

Wrangle Report

1. Introduction

Data preparation is always the hardest part of a data analyst's work flow, in this project, we will use the data wrangling skills to pull real-world data from Twitter, clean it, and do some analysis. We will get the original Twitter data from Twitter user @dog_rates, along with a image prediction dataset, to build our analysis.

WeRateDogs is a popular Twitter hash tag, as the name tells, people rate dogs with a denominator of 10 and the numerator is usually higher than 10 to show how lovely the dog is.

2. Gathering data

The data was gathered from three sources:

- **Enhanced Twitter archive:** This part is kind of on-hand data, stored in the twitter_archive_enhanced.csv
- **Image prediction:** We get the image prediction data from web scraping
- **Twitter API:** This dataset is archived from Twitter's API and parsed from JSON to csv

3. Assessing Data

After assessing the datasets, I summarized several quality and tidiness problems of them:

Quality problems

df_twitter_archive:

remove tweets that has been replay as its not original.

2-remove tweets that has been retweet as its not original.

3-When we back to twitter account"@dog_rates", there are ratings that are incorrect. I ordered the ratings from low to high and looked at the extremes only for incorrect ratings therefore there are likely more than I missed and will be difficult to find them all programmatically. Examples where things may have gone wrong is the use of decimals, or when two instances of numbers separated by a slash are present in 1 text and I assume the first was chosen. Also, there are ratings with decimals such as 13.5/10, 9.5/10 have been incorrectly extracted as 5/10 (in addition to other numbers with decimals such as 11.26 and 11.27). There are instances of 1/2 and 50/50 which are not ratings such signifying "half" which have been considered as ratings. Finally, use of 4/20 and 24/7 has been confused as ratings. For future analysis it could be confusing to interpret unstandardized ratings. It is their gimmick to give dogs a rating of 100% but not all are above 100% so it could be interesting to see what % are below or above 100% and how this changed overtime by calculating a single value for rating.

There are many columns in this dataframe making it hard to read, and some will not be needed for analysis.

4-Correct numerator bigger than 14.

5-Correct denominators not equal 10.

6- Combine rating numerator and rating denominator columns into one column.

7-replace the url from value of source .

8-The name column has many invalid values like , a, an, the.

df_image_predictions:

1- drop duplicate jpg_url.

2-p1,p2 and p3 have inconsistent capital words.

3-p1, p2, and p3 contain underscores instead of spaces in the labels

Tidiness:

Tidy data is a standard way of mapping the meaning of a dataset to its structure. A dataset is messy or tidy depending on how rows, columns and tables are matched up with observations, variables and types. In tidy data:

1- replace the url from value of source.

2-There are multiple dog stages columns present e.g. doggo, pupper, etc. They should be merged into 1 column.

4. Cleaning data

The data quality and tidiness problems mentioned above was cleaned by multiple methods including pandas join, regular expression, combining multiple columns, pandas sub setting, removing missing values and so on.

In the end of this part, I stored the cleaned version of the data into a csv file for future usage.

5. Conclusion

This data wrangling project is a good practice of what I learned in the course, I've struggled to find subtle data problems, and cleaned those hard, time-consuming data quality problems. The hardest part of this project is the dealing with strings by using regular expressions.

In the end, I'm able to build a well-structured wrangling processing notebook with details in every part: exhaustive steps in how I gathered, assessed, and cleaned the data with every problem and its solving process illustrated.

The analyzing and visualizing part will be in another file called act_report.htm