

# Data Wrangling Report

## 1. Data Gathering

### About the Data:

The Dataset wrangled is the Tweet archive of Twitter account WeRateDogs--@dog\_rates. Three (3) different datasets are used.

- 1.The archive dataset which is the main dataset and contains 2356 tweet data.
2. The image predictions dataset which contains images based on image predictions of the twitter archive dataset and it contains the tweetID, image URL, and the image number with the most confident prediction
3. The tweet\_json.txt data from which Retweet count,Favourite Count,Display text range, and full text were extracted

### Gathering the data:

The process used to gather each of the datasets are:

1. Twitter archive file: This was downloaded manually using the link provided by Udacity.
2. Image Predictions: This was downloaded from the file hosted on udacity server using python's requests and locally saved it to image\_predictions.tsv file. The downloaded file was then imported using Pandas as img\_df
3. Tweet\_json : I downloaded the tweet\_json file as provided on the udacity servers.  
I then created a data frame 'tweet\_json\_df' from the the JSON file extracting the tweet\_id,'retweet\_count','favorite\_count','display\_text\_range','full\_text' columns.

## Assessing Data

The data quality issues were assessed based on

1. Completeness
2. Validity
3. Accuracy
4. Consistency

### Quality Issues:

- Inappropriate datatype for tweet\_id across 3 tables-- should be string not int.
- They are a lot of missing data in some features in the df\_archive table - e.g "expanded\_url" column.
- Missing values are also represented as "None" in name,doggo,floofer,pupper,puppo tables of df\_archive tables.
- The df\_archive table contains Retweets while its only original tweets that are needed
- Inappropriate datatype for timestamp column in df\_archive table-- should be datetime instead
  
- There are unnecessary html tags in the 'Source' column instead of just source name eg 'Twitter for iPhone' instead of <a href=""http://twitter.com/download/iphone"" rel=""nofollow"">Twitter for iPhone</a>
- The rating\_numerator of df\_archive table has some erroneously large numbers (e.g. 1176,960,420,182)
- The rating\_denominator column also has other values besides 10(2,150, 170)-- some are due to the text having 2 ratings (like) figures in it as seen from the programmatic assessment of denominators not 10.
- Dog names have inappropriate names eg(very,unacceptable,this,the,such), all starting with lower cases-- (a validity issue)
- Many tweet\_id(s) --186 tweet data rows from the archive table are missing in img\_df (image predictions) table.
- 2 tweet\_id(s) from the archive table are missing in tweet json table.

## Tidiness Issues:

- Dogs kind 4 stages in 4 diffedifferent columns(doggo,floofer,pupper and puppo) and they should be in one column.
- p1, p2, and p3 in the img\_df table are not properly formatted.
- All three tables need to be merged into 1.

## Data Cleaning

-- I made a copy of all 3 tables so as to preserve the original data in cases of mistake.

-- Also Since they all have a quality or tidiness issues one way or the other, and would be needed to be merged into one table, cleaning the 3 datasets was necessary.

-- I performed a programmatic data cleaning process using the Define, COde, Test framework on the copied version of the datasets addressing the issues alighted in my data assesment.

--I thereby merged all the 3 datasets in to 1 using a left join joining on the tweetID.

In [ ]:

## Data Storage

After successfully cleaning and merging all the datasets into one, I saved the clean dataframe to file as twitter\_archive\_master.csv.

In [ ]: