

# Data wrangling report

By: Eslam Ali Abou-shashaa

Dec. 2020

As an assignment for the Udacity Data Analysis Nanodegree; This report illustrates the main steps involved in the data wrangling of Twitter account "WeRateDog"

## Data Gathering

in this step collecting data takes place. For this project, there were three main sources for the data to deal with:

1-Twitter\_archive\_enhanced.csv file this file was delivered by email and downloaded manually to our working directory and then imported into our working environment using pandas function "pd.read\_csv"

2-Image\_prediction.tsv is the second file that has been hosted in a webpage and downloaded from its relevant URL using the requests library get function and pd.read\_csv pandas function . this file encompassed image predictions for the dogs' breeds obtained through a neural network on most of the tweets in the archive file

3- the final dataset was gathered from tweet-json.txt after then use json library to use file and use loop to read line by line in file to extract retweets count and favorite count

## Assess

this step we use two type from Assess

1-visual this type i use google sheet to see all data and select issues

2-programmatic this type use functions like info,describe,value\_counts

function info this function merged information like type of data in columns and number of None in columns

function value\_counts this function appear number of character in columns

function unique this function appear all character in specially columns

from above step we find two type from issues

quality issues like some name don't correct like a

and some data is missing like data in in\_reply\_to\_status\_id ,

in\_reply\_to\_user\_id ,name , retweeted\_status\_id

,retweeted\_status\_user\_id and retweeted\_status\_timestamp

timestamp is object should be data time

most rating is over rating

## cleaning

we should cleaning this data to analysis it

let's cleaning the data

first remove all columns which don't have data like in\_reply\_to\_s

,retweeted\_status\_timestamp to remove this columns we use drop function in pandas

convert timestamp to data time used is operation to\_datetime function in pandas

make new columns to dog type and remove "puppo, doggo , floofer, poppur" columns

make the three table is one table

