

INFO 6210

Data Management and Database

Assignment 3

Changrong Chen 001276880

I chose games as my domain. The original data comes from steam and it has created a PostgreSQL database successfully. In this assignment I need to convert my PostgreSQL database to mongoDB.

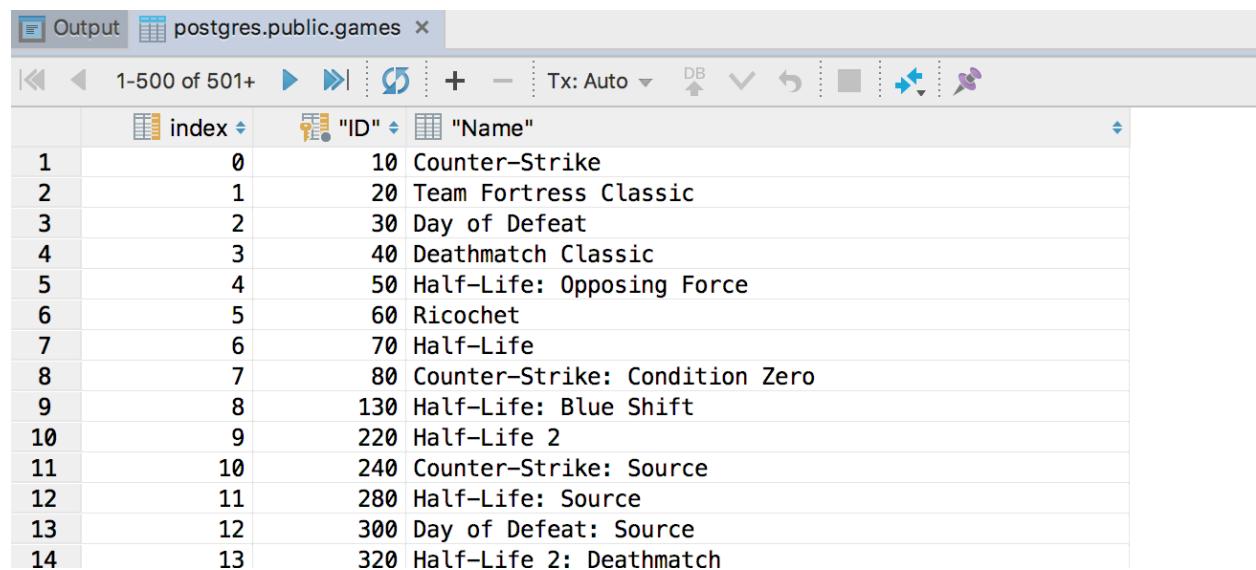
The original database: 2 parts, game & social media data.

Game has 4 tables: games, games_features, games_reviews, games_users.

Social media data has 2 tables: game_tweets, game_hashtags.

The original sample data:

games:



A screenshot of the pgAdmin interface showing the 'games' table in the 'postgres.public' schema. The table has two columns: 'index' and 'Name'. The data consists of 501 rows, with the first 14 rows listed below:

	index	Name
1	0	10 Counter-Strike
2	1	20 Team Fortress Classic
3	2	30 Day of Defeat
4	3	40 Deathmatch Classic
5	4	50 Half-Life: Opposing Force
6	5	60 Ricochet
7	6	70 Half-Life
8	7	80 Counter-Strike: Condition Zero
9	8	130 Half-Life: Blue Shift
10	9	220 Half-Life 2
11	10	240 Counter-Strike: Source
12	11	280 Half-Life: Source
13	12	300 Day of Defeat: Source
14	13	320 Half-Life 2: Deathmatch

games_features:

