

## Wrangle Report :

### WeRateDogs Twitter Archive

---

#### Gathering Data:

---

1. Manually download the *twitter archive enhanced csv* file from Udacity. I create the dataframe named *twit\_arch*.
2. 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* file. Read this file and create a dataframe with tweet ID, retweet count, and favorite count named *df\_cleanjson*.
3. Systemically download *Image\_prediction* file and saved the dataframe named *df\_image*.

#### Assessing Data & Cleaning Data:

---

I assess those three dataframe visually and programmatically for quality and tidiness issues.

##### 1. *twit\_arch* dataframe

###### Quality

- Timestamp should be changed from object datatype into datetime datatype
- Removed the row that *in\_reply\_to\_status\_id* != NaN because that row just retweet the original content.
- Due to text column, there are some *rating\_numerator* should be float such as 75 should be 9.75. I will change *rating\_numerator* datatype from int into float.
- Some *rating\_numerator* and *rating\_denominator* are invalids (i.e. date, 7/11, etc.). I will extract the second xx/xx pattern to *rating\_numerator* and *rating\_denominator*.
- There are too various *rating\_denominator*. I will create the score out of ten column by using *rating\_numerator* divided by *rating\_denominator* and multiply by ten
- Removed duplicated *expanded\_urls* data
- Cleansing the data in *source* column to be more readable

- There are some incorrect name such as 'a', 'an', 'the', etc. I will use string patterns like named \* , and name is \* to find the real dog name. (The pattern This is\*, and , and Meet \* is already used to detect name in this dataframe). If the text does not contain dog names, I will change invalid names into 'None'.
- Remove the column will not be used for analyzing : in\_reply\_to\_status\_id, in\_reply\_to\_user\_id, retweeted\_status\_id, retweeted\_status\_user\_id, retweeted\_status\_timestamp, rating\_numerator, rating\_denominator

### **Tidiness**

- Merge all dog stages (doggo, floofer, pupper, puppo) into one column called 'dog\_stages' and make a 'multiple' stage value for those tweet\_id which contains more than 1 stage.

## **2. df\_cleanjson dataframe**

### **Quality**

- Rename id into tweet\_id

### **Tidiness**

- Merge df\_arch\_copy and df\_cleanjson\_copy and create a new dataframe named *combine*

## **3. df\_image dataframe**

### **Quality**

- Delete the duplicate jpg\_url
- Because I will use only p1 for analysis, I will remove unnecessary columns: img\_num, p1\_conf, p1\_dog, p2, p2\_conf, p2\_dog, p3, p3\_conf, p3\_dog
- Make all dog breeds prediction of p1 into capitalize name
- Rename p1 column to predict\_dogbreed

### **Tidiness**

- Merge combine and df\_image\_copy table and create a new dataframe named *full\_combine*

## Storing Data:

---

After merge all three tables name `full_combine`, I save it into *twitter\_archive\_master.csv* and I will use it for data analysis and visualization.