# Package 'gatherNet'

## Claudia Kann

### October 6, 2022

**Version** 0.1
**Date** October 6, 2022
**Title** A Python Package for Collecting Social Media Data on Online Events
**Author** Claudia Kann (ckann@caltech.edu), Sarah Hashash ( shashash@caltec.edu)
**Maintainer** Claudia Kann (ckann@caltech.edu)
**Depends**
**Description**
**License**
**Repository**
**Date/Publication**

# Contents

# Introduction

This manual will walk you through the basic functionality of the package 'gatherNet' designed to collect Twitter networks of dynamic movements. Throughout, the code included should be sufficient to get you started. This package is created to be used with the Twitter Version 2 API as of September 2022.

# Generating Movements

There are two main ways to import a series of events that make up a movement. The first, which is recommended for beginners is using the excel sheet "movement_template.xlsx" found in the repo. The second is to manually import the information into the movement structure. Regardless, the first step is to instantiate your movement object:

```python
import twitterMovments as tw
movement = tw.TwitterMovment(name,
                    base_directory = '',
                    separator = 'CityTown')
```

| | |
|---|---|
| `name` | name for the movement |
| `base_directory` | path to where all of the data will be stored |
| `separator` | describes how the events are split in the data saving structure, it can be any of the event objects. For geographical differences 'CityTown' makes the most sense. |

## From Excel File

To fill out the excel sheet, you must include the values below. The bolded values are optional:

- events:

    - ID: unique identification string
    - starting_date: beginning date of the event within the movement in the format dd mm yyyy hh:mm:ss
    - ending_date: ending date of the event within the movement in the format dd mm yyyy hh:mm:ss
    - CityTown: City or Town where you want to search for individuals in the nucleus
    - StateTerritory: State or Territory where you want to search for individuals in the nucleus
    - Date: Date of the event in m/dd/yy format
    - BestGuess: best guess of the size of the event, if unknown, write 0
    - **period_start**: start of the movement (when you want the timeline gathering to start) if different from the rest of the events in the movement. In the format dd mm yyyy hh:mm:ss
    - **period_end**: end of the movement (when you want the timeline gathering to end) if different from the rest of the events in the movement. In the format dd mm yyyy hh:mm:ss

- keywords:

    - keywords: list of keywords you want to use to identify the nucleus. Can be written individually or in Twitter accepted format. For example, it is equivalent to have:

        George Floyd
        GeorgeFloyd
        vs.
        (GeorgeFloyd OR (George Floyd))

- time span:

  - start time: begining of the timeline gathering period in dd mm yyyy hh:mm:ss format
  - end time: end of the timeline gathering period in dd mm yyyy hh:mm:ss format

- keys: Twitter key and secret keys, make sure not to share these with other individuals

- location: coordinates to find bounding box for each location, a row needs to be added for each unique CityTown-StateTerritory pair found in the events tab

**upload_from_excel**

Once the excel file is filled in, it should be saved in the directory `dirr` as specified above. The built in function can then be used to populate the movement object:

```
movement.upload_from_excel(path = 'movement_template.xlsx')
```

| path | this should be the path to the excel file. If the path given is not an .xlsx file it will replace it with `base_directory` + `'movement_template.xlsx'`. If the path is `'movement_template.xlsx'`, it also assumes its in the `base_directory` |
| --- | --- |

**upload_from_file_structure**

If the data has already been pulled from Twitter but postprocessing needs to occur, the entire movement may not be necessary. In this case, the function `upload_from_file_structure` can be used in order to populate the movement sufficiently to work with the pulled data

```
movement.upload_from_file_structure()
```

## Manually

The movement can also be added into the object manually. In this case, all of the information must be added through a series of functions.

**TwitterEvent**

TwitterEvent objects must be created for each event within the movement:

```
event1 = tw.TwitterEvent(ID, starting_date, ending_date,
                         CityTown, StateTerritory, Date, BestGuess)
```

| ID | unique identification string |
| --- | --- |
| starting_date | beginning date of the event within the movement in the format dd mm yyyy hh:mm:ss |
| ending_date | ending date of the event within the movement in the format dd mm yyyy hh:mm:ss |
| CityTown | City or Town where you want to search for individuals in the nucleus |
| StateTerritory | State or Territory where you want to search for individuals in the nucleus |
| Date | Date of the event in m/dd/yy format |
| BestGuess | best guess of the size of the event, if unknown, write 0 |

If you want to pull a separate full timeline for individuals associated with this event, that information must be added as:

---
`event1.add_timing(period_start, period_end)`
---

| | |
|---|---|
| `period_start` | start of the movement (when you want the timeline gathering to start) if different from the rest of the events in the movement. In the format dd mm yyyy hh:mm:ss |
| `period_end` | end of the movement (when you want the timeline gathering to end) if different from the rest of the events in the movement. In the format dd mm yyyy hh:mm:ss |