Output postgres.public.games_features x

1-500 of 501+ Tx: Auto DB View Query

	index	"ResponseID"	"ReleaseDate"	"Metacritic"	"RecommendationCount"	"IsFree"	"GenresNonGame"	"GenresIndie"	"GenresAction"	"GenresAdventure"	"GenresCasual"	"GenresStrategic"
1	0	10 Nov 1 2000	88	68991								
2	1	28 Apr 1 1999	0	2439								
3	2	30 May 1 2003	79	2319								
4	3	40 Jun 1 2001	0	888								
5	4	50 Nov 1 1999	0	2934								
6	5	60 Nov 1 2000	0	1965								
7	6	70 Nov 8 1998	96	12486								
8	7	80 Mar 1 2004	65	7667								
9	8	130 Jun 1 2001	71	2219								
10	9	220 Nov 16 2004	96	35792								
11	10	240 Nov 1 2004	88	53931								
12	11	280 Jun 1 2004	0	2547	✓							
13	12	300 Jul 12 2010	80	7185								
14	13	320 Nov 1 2004	0	4328								
15	14	340 Oct 27 2005	0	4352	✓							
16	15	360 May 1 2006	0	864	✓							
17	16	380 Jun 1 2006	87	4437								
18	17	400 Oct 10 2007	90	27535								
19	18	420 Oct 10 2007	90	6184								
20	19	440 Oct 10 2007	92	383949	✓							
21	20	500 Nov 17 2008	89	9980								
22	21	550 Nov 16 2009	89	140726								
23	22	570 Jul 9 2013	90	590480	✓							
24	23	620 Apr 18 2011	95	73128								
25	24	630 Jul 19 2010	77	14625	✓							
26	25	730 Aug 21 2012	83	1427633								
27	26	1002 Oct 12 2005	69	0						✓		
28	27	1200 Mar 14 2006	81	1220								
29	28	1250 May 14 2009	72	44364								

games_reviews:

Output postgres.public.games_reviews x

1-500 of 501+ Tx: Auto DB View Query

	index	"ResponseID"	"AboutText"	"ShortDescrip"
1	0	10	Play the worlds number 1 online action game. Engage in an incredibly realistic bra...	
2	1	20	One of the most popular online action games of all time Team Fortress Classic feat...	
3	2	30	Enlist in an intense brand of Axis vs. Allied teamplay set in the WWII European Th...	
4	3	40	Enjoy fast-paced multiplayer gaming with Deathmatch Classic (a.k.a. DMC). Valves t...	
5	4	50	Return to the Black Mesa Research Facility as one of the military specialists assi...	
6	5	60	A futuristic action game that challenges your agility as well as your aim Ricochet...	
7	6	70	Named Game of the Year by over 50 publications Valves debut title blends action an...	
8	7	80	With its extensive Tour of Duty campaign a near-limitless number of skirmish modes...	
9	9	130	Made by Gearbox Software and originally released in 2001 as an add-on to Half-Life...	
10	10	220	1998. HALF-LIFE sends a shock through the game industry with its combination of po...	#app_220_note_1
11	11	240	THE NEXT INSTALMENT OF THE WORLDS # 1 ONLINE ACTION GAME Counter-Strike: Source b...	Just updated to include player stats achievements new scoreboards and more!
12	12	280	Winner of over 50 Game of the Year awards Half-Life set new standards for action g...	
13	13	300	Day of Defeat offers intense online action gameplay set in Europe during WWII. Ass...	
14	14	320	Fast multiplayer action set in the Half-Life 2 universe! HL2s physics adds a new d...	
15	15	340	Originally planned as a section of the Highway 17 chapter of Half-Life 2 Lost Coas...	
16	16	360	Half-Life Deathmatch: Source is a recreation of the first multiplayer game set in ...	
17	17	380	Half-Life 2 has sold over 4 million copies worldwide and earned over 35 Game of th...	Supports cross-platform Steam Cloud – continue your game on any supported OS. Plus...
18	18	400	Portalâ 4 is a new single player game from Valve. Set in the mysterious Aperture S...	
19	19	420	Half-Lifeâ 2: Episode Two is the second in a trilogy of new games created by Valv...	Supports cross-platform Steam Cloud – continue your game on any supported OS.
20	20	440	The most fun you can have online – PC Gamer Is now FREE! Theres no catch! Play as ...	Nine distinct classes provide a broad range of tactical abilities and personalitie...
21	21	500	From Valve (the creators of Counter-Strike Half-Life and more) comes Left 4 Dead ...	

