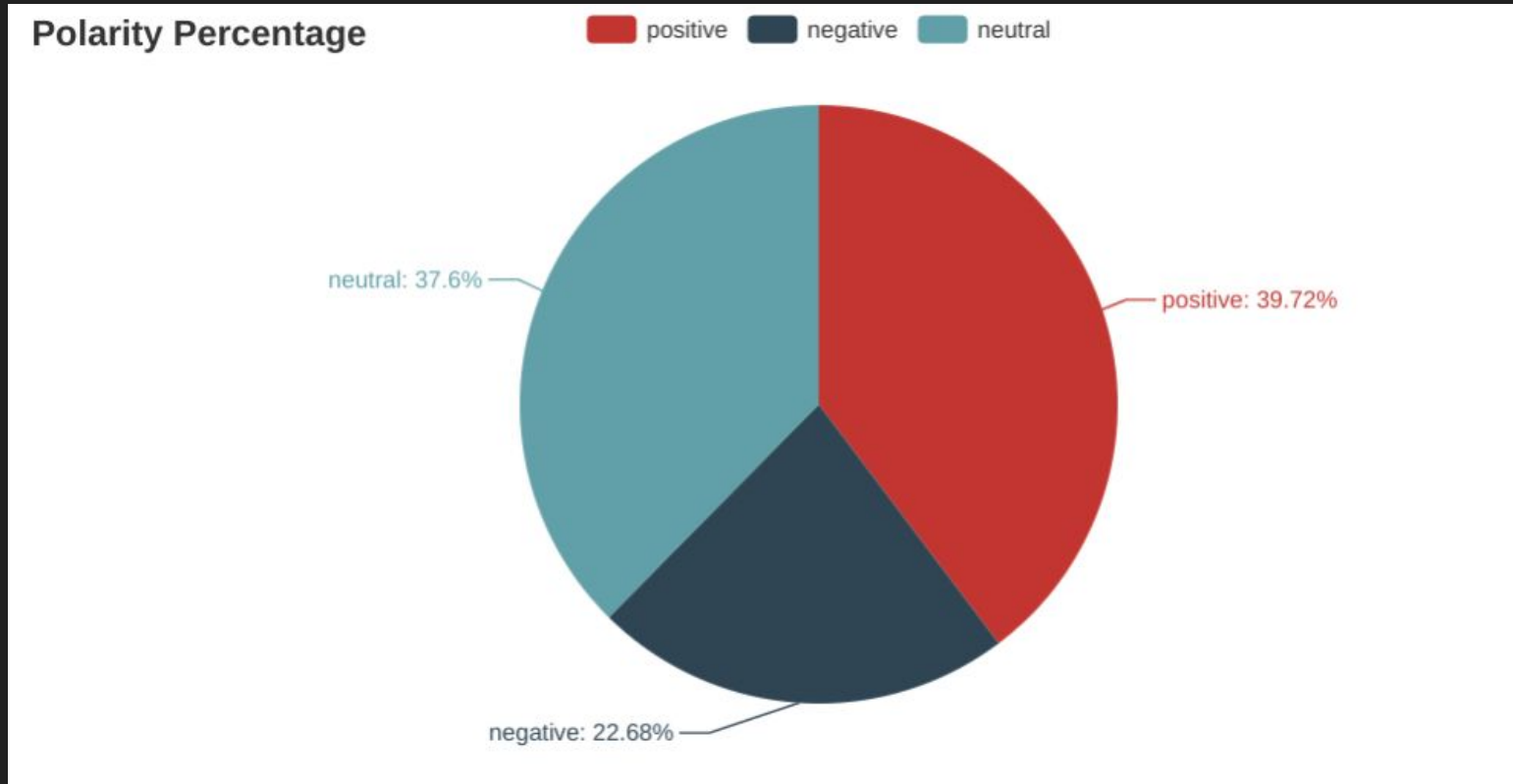# Tweets Analysis with Spark
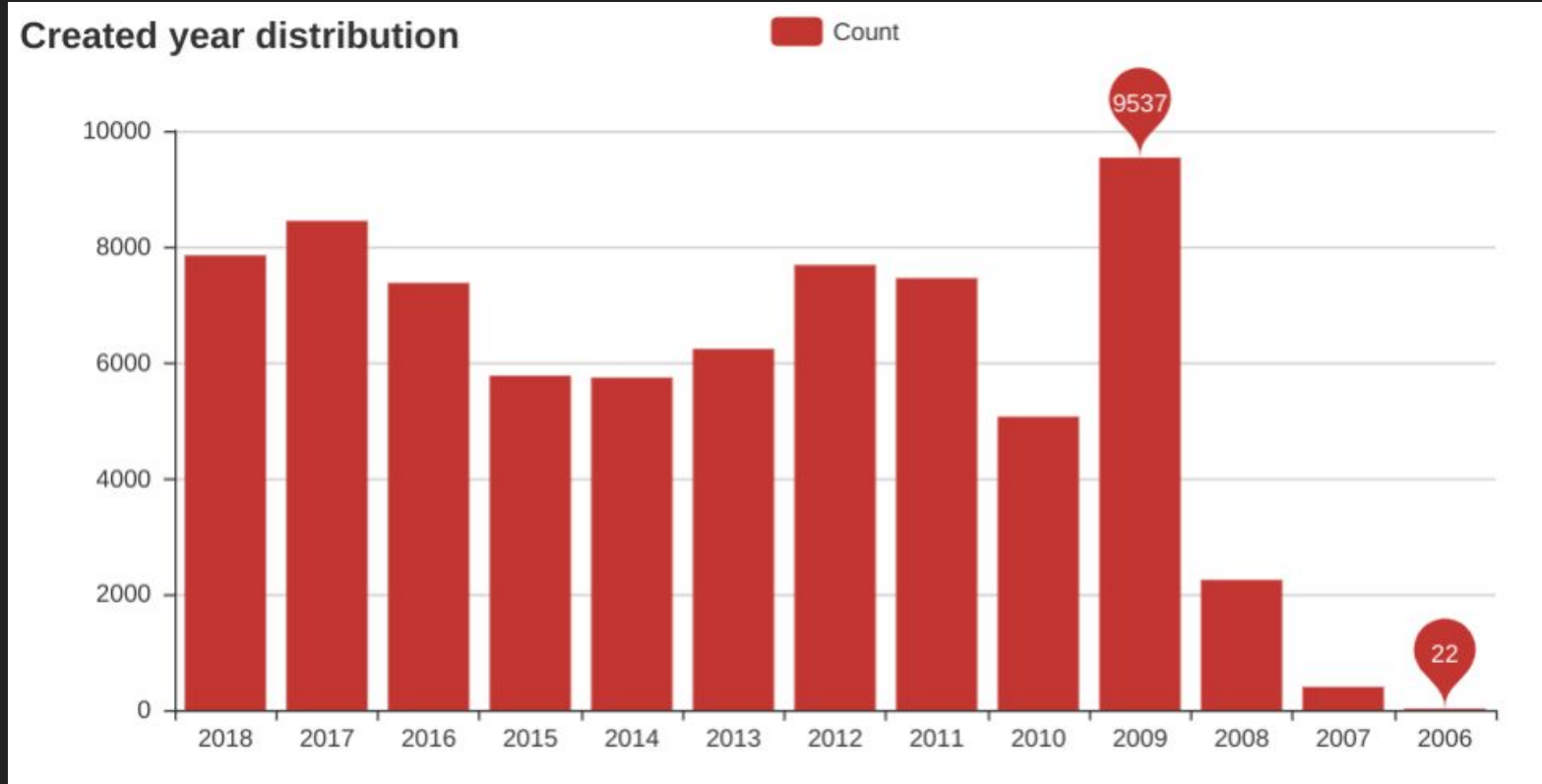
Harshil Patel (8), Matthew Boerner (1), Stephanie Retzke (10) and Yong Zheng (14)

