

PRINCIPLES OF BIG DATA MANAGEMENT
PROJECT – II
SPRING 2017

Submitted by: TEAM 7

Sri Sai Anusha Gandu (sgr43)

Sushma Mitta (smgp6)

Sadanand Kallakuri (sk789)

Abhilash Reddy Gaddam (aggg6)

CONTENTS

Title	Pg. No.
1. Introduction	1
1.1 About Twitter	1
1.2 About the Project	1
2. Requirements	1
2.1 Languages	1
2.2 Software	1
3. Tasks	2
4. Map Reduce	2
5. Extra Requirement	7
6. References	8

1. INTRODUCTION

1.1 About Twitter

Twitter is an online news and social networking service where users post and interact with messages, "tweets," restricted to 140 characters. Registered users can post tweets, but those who are unregistered can only read them. Twitter Inc. is based in San Francisco, California, United States, and has more than 25 offices around the world. Twitter was created in March 2006 by Jack Dorsey, Noah Glass, Biz Stone, and Evan Williams and launched in July, whereby the service rapidly gained worldwide popularity. As of 2016, Twitter had more than 319 million monthly active users.

1.2 About the Project

Here, we have collected the tweets using twitter API through tweepy using the keywords Python, JavaScript and Ruby in JSON (JavaScript Object Notation) format. The tweets then collected have been analyzed and different SQL queries are written to obtain the result.

2. REQUIREMENTS

2.1 Languages

1. Python
2. Scala
3. SQL
4. Java

2.2 Software

1. IntelliJ IDEA 3.4 (IDE)
2. Python 3.6
3. JDK 1.8
4. Scala 2.12.1
5. Spark 2.1
6. Virtual Box (Cloudera)

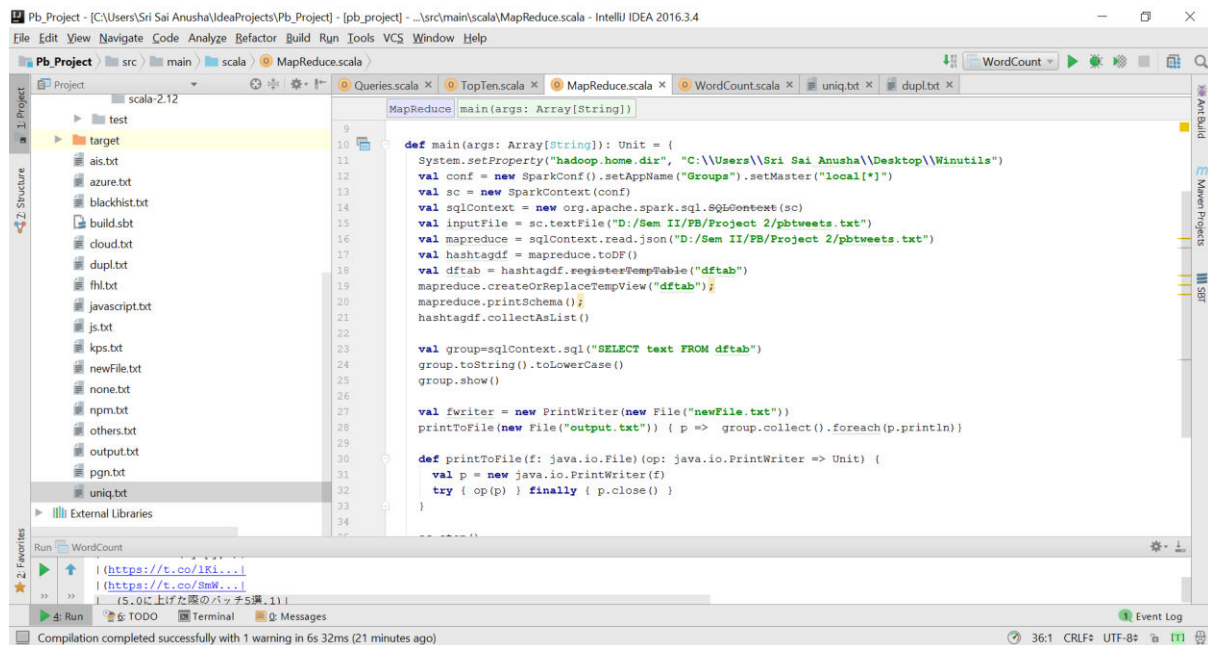
3. TASKS

1. To collect the tweets in JSON (JavaScript Object Notation) format.
2. Implement MapReduce to determine the uniqueness of the dataset.
3. To find the list of words in the tweets' text that are unique and duplicates.
4. Store the lists in two text files: uniq.txt and dupl.txt
5. To print the ratio of number of unique words to the number of duplicate words.

