



## **COLLECTING TWEETS USING TWITTER STREAMING API'S**

Principles of Big Data Management  
(Phase 2 of Project)

### **Team members**

Sindhusha Tiyyagura(16280708)

Pradeepika Kolluru(16283597)

Thoshita Movva (16279838)

### **Instructor**

Dr. PRAVEEN RAO, Ph.D.

## Abstract:

Phase 2 of this project deals with the following requirements:

1. Writing interesting analytical queries on twitter data that we collected.
2. Developing interesting visualizations.

## Title:

Technology in different domains.

## Technologies and Tools used:

1. Hadoop
2. Spark
3. Scala
4. Tableau
5. Java and MapReduce programs

## Queries and Analysis:

### **Query-1 :**

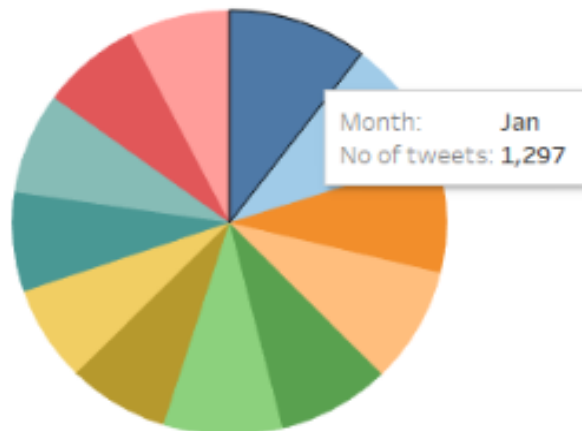
In this query, We found number of users created based on month.

### **Code :**

```
scala> val Query1 = sqlContext.sql("SELECT substring(user.created_at,5,3) as month, count(user.id) from tweetDatatable group by month");
19/05/03 10:08:14 WARN ObjectStore: Failed to get database global_temp, returning NoSuchObjectException
Query1: org.apache.spark.sql.DataFrame = [month: string, count(user.id AS 'id'): bigint]
```

```
scala> Query1.show();
+-----+-----+
|month|count(user.id AS 'id')|
+-----+-----+
|Oct|968|
|Sep|946|
|Dec|946|
|Aug|912|
|May|1054|
|Jun|1110|
|Feb|1207|
|Nov|928|
|Mar|1101|
|Jan|1297|
|Apr|1109|
|Jul|943|
+-----+-----+
```

## Visualization :



## Query-2 :

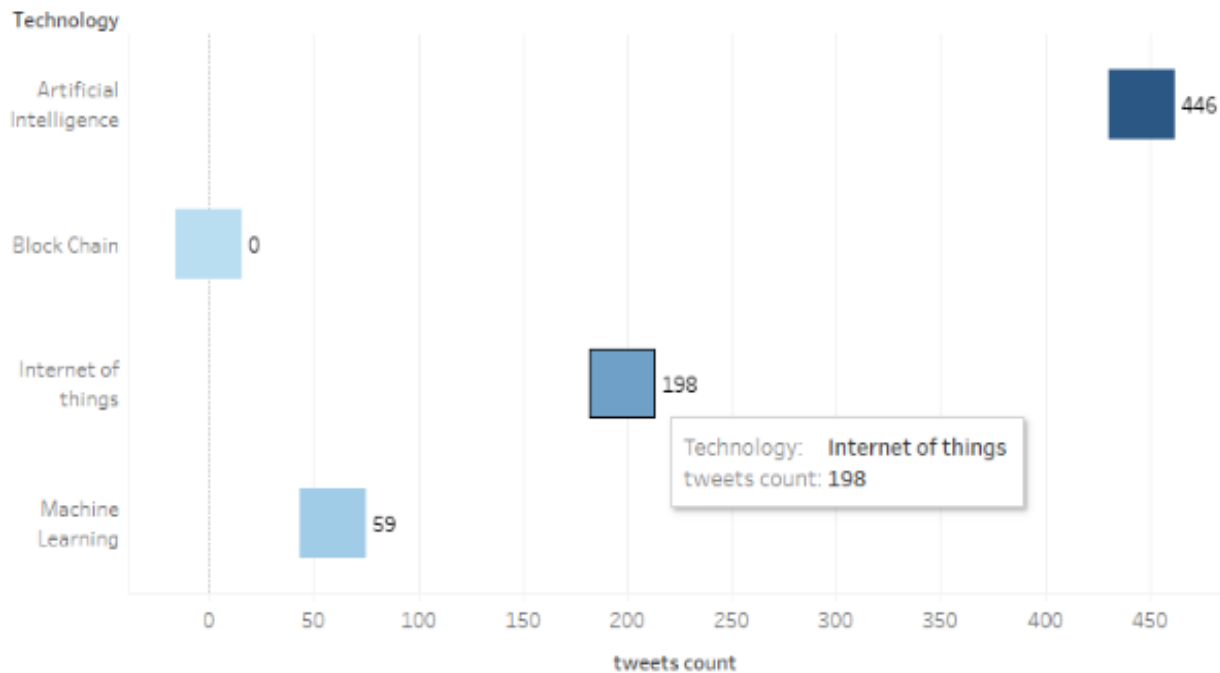
In this query we found the trend in twitter for Artificial Intelligence, Internet of Things, Machine learning, Block chain etc.,