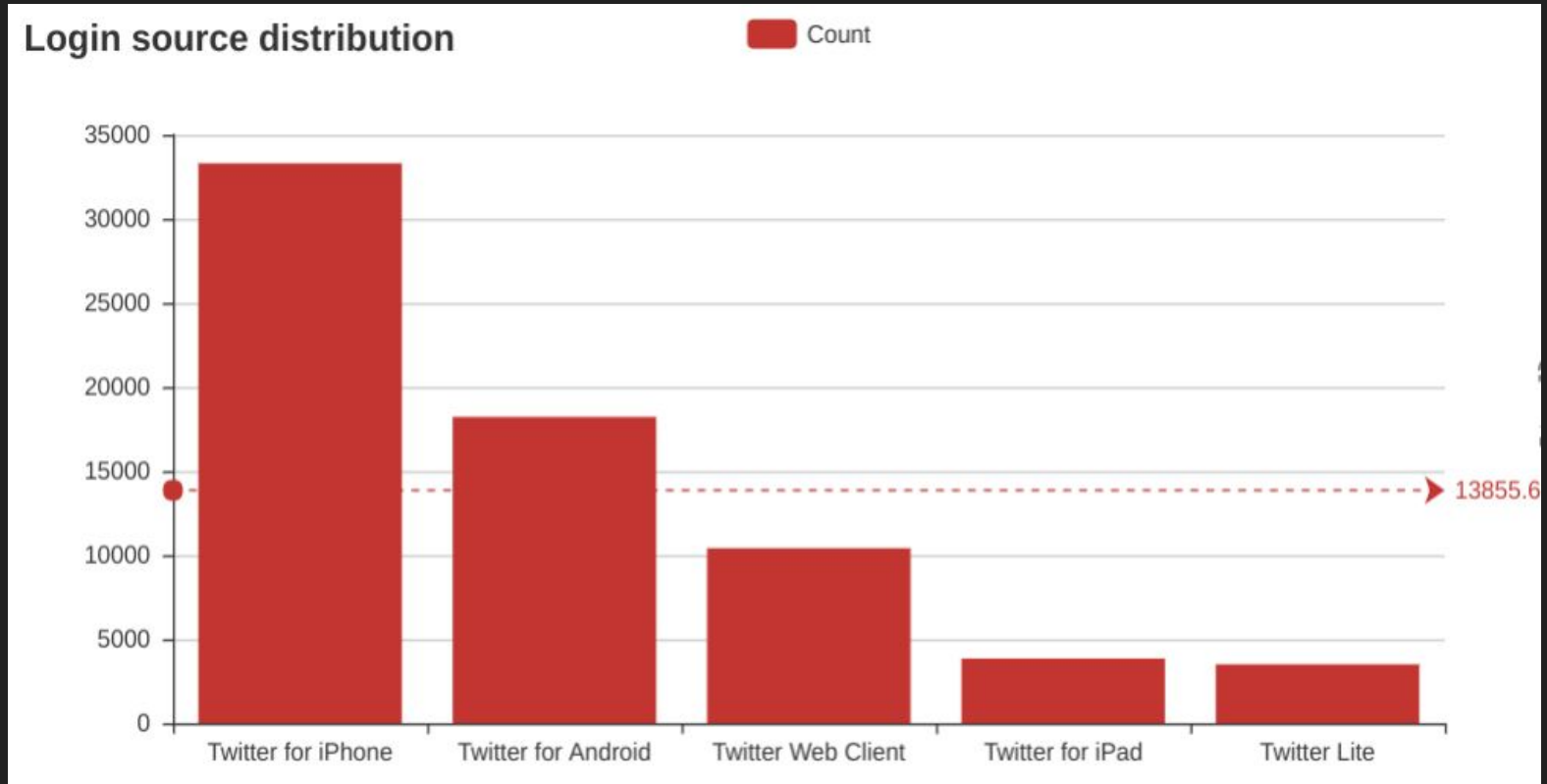Yong

# Increment 1 - Query 1 - Polarity Percentage
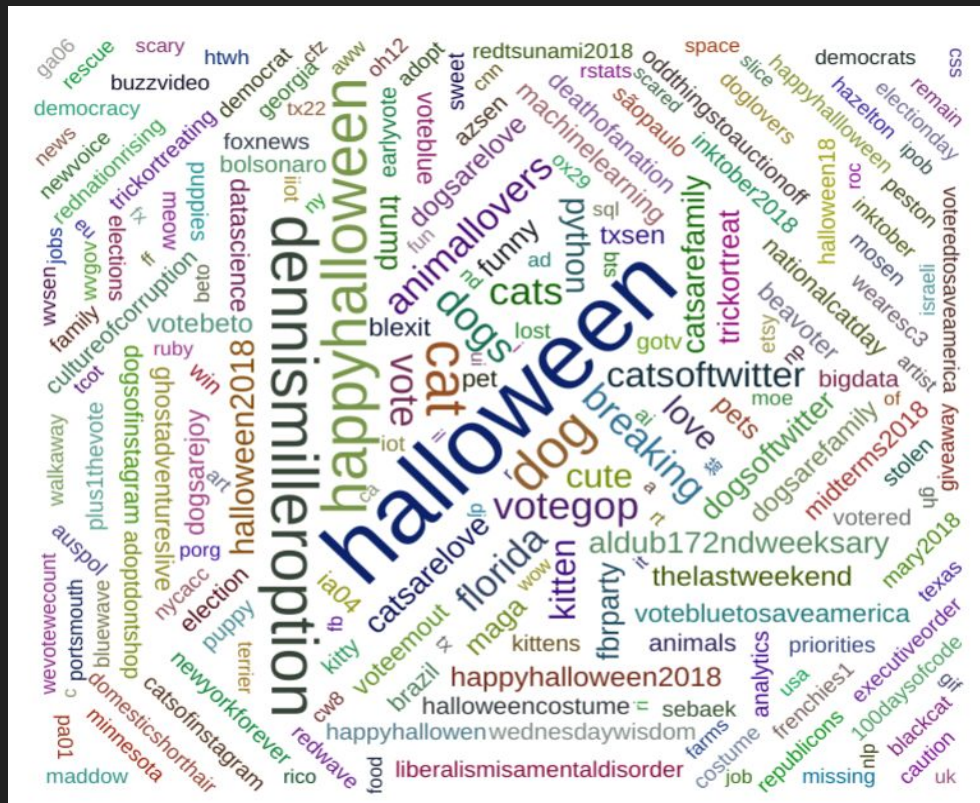
# Increment 1 - Query 2 - Created Year Distribution

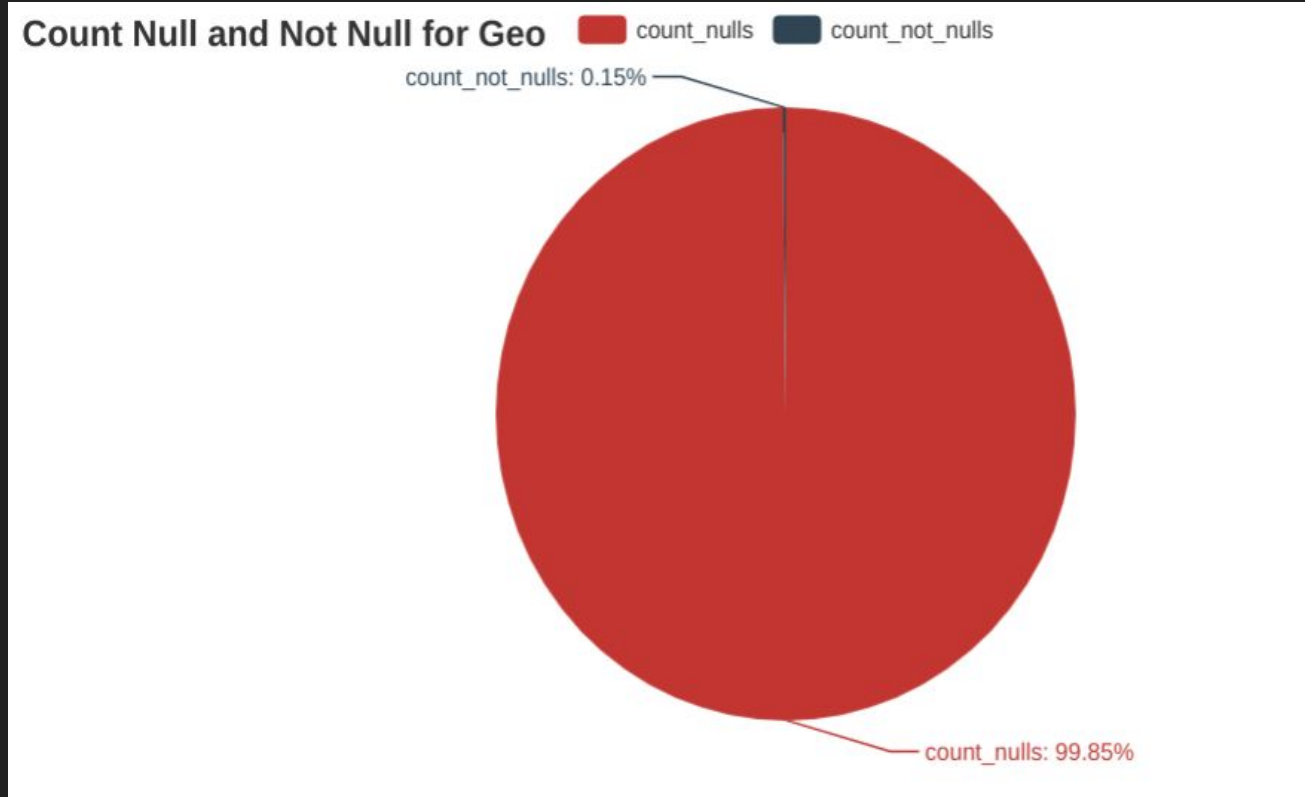# Increment 1 - Query 3 - Login Source Distribution

# Increment 1 - Query 4 - Word Cloud

# Increment 1 - Query 5 - With/Without Geo

# Increment 1 - Query 6 - Country Distribution (Geo)

| | Count |
|---|---|
| USA | 119 |
| Canada | 11 |
| New Zealand/Aotearoa | 3 |
| UK | 3 |
| Brasil | 2 |
| Philippines | 2 |
| ประเทศไทย | 1 |
| Chile | 1 |
| Ireland | 1 |
| Australia | 1 |
| 대한민국 | 1 |
| India | 1 |

# Increment 1 - Query 7 - Country Distribution (Placed)

```
+--------------+-----+
|country_code|count|
+--------------+-----+
|           US| 1147|
|           CA|   49|
|           GB|   47|
|           AU|   17|
|           PH|   10|
|           BR|    6|
|           IN|    6|
|           JP|    5|
|           NZ|    4|
|           MY|    4|
+--------------+-----+
```

# Increment 1 - Query 8 - Long Message Percentage

**Message truncated percentage**

9%

# Increment 1 - Query 9 - With/Without Verified

```
+---------+-----+
|verified|count|
+---------+-----+
|     true| 1221|
|    false|98779|
+---------+-----+
```

```
+---------------+-----+
|default_profile|count|
+---------------+-----+
|           true|  218|
|          false| 1003|
+---------------+-----+
```

Verified Accounts

```
+---------------+-----+
|default_profile|count|
+---------------+-----+
|           true|55598|
|          false|43181|
+---------------+-----+
```

