

## WeRateDogs data wrangling

I divided data wrangling process to 4 stages: Gathering the data, assessing the data, cleaning the data and analyzing the data.

At first stage I started to download all the necessary files. After downloading the "twitter\_archive\_enhanced.csv" csv file manually, I created a folder named "Project" and downloaded "image\_predictions\_tsv" file programmatically using requests library to this folder. Then through the tweepy library I downloaded Twitters' JSON data and created "twitter\_data" file.

Second stage was assessing the data, I have identified several quality and tidiness issues.

### ***Quality issues are:***

- "twitter\_archive" and "image\_prediction" - twitter\_id column should be string
- "twitter\_archive" - timestamp column's data type should be datetime
- "twitter\_archive" - in dog stage columns null values is not correctly displayed
- "twitter\_archive" - Delete retweets
- "twitter\_archive" - some columns miss lot of values (in\_reply\_to\_status\_id, in\_reply\_to\_user\_id, retweeted\_status\_id, retweeted\_status\_timestamp)
- "twitter\_archive" - keep original tweets with images
- "twitter\_archive" - there are inappropriate names in dog name column
- "twitter\_archive" - Unnecessary columns
- "twitter\_archive" text column includes links
- "twitter\_archive" The standard denominator for "rating\_denominator" is 10
- "image\_prediction" - Capitalize the dog names in p1, p2, p3 columns

### ***Tidiness issues are:***

- "twitter\_archive" - dog stages should be under one column
- "Image\_prediction" and "twitter\_data" tables should be merged to "twitter\_archive" table

At the third stage of the wrangling process, I cleaned the data through solving the quality and tidiness issues. Finally, I saved the new data frame as "twitter\_archive\_master.csv" file.

Finally, I started to analysis the last version of the data. First, I looked how the number of retweets changed from the beginning of the account. Then relation between the favorite and retweet counts are displayed with scatter plot. Third insight was to identify the most popular dog breeds in dataset. Next step was to find out the dog breeds which are the most favorable. These results are showed with bar chart. The result of Top five dog names is also visualized with bar chart. At the last stage of analysis most rated dog are identified.

