

DAND\_C9

WeRateDogs

Wrangle\_Report

Data Wrangling Project



**MOSHIRA EZZAT**

**moshiraezzat@pg.cu.edu.eg**

## Introduction

This is a data wrangling analysis for WeRateDogs Twitter Archive, which is a Twitter account for rating the dog's cool images with humorous comments. These ratings almost have denominator as 10, and numerators that are greater than 10.

## Wrangling Process

### i. Gathering

In this data analysis, the data was gathered as three datasets;

- 1- The WeRateDogs Twitter archive, manually downloaded as csv file.
- 2- The tweet image predictions, programmatically downloaded as tsv file.
- 3- Tweets JSON data queried from Twitter API using tweepy

The three datasets were used to create 3 dataframes by pandas;

- 1- archive\_df,
- 2- image\_prediction\_df,
- 3- api\_df

### ii. Assessment

The dataframes were visually and programmatically assessed, to identify quality and tidiness issues summarized in the following table:

<u>Tidiness Issues</u>	<u>Quality Issues</u>
1- Dog stage represented in 4 columns	1- Erroneous dtype of tweet_id columns
2- Three predictions and confident columns	2- Erroneous dtype of timestamp
3- Inconsistent name of id column in api_df	3- Presence of retweets
4- Data separated in three tables	4- Tweets without images
	5- Unnecessary columns
	6- Invalid values in name column
	7- