

# Reporting: wragle\_report

This project uses three datasets: Twitter archive dataset, Additional data from twitter api and image predictions dataset. I used both visual assessment and programmatic assessment to know what issues to clean in all the datasets.

## Twitter archive dataset

The Twitter archive dataset has a total of 17 columns:

tweet\_id, in\_reply\_to\_status\_id, in\_reply\_to\_user\_id, timestamp, source, text, retweeted\_status\_id, retweeted\_status\_user\_id, retweeted\_status\_timestamp, expanded\_urls, rating\_numerator, rating\_denominator, name, doggo, floofer, pupper, puppo.

The quality issues I picked up after assessing were:

1. Mismatched datatypes in timestamp should be datetime not object.
2. Finding incorrect ratings visually and updating them manually.
3. Some columns have null values which are entered as none instead of NaN.
4. In dog names column, some dog names are recorded as None, such, the, quite, an.
5. Some expanded\_urls are duplicated.
6. Removing rows that have non-empty: retweeted\_status\_id, retweeted\_status\_user\_id, and retweeted\_status\_timestamp.
7. Some columns in twitter\_archives\_enhanced\_csv are not relevant to my analysis need to be dropped (e.g retweeted\_status\_id, retweeted\_status\_timestamp e.t.c)

The tidiness issues I picked up after assessing were:

1. In twitter\_archives\_enhanced\_csv Dog stages are spread into different columns (puppo, pupper, doggo, floofer).

## Solving procedures

For quality issues: The dataset had no missing data but had duplicates in expanded\_urls column which I dropped. The timestamp datatype was object instead of datetime, hence the problem was solved by converting the datatype to datetime. In the dataset some ratings were not correct compared to the ratings given in text column. This was solved by visually assessing the incorrect ratings and updating them manually. Some columns are recorded as None where there is a missing value instead of NaN, the problem was solved by replacing the None values with NaN. In the name column some dog names are invalid such as, quite, very, None, e.t.c, the invalid names were replaced by NaN. Irrelevant columns such as in\_reply\_to\_status\_id, in\_reply\_to\_user\_id are dropped. Rows that had non-empty retweeted\_status\_id, retweeted\_status\_timestamp, were dropped since they are retweets and we don't want them to interfere with the findings. For tidiness issues: Some columns; doggo, puppo, floofer and pupper are spread instead of being in one column, hence the columns are melted. Later the dataset was merged to make a master dataset.

### **Image predictions dataset.**

The image prediction dataset has 12 columns: tweet\_id, jpg\_url, img\_num, p1, p1\_conf, p1\_dog, p2, p2\_conf, p2\_dog, p3, p3\_conf, p3\_dog.

The quality issues picked after assessing:

1. Some jpg\_urls are duplicated

The tidiness issues were:

2. Melting the image\_predictions dataframe

### **Solving procedures**

The jpg\_url column had duplicates which were dropped. In tidiness issue, columns p1, p1\_conf, p1\_dog, p2, p2\_conf, p2\_dog, p3, p3\_conf and p3\_dog were melted and the dataframe was later merged with twitter archive. Irrelevant columns were also dropped.

### **Additional twitter data**

This dataset has four columns: tweet\_id, create\_date, retweet\_count, favorite\_count.

The quality issues picked up were:

1. Create\_date is datatype object instead of datetime

The create\_date datatype was object instead of datetime, this was solved by converting the datatype to datetime. Later the dataset was merged with the twitter archive dataset to make a master dataset.

### **Conclusion**

In the wrangling process all the three datasets were cleaned and merged to form a master dataset that contains the relevant columns needed for analysis. Wrangling helped me identify and correct mistakes which may have led to inaccurate findings.