

# Data Wrangling Report

## Project Objectives

- Perform data wrangling on datasets gathered from three different sources which includes: *'twitter-archive-enhanced.csv'*, *'tweet\_json.txt'*, and *'image\_predictions.tsv'*.
- Store and analyze *'twitter\_archive\_master.csv'*, created after wrangling and combining data from our three different data sources.

This report will focus mainly on the Data Wrangling aspect of this Project.

## Gathering Data

- The WeRateDogs Twitter archive (file on hand, manual download of *'twitter-archive-enhanced.csv'*).
- The tweet image predictions (*'image\_predictions.tsv'*). This file was be downloaded programmatically using the Requests library from a provided URL.
- Each tweet's entire set of JSON data (with at minimum tweet ID, retweet count, and favorite count) in a file called *'tweet\_json.txt'* were stored using Twitter API and Python's Tweepy library. Each tweet's JSON data was written to its own line.

## Accessing Data

During Data assessment both visually and programmatically, made the following observation.

Issue	Observation
Data Quality	<ol style="list-style-type: none"><li>1. df_api table - missing values in <i>geo</i> column.</li><li>2. df_api table - <i>possibly_sensitive_appealable</i> column does not contain any data.</li><li>3. df_archive table - does not seem to contain complete data for dog stages.</li><li>4. df_archive table - <i>timestamp</i> column is a string object instead of a datetime object.</li><li>5. df_archive table - missing <i>retweeted_status_id</i> values</li><li>6. df_archive table - missing <i>retweeted_status_user_id</i> values</li></ol>

	<ol style="list-style-type: none"> <li>7. df_img_pred table - column <i>p1</i> values start with varying cases (upper and lower).</li> <li>8. df_img_pred table - column <i>p2</i> values start with varying cases (upper and lower).</li> </ol>
Data Tidiness	<ol style="list-style-type: none"> <li>1. df_api and df_archive both have ratings and tweet URL in the <i>full text</i> column.</li> <li>2. df_api table - the same data in two columns <i>id</i> and <i>id_str</i>.</li> </ol>

## Cleaning Data

After data assessment, I went on to clean each DataFrame one after the other and the cleaning steps used are as followed

### df\_api

1. First created a copy of *df\_api* called *df\_api\_clean*.
2. Dropped redundant columns i.e. columns with over 70% NAN values.
3. Next step to cleaning *df\_api\_clean* was dropping single cardinality variables that do not offer any information to us i.e. columns with only one value.
4. During assessment, found that the "*id*" column is duplicated in *id\_str* so we drop *id\_str* column and rename *id* to *tweet\_id*.
5. Columns *entities*, *user*, and *extended\_entities* contain data in other columns/tables and appear to be very noisy, so dropped them.
6. During Data Tidiness Assessment, found that column *full\_text* contains the ratings and the tweet link alongside the tweet, extracted all three into different columns using Regular Expression and created a function *extract\_multi\_cols* in order to avoid repetition of this task.

### df\_archive

1. First created a copy of *df\_archive* called *df\_archive\_clean*.
2. Again, the first step of cleaning *df\_api\_clean* was to drop redundant columns i.e columns with over 70% NAN values.

3. Convert column *timestamp* to a pandas Datetime object as observed during assessment.
4. Rename column *text* to *full\_text*.
5. Use *extract\_multi\_cols* function to solve the same tidiness issues we saw while cleaning *df\_api*.

### **df\_img\_pred**

1. First, make a copy of *df\_img\_pred* called *df\_pred\_clean*.
2. Convert the dog breed names in columns *p1*, *p2* and *p3* to lowercase.

### Storing Data

Combine our three cleaned Dataframes and store it as a csv file ***twitter\_archive\_master.csv***.