

Assignment 3

Yijun Zhou #52533998

```
In [1]: import pandas as pd
import numpy as np
from sqlalchemy import create_engine
import sqlalchemy
```

Collect Data from Twitter API

```
In [2]: #Keys for Twitter API
CONSUMER_KEY = "ZNLSW0nyWYv5fqEQ3eawyBj8z"
CONSUMER_SECRET = "1k1XyWjeCMqSbGNuvQeYcFj2pYawN6GQNZNfhXP3hGZjvyPzla"
ACCESS_KEY = "3309288092-yxGwoIzFDBRh1jkEgpCgLfWf0MAwLyhQkXGhQgm"
ACCESS_SECRET = "lsioTIWicgaTacW3T6zXpxAd25D41LzdP2tXTMr6EYOCZ"
```

```
In [3]: import tweepy

auth = tweepy.OAuthHandler(CONSUMER_KEY, CONSUMER_SECRET)
auth.set_access_token(ACCESS_KEY, ACCESS_SECRET)

api = tweepy.API(auth)
```

```
In [50]: def collectData(query, item_num):
    data_results = []
    for tweet in tweepy.Cursor(api.search_tweets,
                               q=query,
                               result_type="recent",
                               include_entities=True,
                               lang="en").items(item_num):
        data_results.append(tweet)
    print("List length: " + str(len(data_results)))
    return data_results
```

```
In [52]: data_B1 = collectData("Bridgerton", 2500)

List length: 2500
```

```
In [54]: data_B2 = collectData("bridgerton netflix", 2500) #9:31

List length: 2500
```

```
In [65]: data_B3 = []
for tweet in tweepy.Cursor(api.search_tweets,
                           q='bridgerton',
                           until="2022-01-30",
                           include_entities=True,
                           lang="en").items(2500):
    data_B3.append(tweet)
print("List length: " + str(len(data_B3)))

List length: 2500
```

```
In [56]: data_Q1 = collectData("The Queen's Gambit", 2500)

List length: 2441
```

```
In [58]: data_Q2 = collectData("Beth Harmon", 2500)

List length: 426
```

```
In [80]: data_Q3 = collectData("The Queen's Gambit", 2500)

List length: 2459
```

```
In [88]: data_Q4 = collectData("The Queen's Gambit",2500)
List length:2456

In [66]: data_Bridgerton = set(data_B1 + data_B2 + data_B3)
print(len(data_Bridgerton))

5009

In [90]: data_Queen = set(data_Q1 + data_Q2 + data_Q3 + data_Q4)
print(len(data_Queen))

2866

In [91]: data_all = data_Bridgerton.union(data_Queen)
print(len(data_all))

7867
```

Transform Data into Dataframe

```
In [97]: ##transform functions:

def getTweets(tweet):
    trec = {"created_at": tweet.created_at,#time.mktime(tweet.created_at.timetuple()),
            "tweetid":tweet.id,
            "userid":tweet.user.id,
            "text":tweet.text,
            "in_reply_to_status_id":tweet.in_reply_to_user_id,
            "raw_tweet":tweet._json}
    return trec

def getStaticUserInfo(user):
    trec = {"userid":user.id,
            "username":user.name,
            "screen_name":user.screen_name,
            "userdescription":user.description}
    return trec

def getDynamicUserInfo(user):
    trec = {"created_at":user.created_at,
            "userid":user.id,
            "statuses_count":user.statuses_count,
            "followers_count":user.followers_count}
    return trec

def getEntityInfo(tweet):
    trec = {"tweetid":tweet.id,
            "hashtags":tweet.entities['hashtags']}
    return trec
```

```
In [76]: ##Postgres connection info
url = 'postgresql+psycopg2://{}:{}@{}:{}/{}'.format('postgres','2245','localhost','5432','stats170a')
engine = create_engine(url)
```