EXTRA REQUIREMENT

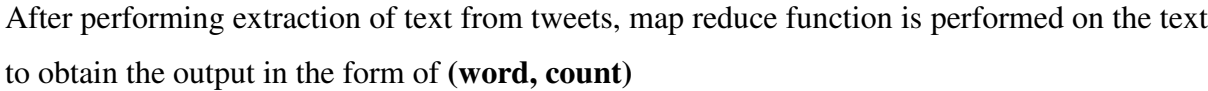
- To propose the best time to post a tweet on twitter.

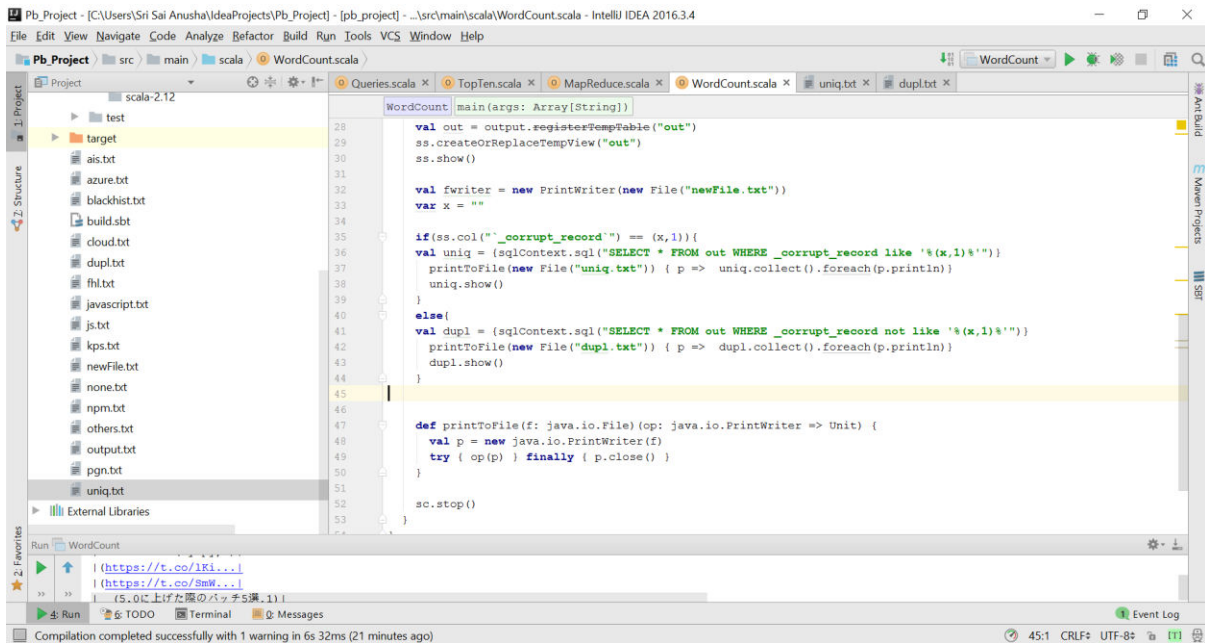
4. MAP REDUCE



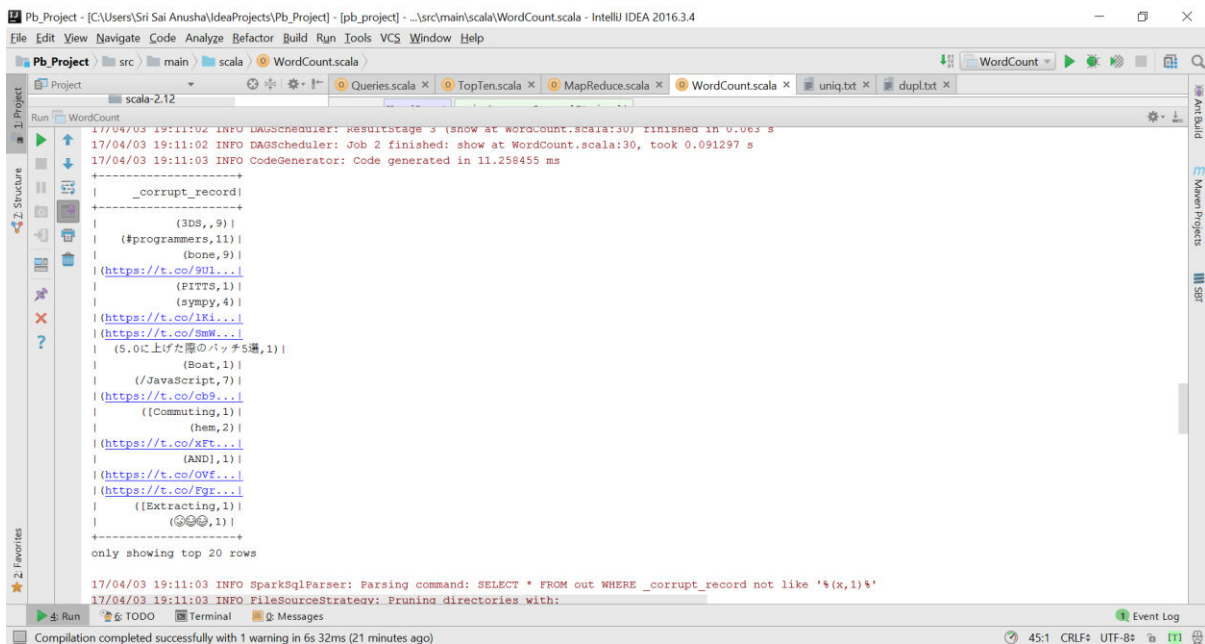
```
def main(args: Array[String]): Unit = {
  10
  11  System.setProperty("hadoop.home.dir", "C:\\Users\\Sri Sai Anusha\\Desktop\\Winutils")
  12  val conf = new SparkConf().setAppName("Groups").setMaster("local[*]")
  13  val sc = new SparkContext(conf)
  14  val sqlContext = new org.apache.spark.sql.SQLContext(sc)
  15  val inputFile = sc.textFile("D:/Sem II/PB/Project 2/pbtweets.txt")
  16  val mapreduce = sqlContext.read.json("D:/Sem II/PB/Project 2/pbtweets.txt")
  17  val hashtagdf = mapreduce.toDF()
  18  val dftab = hashtagdf.registerTempTable("dftab")