## Code :

```
scala> val Query2 = sqlContext.sql("SELECT COUNT(*) AS NumberOfTweets, 'Artificial Intelligence' as Language FROM Technology where text LIKE '%Artificial Intelligence%' or text like '%AI%' UNION SELECT COUNT(*) AS NumberOfTweets, 'Internet of Things' as Language FROM Technology where text LIKE '%Internet of Things%' or text like '%IoT%' UNION SELECT COUNT(*) AS NumberOfTweets, 'Machine Learning' as Language FROM Technology where text LIKE '%Machine Learning%' or text like '%ML%' UNION SELECT COUNT(*) AS NumberOfTweets, 'Block Chain' as Language FROM Technology where text LIKE '%Block Chain%' or text like '%Blockchain%'")
Query2:Query2.show()
```

NumberOfTweets	Language
198	Internet of things
446	Artificial Intell...
0	Block Chain
59	Machine Learning

## Visualization :

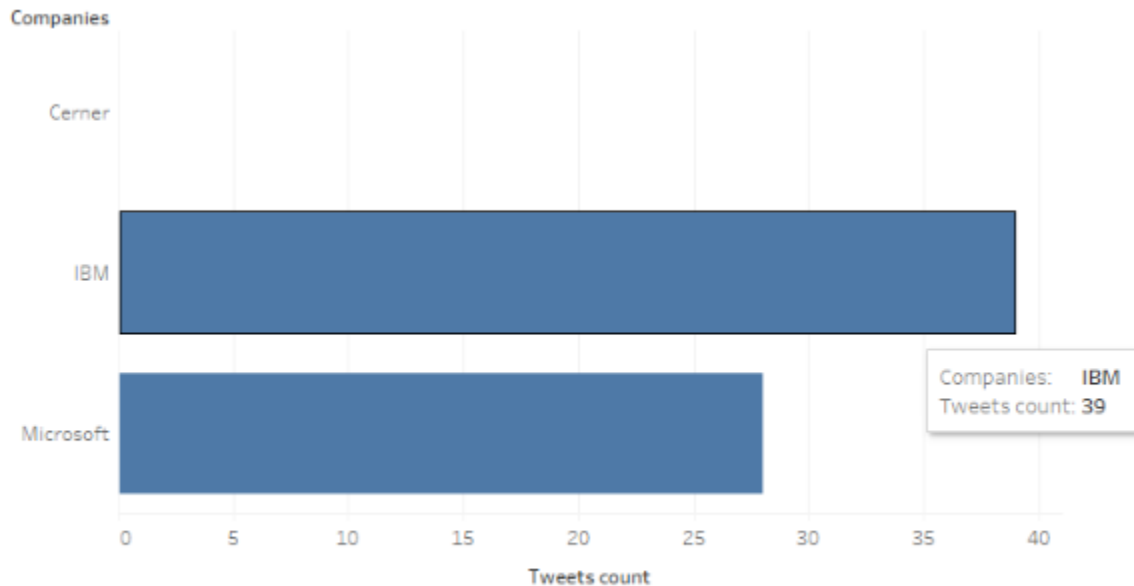


## Query-3 :

In this query we found how Microsoft, Cerner and IBM associates/ employees are actively tweeting about AI.

```
scala> val Query3 = sqlContext.sql("SELECT 'Microsoft' as Company, count(*) as Count from Technology where text like 'Microsoft%' and (text like '%technology%' or text like '%tech%') UNION SELECT 'IBM' as Company, count(*)  
Query3: org.apache.spark.sql.DataFrame = [Company: string, Count: bigint]  
scala> Query3.show()  
+-----+-----+  
| Company|Count|  
+-----+-----+  
|Microsoft| 28|  
|Cerner| 0|  
|IBM| 39|  
+-----+-----+
```

## Visualization :



## Query-4 :

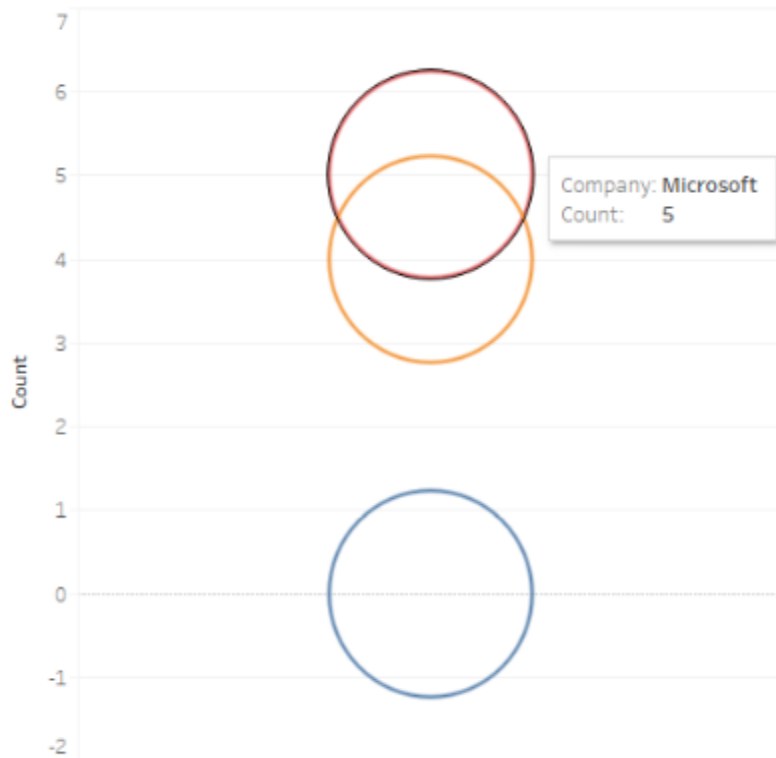
In this query we found number of tweets made by companies like Amazon, Microsoft, IBM on technologies like Artificial Intelligence, Machine Learning, Internet of Things etc.,