BridgertonTweetInfo

```
In [75]: Bridgerton = [getTweets(pt) for pt in data_Bridgerton]

df_Bridgerton = pd.DataFrame(Bridgerton).drop_duplicates(['tweetid'])

print(len(df_Bridgerton))
df_Bridgerton.head()

5009
```

	created_at	tweetid	userid	text	in_reply_to_status_id	raw_tweet
0	2022-02-02 03:52:00+00:00	1488721839887384576	1604249562	RT @Anladybee: Kate Sharma in Bridgerton Seaso...	NaN	{'created_at': 'Wed Feb 02 03:52:00 +0000 2022...}
1	2022-01-28 23:08:15+00:00	1487200881460215808	3316300908	@Christopher_NLC I'm currently reading An Offe...	1.821024e+08	{'created_at': 'Fri Jan 28 23:08:15 +0000 2022...}
2	2022-02-01 02:23:34+00:00	1488337199313551362	1344867125790298112	RT @chararirii: The Bridgerton. https://t.co/m...	NaN	{'created_at': 'Tue Feb 01 02:23:34 +0000 2022...}
3	2022-01-28 20:19:53+00:00	1487158508017516544	1305937228783181824	How rich is the Bridgerton family? *runs to Go...	NaN	{'created_at': 'Fri Jan 28 20:19:53 +0000 2022...}
4	2022-01-28 17:18:36+00:00	1487112889550864384	742703809235804162	everyone involved in bridgerton finds my suffe...	7.427038e+17	{'created_at': 'Fri Jan 28 17:18:36 +0000 2022...}

```
In [79]: ##insert table
df_Bridgerton.to_sql("bridgertontweetinfo", con = engine, if_exists = 'replace', index = False,dtype = {'raw_tweet':sqlalchemy.types.JSON}, schema = 'twitter')
```

Out[79]: 9

QueenTweetInfo

```
In [84]: Queen = [getTweets(pt) for pt in data_Queen]

df_Queen = pd.DataFrame(Queen).drop_duplicates(['tweetid'])

print(len(df_Queen))
df_Queen.head()

2866
```

	created_at	tweetid	userid	text	in_reply_to_status_id	raw_tweet
0	2022-02-02 04:21:03+00:00	1488729150068244480	993748794880671745	@meleesadposts Is that the guy from the Queen'...	3.404476e+09	{"created_at": "Wed Feb 02 04:21:03 +0000 2022..."}
1	2022-01-31 05:16:14+00:00	1488018261627502593	812424309285105664	RT @AcademiaDreams: The Queen's Gambit (2020) ...	NaN	{"created_at": "Mon Jan 31 05:16:14 +0000 2022..."}
2	2022-02-02 04:44:20+00:00	1488735010827870210	1306319934528475136	they should have ended the queen's gambit with...	NaN	{"created_at": "Wed Feb 02 04:44:20 +0000 2022..."}
3	2022-02-02 04:30:45+00:00	1488731593217015810	751022622373191684	RT @tasty: Sit down and watch The Queen's Gamb...	NaN	{"created_at": "Wed Feb 02 04:30:45 +0000 2022..."}
4	2022-02-02 03:27:56+00:00	1488715784578867204	1079563520478453760	@laisassauro the queen's gambit	1.473370e+18	{"created_at": "Wed Feb 02 03:27:56 +0000 2022..."}

```
In [86]: ##insert table
df_Queen.to_sql("queentweetinfo", con=engine, if_exists='replace', index=False, dtype = {'raw_tweet':sqlalchemy.types.JSON}, schema = 'twitter')
```

Out[86]: 866

StaticUserInfo

```
In [98]: Static = [getStaticUserInfo(pt.user) for pt in data_all]

df_Static = pd.DataFrame(Static).drop_duplicates(['userid'])

print(len(df_Static))
df_Static.head()
```

6128

	userid	username	screen_name	userdescription
0	1604249562	-`ami7`-	amikosmos	InternationalPopKSensationSunshineRainbowTradi...
1	3316300908	dark_princess000	darkprincess_00	
2	1344867125790298112	delilah	dgeyoydelilah	
3	1305937228783181824	A Fezco Apologist	sameolegal	manifesting Carter x Taylor 🌟
4	742703809235804162	so you find my smile pleasing	connellschains	i will always seek to make it summer for you

```
In [99]: ##insert table
df_Static.to_sql("staticuserinfo", con=engine, if_exists='append', index=False, schema = 'twitter')
```

Out[99]: 128

DynamicUserInfo

```
In [101...]: Dynamic = [getDynamicUserInfo(pt.user) for pt in data_all]

#get the lastest dynamicinfo in collected dataset
df_Dynamic = pd.DataFrame(Dynamic).sort_values(by='created_at', ascending=False).drop_duplicates(['
```

```
df_Dynamic = df_Dynamic.iloc[:, 1:4]
print(len(df_Dynamic))
df_Dynamic.head()
```

