# PRINCIPLES OF BIG DATA PROJECT

# TRAVEL

**TEAM MEMBERS : Shreyaa Sridhar**

**Vinay Maturi**

**Hari**

# SOURCE CODE

https://drive.google.com/drive/folders/1cj1nqNs5ntlsiEZfSP12TR1WTYtwikCS?usp=sharing

# YOUTUBE LINK

https://youtu.be/-ZggYl9NsL8

# TECHNOLOGIES USED

- Pycharm
- Cloudera
- IntelliJ
- Scala
- Spark
- Apache Hadoop
- Tableau

# PHASE 1

- Collected tweets and saved as a JSON file.
- From the JSON file , extracted the hashtags and URLs and saved as text file.
- Word count on the hashtags and URLs using Apache Hadoop and Apache Spark.

```python
import sys
from tweepy import OAuthHandler
from tweepy import Stream
from tweepy.streaming import StreamListener

consumer_key='cWf9jAnJHeDY5ECjDoKyijl6z'
consumer_secret='Gi1cZa9jUFzIiup0wvTUtykSy0ygCyvdpAed3YvSYZ8RpOKlzH'
access_token='3280449488-PRUJP4mIjXIYFmYdhrlvsM38hqsXBdqkYRqwmUa'
access_token_secret='WN6vpG4EJzzJGLtaXc9Bd5w0rxn9H5sVCpAtTYoWPv2z1'

class StdOutListener(StreamListener):
    def on_data(self,data):
        print(data)
        saveFile=open('TweetsTravel.json','a')
        saveFile.write(data)
        saveFile.write('\n')
        saveFile.close()
        return True

    def on_error(self,status):
        print(status)

if __name__== '__main__':
    x=StdOutListener()
    auth = OAuthHandler(consumer_key, consumer_secret)
    auth.set_access_token(access_token, access_token_secret)

    stream = Stream(auth, x)
    stream.filter(track=['vacation', '#vacation', 'travel', '#travel','vacay',
    '#vacay','nature','#nature','#photography','@lonelyplanet_in •','@NGTIndia •','@travelscenes •','Travel Vibes','@thetblogger
    •','@timestravel','@TravelLeisure •','@travelchannel •','@TripAdvisor •','@nytimestravel •','@usatodaytravel
    •','@NBCNewsTravel •','@GoogleTravel •','@BudgetTravel •','@BBCTravelShow •','@SmarterTravel •','@ST_Travel •','@travelmail
    •','@SFChronTravel •','@TravelMagazine ',
    '#adventure','#TravelSkills','@BBC_Travel','@travelskills','#travelsmart','#tour','#travelblog','#exploring','#travelling','
    @lonelyplanet','@CNTraveler','@travelingspots','@LuxuryTravel77','@easyplanetravel','#familytravel','@NatGeoTravel','#Travel
    Ban','Travel ban','@TravelBan'])
```

```python
import codecs
import json
import sys

def parse_json_tweet(line):
    tweet = json.loads(line)
    htags = [hashtag['text'] for hashtag in tweet['entities']['hashtags']]
    urls = [url['expanded_url'] for url in tweet['entities']['urls']]
    return [htags, urls]
if __name__ == "__main__":
    file_timeordered_json_tweets = codecs.open("TweetsTravel.json", 'r', 'utf-8')
    fout = codecs.open("TravelsOutput.txt", 'w', 'utf-8')
for line in file_timeordered_json_tweets:
    try:
        [htags, urls] = parse_json_tweet(line)
        fout.write(str([htags, urls]) + "\n")
    except:
        pass

file_timeordered_json_tweets.close()
fout.close()
```
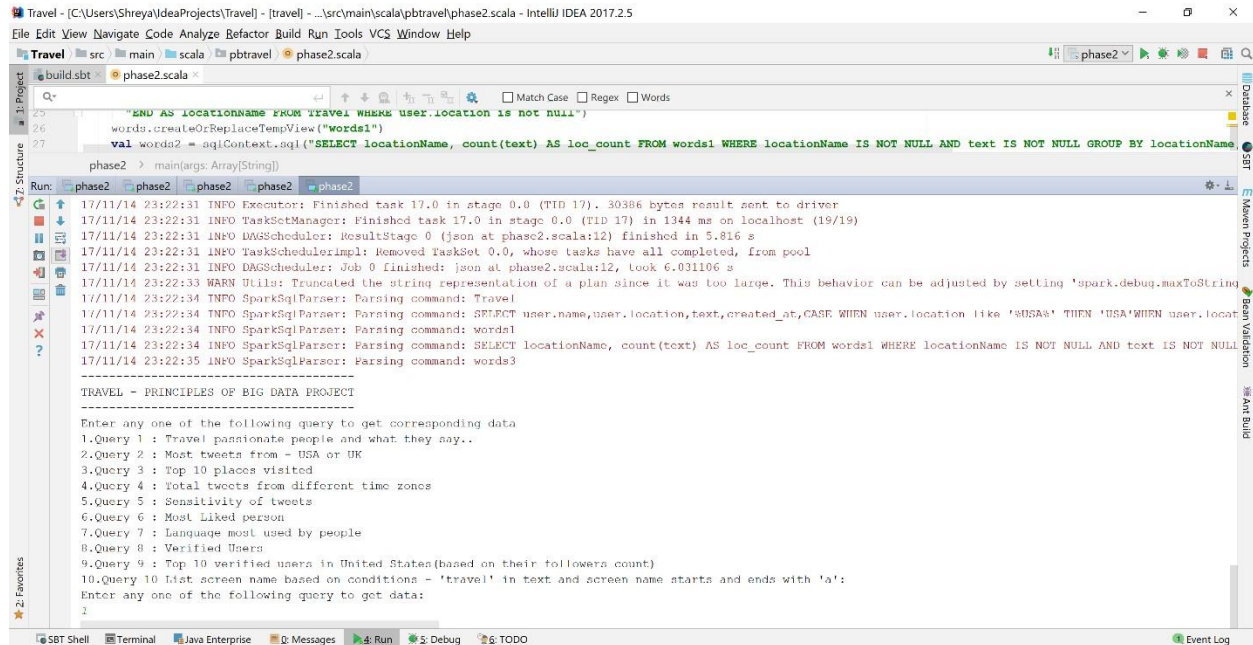
# PHASE 2

