

## Wrangle Report

We were given a tweet archive dataset from Twitter user weratedogs (@dog\_rates). WeRateDogs is an account on Twitter that posts images and comments of different dogs.

The major goal of the project was to go through the entire data wrangling process from gathering to assessing to cleaning data after which the cleaned dataset is used to run an analysis.

The high-level process was involved:

- Data gathering
- Assessing of gathered data
- Cleaning of data based on issues identified
- Storing, analyzing and visualizing the data to get insights

### Data Gathering

This involved gathering data across three sources. Data gathered was in different formats (csv, via a URL and via twitter APIs)

- The first dataset represented the Twitter archive of WeRateDogs. This was provided by Udacity in a csv format. The file named `twitter_archive_enhanced.csv` was read using the `pandas read_csv` function.
- The second dataset represented the tweet image predictions with information on dog breed based on predictions by a neural network. This was also provided by Udacity and pulled from a URL and was read using the `pandas read_csv` function while specifying the separator as `'\t'`
- The third dataset was gathered from Twitter APIs using the Python Tweepy library. The purpose of this was to pull favorite and retweet count.

### Assessing of Gathered Data

In this section, the goal was to look through the various datasets to identify issues present in each of them. For this, a programmatic and visual assessment was run on each dataset. The requirements for this were split into quality and tidiness. For quality, I checked for:

- Completeness: Identify missing data
- Validity: Check that data present are valid and adequately represent their fields
- Accuracy: Besides validity, is the data correct?
- Consistency: Is it standardized across fields?

While for tidiness, I checked that:

- Each variable forms a column
- Each observation forms a row
- Each type of observational unit forms a table

The issues identified in this process are listed below:

## Quality Issues

### *twitter\_archive:*

- Missing data in columns: in\_reply\_to\_status\_id, in\_reply\_to\_user\_id, retweeted\_status\_id, retweeted\_status\_user\_id, retweeted\_status\_timestamp, expanded\_urls
- Dog names: some dogs have 'None' as a name, or 'a', or 'an.'
- Dataset has retweets which will lead to duplicated data resulting in some empty columns (Columns with the the retweeted tag)
- Datatype of tweet\_id is int (on all tables)
- Datatype of timestamp is an object
- Rating\_numerator has values up to 1776
- Rating\_denominator has values above 10
- The source column still has HTML tags

### *image\_predictions:*

- p1, p2 and p3 columns have invalid data
- p1, p2 and p3 columns aren't consistent: It varies between lowercase and sentence case
- Multi word breeds are represented with an underscore in columns p1, p2, p3

### *action\_counts:*

- Missing data in columns

## Tidiness Issues

### *twitter\_archive:*

- Dogoo, Fluffer, Pupper and Puppo all relate to a variable

### *image\_predictions:*

- Should form one observational unit with twitter\_archive

### *action\_counts:*

- Should form one observational unit with twitter\_archive

## Cleaning of data based on issues identified

In this section, I cleaned the data based on the issues I had identified. Steps taken are:

- Copy the independent data sets and merge together to form one table/dataset
- Correct issues in the name column of twitter\_archive
- Delete retweets
- Remove columns with missing data not needed for the analysis. These are: in\_reply\_to\_status\_id, in\_reply\_to\_user\_id, retweeted\_status\_id, retweeted\_status\_user\_id, and retweeted\_status\_timestamp
- Change datatype of tweet\_id from integer to string
- Change datatype of timestamp from object to datetime format
- Melt the various dog types into one column
- Remove all columns which are no longer needed
- Standardize dog ratings for consistency
- Create new breed column using data from the image prediction table

This involved coding and testing after which the data was stored and further analysis run.