6128

Out[101]:

	userid	statuses_count	followers_count
4820	1488584836084342786	46	0
6016	1488577819458342913	7	3
7826	1488575055290380331	8	0
6431	1488571885122056194	7	4
7709	1488476516090626050	1	3

In [103]: `##insert table`
`df_Dynamic.to_sql("dynamicuserinfo", con=engine, if_exists='append', index=False, schema = 'twitter')`

Out[103]: 128

TweetEntity

In [106]:

```
Entity = [getEntityInfo(pt) for pt in data_all]

df_Entity = pd.DataFrame(Entity)

##make the table satisfy 1NF
arr_hashtags = []
for i in df_Entity.hashtags:
    if len(i) == 0:
        arr_hashtags.append(i)
    else:
        temp = []
        for j in range(len(i)):
            temp.append(i[j]['text'])
        arr_hashtags.append(temp)
arr_hashtags

df_Entity_1NF = df_Entity.loc[df_Entity.index.repeat(df_Entity.hashtags.str.len())].assign(hashtags=df_Entity_1NF.rename(columns={'hashtags':'hashtag'}, inplace = True))

print(len(df_Entity_1NF))
df_Entity_1NF.head()
```

2287

Out[106]:

	tweetid	hashtag
8	1488249503865995265	booktwt
8	1488249503865995265	booktwitter
29	1487122067543625731	bridgerton
29	1487122067543625731	ladywhistledown
35	1487512318917361670	bridgerton

In [107]: `##insert table`
`df_Entity_1NF.to_sql("tweetentity", con=engine, if_exists='append', index=False, schema = 'twitter')`

Out[107]: 287

In []:

Part 2

```
1 SET search_path TO twitter;
2
3 CREATE TABLE IF NOT EXISTS BridgertonTweetInfo
4 (
5     created_at          timestamp NOT NULL,
6     tweetid             bigint    NOT NULL,
7     userid              bigint    NOT NULL,
8     text                varchar,
9     in_reply_to_status_id bigint,
10    raw_tweet            jsonb,
11    PRIMARY KEY (tweetid),
12    FOREIGN KEY (userid) REFERENCES StaticUserInfo
13 );
14
15 CREATE TABLE IF NOT EXISTS QueenTweetInfo
16 (
17     created_at          timestamp NOT NULL,
18     tweetid             bigint    NOT NULL,
19     userid              bigint    NOT NULL,
20     text                varchar,
21     in_reply_to_status_id bigint,
22     raw_tweet            jsonb,
23     PRIMARY KEY (tweetid),
24     FOREIGN KEY (userid) REFERENCES StaticUserInfo
25 );
```

```
27 CREATE TABLE IF NOT EXISTS StaticUserInfo
28 (
29     userid               bigint NOT NULL,
30     username             varchar,
31     screen_name          varchar,
32     userdescription      varchar,
33     PRIMARY KEY (userid)
34 );
35
36 CREATE TABLE IF NOT EXISTS DynamicUserInfo
37 (
38     userid               bigint NOT NULL,
39     statuses_count        bigint,
40     followers_count       bigint,
41     PRIMARY KEY (userid),
42     FOREIGN KEY (userid) REFERENCES StaticUserInfo
43 );
44
45 CREATE TABLE IF NOT EXISTS TweetEntity
46 (
47     tweetid              bigint NOT NULL,
48     hashtag              varchar
49 );
```

Part 3

a)

Observation:

I think most of the top words are neutral for both shows. There are many words are common stop words in English and neutral nouns . Other than that, for Bridgerton, there are three positive words detected: “like”, “love” and “kindly”. For the Queen's Gambit, negtive and postive word are both detected, such as “deny” and “like”.

