# Principles of Big Data Management

**Phase 2 Report**

Team :

Chandrasekhar Pentakota (16267781)

Kartheek Katta (16273764)

Jaya Prakash Ravella (16271007)

**Links :**

https://github.com/chanduhub/TwitterProject_CS5540

## Goal:

* To store tweets in Spark SQL and run queries on the data to analyse some interesting things.

* Make use of visualization tools to understand data better.

## Environment Setup:

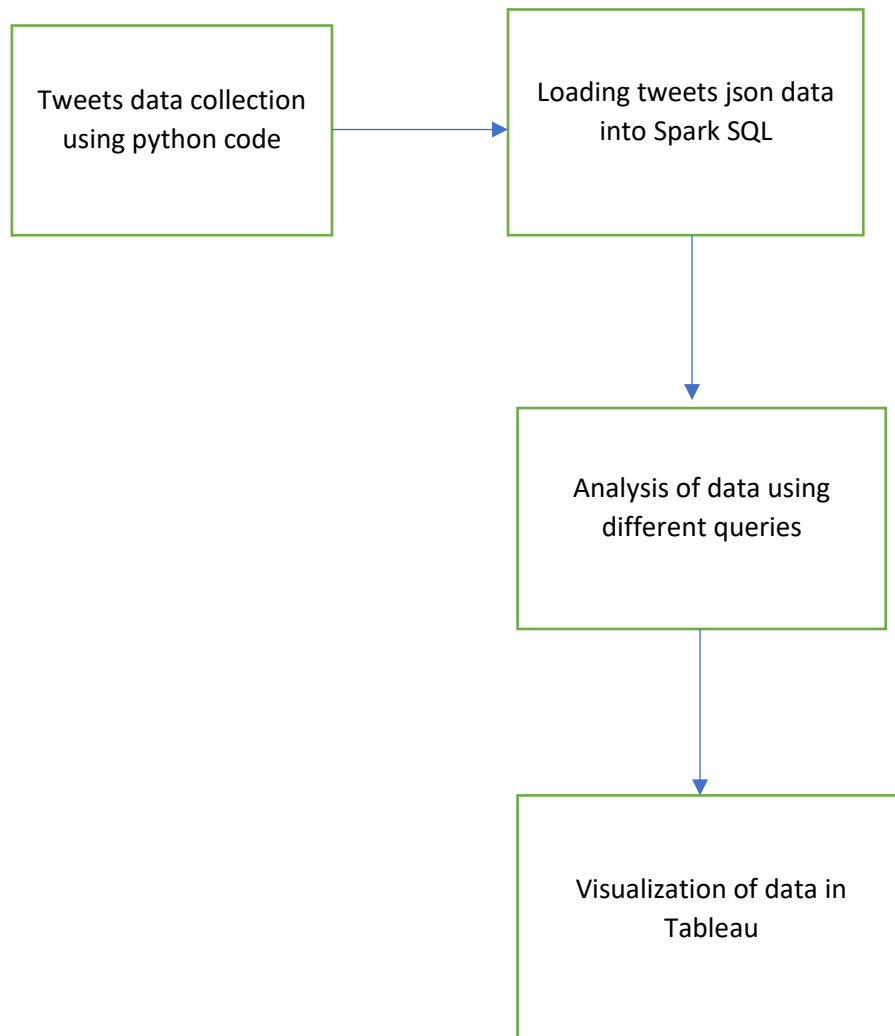We have installed spark in windows. The following steps are used

1.Download spark,unzip and set SPARK_HOME.

2.Download winutils jar and set it as HADOOP_HOME in environmental variables.

3.Install IntelliJ, the build.sbt file in intelliJ downloads all the required libraries for the project.

