

## Report: act\_report

My initial or preliminary work was to gather data from different sources which includes flat file containing WeRateDogs archived tweets.  
Download file programmatically from Udacity server and using twitter API to gather tweets from twitter real time.  
Other resources from Udacity and from internet for supports.

TIDINESS ASSESSMENT: For predictions\_df: columns [ prediction 1, prediction 2, and prediction 3] are untidy twitter\_archive\_df: stage columns [doggo, floofer, pupper, and puppo] are untidy IN predictions\_df table jpg\_url variable should be in twitter\_archive\_df table to satisfy tidiness definition twitter\_archive\_df and predictions\_df tables form two different obseravations units and will be kept seperately IN twitter\_archive\_df table: doggo, floofer, pupper and puppo will be merged into one column tables twitter\_archive\_df and tweet\_in\_json\_df will be under one observational unit There should have been a dataframe for each of these observational units [tweet data, dog data, and image predictions]

QUALITY ASSESSMENT: IN predictions\_df: the prediction number needs to have data type integer One row has total confidence greater than 1 which is not normal The values in p1, p2, and p3 should be properly named Change 'None' to empty cell in doggo,floofer, pupper, puppo then delete There are duplicated jpg\_url IN twitter\_archive\_df: some tweets in archive have missing data for retweet\_count or favorite\_count alot of 0 in the favorite\_count IN twitter\_archive\_df: the dog\_stage columns (doggo, floofer, pupper, and puppo) should have a categorical data type using data type float for rating so that the column can be uniform All dogs should have a name with proper letter casing dropping timestamp column not needed

Cleaning was performed to take care of the following:  
removing retweets and duplicate values in jpg\_url column.  
removing duplicate values in jpg\_url column  
Removing pictures that are not of dogs when at least one value of p1\_dog, p2\_dog, and p3\_dog is true with confidence greater than 1.  
Remove any invalid names and text related to lowercase names that do not have the phrase "We only rate dogs"  
timestamp column was deleted

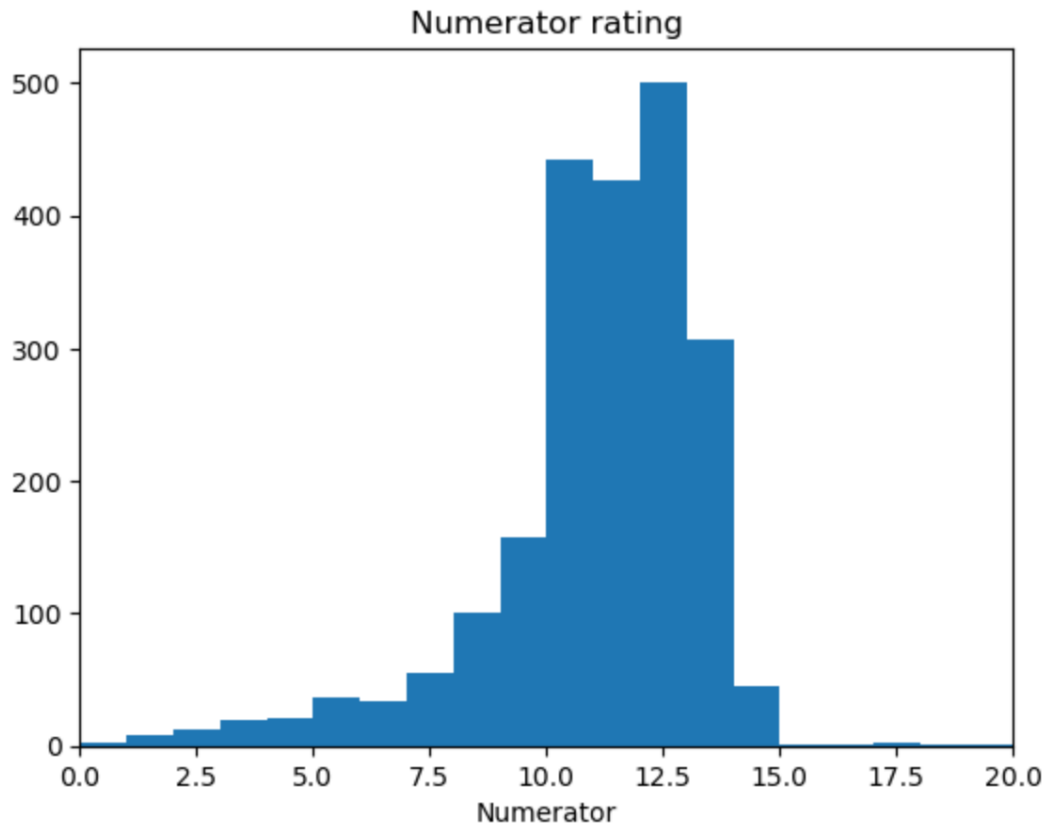
IN predictions\_df: the prediction number was changed integer  
One row has total confidence greater than 1 which is not normal  
The values in p1, p2, and p3 was properly named  
Change 'None' to empty cell in doggo,floofer, pupper, puppo then delete

