

---

# Project 2: Wrangle and Analyze Data

---

## Data Wrangling Report

*Prepared By: Suzan Hamza*

### Introduction

This report briefly describes the efforts used to wrangle *WeRateDogs* Twitter data to create interesting and trustworthy analyses and visualizations.

The following packages (libraries) were installed and imported at the beginning of the `wrangle_act.ipynb` notebook.

- *pandas*
- *NumPy*
- *requests*
- *tweepy*
- *json*
- *matplotlib*
- *seaborn*

### Gathering Data for this Project

Each of the **three pieces of data** as described below were gathered in a Jupyter Notebook titled `wrangle_act.ipynb`:

#### 1) The WeRateDogs Twitter archive.

This file named `twitter_archive_enhanced.csv` was downloaded manually and read into a dataframe using pandas `read_csv()`.

## 2) Tweet Image Predictions File

This file (`image_predictions.tsv`) is hosted on Udacity's servers and was *downloaded programmatically* using the `Requests` library and the following URL: [https://d17h27t6h515a5.cloudfront.net/topher/2017/August/599fd2ad\\_image\\_predictions/image\\_predictions.tsv](https://d17h27t6h515a5.cloudfront.net/topher/2017/August/599fd2ad_image_predictions/image_predictions.tsv)

Then the file was read into a dataframe using pandas `read_csv()`.

## 3) Additional Data via the Twitter API

First, a Twitter developer account was created, then it was used to generate the *Consumer* API keys, and the Access Token and Access Token Secret needed.

Those secret credentials were stored in a separate `.env` file, and were loaded using Python's `dotenv` library.

Each tweet's **retweet count** and **favorite** ("like") **count** was gathered using the **tweet IDs** in the WeRateDogs Twitter archive, by querying the Twitter API for each tweet's JSON data using Python's `Tweepy` library.

Then each tweet's entire set of JSON data was stored in a file called `tweet_json.txt` file. Each tweet's JSON data was written to its own line.

Then this `.txt` file was read line by line into a pandas *DataFrame* with `tweet ID`, `retweet count`, and `favorite count`.

## Assessing Data for this Project

After gathering each of the above pieces of data and loading them into a separate pandas dataframe, first they were assessed *visually* by displaying each dataframe in the notebook and by opening the files in Excel to investigate its contents for **quality** and **tidiness** issues.

Then a *programmatic assessment* was conducted using pandas methods, like `info`, `head`, `describe` and `value_counts`.

The following **quality** and **tidiness** issues were detected, including several issues that did not satisfy the Project Motivation:

## Quality Issues

### archive\_df table

1. **78** entries are **replies** and not original tweets (**in\_reply\_to\_status\_id** and **in\_reply\_to\_user\_id** has values).
2. **181** entries are **retweets** and not original tweets (**retweeted\_status\_id**, **retweeted\_status\_user\_id** and **retweeted\_status\_timestamp** has values).
3. **in\_reply\_to\_status\_id**, **in\_reply\_to\_user\_id**, **retweeted\_status\_id**, **retweeted\_status\_user\_id** and **retweeted\_status\_timestamp** columns are not useful for analysis and should be removed.
4. **59** entries do not contain images (**expanded\_urls** is null).
5. **tweet\_id** is integer instead of string (object).
6. **timestamp** is a string and not datetime.
7. **rating\_numerator** is integer instead of float.
8. **55** dog **names** incorrectly extracted as **a**, because the tweet was in the format *'This is a ...'* instead of *'This is (Dog Name) ...'*.
9. **745** dog **names** incorrectly extracted as **None**, because the tweet was not in the format *'This is (Dog Name) ...'*.
10. **tweet\_id = 810984652412424192** does not include a dog rating, the text= 'Meet Sam. She smiles 24/7 & secretly aspires to be a reindeer.' and the rating was incorrectly extracted as 24/7.
11. **tweet\_id = 666287406224695296** rating was incorrectly extracted as 1/2 while it should be 9/10. text = 'This is an Albanian 3 1/2 legged Episcopalian. Loves well-polished hardwood flooring. 9/10'.
12. **6** entries in the **rating\_numerator** column were incorrectly extracted, because the original rating had decimal points, so only the number after the decimal point was extracted.
13. Other incorrectly extracted ratings for **tweet\_ids:** **775096608509886464**, **740373189193256964**, **716439118184652801**.

#### img\_pred\_df table

1. `tweet_id` is integer instead of string (object).
2. Only **2075** image predictions are available, which indicates that another **281** tweets have no images and should be excluded from the `archive_df` table.

#### tweepy\_df table

1. `tweet_id` is integer instead of string (object).

### Tidiness Issues

#### archive\_df table

1. Dog Stages variables split into four columns (`doggo`, `floofer`, `pupper` and `puppo`) instead of one.
2. A single observational unit (Tweet information) is stored in multiple tables (`archive_df` and `tweepy_df`).

## Cleaning Data for this Project

Each of the issues documented in the assessment phase were cleaned as follows:

### 1) Missing Data

Missing data or data that does not satisfy the Project Motivation was cleaned first.

- Removed all **replies** from archive table; entries where `in_reply_to_status_id` and `in_reply_to_user_id` are not null.
- Removed all **Retweets** from archive table; entries where `retweeted_status_id`, `retweeted_status_user_id` and `retweeted_status_timestamp` are not null.
- Dropped `in_reply_to_status_id`, `in_reply_to_user_id`, `retweeted_status_id`, `retweeted_status_user_id` and `retweeted_status_timestamp` columns from archive table, because they will not be useful to our analysis.
- Dropped entries where `expanded_urls` is null from archive table, because that means that these tweets do not contain any images.

- Removed entries from archive table with `tweet_id` that did not exist in `img_pred` table, because that means that these tweets do not contain any images.

## 2) Tidiness Issues

- Combined the four Dog Stages columns (`doggo`, `floofer`, `pupper` and `puppo`) into one column named `dog_stages`.
- Merged the `retweet_count` and `favorite_count` columns with the `archive` dataframe.

## 3) Remaining Quality Issues

- Converted `tweet_id` data type from integer to string in both `archive` and `img_pred` tables.
- Converted `timestamp` from string data type to datetime in `archive` table.
- Converted `rating_numerator` from integer data type to float in `archive` table.
- Corrected some of the dog names that were incorrectly extracted as `a` or `an` or `None` programmatically.
- Replaced the remaining `a`, `an` and `None` name values with `NaN`.
- Replaced the `None` values in the `dog_stages` column with `NaN`.
- Inspected the tweet's text for `tweet_id = 666287406224695296` and corrected the rating from 1/2 to 9/10 manually, since it is a one off occurrence.
- Extracted the correct `rating_numerator` from tweet's text using regular expressions, for tweets that had decimal ratings.
- Inspected the tweet's text for `tweet_ids: 740373189193256964` and `716439118184652801` and corrected their ratings manually.
- 

*Note:* The rating for `tweet_id = 810984652412424192` was not found and could not be corrected, but that issue will not affect our analysis.

The result was stored in two high quality and tidy master pandas DataFrames, `archive_clean` and `img_pred_clean`.

## Storing Data for this Project

Finally, the two clean DataFrame(s) were stored in a **CSV file** with the main one named `twitter_archive_master.csv` and the other named `image_predictions_clean.csv` for image predictions data.