Code:

```
1 --For Bridgerton:
2 SELECT word,
3     count(1)::float/(select count(1) from BridgertonTweetInfo)::float as occurrence_rate
4 FROM BridgertonTweetInfo
5 ,UNNEST(regexp_split_to_array(
6     left(text,position('https://' in text)-1) --ignore the urls in tweet
7     , '[^A-Za-z0-9@#]') -- split by special characters
8 ) AS word
9 WHERE word NOT IN ('', 'RT')
10    AND length(word)>=3 --only keep sequence of 3 or more alphanumeric characters
11 GROUP BY word
12 ORDER BY occurrence_rate DESC
13 LIMIT 100;
14
15 --For the Queen's Gambit
16 SELECT word,
17     count(1)::float/(select count(1) from QueenTweetInfo)::float as occurrence_rate
18 FROM QueenTweetInfo
19 ,UNNEST(regexp_split_to_array(
20     left(text,position('https://' in text)-1) --ignore the urls in tweet
21     , '[^A-Za-z0-9@#]') -- split by special characters
22 ) AS word
23 WHERE word NOT IN ('', 'RT')
24    AND length(word)>=3 --only keep sequence of 3 or more alphanumeric characters
25 GROUP BY word
26 ORDER BY occurrence_rate DESC
27 LIMIT 100;
```

Output for Bridgerton :

	word text	occurrence_rate double precision		word text	occurrence_rate double precision		word text	occurrence_rate double precision	
1	Bridgerton	0.43960870433220206		21	like	0.05210620882411659	41	Atlanta	0.03952884807346776
2	the	0.38151327610301455		22	Kate	0.04791375524056698	42	recaps	0.03952884807346776
3	for	0.18925933320023958		23	was	0.0467159113595528	43	Anthony	0.038929926132960674
4	and	0.1780794569774406		24	coming	0.046316630065881416	44	stop	0.03793172289878219
5	#Bridgerton	0.12377720103813136		25	#bridgerton	0.045717708125374325	45	Netflix	0.03693351966460371
6	bridgerton	0.12317827909762427		26	that	0.04451986424436015	46	just	0.03673387901776802
7	Nicola	0.11399480934318228		27	can	0.043322020363345974	47	but	0.03673387901776802
8	her	0.11339588740267519		28	looking	0.042124176482331806	48	your	0.03573567578358954
9	Coughlan	0.11060091834697544		29	Season	0.04192453583549611	49	BRIDGERTON	0.03533639448991815
10	you	0.09502894789379118		30	are	0.04132561389498902	50	people	0.03473747254941106
11	star	0.09023757236973448		31	@lmnicjuarez	0.04032741066081054	51	see	0.033140347374725494
12	The	0.08704332202036334		32	following	0.04032741066081054	52	Sharma	0.0329407067278898
13	season	0.08464763425833499		33	writers	0.04012777001397484	53	opinions	0.03194250349371132
14	body	0.07306847674186465		34	Snowfall	0.04012777001397484	54	commenting	0.030744659612697144
15	this	0.06787781992413655		35	@vulture	0.04012777001397484	55	have	0.030744659612697144
16	@chaoticguitar	0.06687961668995808		36	Bridge	0.03992812936713915	56	has	0.03054501896586145
17	about	0.06508285086843682		37	culture	0.039728488720303455	57	not	0.03014573767219006
18	@NoContxtBR	0.06288680375324417		38	assign	0.03952884807346776	58	would	0.029746456378518665
19	series	0.05709722499500899		39	Ackman	0.03952884807346776	59	@charariril	0.02954681573168297
20	with	0.05350369335196646		40	Bill	0.03952884807346776	60	from	0.029147534438011578
61	2022	0.028548612497504493		word text	occurrence_rate double precision				
62	all	0.02715112796965462		81	everyone	0.021561189858255142			
63	one	0.026552206029147533		82	characters	0.021361549211419445			
64	love	0.02635256538231184		83	kept	0.020962267917748054			
65	had	0.02635256538231184		84	will	0.020962267917748054			
66	out	0.026152924735476143		85	2020	0.020962267917748054			
67	THE	0.025753643441804752		86	yourself	0.0207626227270912357			
68	fans	0.025554002794969055		87	Her	0.020562986624076664			
69	being	0.02535436214813336		88	@BuzzFeed	0.020562986624076664			
70	@Anladybee	0.024555799560790577		89	month	0.020562986624076664			
71	BALL	0.024156518267119186		90	Daphne	0.020363345977240967			
72	days	0.024156518267119186		91	kindly	0.02016370533040527			
73	she	0.02395687762028349		92	when	0.019964064683569576			
74	get	0.023158315032940708		93	@KennethDredd	0.01976442403673388			
75	until	0.02295867438610501		94	SHUT	0.01976442403673388			
76	wait	0.02275903739269314		95	appreciate	0.01976442403673388			
77	Julia	0.02275903739269314		96	Hell	0.01976442403673388			
78	real	0.021760830505090836		97	his	0.019564783389898182			
79	two	0.021760830505090836		98	Body	0.019564783389898182			
80	time	0.021760830505090836		99	Quinn	0.019564783389898182			
				100	finally	0.018965861449391094			