  19  mapreduce.createOrReplaceTempView("dftab")
  20  mapreduce.printSchema()
  21  hashtagdf.collectAsList()
  22
  23  val group=sqlContext.sql("SELECT text FROM dftab")
  24  group.toString().toLowerCase()
  25  group.show()
  26
  27  val fwriter = new PrintWriter(new File("newFile.txt"))
  28  printToFile(new File("output.txt")) { p => group.collect().foreach(p.println) }
  29
  30  def printToFile(f: java.io.File) (op: java.io.PrintWriter => Unit) {
  31    val p = new java.io.PrintWriter(f)
  32    try { op(p) } finally { p.close() }
  33  }
  34
  35  printToFile(new File("output.txt")) { p => group.collect().foreach(p.println) }
```

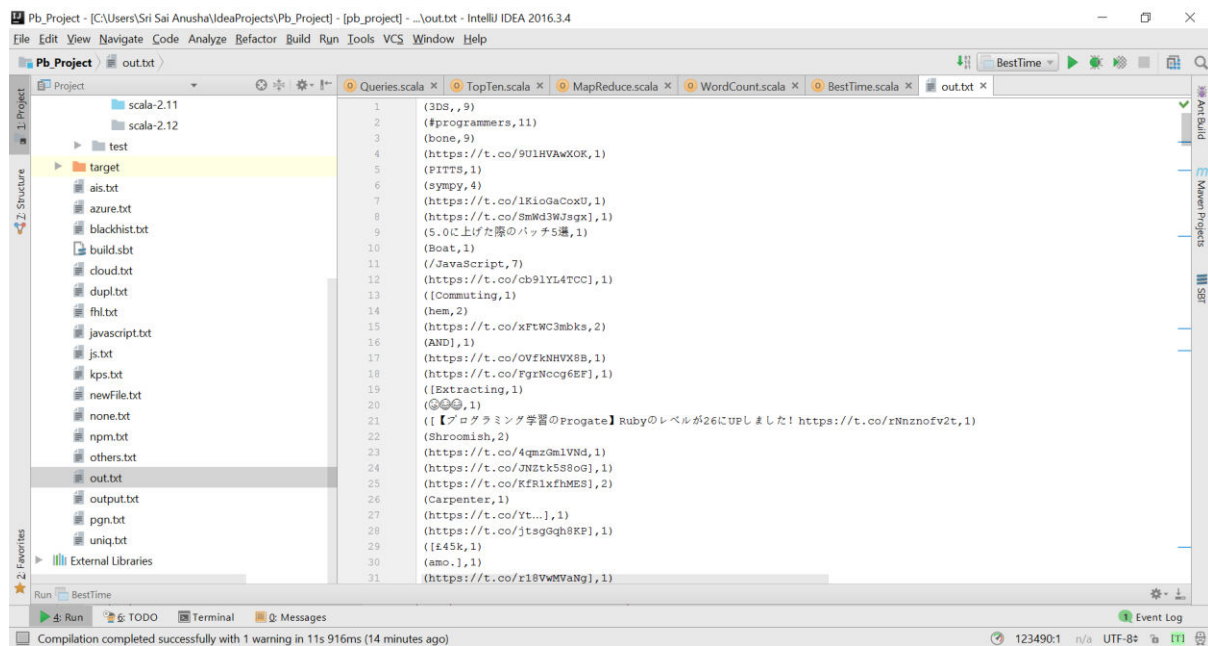
This bit of code extracts the text part from the tweets collected and saves it to a new text file named **output.txt**



