

# DATA WRANGLING REPORT

## Gathering the data

- Udacity provided the csv file 'twitter\_master\_enhanced.csv' which was loaded into a dataframe called 'archive'.
- Using the Requests library, 'image\_predictions.tsv' was downloaded and read into a dataframe called 'df\_image'
- The json text file provided was used in which the retweets counts and favorite count were gathered and called into a dataframe called 'df\_json'

## Assessing the data

### Quality issues

- The in\_reply\_to\_status, in\_reply\_to\_user\_id, retweeted\_status\_id, retweeted\_status\_user\_id, retweeted\_status\_timestamp columns have empty values.
- Drop the unnecessary columns.
- The datatype of the 'tweet\_id' is in integer, and should be converted to string (for all the 3 dataset)
- The source text is embedded in the source url.
- The timestamp is in object datatype, and should be converted to datetime.
- The rating numerator and denominator column have some wrong values.
- Some of the names 'none', 'a', 'an', 'this' etc. are not dog names.
- The dog breeds columns (p1, p2 and p3) have inconsistent case sensitivities.

### Tidiness issues

- The dogger, floofer, pupper and puppo columns should all be under one column as it represents only one variable (dog stage).
- Merge the archive data, df\_image and df\_json into one master data dataset.

## **Cleaning the data**

- Made copies of the original dataset.
- Removed rows with 'retweeted\_status\_id'.
- Dropped the unnecessary columns.
- From the archive data, the datatype of the tweet\_id was converted to string.
- The text is embedded in the source url was extracted.
- The timestamp was converted to datetime.
- The ratings were compared with the texts and wrong values were replaced. Also, ratings with decimals that were incorrectly extracted were fixed.
- Some of the unusual dog names were dropped and 'None' changed to unknown.
- The dog breed columns were standardized by changing the values to lowercase.

## **Storing the data**

- The clean data was saved into a csv file 'twitter\_archive\_master.csv'.