## ANALYTICAL QUERIES



### 1. Travel passionate people and what they say..

**SELECT user.name, text FROM Travel WHERE text LIKE '%travel%'**

```
17/11/14 19:41:52 INFO DAGScheduler: Job 1 finished: show at phase2.scala:49, took 0.135430 s
+-------------------+-------------------+
|               name|               text|
+-------------------+-------------------+
|    Alexandra Weiser|It is really some...|
|       Ania_Travels|I loved walking a...|
|       Realetybytes|RT @starcrosswolf...|
|     Salena Salazar|@ninjasexparty wi...|
|        Jobs at CWT|Want to work at C...|
|        lancy ✈👀😩💯|RT @MrGoodBeard_:...|
|    Top 10 Travel ✈ |#... https://t.co...|
|                eba|RT @Blackenedruin...|
|    Cristian Ayalað|Nothing to say......|
|      Steve Collins|RT @kylegriffin1:...|
|     Kamelia Britton|RT @nomadbytrade1...|
|      🐍 Ṣℂ𝑜𝑟𝑝𝑖𝑜𝑛 🐍 |RT @MK7786: Major...|
|        zmoneygrip|RT @SeattleTam: @...|
|           Zoe King|RT @NEI: Nuclear ...|
|        Sofea Varsy|RT @MrGoodBeard_:...|
|     Marisol Becerra|Want to travel to...|
|               gina|RT @starcrosswolf...|
|       Philip Rocco|RT @donmoyn: As t...|
|           EDMit it|That moment when ...|
|Intelligent Traveler|RT @Sat_Sightseer...|
+-------------------+-------------------+
only showing top 20 rows
```