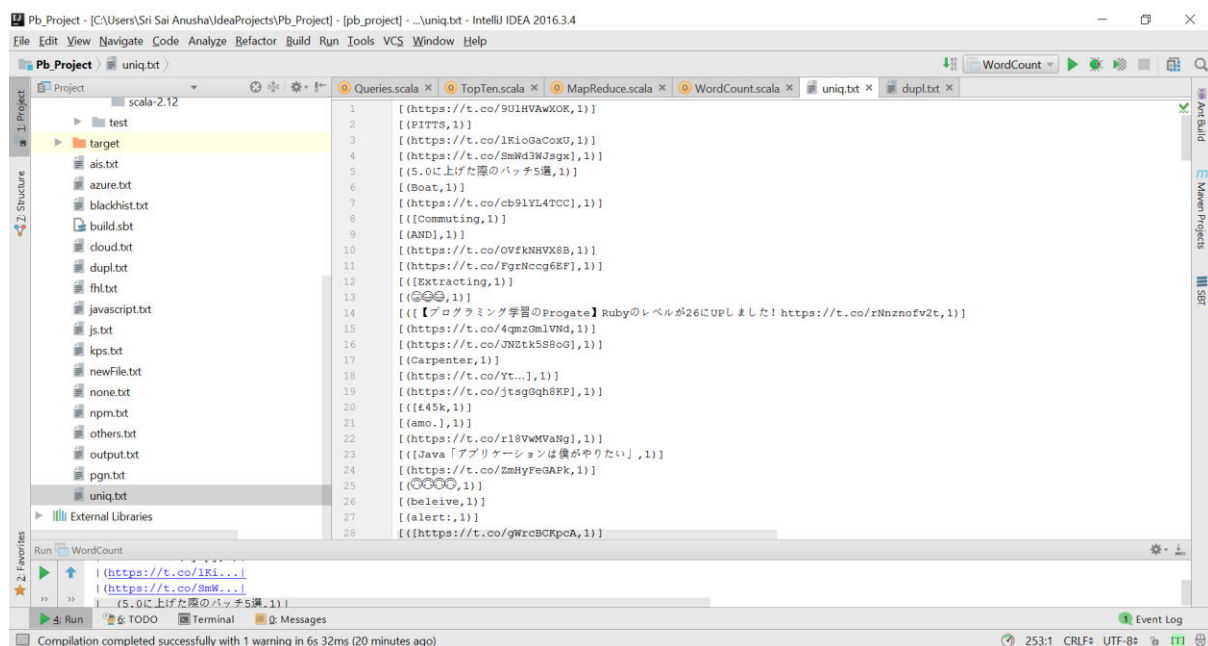


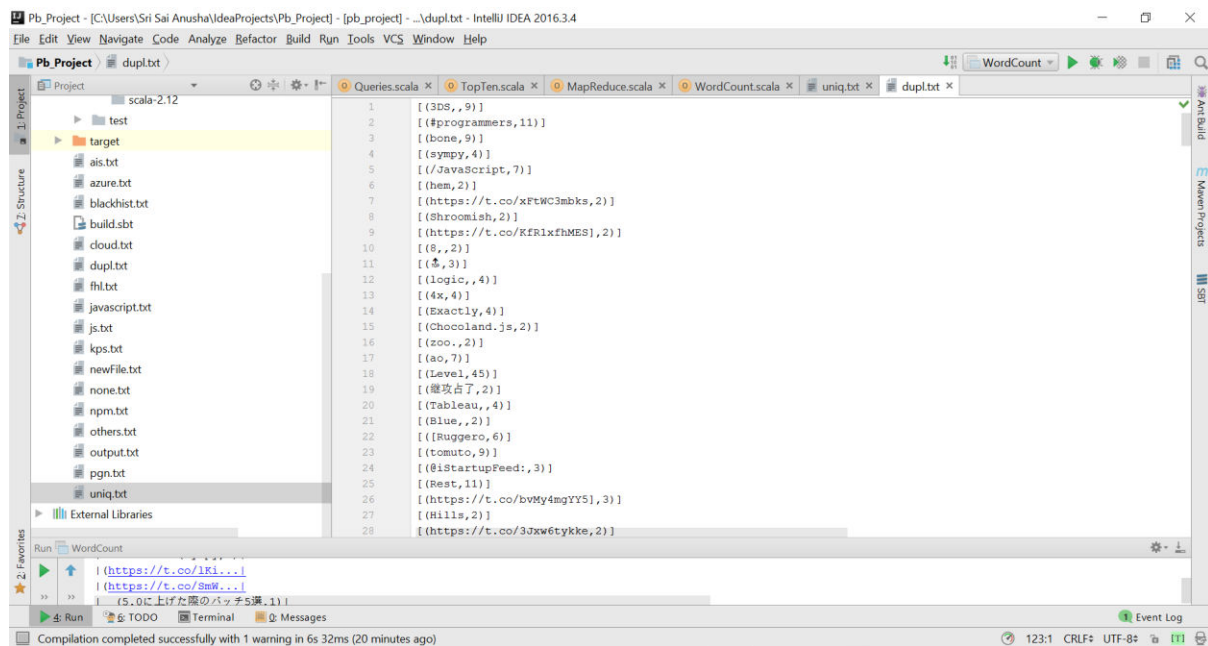
The output obtained through map reduce is shown below as a table and text file.