**TwitterKeyPair**

For each set of Twitter keys, a TwitterKeyPair object must be created:

---
`key1 = tw.TwitterKeyPair(key, secret)`
---

| | |
|---|---|
| `key` | Twitter developer key as a string |
| `secret` | Twitter developer secret key as a string |

More information for getting these credentials can be found in the Twitter Documentation

**add_event**

---
`movement.add_event(events)`
---

| | |
|---|---|
| `events` | either a TwitterEvent object or a list of TwitterEvent objects |

**add_key**

---
`movement.add_key(keypair)`
---

| | |
|---|---|
| `keypair` | either a TwitterKeyPair object or a list of TwitterKeyPair objects |

**add_keyWords**

---
`movement.add_keyWords(words)`
---

| | |
|---|---|
| `words` | list of keywords you want to use to identify the nucleus. Can be written individually or in Twitter accepted format. For example, it is equivalent to have: |

<div align="center">

George Floyd

GeorgeFloyd

vs.

(GeorgeFloyd OR (George Floyd))

</div>

**add_timing**

---
`movement.add_timing(start, end)`
---

| | |
|---|---|
| `start` | begining of the timeline gathering period in dd mm yyyy hh:mm:ss format |
| `end` | end of the timeline gathering period in dd mm yyyy hh:mm:ss format |

**add_location**

---

`movement.add_location(CityTown, StateTerritory, west, south, east, north)`

---

| | |
|---|---|
| `CityTown` | list of CityTown entries that occur in the events |
| `StateTerritory` | list of StateTerritories associate with CityTown list |
| `west` | list of west latitude for bounding box of CityTowns |
| `south` | list of south longitude for bounding box of CityTowns |
| `east` | list of east latitude for bounding box of CityTowns |
| `north` | list of north longitude for bounding box of CityTowns |

**print_protests**

---

`protest_ids = movement.print_protests()`

---

| | |
|---|---|
| `protest_ids` | a dataframe of the events with relevant information. Run after generating the movement in order to remove incomplete entries and create a check of whether the movement is presenting as expected. |

# Pulling Data from Social Media

Once the movement is created, the user can begin pulling the data from Twitter. The main function used here is `get_tweets`, the rest are used within it but can be accessed by the precocious user. In order to use any of these functions you must have a movement object, the function is then called as `movement.pull.function()`.

## Main Functions

**get_tweets**

---

```
movement.pull.get_tweets(movement,
          types = [],
          max_results = 500,
          tweets_per_file = 1000,
          expansions = ['author_id', 'in_reply_to_user_id'],
          tweetfields = ['author_id','created_at', 'geo',
          'entities','public_metrics', 'text','referenced_tweets'])
```

---

| | |
|---|---|
| `movement` | TwitterMovement object types: a list of the types of Tweets and users you want to collect, options are: ['Nucleus', 'NucleusTimeline', 'Echos', 'EchosTimeline', 'Influences', 'InfluencesTimeline']. If left as an empty list, all will be evaluated. In order to run the Timeline versions, all events are checked to make sure they have the base version. For 'Echos' and 'Influences' it checks that a nucleus exists. Everything but the 'Nucleus' can restart after being interrupted with minimal redundancy. |
| `max_results` | number of tweets to attempt to pull in each query, must be an integer between 1 and 500 |
| `tweets_per_file` | number of tweets to save per file before beginning a new file |

| | | |
|---|---|---|
| expansions | refers to which of the tweetfields the user would like more information on, taken from Twitter Documentation are: | |
| | author_id | *Returns a user object representing the Tweet's author* |
| | referenced_tweets.id | *Returns a Tweet object that this Tweet is referencing (either as a Retweet, Quoted Tweet, or reply)* |
| | in_reply_to_user_id | *Returns a user object representing the Tweet author this requested Tweet is a reply of* |
| | attachments.media _keys | *Returns a media object representing the images, videos, GIFs included in the Tweet* |
| | attachments.poll_ids | *Returns a poll object containing metadata for the poll included in the Tweet* |
| | geo.place_id | *Returns a place object containing metadata for the location tagged in the Tweet* |
| | entities.mentions. username | *Returns a user object for the user mentioned in the Tweet* |
| | referenced_tweets. id.author_id | *Returns a user object for the author of the referenced Tweet* |
| tweetfields | the values within each tweet to be returned from each call are taken from Twitter Documentation: | |
| | id (default) | *The unique identifier of the requested Tweet.* |
| | text (default) | *The actual UTF-8 text of the Tweet. See twitter-text for details on what characters are currently considered valid.* |
| | attachments | *Specifies the type of attachments (if any) present in this Tweet.* |
| | author_id | *The unique identifier of the User who posted this Tweet.* |
| | context_annotations | *Contains context annotations for the Tweet.* |
| | conversation_id | *The Tweet ID of the original Tweet of the conversation (which includes direct replies, replies of replies).* |
| | created_at | *Creation time of the Tweet.* |
| | entities | *Entities which have been parsed out of the text of the Tweet. Additionally see entities in Twitter Objects.* |
| | geo | *Contains details about the location tagged by the user in this Tweet, if they specified one.* |
| | in_reply_to_user_id | *If the represented Tweet is a reply, this field will contain the original Tweet's author ID. This will not necessarily always be the user directly mentioned in the Tweet.* |
| | lang | *Language of the Tweet, if detected by Twitter. Returned as a BCP47 language tag.* |
| | non_public_metrics | *Non-public engagement metrics for the Tweet at the time of the request.* |
| | organic_metrics | *Engagement metrics, tracked in an organic context, for the Tweet at the time of the request.* |