2. **Most tweets from USA or UK**

```
SELECT user.name,user.location,text,created_at," +
  "CASE WHEN user.location like '%USA%' THEN 'USA'" +
  "WHEN user.location like '%States%' THEN 'USA'" +
  "WHEN user.location like '%Kingdom%' THEN 'UK'" +
  "WHEN user.location like '%Scotland%' THEN 'UK'" +
  "WHEN user.location like '%England%' THEN 'UK'" +
  "WHEN user.location like '%Ireland%' THEN 'UK'" +
  "WHEN user.location like '%Wales%' THEN 'UK'" +
  "WHEN user.location like '%Europe%' THEN 'UK'" +
  "END AS locationName FROM Travel WHERE user.location is not null
```

```
SELECT locationName, count(text) AS loc_count FROM words1 WHERE
locationName IS NOT NULL AND text IS NOT NULL GROUP BY locationName,text
ORDER BY COUNT(text) DESC
```

```
SELECT locationName, SUM(loc_count) AS No_of_tweets FROM words3 WHERE
locationName IS NOT
NULL AND loc_count IS NOT NULL GROUP BY locationName ORDER BY
SUM(loc_count) DESC
```

```
17/11/14 19:11:25 INFO CodeGenerator: Code generated in 6.674077 ms
+-----------+-----------+
|locationName|No_of_tweets|
+-----------+-----------+
|        USA|       9713|
|         UK|       3595|
+-----------+-----------+
```

### 3. Top 10 places visited

**SELECT DISTINCT user.location FROM Travel WHERE user.location IS NOT NULL LIMIT 20**

```
17/11/14 19:36:09 INFO CodeGenerator: Code generated in 13.278126 ms
+--------------------+
|            location|
+--------------------+
|  San Diego, CA, USA|
|         Tacoma, WA|
|         Chicago, IL|
|        Midland, MI|
|      Minnesota, USA|
|In nearly 150 cou...|
|              2 ▯ ▯ ▯|
|        San Jose, CA|
|Bay Area and San ...|
|People's Republic...|
|    Juba,South Sudan|
|         Chicago, IL|
|          Australia|
|            sandton |
|Bff: Cleopatria S...|
|        Newtown, PA|
|       United States|
|    Alexandria, Egypt|
|            (Madrid)|
|  San Diego, CA, USA|
+--------------------+
only showing top 20 rows
```

### 4. Total tweets from different time zones

**SELECT user.time_zone,count(text) AS Total FROM Travel WHERE user.time_zone IS NOT NULL GROUP BY user.time_zone ORDER BY Total DESC LIMIT 10**

```
17/11/14 19:57:19 INFO CodeGenerator: Code generated in 17.096324 ms
+--------------------+-----+
|           time_zone|Total|
+--------------------+-----+
|Pacific Time (US ...|17631|
|Eastern Time (US ...|11352|
|Central Time (US ...| 6102|
|              London| 3608|
|               Quito| 1776|
|Atlantic Time (Ca...| 1565|
|Mountain Time (US...| 1436|
|             Arizona| 1079|
|           Amsterdam| 1027|
|              Athens|  704|
+--------------------+-----+
```

5. **Sensitivity of tweets**

**SELECT possibly_sensitive,count(*) AS value  FROM Travel GROUP BY possibly_sensitive**

```
17/11/14 20:08:45 INFO CodeGenerator: Code generated in 7.328918 ms
+-----------------+-----+
|possibly_sensitive|value|
+-----------------+-----+
|             null|46852|
|             true|  802|
|            false|58055|
+-----------------+-----+
```

6. **Most Liked person**

**SELECT user.screen_name AS Screen_Name,user.favourites_count AS Favorites FROM Travel ORDER BY Favorites DESC LIMIT 3**

```
+-----------+---------+
|Screen_Name|Favorites|
+-----------+---------+
|    TheBoydP|  1013362|
|  paoloigna1|   730579|
|  paoloigna1|   730570|
+-----------+---------+
```