Unverified Accounts

# Increment 1 - Query 10 - Friends VS. Followers

```
# Users count where the number of friends_count is greater than followers_count
data_df.where('user.friends_count > user.followers_count').count()
```

61730

```
# Users count where the number of friends_count is less than followers_count
data_df.where('user.friends_count < user.followers_count').count()
```

37952

```
# Users count where the number of friends_count is the same as followers_count
data_df.where('user.friends_count = user.followers_count').count()
```

318

# Increment 2 - ML - Features

```python
# Generate new data list for ML
ml_data_list = []
index = 0

for row in data_list:
    temp_data_list = []
    if row[0] == 1:
        temp_data_list.append('truncated')
    if row[1] == 1:
        temp_data_list.append('verified')
    if row[2] == 1:
        temp_data_list.append('geo_enabled')
    if row[3] == 1:
        temp_data_list.append('profile_background_tile')
    if row[4] == 1:
        temp_data_list.append('profile_use_background_image')
    if row[5] == 1:
        temp_data_list.append('default_profile')
    if temp_data_list:
        ml_data_list.append((index, temp_data_list))
        index += 1
```

# Increment 2 - ML - Model

```python
# Generate analysis dataframe
analysis_df = spark.createDataFrame(ml_data_list, ['id', 'items'])
```

```python
# Create the model and fit the dataframe into the model
fp_growth = FPGrowth(itemsCol="items", minSupport=0.50, minConfidence=0.6)
model = fpGrowth.fit(analysis_df)
```

```python
# Display frequent itemsets
model.freqItemsets.show(truncate=False)
```

```
+------------------------------------------------+-----+
|items                                           |freq |
+------------------------------------------------+-----+
|[profile_use_background_image]                  |80317|
|[default_profile]                               |55816|
|[default_profile, profile_use_background_image] |55816|
+------------------------------------------------+-----+
```

# Increment 2 - ML - Result

```python
# Generate analysis dataframe
analysis_df = spark.createDataFrame(ml_data_list, ['id', 'items'])
```

```python
# Create the model and fit the dataframe into the model
fp_growth = FPGrowth(itemsCol="items", minSupport=0.50, minConfidence=0.6)
model = fpGrowth.fit(analysis_df)
```

```python
# Display frequent itemsets
model.freqItemsets.show(truncate=False)
```

```
+------------------------------------------------------+-----+
|items                                                 |freq |
+------------------------------------------------------+-----+
|[profile_use_background_image]                        |80317|
|[default_profile]                                     |55816|
|[default_profile, profile_use_background_image]       |55816|
+------------------------------------------------------+-----+
```

# Increment 2 - ML - Recommendations

```
# Display prediction
model.transform(analysis_df).show(tdruncate=False)
```

```
+---+----------------------------------------------------------------+----------------+
|id |items                                                           |prediction      |
+---+----------------------------------------------------------------+----------------+
|0  |[profile_use_background_image, default_profile]                 |[]              |
|1  |[geo_enabled, profile_use_background_image, default_profile]    |[]              |
|2  |[profile_use_background_image]                                  |[default_profile]|
|3  |[truncated, profile_use_background_image, default_profile]      |[]              |
|4  |[profile_use_background_image, default_profile]                 |[]              |
|5  |[profile_use_background_image, default_profile]                 |[]              |
|6  |[profile_use_background_image, default_profile]                 |[]              |
|7  |[geo_enabled, profile_use_background_image]                     |[default_profile]|
|8  |[geo_enabled, profile_use_background_image, default_profile]    |[]              |
|9  |[profile_use_background_image, default_profile]                 |[]              |
|10 |[profile_use_background_image, default_profile]                 |[]              |
|11 |[truncated, geo_enabled, profile_use_background_image, default_profile]|[]       |
|12 |[truncated, profile_use_background_image]                       |[default_profile]|
|13 |[geo_enabled, profile_background_tile, profile_use_background_image]|[default_profile]|
|14 |[profile_use_background_image, default_profile]                 |[]              |
|15 |[profile_background_tile, profile_use_background_image]         |[default_profile]|
|16 |[profile_use_background_image, default_profile]                 |[]              |
|17 |[profile_use_background_image, default_profile]                 |[]              |
|18 |[profile_use_background_image, default_profile]                 |[]              |
|19 |[profile_use_background_image, default_profile]                 |[]              |
+---+----------------------------------------------------------------+----------------+
only showing top 20 rows
```

# Increment 2 - Graph - Graph Creation

```python
# Create edges
e = data_df.select(
    col("user.id").alias("src"),
    col("retweeted_status.user.id").alias("dst"),
    lit("retweet").alias("relationship")
).where(
    col("retweeted_status.user.id").isNotNull()
).distinct()
```

```python
# Create vertices
v = data_df.select(
    col("user.id"),
    col("user.screen_name")
).union(
    data_df.select(
        col("retweeted_status.user.id"),
        col("retweeted_status.user.name"),
    ).where(
        col("retweeted_status.user.id").isNotNull()
    )
).distinct()
```

```python
# Create graph
g = GraphFrame(v, e)
```

# Increment 2 - Graph - In-degree

```
# Get top 10 id with highest in degree
g.inDegrees.sort("inDegree", ascending=False).show(10, False)
```

```
+--------------------+--------+
|id                  |inDegree|
+--------------------+--------+
|1604444052          |2201    |
|487297085           |1979    |
|130496027           |1579    |
|4196983835          |1359    |
|100844048791572O708 |1319    |
|16989178            |1110    |
|19697415            |1070    |
|150078976           |1056    |
|2828212668          |984     |
|3267456386          |930     |
+--------------------+--------+
only showing top 10 rows
```

# Increment 2 - Graph - In-degree (sorted)

```
# Get top 10 id with lowest in degree
g.inDegrees.sort("inDegree").show(10, False)
```

```
+------------------+--------+
|id                |inDegree|
+------------------+--------+
|3358687222        |1       |
|18611344          |1       |
|2859622263        |1       |
|20813564          |1       |
|179176923         |1       |
|948171093818343425|1       |
|540321025         |1       |
|72954856          |1       |
|95755482          |1       |
|135138678         |1       |
+------------------+--------+
only showing top 10 rows
```

# Increment 2 - Graph - Out-degree

```
# Get top 10 id with highest out degree
g.outDegrees.sort("outDegree", ascending=False).show(10, False)
```

```
+--------------------+---------+
|id                  |outDegree|
+--------------------+---------+
|988374538491777025  |48       |
|1053011875007483905 |24       |
|824216698551209984  |18       |
|822210115784806400  |18       |
|80602426            |17       |
|2905278738          |16       |
|62147616            |16       |
|4053477192          |16       |
|954454874665771013  |16       |
|2874358611          |15       |
+--------------------+---------+
only showing top 10 rows
```

# Increment 2 - Graph - Shortest Paths  - Less Popular Landmarks

```python
# Compute shortest paths from each vertex to the given set of landmark vertices
# For one id with low number of in degree
g.shortestPaths(
    landmarks=["3358687222"]
).select("id", "distances").where(
    size(col("distances")) > 0
).show(10, False)
```

```
+------------------+------------------+
|id                |distances         |
+------------------+------------------+
|885245073813798912|[3358687222 -> 1]|
|3358687222        |[3358687222 -> 0]|
+------------------+------------------+
```

# Increment 2 - Graph - Shortest Paths - Popular Landmarks

```python
# Compute shortest paths from each vertex to the given set of landmark vertices
# For one id with high number of in degree
g.shortestPaths(
    landmarks=["150078976"]
).select("id", "distances").where(
    size(col("distances")) > 0
).show(10, False)
```

```
+------------------+----------------+
|id                |distances       |
+------------------+----------------+
|740005760801726465|[150078976 -> 1]|
|334776400         |[150078976 -> 1]|
|849374799868637185|[150078976 -> 1]|
|935707527580397569|[150078976 -> 1]|
|3424129696        |[150078976 -> 1]|
|827544166624325632|[150078976 -> 1]|
|303950400         |[150078976 -> 1]|
|498077600         |[150078976 -> 2]|
|241941600         |[150078976 -> 1]|
|1566650000        |[150078976 -> 1]|
+------------------+----------------+
only showing top 10 rows
```

# Increment 2 - Graph - Shortest Paths - Triangle Count

```
# Computes the number of triangles passing through each vertex
g.triangleCount().select("id", "count").where(col("count") > 0).show(10, False)

+-------------------+-----+
|id                 |count|
+-------------------+-----+
|908149255654854656 |104  |
|883855148136763392 |4    |
|883855148136763392 |4    |
|1003631763133001729|4    |
|719662077934243840 |12   |
|197111814          |4    |
|862106769862078464 |2    |
|909519301299892225 |40   |
|909519301299892225 |40   |
|3100613023         |4    |
+-------------------+-----+
only showing top 10 rows
```

Harshil

# Model



python script to download tweet

Pyspark ML

model creation

tweet analysis

100K tweet in json format

# Get top 10 user which have most favourites count

```
C:\Users\harsh\AppData\Local\Programs\Python\Python36-32\pyth
-------------favourites_count-----
meh
1073419
Brian D.
965279
EnigMAA
851456
Madana Bhat-Khandige
823206
Monica Cates
821897
DerekPlatt
815322
rebecca lauren
812977
Brennen Burleson
800649
Jeanette Baratta
799905
Grand Moff Snarkin, Surefire Intelligence CFO
772169
```

# Find the tracks in the tweet and shows its distribution

# Show tweets generated from united state and other than united states

# Get the count of the user which have account verified

```
--------account verified-----------
1319
```

Get top 10 tweet that are retweeted the most

```
-------tweet that are retweeted the most-------------
Rylie Geraci
633089
c h i e f
178413
BAYU ARISANDY
125011
monika bielskyte
125010
Rihanna
125009
Nasir Shakur
125007
shay
125007
RichFanAcc
118884
SaltSaltSalt
110527
rory miller
107464
```

# Get top 10 user which have most followers count

```
------followers_count-------
Donald J. Trump
55525198
President Trump
24414561
The Economist
23428140
Reuters Top News
19957413
The White House
17502570
The Washington Post
13010477
The Washington Post
13010416
China Xinhua News
11550188
Jimmy Kimmel
11427575
HuffPost
11389165
---------------------------
```

# Get top 10 user which have most friends count

```
-------friends_count---------
Ed Krassenstein
641614
Travel
572779
Jeffrey Levin
428351
Jeffrey Levin
428351
Music Lovers Fans♫
306347
Ⓜ MonsterFunder
305853
Tina Stull
236552
Reg Saddler
229134
Reg Saddler
229134
Dorian Sage ©
227725
--------------------------
```

# Get percentage of tweet based on negative, positive, neutral

```
--------negative postive tweet-------------
negative:27.449725502744972
positive:34.2296577034229666
neutral:38.31961680383196
```

# Find gender of the user based on their name

# Get the count of tweets

```
------count-----------
100000
```

# HashingTF + IDF + Logistic Regression

```
+---+-------------------+------+-------------------+-------------------+-------------------+-----+
|_c0|               text|target|              words|                 tf|           features|label|
+---+-------------------+------+-------------------+-------------------+-------------------+-----+
|  0|RT Trump�s press ...|     1|[rt, trump�s, pre...|(65536,[312,7752,...|(65536,[312,7752,...|  0.0|
|  1|RT What is tweeti...|     0|[rt, what, is, tw...|(65536,[12716,158...|(65536,[12716,158...|  1.0|
|  2|RT Go to a hand r...|     0|[rt, go, to, a, h...|(65536,[1038,1354...|(65536,[1038,1354...|  1.0|
|  3|RT I�m moving to ...|     1|[rt, i�m, moving,...|(65536,[1197,8436...|(65536,[1197,8436...|  0.0|
|  4|RT Little-known f...|     1|[rt, little-known...|(65536,[19996,290...|(65536,[19996,290...|  0.0|
|  5|RT FACT funded th...|     0|[rt, fact, funded...|(65536,[20464,217...|(65536,[20464,217...|  1.0|
|  6|RT Trump's firing...|     1|[rt, trump's, fir...|(65536,[9639,1553...|(65536,[9639,1553...|  0.0|
|  7|RT _hooty _regula...|     1|[rt, _hooty, _reg...|(65536,[3811,7612...|(65536,[3811,7612...|  0.0|
+---+-------------------+------+-------------------+-------------------+-------------------+-----+
only showing top 8 rows

18/11/28 19:51:04 WARN BLAS: Failed to load implementation from: com.github.fommil.netlib.NativeSystemBLAS
18/11/28 19:51:04 WARN BLAS: Failed to load implementation from: com.github.fommil.netlib.NativeRefBLAS
0.9415204678362573
```

## CountVectorizer + IDF + Logistic Regression

```
Accuracy Score: 0.9181
ROC-AUC: 0.9544

Process finished with exit code 0
```

