

# Report for Data gathering, assessing, cleaning, and storing parts of twitter project

## Data gathering:

- **api\_df:**
  - Collected data from twitter archive using twitter API and saved the data json information in a text file.
  - Looped over the data text file line by line, read each line using `json.loads()` function, extracted the variables I'm interested in (tweet\_id, retweet\_count, favorite\_count, user\_followers\_count) and saved them with their values as key value pairs in a dictionary object. All dictionary objects were stored in an array. The final result is a list of dictionaries where each dictionary represents a tweet.
  - Converted the list of dictionaries into a dataframe using pandas `DataFrame` function and named it `api_df`.
  - Note: my twitter developer account wasn't activated by the time I completed this project so I used the code and text file provided by udacity mentors.
- **archive\_df:**
  - read `twitter-archive-enhanced.csv` file into a dataframe named `archive_df`.
- **image\_predictions\_df:**
  - read `image-predictions.tsv` file into a dataframe named `image_predictions_df` specifying the separator as tab separator (`\t`).

## Data assessment:

- **api\_df:**
  - Displayed the table head and info to inspect variables datatypes and values.
  - **Quality Issues:**
    - `tweet_id` should be string not int
- **archive\_df:**

- Displayed tab head and info to inspect variable datatypes and values to find quality and tidiness issues.
- Displayed values of certain variables as name, source, and rating denominator to further inspect them and find null and wrong values.
- **Quality Issues:**
  - some names are uppercase and some are lowercase.
  - some dogs have none as names.
  - tweet\_id should be string not int.
  - source column contains whole html element instead of source name only.
  - tweets with expanded\_urls value of null should be removed.
  - rating\_numerator, rating\_denominator are int when they should be floats.
  - timestamp is string not datetime.
  - retweeted\_status\_id, retweeted\_status\_user\_id, retweeted\_status\_timestamp all mostly nulls so should be removed.
  - in\_reply\_to\_status\_id, in\_reply\_to\_user\_id are mostly nulls so should be removed.
- **Tidiness Issues:**
  - four dogs columns should be combined into one.
- **image\_predictions\_df:**
  - Displayed tab head and info to inspect variable datatypes and values to find quality and tidiness issues.
  - **Quality Issues:**
    - tweet\_id should be string not int.
    - rename jpg\_url column name to image\_url
- **Tidiness Issue:** api\_df, archive\_df, image\_predictions\_df should all be merged into one master table named twitter\_archive\_master.

## Data cleaning:

- Followed the define-code-test pattern to define my solution and code it then test the output to check if it aligns with my solution definition or not.
- **archive\_df:**
  - **Quality Issues:**
    - Changed tweet\_id from int datatype to string to eliminate any accidental ordering or misunderstanding that any order is implied through the id.
    - make all dog names lowercase
    - change all dog names that are none to np.nan.
    - change all denominator values to 10.
    - make timestamp datatype datetime instead of string
    - make source column contain only the source name instead of the whole html element.
    - remove retweeted\_status\_id, retweeted\_status\_user\_id, retweeted\_status\_timestamp, in\_reply\_to\_status\_id, in\_reply\_to\_user\_id columns
    - drop tweets with null expanded\_urls values.
    - changed rating\_denominator and rating\_numerator datatypes from int to float.
  - **Tidiness Issues:**
    - combine 4 dogs columns [doggo, floofer, pupper, puppo] into one column named dog\_stage
- **image\_predictions\_df:**
  - **Quality Issues:**
    - Changed tweet\_id from int.
    - Renamed jpg\_url column to image\_url to describe better the content of the column (url of the image being classified by the network).
- **api\_df:**
  - **Quality Issues:**

- Changed tweet\_id from int.

### **Data storing:**

- **twitter\_archive\_master\_df:**
  - Stored master table into csv file named twitter\_archive\_master.csv.
- **archive\_df:**
  - Stored clean version of archive\_df table into csv file named twitter\_archive\_clean.csv.
- **image\_predictions\_df:**
  - Stored clean version of image\_predictions\_df table into csv file named image\_predictions\_clean.csv.