

## Wrangle Report –

For this analysis, I wrangled three different datasets. The first was a dataset I downloaded using the requests library from [https://d17h27t6h515a5.cloudfront.net/topher/2017/August/599fd2ad\\_image-predictions/image-predictions.tsv](https://d17h27t6h515a5.cloudfront.net/topher/2017/August/599fd2ad_image-predictions/image-predictions.tsv). This was the dataset which included the breed predictions from the images on the weratedogs tweets.

The second dataset used was provided in the project materials and was read into a dataframe from the csv file. This was the twitter archive dataset.

2. The third dataset used was a dataset created from api requests from the Twitter API developer portal. To create this dataset, I created a twitter developer account and used the Tweepy library to make api "df\_twitter\_archive" columns "doggo", "floofer", "pupper", and "puppo" have "None" listed for all values. These columns hold no value and should just be dropped. (Tidiness issue because columns should store actual/useful variables and these columns store nothing).

requests of all the tweet ids from the other two dataframes.

Due to the request limit from the twitter api, I made a function which paused the requests for fifteen minutes every 900 requests.

I saved these requests in a json format in a .txt file, then created another function which read the file line by line and parsed the json format to extract only the tweet ids, favorite counts, and retweet counts for the dataframe.

I used both visual and programmatic assessment to see what needed to be cleaned in the datasets. I found the following issues and addressed each of them:

1. In "df\_twitter\_archive", tweet ids were listed in descending order (other two dataframes are listed in tweet id ascending order). I fixed this by changing the order to ascending to match the other two dataframes.
2. In "df\_image\_predictions", breed names were not uniform in their casing (some start with uppercase, some lowercase). I used the str.lower() method to change all of the casing to lowercase to be uniform.
3. NaN values in "df\_twitter\_archive" dataframe. I ended up dropping several columns including "in\_reply\_to\_status\_id", "in\_reply\_to\_user\_id", "retweeted\_status\_id", "retweeted\_status\_user\_id", and "retweeted\_status\_timestamp." These columns were ultimately dropped because the vast majority of the rows were NaN, so they were not useful for further analysis across the datasets.
4. Removed ids missing from "api\_tweets\_df" from "df\_image\_predictions." I trimmed this dataset down to match the tweet\_ids in the dataset requested from the twitter api.

5. Removed ids missing from "api\_tweets\_df" from "df\_twitter\_archive." I trimmed this dataset down to match the tweet\_ids in the dataset requested from the twitter api.
6. Column names in "df\_image\_predictions" were not descriptive. I used the project motivation documentation to find the meanings of the column headers (such as p1, p1\_conf, etc) and renamed the columns to be more descriptive so they were understandable to the end user.
7. In "df\_twitter\_archive", "Timestamp" column was just a string/object, not a datetime type. I used the pandas to\_datetime method to convert this into a datetime format instead of just a string.
8. In "df\_twitter\_archive" "name" column - several values were obviously not names. I manually looked through the unique values in the name column and removed any values which were obviously not names (such as "a", "an", "unacceptable", etc). I replaced these values in the name column with the string "NOT FOUND," as the name was unable to be parsed correctly.

Tidiness issues:

1. In "df\_twitter\_archive," I separated the datetime object into two separate columns – date and time.
2. In "df\_twitter\_archive" columns "doggo", "floofer", "pupper", and "puppo" had "None" listed for all values. These columns hold no value and were dropped completely. (Tidiness issue because columns should store actual/useful variables and these columns store nothing).
3. I separated the df\_twitter\_archive into the following datasets – "tweet\_times" with the id, posted date, and posted time, "tweet\_ratings" with the id, rating numerator, and rating denominator, "tweet\_text" with the id, text, and dog name, and "tweet\_urls" with the id, source, and expanded urls.

I also created a cleaned master dataframe and saved it to a separate csv file, twitter\_archive\_master.csv

This concluded my wrangling and cleaning efforts for these datasets.