

# Data Wrangling Report

## Project objectives

The project main objectives were:

- Perform data wrangling (gathering, assessing and cleaning) on the provided sources of data.
- Store, analyze, and visualize the wrangled data.

## Step 1: Gathering Data

In this phase, the three pieces of data were gathered and represented as pandas dataframes:

- The WeRateDogs Twitter archive (file on hand, manual download of 'twitter-archiveenhanced.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.

## Step 2 and 3: Assessing and Cleaning Data

While working with data, a number of observations were made. In the below table there are the observations along with actions taken in the Cleaning Step.

Dataset	Observation	Solution
Twitter Archive	<p><b>Quality</b></p> <ul style="list-style-type: none"><li>• some of the column names are not meaningful (like: retweet, reply to the original tweet, .. )</li><li>• timestamp column is str instead of datetime</li></ul> <p><b>Tidyness (structure)</b></p> <ul style="list-style-type: none"><li>• more than one stage is filled for a particular dog</li><li>• "source" and "expanded_urls" have several informations inside them.</li><li>• doggo, floofer, pupper and puppo refer to the same dog stage.</li><li>• rating_numerators should be of type float and</li></ul>	<ul style="list-style-type: none"><li>- Replaced Non values with np.nan.</li><li>- Removed the rating score and tweet link from the tweets text column using RegEx and pandas extract method.</li><li>- Converted timestamp to datetime data type using pandas to_datetime function.</li><li>- Removed retweets rows from data.</li><li>- Removed replys rows from data.</li><li>- Extracted the rating score correctly and converted it to float.</li><li>- Removed any rows with denominator more than 10.</li></ul>

	are not always correctly extracted.	<ul style="list-style-type: none"> <li>- Removed rows with missing expanded urls as they are not valid data.</li> <li>- Replaced None and invalid names with np.nan.</li> </ul>
Image Predictions	<b>Quality</b> <ul style="list-style-type: none"> <li>• there are 2075 tweet id, and the archive dataset has a total of 2356 ids which means 281 IDs are missing.</li> <li>• columns name are not the best thing.</li> <li>• p1, p2, and p3 contain underscores instead of spaces in the labels.</li> <li>• img_num is not needed.</li> </ul>	<ul style="list-style-type: none"> <li>- Removed other columns</li> <li>- Renamed it to match the other 2 datasets</li> </ul>
JSON File	<b>Quality</b> <ul style="list-style-type: none"> <li>• the original twitter_arch has 2356 tweet_id and JSON file 2354 ( number of missing IDs = <math>2356 - 2354 = 2</math>)</li> </ul>	
All	All datasets should be combined into 1 dataset only	Combined all the 3 datasets into one pandas df

## Result

A combined data set with all needed information was stored in CSV file called “[twitter\\_archive\\_master\\_new.csv](#)”