## Code :

```
scala> val Query4 = sqlContext.sql("SELECT 'Microsoft' as Company, count(*) as Count From Technology where text like 'Microsoft%' and (text like '%AI%' or text like '%IoT%' or text like '%MLN%') UNION SELECT 'IBM' as Company, count(*) as Count From Technology where text like 'IBM%' and (text like '%AI%' or text like '%IoT%' or text like '%MLN%')")
Query4: org.apache.spark.sql.DataFrame = [Company: string, Count: bigint]

scala> Query4.show()
+-----+
| Company|Count|
+-----+
| Amazon|    0|
| IBM   |    4|
| Microsoft|  5|
+-----+
```

## Visualization :



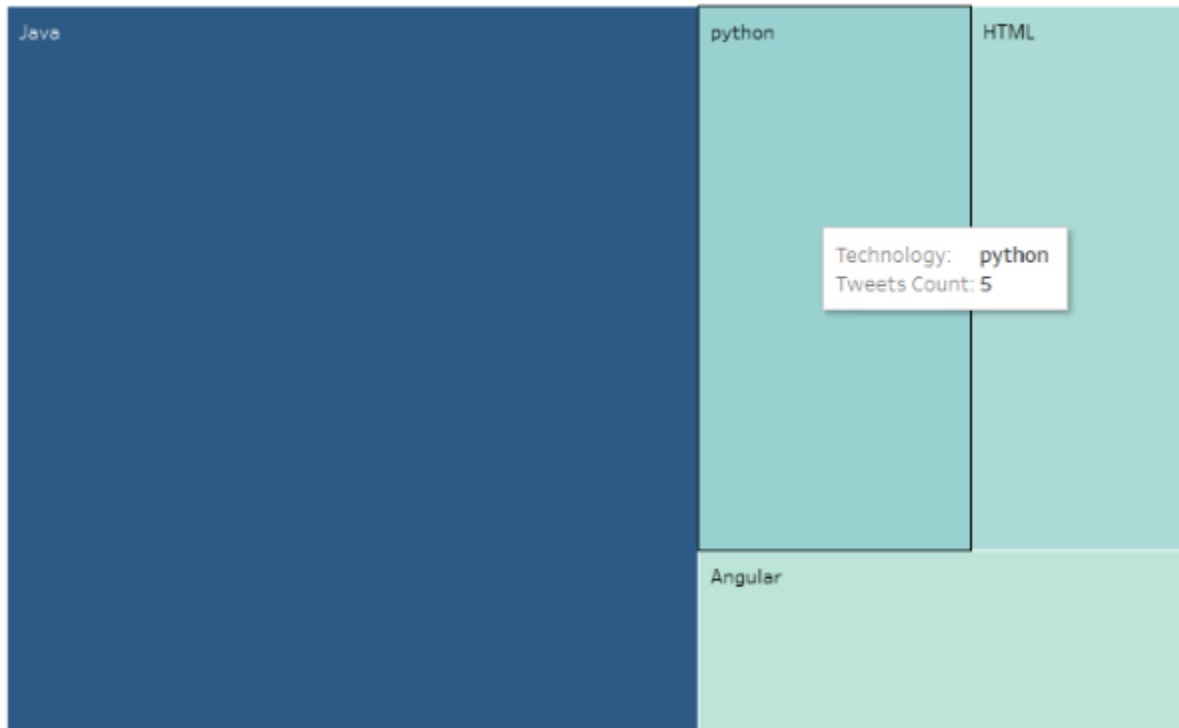
## Query-5 :

In this query we found number of tweets made on different technologies like Java, Python, HTML, Angular etc.,

## Code :

```
scala> val Query5 = sqlContext.sql("SELECT COUNT(*) AS NumberOfTweets, 'HTML' as Language FROM Technology where text LIKE 'HTML%' UNION SELECT COUNT(*) AS NumberOfTweets, 'Java' as Language FROM Technology where text LIKE '%Java%'")
Query5: org.apache.spark.sql.DataFrame = [NumberOfTweets: bigint, Language: string]
scala> Query5.show()
+-----+
|NumberOfTweets|Language|
+-----+
|17|Java|
|5|python|
|3|Angular|
|4|HTML|
+-----+
```

## Visualization :



## Query-6 :

In this query we counted number of tweets made on technologies according to days.