Output for the Queen's Gambit :

	word text	occurrence_rate double precision		word text	occurrence_rate double precision		word text	occurrence_rate double precision	
1	The	0.7107466852756456		21	harmon	0.09909281228192603	41	post	0.06175854849965108
2	Queen	0.6577110956036287		22	first	0.09769713886950454	42	appreciation	0.06071179344033496
3	Gambit	0.6273551988834613		23	never	0.09665038381018842	43	face	0.06001395673412421
4	the	0.4124214933705513		24	woman	0.0942079553384508	44	has	0.058967201674808095
5	Netflix	0.27983251919050944		25	gambit	0.0928122819260293	45	@chesscom	0.052686671318911374
6	for	0.2198185624563852		26	claimed	0.09176552686671319	46	with	0.04989532449406839
7	defamation	0.20341939986043267		27	named	0.09106769016050244	47	ruled	0.04954640614096301
8	chess	0.16922540125610608		28	@ajplus	0.09071877180739707	48	over	0.04884856943475227
9	and	0.15561758548499652		29	fac	0.0903698534542917	49	MUST	0.0471039776692254
10	after	0.15073272854152128		30	lawsuit	0.08653175157013258	50	fal	0.04675505931612003
11	2020	0.12630844382414516		31	judge	0.07920446615491974	51	Beth	0.044661549197487785
12	@AcademiaDreams	0.12037683182135381		32	dismiss	0.07641311933007676	52	Nona	0.03977669225401256
13	suing	0.10711793440334962		33	motion	0.07257501744591766	53	Harmon	0.034891835310537335
14	she	0.10432658757850663		34	#SeasonOfTheLost	0.07013258897418004	54	her	0.032449406838799724
15	queen	0.10432658757850663		35	#Destiny2Art	0.07013258897418004	55	case	0.0314026517794836
16	that	0.10293091416608513		36	denied	0.07013258897418004	56	can	0.031053733426378228
17	grandmaster	0.10118632240055828		37	@atmosfioric	0.06978367062107467	57	you	0.027564549895324496
18	this	0.10013956734124214		38	look	0.06385205861828332	58	from	0.025819958129797628
19	show	0.09979064898813678		39	was	0.06350314026517795	59	line	0.025819958129797628
20	beth	0.09909281228192603		40	@everyhouranya	0.06245638520586183	60	like	0.025471039776692253

	word text	occurrence_rate double precision	word text	occurrence_rate double precision
61	watch	0.02512212142358688	81	script
62	Judge	0.024773203070481507	82	misrepresented
63	series	0.023028611304954642	83	based
64	just	0.0209351011863224	84	Break
65	@moviemenes	0.0209351011863224	85	written
66	Gaprindashvili	0.020586182833217028	86	Surprised
67	Good	0.020237264480111653	87	boo
68	have	0.019539427773900907	88	gambi
69	Soviet	0.01884159106769016	89	legendary
70	are	0.01884159106769016	90	Picky
71	but	0.018143754361479414	91	@Thompson
72	Blinders	0.018143754361479414	92	Wife
73	@JenShahade	0.018143754361479414	93	HTGAWM
74	Gambi	0.018143754361479414	94	Dexter
75	down	0.018143754361479414	95	Mentalist
76	Prison	0.01779483600837404	96	Diaries
77	man	0.01779483600837404	97	Castle
78	Chess	0.01779483600837404	98	Blacklist
79	@vibes	0.017445917655268667	99	Vampire
80	amp	0.017445917655268667	100	Mad

b)

```

1 SELECT count(distinct hashtag) AS diff_tag_num
2 FROM TweetEntity;

