

# Udacity

## Data Analysis Nanodegree

### Wrangle Report

By: Mohamed Gabr

#### Introduction

This report is based on the data wrangling project created by Udacity for the Data Analysis Nanodegree program. The project's main task is to wrangle and analyze the tweets archive of the @dog\_rates Twitter account which goes under the name WeRateDogs.

Several datasets were gathered in the process and after assessing them I found several issues that needed some adjustments.

#### Quality issues:

##### **twitter-archive-enhanced.csv**

- rating\_denominator column has values other than 10
- Dog names in name column have (none, the, an, a, this, None) as values
- there are 181 different retweets which are not original dog ratings
- (retweeted\_status\_id, retweeted\_status\_user\_id, retweeted\_status\_timestamp) columns should be dropped after removing the retweets
- There are 78 replies to original dog rating tweets
- (in\_reply\_to\_status, in\_reply\_to\_user\_id) columns should be dropped after the replies are removed
- timestamp is a string not datetime
- source column is not informative

##### **tweet-json.txt**

- (friends\_count, source, retweeted\_status, url) columns should be dropped

#### Tidiness issues:

- timestamp consists of one column only (date + time) and not two separate columns
- twitter-archive-enhanced.csv and image-predictions.tsv are not merged together in one dataframe or csv file

After going through the issues mentioned above, I created a final dataset which included the modifications of the issues. This dataset was exported as 'twitter\_archive\_master.csv' and was a result of merging the three datasets together. Those datasets are: twitter-archive-enhanced.csv, tweet-json.txt, and image-predictions.tsv.