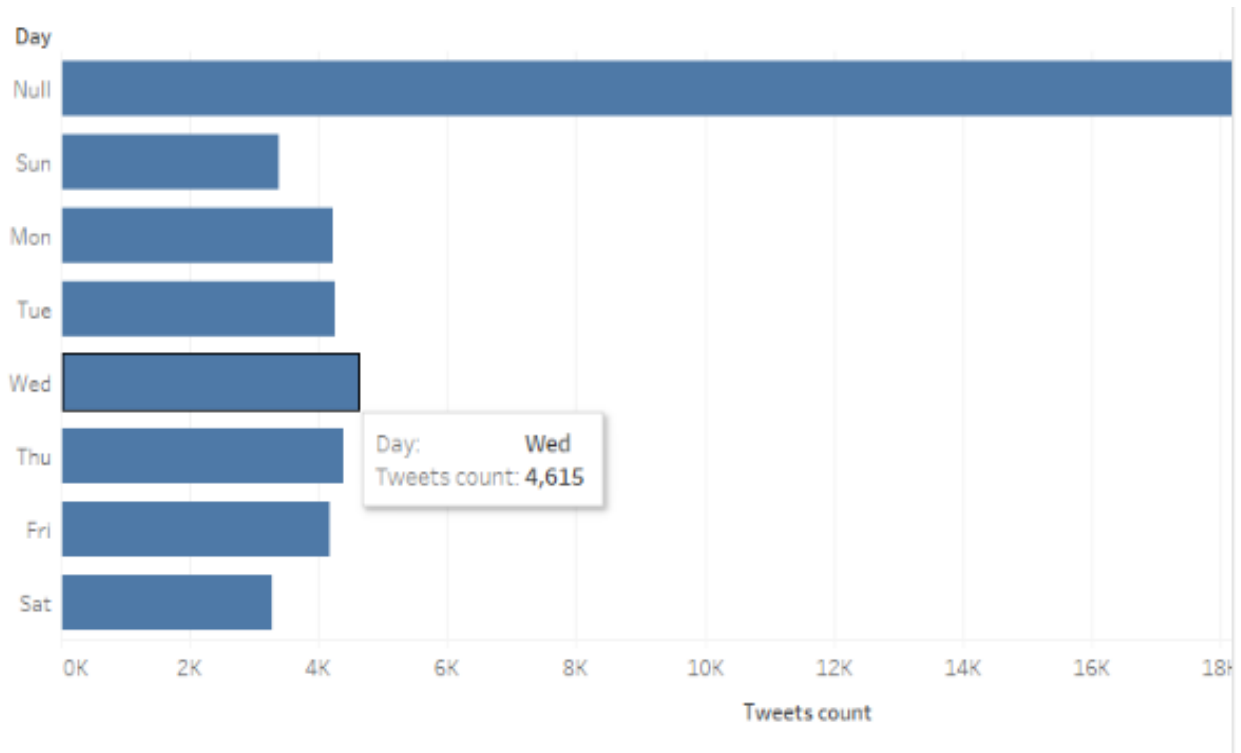
## Code :

```
scala> val Query6=sqlContext.sql("SELECT substring(user.created_at,1,3) as day,count(*) as count from technology group by day");
19/05/05 08:36:55 WARN ObjectStore: Failed to get database global_temp, returning NoSuchObjectException
Query6: org.apache.spark.sql.DataFrame = [day: string, count: bigint]
```

```
scala> Query6.show();
```

```
+---+-----+
| day|count|
+---+-----+
| Sun| 3399|
| null|22095|
| Mon| 4218|
| Thu| 4397|
| Sat| 3277|
| Wed| 4615|
| Tue| 4250|
| Fri| 4184|
+---+-----+
```

## Visualization :



## Query-7 :

In this query, we found different languages used to tweet about technologies and their count respectively.

## Code :

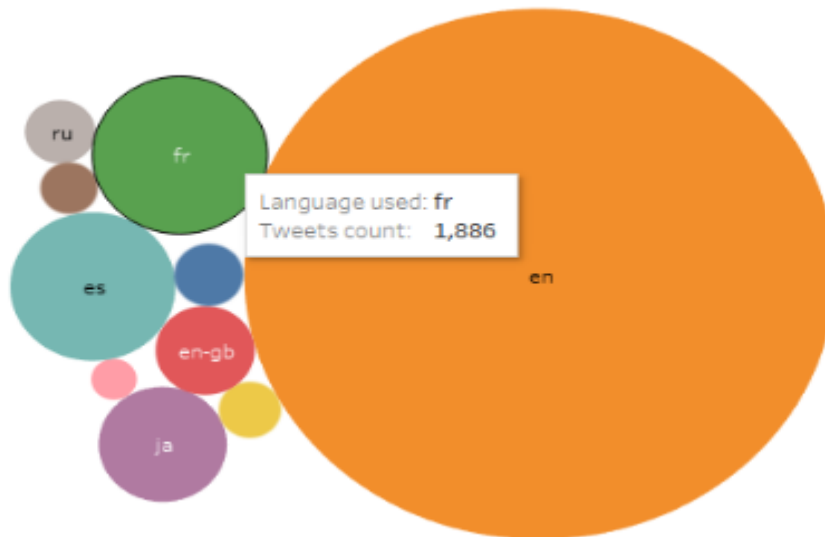
```
scala> val Query5 = sqlContext.sql("SELECT user.lang, count(*) AS count FROM technology WHERE lang<>'null' GROUP BY user.lang ORDER BY count DESC LIMIT 10");  
Query5: org.apache.spark.sql.DataFrame = [lang: string, count: bigint]
```

```
scala> Query5.show();
```

lang	count
en	21412
fr	1886
es	1677
ja	1007
en-gb	602
ru	307
de	295
it	239
pt	203
nl	128



## Visualization :



## Query-8 :

In this query, we found number of tweets made on technologies based on user names.