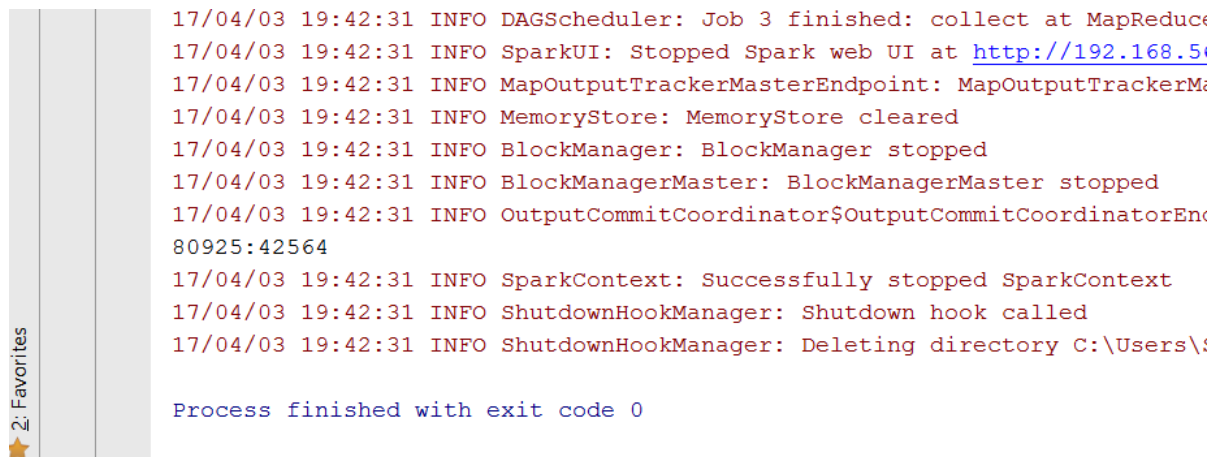


This output file is an aggregation of unique words and duplicate words in the text. These files are now segregated into uniq.txt and dupl.txt



