# Wrangle Report

# Analysis of Twitter user "WeRateDogs"

## Introduction

This report briefly describes the data wrangling efforts exerted in this project.

The dataset that you will be wrangling (and analyzing and visualizing) is the tweet archive of Twitter user @dog_rates, also known as WeRateDogs. WeRateDogs is a Twitter account that rates people's dogs with a humorous comment about the dog. These ratings almost always have a denominator of 10. The numerators, though? Almost always greater than 10. 11/10, 12/10, 13/10, etc. Why? Because "they're good dogs Brent." WeRateDogs has over 4 million followers and has received international media coverage.

The entirety of this project was implemented using Python language, Pandas, Matplotlib, Seaborn, Requests, Tweepy, Json, and NumPy.

**Data wrangling, which consists of:**

- **Gathering Data**
  The data used was gathered from three different sources:
  1. **Enhanced Twitter Archive dataset** was a file on hand. The file was stored in the CSV file.



  2. **Image Predictions File dataset** was retrieved from Udacity's servers through the Requests library. The file was stored in the TSV file.

3. **Data via the Twitter API dataset** was retrieved from Tweeter's API call, using the Tweepy library. Retrieved data was stored in JSON format.

```
{'created_at': 'Tue Aug 01 16:23:56 +0000 2017',
 'id': 892420643555336193,
 'id_str': '892420643555336193',
 'full_text': "This is Phineas. He's a mystical boy. Only ever appears in the hole of a donut. 13/10 https://t.co/MgUWQ76dJU",
 'truncated': False,
 'display_text_range': [0, 85],
 'entities': {'hashtags': [],
  'symbols': [],
  'user_mentions': [],
  'urls': [],
  'media': [{'id': 892420639486877696,
    'id_str': '892420639486877696',
    'indices': [86, 109],
    'media_url': 'http://pbs.twimg.com/media/DGKD1-bXoAAIAUK.jpg',
    'media_url_https': 'https://pbs.twimg.com/media/DGKD1-bXoAAIAUK.jpg',
    'url': 'https://t.co/MgUWQ76dJU',
    'display_url': 'pic.twitter.com/MgUWQ76dJU',
    'expanded_url': 'https://twitter.com/dog_rates/status/892420643555336193/photo/1',
    'type': 'photo',
    'sizes': {'large': {'w': 540, 'h': 528, 'resize': 'fit'},
     'thumb': {'w': 150, 'h': 150, 'resize': 'crop'},
     'small': {'w': 540, 'h': 528, 'resize': 'fit'},
     'medium': {'w': 540, 'h': 528, 'resize': 'fit'}}}]},
 'extended_entities': {'media': [{'id': 892420639486877696,
    'id_str': '892420639486877696',
    'indices': [86, 109],
    'media_url': 'http://pbs.twimg.com/media/DGKD1-bXoAAIAUK.jpg',
    'media_url_https': 'https://pbs.twimg.com/media/DGKD1-bXoAAIAUK.jpg',
    'url': 'https://t.co/MgUWQ76dJU',
    'display_url': 'pic.twitter.com/MgUWQ76dJU',
    'expanded_url': 'https://twitter.com/dog_rates/status/892420643555336193/photo/1',
    'type': 'photo',
    'sizes': {'large': {'w': 540, 'h': 528, 'resize': 'fit'},
     'thumb': {'w': 150, 'h': 150, 'resize': 'crop'},
     'small': {'w': 540, 'h': 528, 'resize': 'fit'},
```

- **Assessing Data**

  After gathering each of the above pieces of data, I assessed them visually and programmatically for quality and tidiness issues.

  1. **Tidiness Issues**
     - doggo, floofer, pepper, and puppo columns should be one column (merge columns).
     - twitter_archive, image_predictions, and tweet_data Dataframes should be part of one Dataframe (merge dataframes).
  2. **Quality Issues**
     - Delete retweets by filtering the NaN of retweeted_status_user_id.
     - Delete in_reply_to_status_id, retweeted_status_id, in_reply_to_user_id retweeted_status_user_id, retweeted_status_timestamp columns.
     - Correcting data type in tweet_id (from int into string).
     - Correcting data type in timestamp (from string into datetime).
     - Delete rows with jpg_url have missing.
     - Create 1 column for dog image prediction and 1 column for dog image prediction confidence.
     - Delete p1, p1_conf, p1_dog, p2, p2_conf, p2_dog, p3, p3_conf, p3_dog, img_num columns.
     - Convert underscore to space and convert lowercase to uppercase in prediction_dog.

- **Cleaning Data**

  The previous issues I cleaned as appropriate resulting in high quality and tidy master pandas DataFrame.