# Udacity Data Analyst Nanodegree

## Wrangle and Analyze Data Project

## Project Details

The tasks in this project are as follows:

- Data wrangling, which consists of:
    - Gathering data
    - Assessing data
    - Cleaning data
- Storing, analyzing, and visualizing the wrangled data

## Gathering Data:

The data for this project consist on three different dataset that were obtained as following:

- **The Twitter archive file:**
  The (twitter_archive_enhanced.csv) file was provided by Udacity, I downloaded it manually and stored it as a DataFrame.

- **The tweet image predictions file:**
  This file (image_predictions.tsv) is hosted on Udacity's servers; I downloaded it programmatically using the Requests library and stored it as a DataFrame.

- **Data from Twitter:**
  The (tweet_json.txt) text File with jSON structure, provided by Udacity, I read the tweet's JSON data from this file line by line into a list of dictionaries then create a DataFrame from this list.

## Assessing data:

- **Visually:**
  By checking the Twitter archive CSV file in Excel.

- **Programmatically**:
  By using different methods:
    - info()
    - duplicated()
    - value_counts()
    - describe()
    - sort_values
    - head()

I found the following issues:

| • **Quality issues:** | |
|---|---|
| (twitter_archive) table:<br>- tweet_id is an int not string.<br>- Timestamp is an object not data time.<br>- Retweets columns are not necessary, as this project concerns wrangling and exploring original rating.<br>- By manually checking the top 10 rating_numerator with Rating_denominator =10, I found that:<br>   - 1 not dog image.<br>   - 3 without image.<br>   - 6 extraction issues, ex:the actual rating is "11.27/10" while the rating in the set is "27/10 | - Some columns will not be used for analysis.<br>- 'source' column contain HTML tags.<br>- Rating_denominator not equal to 10 in 23 rows:<br>   - 1 row with denominator =0.<br>   - 2 rows < 10.<br>   - 20 > 10, the images in most of them are contains more than 1 dog. |
| (tweet_json) table:<br>  - tweet_id is an int not string. | (image_predictions) table:<br>  - tweet_id is an int not string. |
| • **Tidiness issues:** | |
| (twitter_archive) table:<br>- Dog Stages represented in four columns (doggo, floof, pupper, and puppo). | |
| (image_predictions) table:<br>- The table should be part of the master table. | |
| (image_predictions) table:<br>- The table must be merged with the other tables. | |

## Cleaning data:

- First of all I create a copy of the three DataFrames to keep the originals.
- Example for cleaning:
  - Remove the HTML tags by replacing them with readable sources.
  - Drop the row with rating_denominator values not equal to 10.
  - Drop the rows with rating_numerator > 17.
  - Create column for dog Stages and drop columns :(doggo,floofer,pupper,puppo).
  - Merge the three tables.
- For each issue described in the assessing section I followed the approach of  Define, Code and Test.
- I used different methods
  - astype()
  - type()
  - to_datetime()
  - isnull()
  - value_counts()
  - list()
  - drop()
  - replace()
  - sum()
  - sort_values()
  - extract()
  - head()
  - merge()