

# Wrangle Act Report

August 23,2018

## Data Gathering:

Data is gathered from three sources 'twitter\_archive\_enhanced.csv', 'image\_predictions.tsv' and 'tweet\_json.txt'.

First data was collected from the "twitter-archive-enhanced.csv" file which was in the same directory in which project notebook was located. The csv file was imported into pandas dataframe. The dataframe was named "df"

Second data was extracted programmatically from a URL:

[https://d17h27t6h515a5.cloudfront.net/topher/2017/August/599fd2ad\\_imagepredictions/image-predictions.tsv](https://d17h27t6h515a5.cloudfront.net/topher/2017/August/599fd2ad_imagepredictions/image-predictions.tsv). Python's request library was used to extract data from this URL. This file was imported as a dataframe in pandas using tab as the separator. The dataframe was named "img\_pred".

The third data was extracted from Twitter API using python's tweepy library. I needed to extract id, the favourites and retweet counts for each tweet. This data was then saved as a JSON file using UTF-8 encoding.

## Data Assessing:

Data was assessed visually as well as programmatically. Through pandas info(), head(), describe() methods I was able to detect some quality and tidiness issues like many tweets were retweeted, dataframe contained the faulty names, there were several empty values in in\_reply\_to\_status, in\_reply\_to\_user\_id, retweeted\_status\_id, retweeted\_status\_user\_id, retweeted\_status\_timestamp.

In several columns, the null values are treated as non-null values. Some entries contain "Nan" as string.

Some of the ratings did not look right because the expected value for numerator and denominator was around 10, but there were many values above 100 also.

Dog "stage" variable were in four columns: "doggo, floofer, pupper, puppo" individually. That could be condensed in one column.

## Data Cleaning:

The copy of all the data frame were created as "df\_clean", "img\_pred\_clean", tweet\_info\_clean. Retweets were removed and tweets which did not include images were also removed because those tweets were not dog ratings.

dog\_stage column was created which showed the type of dog(dog stages).

Duplicate jpg\_url in "img\_pred\_clean" dataframe was dropped . The text in 'img\_pred' was made consistent and pretty.

Incorrect dog names in "df\_clean" were changed. tweet without rating was removed. The value for rating numerator and denominator was corrected. The column "source" was made in more readable categories.

Datatypes of timestamp,dog\_stage and tweet\_id, in\_reply\_to\_status\_id,in\_reply\_to\_user\_id was changed to datetime, categorical, and to strings respectively.

After doing all the process of cleaning all the dataframe was merged in single data frame named "df\_merge" and was stored in a csv file named "df\_master.csv" and then dataset was visualized in different forms