

# Data Wrangling Report prepared by Caleb Omariba

## Udacity Data wrangling of datasets from “WeRateDogs” tweets.

### INTRODUCTION

The WeRateDogs which is the source of our data is a Twitter account that rates people's dogs with humorous comments about the dog. The account was started in 2015 by a college student Matt Nelson. This account has attracted 8.8 million+ followers ever since.

### Data wrangling steps in this project:

- Step 1: Gathering data
- Step 2: Assessing data
- Step 3: Cleaning data

### Gathering data.

The project required me to gather data from three different sources.

- a) I downloaded one dataset given by Udacity i.e ([twitter-archive-enhanced.csv](#)) which I successfully loaded. Which i named (twits\_arc\_df)
- b) Programmatically downloaded (`image_predictions.tsv`) using the [Requests](#) library. Which i named (predictions\_df )
- c) I used [Tweepy](#), a Python library, to query Twitter's API for WeRateDogs Twitter data whereby i create a table called (api\_tweets\_df)

### Assessing data

In the assessment process, I found some issues and documented it in two parts, the **quality issues** and the **tidiness issues**

## Quality issues observed.

twits\_arc\_df

- **in\_reply\_to\_status\_id, in\_reply\_to\_user\_id, retweeted\_status\_id, retweeted\_status\_user\_id, retweeted\_status\_timestamp** variables have a lot of missing data and, moreover, we do not need them for the analysis.
- The **timestamp** column datatype is wrongly string data type.
- The standard **rating\_numerators** range from 0 to a maximum of 15 the rest of the values are considered as outliers.
- 10 is the standard **rating\_denominator** that is used in the WeRateDogs handle. The rest of the values are incorrect.
- We can observe that the **majority of Dog names start with a capital letter** so those that do not are wrong.

predictions\_df

- The **column names in the predictions\_df are not meaningful** to be understood clearly.
- we have **66 duplicates in the jpg\_url column**
- The **dog breeds' names have some names beginning with uppercase letters while others begin with lowercase letters** which is not consistent.

## Tidiness Issues.

twits\_arc\_df

- we have **4 columns for the dog stages which is wrong.**

api\_tweets\_df

- *The api\_tweets\_df table should be combined with the twits\_arc\_df table.*

## Cleaning data

1. Dropped all columns with a lot of missing data which will not be helpful in the analysis process. I used the drop() method to remove the specified redundant values.
2. Changed the timestamp to datetime type data type. I used the pd.to\_datetime() function to the strings to datetime data type.
3. Drop rating\_numerator that are >15. Since the majority values of the rating\_numerator were less than 15 the ones more than 15 were considered as outliers.
4. Dropped the rating\_denominator column since it will not be useful because it is a constant that is all the denominator values are to be 10 though some are not 10.

5. Converted the first letter of the dog names to capital for all the dogs names to make them consistent.

6. Renamed the column titles of the predictions\_df to meaningful titles.

7. Dropped the duplicate values from the jpg\_url column in the predictions\_df.

8. Converted all dog names to begin with uppercase letters in the predictions\_df.

9. When handling the tidiness issues I used the melt function Melting the 4 dog stages into one column.

10. Merged twits\_arc\_df and api\_tweets\_df tables together using the merge function and finally I joined the two with the predictions\_df to form the final master dataset table that I used for analysis and visualizations.