4.Install Tableau Public software for visualization.

```
build.sbt ×    sparkTwetts.scala ×
1      name := "PB Project"
2
3      version := "0.1"
4
5      scalaVersion := "2.11.8"
6
7      // https://mvnrepository.com/artifact/org.apache.spark/spark-core
8      libraryDependencies += "org.apache.spark" %% "spark-core" % "2.3.2"
9
10
11     // https://mvnrepository.com/artifact/org.apache.spark/spark-core
12     libraryDependencies += "org.apache.spark" %% "spark-sql" % "2.3.2"
13
```
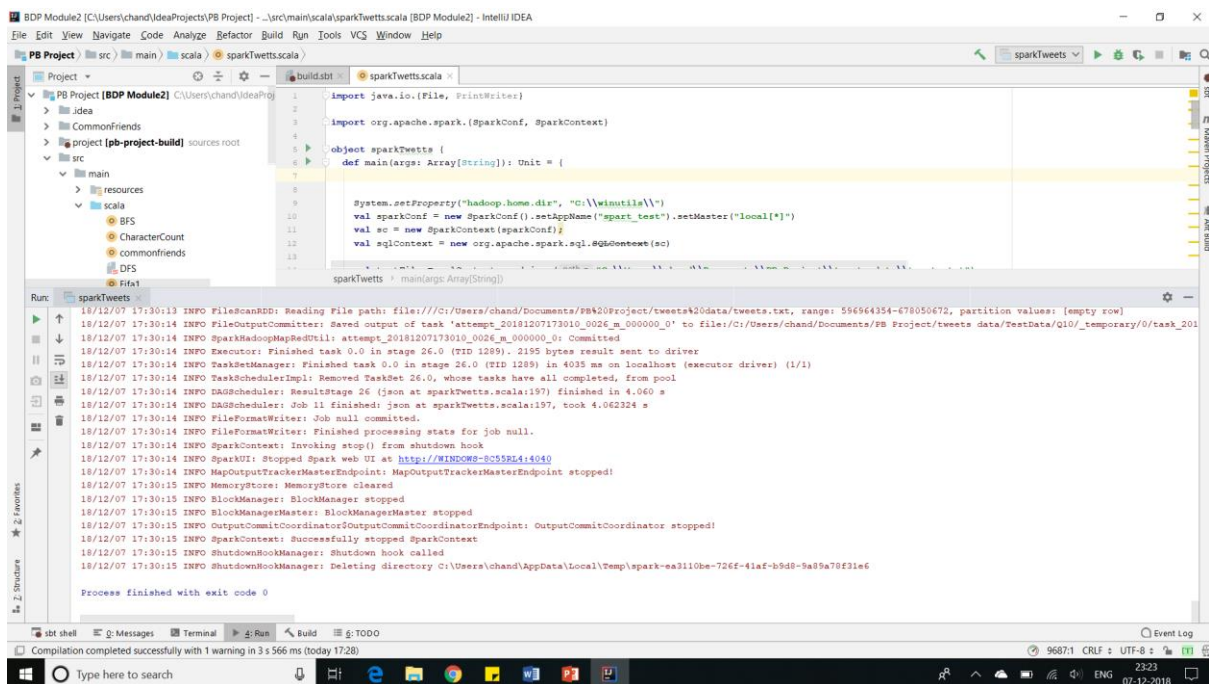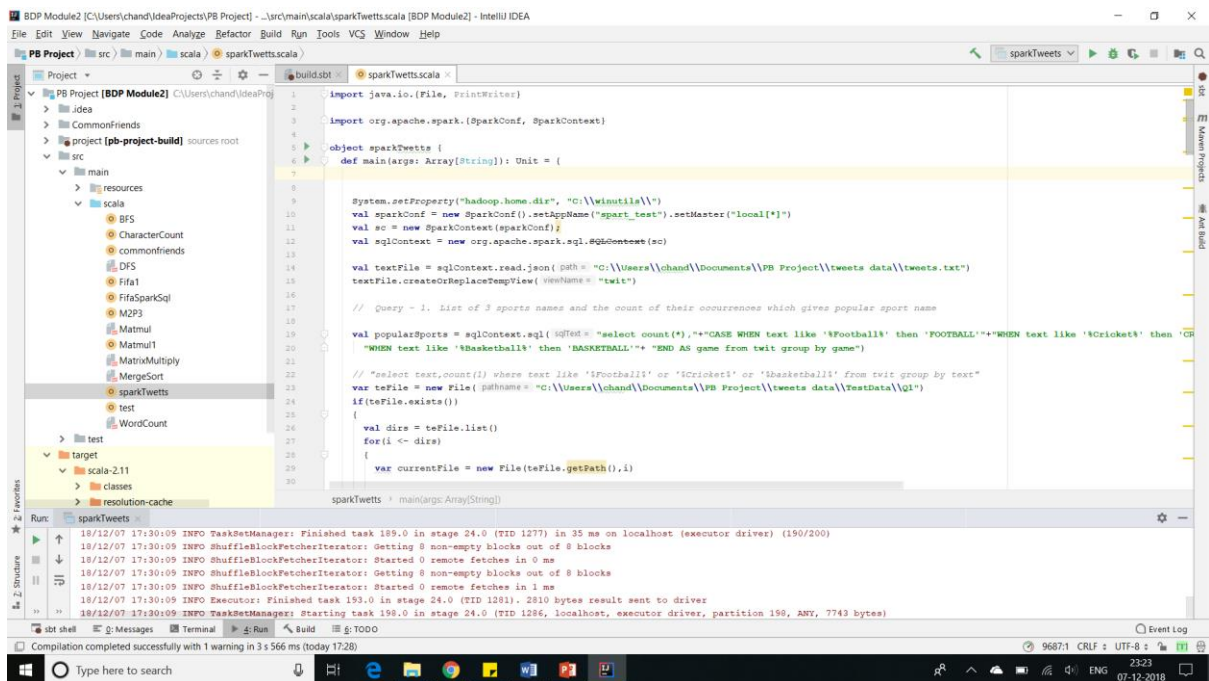
**Design:**

In phase1, we have collected twitter data by using the code which is of size 650 MB. The total tweets collected are 100K+. In this phase, we have done analysis on this data using spark SQL.

```
┌─────────────────────┐        ┌─────────────────────┐
│  Tweets data        │───────▶│  Loading tweets     │
│  collection         │        │  json data          │
│  using python code  │        │  into Spark SQL     │
└─────────────────────┘        └──────────┬──────────┘
                                          │
                                          ▼
                               ┌─────────────────────┐
                               │  Analysis of data   │
                               │  using              │
                               │  different queries  │
                               └──────────┬──────────┘
                                          │
                                          ▼
                               ┌─────────────────────┐
                               │  Visualization of   │
                               │  data in            │
                               │  Tableau            │
                               └─────────────────────┘
```

**Implementation:**

The main code is written in scala spark and the queries are written in Spark SQL. Output data generated from this code is fed into the Tableau to get insights from that data.
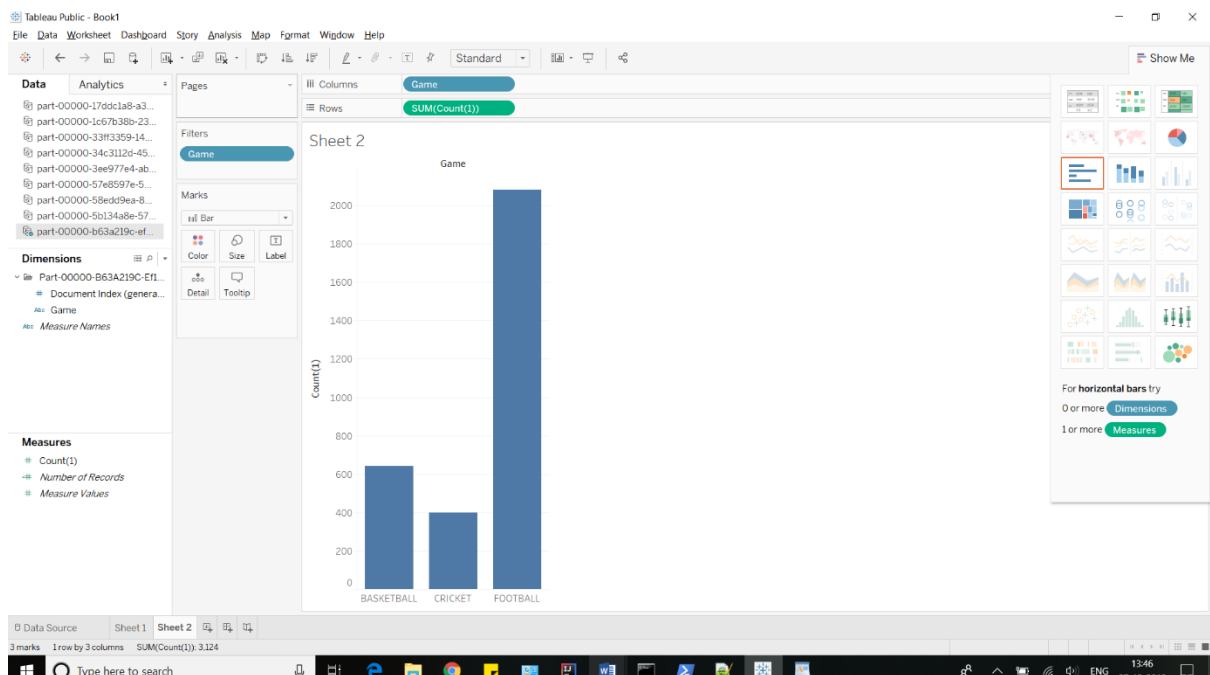
Process finished with the success code.

# Queries and output :

**Query1 :** List of 3 sports names and the count of their occurrences which gives popular sport name among them

```
//   Query - 1.  List of 3 sports names and the count of their occurrences which gives popular sport name

val popularSports = sqlContext.sql( sqlText = "select count(*),"+"CASE WHEN text like '%Football%' then 'FOOTBALL'"+"WHEN text like '%Cricket%' then 'CRICKET'"+
    "WHEN text like '%Basketball%' then 'BASKETBALL'"+ "END AS game from twit group by game")
```
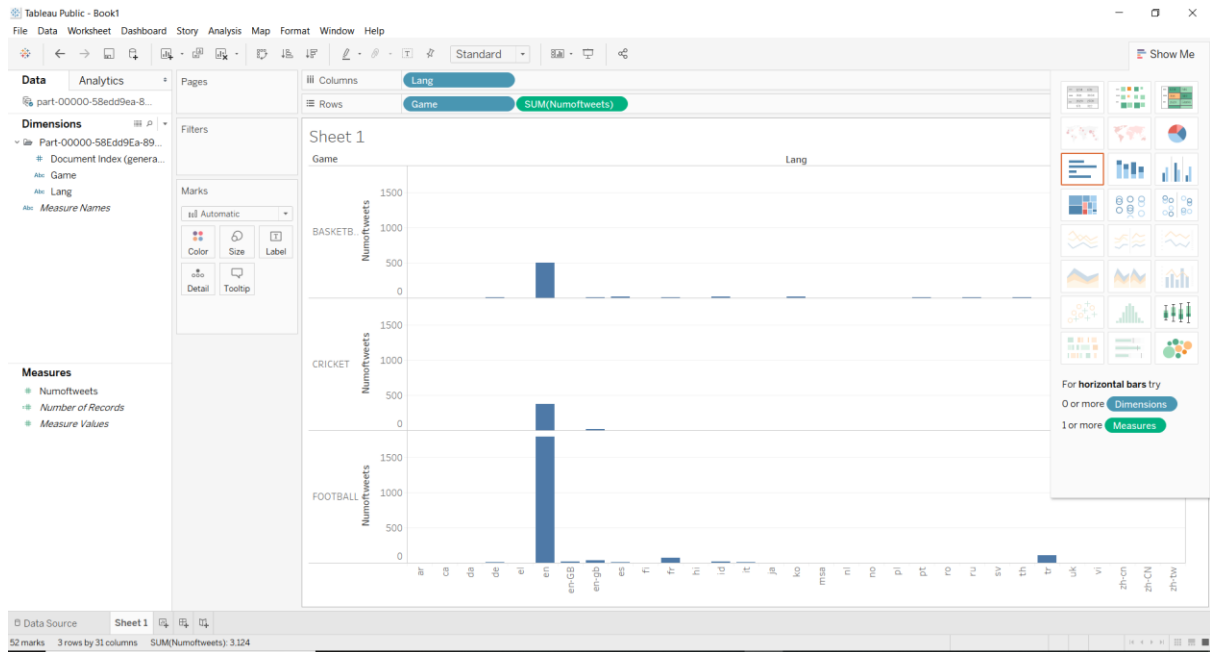


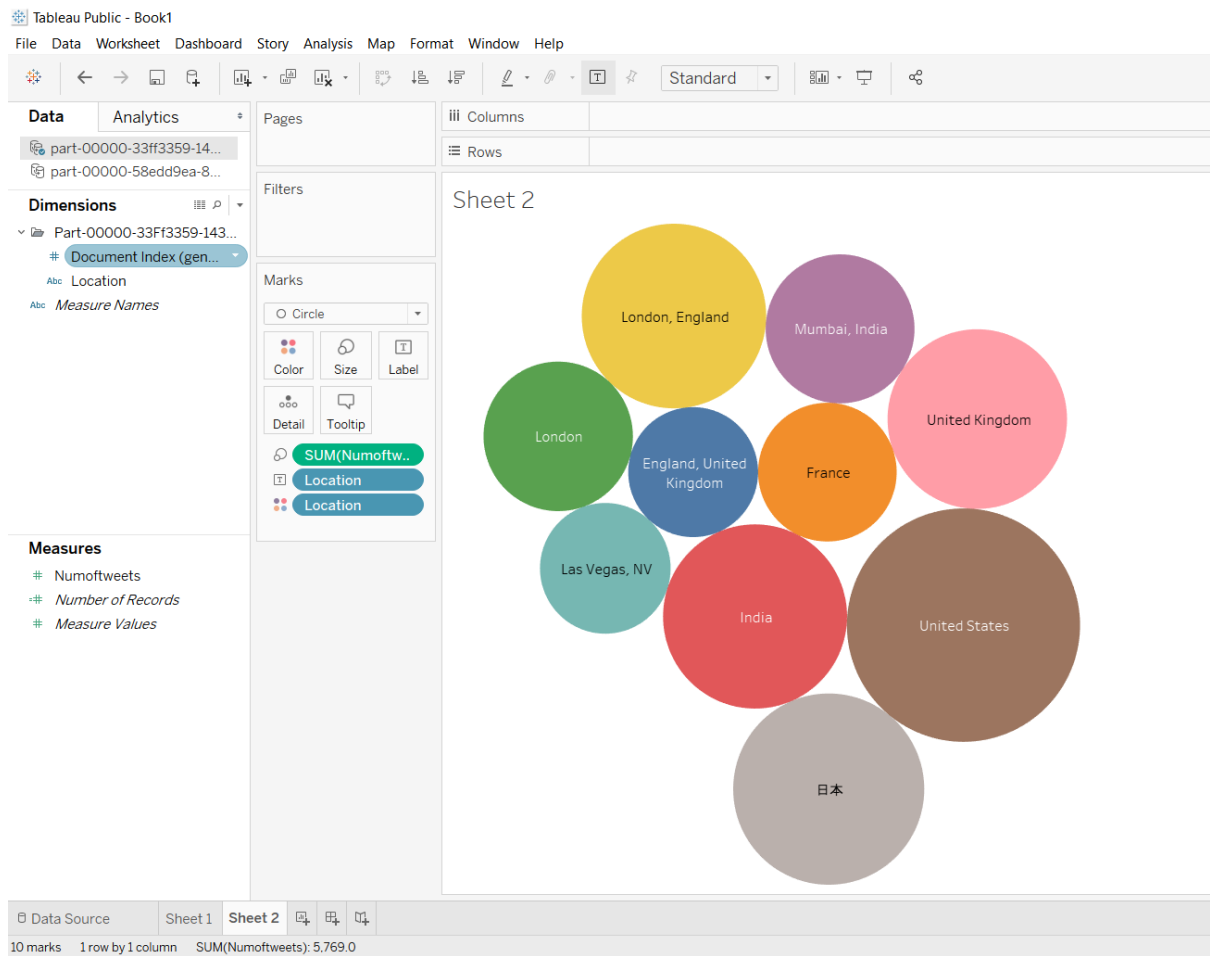**Query 2:** Favourite sports of people speaking different languages

```
// Query 2 - Favourite sports of people speaking different languages
val languague = sqlContext.sql( sqlText = "select user.lang,count(*) as numoftweets,"+"CASE WHEN text like '%Football%' then 'FOOTBALL'"+"WHEN text like '%Cricket%' then 'CRICKET'"+
    "WHEN text like '%Basketball%' then 'BASKETBALL' ELSE 'NogivenSport'"+ "END AS game from twit where text is not null group by game,user.lang order by game")
```

