

The project was a to take a twitter data from a twitter account by the name @WeRateDogs and do the data wrangling steps to gather, assess, clean the data to prepare it to store it and visualize it. the steps taken to prepare this data is:

## quality assures

- renaming columns (p1, p1\_conf, p1\_dog, p2, p2\_conf, p2\_dog, p3, p3\_conf, p3\_dog) to (algo\_A\_prediction, algo\_A\_confident, algo\_A\_dog, algo\_B\_prediction, algo\_B\_confident, algo\_B\_dog, algo\_C\_prediction, algo\_C\_confident, algo\_C\_dog)
- Delete columns from df1(in\_reply\_to\_status\_id, in\_reply\_to\_user\_id, source, text, retweeted\_status\_id, retweeted\_status\_user\_id, retweeted\_status\_timestamp, expanded\_urls), from df2(jpg\_url, img\_num), from df3('contributors', 'coordinates', 'created\_at', 'display\_text\_range', 'entities', 'extended\_entities', 'favorited', 'full\_text', 'geo', 'id\_str', 'in\_reply\_to\_screen\_name', 'in\_reply\_to\_status\_id', 'in\_reply\_to\_status\_id\_str', 'in\_reply\_to\_user\_id', 'in\_reply\_to\_user\_id\_str', 'is\_quote\_status', 'lang', 'place', 'possibly\_sensitive', 'possibly\_sensitive\_appealable', 'quoted\_status', 'quoted\_status\_id', 'quoted\_status\_id\_str', 'retweeted', 'retweeted\_status', 'source', 'truncated', 'user') by using the drop() method in the pandas library
- dropping all the rows with rating\_denominator > 10 or < 10 by making a query that gets all the rows with denominator bigger than 10 or smaller than 10
- change the id type from int to string in all the dfs by using the astype() method with parameter 'str' for string
- change the timestamp to date by using the pandas method to\_datetime().
- divide all the rating\_numerator by 10
- change the wrong predicted names of dogs (a , an , the) to null by using the replace method
- change the name of column id in the df3 to tweet\_id to match the name with df1, df2.
- I must chose 1 algorithm since having 3 is inassure, depending on the mean of the confidence of the algorithm
- drop rows that aren't dogs according to our chose algorithm

## tidyness

- type of dogs should be in 1 column.
- merge all dfs to one df

After doing this the data was clean and I was able to store the data in my machine for future use or sharing it with other