games_users:

Output postgres.public.games_users x

1-500 of 501+ Tx: Auto DB View Query

	index	"ID"	user_id	games
1	0	377160	151603712	Fallout 4
2	2	17390	151603712	Spore
3	4	550	151603712	Left 4 Dead 2
4	6	339800	151603712	HuniePop
5	8	238960	151603712	Path of Exile
6	10	367450	151603712	Poly Bridge
7	12	500	151603712	Left 4 Dead
8	14	440	151603712	Team Fortress 2
9	16	203160	151603712	Tomb Raider
10	18	237990	151603712	The Banner Saga
11	20	8870	151603712	BioShock Infinite
12	22	22370	151603712	Fallout 3 – Game of the Year Edition
13	24	12210	151603712	Grand Theft Auto IV
14	26	200210	151603712	Realm of the Mad God
15	28	298160	151603712	Eldevin
16	30	570	151603712	Dota 2
17	32	7670	151603712	BioShock
18	34	301520	151603712	Robocraft
19	36	4000	151603712	Garry's Mod
20	38	250260	151603712	Jazzpunk
21	40	108710	151603712	Alan Wake

game_tweets:

game_hashtags:

0	12	VarietyStreaming	Counter-Strike
1	15	_Ø¹Ø¬Ø“__Ø¬Ø“_Ø§Ø°Ø§	Counter-Strike
2	16	VarietyStreaming	Counter-Strike
3	43	BullyHunters	Counter-Strike
4	57	CSGO	Counter-Strike
5	76	CSGO	Counter-Strike
6	83	FAv	Counter-Strike
7	84	FAv	Counter-Strike
8	85	ThrowbackThursday	Counter-Strike
9	88	FAv	Counter-Strike
10	89	FAv	Counter-Strike
11	101	Infinite	Counter-Strike

In order to convert SQL database to NoSQL database, we need to denormalization table.

So, games, games_features, games_reviews, games_users will be denormalized to one table; game_tweets, game_hashtags will be denormalized to one table.

In SQL database, use SQL syntax to do social media denormalization:

Sample code:

Then convert it to a csv file and then use mongoimport command to import data to mongoDB.

Sample code:

```
In [8]: df_k.to_csv('full_tweets_data.csv', encoding='utf-8', index=False)

In [11]: twitter_file = pd.read_csv('/Users/chen/Desktop/data/full_tweets_data.csv', encoding='latin1')

In [12]: twitter_file.duplicated().sum()

Out[12]: 0

In [13]: twitter_csv_path='/Users/chen/Desktop/data/full_tweets_data.csv'

In [14]: db_name='local'
        coll_name1='game_twitter'

In [15]: mongoimport(twitter_csv_path, db_name, coll_name1)
          /Users/chen/anaconda3/lib/python3.6/site-packages/ipykernel_launcher.py:11: DeprecationWarning: remove is deprecated.
          Use delete_one or delete_many instead.
          # This is added back by InteractiveShellApp.init_path()
          /Users/chen/anaconda3/lib/python3.6/site-packages/ipykernel_launcher.py:12: DeprecationWarning: insert is deprecated.
          Use insert_one or insert_many instead.
          if sys.path[0] == '':
Out[15]: 3261
```

In SQL database, use SQL syntax to do game data denormalization:

Sample code:

First, denormalization 4 table to one table.

```
In [5]: import psycopg2
conn = psycopg2.connect(database='postgres', user='postgres', password='')
df1=pd.read_sql('''
SELECT A."Name", B.* , C."AboutText",C."ShortDescrip", C."Reviews", D.user_id
FROM
(
    (SELECT "ID","Name" FROM games) A
    FULL OUTER JOIN
    (SELECT * FROM games_features) B on A."ID" = B."ResponseID"
    FULL OUTER JOIN
    (SELECT * FROM games_reviews) C on A."ID" = C."ResponseID"
    FULL OUTER JOIN
    (SELECT * FROM games_users) D on A."ID" = D."ID"
);'', conn=conn)
```