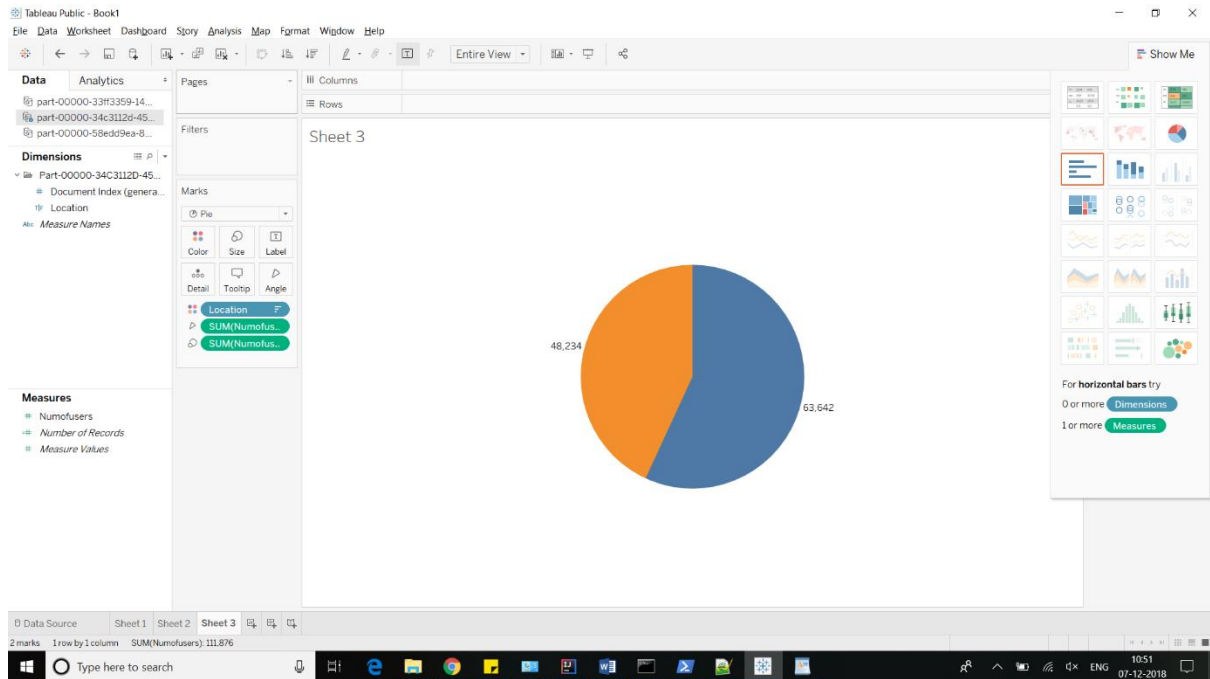## Query 3: More tweets coming from locations

```
// Query 3 - More tweets coming from locations
val moretweetLoc = sqlContext.sql( sqlText = "select user.location as location,count(*) as numoftweets  from twit where user.location is not null " +
  "group by user.location " + "order by count(1) desc limit 10")
```

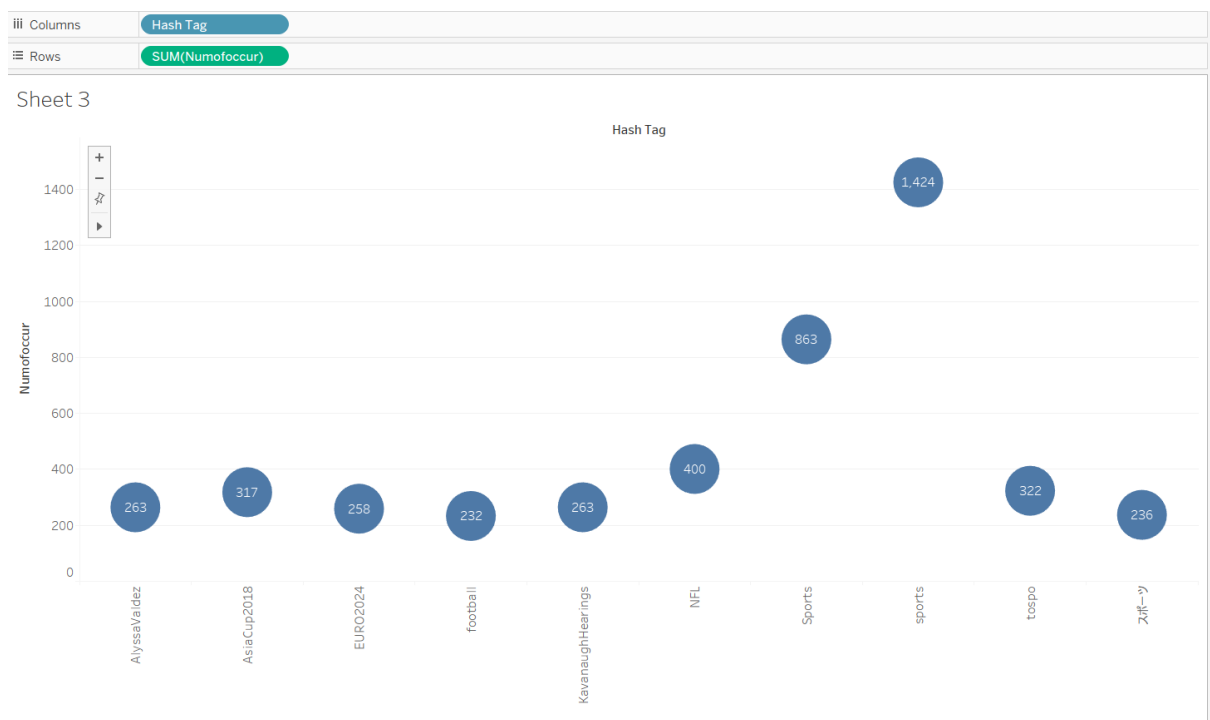## Query 4: Users who enabled geolocation

```
//Query 4- Users who enabled geolocation

val geoEnabled = sqlContext.sql( sqlText = "select user.geo_enabled as location,count(*) as numofusers  from twit " +
    "where user.geo_enabled is not null group by user.geo_enabled")
```

