

## Data Wrangling Report

### Investigate twitter\_archive\_enhanced.csv

The dataframe "tae" is created from twitter\_archive\_enhanced.csv which is merged data from twitter\_archive\_enhanced.csv, tweet\_json.txt .

tweet\_json.txt , twitter\_archive\_enhanced.csv , image\_prediction.tsv are first visually analyzed by opening as csv file on google sheets.

Dataframe tae is investigated programmatically using below mentioned panda methods.

Tae.info() displays

- in\_reply\_to\_status\_id and in\_reply\_to\_user\_id are floats.
- Timestamp is an object.
- There are 181 values in retweet\_status\_id and retweet\_status\_user\_id columns.
- The columns retweeted\_status\_id and retweeted\_status\_user\_id exist which are not needed for analysis.

Tae.describe() displays numeric columns. tae.tail(),tae.sample().tae.info() used to take the first look at all the columns info and sample values.

Min max rating denominator and min max rating numerator are evaluated

To help understand range on ratings. Ratings are needed to be calculated in further wrangling efforts. The value counts in ratings numerator and denominator are analyzed for finding the range. The Zero value in numerator and denominator was not detected in min max analysis. The value\_counts showed there are few records with 0 value. These rows will need attention.

Below quality issues are detected in twitter\_archive\_enhanced.csv

**Issue 1. Retweet id is present for many rows**

**Issue 2. Ratings are accessed by programmatic addition of column ratings and fixing the rows with 0 denominator values**

**Issue 3. 181 are retweet**

**Issue 4. in\_reply\_to\_status\_id is float and has nulls**

**Issue 5. in\_reply\_to\_user\_id is float and has nulls**

**Issue 6. retweeted\_status\_id,retweeted\_status\_user\_id are not needed for analysis.**

**Issue 7. Names with value 'None' or 'a'**

**Issue 8. tweet\_id is defined as object and int64**

**Issue 9. timestamp is object**

**Issue 10 expanded\_url is null**

### **Investigate image\_prediction**

Describe(),info(),value\_counts(),sample(),head(),tail() are used for general programmatic assessment of columns,data type and values.

1. P1\_conf > p2\_conf > p3\_conf shows that p1 is most predictable value

### **Investigate jasondf created from tweet\_json.txt**

Describe(),info(),value\_counts(),sample(),head(),tail() are used for general programmatic assessment of columns,data type and values.

1. tweet\_id is defined as object and int64 in twitter\_archive\_enhanced.csv & image prediction

2. retweet\_count & favorite\_count are objects in jason df.

Below tidiness issues are noted from investigating above 3 dataframes:

1. Columns doggo ,floofer, pupper, puppo need to be under one column heading Doggo Lingo

2. Column retweet\_count , favorite\_count of jtweet\_jason\_df need to be with columns of twitter\_archive\_enhanced.csv as it makes more relevance along with other variables

3. Image and highest prediction need to be with columns of twitter\_archive\_enhanced.csv as it makes more relevance along with other variables

Cleaning Efforts:

- Drop rows with retweeted\_status\_id <> 0
- Remove null from in\_reply\_to\_status\_id and convert to int
- Remove null from in\_reply\_to\_user\_id and convert to int
- Convert timestamp from object to datetime
- Change name = 'none' and 'a' to NoNameDogs
- Drop row with rating denominator min = 0 and numerator min = 0. This decision is made to avoid math error issue for calculating ratings\_round = numerator/denominator
- Calculate and Round the ratings: A New column is added ratings\_round to store the calculated ratings.
- Replace nulls and combine the columns doggo,floofer,pupper,puppo

- Fix the data type of jasondf columns and merge columns retweet\_count , favorite\_count of jasondf with twitter\_archive\_enhanced on tweet\_id as it makes more relevance along with other variables.
- join image\_prediction with tae\_merge for image\_url and p1