

# Data Wrangling Report:

---

## About Dataset:

The Dataset I wrangled is a twitter archive of a user known as **WeRateDogs** who rates the dogs with a amusing comment. They dont rate dogs with a specific fixed denominator , most of the times the rating numerator is more than denominator. this data consists of 2356 rows and it consists of all the tweet data from November 2015 to August 2017.

Based on the above images data we have another Dataset which contains the image predictions along with tweet ID and image URL and image number. It has top three predictions of dog breeds along with respective confidence in those predictions for every tweet.

## Data Gathering:

### Gathering from a csv file:

I just imported the **csv** file which is a twitter archive of WeRateDogs.I manually downloaded it from udacity projects lesson.The name of the file is **twitter\_archive\_enhanced.csv**.

### Gathering from a server:

I downloaded the image predictions tsv file from the server using python **requests** library .the link to file hosted in Udacity server is [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) I loaded the data from the tsv file into a dataframe using **read\_csv()** and **sep** attribute of that function.The name of this file is **image\_predictions.tsv**

### Gathering using Twitter API:

I gathered the required twitter data from **Twitter API** using **tweepy** access library in python and the tweet ID's from other the **twitter\_archive\_enhanced.csv** . By default it is JSON data and stored it in **tweet\_json.txt** using **json** library in python. and extracted the attributes(retweet\_count, favorite\_count and tweet ID) data of all the tweets from their respective tweet objects.

## Data Assessment:

### Visual Assessment:

I went through the whole DataFrame visually and found some issues in some datasets.

Here are the issues I found in each dataset visually:

#### **twitter\_enhanced:**

- None values in **name** column.
- hyperlinks in source column.
- None values in all the **doggo, floofer, pupper, puppo** columns.
- Multiple columns for dogstage\_\_(tidyness issue)\_\_

- Retweets are Present in the Dataset.

### image\_pred:

- Multiple DogBreeds are possible\_\_(tidyness issue)\_\_
- Unnecessary columns for analysis
- Dog breed names starting with small case alphabets

### Programatic Assessment:

I assessed all the dataframes using some pandas and numpy functions and found some more issues in the datasets.

Here are the issues I found in each dataset programatically:

### twitter\_enhanced:

- Erroneous names in **name** column.
- incorrect Numerators and Denominators for some of the tweets
- Datatypes of **tweet\_id,in\_reply\_to\_status\_id, in\_reply\_to\_user\_id, retweeted\_status\_id, retweeted\_status\_user\_id**.

### image\_pred:

- Datatype of **tweet\_id** is Int

## Cleaning Dataset:

As a first step I made copies of all the three original datasets so as to enable us to access the original datasets wherever required.I followed a organised **Define,Code and Test** programmatical clean process.

## Storing Data:

I joined all the dataframes performing "**Inner Join**" on tweet\_id programmatically using **merge()** function.Then, I stored the cleaned dataframe in a csv file named **twitter\_archive\_master.csv**