

WeRateDogs Dataset

# Data Wrangling Report

This report documents the wrangle efforts put into preparing the dataset for analysis. The data wrangling procedure includes the following steps:

1. Data Gathering
2. Data Assessment
3. Data Cleaning

## 1. Data Gathering

In this project, data was gathered from three sources in three different formats. The data gathered are as follows:

- a. WeRateDogs Twitter Archive: This file was provided in CSV format, which was downloaded manually from the Udacity classroom.
- b. Image Predictions: This file is a TSV file that is hosted on Udacity server. The file was downloaded programmatically using *requests* package and stored in my working directory from where the file was loaded into the notebook.

- c. Tweets: To retrieve tweets from Twitter API, the Tweepy package was used to get the status of each tweet by passing the *tweet\_id* as a parameter to *get\_status* method. The *json* package was used to read the response, which was in JSON format and the response was writing to a text file *tweet\_json.txt*.

Finally, the text file was read and the *retweet\_count* and *favourite\_count* were retrieve for each available tweet. Note some tweet were not available.

## 2. Data Assessment

In this step, I assessed the data visually and programmatically. I used Microsoft Office Excel application for the visual assessment while pandas DataFrame methods were used to assess the data programmatically, including *head()*, *info()*, *describe()*. The following data quality and tidiness issues were identified:

### Quality Issues

#### *twitter\_archive Table:*

- i. Only original rating with images are needed but the table includes retweets and replies ratings.
- ii. The following columns: *in\_reply\_to\_status\_id*, *in\_reply\_to\_user\_id*, *retweeted\_status\_id*, *retweeted\_status\_user\_id*, and *retweeted\_status\_timestamp* are not needed for the analysis.
- iii. *timestamp* column datatype is string instead of datetime.
- iv. The datatype for the *rating\_numerator* and the *rating\_denominator* is int instead of float.
- v. Ratings with decimal points are presented wrongly
- vi. Wrong rating for rows with multiple rating like text
- vii. The *source* column contains html markup

#### *image\_predictions Table:*

- viii. Only interested in the highest confidence level p1 that are dogs
- ix. Inconsistent captilization in the p1 => breed column

### **Tidiness Issues**

- x. The doggo, floofer, pupper, puppo columns in the twitter\_archive represent different dog stages and should be into a single column.
- xi. The three tables should be combined into one table

## **3. Data Cleaning**

In this step, I created a copy of the three data frames and all identify quality and tidiness issues were resolved using the Define-Code-Test framework. Describe briefly below is the action taken to resolve each issue.

### **Quality Issues**

#### **twitter\_archive Table:**

- i. Dropped retweets, replies and rows without an image.
- ii. Dropped unwanted columns.
- iii. Converted the datatype of timestamp column to datetime.
- iv. Converted the datatype of the rating\_numerator and the rating\_denominator to float.
- v. Corrected ratings for rows with fractional rating numerator.
- vi. Corrected rating for rows with multiple rating like text
- vii. Extracted only the content of the anchor tag in the source column.

#### **image\_predictions Table:**

- viii. Filtered for only rows which are dogs in the p1 predictions.
- ix. Updated capitalisation of the p1 column to first letter capital using the capitalize method of the pandas Series.

### **Tidiness Issues**

- x. Combined the doggo, floofer, pupper, puppo columns in one column using the melt method. Afterwards the variable column was dropped and the value column sorted and duplicate removed. (Note dogs with multiple dog stages, only one was selected as the dog stage).
- xi. The three tables were merged into one table using the DataFrame merge method.