

# Project#4 - Wrangle and Analyze Data

**Developed By:** Rakan210 on June 17, 2019 - as part of Udacity Data Analyst Nanodegree program (DAND)

In this project, all the analysts were requested to unleash their analytic skills by walking through all data analysis process [Gather, Assess, Clean, Analyze, and Report]. The focus is at the data wrangling part, which includes gathering, assessing and cleaning the data. Below are the detailed requirements:

## Part#1: Gathering Data:-

This project required 3 sources of data to be gathered:

- The WeRateDogs Twitter archive, which is provided in the following link: [[twitter\\_archive\\_enhanced.csv](#)].
- The tweet image predictions, i.e., what breed of dog (or other object, animal, etc.) is present in each tweet according to a neural network. This file (image\_predictions.tsv) is hosted on Udacity's servers and should be downloaded programmatically using the Requests library and the following [URL](#).
- Each tweet's retweet count and favorite ("like") count at minimum, and any additional data you find interesting. Using the tweet IDs in the WeRateDogs Twitter archive, query the Twitter API for each tweet's JSON data using Python's Tweepy library and store each tweet's entire set of JSON data in a file called tweet\_json.txt. Each tweet's JSON data should be written to its own line. Then read this .txt file line by line into a pandas DataFrame with (at minimum) tweet ID, retweet count, and favorite count. Note: do not include your Twitter API keys, secrets, and tokens in your project submission.

## Part#2: Assessing Data:-

After gathering each of the above pieces of data, assess them visually and programmatically for quality and tidiness issues. Detect and document at least eight (8) quality issues and two (2) tidiness issues in your wrangle\_act.ipynb Jupyter Notebook.

Quality issues: (completeness, validity, accuracy, and consistency issues)

**Note: Visual assessment (VA) & Programmatic assessment (PA)**

**df\_archive:**

1. (PA) Change datatype of [timestamp] to datetime64 and [rating\_numerator] to float.
2. (VA) Rename column labels to a descriptive ones.
3. (VA) Change all 'None' values should be changed to NP.NaN.
4. (PA) Delete all Retweets/Replies records since they shouldn't be considered. ([retweeted\_status] & [in\_reply\_to] columns)

5. (VA) Modify some records that have inaccurate value in [rating\_denominator] as it reasonably to be fixed to or be multiple of "10" ( rating\_denominator %10 = 0).
6. (PA) Modify some records that have inaccurate value in [rating\_numerator] due to the existing of the dicimal.
7. (VA) Replace invalid dogs' names with NP.NaN.
8. (PA) Not all tweets have predictions photo (Archive Tweets "df\_archive": 2356 rows , Predection Images "df\_images": 2075 rows), so only tweet with images will be considered.
9. (PA) Delete all records that missing values in [expanded\_url] column as they indicate ratings without an image.

### Tidiness issues: (structural issues)

#### **df\_archive:**

1. Combine dog stages (doggo, floofer, pupper and puppo) into one column called [dog\_stage].
2. Combine rating columns into one column called [dog\_rating]

#### **df\_master:**

3. Combine all DataFrames into one master DataFrame called [df\_master]. Its content should be stored in twitter\_archive\_master.csv.

### Part#3: Cleaning Data:-

Detailed step-by-step cleaning process was done to overcome all the mentioned issues above in ` wrangle\_act.ipynb`. In general, each data frame was handled and cleaned by its own before merging them into a master data frame. After that, another cleaning process was done on the master data frame to finally polish the data and start the analysis and reporting stages.