## Code :

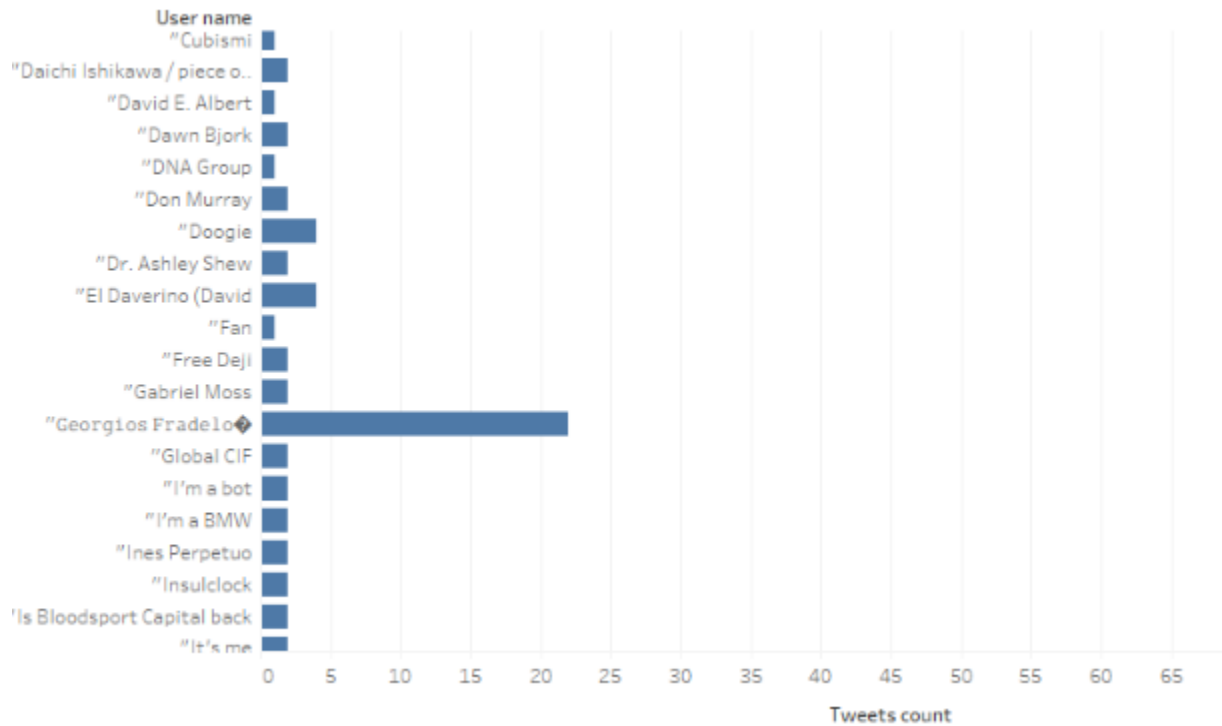
```
scala> val Query8 = sqlContext.sql("SELECT count(*) as count, user.name from Technology where user.name is not null group by user.name order by count desc");
Query8: org.apache.spark.sql.DataFrame = [count: bigint, name: string]
```

```
scala> Query8.show();
```

count	name
107	Abhishek
85	TendenciasTech
82	SkyFree Marketing
44	Lauro Espinoza Creel
40	???? ???????????
32	Money Making Arti...
32	startupcrunch
31	global-tech-news.net
28	Institute High Tech
28	ZoeGeop Technolo...
28	QCS Tech Reviews
28	Podcast Listeners
27	Technyc
26	Business News
25	The Technology?
22	????????? ...
20	Machine Learning
20	Eva Prokop
20	Tech Digi
20	GoldenHarvest

only showing top 20 rows

## Visualization :



## Query-9 :

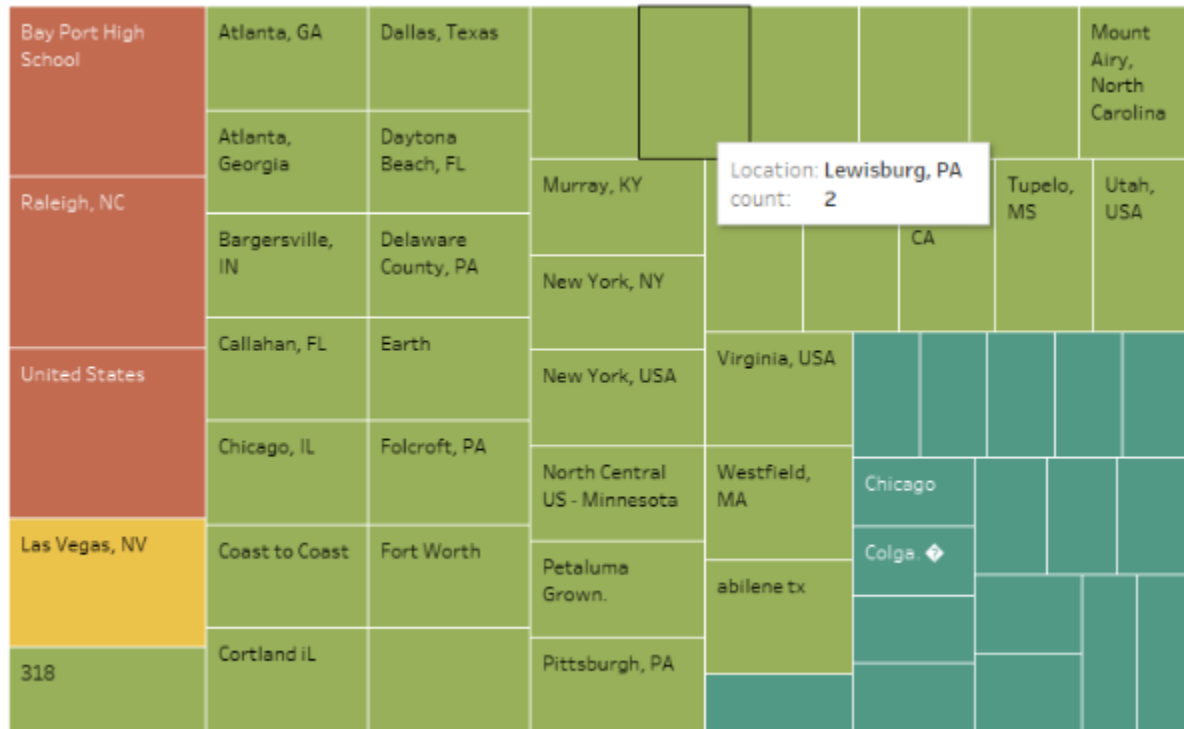
In this query, we found number of tweets made on technologies based on the location.

