

Data Challenge

Instructions

You will find a database of a representative sample of search and share events for the month of June for GIF Keyboard (US-only). Please analyze this data, provide written responses to questions 1-4 below, and prepare to discuss the data during the on-site interviews.

1. "Warmup" metrics:

- Compute the aggregate "share rate": Ratio of shares to searches
- Compute the share rate for the top 1000 most-searched terms
- How many unique search terms are there? What does the distribution look like (consider the count per search term; plot or describe the results)?
- How many unique gifs are shared? What does the distribution look like?

2. Advanced metrics:

- A "search session" is defined as a series of consecutive searches by a user, close in temporal proximity (typically on the order of seconds), that may or may not result in a share. Given the available data, how would you compute the average duration of a search session? Bonus: Estimate the average duration of a search session, and compute the ratio of shares to search sessions.

3. Experimental design (hint: A/B testing):

- Suppose you've built a new model for the search engine that you'd like to test against the current model.
 - * How would you measure the effectiveness of the new model?
 - * What additional data would you need in the dataset to compare the effectiveness of the new vs. current models?
 - * How would you determine the statistical significance of the results?

4. Predictive models:

- There is sufficient data in the dataset to build a crude "search suggestions" model: Given a specific search term, return a list of related search terms. Describe how you might build such a model (describe the data you would use, techniques, inferences, etc.; you do not need to implement this).

5. Insights and discussion points for on-site:

- This section is not needed for the written responses. For the on-site interviews, please take another look at the dataset, and be prepared to share and discuss additional insights you can extract from the dataset, or useful models or features you can potentially build from the dataset. Example topics might include semantic analysis of searches, graph analysis of shared GIFs, etc.

Database Info:

You will have read-write access to the "public" schema. You may create temp tables as necessary. The table with relevant data is "ios_events". The data may be noisy, but we tried to limit the dataset to the most relevant data.

Notes on relevant columns:

timestamp - time of user action

keyboardid - unique anonymous user id

eventname - textsearch = browse or search; share = share

tags - primary search term(s)

riffid - id of the gif shared (gif metadata is not provided)

tag1/tag2 - may be relevant if autosuggest or filters are used, but this data is noisy

actions - app behavior triggering the search

component - sub-component of app

viewindex - position of the gif shared (or -1 if unknown/unavailable)

In [1]:

```
1 import sys
2 import json
3 import psycopg2
4 import pandas as pd
5 import numpy as np
6 import pprint
7 import matplotlib.pyplot as plt
8 import datetime
```

```
In [2]: 1 def connect_db(config_file):
2         with open(config_file) as f:
3             config = json.load(f)
4
5         connect_str = "host={} port={} dbname={} user={} password={} \
6         ".format(config['host'], config['port'], config['database'], config
7
8         connection = None
9         try:
10            connection = psycopg2.connect(connect_str)
11            print("Database connected!")
12        except:
13            print("Unable to connect to the database.")
14
15        return connection
```

```
In [7]: 1 conn = connect_db('config.json')
```

Database connected!

```
In [ ]: 1 """df = pd.DataFrame()
2 for chunk in pd.read_sql('SELECT * FROM IOS_EVENTS', con = conn, chunks:
3     df = df.append(chunk)
4 #df.head(5)
5 """
```

```
In [53]: 1 conn.close()
```

```
In [9]: 1 def execute_sql(sql):
2         global conn
3         if conn.closed:
4             conn = connect_db('config.json')
5         cursor = conn.cursor()
6         cursor.execute(sql)
7         rows = cursor.fetchall()
8         return rows
```

```
In [54]: 1 sql = 'SELECT * FROM IOS_EVENTS limit 5'
2 rows = execute_sql(sql)
3 for row in rows:
4     print(row)
```