IN twitter\_archive\_df: some tweets in archive have missing data for retweet\_count or favorite\_count  
alot of 0 in the favorite\_count so retweet\_count and favorite\_count was removed from the datafame.

rating number column was converted to float so that the column can be uniform  
All dog names were converted to proper letter casing  
dropping timestamp column because it was not needed

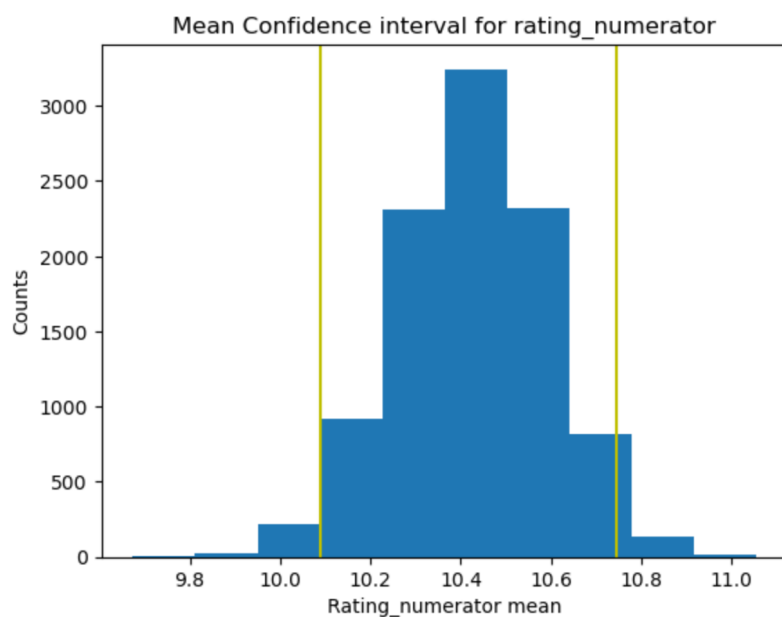
It was obvious a lot of names are in between 5 and 15 in their Numerator rating as shown below Because of the able we can see that the rating numerator histogram is skewed to the left.

```
In [73]: 1 bins = np.arange(0, twitter_archive_df_clean['rating_numerator'].max()+1, 1)
2 plt.hist(data=twitter_archive_df_clean, x='rating_numerator', bins=bins)
3 plt.title('Numerator rating')
4 plt.xlabel('Numerator')
5 plt.xlim(0, 20);
```



By defining confidence interval mean of 95% accuracy. I can confidently say the mean value of the rating\_numerator is inbetween 10.09, 10.75.

```
In [77]: 1 fig, ax = plt.subplots()
2 left, right = np.percentile(sample_numerator_rating_means, 2.5), np.percentile(sample_numerator_rating_means, 97.5),
3 plt.hist(sample_numerator_rating_means);
4 plt.title('Mean Confidence interval for rating_numerator')
5 ax.set_xlabel('Rating_numerator mean');
6 ax.set_ylabel('Counts');
7 plt.axvline(left, color='y');
8 plt.axvline(right, color='y');
```



And finally, the clean data set was saved as twitter\_archive\_clean\_master.csv

