

## Sources

I gathered data from a given CSV, a website, and Twitter's API. I used Tweepy to access the API and gather the JSON data for the tweets. I read the JSON data in a text file, then loaded what I needed into a pandas dataframe.

## Process

I audited the data by checking datatypes, value counts, number of non-null entries, and numeric summaries. Data for some of a few tweets. I combined (that is, inner joined) all three tables because each column is a feature of its tweet. I reshaped the dog\_stages (e.g., floofer, puppo, pupper, doggo) into a single column rather than multiple columns. I converted several columns to new datatypes:

- @archive\_df timestamp column its dtype converted from string to datetime
- @archive\_df name column its missing values converted from None to NaN
- @archive\_df doggo column its missing values converted from None to NaN
- @archive\_df floofer column its missing values converted from None to NaN
- @archive\_df pupper column its missing values converted from None to NaN
- @archive\_df puppo column its missing values converted from None to NaN
- @archive\_df text column is cut off with ellipses
- @tweets\_df tweet\_id column needs to convert its dtype from string to int

For each fixed issue, I identified the issue, stated my intention, then tested to ensure that I enacted my intention.

## Storage

I stored the data in two CSVs, one for the tweet data and one for the image predictions.

## Remaining Issues

- I noticed a strange interaction where a tweet's text was cutoff in the dataframe.
- I decided not to change the numerators or denominators for the dog ratings because they probably were not mistakes.
- Some of the variables did not seem necessary or useful (e.g., *in\_reply\_to\_status\_id*, *in\_reply\_to\_user\_id*, *retweeted\_status\_id*, *retweeted\_user\_id*).