Database connected!

```
(datetime.datetime(2017, 6, 1, 0, 4, 50), 'MjA0Mzg5MjY', 'textsearch', 's
earch', 'party', '0', 'iOS_GK_3.5', '', 'iPhone7,1|10.3.1', 'Part', 'en_U
S', 'party', '', 'autosuggest', 'keyboard', -1)
(datetime.datetime(2017, 6, 1, 0, 30, 33), 'NjA0MjI5OTQ', 'textsearch',
'search', 'Fucking', '0', 'iOS_GK_3.5', '', 'iPhone8,1|10.3.2', '', 'en_U
S', 'Fucking', '', 'search_cancel', 'keyboard', -1)
(datetime.datetime(2017, 6, 1, 3, 20, 18), 'MjEzNTcwNA', 'textsearch', 's
earch', 'cocaine', '0', 'iOS_GK_3.5', '', 'iPhone8,1|10.3.2', 'cocaine',
'en_US', 'cocaine', '', 'autosuggest', 'messages', -1)
(datetime.datetime(2017, 6, 1, 3, 42, 24), 'OTg0NTY0NTg', 'textsearch',
'search', 'freshprince', '0', 'iOS_GK_3.5', '', 'iPhone7,2|10.2', 'Fres
h', 'en_US', 'freshprince', '', 'autosuggest', 'keyboard', -1)
(datetime.datetime(2017, 6, 1, 5, 43, 39), 'OTg3NTU1NjY', 'textsearch',
'search', 'congratulations', '0', 'iOS_GK_3.5', '', 'iPhone9,4|10.3.2',
'Congr', 'en_US', 'congratulations', '', 'autosuggest', 'keyboard', -1)
```

Challenge Question - 1

1. "Warmup" metrics:

- Compute the aggregate "share rate": Ratio of shares to searches
- Compute the share rate for the top 1000 most-searched terms
- How many unique search terms are there? What does the distribution look like (consider the count per search term; plot or describe the results)?
- How many unique gifs are shared? What does the distribution look like?

Compute the aggregate "share rate": Ratio of shares to searches

```
In [41]: 1 share_rate_sql = "SELECT TSHARE*1.0/TSEARCH AS SHARE_RATIO \
2     FROM \
3     ( \
4         SELECT 'TOTAL' AS TOTAL, \
5         SUM(CASE WHEN EVENTNAME='textsearch' THEN 1 END) AS TSEARCH,
6         SUM(CASE WHEN EVENTNAME='share' THEN 1 END) AS TSHARE \
7         FROM IOS_EVENTS \
8         GROUP BY 1 \
9     )"
10 results = execute_sql(share_rate_sql)
11 print('The ratio of shares to searches is '+str(results[0][0])[0:4] +'.')
```

The ratio of shares to searches is 0.65.

Compute the share rate for the top 1000 most-searched terms

```

In [58]: 1 top_1000_searched_terms_sql = "SELECT \
2         AVG(TSHARE * 1.0 / TSEARCH) AS TAGS_SHARE_RATIO \
3         FROM \
4         ( \
5             SELECT TAGS AS TAGS, \
6                 SUM(CASE WHEN EVENTNAME='textsearch' THEN 1 END) AS TSEARCH, \
7                 SUM(CASE WHEN EVENTNAME='share' THEN 1 END) AS TSHARE \
8         FROM \
9         ( \
10            SELECT TAGS, EVENTNAME \
11            FROM IOS_EVENTS \
12            WHERE TAGS IN ( SELECT TAGS FROM ( SELECT TAGS, COUNT(*) \
13              WHERE TRIM(TAGS) <> '' GROUP BY TAGS ORDER BY 2 DESC LIMIT 1000 ) ) \
14            ) \
15            GROUP BY 1 ) "
16 results = execute_sql(top_1000_searched_terms_sql)
17 print('The ratio of shares to searches on top 1000 most searched terms is: ', results[0][0])

```

The ratio of shares to searches on top 1000 most searched terms is 0.74.

How many unique search terms are there? What does the distribution look like (consider the count per search term; plot or describe the results)?