## Query5: Top Hashtags in sports tweets

```
// Query 5- Top hashtags
val topHashtags = sqlContext.sql( sqlText = "select text as hashTag,count(*) as numofoccur from hashtags " +
    "group by text order by count(1) desc limit 10")
```
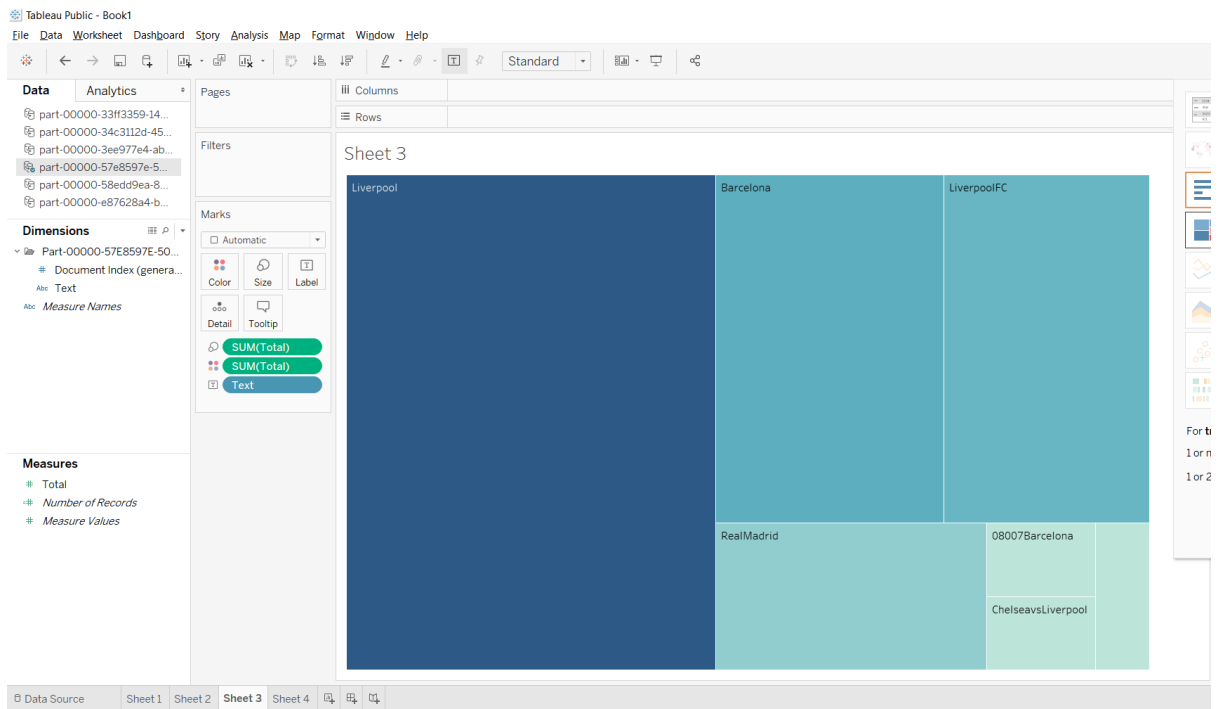
## Query 6: Football related, which is famous among Liverpool, RealMadrid and Barcelona

```
// Query 6 - Football related, which is famous among Liverpool, RealMadrid and Barcelona

val football = sqlContext.sql( sqlText = "select text,count(*) as total from hashtags where text like '%Liverpool%' or text like '%RealMadrid%' or text like '%Barcelona%' " +
    "group by text order by count(1) desc limit 10")
```
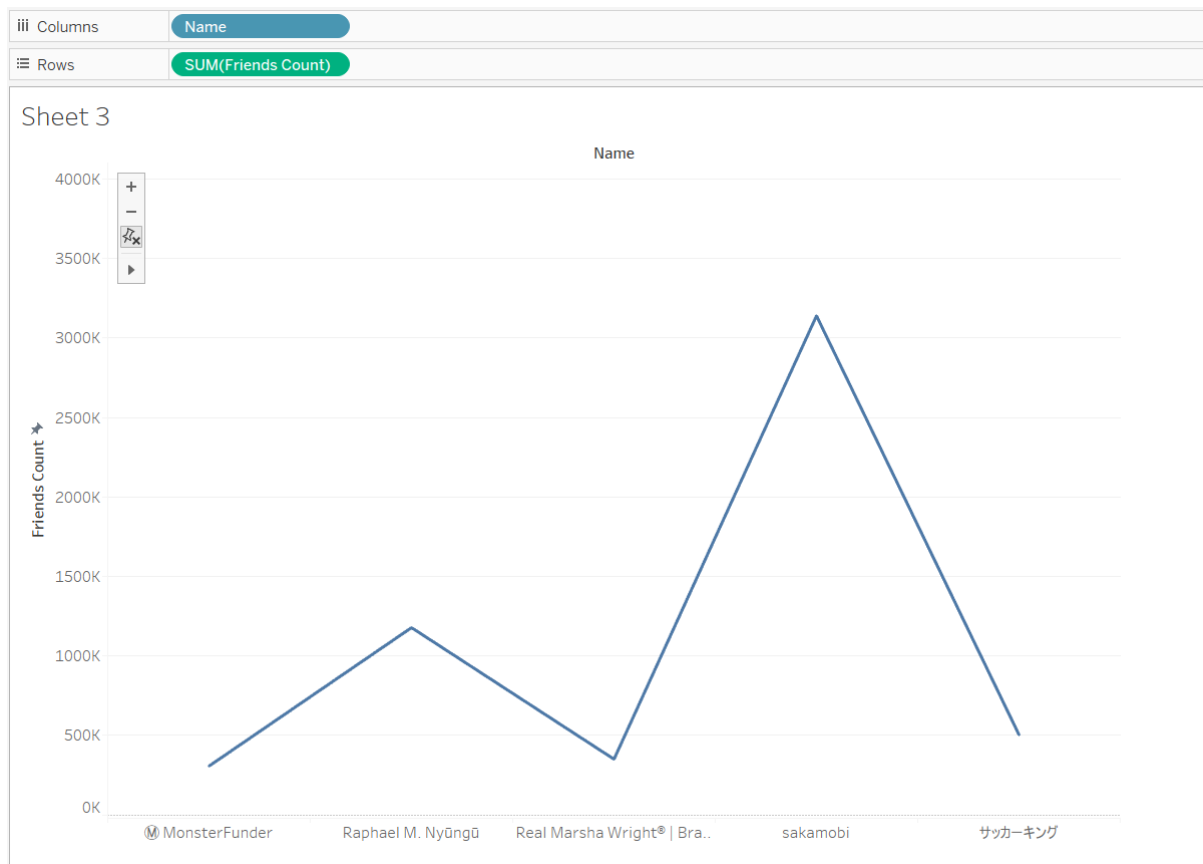


