# REPORT: WRANGLE AND ANALYSE DATA PROJECT



*Image via Boston Magazine*

This report comprises of the summary of the data wrangling of WeRateDogs Twitter archive Data.

I worked with three datasets for this project, they include:

A file on hand provided by Udacity, which I downloaded manually, twitter_archive_enhanced.csv. It contains basic tweet data of 2356 tweets from November 15th 2015 to August 1st 2017.

A file hosted on the Udacity server that I downloaded programmatically, image_predictions.tsv. It contains a table full of image predictions (the top three only) alongside each tweet ID, image URL, and the image number that corresponded to the most confident prediction.

A file named tweet_json.txt containing data including but not limited to tweet IDs, retweet count and favorite count. This was gotten using Python's Tweepy library to access Twitter's API.

While assessing this datasets, I came across eleven quality issues and two tidiness issues, which I cleaned up and then I merged the dataset into one master dataframe.

## Quality issues

1.archive missing values in in_reply_to_status_id, in_reply_to_user_id, retweeted_status_id, retweeted_status_user_id, retweeted_status_timestamp

2.archive timestamp is object datatype instead of datetime

3.archive there are 181 retweeted tweets

4.json_data missing values across multiple columns

5.json_data user column is a duplicate of id and id_str

6.json_data id and id_str columns have the same values

7.json_data created_at and archive table timestamp have the same values but different titles
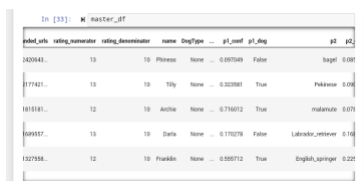
8.json_data has three columns with no values

9.json_data id should be tweet_id

**Tidiness issues**

1.archive table doggo, floofer, pupper and puppo columns should be one column

2.predictions table tweet_id column is arranged serially



*Master dataframe*

# Insights

1.Favorite count has a minimum of 0, a maximum of 144914 and a mean of 6704

2.Retweet count has a minimun of 1, a maximum of 70811 and a mean of 2325

3.The master dataset has a total colunm number of 22 which includes tweet_id, timestamp, source, text, expanded_urls, rating_numerator, rating_denominator, name, dogtype, jpg_url, img_num, p1, p1_conf, p1_dog, p2, p2_conf, p2_dog, p3, p3_conf, p3_dog, favorite_count, retweet_count.

4.P2 has the highest True count next to P1 and then P3 5.The most occurring rating numerator is 12 6.Using df.corr() to clarify the result of the scatterplot, it shows that retweet count and favorite count having a correlation of 0.815396 which is greater than zero confirms it has a positive correlation.

# Visualization



*Linear correlation between Retweet count and Favorite count*

This project was all together engaging and fun.