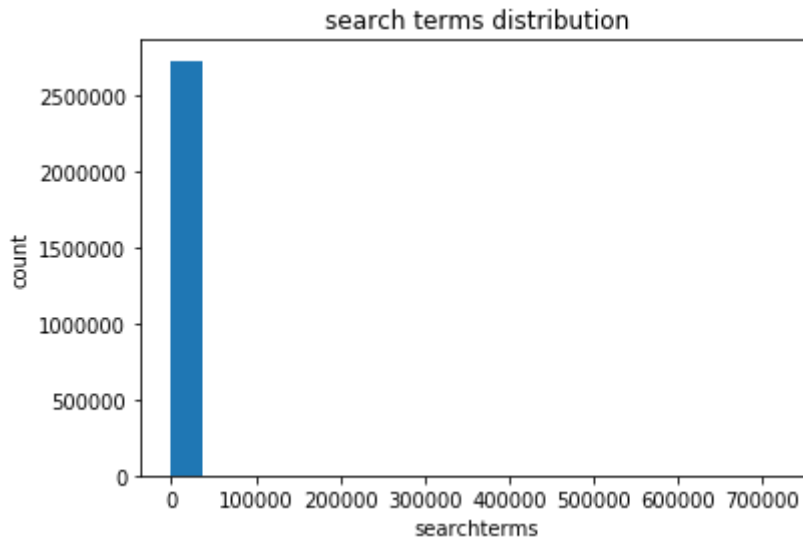
```

In [59]: 1 number_of_unique_searchterms_sql = "SELECT COUNT(DISTINCT TAGS) FROM IOS_EVENTS"
2 results = execute_sql(number_of_unique_searchterms_sql)
3 print('The number of total unique search terms are '+str(results[0][0]))

```

The number of total unique search terms are 2729143.

```
In [17]: 1 unique_searchterms_sql = "SELECT DISTINCT TAGS, COUNT(*) FROM IOS_EVENTS
2 results = execute_sql(unique_searchterms_sql)
3
4 data = []
5 for row in results:
6     data.append(row[1])
7
8 plt.hist(data, bins = 20)
9 plt.xlabel('searchterms')
10 plt.ylabel('count')
11 plt.title('search terms distribution')
12 plt.show()
```



Analysis Description:

1. *Gif for agreement:* Looks like people would want to use gifs to say yes or agree to something. 2. *Gif for happy birthday:* people would want to use gif for occations like happy birthday and to say good night and all.

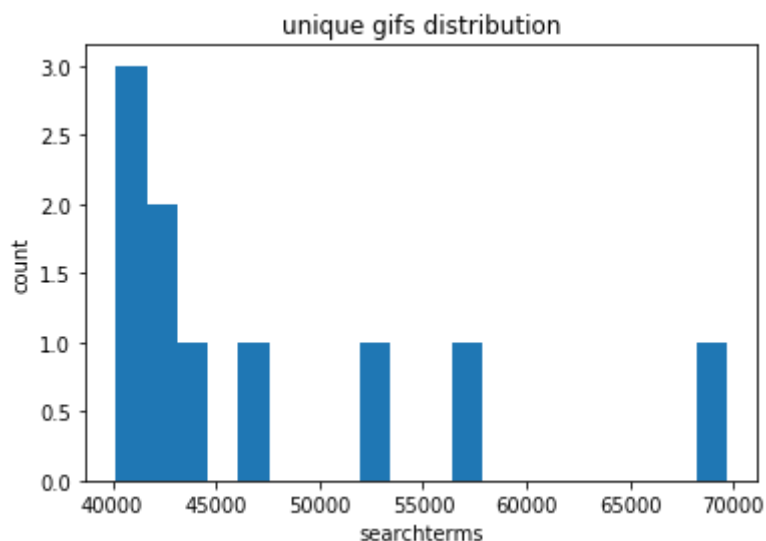
How many unique gifs are shared? What does the distribution look like?

```
In [61]: 1 number_of_unique_gifs_shared_sql = "SELECT COUNT(DISTINCT RIFFID) FROM 1
2 results = execute_sql(number_of_unique_searchterms_sql)
3 print('The number of total unique gifs shared '+str(results[0][0]) +'.')
```