| | | |
|---|---|---|
| | possibly_sensitive | *This field only surfaces when a Tweet contains a link. The meaning of the field doesn't pertain to the Tweet content itbut instead it is an indicator that the URL contained in the Tweet may contain content or media identified as sensitive content.* |
| | promoted_metrics | *Engagement metrics, tracked in a promoted context, for the Tweet at the time of the request.* |
| | public_metrics | *Public engagement metrics for the Tweet at the time of the request* |
| | referenced_tweets | *A list of Tweets this Tweet refers to. For example, if the parent Tweet is a Retweet, a Retweet with comment (also known as Quoted Tweet) or a Reply, it will include the related Tweet referenced to by its parent.* |
| | reply_settings | *Shows you who can reply to a given Tweet. Fields returned are "everyone", "mentioned_users", and "followers".* |
| | source | *The name of the app the user Tweeted from.* |
| | withheld | *When present, contains withholding details for withheld content.* |

## Manually

All of these functions use the inputs in the Main Function section, refer above for help.

### version_2_setup

```
APIs, PARAMS = movement.pull.version_2_setup(movement,
                    max_results = 500,
                    expansions = ['author_id', 'in_reply_to_user_id'],
                    tweetfields = ['author_id','created_at', 'geo',
                                'entities','public_metrics', 'text',
                                'referenced_tweets'])
```

| | |
|---|---|
| APIs | list of connections to Twitter using TwitterAPI and the key and secret keys provided |
| PARAMS | dictionary of parameters that are fed to |

### get_nucleus_users

```
movement.pull.get_nucleus_users(movement,
                            PARAMS,
                            APIs,
                            tweets_per_file = 1000)
```

### get_nucleus_timeline

```
movement.pull.get_nucleus_timeline(movement,
                            PARAMS,
                            APIs,
                            tweets_per_file = 1000)
```

**get_echo_users**

```
movement.pull.get_echo_users(movement,
                             PARAMS,
                             APIs,
                             tweets_per_file = 1000):
```

**get_echo_timeline**

```
movement.pull.get_echo_timeline(movement,
                                PARAMS,
                                APIs,
                                tweets_per_file = 1000):
```

**get_influence_users**

```
movement.pull.get_influence_users(movement,
                                  PARAMS,
                                  APIs,
                                  tweets_per_file = 1000):
```

**get_influence_timeline**

```
movement.pull.get_influence_timeline(movement,
                                     PARAMS,
                                     APIs,
                                     tweets_per_file = 1000):
```

**get_ids**

```
movement.pull.get_ids(path)
```

**pull_tweets**

```
movement.pull.pull_tweets(PARAMS,
                          datestart,
                          APIs,
                          api_i,
                          file_prefix,
                          tweets_per_file,
                          echo = False,
                          influence = False):
```

# Checking Data

**number_of_tweets**

```
movement.check.number_of_tweets(events,
                                users,
                                base_directory = "",
                                separator = 'CityTown',
                                timeline = False)
```

| | |
|---|---|
| `events` | a list of the events that you want to count the tweets in |
| `users` | the types of users you want to count tweets from, options are ['Nucleus', 'Echo', 'Influence'] |
| `base_directory` | where to begin searching for the data |
| `separator` | describes how the events are split in the data saving structure, it can be any of the event objects. For geographical differences 'CityTown' makes the most sense. |
| `Timeline` | Boolean of whether to count tweets in the timeline or not, if not, only Nucleaus original tweets will be counted |

## Reading Data

**read_tweets**

```
movement.read.read_tweets(movement, users = "Nucleus")
```

| | |
|---|---|
| `movement` | TwitterMovement object |
| `users` | which tweets do you want to look at, can be a list or an individual, options are: "Nucleus", "NucleusTimeline", "EchoTimeline", "InflunceTimeline" |

A csv of all the relevant tweets will be added to the `movement.base_directory`