# Graph

```
g.degrees.show(5)
+------------------+------+
|                id|degree|
+------------------+------+
|         133938408|    13|
|        2896675593|     1|
|823651149823688705|     4|
|883855148136763392|     2|
|818632400641097733|     1|
+------------------+------+
only showing top 5 rows
```

```
+--------------------+--------------------+------------+
|                 src|                 dst|relationship|
+--------------------+--------------------+------------+
|  963967475875500033|           395674143|     retweet|
|           320908938| 1034665733005959168|     retweet|
|  789588266198700036|           232901331|     retweet|
|           385689121|          4731701624|     retweet|
|  879531117811965953|  988960980548931585|     retweet|
+--------------------+--------------------+------------+
only showing top 5 rows
```

```
+------------------+-------------+
|               id|  screen_name|
+------------------+-------------+
|         210185175|  emmaburnsapp|
|         822926826|     kjtc1979|
|943892654160531457|    aesopgalt|
|767072603278159872|m8nkey2chm00n|
|         215441260|      LPWhitt|
+------------------+-------------+
only showing top 5 rows
```

```
g.inDegrees.sort("inDegree").show(5, False)
```

```
+----------------------------+---------+
|id                          |inDegree |
+----------------------------+---------+
|15469000                    |1        |
|577432615                   |1        |
|864068108                   |1        |
|10385129071120779264        |1        |
|53190110                    |1        |
+----------------------------+---------+
only showing top 5 rows
```

```
g.inDegrees.sort("inDegree", ascending=False).show(5, False)
```

```
+----------------------------------+--------+
|id                                |inDegree|
+----------------------------------+--------+
|25073877                          |105     |
|21728303                          |44      |
|822215673812119553                |34      |
|18643437                          |32      |
|91386979                          |31      |
+----------------------------------+--------+
only showing top 5 rows
```

```
g.shortestPaths(landmarks=["25073877"]).select("id", "distances").where(size(col("distances")) > 0).show(10, False)
```

```
+--------------------+----------------------+
|id                  |distances             |
+--------------------+----------------------+
|95460995308755356   |Map(25073877 -> 1)|
|733279356492079106  |Map(25073877 -> 1)|
|877879587493085184  |Map(25073877 -> 1)|
|43870406            |Map(25073877 -> 1)|
|977695071096262658  |Map(25073877 -> 1)|
|2880381305          |Map(25073877 -> 1)|
|2370953706          |Map(25073877 -> 1)|
|976262602421448705  |Map(25073877 -> 1)|
|227835612           |Map(25073877 -> 1)|
|2149493108          |Map(25073877 -> 1)|
+--------------------+----------------------+
only showing top 10 rows
```

Stephanie

# Increment 1

Query 1:  Count how many tweets were favorited and those that were not

```
//1
val favTweets = sqlContext
  .sql( sqlText = "SELECT favorited, count(*) as count FROM EntertainmentTable where favorited is not null  group by favorited order by count desc limit 10"
favTweets.show
```

```
+---------+-----+
|favorited|count|
+---------+-----+
|    false| 1000|
+---------+-----+
```

# Increment 1

Query 2:  Count how many tweets were retweeted and those that were not

```
//2
val reTweets = sqlContext
  .sql( sqlText = "SELECT retweeted, count(*) as count FROM EntertainmentTable where retweeted is not null  group by retweeted order by count desc limit 10")
reTweets.show
```

```
+----------+-----+
|retweeted|count|
+----------+-----+
|     false| 1000|
+----------+-----+
```

# Increment 1

Query 3:  Count how many tweets were quotes and those that were not

```
//3
val quoteTweets = sqlContext
  .sql( sqlText = "SELECT is_quote_status, count(*) as count FROM EntertainmentTable where is_quote_status is not null     group by is_quote_status order by count
quoteTweets.show
```

```
+----------------+-----+
|is_quote_status|count|
+----------------+-----+
|           false|  802|
|            true|  198|
+----------------+-----+
```

# Increment 1

Query 4:  See what filter levels people use and how many

```
//4
val filterTweets = sqlContext
  .sql( sqlText = "SELECT filter_level, count(*) as count FROM EntertainmentTable where filter_level is not null group by filter_level order by count desc limit
filterTweets.show
```

```
+------------+-----+
|filter_level|count|
+------------+-----+
|         low| 1000|
+------------+-----+
```

# Increment 1

Query 5:  Count how many tweets were truncated and those that were not

```
//5
val truncatedTweets = sqlContext
  .sql( sqlText = "SELECT truncated, count(*) as count FROM EntertainmentTable where truncated is not null group by truncated order by count desc limit 10")
truncatedTweets.show
```

```
+---------+-----+
|truncated|count|
+---------+-----+
|    false|  905|
|     true|   95|
+---------+-----+
```

# Increment 1

Query 6:  Count how many tweets were quotes and truncated

```
//6
val tqTweets = sqlContext
  .sql( sqlText = "SELECT id, truncated, is_quote_status, count(*) as count FROM EntertainmentTable where truncated = true AND is_quote_status = true  group by id
tqTweets.show
```

```
+-------------------+---------+---------------+-----+
|                 id|truncated|is_quote_status|count|
+-------------------+---------+---------------+-----+
|1057777525441679361|     true|           true|    1|
|1057777569389436928|     true|           true|    1|
|1057772296055119872|     true|           true|    1|
|1057777564482240516|     true|           true|    1|
|1057773221135191553|     true|           true|    1|
|1057777304359686144|     true|           true|    1|
|1057777353340928000|     true|           true|    1|
|1057777303285989376|     true|           true|    1|
+-------------------+---------+---------------+-----+
```

# Increment 1

Query 7: Count how many tweets were created at the same time

```
//7
val createdTweets = sqlContext
  .sql( sqlText = "SELECT created_at, count(*) as count FROM EntertainmentTable where created_at is not null group by created_at order by count desc limit 15")
createdTweets.show
```

```
+--------------------+-----+
|          created_at|count|
+--------------------+-----+
|Wed Oct 31 23:32:...|   24|
|Wed Oct 31 23:33:...|   23|
|Wed Oct 31 23:33:...|   21|
|Wed Oct 31 23:32:...|   21|
|Wed Oct 31 23:33:...|   20|
|Wed Oct 31 23:32:...|   20|
|Wed Oct 31 23:32:...|   19|
|Wed Oct 31 23:33:...|   19|
|Wed Oct 31 23:32:...|   19|
|Wed Oct 31 23:32:...|   18|
|Wed Oct 31 23:32:...|   18|
|Wed Oct 31 23:32:...|   18|
|Wed Oct 31 23:32:...|   18|
|Wed Oct 31 23:32:...|   17|
|Wed Oct 31 23:33:...|   17|
+--------------------+-----+
```