The number of total unique gifs shared2729143.

```
In [20]: 1 unique_gifs_shared_sql = "SELECT RIFFID, COUNT(*) FROM IOS_EVENTS WHERE
2 results = execute_sql(unique_gifs_shared_sql)
3 pprint.pprint(results[0:1000])
4
5 data = []
6 for row in results:
7     data.append(row[1])
8
9 plt.hist(data, bins = 20)
10 plt.xlabel('searchterms')
11 plt.ylabel('count')
12 plt.title('unique gifs distribution')
13 plt.show()
```

```
[('7513882', 69662),
 ('7212866', 57450),
 ('5488810', 52495),
 ('5795910', 46818),
 ('5663730', 43491),
 ('5943705', 42212),
 ('4790020', 42092),
 ('3850980', 41186),
 ('4988274', 40465),
 ('5751664', 40172)]
```



2. Advanced metrics:

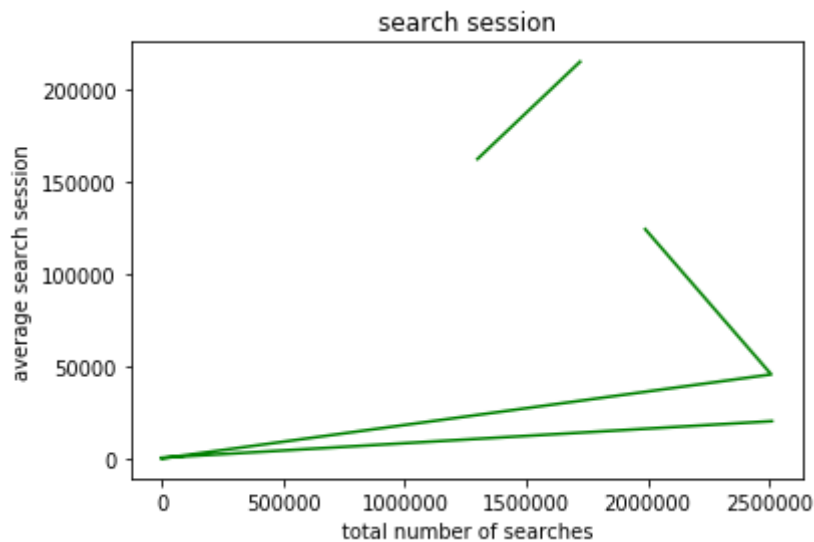
A "search session" is defined as a series of consecutive searches by a user, close in temporal proximity (typically on the order of seconds), that may or may not result in a share. Given the available data, how would you compute the average duration of a search session? Bonus: Estimate the average duration of a search session, and compute the ratio of shares to search sessions.

Average Duration of Search session

```
In [27]: 1 search_session_sql = "SELECT KEYBOARDID, COUNT(*) AS TOTAL_NUMBER_OF_SEA
2          SUM(extract(epoch from duration)) AS TOTAL_SEARCH_DURATION, \
3          AVG(extract(epoch from duration)) AS AVERAGE_SEARCH_DURATION \
4          FROM \
5          (SELECT KEYBOARDID, \
6           TIMESTAMP AS T, \
7           (TIMESTAMP - LAG(TIMESTAMP) OVER (PARTITION BY KEYBOARDID OF
8           FROM IOS_EVENTS WHERE EVENTNAME = 'textsearch' \
9           ) GROUP BY KEYBOARDID LIMIT 10"
```

```
In [29]: 1 x = []
2 y = []
3 for row in execute_sql(search_session_sql):
4     print(row)
5     x.append(row[2])
6     y.append(row[3])
7
8 plt.plot(x, y, 'g')
9 plt.xlabel('total number of searches')
10 plt.ylabel('average search session')
11 plt.title('search session')
12 plt.show()
```

```
('MTA0MDA0Mjc0', 17, 1990390.0, 124399.375)
('MTA0MDA0NTEy', 56, 2509166.0, 45621.2)
('MTA0MDA0NjU1', 2, 8.0, 8.0)
('MTA0MDA1MzAw', 1, None, None)
('MTA0MDA1NzEx', 1, None, None)
('MTA0MDA1Nzg1', 2, 389.0, 389.0)
('MTA0MDA2MDg0', 125, 2511861.0, 20256.9435483871)
('MTA0MDA2MjUz', 1, None, None)
('MTA0MDA2NjE5', 9, 1299906.0, 162488.25)
('MTA0MDA3MjE3', 9, 1722138.0, 215267.25)
```



Ratio of Shares to Searches

3. Experimental design (hint: A/B testing):

- Suppose you've built a new model for the search engine that you'd like to test against the current model.
 - How would you measure the effectiveness of the new model?
 - What additional data would you need in the dataset to compare the effectiveness of the new vs. current models?
 - How would you determine the statistical significance of the results?

