |
| Major Lang | Java |
| GitHub | https://github.com/SeokJune |
| Part | - MariaDB 환경 구축 및 관리<br>- Hadoop(Map)구현<br>- 문서 작성 및 수정 |

| | |
|---|---|
| Name | Lee YunHyuck |
| Student ID | 20165062 |
| Cell Phone | 010-4220-5134 |
| E-mail | leeyh5134@naver.com |
| Major Lang | Python |
| GitHub | https://github.com/yunhyuck |
| Part | - Hadoop3 환경 구축<br>- Hadoop, DB 연동 구현<br>- Sqoop 환경 구축<br>- Hadoop(Reduce)구현 |

| | |
|---|---|
| Name | Bae InGyu |
| Student ID | 20165073 |
| Cell Phone | 010-4679-4968 |
| E-mail | happykkk789@naver.com |
| Major Lang | Python |
| GitHub | https://github.com/BaeInGyu |
| Part | - Python, DB 연동 구현<br>- Visualization 구현 |

| | |
|---|---|
| Name | Seo JaeIck |
| Student ID | 20144773 |
| Cell Phone | 010-2460-7617 |
| E-mail | nero8879@naver.com |
| Major Lang | Python |
| GitHub | https://github.com/nero8879 |
| Part | - Twitter API 구현<br>- Visualization 구현 |

# T-SA: Purpose of Development

Twitter Keyword Search API based Tweet Analysis

대한민국 지역 및 특정 기간에 사용된 키워드 트렌드 분석

특정 인물의 트윗 스타일 분석

# T-SA: Development Environment

Twitter Keyword Search API based Tweet Analysis



| ubuntu | python | MariaDB | hadoop | OpenJDK | eclipse |
|--------|--------|---------|--------|---------|---------|
| 18.04.2 LTS | 3.6.8 | 10.1.38 | 3.2.0 | 1.8.0_191 | 2019-03(4.11) |

ubuntu

**18.04.2 LTS**

Ubuntu is an _open source software operating system_ that runs from the desktop, to the cloud, to all your internet connected things.

Ubuntu Site:
- https://www.ubuntu.com/

Python features a *dynamic type system* and *automatic memory management*.
It supports multiple programming paradigms, including *object-oriented*, *functional* and *procedural*.

3.6.8

Python Stie:
- https://www.python.org/

민형기, 파이썬으로 데이터 주무르기, 2017.12.29, 비제이퍼블릭

파이썬으로 데이터 주무르기 저자의 블로그 중 파이썬 목록
- https://pinkwink.kr/category/Software/Python

MariaDB

10.1.38

MariaDB is an open source *relational database management system (RDBMS)*. Based on the same source code as MySQL, follow the *GPL v2 license*.

MariaDB Stie:
- https://mariadb.com/kb/ko/mariadb

Hadoop software library is a framework that allows for the _distributed processing of large data sets_ across clusters of computers using simple programming models.

Hadoop Site:
- https://hadoop.apache.org/

정재화, 시작하세요! 하둡 프로그래밍 빅데이터 분석을 위한 하둡 기초부터 YARN까지 [개정2판], 2016.05.13, 위키북스

hadoop

3.2.0

## Using Hadoop Ecosysytem



Sqoop
 - RDBMS에서 데이터 수집
 - 배치 처리 후 RDBMS에 데이터 저장

HDFS
 - 분산 데이터 저장

Map/Reduce
 - 분산 데이터 배치 처리

OpenJDK is a _free and open-source_ implementation of the Java Platform, Standard Edition. also produces the _virtual machine_, the _Java Class Library_, the _Java compiler_ and etc.

OpenJDK Site:
- https://openjdk.java.net/

**OpenJDK**

1.8.0_191

# T-SA: Development Environment

Eclipse is an *integrated development environment(IDE)* used in computer programming, is the *most widely used Java IDE*, may also be used to develop applications in other programming languages via various plug-ins

Eclipse Site:
- https://www.eclipse.org/

2019-03 (4.11)

## Twitter API

5.0

Twitter API furnish developer with *publish and analyze of Tweets*, *optimize ads*, and *create unique customer experiences*.

Twitter Developer Site:
- https://developer.twitter.com

Tweepy Site:
- http://www.tweepy.org

Twitter Analysis Site:
- http://tweetrend.com/
- https://foller.me/

## Importing Twitter API Key



Get Started 를 누르면 아래와 같이 Twitter API를 사용하기 위한 인증키 발급을 받을 수 있는 목록을 받을 수 있다.
Create an app 을 제외한 나머지는 유료 이므로 무료로 사용하기 위한 인증키를 발급 받는다.

