

For Data Wrangling of WeRateDogs, I first create Developer Account on Twitter to fetch Twitter Archive Tweets Data. Then I checked all libraries present with my Jupyter.

### **Gathering Data:**

I installed all required libraries and read Twitter archive CSV file, tweet image predictions TSV file. Create a logic to fetch tweets' data from Twitter using API with start and end time.

I then created list of dictionaries to read tweet's JSON data line by line and later convert to a DataFrame.

### **Assessing Data:**

I followed below checks on the data gathered:

- 1) If there are any records in tweets\_df which are retweets
- 2) If there are any records in tweets\_df whose corresponding record with same tweet\_id is missing in img\_df table
- 3) If there are any records in tweets\_df whose corresponding record with same tweet\_id is missing in status\_df table
- 4) Sort by rating\_denominator values
- 5) Sort by rating\_numerator values
- 6) Sort by names values

After Analysing, I got the following Quality and Tidy issues:

#### **Quality:**

- 1) contains retweets and therefore, duplicates
- 2) many tweet\_id(s) of tweets\_df table are missing in img\_df (image predictions) table
- 3) many tweet\_id(s) of tweets\_df table are missing in status\_df (Twitter Archive) table
- 4) erroneous datatypes (in\_reply\_to\_status\_id, in\_reply\_to\_user\_id and timestamp columns)
- 5) unnecessary html tags in source column in place of utility name e.g. <a href=""http://twitter.com/download/iphone"" rel=""nofollow"">Twitter for iPhone</a>
- 6) some records have more than one dog stage
- 7) erroneous dog names starting with lowercase characters (e.g. a, an, actually, by)
- 8) extract dog breed from predetection data into img\_df

#### **Tidiness:**

- 1) tweets\_df table without any duplicates (i.e. retweets) have empty retweeted\_status\_id, retweeted\_status\_user\_id and retweeted\_status\_timestamp columns, which can be dropped

- 2) doggo, floofer, pupper and puppo columns should be merged into one column named "stage"
- 3) merge all 3 dataframes to get single master file in archive\_clean
- 4) Dropping all extra columns which are not useful

### **Cleaning Data:**

Take a copy of tweets\_df on which the cleaning tasks will be performed and start removing the quality and tidiness issues. After cleaning and dropping the unrequired column, I merged the dataframe into single dataframe for further storing the file.

### **Storing Data:**

Dataframe is stored in twitter\_archive\_master.csv as a single file for further Analysis and Visualisation.

### **Analysing and Visualizing Data:**

For analysis, I created a copy of the cleaned twitter archive data. Following are the analysis and visualisations done:

- 1) Analysis of rating of dogs
- 2) Most used Twitter source
- 3) Analysis of famous Dog's Names