

# Wrangle Report

## WeRateDogs (Twitter Archive)

### CONTENT

#### Gather Data

1. Enhanced Twitter Archive
2. Image Predictions File
3. Download Tweet JSON Data

#### Assess Data

1. Twitter Archive
2. Image Predictions
3. Tweet\_Data

#### Clean Data

1. Twitter Archive
2. Image Predictions
3. Tweet\_Data

---

#### Gather Data

I gathered data from (**3 sources**), stored in separate files:

1. WeRateDogs (Twitter Enhanced archive), manually (downloaded from the Udacity servers).
2. The image predictions file, programmatically downloaded from the Udacity servers.
3. The entire set of each tweets' JSON data, downloaded by querying the Twitter API using the Tweepy library. The favourite\_count and retweet\_count was extracted programmatically from this file. I loaded the 3 raw data files into separate tables: archive, predictions and json\_data.

#### Assess Data & Clean Data

Quality 1. Twitter Archive: -

- 181 retweets
- 78 reply tweets
- 2297 tweets with expanded\_urls
- timestamp column is in string format
- NOT a valid name 109 tweets
- 17 tweets with rating\_denominator NOT equal to 10

Quality 2. image prediction: -

- 2075 image predictions, 281 less than

Quality 3. Json: -

- TweetErrors

tidiness 1. archive: -

- There are 4 columns for dog stages
- Key Points indicates
- rating\_denominators

tidiness 2. image predictions: -

- p1, p2 and p3 contain the same type of data

tidiness 3. json\_data

- combine json\_data with archive table
-