## Importing Twitter API Key

**Tell us how this app will be used** (required)

This field is only visible to Twitter employees. Help us understand how your app will be used. What will it enable you and your customers to do?

> 이 앱의 사용방법은 사용자의 키워드를 분석하여, 해당 키워드에 대해 분석을 통해 얻을 수 있는 정보들에 대해 시각화 하는 것에 목적이 있 습니다.

⚠ **Must be 100 characters or longer**

Minimum characters: **100**

앱의 이름, 앱의 설명, 사용하는 주소, 앱의 사용 방법에 대한 필수적인 요소를 작성합니다.

## Importing Twitter API Key



개인 정보의 Apps를 통해 본인이 사전에 작성한 제목을 통해API 키가 발급이 된 것을 확인 할 수 있다.

# T-SA: Program Flowchart_W05
Twitter Keyword Search API based Tweet Analysis



TwitterAPI: To import data from Twitter

Python: Provides tweepy which is twitterAPI, Visualization of Data

MariaDB: Open source R-DBMS, Based on the same source as MySQL

Hadoop: Distributed storge and Processing of big data, Pseudo-distributed

Sqoop: For BigData Transfers between Hadoop and MariaDB

① Twitter API를 이용한 데이터 크롤링

② 크롤링 된 데이터를 MariaDB에 저장

③ Sqoop을 이용해 HDFS에 분산 저장 처리

④ Map/Reduce를 통한 분산 데이터 배치 처리

⑤ Sqoop을 이용해 MariaDB에 저장

⑥ 저장된 데이터를 Python에 로드 및
　시각화 라이브러리를 이용한 데이터 시각화

# T-SA: Development Schedule

Twitter Keyword Search API based Tweet Analysis

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Conference | | | | | | | | | | | | | | | |
| Create document and PPT | | | | | | | | | | | | | | | |
| Create class of Twitter API in Python | | | | | | | | | | | | | | | |
| Create control class of MariaDB in Python | | | | | | | | | | | | | | | |
| Tasking of HDFS | | | | | | | | | | | | | | | |
| Tasking of Map/Reduce | | | | | | | | | | | | | | | |
| Normalized data save to MariaDB | | | | | | | | | | | | | | | |
| Visualization through by Python | | | | | | | | | | | | | | | |
| Test and Modification of T-SA | | | | | | | | | | | | | | | |

# T-SA: Weekly Progress_W05

Twitter Keyword Search API based Tweet Analysis

| Lee SeokJune | Lee YunHyuck | Bae InGyu | Seo JaeIck |
|---|---|---|---|
| 문서 작성 및 수정<br><br>발표 준비 | Map/Reduce 구현<br><br>DB, Hadoop 연동 | Python, MariaDB 의<br>DML(Insert, RowCheck)<br>구현 및 전체 오류 수정 작업 | ? |

# Hadoop
# (+Sqoop)

## HDFS에 저장된 데이터 확인

```
vi@vi:~$ hdfs dfs -cat /user/yunhyuck/test2/part-m-00000 | head -10
100,P1,Computer,C413,A
100,P1,Computer,E412,A
200,P2,Electric,C123,B
300,P3,Computer,C312,A
300,P3,Computer,C324,C
300,P3,Computer,C413,A
400,P1,Computer,C312,A
400,P1,Computer,C324,A
400,P1,Computer,C413,B
400,P1,Computer,E412,C
```