7. **Language most used by people**

**SELECT user.lang,count(*) AS lang_count  FROM Travel WHERE user.lang IS NOT NULL GROUP BY user.lang ORDER BY lang_count DESC LIMIT 1**

```
17/11/14 20:19:36 INFO CodeGenerator: Code generated in 7.336801 ms
+----+----------+
|lang|lang_count|
+----+----------+
|  en|     88144|
+----+----------+
```

## 8. Verified Users

**SELECT count(user.verified) AS Verified_Users FROM Travel WHERE user.verified=true**

```
17/11/14 21:16:59 INFO CodeGenerator: Code generated in 8.932625 ms
+-------------+
|Verified_Users|
+-------------+
|         1012|
+-------------+
```

## 9. Top 10 verified users in United States(based on their followers count)

**SELECT user.name AS Name FROM Travel WHERE user.verified=true AND user.location='United States' ORDER BY user.followers_count DESC LIMIT 10**

```
17/11/14 22:19:23 INFO CodeGenerator: Code generated in 6.981555 ms
+----------------+
|            Name|
+----------------+
|             NWS|
|  beIN SPORTS USA|
|     Clark Howard|
|     Clark Howard|
|   The Wirecutter|
|          AskTSA|
|     miki howard|
|   Liberty Travel|
|   Monsanto BioAg|
|Kaiser Permanente|
+----------------+
```

