

DATA WRANGLING PROJECT

We Rate Dogs Data

I started by importing the necessary libraries like pandas, numpy, seaborn, json, requests and tweepy.

Gathering Data Phase:

My project began with a manual download of the 'twitter-archive-enhanced.csv' file. Then I made a folder called 'image_predictions' and used the requests library to retrieve 'image-predictions.tsv' from Udacity's server. I then saved it as image-predictions.tsv.

'twitter_data' was produced by utilizing the tweepy package to retrieve and download Twitter's JSON data. To acquire the JSON data for each tweet, I first extracted a list of tweet IDs from the 'twitter-archive-enhanced.csv' file, then looped over each ID and queried Twitter's API using the ID.

I then saved the information in a text file called 'tweet-json.txt,' with each tweet's data typed on a separate line. After the query was finished and all of the data was saved to the text file, I read the text file line by line and used the json library to get the information for each tweet (tweet ID, retweet count, like count, and followers count) and appended it to an empty list.

Finally, I saved the list of dictionaries in the 'twitter data' folder as a pandas DataFrame.

Accessing and Cleaning Data:

The three tables were found to have some quality and tidiness concerns. The table below lists the issues and their remedies in detail:

Twitter Archive Dataframe

PROBLEMS	CLEAN(SOLUTIONS)
<ul style="list-style-type: none">Keep original ratings with images(no retweets).	<ul style="list-style-type: none">Delete retweets by filtering the NaN of retweeted_status_user_id
<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,	<ul style="list-style-type: none">Convert tweet_id to strConvert timestamp to datetimeConvert source to category datatype

retweeted_status_timestamp, doggo, floofer, pupper, and puppo)	
<ul style="list-style-type: none"> Name more than 745 records do not contain a valid name 	<ul style="list-style-type: none"> All names should start with a capital letter.
<ul style="list-style-type: none"> Source column is in HTML-formatted string in the source column, not a normal string 	<ul style="list-style-type: none"> Extract HTML values from source
<ul style="list-style-type: none"> doggo, floofer, pupper, puppo columns contain 'None' value. 	<ul style="list-style-type: none"> Nan should be used, I replaced it with Nan values
<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. Then to Nan values
<ul style="list-style-type: none"> Drop columns that is not needed for our analysis 	<ul style="list-style-type: none"> Drop columns

Image Prediction Dataframe

PROBLEMS	CLEAN(SOLUTIONS)
<ul style="list-style-type: none"> p1, p2, and p3 are inconsistent in a way capital and small letters are used in values. 	<ul style="list-style-type: none"> ○ Capitalize the first letter of all items in p1, p2, and p3 columns.
<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 Dataframe:

PROBLEMS	CLEAN(SOLUTIONS)
----------	------------------

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

TIDINESS

PROBLEMS	CLEAN(SOLUTIONS)
<ul style="list-style-type: none"> ▪ Image predictions table should be added to twitter archive table ▪ Twitter data api columns(retweet_count, favorite_count, followers_count) should be added to twitter archive table. 	<ul style="list-style-type: none"> ▪ Merge twitter api clean dataframe and image prediction into twitter archive table.

Storing Cleaned Data:

The data is now clean and ready to be analyzed. The master table was saved as twitter archive master.csv.