

# #Charlottesville on Twitter

Vincent La

August 17, 2017

## Data Description

### Source

This dataset is a collection of tweets taken from the Twitter Streaming API. This is a continuous collection of tweets mentioning "Charlottesville" or using the hashtag #Charlottesville. However, a small amount of tweets were lost due to programming bugs and intermittent connection failures.

### License

I am distributing this dataset under the terms of the CC BY-SA 4.0. Furthermore, Twitter also requests that usage of this data abide by the Twitter Developer Agreement. Most notably, you should display individual tweets in accordance with Twitter's display policy.

### The Files

Each file contains a random sample of 50,000 tweets (in accordance with the Twitter Developer Agreement) from each day. It should be noted that due to programming bugs and intermittent connection failures, a small number of tweets were not collected. Therefore, these samples may potentially be less than truly random. Furthermore, because I started collecting data on August 15, that day's sample only includes tweets after 9PM Eastern Time.

### Attributes

Since the vast majority of attributes are unmodified and self-explanatory, I'm only going to describe the less obvious ones and the attributes I either created or cleaned (there's only two). For the rest, I will describe what attributes from the Twitter API they came from. An overview of tweet attributes can be found here on Twitter's website.

| Attribute            | Source                    | Description   |
|----------------------|---------------------------|---|
| id                   | Unmodified                | Integer corresponding to the Tweet ID   |
| user_id              | user -> 'id'              | Twitter user name   |
| user_name            | user -> 'name'            |   |
| screen_name          | user -> 'screen_name'     |   |
| user_statuses_count  | user -> 'statuses_count'  |   |
| user_favorites_count | user -> 'favorites_count' |   |
| friends_count        | user -> 'friends_count'   |   |
| followers_count      | user -> 'followers_count' |   |
| user_description     | user -> 'description'     | How the user chooses to describe them self  |
| user_location        | user -> 'location'        | <b>Note:</b> Twitter places no restrictions on what users can enter as their location |

|                               |  |   |
|-------------------------------|--|---|
| user_time_zone                | user -> 'time_zone'                          |   |
| user_profile_text_color       | user -> 'profile_text_color'                 |   |
| user_profile_background_color | user -> 'profile_background_color'           |   |
| full_text                     | Either text or extended_tweet -> 'full_text' |   |
| created_at                    | Unmodified                                   | UTC timestamp of when Tweet was posted. For reference, Eastern Standard Time is 4 hours behind UTC. |
| is_retweet                    |  | A binary variable I created to help me subset the data. Not very useful outside of that.            |
| retweeted_status_text         | retweeted_status -> 'text'                   |   |
| retweeted_status_id           | retweeted_status -> 'id'                     |   |
| quoted_status_text            | quoted_status -> 'text'                      | The text of the tweet that this status referenced (if applicable)                                   |
| quoted_status_id              | quoted_status -> 'id'                        |   |
| in_reply_to_status_id         | Unmodified                                   |   |
| in_reply_to_user_id           | Unmodified                                   |   |
| hashtags                      | entities -> 'hashtags'                       | I used a Postgres function to flatten out the JSON array which contained the list of hashtags.      |

---