## Yarn 실행

yarn jar /home/vi/hadoop/jar/Wordcount.jar KeywordCount /user/yunhyuck/test2/part-m-00000 output

```
2019-04-07 19:15:05,533 INFO mapreduce.Job: The url to track the job: http://vi:8088/proxy/application_1554631447670_0001/
2019-04-07 19:15:05,534 INFO mapreduce.Job: Running job: job_1554631447670_0001
2019-04-07 19:15:12,640 INFO mapreduce.Job: Job job_1554631447670_0001 running in uber mode : false
2019-04-07 19:15:12,642 INFO mapreduce.Job:   map 0% reduce 0%
2019-04-07 19:15:16,722 INFO mapreduce.Job:   map 100% reduce 0%
2019-04-07 19:15:22,769 INFO mapreduce.Job:   map 100% reduce 100%
2019-04-07 19:15:22,785 INFO mapreduce.Job: Job job_1554631447670_0001 completed successfully
2019-04-07 19:15:22,873 INFO mapreduce.Job: Counters: 54
        File System Counters
                FILE: Number of bytes read=536
                FILE: Number of bytes written=444121
                FILE: Number of read operations=0
                FILE: Number of large read operations=0
                FILE: Number of write operations=0
                HDFS: Number of bytes read=349
                HDFS: Number of bytes written=108
                HDFS: Number of read operations=8
                HDFS: Number of large read operations=0
                HDFS: Number of write operations=2
                HDFS: Number of bytes read erasure-coded=0
        Job Counters
                Launched map tasks=1
                Launched reduce tasks=1
                Data-local map tasks=1
                Total time spent by all maps in occupied slots (ms)=1996
                Total time spent by all reduces in occupied slots (ms)=2256
                Total time spent by all map tasks (ms)=1996
                Total time spent by all reduce tasks (ms)=2256
                Total vcore-milliseconds taken by all map tasks=1996
                Total vcore-milliseconds taken by all reduce tasks=2256
                Total megabyte-milliseconds taken by all map tasks=2043904
                Total megabyte-milliseconds taken by all reduce tasks=2310144
```

```
        Map-Reduce Framework
                Map input records=10
                Map output records=50
                Map output bytes=430
                Map output materialized bytes=536
                Input split bytes=119
                Combine input records=0
                Combine output records=0
                Reduce input groups=17
                Reduce shuffle bytes=536
                Reduce input records=50
                Reduce output records=17
                Spilled Records=100
                Shuffled Maps =1
                Failed Shuffles=0
                Merged Map outputs=1
                GC time elapsed (ms)=87
                CPU time spent (ms)=1060
                Physical memory (bytes) snapshot=501370880
                Virtual memory (bytes) snapshot=5316816896
                Total committed heap usage (bytes)=457703424
                Peak Map Physical memory (bytes)=281010176
                Peak Map Virtual memory (bytes)=2656210944
                Peak Reduce Physical memory (bytes)=220360704
                Peak Reduce Virtual memory (bytes)=2660605952
        Shuffle Errors
                BAD_ID=0
                CONNECTION=0
                IO_ERROR=0
                WRONG_LENGTH=0
                WRONG_MAP=0
                WRONG_REDUCE=0
        File Input Format Counters
                Bytes Read=230
        File Output Format Counters
                Bytes Written=108
```

## Yarn 상세 결과 확인

```
vi@vi:~$ hdfs dfs -cat output3/part-r-00000 | head -10
100        2
200        1
300        3
400        4
A          6
B          2
C          2
C123       1
C312       2
C324       2
```

## Sqoop Job 실행

sqoop export --connect jdbc:mysql://localhost/mysql --username root -P --table test --export-dir
hdfs://localhost:9000/user/vi/output3/part-r-00000 --columns a,b --input-fields-terminated-by "₩t"

```
2019-04-08 14:46:23,438 INFO mapreduce.Job:  map 0% reduce 0%
2019-04-08 14:46:34,607 INFO mapreduce.Job:  map 100% reduce 0%
2019-04-08 14:46:34,646 INFO mapreduce.Job: Job job_1554733044141_0004 completed s
uccessfully
2019-04-08 14:46:34,762 INFO mapreduce.Job: Counters: 33
        File System Counters
                FILE: Number of bytes read=0
                FILE: Number of bytes written=918424
                FILE: Number of read operations=0
                FILE: Number of large read operations=0
                FILE: Number of write operations=0
                HDFS: Number of bytes read=806
                HDFS: Number of bytes written=0
                HDFS: Number of read operations=16
                HDFS: Number of large read operations=0
                HDFS: Number of write operations=0
                HDFS: Number of bytes read erasure-coded=0
        Job Counters
                Launched map tasks=4
                Data-local map tasks=4
                Total time spent by all maps in occupied slots (ms)=34796
                Total time spent by all reduces in occupied slots (ms)=0
                Total time spent by all map tasks (ms)=34796
                Total vcore-milliseconds taken by all map tasks=34796
                Total megabyte-milliseconds taken by all map tasks=35631104
```

```
        Map-Reduce Framework
                Map input records=17
                Map output records=17
                Input split bytes=524
                Spilled Records=0
                Failed Shuffles=0
                Merged Map outputs=0
                GC time elapsed (ms)=1056
                CPU time spent (ms)=5470
                Physical memory (bytes) snapshot=986075136
                Virtual memory (bytes) snapshot=10641772544
                Total committed heap usage (bytes)=745537536
                Peak Map Physical memory (bytes)=247001088
                Peak Map Virtual memory (bytes)=2662076416
        File Input Format Counters
                Bytes Read=0
        File Output Format Counters
                Bytes Written=0
2019-04-08 14:46:34,767 INFO mapreduce.ExportJobBase: Transferred 806 bytes in 20.
8141 seconds (38.7237 bytes/sec)
2019-04-08 14:46:34,770 INFO mapreduce.ExportJobBase: Exported 17 records.
```