## Code :

```
scala> val Query9=sqlContext.sql("SELECT user.location,count(text) as count FROM technology WHERE place.country='United States' AND user.location is not null GROUP BY user.location ORDER BY count")
Query9: org.apache.spark.sql.DataFrame = [location: string, count: bigint]

scala> Query9.show();
+-----+-----+
|location|count|
+-----+-----+
|Bay Port High School|4|
|United States|4|
|Raleigh, NC|4|
|Las Vegas, NV|3|
|Delaware County, PA|2|
|Mount Airy, North...|2|
|Knoxville, TN|2|
|North Central US ...|2|
|Atlanta, Georgia|2|
|Atlanta, GA|2|
|Cortland IL|2|
|Rancho Mirage, CA|2|
|Earth|2|
|Callahan, FL|2|
|Queens, NY|2|
|HARRISONBURG|2|
|Louisiana, USA|2|
|Utah, USA|2|
|Bangorsville, IN|2|
|Westfield, MA|2|
+-----+-----+
only showing top 20 rows
```

## Visualization :



## Query-10 :

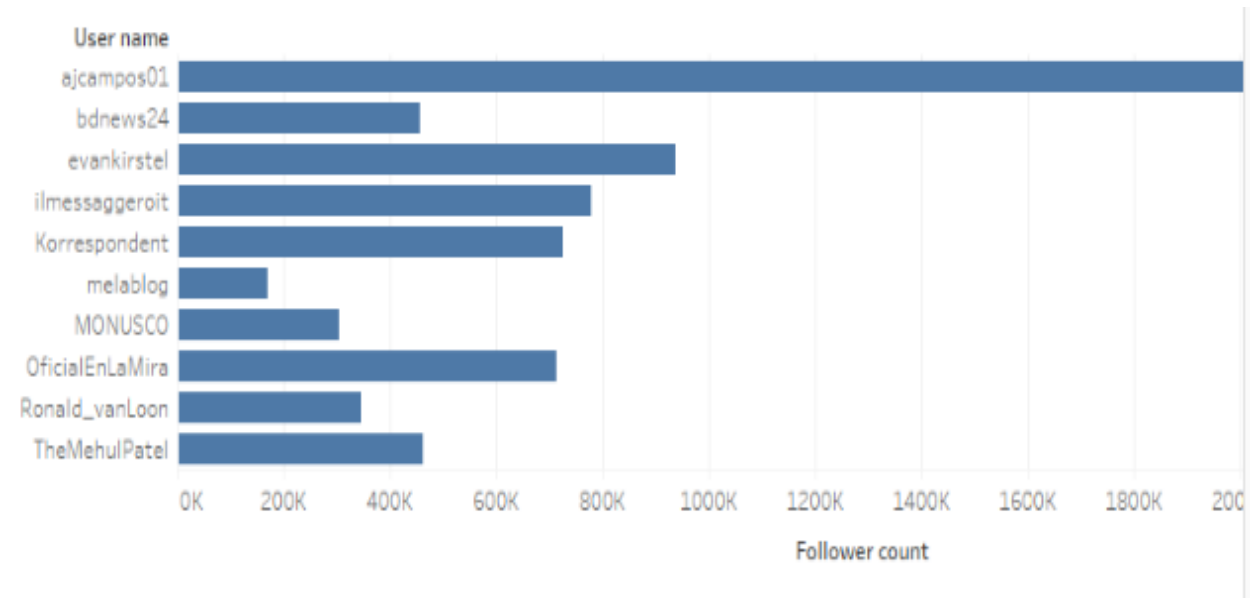
In this query, we found followers count for the non-verified accounts.

## Code :

```
scala> val Query10 = sqlContext.sql("SELECT user.screen_name,user.followers_count FROM Technology WHERE user.verified = false ORDER BY user.followers_count DESC LIMIT 20")
Query10: org.apache.spark.sql.DataFrame = [screen_name: string, followers_count: bigint]

scala> Query10.show();
+-----+-----+
| screen_name | followers_count |
+-----+-----+
| ajcampos01  | 1227793         |
| ajcampos01  | 1227793         |
| ilmessaggeroit | 389138         |
| ilmessaggeroit | 389138         |
| Korrespondent | 363049         |
| Korrespondent | 363049         |
| OfficialEnlaWira | 357310         |
| OfficialEnlaWira | 357310         |
| MONUSCO     | 303123         |
| evankirstel  | 234363         |
| evankirstel  | 234362         |
| evankirstel  | 234362         |
| evankirstel  | 234360         |
| TheMehuIPatel | 231373         |
| TheMehuIPatel | 231373         |
| bdnews24     | 227957         |
| bdnews24     | 227957         |
| Ronald_vanLoon | 172560         |
| Ronald_vanLoon | 172560         |
| me1ablog     | 171501         |
+-----+-----+
```

## Visualization :



## Query-11 :

In this query, we found the followers count for verified accounts.

