

# WeRateDogs - Wrangle Report

By Ugochukwu Uche-Agada

July, 2022

This is a report based on the data wrangling activity carried out on the Tweet archives of WeRateDogs Twitter account. It covers the 3 actions in data wrangling: Gather, Assess and Clean.

## Data Gathering

Data required for analysis sometimes is not always available at one location or data set. In this case, data had to be sourced from three separate locations and with three different methods of gathering.

1. From file on hand

The file **twitter\_archive\_enhanced.csv** was already provided for so all that was done was to load the csv file into a dataframe.

2. Using the Requests Library

The **image-predictions.tsv** was downloaded programmatically using the requests library and a link provided by Udacity.

3. Using Python's Tweepy library

The last of the three files was downloaded using Tweepy library. This was made possible by the elevated access granted to my Twitter Developer Account. The **tweet\_json.txt** was downloaded here.

## Assess Data

When assessing, you're inspecting your data set for two things:

Data quality issues: Data that has quality issues have issues with content like missing, duplicate, or incorrect data. This is called dirty data.

Lack of tidiness: Data that has specific structural issues that slow you down when cleaning and analyzing, visualizing, or modeling your data later.

## **Source: Udacity**

After using visually and programmatically assessing the datasets, **quality** and **tidiness** issues were identified.

### **Quality Issues**

#### **tweet\_arch DataFrame**

1. row values in retweeted\_status\_id, retweeted\_status\_user\_id, retweeted\_status\_timestamp are not needed in the analysis
2. values in in\_reply\_to\_status\_id and in\_reply\_to\_user\_id are not needed in the analysis
3. timestamp has wrong datatypes
4. There are outliers in the numerator and denominator columns
5. Identical values in expanded\_url
6. Source column is noisy since it can be shortened.
7. Incorrect values in name column

#### **image\_predict DataFrame**

8. Inconsistent names in P1, P2 and P3

#### **tweet\_df DataFrame**

9. wrong column datatypes for retweets and favorites

### **Tidiness issues**

1. Retweets are not needed since only original tweets are needed.
2. Replies are not needed since only original tweets are needed.
3. From the Tweet\_arch data frame, Puppo, Floofer, pupper and doggo should be values on a single column
4. Timestamp should be separated to capture day, month and year.
5. 'img\_num', 'p1', 'p1\_conf', 'p1\_dog', 'p2', 'p2\_conf', 'p2\_dog', 'p3', 'p3\_conf', 'p3\_dog' does not fulfil requirements of being different columns.
6. tweets\_df should be merged with tweet\_arch data frame

## Clean Data

All the issues detected while assessing the data was cleaned using code and tested afterwards

The following actions were taken to clean the quality issues:

## Quality Issues Cleaning

### tweet\_arch DataFrame

1. row values in retweeted\_status\_id, retweeted\_status\_user\_id, retweeted\_status\_timestamp are not needed in the analysis
2. values in in\_reply\_to\_status\_id and in\_reply\_to\_user\_id are not needed in the analysis
3. timestamp has wrong datatypes
4. There are outliers in the numerator and denominator columns
5. Identical values in expanded\_url
6. Source column is noisy since it can be shortened.
7. Incorrect values in name column

### image\_predict DataFrame

8. Inconsistent names in P1, P2 and P3

### tweet\_df DataFrame

9. wrong column datatypes for retweets and favorites

## Tidiness issues Cleaning

1. Columns relating to retweets were dropped .
2. Columns relating to replies were dropped.
3. Puppo,Floofer,pupper and doggo were put in a single column **dog\_stage** using the melt() function
4. Timestamp was separated to capture day, month and year.
5. 'img\_num','p1','p1\_conf','p1\_dog','p2','p2\_conf','p2\_dog','p3','p3\_conf','p3\_dog' does not fulfil requirements of being different columns hence the named columns were melted into the **breed** and **confidence** columns.
6. The three data frames were merged into one master data.

At the end of this wrangling process, the final data set was saved in the project folder ready for analysing.