```

diff_tag_num
bigint

1 338

c)

```

1 SELECT collection,
2     count(hashtag)::float/count(distinct userid)::float AS avg_tag_num
3 FROM(
4     SELECT e.tweetid,
5         hashtag,
6         COALESCE(b.userid, q.userid) AS userid,
7         CASE WHEN b.tweetid is not null then 'Bridgerton'
8             ELSE 'the Queens Gambit'END AS collection
9     FROM TweetEntity e
10    LEFT JOIN BridgertonTweetInfo b
11      ON e.tweetid = b.tweetid
12    LEFT JOIN QueenTweetInfo q
13      ON e.tweetid = q.tweetid
14 ) t
15 GROUP BY collection;

```

	collection text	avg_tag_num double precision
1	Bridgerton	2.58179012345679
2	the Queens Gambit	1.9935064935064934

d)

```

1  SELECT hashtag
2  FROM (
3      SELECT hashtag,
4          count(distinct b.tweetid )   as b_count,
5          count(distinct q.tweetid )   as q_count
6  FROM TweetEntity e
7  LEFT JOIN BridgertonTweetInfo b
8      ON e.tweetid = b.tweetid
9  LEFT JOIN QueenTweetInfo q
10     ON e.tweetid = q.tweetid
11    group by hashtag
12  ) t
13 WHERE b_count >0 and q_count>0

```

hashtag	character varying
1	Louies
2	BestFanArmy
3	podcasts
4	Netflix
5	tvtime
6	News
7	BestCoverSong
8	BestMusicVideo
9	Butter
10	HappierThanEver
11	iHeartAwards
12	NEWS
13	bancodeseries
14	ShawnMendes

e)

```

1  SELECT hashtag,b_count,q_count
2  FROM (
3      SELECT hashtag,
4          count(distinct b.tweetid )   as b_count,
5          count(distinct q.tweetid )   as q_count
6  FROM TweetEntity e
7  LEFT JOIN BridgertonTweetInfo b
8      ON e.tweetid = b.tweetid
9  LEFT JOIN QueenTweetInfo q
10     ON e.tweetid = q.tweetid
11    group by hashtag
12  ) t
13 WHERE b_count >0 and q_count>0
14 ORDER BY b_count+q_count DESC
15 LIMIT 10

```

hashtag	character varying	b_count	bigint	q_count	bigint
1	Netflix	17		20	
2	iHeartAwards	21		5	
3	BestFanArmy	14		3	
4	tvtime	4		12	
5	Louies	9		3	
6	HappierThanEver	7		1	
7	bancodeseries	1		7	
8	podcasts	7		1	
9	BestCoverSong	7		1	
10	News	2		2	

f)

```

1  SELECT b.userid,username,screen_name
2  FROM BridgertonTweetInfo b
3  LEFT JOIN StaticUserInfo s
4      ON b.userid = s.userid
5  WHERE b.userid in (
6      SELECT userid
7          FROM QueenTweetInfo
8      )

```

userid	username	screen_name	userid	username	screen_name	userid	username	screen_name
1	225187661	GIGICAPONEPR	23	111556701	Daily Mail Celebrity	45	787546010	Berkley Bear
2	44686891	Barbie Williams	24	1561416842	gaby	46	1190934402584338400	Kisahki99454428
3	613385915	Lauri Donahue	25	280249909	angie *; la	47	1221366007052324600	World_News_eng
4	14335144	Digital Spy	26	1196845795703742500	Internewscast Media	48	1257771553988755500	jessfrancis
5	283604227	Andy Vermaut	27	4725957312	Radio TFI (Home of The Taxi Stand Hour)	49	4729957312	TheRadioTFI
6	1913959446	9 Breaking News	28	1326561978308059100	NewsPlayer+	50	1420660507	Woody
7	1257771553988755500	Jude 51 days	29	821957852	kara aielo	51	1420660507	Knewz_Currently
8	1257771553988755500	Jude 51 days	30	456192512	jess	52	759655316946554900	Lydia Massiah
9	171447719	DS Breaking News	31	577329310	News 100 Media	53	2971469687	Rachelsmith999
10	2544099769	██████████	32	1448332054186299400	chrts ³⁹	54	17434245	maxvaldes
11	14173315	NBC News	33	225187661	GIGICAPONEPR	55	1257771553988755500	jessfrancis
12	577329310	News 100 Media	34	1420660507	Woody	56	1362177019	N A
13	999907020	OFFICIALGIGICAPONE	35	1313269879	The 206geek	57	1420660507	Woody
14	280617017	50 Shades of Rein	36	1257771553988755500	Jude 51 days	58	283604227	Andy Vermaut
15	1157901558365667300	Amstel News	37	225187661	GIGICAPONEPR	59	456192512	jess
16	11156701	Daily Mail Celebrity	38	1435638352204554200	ReportWire	60	1257771553988755500	jessfrancis
17	456192512	jess	39	49395291	Charles Myrick	61	1257771553988755500	jessfrancis
18	1458943296118612000	Georgina Wiggins	40	1257771553988755500	Jude 51 days	62	999907020	OFFICIALGIGICAPONE
19	456192512	jess	41	277600246	Gossip Bucket	63	230420580	rizkafaribila
20	283604227	Andy Vermaut	42	273353243	Daily News Kit	64	1007725923720835100	NewsfeedsMedia
21	117139256404082700	eyien · DOYOUNG DAY♥	43	4725957312	Radio TFI (Home of The Taxi Stand Hour)	65	1898672911	repeatinglitanies
22	143425578462847000	Katherine	44	456192512	jess	66	129313311087645400	TexasDem2

