# Wrangle Act

This project set out to extract data from the Weratedogs twitter database, clean the data and derive useful insights from the data. In the project, three sets of data were used for the analysis, two out of the three were provided by the Udacity instructors, while I had to obtain the third from the twitter database using their API.

## Obtaining the first data:

The first data, labelled twitter_enhanced_data was pretty straightforward to obtain as it was provided directly. All I did was download the file and read it into a pandas dataframe.

## Obtaining the second data:

For the second data, I was provided the url, and I had to make use of the request library in python. After obtaining the file which was stored temporarily in my system memory, I wrote it into a file, after which I read it into a pandas dataframe. The syntax and the codes for this are represented in a snapshot below:



## Obtaining the third data:

For the third dataset, I used the API for twitter to obtain the data. In order to get access to the API, I applied for a twitter developer account, after which I was granted access in about two weeks.

I pip installed tweepy (a library for querying twitter API). The syntax for the codes I used for this are displayed below.

The codes took roughly 35mins to download the contents and write to a twitter_json.txt file.

**Creating a Dataframe from twitter_json.txt file:**

After writing the contents of each tweet_id to a txt file, I converted the file to a pandas dataframe. I made use of the json library method (json.dumps) in order to convert the json.txt file to a json object. I did this so that I could access the contents I needed in the form of a python dictionary.

The additional data I gathered from the json file include:

- 'tweet_id':[],
- 'date':[],
- 'retweet_count':[],
- 'favorite_count':[],
- 'full_text':[],
- 'followers_count':[]

**Data Cleaning**

After obtaining the data, I performed some data wrangling processes. 11 Quality issues and 2 Tidiness issues were identified and fixed using the **Define, Code and Test** format.

Notably, for the tweet_ids that had erroneous numerator and denominator ratings, I went to the werate_dogs twitter database to obtain the correct ratings for those posts. I did this so I could have a dataset that was near perfect for the exploratory data analysis.

After cleaning the datasets, I set to explore the data to generate insights about the data.