

DATA WRANGLING PROJECT

We Rate Dogs Data.

Gathering Data Phase:

I started my project by downloading the 'twitter-archive-enhanced.csv' file manually. Then, created a folder named 'image_predictions' before I downloaded 'image-predictions.tsv' programmatically from Udacity's server using the requests library. Next I wrote it into image_predictions.tsv

'twitter_data' was created by accessing and downloading Twitter's JSON data using the tweepy library. Firstly, I extracted a list of tweet ID from the 'twitter-archive-enhanced.csv' file, then looped through each ID and query Twitter's API with the ID to get each tweet's JSON data. Subsequently, I recorded the data in a text file named 'tweet-json.txt', with each tweet's data written in a new line. After the query was completed and all the data was written in the text file, I read the text file line by line, obtained each tweet's information (tweet ID, retweet count, favorite count, and followers count) using the json library, and appended the information into an empty list.

Finally, I convert the list of dictionaries to a pandas DataFrame and saved it into 'twitter_data'

Assessing and Cleaning Data:

Some quality and tidiness issues were identified for the three tables.
Details of the issues identified and solutions are in the table below:

QUALITY

TWITTER ARCHIVE TABLE:

ISSUES	SOLUTIONS
<ul style="list-style-type: none">Keep original ratings (no retweets) that have images	<ul style="list-style-type: none">Delete retweets by filtering the NaN of retweeted_status_user_id
<ul style="list-style-type: none">drop columns not needed for our analysis	<ul style="list-style-type: none">drop columns
<ul style="list-style-type: none">Erroneous datatypes in these columns (tweet_id, in_reply_to_status_id, in_reply_to_user_id, timestamp, retweeted_status_id, source, retweeted_status_user_id, retweeted_status_timestamp, doggo, floofer, pupper, and puppo)	<ul style="list-style-type: none">Convert tweet_id to str from twitter_archiveConvert timestamp to datetimeconvert source to category datatype

<ul style="list-style-type: none"> Missing values in 'name' and dog stages represented as 'None' 	<ul style="list-style-type: none"> Change missing values in dog name to unnamed.
<ul style="list-style-type: none"> Some records have more than one dog stage 	<ul style="list-style-type: none"> Separate the dog stages to know which records have more than one dog stage.
<ul style="list-style-type: none"> Source column is in HTML-formatted string, not a normal string. 	<ul style="list-style-type: none"> Extract HTML values from source
<ul style="list-style-type: none"> Error in dog names (e.g a,an,actually) are not a dog's name. 	<ul style="list-style-type: none"> Change error name in dog name to None.
<ul style="list-style-type: none"> Some values in rating_numerator not showing proper float values. 	<ul style="list-style-type: none"> Spot those records and confirm changes made.
<ul style="list-style-type: none"> Text column includes a text and a short link. 	<ul style="list-style-type: none"> Remove hyperlinks in tweets.

IMAGE PREDICTION TABLE:

ISSUES	SOLUTIONS
<ul style="list-style-type: none"> Erroneous datatype (tweet_id) 	<ul style="list-style-type: none"> Convert tweet_id to str
<ul style="list-style-type: none"> Missing images (only 2075 counts out of possible 2356) 	<ul style="list-style-type: none"> Drop rows with missing images

TWITTER API TABLE:

ISSUES	SOLUTIONS
<ul style="list-style-type: none"> Erroneous datatype (tweet_id) 	<ul style="list-style-type: none"> Convert tweet_id to str

TIDINESS

TWITTER ARCHIVE TABLE:

ISSUES	SOLUTIONS
<ul style="list-style-type: none">doggo, floofer, pupper and puppo columns in twitter_archive table should be merged into one column named "dog_stage"	<ul style="list-style-type: none">Merge columns into one column named 'dog_stage'

TWITTER API TABLE:

ISSUES	SOLUTIONS
<ul style="list-style-type: none">twitter api table columns(retweet_count, favorite_count, followers_count)	<ul style="list-style-type: none">Merge table with twitter archive table.

IMAGE PREDICTION TABLE:

ISSUES	SOLUTIONS
<ul style="list-style-type: none">Image predictions table should be added to twitter archive table	<ul style="list-style-type: none">Merge table with twitter archive table.

Storing Cleaned Data

Now the data set is clean and ready for analysis. I saved the master table to twitter_archive_master.csv

Then I started my analysis.