g)

```
1 SELECT s.userid,  
2     username,  
3     screen_name,  
4     userdescription,  
5     statuses_count,  
6     followers_count  
7 FROM StaticUserInfo s  
8 LEFT JOIN DynamicUserInfo d  
9     ON s.userid = d.userid  
10 ORDER BY d.statuses_count DESC  
11 LIMIT 10;
```

	userid bigint	username character varying	screen_name character varying	userdescription character varying	statuses_count bigint	followers_count bigint
1	289148078	♥...david_jones @MrDJones	MrDJones	You've come to right place. There's NEW...	2582297	5019
2	1420660507	Woody	Knewz_Currently	Love my Family, my Harley, my Jeeps, b...	2551391	1570
3	275942684	rashid al dosari	rashidaldosari	I listen to and respect others views and I...	2272146	2314
4	787546010	Berkley Bear	BerkleyBearNews	Bringing the news that is important, that...	1932470	1382
5	40173650	Global News Report	robinsonswire	Flying the web for trusted world news re...	1881016	23900
6	1596829382	Zyite.news	ZyiteGadgets	Zyite is a digital media provider	1742740	3216
7	49404200	Charles Myrick -CEO	medicinehelp	American Consultants Rx is focused on ...	1696332	1358
8	2544099769	🇨🇴Juan Carlos	falconhamada_90	Official Account	1629003	2343
9	1221366000752324600	World News	world_news_eng	For All English Users	1590967	2444
10	3534222021	NewsOnePlace.com	newsoneplace	https://t.co/WWMmlkl7DB is a news ag...	1586895	998

h)

```
1 CREATE INDEX gin_text_b ON BridgertonTweetInfo USING gin(to_tsvector('english', text));
2 CREATE INDEX gin_text_q ON QueenTweetInfo USING gin(to_tsvector('english', text));
```

i)

Finding :

I filter the tweets about Brigerton that have been liked by at least 100 Twitter users. And I choose my keyword by looking at the most frequent words in those “liked” tweets. I found people are talking about the new season of Brigerton which will come in February and the final ball in Brigerton. About 30% tweets that I collected are taking about this topic.

```
1 --Find the keywords using favorite count
2 SELECT lower(word),count(1) as word_count
3 FROM BridgertonTweetInfo,UNNEST(regexp_split_to_array(
4     left(text,position('https://in text)-1) --ignore the urls in tweet
5     , '[^A-Za-z0-9@]') -- split by special characters
6 ) AS word
7 WHERE (raw_tweet->'favorite_count'):: text::int >= 100
8     AND length(word)>=3 --ingore word shorter than 3 characters
9 GROUP BY lower(word)
10 ORDER BY word_count desc
11
12 --Search with keywords (february,month,ball,final,coming,next,season)
13 SELECT tweetid,text
14 FROM BridgertonTweetInfo
15 WHERE lower(text) SIMILAR TO '%(february|month|ball|final|next|season|coming)%'
16
17 --Check the porportion of tweets with keywords
18 SELECT count(k.tweetid) as ktweets_num,
19     count(b.tweetid)      as total_tweets_num,
20     count(k.tweetid)::float/count(b.tweetid)::float as rate
21 FROM BridgertonTweetInfo b
22 LEFT JOIN (
23     SELECT tweetid,text
24     FROM BridgertonTweetInfo
25     WHERE lower(text) SIMILAR TO '%(february|month|ball|final|next|season|coming)%'
26 ) k
27     ON b.tweetid = k.tweetid
```