10. **List the User's Screen Name, followers count , friends count and favourites count when it satisfies two conditions : 'Travel' in the text and screen name starts with 'a'.**
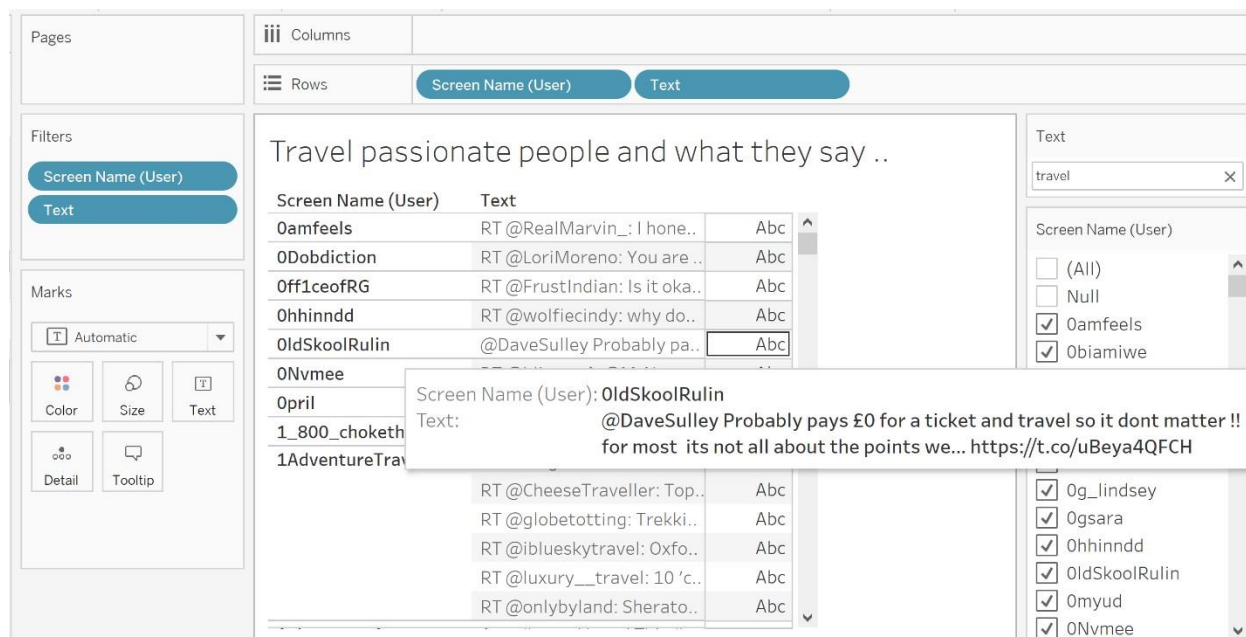
**SELECT user.screen_name, user.followers_count,user.favourites_count, user.friends_count FROM Travel WHERE user.screen_name LIKE 'a%' AND text LIKE '%travel%'**

```
|    screen_name|followers_count|favourites_count|friends_count|
+---------------+---------------+----------------+-------------+
|aldub_sherry03|             37|              24|           60|
|      afhreaves|            51|             662|          256|
|  amazintravels|           518|               2|         1884|
|        azimgpj|           581|            7477|         1243|
|        azimgpj|           581|            7478|         1243|
|        azimgpj|           581|            7478|         1243|
|     appertunity|          2756|           12281|         2720|
|      anymattoni|          2260|           32350|         4669|