In [10]: df1

	Name	index	ResponseID	ReleaseDate	Metacritic	RecommendationCount	IsFree	GenreisNonGame	GenreisIndie	GenreisAction	GenreisE...
0	Fallout 4	6635.0	377160.0	Nov 9 2015	84.0	72929.0	False	False	False	False	False ...
1	Left 4 Dead 2	21.0	550.0	Nov 16 2009	89.0	140726.0	False	False	False	True	True ...
2	HuniePop	5043.0	339800.0	Jan 19 2015	0.0	9827.0	False	False	True	False	False ...
3	Path of Exile	2147.0	238960.0	Oct 23 2013	86.0	33124.0	True	False	True	True	True ...

Then use group command to correct the 1 to n relationship:

Then convert it to a csv file and then use mongoimport command to import data to MongoDB.

```
In [20]: gamedata_csv_path='/Users/chen/Desktop/steam-data/datafile/data1/fulldata2.csv'
db_name='local'
coll_name3='games_data'

In [21]: mongoimport(gamedata_csv_path, db_name, coll_name3)

Out[21]: 13208
```

Sample data in MongoDB:

games:

The screenshot shows the MongoDB Compass interface connected to a local host at port 27017. The database is 'local' and the collection is 'games_data'. The interface displays 13.2k documents with a total size of 26.4MB and an average size of 2.0KB. There is one index with a total size of 128.0KB and an average size of 128.0KB. The 'Documents' tab is selected, showing two document snippets. The first snippet is a game entry for 'The Count of Monster Disco' with various metadata fields like _id, Name, ResponseID, ReleaseDate, Metacritic, etc. The second snippet is another game entry with similar fields. The left sidebar shows other databases and collections like admin, config, local, Q9_hashtags, Q9_mentions, game_twitter, startup_log, and games_data.

Social Media data:

Question Sample Code:

Answer Questions

```
In [16]: from pymongo import MongoClient
client = MongoClient('localhost', 27017)

In [17]: db = client.local
collection=db.game_twitter

In [43]: pprint.pprint(collection.find_one({"game": "Dota 2"}))

{'Media': None,
 'Urls': "[{'url': 'https://t.co/u8r04MECZc', 'expanded_url': ''}
           '\"https://www.twitch.tv/sagembblack', 'display_url': ''
           '\"twitch.tv/sagembblack', 'indices': [84, 107]}]", "
 '_id': ObjectId('5ad55a3f4ale94208242add2'),
 'array_agg': "'VarietyStreaming']",
 'created_at': '2018-04-13 22:30:06',
 'game': 'Dota 2',
 'id': 961,
 'retweet_count': 1,
 'tweets': 'RT @Ashtheimpailer: Hi Guys! \n'
           '\n'
           'I am now live playing Counter Strike!\n'
           '\n'
           'Come say hi!\n'
           '\n'
           'https://t.co/u8r04MECZc\n'
           '\n'
           '#varityStreaming \n'
           '#SmallStreamÃ¢Ã\x80Ã¡',
 'user_id': 2910985259,
 'user_location': 'Living with @kyroskoh (SG)'}]
```

```
In [21]: #Question1 What are tags are associated with a person, place or thing
list(db.game_twitter.aggregate([
    {"$match": {"game": "Dota 2"}},
    {"$group": {"_id": "$array_agg", "count": {"$sum": 1}}},
    {"$sort": {"count": -1}}]))
[{"_id": "[DanaDurnford', 'BanTokyo2020', 'StopTokyo2020]", "count": 13},
 {"_id": "[quake]", "count": 12},
 {"_id": "[Divinity]", "count": 12},
 {"_id": "[Call]", "count": 9},
 {"_id": "[Fallout]", "count": 9},
 {"_id": "[Facebook]", "count": 8},
 {"_id": "[Fallout", "ThrowbackThursday]", "count": 7},
 {"_id": "[Divinity', 'PS4]", "count": 6},
 {"_id": "[DivinityOriginalSin2', 'XboxOne', 'Playstation4]", "count": 6},
 {"_id": "[Dota]", "count": 6},
 {"_id": "[Dota', 'Steam]", "count": 6},
 {"_id": "[Quake', 'Quake', 'DadJokes]", "count": 6},
 {"_id": "[Streamhype', 'supportsmallstreamers', 'Call]", "count": 5},
 {"_id": "[Left', 'Dead]", "count": 5},
 {"_id": "[QuakeChampions]", "count": 4},
 {"_id": "[Repost', 'BandaiNamco]", "count": 4},
 {"_id": "[FAV]", "count": 4},
 {"_id": "[Repost', 'BandaiNamco', 'PS4]", "count": 4},
 {"_id": "[dota', 'meepo', 'techies', 'SupportSmallerStreamers]", "count": 4}]

In [22]: list(db.game_twitter.aggregate([
    {"$match": {"game": "Left 4 Dead"}},
    {"$group": {"_id": "$array_agg", "count": {"$sum": 1}}},
    {"$sort": {"count": -1}})))
Out[22]: [{"_id": "[DOTA]", "count": 56},
 {"_id": "[earthquake', 'quake]", "count": 20},
 {"_id": "[Divinity', 'XboxOne]", "count": 13},
 {"_id": "[DanaDurnford', 'BanTokyo2020', 'StopTokyo2020]", "count": 13},
 {"_id": "[quake]", "count": 12},
 {"_id": "[Divinity]", "count": 12},
 {"_id": "[Call]", "count": 9},
 {"_id": "[Fallout]", "count": 9},
 {"_id": "[Facebook]", "count": 8},
 {"_id": "[DivinityOriginalSin2', 'XboxOne', 'Playstation4]", "count": 7},
 {"_id": "[Fallout", "ThrowbackThursday]", "count": 7},
 {"_id": "[Divinity', 'PS4]", "count": 6},
 {"_id": "[Dota]", "count": 6},
 {"_id": "[Quake', 'Quake', 'DadJokes]", "count": 6},
 {"_id": "[Repost', 'BandaiNamco', 'PS4]", "count": 5},
 {"_id": "[Dota', 'Steam]", "count": 5},
 {"_id": "[Streamhype', 'supportsmallstreamers', 'Call]", "count": 5},
```

```
In [28]: #Question2 What social media users are like other social media users in your domain?
match = {
    'tweets': {"$regex": "Counter Strike"}
}

group = {
    '_id': "$user_id"
}

ret2 = db.game_twitter.aggregate([
    {"$match": match},
