

# Data Wrangling for WeRateDogs Datasets

Author: Xueyun Zhou

This report summarizes the processes of data wrangling for the three WeRateDogs Datasets.

## 1. Gathering data

Three pieces of data as follows are gathered:

- 1) the WeRateDogs Twitter archive, saved as `df_enhanced`;
- 2) the tweet image predictions, saved as `df_image`;
- 3) each tweet's retweet count and favorite ("like") count from the Twitter API, saved as `df_json`.

## 2. Assessing data

### 2.1 Quality issues

#### 2.1.1 `df_enhanced`

- (1) Some "rating\_denominator"s are not 10.
- (2) The values for "doggo", "floofer", "pupper", and "puppo" are not the Boolean type but showed as the column names themselves.
- (3) There are observations whose "retweeted\_status\_id" are not NaN, so they are retweets.
- (4) The data type of "tweet\_id", "rating\_numerator", and "rating\_denominator" is int64 and that of "timestamp" is str.
- (5) Some values of "rating\_numerator" are unusual (420, 666, 960, 1776, etc.).
- (6) There are zeros in "rating\_denominator".
- (7) There are too much missing data for the "doggo", "floofer", "pupper", and "puppo" columns.

### 2.1.2 df\_image

- (8) The values of "img\_num" for some observations do not represent the images with the highest confidence.
- (9) The data type of "tweet\_id" is int64.
- (10) The "img\_num" of some observations is 4, while there is no data for the 4th image.
- (11) There are 2356 observations in df\_enhanced, but 2075 in df\_image, so there are tweets in df\_enhanced which do not have images.

### 2.1.3 df\_json

- (12) The data type of "id" is int64.

## 2.2 Tidiness issues

- (13) The column name of the tweet ID in df\_json is different from that in df\_enhanced and df\_image.
- (14) There are too many unrelated columns collected in df\_json.
- (15) The "source" column is in a html formatting.
- (16) The columns "in\_reply\_to\_status\_id" and "in\_reply\_to\_user\_id" in df\_enhanced are useless. The columns "retweeted\_status\_id", "retweeted\_status\_user\_id", and "retweeted\_status\_timestamp" are useless after the retweets are deleted.
- (17) The three datasets should be merged.

## 3. Cleaning data

### 3.1 Issue (13), (14)

Remove all columns for df\_json\_copy except "id", "favorite\_count", and "retweet\_count".  
Change the column name of "id" to "tweet\_id".

### **3.2 Issue (3)**

Delete retweets.

### **3.3 Issue (4), (9), (12)**

Change the data type of "tweet\_id" in three datasets to string.

Change the data type of "timestamp" to pandas datetime.

Change the data type of "rating\_numerator" and "rating\_denominator" to float.

### **3.4 Issue (15)**

Delete the html contents in "source".

### **3.5 Issue (16)**

Delete the columns "in\_reply\_to\_status\_id", "in\_reply\_to\_user\_id", "retweeted\_status\_id", "retweeted\_status\_user\_id", and "retweeted\_status\_timestamp".

### **3.6 Issue (2), (7)**

Gather information of the stage from "text".

Reassign the stage to "doggo", "floofer", "pupper", "puppo" as Boolean type.

### **3.7 Issue (1), (5), (6)**

Check the unusual rating\_denominator and rating\_numerator, and update based on "text".

### **3.8 Issue (8), (10)**

Create a column in df\_image\_copy for the dog type based on confidence columns and "is dog" columns.

Delete useless columns.

### **3.9 Issue (11), (17)**

Merge 3 datasets. The `df_enhanced_copy` and `df_image_copy` should be merged by inner to remove tweets without images.

After cleaning and merging, the clean dataset is ready for storing, analysing, and visualizing.