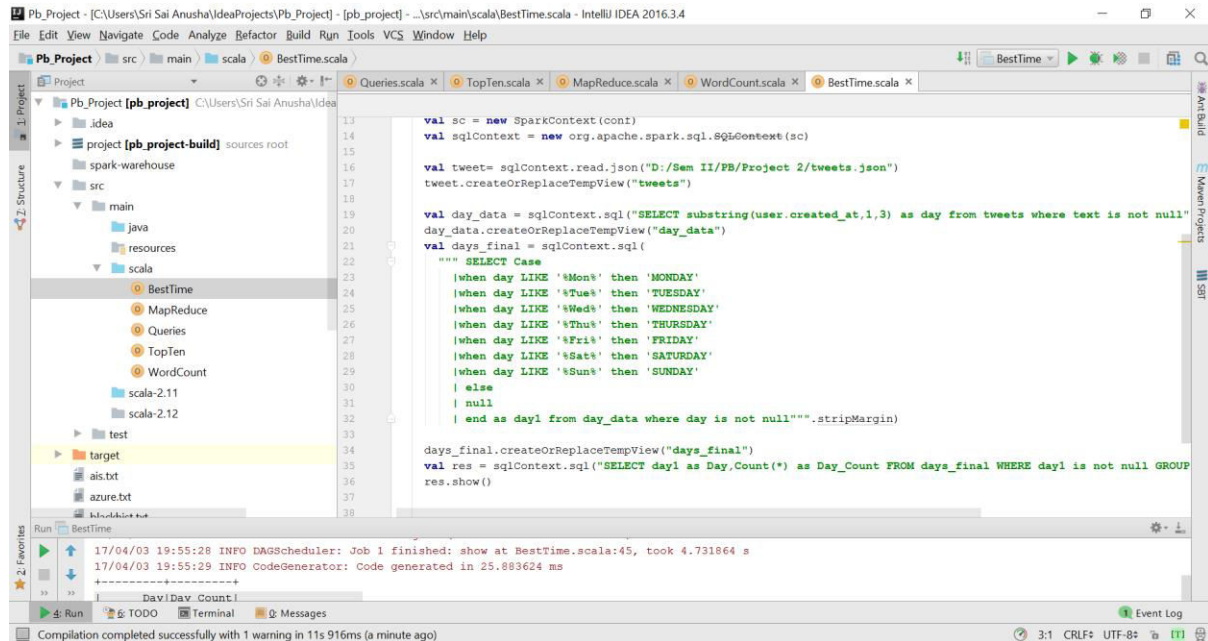


The ratio of the number of unique words to the number of duplicate words is shown below.



5. EXTRA REQUIREMENT

Based on the tweets JSON format, the metric to calculate the best time is the day of the week. According to the proposed metric, the output is as shown below.



```
13 val sc = new SparkContext(conf)
14 val sqlContext = new org.apache.spark.sql.SQLContext(sc)
15
16 val tweet = sqlContext.read.json("D:/Sem II/PB/Project 2/tweets.json")
17 tweet.createOrReplaceTempView("tweets")
18
19 val day_data = sqlContext.sql("SELECT substring(user.created_at,1,3) as day from tweets where text is not null")
20 day_data.createOrReplaceTempView("day_data")
21 val days_final = sqlContext.sql(
22     """ SELECT Case
23         |when day LIKE '%Mon%' then 'MONDAY'
24         |when day LIKE '%Tue%' then 'TUESDAY'
25         |when day LIKE '%Wed%' then 'WEDNESDAY'
26         |when day LIKE '%Thur%' then 'THURSDAY'
27         |when day LIKE '%Fri%' then 'FRIDAY'
28         |when day LIKE '%Sat%' then 'SATURDAY'
29         |when day LIKE '%Sun%' then 'SUNDAY'
30         | else
31         | null
32         | end as day1 from day_data where day is not null""").stripMargin)
33
34 days_final.createOrReplaceTempView("days_final")
35 val res = sqlContext.sql("SELECT day1 as Day,Count(*) as Day_Count FROM days_final WHERE day1 is not null GROUP BY day1")
36 res.show()
37
38
```

Run BestTime

17/04/03 19:55:28 INFO DAGScheduler: Job 1 finished: show at BestTime.scala:45, took 4.731864 s
17/04/03 19:55:29 INFO CodeGenerator: Code generated in 25.883624 ms

Day	Day_Count
MONDAY	9623
TUESDAY	8943
SUNDAY	8926
WEDNESDAY	8178
THURSDAY	7885
FRIDAY	7881
SATURDAY	7289

17/04/03 19:55:29 INFO SparkContext: Invoking stop() from shutdown hook

According to the output thus obtained, the best day of the week to tweet is MONDAY and the least number of tweets are observed on SATURDAY

6. REFERENCES

- <https://twitter.com/>
- <http://stackoverflow.com/questions/>
- <https://www.jetbrains.com/idea/>
- docs.scala-lang.org
- alvinalexander.com
- www.tutorialspoint.com