|     andresweep1|           101|               0|            1|
|    aiuabdullah_|          1063|            6987|          267|
|africanesque79|          1899|            4534|         3266|
|alihunterwadas|           168|            3422|          103|
|           awtyk|           469|             792|          251|
|africanesque79|          1899|            4535|         3266|
|    asropuli1984|             0|               0|           48|
|     abbacharice|          1690|             983|          831|
|    azarialaster|           246|             520|          199|
|       anuj_sena|          1865|           12311|          246|
|     ash_slays13|           299|            7847|          382|
|    anuragmodi31|            42|            1648|          232|
+---------------+---------------+----------------+-------------+
only showing top 20 rows
```
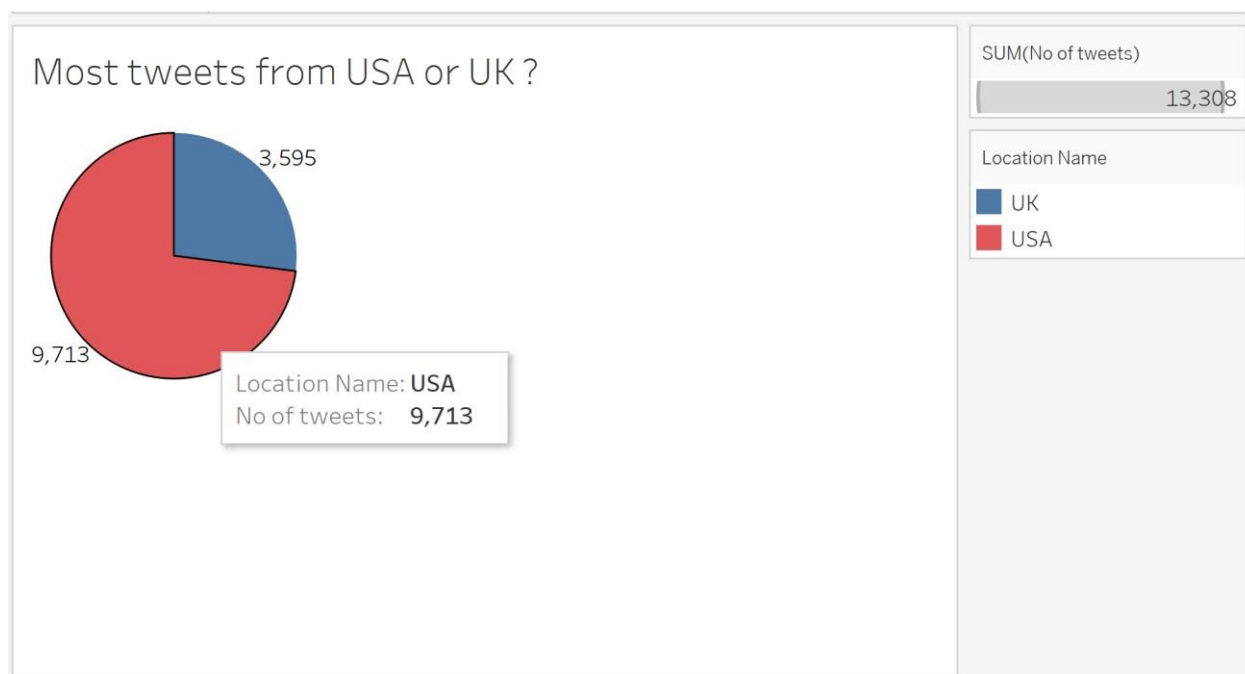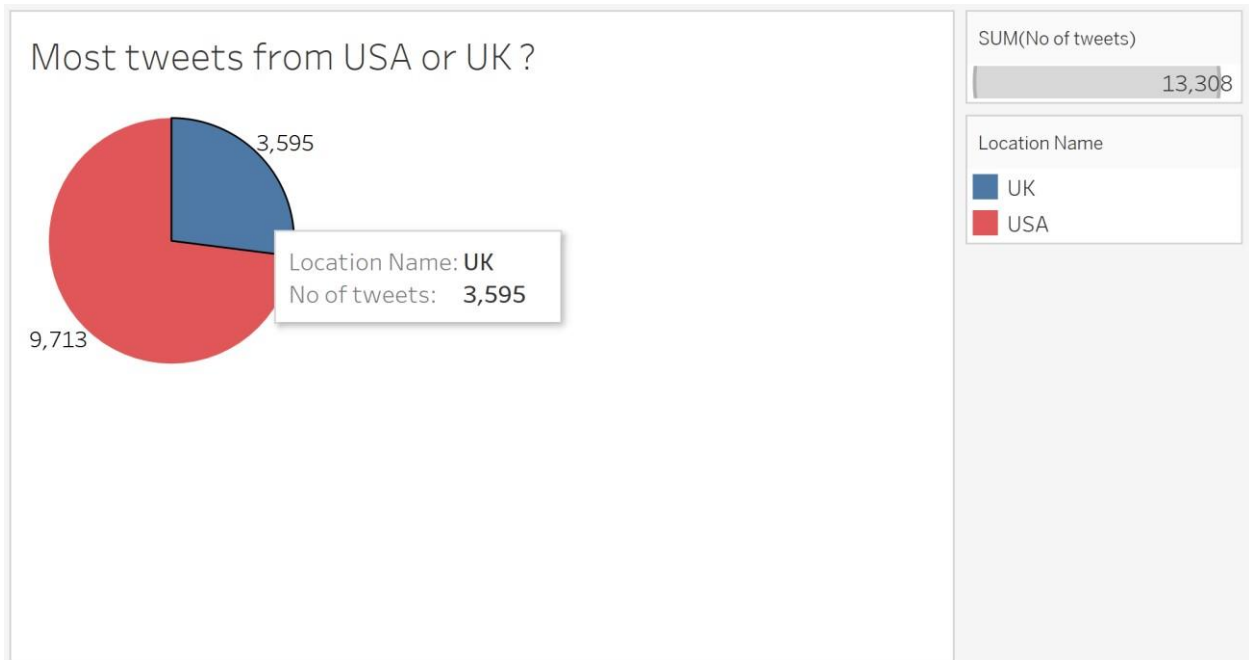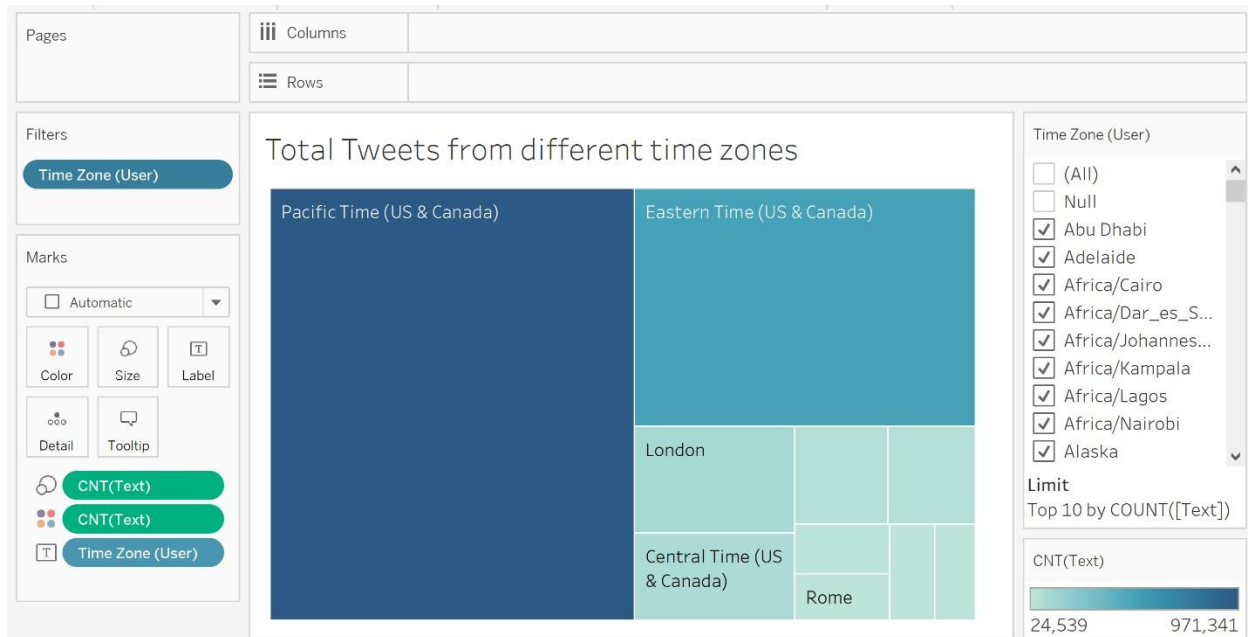
# VISUALIZATION