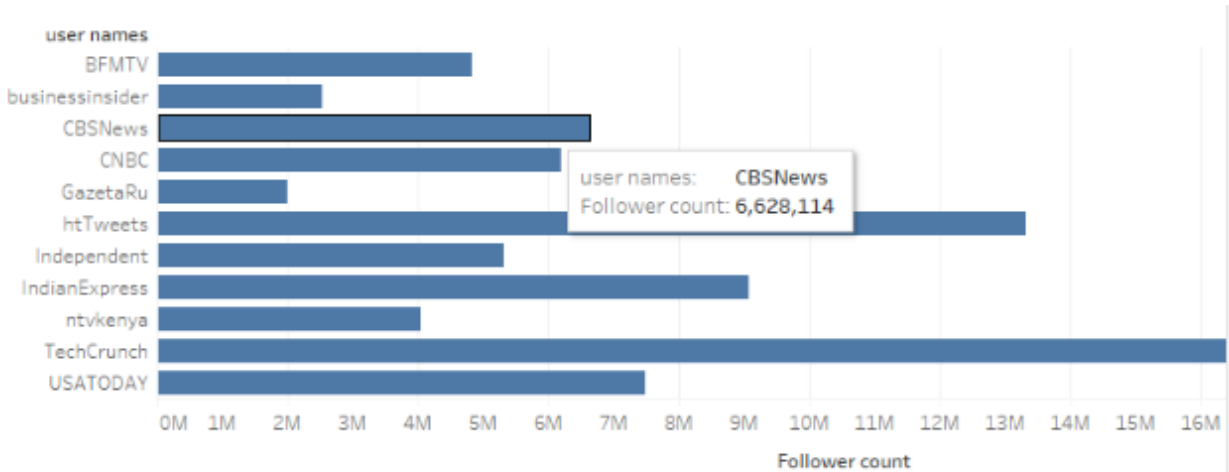
## Code :

```
scala> val Query10 = sqlContext.sql("SELECT user.screen_name,user.followers_count FROM Technology WHERE user.verified = true ORDER BY user.followers_count DESC LIMIT 20")
Query10: org.apache.spark.sql.DataFrame = [screen_name: string, followers_count: bigint]
```

```
scala> Query10.show();
```

screen_name	followers_count
TechCrunch	10031055
TechCrunch	10031055
htTweets	6654072
htTweets	6654072
CBSNews	6628114
USATODAY	3741062
USATODAY	3741062
CNBC	3093696
CNBC	3093696
IndianExpress	3027055
IndianExpress	3026993
IndianExpress	3026993
Independent	2660160
Independent	2660142
businessinsider	2547083
BFMTV	2415473
BFMTV	2415465
ntvkenya	2027313
ntvkenya	2027313
GazetaRu	1995501

## Visualization :



## Query-12 :

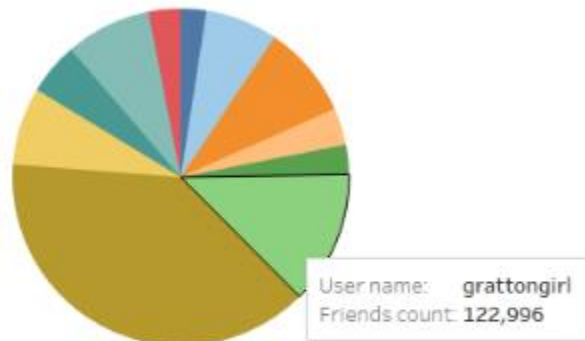
In this query, we found the friends count for verified account

## Code :

```
scala> val Query10 = sqlContext.sql("SELECT user.screen_name,user.friends_count FROM Technology WHERE user.verified = true ORDER BY user.friends_count DESC LIMIT 20");
19/05/06 08:12:50 WARN ObjectStore: Failed to get database global_temp, returning NoSuchObjectException
Query10: org.apache.spark.sql.DataFrame = [screen_name: string, friends_count: bigint]
```

```
scala> Query10.show();
+-----+-----+
| screen_name | friends_count |
+-----+-----+
| HerbertRSim | 185836         |
| HerbertRSim | 185836         |
| grattongirl | 41002          |
| grattongirl | 40997          |
| grattongirl | 40997          |
| nikkeibpITpro | 39535         |
| nikkeibpITpro | 39535         |
| Kevin_Jackson | 35872         |
| Kevin_Jackson | 35872         |
| frenchweb | 34061          |
| DavidPapp | 33418          |
| DavidPapp | 33418          |
| ZAGrrl | 29795          |
| derStandardat | 28705         |
| derStandardat | 28705         |
| derStandardat | 28705         |
| gratonboy | 27341          |
| Lenovodc | 24062          |
| Lenovodc | 24062          |
| cfarivar | 23557          |
+-----+-----+
```

## Visualization :



## Query-13 :

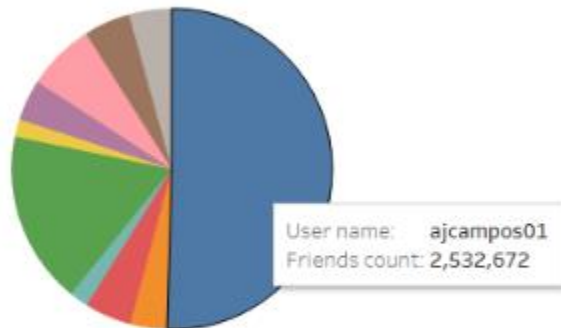
In this query, we found the friends count for the non-verified accounts.

