

# 2021

## Human Besties



**WeRateDogs™ (author)** ✓  
@dog\_rates

Following

This is Stuart. He's sporting his favorite fanny pack. Secretly filled with bones only. 13/10 puppered puppo [#BarkWeek](#)



1:02 PM - 24 Jul 2017

Salma Gouda  
Udacity Projects  
1/6/2021

# **Data Wrangling Report**

By Salma Ahmed Gouda

1/6/21

This report is presented as an assignment for the Udacity Data Analyst Nanodegree. In this report, the main steps involved in the data wrangling of the Twitter account "WeRateDogs" are illustrated.

## **Data Gathering**

In this first step, we collect data from different sources through various methods. For this project, there were three main sources for the data required:

1. Twitter\_archive\_enhanced.csv file, the first file was provided in the Udacity project classroom and downloaded manually to our working directory and then imported into our working environment using Pandas function "pd.read\_csv".
2. Image\_predictions.tsv file, the second file that has been hosted on a webpage and downloaded from its relevant URL using the Requests library "get" function and Pandas function "pd.read\_csv". This file contained image predictions for the dog breeds obtained through a neural network on most of the tweets in the twitter archive file.
3. The third and final dataset was gathered from Twitter REST API via the Tweepy library by querying the API to obtain extra information pertinent to the tweets ids in the first file, e.g. retweet count and favourite count aspects. In case of students having troubles establishing their Twitter developer account, the code and file containing information was provided for them in the Udacity project classroom.

## Data Assessing

In this second step, we investigate our imported datasets both visually and programmatically for quality and tidiness issues.

### *Visual Assessment:*

The visual assessment is done on a spreadsheet application or in a Jupiter notebook to find some issues that can be easily detected with the eye.

### *Programmatic Assessment:*

The programmatic assessment is done in a Jupiter notebook using Pandas functions:

- **head()**
- **info()**
- **columns()**
- **value\_counts()**
- **sample()**
- **describe()**

While assessing the data, I found the following issues:

### ***Tidiness Issues:***

- All tables should be part of one master dataset then duplicate columns should be removed
- a. archive\_df
- values are column names (doggo, floofer, pupper, puppo)
  - timestamp column should be date and time columns

### ***Quality Issues:***

#### a. archive\_df

- retweets (3 tables)
- tweet\_id dtype should be a string
- inaccurate name data
- inaccurate dog ratings (numerator and denominator data)
- some columns aren't useful for analysis

#### b. image\_predictions

- 66 duplicated urls
- Various columns for dog breed predictions and their confidence levels
- tweets with no images
- some columns aren't useful for analysis

Note: There're more issues with the datasets, but those are the ones I decided to fix according to the project rubric

## **Data Cleaning**

In this third and last step in data wrangling, we define the issues we found, build a code to fix it, and then test if our code worked. We should end up with a clean dataset for further analysis and visualization.

Before cleaning, we should first make copied of the datasets to work with.

The issues addressed in the assessment process were fixed as follows:

Tidiness Issue	Solution
<ul style="list-style-type: none"><li>• All tables should be part of one master dataset then duplicate columns should be removed</li></ul>	<ul style="list-style-type: none"><li>• Merge the 3 datasets into one using "concat" function then remove duplicate columns</li></ul>
<ul style="list-style-type: none"><li>• Values are column names (doggo, floofer, pupper, puppo)</li></ul>	<ul style="list-style-type: none"><li>• Create one column to include the values (doggo, floofer, pupper, puppo) using "extract" function, then drop the values columns using "drop" function</li></ul>
<ul style="list-style-type: none"><li>• "timestamp" column should be date and time columns</li></ul>	<ul style="list-style-type: none"><li>• Convert timestamp column into datetime format using "to_datetime" function then convert timestamp column into date and time columns</li></ul>

Quality Issue	Solution
<ul style="list-style-type: none"> <li>Retweets (3 tables)</li> </ul>	<ul style="list-style-type: none"> <li>Keep only original tweets (delete all retweets) by filtering NaN values in the <code>retweeted_status_user_id</code> column</li> </ul>
<ul style="list-style-type: none"> <li>"Tweet_id" dtype should be a string</li> </ul>	<ul style="list-style-type: none"> <li>Change <code>tweet_id</code> dtype from <i>int</i> to <i>str</i> using "astype" function</li> </ul>
<ul style="list-style-type: none"> <li>Inaccurate name data</li> </ul>	<ul style="list-style-type: none"> <li>Change inaccurate name data to None using "replace" function</li> </ul>
<ul style="list-style-type: none"> <li>Inaccurate dog ratings (numerator and denominator data)</li> </ul>	<ul style="list-style-type: none"> <li>Fix dog ratings (numerator and denominator issues) both manually and programmatically then merge both columns into one "rating" column</li> </ul>
<ul style="list-style-type: none"> <li>66 duplicated "<i>jpg_urls</i>"</li> </ul>	<ul style="list-style-type: none"> <li>Drop 66 duplicated "<i>jpg_url</i>" using "drop" function</li> </ul>
<ul style="list-style-type: none"> <li>Various columns for dog breed predictions and their confidence levels</li> </ul>	<ul style="list-style-type: none"> <li>Create 1 column for "<i>dog_breed</i>" and 1 column for "<i>conf_lvl</i>"</li> </ul>
<ul style="list-style-type: none"> <li>Tweets with no images</li> </ul>	<ul style="list-style-type: none"> <li>Remove tweets with no images</li> </ul>
<ul style="list-style-type: none"> <li>Some columns are not useful for analysis</li> </ul>	<ul style="list-style-type: none"> <li>Drop all columns not useful for analysis using "drop" function</li> </ul>