

# Wrangling data for Wrangling and Analyze Data project

By Jennifer Powell

## Gather

The first step to wrangle the data for this project was to gather the 3 files required for this project:

1. WeRateDogs Twitter archive (twitter-archive-enhanced.csv), which I downloaded from udacity, and read into a dataframe.
2. The tweet image predictions file (image\_predictions.tsv), which was downloaded programmatically and read into a dataframe
3. The retweet and favorite count data, which was queried from Twitter and acquired programmatically using the Twitter API. The tweet\_id from the archive file was used to obtain JSON data from Python's Tweepy library. The resulting data was stored one line at a time in a text file: tweet\_json.txt. This was then read into a 3<sup>rd</sup> dataframe.
4. I decided to obtain a 4<sup>th</sup> file ( language-codes\_csv.csv) to use for decoding the 2 char code in the language column of the image\_predictions table to a readable text. I downloaded this file programmatically from <https://datahub.io/core/language-codes> and read it into a 4<sup>th</sup> dataframe.

## Assess

*Detect and document at least **eight (8) quality issues** and **two (2) tidiness issues***

The next step was to assess the data in the 4 dataframes. I went through each dataframe and looked at the quality of the data (missing data, garbage data, duplicates). These were the quality issues found in each of the datasets:

### Quality

#### twitter\_df table

1. erroneous datatypes - id, retweet\_count and favorite\_count should be int; timestamp should be datetime
2. language abbrev should be full name, not cryptic
3. many columns have only 1 value or no data (contributors, coordinates, favorited, geo, place, possibly\_sensitive, possibly\_sensitive\_appealable)

#### img\_pred\_df table

1. some of the p1 values look to be garbage (eg crash helmet, water bottle)
2. many records are not dogs (p1\_dog=False)
3. some records are retweets

#### wrd\_df table

1. Some names look to be garbage (eg a, an) clean up and set to 'None'?
2. some records contain retweets
3. Erroneous datatypes (source column, stage should be category, timestamp should be datetime)
4. a handful of denominators aren't 10
5. missing images
6. source column contains extra info, only 4 values
7. stages of dog: doggo, pupper, puppo, and floof(er) some have multiple ratings for same dog
8. missing floof stage, only checked for floofer

## Tidiness

Column headers are values, not variable names

- wrd\_df table -doggo, floofer, pupper, puppo columns contain only 4 values in 4 columns

A single observational unit is stored in multiple tables.

- The source and text of the tweet are stored in both wrd\_df and twitter\_df. We only want the twitter\_df favorite\_count and retweet\_count so can combine with wrd\_df and drop rest of columns in twitter\_clean
- jpg\_url is repeated in img\_pred\_df. Only need p1 and p1\_conf, can combine these with wrd\_df and drop rest of columns

## Clean

1. Language abbreviations in twitter\_df were replaced using a lookup table.
2. Determined records with missing images and dropped them from wrd\_df.
3. Removed records from wrd\_df that are retweets.
4. One character names in wrd\_df set to 'None'.
5. Numerators can be greater than 10, but decided to make denominators all 10 since so few.
6. Created subsets of **twitter\_df**: 'id','retweet\_count','favorite\_count','lang' and **wrd\_df**: 'tweet\_id', 'source', 'timestamp','expanded\_urls', 'text', 'rating\_numerator', 'rating\_denominator','name','doggo', 'floofer', 'pupper', 'puppo' dropping unused columns.
7. Combined the doggo, floofer, pupper, puppo columns into 1 categorical **stage** column.
8. Convert wrd\_df.source to category type with 4 values (iPhone,Vine,WebClient, TweetDeck).
9. Deleted rows which weren't dogs (ie p1\_dog == False).
10. Converted twitter\_df.retweet\_count and favorite\_count to integer.
11. Converted wrd\_df.timestamp to datetime.
12. Combined the dataframes into one master.