

ASSIGNMENT 21.1

Input Dataset:

Tweets – json file

```
tweets x
1
2 "filter_level": "low",
3 "retweeted": false,
4 "in_reply_to_screen_name": "FilmFan",
5 "truncated": false,
6 "lang": "en",
7 "in_reply_to_status_id_str": null,
8 "id": 689085590822891521,
9 "in_reply_to_user_id_str": "6048122",
10 "timestamp_ms": "1453125782100",
11 "in_reply_to_status_id": null,
12 "created_at": "Mon Jan 18 14:03:02 +0000 2016",
13 "favorite_count": 0,
14 "place": null,
15 "coordinates": null,
16 "text": "@FilmFan hey its time for you guys follow @acadgild To #AchieveMore and participate in contest Win Rs.500 worth vouchers",
17 "contributors": null,
18 "geo": null,
19 "entities": {
20   "symbols": [],
21   "urls": [],
22   "hashtags": [{
23     "text": "AchieveMore",
24     "indices": [56, 68]
25   }],
26   "user_mentions": [{
27     "id": 6048122,
28     "name": "Tanya",
29     "indices": [0, 8],
30     "screen_name": "FilmFan",
31     "id_str": "6048122"
32   }, {
33     "id": 2649945906,
34     "name": "ACADGILD",
35     "indices": [42, 51],
36     "screen_name": "acadgild",
37     "id_str": "2649945906"
38   }]
39 },
40 "is_quote_status": false,
41 "source": "<a href='\"https://about.twitter.com/products/tweetdeck\"' rel='\"nofollow\"'>TweetDeck</a>",
42 "favorited": false,
43 "in_reply_to_user_id": 6048122,
44 "retweet_count": 0,
45 "id_str": "689085590822891521",
46 "user": {
47   "location": "India ",
48   "default_profile": false,
49   "profile_background_tile": false,
50   "statuses_count": 86548,
51   "lang": "en",
52   "profile_link_color": "94D487",
53   "profile_banner_url": "https://pbs.twimg.com/profile_banners/197865769/1436198000",
54   "id": 197865769,
55   "following": null,
56   "protected": false,
57   "favourites_count": 1002,
58   "profile_text_color": "000000",
59   "verified": false,
60   "description": "Proud Indian, Digital Marketing Consultant,Traveler, Foodie, Adventurer, Data Architect, Movie Lover, Namo Fan",
61   "contributors_enabled": false,
62   "profile_sidebar_border_color": "000000",
63   "name": "Babubali",
64   "profile_background_color": "000000",
65   "created_at": "Sat Oct 02 17:41:02 +0000 2010",
66   "default_profile_image": false,
67   "followers_count": 4467,
68   "profile_image_url_https": "https://pbs.twimg.com/profile_images/664486535040000000/GOiDUiuK_normal.jpg",
69   "geo_enabled": true,
70   "profile_background_image_url": "http://abs.twimg.com/images/themes/theme1/bg.png",
```

```

71 "profile_background_image_url_https": "https://abs.twimg.com/images/themes/theme1/bg.png",
72 "follow_request_sent": null,
73 "url": null,
74 "utc_offset": 19800,
75 "time_zone": "Chennai",
76 "notifications": null,
77 "profile_use_background_image": false,
78 "friends_count": 810,
79 "profile_sidebar_fill_color": "000000",
80 "screen_name": "Ashok_Uppuluri",
81 "id_str": "197865769",
82 "profile_image_url": "http://pbs.twimg.com/profile_images/664486535040000000/GOiDUiuK_normal.jpg",
83 "listed_count": 50,
84 "is_translator": false
85 }
86

```

Here is the Spark code snippet to count popular hash tags using Spark SQL:

```

// import required Spark packages
import org.apache.spark.sql.SparkSession

object Assignment21_1 {

  def main(args: Array[String]): Unit = {

    // create a SparkSession object that can be used to create various contexts of Spark such as sqlContext
    val spark = SparkSession
      .builder()
      .config("spark.sql.warehouse.dir", "file:///c:/tmp/spark-warehouse")
      .master("local[*]")
      .getOrCreate()

    // read json file and convert the data into a temporary table – ‘tweets’

    val tweets = spark.read.json("file:///E:/Acadgild/Session 21/assignment 21.1/ tweets")
    .registerTempTable("tweets")

    // SQL query to fetch tweet id and hashtags present in tweets and convert it into a table – ‘hashtags’

    val hashtags = spark.sql("select id as id, entities.hashtags.text as words from tweets")
    .registerTempTable("hashtags")
  }
}

```

```
// SQL query to separate only the hashtags and store them into a temporary table – ‘hashtag_word’
val hashtag_word = spark.sql("select id as id, hashtag from hashtags LATERAL VIEW
explode(words) w as hashtag").registerTempTable("hashtag_word")
```

```
// SQL query to get hashtags and their count. Print result to the console
```

```
val popular_hashtags = spark.sql("select hashtag, count(hashtag) as cnt from hashtag_word
group by hashtag order by cnt desc").show

}

}
```

```
Assignment21_1.scala x
1 package org.spark_samples
2
3 import org.apache.spark.sql.SparkSession
4
5 object Assignment21_1 {
6
7   def main(args: Array[String]) {
8
9     val spark = SparkSession
10      .builder()
11      .config("spark.sql.warehouse.dir", "file:///c:/tmp/spark-warehouse")
12      .master("local[*]")
13      .getOrCreate()
14
15     val tweets = spark.read.json("file:///E:\\Acadgild\\Session 21\\assignment 21.1\\tweets").registerTempTable("tweets")
16
17     val hashtags = spark.sql("select id as id,entities.hashtags.text as words from tweets").registerTempTable("hashtags")
18
19     val hashtag_word = spark.sql("select id as id,hashtag from hashtags LATERAL VIEW explode(words) w as hashtag").registerTempTable("hashtag_word")
20
21     val popular_hashtags = spark.sql("select hashtag, count(hashtag) as cnt from hashtag_word group by hashtag order by cnt desc").show
22   }
23
24 }
```

Output:

```
+-----+---+
|  hashtag|cnt|
+-----+---+
|AchieveMore| 1|
+-----+---+
```

To conclude, there is only one hashtag in the given json file of tweets, ‘AchieveMore’ which has occurred once.