	lower text	word_count bigint	tweetid bigint	text text
1	bridgerton	25	1	1488721839887384600 RT @Anladyybee: Kate Sharma in Bridgerton Season 2 https://t.co/PtBEllleGbT
2	the	12	2	1488505691346436000 RT @chaoticguitar: FINAL BALL OF BRIDGERTON S2 COULD BE A MASQUERADE BALL IM FREAKING OUT CHRIS VAN DUSEN SKSKSKSKSS https://t.co/wkmAbaJ...
3	february	4	3	1487091052821446700 RT @KennethDredd: Bill Ackman 2020: "Hell is coming"
4	and	4	4	1488505741094781000 RT @chaoticguitar: FINAL BALL OF BRIDGERTON S2 COULD BE A MASQUERADE BALL IM FREAKING OUT CHRIS VAN DUSEN SKSKSKSKSS https://t.co/wkmAbaJ...
5	month	4	5	1488367921675993000 RT @BuzzFeed: Bridgerton star Nicola Coughlan would kindly appreciate it if you kept your opinions on her body to yourself. https://t.co/fW...
6	ball	4	6	1488607929825640400 RT @NeenaQueen20: February 1. That means we get Bridgerton next month. We get to see the most anticipated book of the series happen NEXT MO...
7	she	4	7	1488546535721869300 RT @chaoticguitar: Good morning! Happy February and 52 days to #Bridgerton 🎉 + ❤️ The month of love is upon us weeeeee! https://t.co/nJxGmETaMK
8	her	4	8	1488512557598052400 "I'm not a body positivity activist, I'm an actor."
9	his	3	9	1487483267041398800 RT @kathonyshbee: No idea if anybody here cares about this modern Kathony AU, but chapter 3 of 'Tis The Damn Season is up! No smut in this o...
10	kate	3	10	1488440672025776000 I finished bridgerton now i'm waiting on season 2 🌟 iykyk
11	final	3	11	1487343238281314300 RT @Anladyybee: Kate Sharma in Bridgerton Season 2 https://t.co/PtBEllleGbT
12	was	3	12	1487170303948165000 @brauhala hear me out...for your next "Broadway Loves" series...The (Unofficial) Bridgerton Musical by @abigailbarloww... https://t.co/h1tjkx2dM
13	next	3	13	1487285813872365600 Bridgerton season 2 when
14	star	3	14	1487169778720641000 RT @menacetanthony: Can't wait to see more of Shelley Conn as Mary Sharma in season 2 of Bridgerton! Look how lovely she is ❤️ ☺
15	can	3	15	1487022402588577800 Nicola Coughlan teases steamier scenes between Penelope and Colin in #Bridgerton season 2:
16	out	3	16	1488169675452260400 RT @NoContxtBR: 53 days until we finally see the real star of #bridgerton https://t.co/ZtrHseq4zI
17	season	3	17	1488560141238734800 that's great bridgerton teaser next tho? https://t.co/ftZgptOjW8
18	happy	2	18	148838743301299200 RT @BuzzFeed: Bridgerton star Nicola Coughlan would kindly appreciate it if you kept your opinions on her body to yourself. https://t.co/fW...
19	relaxed	2	19	1488444323779724000 RT @NoContxtBR: 53 days until we finally see the real star of #bridgerton https://t.co/ZtrHseq4zI
20	coming	2	20	1488250641424261000 RT @CNN: 'Bridgerton' star Nicola Coughlan has asked people to stop commenting on her body -- the latest celeb to speak out about bodyshami...
21	confirming	2	21	1487052284043350000 I've starting watching Bridgerton, and watching Simon tell Daphne that they were never really friends is just watch... https://t.co/AKoOUHJ0og
		22	22	1487445515478388700 RT @Anladyybee: Kate Sharma in Bridgerton Season 2 https://t.co/PtBEllleGbT
	ktweets_num bigint	total_tweets_num bigint	rate double precision	
21	1	1593	5009	0.31802755040926334