

## Report Data

the wrangling of datasets from “WeRateDogs” tweets.

Data wrangling procedure has been completed on the datasets obtained from “WeRateDogs”. There have been 3 main stages:

1. Data gathering;
2. Data assessment;
3. Data cleaning.

Data has been gathered from different resources such as provided files, a downloaded files via the internet using the provided link and via tweepy API. As a result 3 datasets (master\_df, images, api\_dfi) have been collected for further assessment and cleaning.

Data assessment was conducted against Quality and Tidness issues. Standard python methods and functions were used for data assessment such as .head(), .value\_countes(), .sample(), .describe(), .info() and etc.

The following problems were identified:

### **Quality issues:** of archive\_df

table column source in master\_df was too long for such source information in\_reply\_to\_status\_id, in\_reply\_to\_user\_id, retweeted\_status\_id, retweeted\_status\_user\_id, retweeted\_status\_timestamp variables have a lot of missing data and, moreover, we do not need them for the analysis rating\_denominator is 10 in 2333 cases out 2356 cases, consider rating\_denominator to be 10 rating\_numerator in most cases is in between 0 and 15, the rest consider as outliers variable name has 745 None values and 55 "a" values in timestamp +0000 is redundant information name variable has some entries starting with low case letters(example: An). Is An a dog name? It occurs 7 time in the dataset, though. expanded\_urls contains link which is not valid, possibly because they are expired remove rows in retweet\_count and favorite\_count with missing values

**prediction table** variable img\_num is not needed to change variable types to the appropriate, where it is needed prediction table is missing one important variable which would show if the picture truly contains the breed of dog or not like it is shown above. one picture has Straus on it and algorithm classified it as it is not a breed of dog, but another picture has a dog on it; however, it was misclassified as not a breed of dog

**Tidiness issues** master\_df table Variables doggo, floofer, pupper and puppo in one column tables wt\_arc and tweets\_api form one observational unit prediction table  
jpg\_url variable should be in tw\_arc table to satisfy tidiness definition  
tw\_arc and images tables form two different observations units and will be kept separately

One important issue came up during the assessment of the prediction dataset. This is a lack of the boolean variable which confirms that the pictures either have or do not have a breed of dog.

As it was shown manually that a picture of Straus was correctly classified as False meaning that it is not a breed of dog (p1\_dog=p2\_dog=p3\_dog=False), but there is no indication that the picture is not a breed of a dog. On the other hand, we have a breed of dog, which was two times misclassified p1\_dog=p3\_dog=False and one time correctly classified as a breed of dog p2\_dog=True. Again, without manual confirmation that this is a breed of dog or not, we cannot assess the correctness of the algorithm. We also discovered that the source variable takes only for 4 values which can be changed to more readable values such as 'Twitter for iPhone', 'Vine', 'Twitter Web Client', and 'TweetDeck'. Also, one can note that the expended\_urls variable has links that are broken; hence, this variable was dropped. The timestamp was cleaned by removing "+0000". In addition, to satisfy the tidiness definition tw\_arc and tweets\_api datasets were merged together. To satisfy the tidiness definition jpg\_url variable was moved from the prediction table to the tw\_arc table and after that prediction, the table is ignored due to the fact which was described above. During the cleaning stage for each step cleaning procedure was documented as "Define", code was developed and tested. Because data wrangling is an iterative process, some outliers were found and cleaned during the Analyzing and Visualizing Data stage, when histograms and scatter plots were drawn. Finally, the dataset was stored as .csv file: twitter\_archive\_master.csv.