

# **Data Wrangling Report**

*By Amr Mohamed Helal*

*August 2020*

This report illustrates the main data analytics process steps to wrangle Twitter "We Rate Dogs" data.

## **Data Gathering:**

In this step we collected data from different resources:

1. twitter-archive-enhanced.csv file which is downloaded manually and then imported to our project notebook using "pandas.read\_csv()".
2. image-predictions.tsv file which is exported from a given web page link using "requests.get()" and then imported to our notebook using "pandas.read\_csv()".
3. tweet-json file which is extracted using tweepy API and then imported at our notebook using "json.loads()".

## **Data Assessment:**

In this step we performed both quality and tidiness assessment on our 3 data frames visually and programmatically as follows:

1. Quality Assessment:

### **a. archive\_df**

1. Should include only original ratings tweets (remove retweets and replies)
2. Missing values at "expanded\_urls" are for tweets without photos (could be dropped safely)
3. Most of records are null at these columns "in\_reply\_to\_status\_id, in\_reply\_to\_user\_id, retweeted\_status\_id, retweeted\_status\_user\_id, retweeted\_status\_timestamp" (could be dropped)
4. Datatype of "tweet\_id" is int (should be str format)

### **b. image\_predictions\_df**

1. Column name of "p1, p1\_conf, p1\_dog ..." are not descriptive
2. Datatype of "tweet\_id" is int (should be str format)
3. Tweets with no images should be dropped (total tweets at archive\_df more than images at image\_predictions\_df)

### **c. api\_df**

1. Most of records are null at these columns "contributors, coordinates, geo, in\_reply\_to\_screen\_name, in\_reply\_to\_status\_id, in\_reply\_to\_status\_id\_str, in\_reply\_to\_user\_id, in\_reply\_to\_user\_id\_str, place" (could be dropped)

## 2. Tidiness Assessment:

### a. **archive\_df**

1. Column names are values for "doggo, floofer, pupper, puppo"

### b. **api\_df**

2. Multiple variables are listed in one column "display\_text\_range" (should be splitted into 2 columns "display\_text\_min & display\_text\_max")

## **Data Cleaning:**

In this step we proceeded with Define, Code, Test process to define each solution we worked on to solve issues discovered in assessment step with 3 outcome data frames:

1. archive\_df\_clean
2. image\_predictions\_df\_clean
3. archive\_df\_clean

## **Data Storing:**

In this step we stored the final cleaned data frames in CSV files format as below:

1. Twitter\_archive\_master.csv
2. Twitter\_image\_predictions\_master.csv
3. Twitter\_api\_master.csv