# Increment 1

Query 8:  Find the max replies a tweet had at a certain time

```
val replyCountTweets = sqlContext
  .sql( sqlText = "SELECT created_at, max(reply_count) as max_reply FROM EntertainmentTable where created_at is not null group by created_at order by max_reply de
replyCountTweets.show
```

```
+--------------------+---------+
|          created_at|max_reply|
+--------------------+---------+
|Wed Oct 31 23:32:...|        0|
|Wed Oct 31 23:32:...|        0|
|Wed Oct 31 23:32:...|        0|
|Wed Oct 31 23:33:...|        0|
|Wed Oct 31 23:32:...|        0|
|Wed Oct 31 23:32:...|        0|
|Wed Oct 31 23:32:...|        0|
|Wed Oct 31 23:32:...|        0|
|Wed Oct 31 23:32:...|        0|
|Wed Oct 31 23:33:...|        0|
|Wed Oct 31 23:33:...|        0|
|Wed Oct 31 23:32:...|        0|
|Wed Oct 31 23:32:...|        0|
|Wed Oct 31 23:32:...|        0|
|Wed Oct 31 23:33:...|        0|
|Wed Oct 31 23:32:...|        0|
+--------------------+---------+
```

# Increment 1

Query 9:  Find the max quotes a tweet had at a certain time

```
val quoteCountTweets = sqlContext
    .sql( sqlText = "SELECT created_at, max(quote_count) as max_quote FROM EntertainmentTable where created_at is not null group by created_at order by max_quote d
quoteCountTweets.show
```

```
+--------------------+---------+
|          created_at|max_quote|
+--------------------+---------+
|Wed Oct 31 23:32:...|        0|
|Wed Oct 31 23:32:...|        0|
|Wed Oct 31 23:32:...|        0|
|Wed Oct 31 23:32:...|        0|
|Wed Oct 31 23:33:...|        0|
|Wed Oct 31 23:32:...|        0|
|Wed Oct 31 23:32:...|        0|
|Wed Oct 31 23:32:...|        0|
|Wed Oct 31 23:32:...|        0|
|Wed Oct 31 23:33:...|        0|
|Wed Oct 31 23:33:...|        0|
|Wed Oct 31 23:32:...|        0|
|Wed Oct 31 23:32:...|        0|
|Wed Oct 31 23:33:...|        0|
|Wed Oct 31 23:32:...|        0|
+--------------------+---------+
```

# Increment 1

Query 10:  Count how many users are protected and those that are not

```
//10
val protectedTweets = sqlContext
  .sql( sqlText = "SELECT user.protected, count(*) as count FROM EntertainmentTable where user.protected = false group by user.protected order by count desc limi
protectedTweets.show
```

```
+---------+-----+
|protected|count|
+---------+-----+
|    false| 1000|
+---------+-----+
```

# Increment 2 - Machine Learning

```scala
//We'll split the set into training and test data
val Array(trainingData, testData) = child.randomSplit(Array(0.8, 0.2))

val labelColumn = "id"

//We define two StringIndexers for the categorical variables

val countryIndexer = new StringIndexer()
  .setInputCol("lang")
  .setOutputCol("replyIndex")

//We define the assembler to collect the columns into a new column with a single vector - "features"
val assembler = new VectorAssembler()
  .setInputCols(Array("reply_count", "replyIndex"))
  .setOutputCol("features")

//For the regression we'll use the Gradient-boosted tree estimator
val gbt = new GBTRegressor()
  .setLabelCol(labelColumn)
  .setFeaturesCol("features")
  .setPredictionCol("Predicted " + labelColumn)
  .setMaxIter(50).setMaxBins(100)
```

# Increment 2 - Machine Learning

```scala
//Construct the pipeline
val pipeline = new Pipeline().setStages(stages)

//We fit our DataFrame into the pipeline to generate a model
val model = pipeline.fit(trainingData)

//We'll make predictions using the model and the test data
val predictions = model.transform(testData)
predictions.show()
```

```
-+-----------------+--------------------+---------------+------------+------------+--------------------+----+----------+---------+--------------------+
t|               id|in_reply_to_screen_name|is_quote_status|reply_count|retweeted|                text|lang|replyIndex| features|        Predicted id|
-+-----------------+--------------------+---------------+------------+------------+--------------------+----+----------+---------+--------------------+
0|1057777287033171975|       VonnieCalland|          false|          0|     false|@VonnieCalland Th...|  en|       0.0|(2,[],[])|1.057777426920863...|
0|1057777351533228032|        molly_knight|          false|          0|     false|@molly_knight My ...|  en|       0.0|(2,[],[])|1.057777426920863...|
0|1057777359221374982|       WhenWeAllVote|          false|          0|     false|@WhenWeAllVote @M...|  en|       0.0|(2,[],[])|1.057777426920863...|
0|1057777379047690240|     realDonaldTrump|          false|          0|     false|@realDonaldTrump ...|  en|       0.0|(2,[],[])|1.057777426920863...|
0|1057777436409163776|             Twitter|           true|          0|     false|@Twitter u r not ...|  en|       0.0|(2,[],[])|1.057777426920863...|
0|1057777441501052928|          littlefonty|          false|          0|     false|@littlefonty @tay...|  en|       0.0|(2,[],[])|1.057777426920863...|
0|1057777461021290497|         JulieAnnLily|          false|          0|     false|@JulieAnnLily @Le...|  en|       0.0|(2,[],[])|1.057777426920863...|
0|1057777534924918786|         Need2Impeach|          false|          0|     false|@Need2Impeach @Da...|  en|       0.0|(2,[],[])|1.057777426920863...|
0|1057777548827551744|         Attractivepup|          false|          0|     false|@Attractivepup He...|  en|       0.0|(2,[],[])|1.057777426920863...|
0|1057777553312632832|             tomezine|          false|          0|     false|@tomezine @Verita...|  en|       0.0|(2,[],[])|1.057777426920863...|
0|1057777572908433408|          KevinMKruse|          false|          0|     false|@KevinMKruse My f...|  en|       0.0|(2,[],[])|1.057777426920863...|
-+-----------------+--------------------+---------------+------------+------------+--------------------+----+----------+---------+--------------------+
```

