

DATA WRANGLING PROJECT

BY MUTHOLIB YUSIRA

INTRODUCTION

The dataset that is wrangled is the tweet archive of Twitter user @dog_rates, also known as WeRateDogs. WeRateDogs is a Twitter account that rates people's dogs with a humorous comment about the dog. The goal is to wrangle WeRateDogs Twitter data to create interesting and trustworthy analyses and visualizations. The Twitter archive is great, but it only contains very basic tweet information. Additional gathering from other sources, then assessing and cleaning will be required for "Wow!"-worthy analyses and visualizations.

DATA GATHERING

I started my project by gathering data from 3 different sources:

- The WeRateDogs Twitter archive. The `twitter_archive_enhanced.csv` file was provided to me by Udacity. This archive contains basic tweet data (tweet ID, timestamp, text, etc.) for all 5000+ of their tweets. I downloaded the 'twitter-archive-enhanced.csv' file manually into the workspace and loaded it into a dataframe titled, 'twitter_archive'
- The tweet image predictions, i.e., what breed of dog (or another object, animal, etc.) is present in each tweet according to a neural network. This file was provided to me by Udacity. I downloaded the 'image-predictions.tsv' programmatically from Udacity's servers using the requests library. Then I read it into 'image_predictions' Pandas dataframe.
- I accessed the third data using the 'tweet-json.txt' file provided by Udacity as an alternative. I read the text file line by line, obtained each tweet's information (tweet ID, retweet count, and favorite 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

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

Quality Issues:

Twitter_archive Dataset:

Issues	Solutions
<ul style="list-style-type: none">Rows with 'retweeted_status_id', 'retweeted_status_user_id', 'retweeted_status_timestamp', 'in_reply_to_user_id', 'in_reply_to_status_id' are retweets and not original ratings	<ul style="list-style-type: none">Keep only original ratings and no retweets, by deleting rows where the values in these columns are not null
<ul style="list-style-type: none">Dataset includes columns not needed for the analysis	<ul style="list-style-type: none">Drop the columns that are not needed
<ul style="list-style-type: none">Erroneous datatypes(tweet_id, timestamp, source)	<ul style="list-style-type: none">Convert columns to correct datatype
<ul style="list-style-type: none">Some rows in the 'rating_denominator' column are not equal to 10	<ul style="list-style-type: none">Drop rows where the rating_denominator is not equal to 10
<ul style="list-style-type: none">The string in the 'source' column is written in HTML format	<ul style="list-style-type: none">Convert string in 'source' column to normal format by removing the parts of the text with HTML format
<ul style="list-style-type: none">Rows in the 'text' column end with short links	<ul style="list-style-type: none">Remove the unnecessary links from the ends of text in the 'text' column
<ul style="list-style-type: none">'text' column includes ratings again, and some do not correspond to the 'rating_numerator' column	<ul style="list-style-type: none">Correct the incorrectly extracted 'rating_numerator' column and remove ratings from the 'text' column
<ul style="list-style-type: none">Null values in the 'expanded_urls' column	<ul style="list-style-type: none">Remove rows with null 'expanded_urls'
<ul style="list-style-type: none">Some urls in the 'expanded_urls' column are written multiple times in the same row	<ul style="list-style-type: none">Remove duplicated urls in the individual rows in the 'expanded_urls' column
<ul style="list-style-type: none">Incorrect dog names e.g 'a', 'an', 'not', 'very', 'the', etc	<ul style="list-style-type: none">Change dog names that are not proper nouns to None

Image_predictions and Twitter_data Dataset:

Issues	Solutions
<ul style="list-style-type: none">Erroneous datatypes(tweet_id)	<ul style="list-style-type: none">Convert 'tweet_id' to string

Tidiness Issues:

Twitter_archive Dataset:

Issues	Solutions
<ul style="list-style-type: none">4 different columns for the 4 dog stages (doggo, floofer, pupper, and puppo)	<ul style="list-style-type: none">Combine all the 4 columns into a single column and drop the other columns not needed

Image_predictions and Twitter_data Dataset:

Issues	Solutions
<ul style="list-style-type: none">Multiple columns for image predictions and confidence intervals	<ul style="list-style-type: none">Use a function to select the first true image predictions and confidence intervals, and deleting the rest of the columns
<ul style="list-style-type: none">Columns should be part of the twitter_archive dataset	<ul style="list-style-type: none">Merge the image_predictions table to the twitter_archive_clean table