## Query7: Users who got more friends who are tweeting about sports

```
// Query 7 - Top Friends
val topfriends= sqlContext.sql( sqlText = "select user.name,user.friends_count  from twit order by user.friends_count desc LIMIT 10")
```
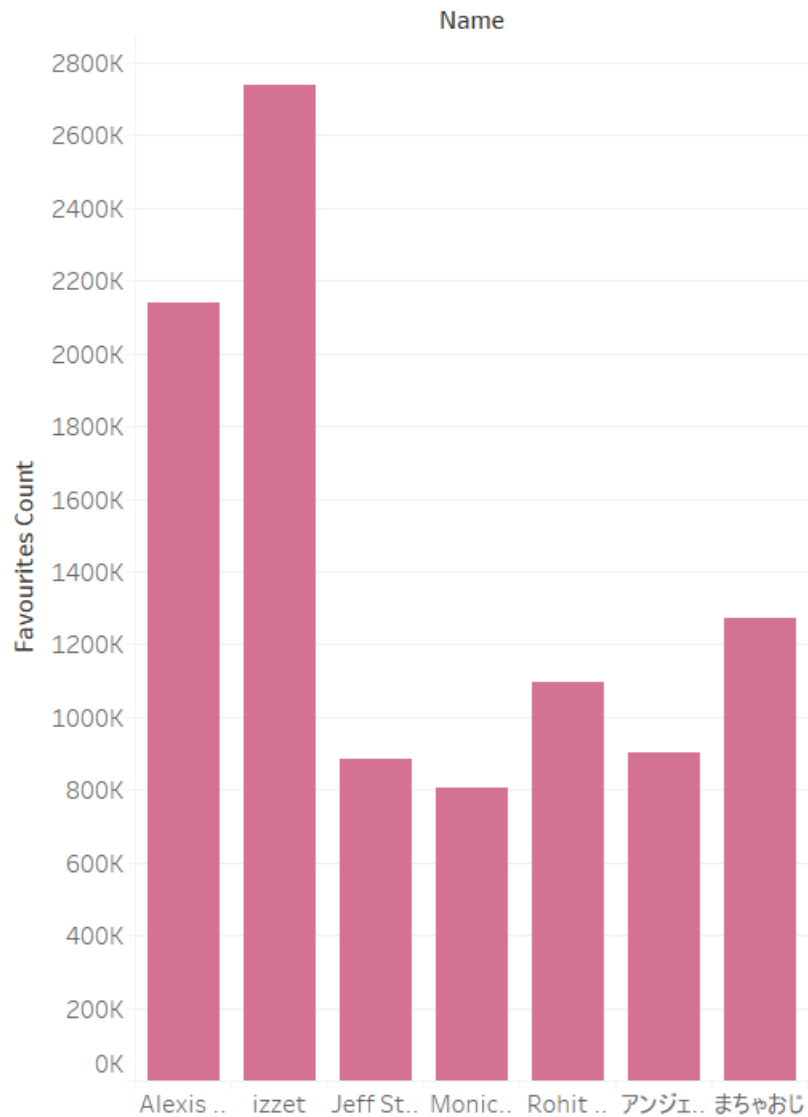
Sheet 3



