

# WeRateDogs Twitter Data Wrangling Report

Three different sets of data all gotten from weratedogs account were used for this analysis. The first was the enhanced tweet archive which contained 2356 samples of data. The second is image prediction dataset that had 2075 samples. In this report, I will refer to them as archive dataset, tweet dataset and prediction dataset.

After careful assessment using both visual and programmatic approach, data quality and tidiness issues are discovered within the datasets. Under quality issues, we have:

## Quality issues

- Missing records (2327 instead of 2356) in tweet dataset
- Non original tweets and Unwanted columns
- Wrong column data type
- 9/11 event mistaken as rating (archive dataset)
- Non 10 rating denominator value (archive)
- Incorrect value for dog stages (archive)
- Underscore used instead of space to separate words (p1, p2 and p3 columns) (prediction)
- No uniform character case in p1, p2 and p3 values (prediction)
- Erroneous datatypes (p1\_dog, p2\_dog, p3\_dog) (prediction)

## Tidiness issues

- Same column bearing different title in another table (timestamp/created\_at, id/tweet\_id, text/full\_text)
- Tweet text duplicated in tweet and archive table
- Dog stages in four columns instead of one column in archive table

I started off by trying to resolve all the missing data issues. Out of 2356 tweet ids residing in archive dataset, I was able to fetch 2327 tweets using tweepy. Next was the tidiness issue. Using df.rename pandas method, I changed the column names timestamp to created, tweet\_id to id, text to full\_text all in archive dataset. I also changed column name tweet\_id to id in prediction dataset. Next, I dropped the text (now full\_text) in the tweet dataset. This is to avoid repetition of columns in multiple datasets (tweet and archive).

For the dog stages issue, I used pandas melt function to merge the four columns (flooper, pupper, doggo and puppo) into two columns titled stage and staged\_value. Dropped stage and renamed stage\_value to dog\_stage which contained string values of stages.

One of the rules of data tidiness is to have one variable in one column. Since entities and extended entities column in tweet dataset has multiples and weren't needed in my future analysis, I dropped them.

## Resolving Quality issues

**Tweet dataset:** The first issue I tackled was converting created\_at column dtype to datetime using pandas to\_datetime() method. Then I used df.drop function to drop columns that had entire value set to null (geo, coordinates, contributors). I dropped id\_str column since it was

the same as id only of type object. Then, I removed all retweets, quotes and reply type samples by dropping samples with non null in\_reply\_to\_status\_id and quoted\_status\_id. Lastly I selected these columns; created\_at, id, retweet\_count, favorite\_count and full\_text and dropped the rest.

**Archive dataset:** Dropped retweets, quotes and replies by deleting data samples that had non null in\_reply\_to\_status\_id and retweeted\_status\_user\_id column values. in\_reply\_to\_status\_id, in\_reply\_to\_user\_id, retweeted\_status\_id, retweeted\_status\_user\_id and retweeted\_status\_timestamp columns was also dropped since it wont be contributing to analysis. Next created\_at column (formerly timestamp) was converted to dtype datetime. Using regex, I extracted rating numerator and denominator from full\_text thereby resolving the 9/11 issue. I then dropped all non 10 rating denominator since they were basically outliers. Lastly, I applied regex on the full\_text column and extracted any occurrence of either of the four known dog stages and filled in the dog\_stage column where found or the string *unknown* otherwise.

**Prediction dataset:** I used df.column\_name.str.replace method to replace all underscores in p1, p2 and p3 columns with space and changed the character case to lower. lastly , I converted the p1\_dog, p2\_dog and p3\_dog column to dtype categorical.