How would you measure the effectiveness of the new model?

ANSWER: By testing the new model against test set. The error rate should be very less. The less error rate, more effectiveness of the new model

What additional data would you need in the dataset to compare the effectiveness of the new vs. current models?

ANSWER: Maybe the input entries by the user as the TAGS currently has fully typed words.

How would you determine the statistical significance of the results?

ANSWER: Test validation should have less error.

4. Predictive models:

- There is sufficient data in the dataset to build a crude "search suggestions" model: Given a specific search term, return a list of related search terms. Describe how you might build such a model (describe the data you would use, techniques, inferences, etc.; you do not need to implement this).

```
In [63]: 1 def get_related_search_terms(search_term):
          2     related_search_terms_sql = "SELECT TAGS FROM (\
          3     SELECT TAGS, COUNT(*) FROM IOS_EVENTS WHERE TAGS LIKE '" + search_term + "%'
          4     GROUP BY TAGS ORDER BY 2 DESC LIMIT 10 )"
          5     results = execute_sql(related_search_terms_sql)
          6     for row in results:
          7         print(row[0])
```

```
In [64]: 1 get_related_search_terms('happ')
```

```
happybirthday
happy
happy birthday
happy fathers day
happy dance
happydance
happyanniversary
happy Birthday
happyfriday
happygilmore
```

Predictive Model Description:

We can probably use TD-IDF (term frequency-inverse document frequency) approach using tags and tag1 features. This will give the important words used in the search terms.

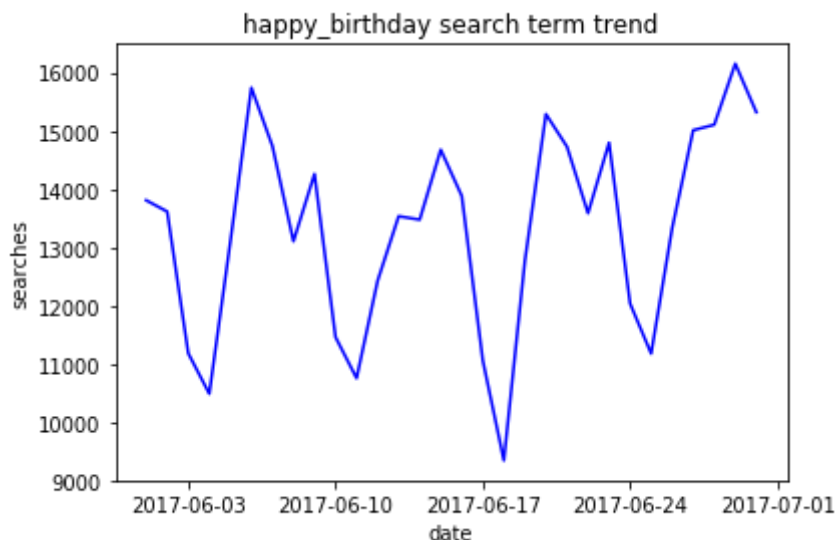
5. Insights and discussion points for on-site:

- This section is not needed for the written responses. For the on-site interviews, please take another look at the dataset, and be prepared to share and discuss additional insights you can extract from the dataset, or useful models or features you can potentially build from the dataset. Example topics might include semantic analysis of searches, graph analysis of shared GIFs, etc.