## MariaDB 결과 확인

```
MariaDB [mysql]> select *from test;
+----------+------+
| a        | b    |
+----------+------+
| B        |    2 |
| C        |    2 |
| C123     |    1 |
| C312     |    2 |
| C324     |    2 |
| Electric |    1 |
| P1       |    6 |
| P2       |    1 |
| P3       |    3 |
| C413     |    3 |
| Computer |    9 |
| E412     |    2 |
| 100      |    2 |
| 200      |    1 |
| 300      |    3 |
| 400      |    4 |
| A        |    6 |
+----------+------+
17 rows in set (0.00 sec)
```

# Python(DB)

## DBModule.getRowByCheck

```python
# 해당 테이블에 레코드 존재파악 함수·························
def getRowByCheck(self,table) :

    # 데이터가 없으면 false 있으면 True
    if bool(self.selectDB(table)) == True :

        return True

    else :

        return False
```
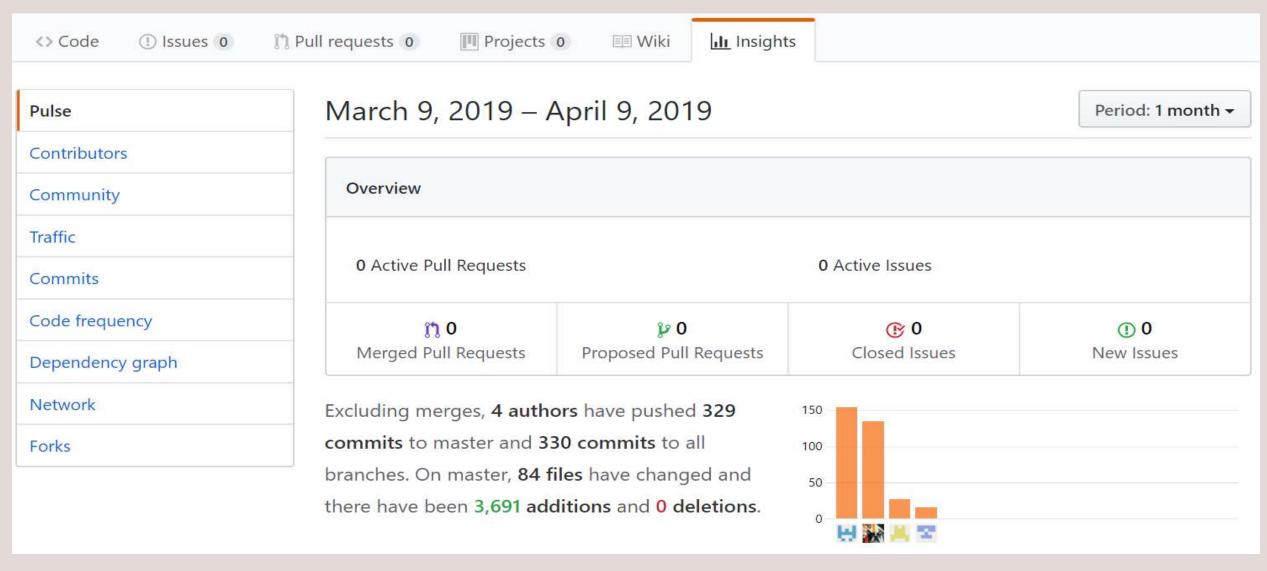
```
1. Keyword Search
2. User Search
3. Visualization
4. Exit
Choice Number:
1
작업할 테이블명 입력: Student
========테이블에 데이터가 존재 합니다.=========
========기존에 있는 데이터 삭제시작========
삭제완료
```

## DBModule.insertDB

```python
# 테이블의 데이터 삽입함수---------------------------------------------------------
def insertDB (self,table,values) :
        try :

                # MariaDB연결 및 Cursor생성
                conn, curs = self.dbConnect()

                # Data삽입
                self.values = values
                sql = "insert into "+table.strip()+values
                curs.execute(sql)
                conn.commit()

                print("삽입완료")

        except :
                print("삽입실패")

        finally :
                # Cursor종료 및 MariaDB연결종료
                self.dbClose()
```

# T-SA: Github_W05

Twitter Keyword Search API based Tweet Analysis

Project Github URL: https://github.com/SeokJune/BigData_VI_T-SA/

# Q & A

Thank you.