**Query 8:** Top 10 Users who are actively liking tweets

```
// Query 8 - Top 10 Users Who are actively liking tweets.
val likingUsers = sqlContext.sql( sqlText = "select user.name,user.favourites_count from twit order by user.favourites_count desc limit 10")
```
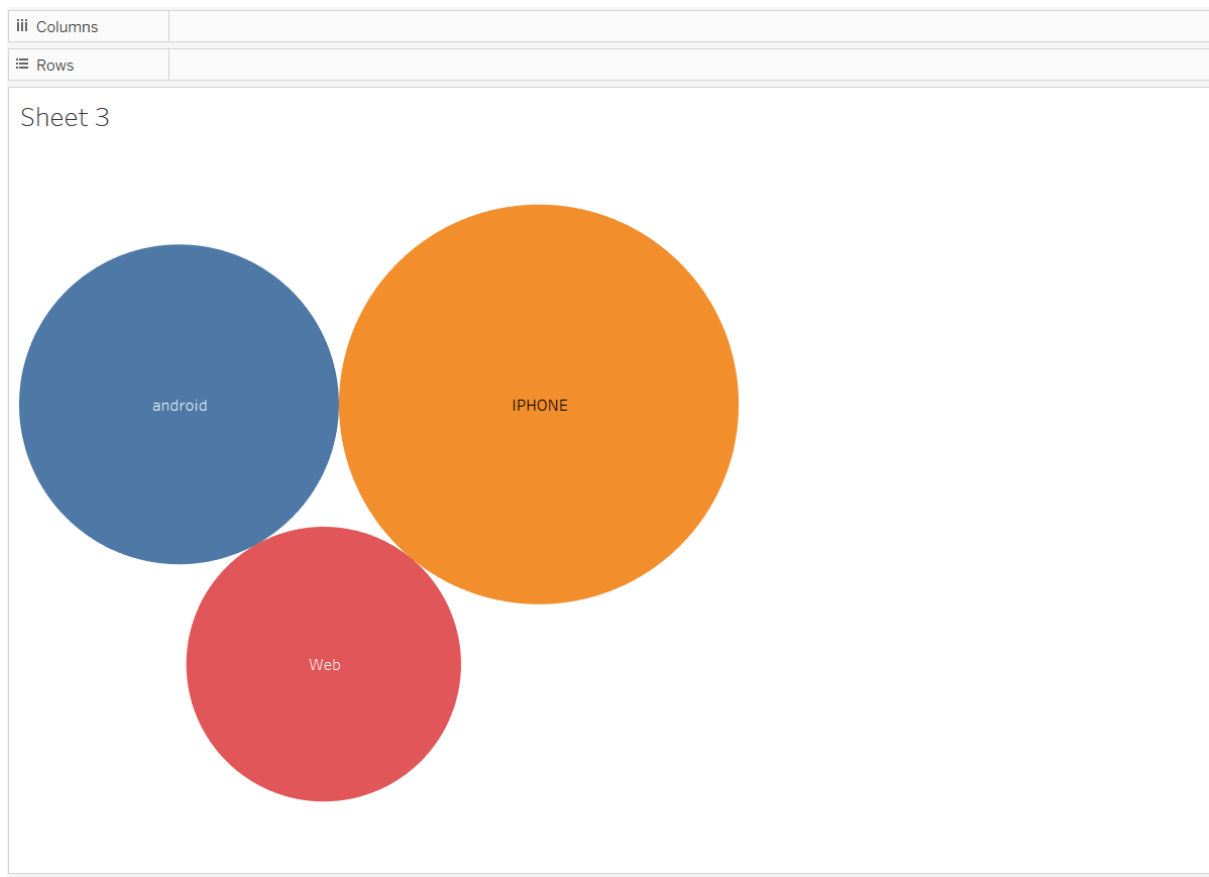
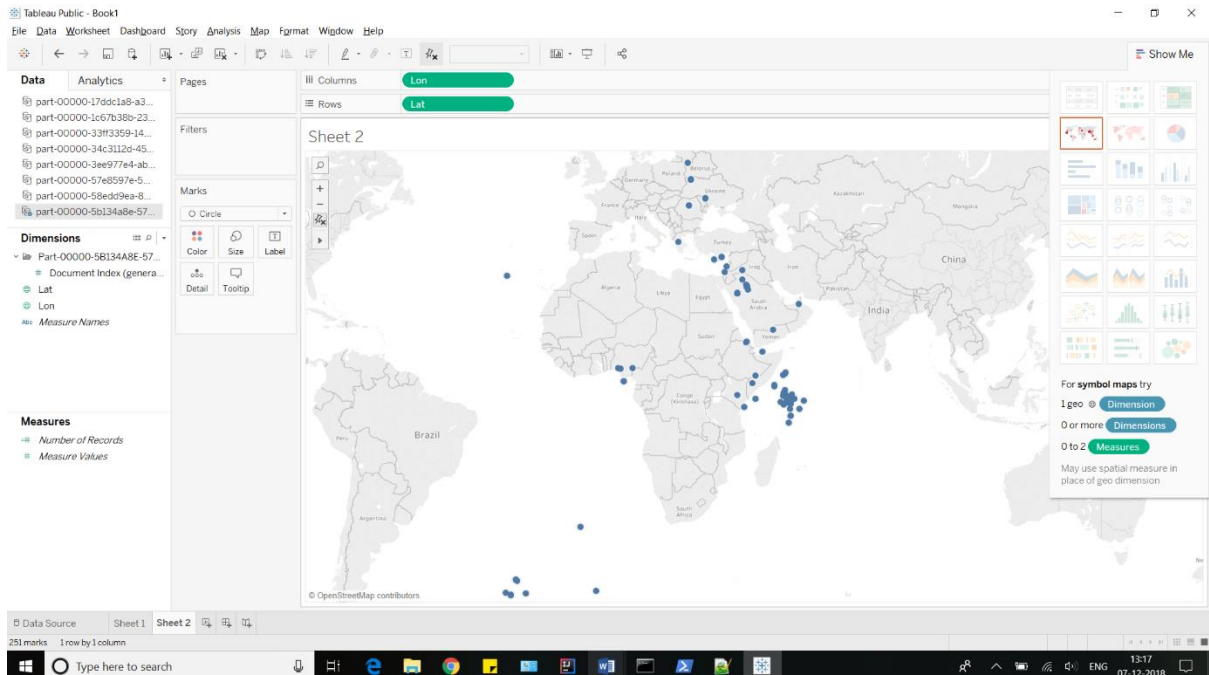**Query 9:** Top 10 Users Who are actively liking tweets.

```
//Query 9 - Top 10 Users Who are actively liking tweets.
val topSources = sqlContext.sql( sqlText = "select count(*) as total,"+
  "CASE WHEN source like '%iphone%' then 'IPHONE'"+"WHEN source like '%android%' then 'android'"+ "WHEN source like '%Web%' then 'Web'"+ "END AS source " +
  "from twit group by source order by count(1) desc limit 3")
```