Below are some of the example graphs and insights from the data

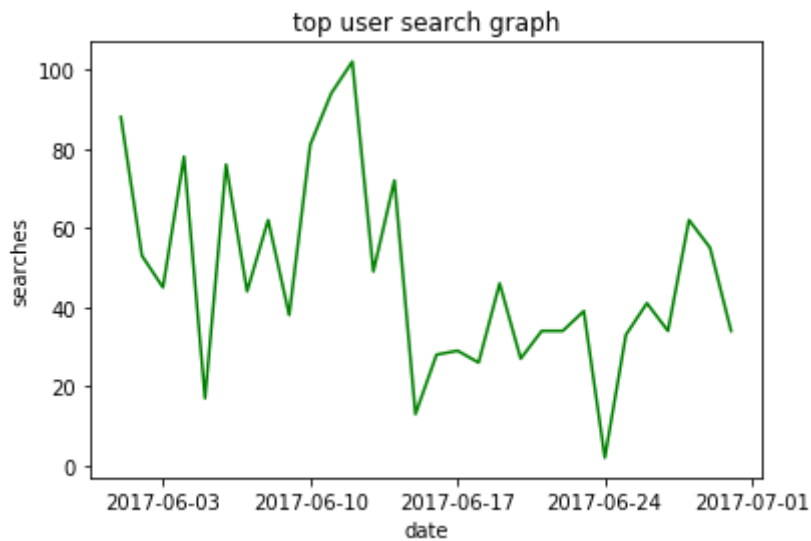
happy birthday search term usage

```
In [31]: 1 happy_birthday_sql = "SELECT TIMESTAMP::DATE , count(*) FROM IOS_EVENTS
2           WHERE TAGS = 'happybirthday' GROUP BY TIMESTAMP::DATE"
3 x = []
4 y = []
5 for row in execute_sql(happy_birthday_sql):
6     x.append(row[0])
7     y.append(row[1])
8 plt.plot(x, y, 'b')
9 plt.xlabel('date')
10 plt.ylabel('searches')
11 plt.title('happy_birthday search term trend')
12 plt.show()
```



top user search activity

```
In [32]: 1 top_user_sql = "SELECT TIMESTAMP::Date, COUNT(*) FROM IOS_EVENTS WHERE I
2 AND KEYBOARDID = 'MTA0NTA4MzU2' GROUP BY TIMESTAMP::Date ORDER BY 1"
3 x = []
4 y = []
5 for row in execute_sql(top_user_sql):
6     x.append(row[0])
7     y.append(row[1])
8
9 plt.plot(x, y, 'g')
10 plt.xlabel('date')
11 plt.ylabel('searches')
12 plt.title('top user search graph')
13 plt.show()
```

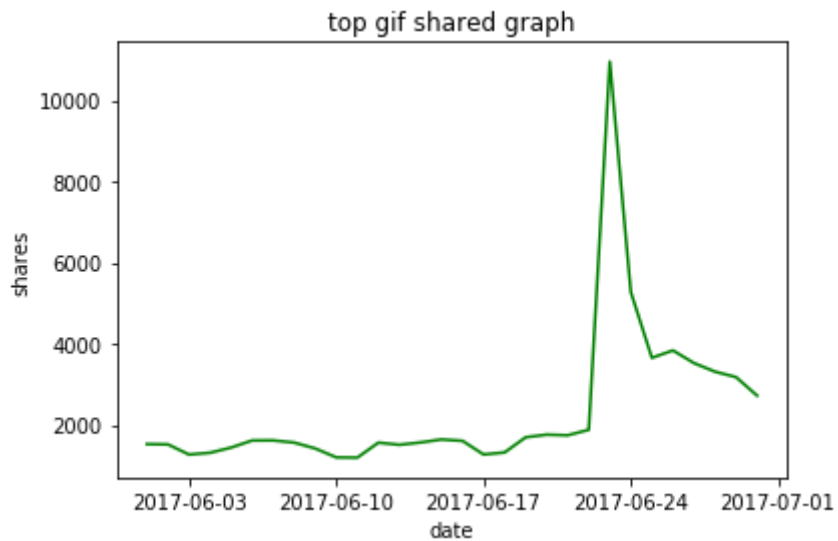


top gif shared trending graph

```

In [33]: 1 top_gif_shared_sql = "SELECT TIMESTAMP::DATE, COUNT(*) FROM IOS_EVENTS \
2 WHERE RIFFID = '7513882' GROUP BY TIMESTAMP::DATE ORDER BY 1"
3 x = []
4 y = []
5 for row in execute_sql(top_gif_shared_sql):
6     x.append(row[0])
7     y.append(row[1])
8
9 plt.plot(x, y, 'g')
10 plt.xlabel('date')
11 plt.ylabel('shares')
12 plt.title('top gif shared graph')
13 plt.show()