## Code :

```
scala> val Query10 = sqlContext.sql("SELECT user.screen_name,user.friends_count FROM Technology WHERE user.verified = false ORDER BY user.friends_count DESC LIMIT 20");
Query10: org.apache.spark.sql.DataFrame = [screen_name: string, friends_count: bigint]
```

```
scala> Query10.show();
+-----+-----+
| screen_name | friends_count |
+-----+-----+
| ajcampos01  | 1266336       |
| ajcampos01  | 1266336       |
| evankirstel | 218615        |
| evankirstel | 218615        |
| evankirstel | 218613        |
| evankirstel | 218613        |
| Ronald_vanLoon | 171635      |
| Ronald_vanLoon | 171635      |
| BSkylstad   | 118228       |
| BSkylstad   | 118228       |
| stojkovic_alex | 117484      |
| stojkovic_alex | 117484      |
| The_News_DIVA | 108282       |
| The_News_DIVA | 108282       |
| ipfconline1 | 103759       |
| ipfconline1 | 103759       |
| chrisifg    | 96121        |
| bdnews24    | 93101        |
| bdnews24    | 93101        |
| GotStockTips | 88750        |
+-----+-----+
```

## Visualization :



## Query-14 :

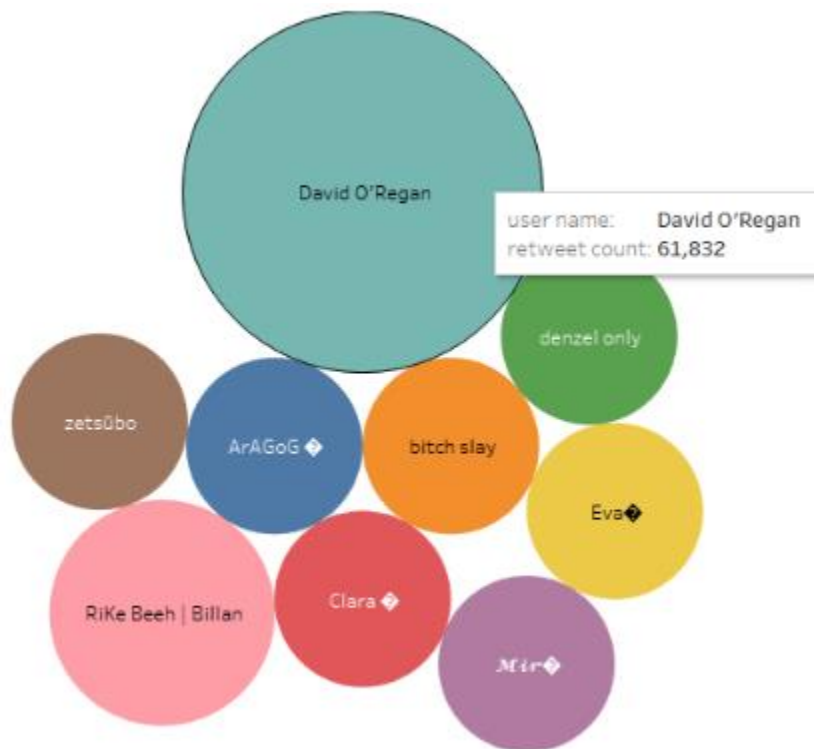
In this query, we found the number of retweets on a particular tweet for a technology based on user.

## Code :

```
scala> val Query10 = sqlContext.sql("SELECT user.name ,retweeted_status.retweet_count AS Retweet_Count FROM technology WHERE retweeted_status.retweet_count IS NOT NULL")
Query10: org.apache.spark.sql.DataFrame = [name: string, Retweet_Count: bigint]

scala> Query10.show();
+-----+-----+
| name | Retweet_Count |
+-----+-----+
| David O'Regan | 30916 | |
| David O'Regan | 30916 |
| RiKe Beeh | Billan | 24216 |
| Eva? | 14820 |
| ArAGOG ? | 14818 |
| zets?bo | 14817 |
| bitch slay | 14816 |
| Clara ? | 14815 |
| denzel only | 14814 |
| ??? | 14813 |
+-----+-----+
```

## Visualization :



## Query-15:

In this query, we found technologies in different fields like gaming, entertainment, movies ets.,

## Code :

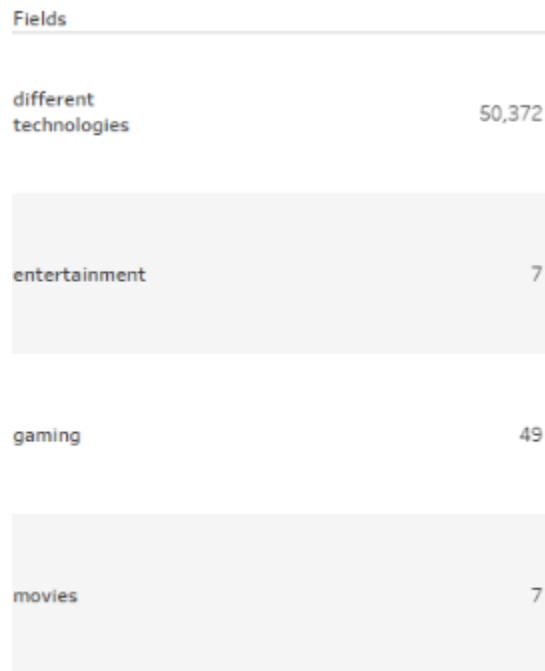
```
scala> val Query9=sqlContext.sql("select count(*) as count,q.text from (select case when text like '%gaming%' then 'gaming' when text like '%entertainment%' then 'entertainment' when text like '%movies%' then 'movies' else 'different technologies' end as text) q")
Query9: org.apache.spark.sql.DataFrame = [count: bigint, text: string]
```

```
scala> Query9.show();
```

count	text
7	movies
7	entertainment
50372	different technologies
49	gaming



## Visualization :



## Testing

### Manual testing :

We tried to test the results manually by taking the text from the collected tweets and finding the tweets using twitter search engine.

For instance consider the text “*South Korean tech firms Netmarble and Kakao as well as private equity fund MBK Partners submitted initial bids*” which is taken from the collected tweets and searched it manually in twitter to find the tweet.