    {"$group": group}
])

list(ret2)
Out[28]: [{"_id": 138156328}, {"_id": 2910985258}]

In [23]: #Question3 What people, places or things are popular in your domain
list(db.game_twitter.aggregate([
    {"$group": {"_id": "$array_agg", "count": {"$sum": 1}}},
    {"$sort": {"count": -1}})))
Out[23]: [{"_id": "[DOTA]", "count": 504},
 {"_id": "[earthquake', 'quake]", "count": 139},
 {"_id": "[quake]", "count": 108},
 {"_id": "[Divinity]", "count": 103},
 {"_id": "[Divinity', 'XboxOne]", "count": 95},
 {"_id": "[Fallout]", "count": 81},
 {"_id": "[Call]", "count": 81},
 {"_id": "[DanaDurnford', 'BanTokyo2020', 'StopTokyo2020]", "count": 79},
 {"_id": "[Facebook]", "count": 72},
 {"_id": "[Dota]", "count": 54},
 {"_id": "[Fallout", "ThrowbackThursday]", "count": 46},
 {"_id": "[Divinity', 'PS4]", "count": 46},
 {"_id": "[Dota', 'Steam]", "count": 43},
 {"_id": "[DivinityOriginalSin2', 'XboxOne', 'Playstation4]", "count": 41},
 {"_id": "[quake', 'earthquake]", "count": 41},
 {"_id": "[Quake', 'Quake', 'DadJokes]", "count": 40},
 {"_id": "[QuakeChampions]", "count": 36},
 {"_id": "[FAV]", "count": 36}, ...]
```

```
In [33]: #Question4 What people, places or things are trending in your domain? (A trend is popularity over time.)  
list(db.game_twitter.aggregate([  
    {"$sort": {"created_at": 1}}]))  
[{"_id": ObjectId('5ad55a3f4ale94208242aefc'),  
    "id": 885,  
    "created_at": "2018-04-03 12:00:03",  
    "tweets": "Double Cross @arianstudios #Divinity Original Sin 2 https://t.co/t6sWDWNIHn",  
    "user_id": 754208650319630337,  
    "user_locatiom": 'Japan',  
    "retweet_count": 0,  
    "Urls": "[{'url': 'https://t.co/t6sWDWNIHn', 'expanded_url': 'http://youtu.be/7P7gbUfBJmw?a', 'display_url': 'youtu.be/7P7gbUfBJmw?a', 'indices': [53, 76]}]",  
    "Media": None,  
    "game": 'Left 4 Dead',  
    "array_agg": "'[Divinity']"},  
    {"_id": ObjectId('5ad55a3f4ale94208242aefd'),  
    "id": 1360,  
    "created_at": "2018-04-03 12:00:03",  
    "tweets": "Double Cross @arianstudios #Divinity Original Sin 2 https://t.co/t6sWDWNIHn",  
    "user_id": 754208650319630337,  
    "user_locatiom": 'Japan',  
    "retweet_count": 0,  
    "Urls": "[{'url': 'https://t.co/t6sWDWNIHn', 'expanded_url': 'http://youtu.be/7P7gbUfBJmw?a', 'display_url': 'youtu.be/7P7gbUfBJmw?a', 'indices': [53, 76]}]"}]
```

```
In [80]: for doc in collection.find({}, {"game":1, "_id":0, "created_at":1}).sort('created_at', pymongo.ASCENDING):
    print(doc)
```

```
['created_at': '2018-04-03 12:00:03', 'game': 'Counter-Strike'},
 {'created_at': '2018-04-03 12:00:03', 'game': 'Left 4 Dead'},
 {'created_at': '2018-04-03 12:00:03', 'game': 'Dota 2'},
 {'created_at': '2018-04-03 12:00:03', 'game': 'Empire: Total War'},
 {'created_at': '2018-04-03 15:44:53', 'game': 'Counter-Strike'},
 {'created_at': '2018-04-03 15:44:53', 'game': 'Left 4 Dead'},
 {'created_at': '2018-04-03 15:44:53', 'game': 'Dota 2'},
 {'created_at': '2018-04-03 15:44:53', 'game': 'Fallout 3'},
 {'created_at': '2018-04-03 15:44:53', 'game': 'Empire: Total War'},
 {'created_at': '2018-04-03 15:44:53', 'game': 'Call of Duty'},
 {'created_at': '2018-04-03 15:44:53', 'game': 'NBA 2K13'},
 {'created_at': '2018-04-03 15:44:53', 'game': 'Divinity 2'},
 {'created_at': '2018-04-03 15:44:53', 'game': 'Quake'},
 {'created_at': '2018-04-03 17:50:05', 'game': 'Counter-Strike'},