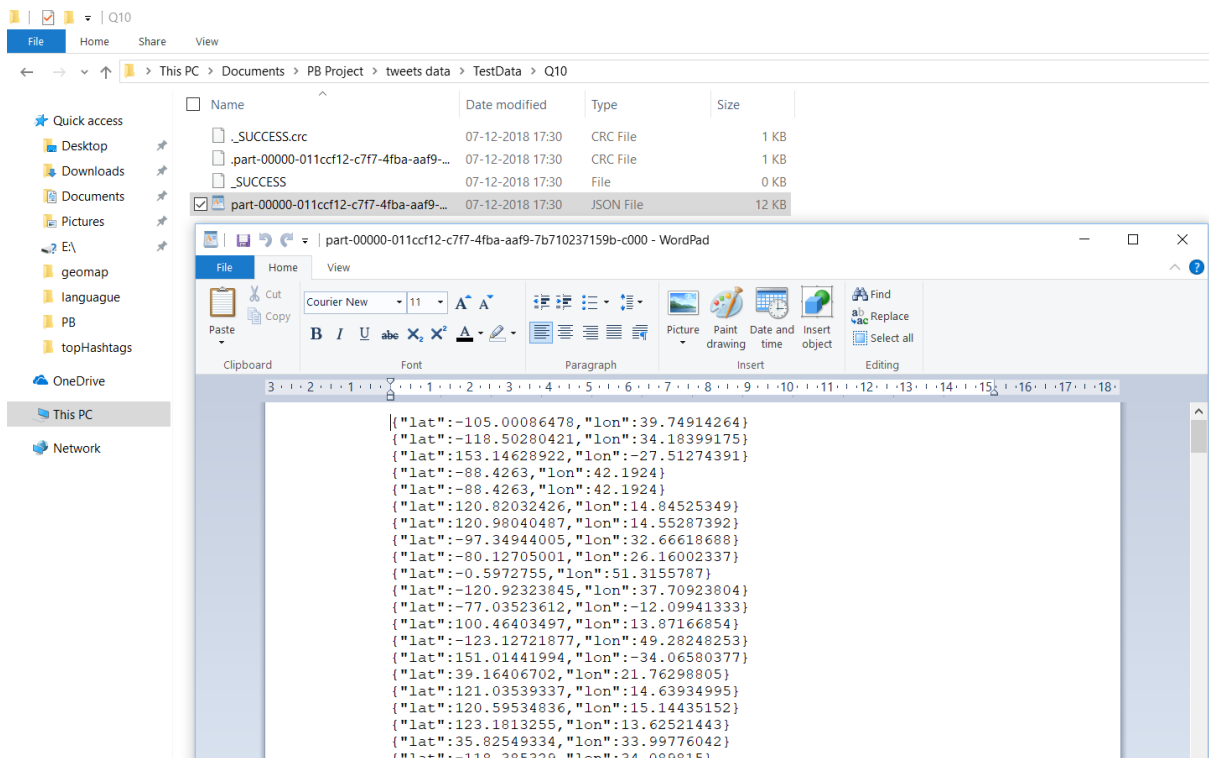| iii Columns | |
| --- | --- |
| ☰ Rows | |

Sheet 3



## Query10: Location of sports from where tweets are coming

```
// Query 10 - Location of Tweets on Map
val geomap=sqlContext.sql( sqlText = "SELECT coordinates.coordinates[0] as lat,coordinates.coordinates[1] as lon FROM twit where coordinates is not null")
```
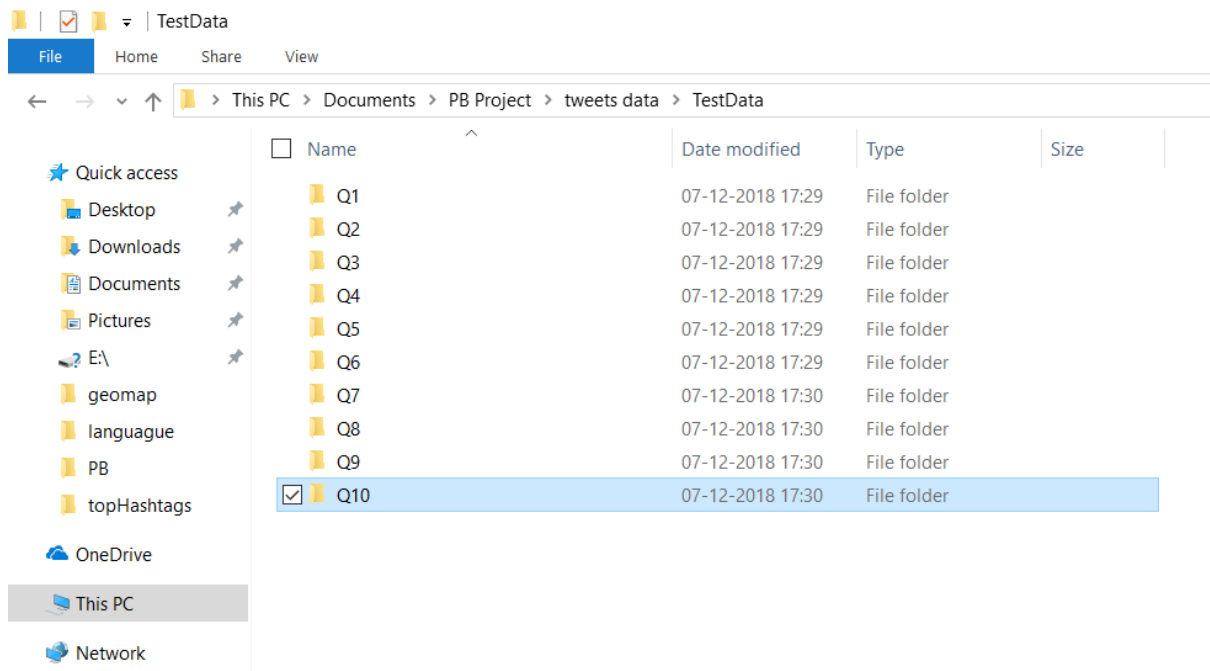
## Testing:

Generated output from the queries are written to the json files to local directory and done manual testing on the output fields whether the output is generated or not.

{"lat":-105.00086478,"lon":39.74914264}
{"lat":-118.50280421,"lon":34.18399175}
{"lat":153.14628922,"lon":-27.51274391}
{"lat":-88.4263,"lon":42.1924}
{"lat":-88.4263,"lon":42.1924}
{"lat":120.82032426,"lon":14.84525349}
{"lat":120.98040487,"lon":14.55287392}
{"lat":-97.34944005,"lon":32.66618688}
{"lat":-80.12705001,"lon":26.16002337}
{"lat":-0.5972755,"lon":51.3155787}
{"lat":-120.92323845,"lon":37.70923804}
{"lat":-77.03523612,"lon":-12.09941333}
{"lat":100.46403497,"lon":13.87166854}
{"lat":-123.12721877,"lon":49.28248253}
{"lat":151.01441994,"lon":-34.06580377}
{"lat":39.16406702,"lon":21.76298805}
{"lat":121.03539337,"lon":14.63934995}
{"lat":120.59534836,"lon":15.14435152}
{"lat":123.1813255,"lon":13.62521443}
{"lat":35.82549334,"lon":33.99776042}
{"lat":-118.385329,"lon":34.089815}

Output folders that are generated from the queries.



**References:**

https://developer.twitter.com/en/docs/tweets/data-dictionary/overview/tweet-object

# Thank You