# Increment 2 - Machine Learning

```scala
//This will evaluate the error/deviation of the regression using the Root Mean Squared deviation
val evaluator = new RegressionEvaluator()
  .setLabelCol(labelColumn)
  .setPredictionCol("Predicted " + labelColumn)
  .setMetricName("rmse")

//We compute the error using the evaluator
val error = evaluator.evaluate(predictions)


println("The Root Mean Square Deviation error: " + error + "\n")
```

```
The Root Mean Square Deviation error: 9.363686403376337E10
```

# Increment 2 - GraphFrame

```scala
val verticesTweets = df
  .select( col = "id", cols = "user.screen_name" )
  .where( conditionExpr = "user.screen_name is not null " +
    "and id is not null " )
```

```scala
val edgesPrototype = df
  .select( col = ("id"), cols = "in_reply_to_screen_name", "in_reply_to_user_id" )
  .where( conditionExpr = "in_reply_to_screen_name is not null " +
    "and in_reply_to_user_id is not null " +
    "and id is not null " )

val e = edgesPrototype.withColumnRenamed( existingName = "id", newName = "src")
val ed = e.withColumnRenamed( existingName = "in_reply_to_screen_name", newName = "dst")
val edgesTweets = ed.withColumnRenamed( existingName = "in_reply_to_user_id", newName = "relationship")
```

```scala
val tweetGraph = GraphFrame(verticesTweets,edgesTweets)
```

# Increment 2 - GraphFrame

```
val triCount = tweetGraph.triangleCount.run()
triCount.select( col = "id", cols = "count").show()
```

```
+-------------------+-----+
|                 id|count|
+-------------------+-----+
|1057777280309706752|    0|
|1057777336354041856|    0|
|1057777348135829504|    0|
|1057777441551409152|    0|
|1057777569389436928|    0|
|1057777361406455809|    0|
|1057777400145174529|    0|
|1057777419015208960|    0|
|1057777480717819904|    0|
|1057777283572805632|    0|
|1057777324941369344|    0|
|1057777335825506305|    0|
|1057777386945691648|    0|
|1057777468797583360|    0|
|1057777536250400768|    0|
|1057777562359906304|    0|
|1057777567019778048|    0|
|1057777280079065088|    0|
|1057777357656899584|    0|
|1057777409607589888|    0|
+-------------------+-----+
```

# Increment 2 - GraphFrame

```
println("Out Degrees: ")
tweetGraph.outDegrees.sort( sortCol = "outDegree").show()
```

```
+-------------------+---------+
|                 id|outDegree|
+-------------------+---------+
|1057777441551409152|        1|
|1057777419015208960|        1|
|1057777283572805632|        1|
|1057777338564440064|        1|
|1057777395900526592|        1|
|1057777384663977984|        1|
|1057777482651201536|        1|
|1057777436409163776|        1|
|1057777401244106753|        1|
|1057777461021290497|        1|
|1057777467455234048|        1|
|1057777286970142720|        1|
|1057777288433958912|        1|
|1057777372068483074|        1|
|1057777557553250306|        1|
|1057777304359686144|        1|
|1057777546421587968|        1|
|1057777543674322946|        1|
|1057777379865755648|        1|
|1057777351533228032|        1|
+-------------------+---------+
```

# Increment 2 - GraphFrame

```
println("In Degrees: ")
tweetGraph.inDegrees.sort( sortCol = "inDegree").show()
```

```
+--------------+--------+
|            id|inDegree|
+--------------+--------+
|      Annaleen|       1|
|       adnilxa|       1|
|    GraceOM1967|      1|
|     kzannarbor|      1|
|   FlatEarthGang|     1|
|      GROGParty|      1|
|   MakedaIsRight|     1|
|    TrollTerrific|    1|
|     sallyacb275|     1|
|schneiderleonid|      1|
|       LouDobbs|      1|
|jjbittenbinders|      1|
|       natvanlis|      1|
|        AyyBates|      1|
|    WhenWeAllVote|    1|
|       _BenMonroe_|    1|
|   grantwarkentin|    1|
|    SoulSolaris23|     1|
|    molly_knight|      1|
|   HawaiianTrash_|     1|
+--------------+--------+
```

Shaun

# Increment 1 - Query 1 - Favorites All Tweet

# Increment 1 - Query 2 - Tweets Containing Trump

# Increment 1 - Query 3 - Tweets from Northern Hemosphere

# Increment 1 - Query 4 - Tweets from Southern Hemisphere

# Increment 1 - Query 5 - Contain "the"

# Increment 1 - Query 6 - Average Lat.

# Increment 1 - Query 7 - Average Longitude

# Increment 1 - Query 8 - Contain "User" and "Generated"

# Increment 1 - Query 9 - Total Retweets

# Increment 1 - Query 10 - Longest Tweet

# Increment 2 - Code

```
val tweets = spark.read.format( source = "csv").option("header", "true").load( path = "C:\\Users\\calcalocalo\\Documents\\food
tweets.createOrReplaceTempView( viewName = "tweet")
val v = spark.sql( sqlText = "select id, text, retweetCount from tweet where replyToSID <> 'NA'")
val e = spark.sql( sqlText = "select id as src, replyToSID as dst from tweet where replyToSID <> 'NA'")

val g = GraphFrame(v, e)
g.vertices.show()
g.edges.show()
```

# Increment 2 - Edges

```
+------------------+------------------+------------+
|                id|              text|retweetCount|
+------------------+------------------+------------+
|                NA|             FALSE|        TRUE|
|                NA|             FALSE|       FALSE|
|                NA|             FALSE|       FALSE|
|                NA|             FALSE|       FALSE|
|                NA|             FALSE|        TRUE|
|                NA|             FALSE|       FALSE|
|                NA|             FALSE|        TRUE|
|                NA|             FALSE|       FALSE|
|1067919125257818112|@nova_meat @Moeda...|           0|
|                NA|             FALSE|        TRUE|
|                NA|             FALSE|       FALSE|
|                NA|             FALSE|       FALSE|
|                NA|             FALSE|       FALSE|
|                NA|             FALSE|        TRUE|
|                NA|"""But #Daniel #r...|firstseekHim|
|                NA|             FALSE|       FALSE|
|                NA|             FALSE|        TRUE|
|                NA|             FALSE|       FALSE|
|                NA|             FALSE|        TRUE|
|                NA|             FALSE|        TRUE|
+------------------+------------------+------------+
only showing top 20 rows
```

# Increment 2 - Vertices