## QUERY 1 :



## QUERY 2 :

## Most tweets from USA or UK ?

3,595

Location Name: UK
No of tweets:    3,595

9,713

SUM(No of tweets)
                    13,308

Location Name
■ UK
■ USA

**QUERY 3 :**

Filters
Country

Marks
○ Automatic ▼

Color | Size | Label
Detail | Tooltip

Country
Country
CNT(Country)

## Top 5 Countries from where the tweets originate

© OpenStreetMap contributors

Country
■ Canada
■ India
■ Italy
■ United Kingdom
■ United States

## QUERY 4 :



## QUERY 5 :

## QUERY 6 :



## QUERY 7:

**QUERY 8 :**

## Pages

## Columns

## Rows
Verified (User)

## Filters

### Verified Users ▾

| Verified (Us.. | |
|---|---|
| Null | · |
| False | ▮ |
| True | · |

Verified (User):
Count of Verified (User): **0**

## Marks

☐ Automatic ▾

Color    Size    Label

Detail    Tooltip

Verified (User)

CNT(Verified (Us..

**CNT(Verified (User))**

| | |
|---|---|
| · | 0 |
| ☐ | 1,000,000 |
| ☐ | 2,000,000 |
| ☐ | 3,000,000 |
| ▣ | 3,930,380 |

**Verified (User)**

| ▮ | Null |
|---|---|
| ▮ | False |
| ▮ | True |

---

## QUERY 9 :

## Pages

## Columns

## Rows

## Filters
Verified (User): True

Location (User)

### Top 10 verified users in United States based on their follower's count

NWS

Name (User):       **Ronny Pascale**
Followers Count (User) per user: **45,568**

## Marks

○ Circle ▾

Color    Size    Label

Detail    Tooltip

SUM(Followers C..

Name (User)

Name (User)

**Verified (User)**

☐ (All)
☐ Null
☐ False
☑ True

**Location (User)**

United States   ✕

**Name (User)**

| ▮ | Kaiser Permanente |
|---|---|
| ▮ | Monsanto BioAg |
| | ASK USA |
| ▮ | Ronny Pascale |
| ▮ | The Wirecutter |

**QUERY 10 :**



**EXTRA VISUALIZATIONS :**

**Top 15 Hashtags**

**Top 10 Users using their Follower's count**