

# wrangle\_report

January 28, 2023

## 1 Data wrangling report of tweet archive of Twitter user WeRateDogs (@dog\_rates).

### 1.1 Introduction

The dataset that we are wrangling (and analyzing and visualizing) 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. These ratings almost always have a denominator of 10. The numerators, though? Almost always greater than 10. 11/10, 12/10, 13/10, etc. Why? Because "they're good dogs Brent." WeRateDogs has over 4 million followers and has received international media coverage.

### 1.2 Data Gathering

We are using three datasets with different file formats and gathering techniques.

1. The WeRateDogs Twitter archive (twitter\_archive\_enhanced.csv) is downloaded from udacity resources and uploaded directly to the workspace environment
2. The tweet image predictions (image\_predictions.tsv) is downloaded programmatically using the Requests library and the given URL
3. Additional data from the Twitter API (tweet\_json.txt) is downloaded either by using the Tweepy library or reading the already given file from udacity

### 1.3 Data Assessing

All the three datasets were assessed using:

#### 1.3.1 1. Visual assessment:

Where we can explore the entire dataset

#### 1.3.2 2. Programmatic assessment:

Some pandas functions were used to help assessing the datasets: [.head(), .tail(), .sample(5), .info(), .describe(), .duplicated(), .value\_counts(), .sort\_values(), .isnull()]

After assessing the datasets, some quality and tidiness issues were found:

### 1.3.3 Quality issues

Completeness, validity, accuracy, consistency (content issues)

- The original tweets only are needed, no retweets
- Many un-needed columns
- ['rating\_numerator', 'rating\_denominator'] columns type is int
- duplicated jpg\_url
- naming issues
- timestamp type is object instead of being datetime
- the extracted numerators did not contain the decimals values
- missing [rating] column that divides the numerator by the denominator
- the tweet\_id type is int instead of being object or string

### 1.3.4 Tidiness issues

rows, columns, and tables (structural issues)

- Many datatypes (doggo, floofer, pupper and puppo columns)
- day, month, year are in one timestamp column
- many dog type and confidence level columns
- twitter\_archive\_clean, image\_predictions\_clean, and twitter\_api\_clean must be combined

## 1.4 Data Cleaning

Here, all the discovered issues were solved.

### 1.4.1 Quality Issues

- Delete retweets by filtering the null values of retweeted\_status\_user\_id to keep the original tweets only.
- Delete the unneeded columns
- Convert ['rating\_numerator', 'rating\_denominator'] columns type from int to float
- Delete the duplicated jpg\_url
- Correct naming issues
- Convert timestamp to datetime
- Correct decimal values in the numerator of the rating
- Create a new [rating] column by deviding the numerator by the denominator
- convert tweet\_id type to string

#### 1.4.2 Tidiness Issues

- Melt [doggo, floofer, pupper, puppo] columns to dogs and dogs\_stage column
- Extract [year, month and day] values to new columns
- Keep only the dog type with the max confidence level
- merge twitter\_archive\_clean and image\_predictions\_clean
- merge df\_twitter1 and twitter\_api\_clean

#### 1.4.3 Storing Data

Save the gathered, assessed, and cleaned master dataset to a CSV file named "twitter\_archive\_master.csv".