```

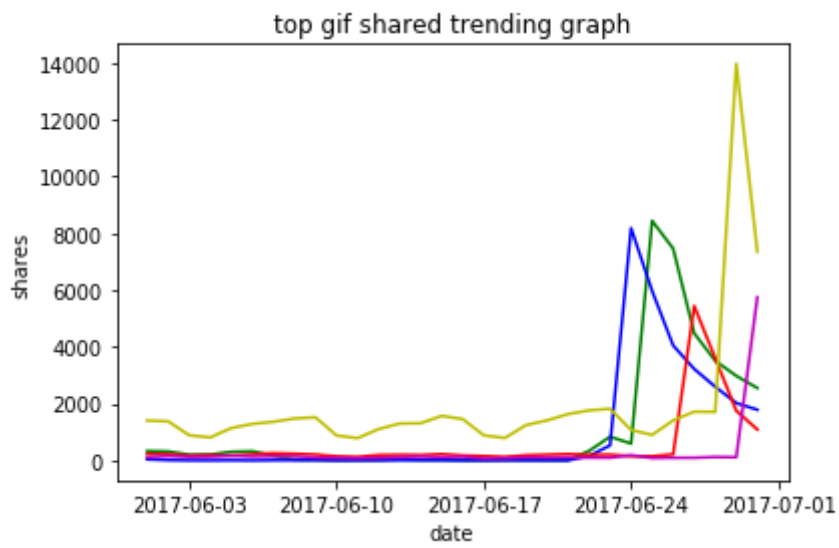


top trending gifs recently

```

1 top_trending_gifs_sql = \
2 "SELECT TIMESTAMP::DATE, RIFFID, COUNT(*) FROM IOS_EVENTS WHERE RIFFID = 1
3 SELECT RIFFID FROM \
4 ( \
5     SELECT RIFFID, _DATETIME, COUNT, (_DATETIME - LAG(_DATETIME) OVER (P
6     (COUNT - LAG(COUNT) OVER (PARTITION BY RIFFID ORDER BY _DATETIME ))
7     FROM \
8     (SELECT RIFFID, TIMESTAMP::DATE AS _DATETIME, COUNT(*) AS COUNT FROM
9 ) \
10 WHERE SHARE_DIFF IS NOT NULL AND SHARE_DIFF > 2000 AND _DATETIME > '2017
11 ORDER BY SHARE_DIFF DESC LIMIT 5 ) GROUP BY TIMESTAMP::DATE, RIFFID ORDER
12 x1 = []
13 x2 = []
14 x3 = []
15 x4 = []
16 x5 = []
17
18 y1 = []
19 y2 = []
20 y3 = []
21 y4 = []
22 y5 = []
23
24 riffids = []
25 for row in execute_sql(top_trending_gifs_sql):
26     if row[1] not in riffids:
27         riffids.append(row[1])
28
29 for row in execute_sql(top_trending_gifs_sql):
30     if riffids[0] == row[1]:
31         x1.append(row[0])
32         y1.append(row[2])
33     if riffids[1] == row[1]:
34         x2.append(row[0])
35         y2.append(row[2])
36     if riffids[2] == row[1]:
37         x3.append(row[0])
38         y3.append(row[2])
39     if riffids[3] == row[1]:
40         x4.append(row[0])
41         y4.append(row[2])
42     if riffids[4] == row[1]:
43         x5.append(row[0])
44         y5.append(row[2])
45
46 plt.plot(x1, y1, 'g')
47 plt.plot(x2, y2, 'b')
48 plt.plot(x3, y3, 'r')
49 plt.plot(x4, y4, 'y')
50 plt.plot(x5, y5, 'm')
51
52 plt.xlabel('date')
53 plt.ylabel('shares')
54 plt.title('top gif shared trending graph')
55 plt.show()

```



In []: 1

In []: 1 conn.close()