 {'created_at': '2018-04-03 17:50:05', 'game': 'Left 4 Dead'},
 {'created_at': '2018-04-03 17:50:05', 'game': 'Dota 2'},
 {'created_at': '2018-04-03 17:50:05', 'game': 'Fallout 3'},
 {'created_at': '2018-04-03 17:50:05', 'game': 'Empire: Total War'},
 {'created_at': '2018-04-03 17:50:05', 'game': 'Call of Duty'}
```

```
In [84]: list(db.game_twitter.aggregate([
    {"$match": {"game": "Left 4 Dead"}},
    {"$sort": {"created_at": 1}}]))
```

```
Out[84]: [ {_id: ObjectId('5ad55a3f4ale94208242ad91'),
  'id': 885,
  'created_at': '2018-04-03 12:00:03',
  'tweets': 'Double Cross @larianstudios #Divinity Original Sin 2 https://t.co/t6sWDWNiHn',
  'user_id': 754208650319630337,
  'user_location': 'Japan',
  'retweet_count': 0,
  'Urls': "[{'url': 'https://t.co/t6sWDWNiHn', 'expanded_url': 'http://youtu.be/7P7gbUfBJmw?a', 'display_url': 'youtu.be/7P7gbUfBJmw?a', 'indices': [53, 76]}]",
  'Media': None,
  'game': 'Left 4 Dead',
  'array_agg': ["'Divinity'"]},
  {_id: ObjectId('5ad55a3f4ale94208242ad90'),
  'id': 884,
  'created_at': '2018-04-03 15:44:53',
  'tweets': "Je viens de m'acheter #Divinity Original Sin 2. J'ai fait le premier sur PS4 et j'ai ador\u00e2\x83\u00a0... j'esperre que je vais\u00e2\x83\u00a0x80\u00a9 https://t.co/GEpMQNxzxK",
  'user_id': 2706780463,
  'user_location': 'Nowhere',
```

```
In [92]: #Question5 use case
for doc in collection.find({}, {"game":1, "_id":0, "tweets":1}).sort('game', pymongo.ASCENDING):
    print(doc)
```

```
In [105]: list(collection.find({"game": "Empire: Total War"}).limit(10))
```

```
In [106]: #return this cursor with a limit of 5 and skip of 20 applied
list(collection.find()[20:25])

Out[106]: [{"_id": ObjectId('5ad55a3f4ale94208242ab0f'),
   'id': 110,
   'created_at': '2018-04-10 12:31:23',
   'tweets': 'Trump=Attorney-Client Privilege DEAD \n#Facebook=Sells 2B Users Info & Peddles 2 Doctors\nCA Legislature=Ban FreeSpeech https://t.co/NsX3L98Xdc',
   'user_id': 582140844,
   'user_location': 'USA',
   'retweet_count': 2,
   'Urls': "[{'url': 'https://t.co/NsX3L98Xdc', 'expanded_url': 'https://twitter.com/i/web/status/983683862965489665',
   'display_url': 'twitter.com/i/web/status/983683862965489665', 'indices': [121, 144]}]",
   'Media': None,
   'game': 'Counter-Strike',
   'array_agg': "[ 'Facebook' ]",
   {'_id': ObjectId('5ad55a3f4ale94208242ab10'),
   'id': 111,
   'created_at': '2018-04-10 12:31:02',
   'tweets': 'Trump=Attorney-Client Privilege DEAD \n#Facebook=Sells 2B Users Info & Peddles 2 Doctors\nCA Legislature=Ban FreeSpeech https://t.co/92VNB0qS0',
   'user_id': 582140844, ...}

In [99]: list(collection.find({
   "created_at": {
      "$gte": "2018-04-10 12:00:06",
      "$lte": "2018-04-10 12:30:06"
   }
}))
```

```
Out[99]: [{"_id": ObjectId('5ad55a3f4ale94208242ab8b'),
   'id': 239,
   'created_at': '2018-04-10 12:01:26',
   'tweets': 'Fallout 3 in Fallout 4 mod cancelled due to voice acting copyright https://t.co/w3QtPSttd7 #Fallout #Fallout',
   'user_id': 901492458307112962,
   'user_location': None,
   'retweet_count': 0,
   'Urls': "[{'url': 'https://t.co/w3QtPSttd7', 'expanded_url': 'https://capitalwasteland.com/#cancellation-annoucement-section', 'display_url': 'capitalwasteland.com/#cancellation-Announcement', 'indices': [67, 90]}]",
   'Media': None,
   'game': 'Counter-Strike',
   'array_agg': "[ 'Fallout', 'Fallout' ]",
   {'_id': ObjectId('5ad55a3f4ale94208242abd8'), ...}]
```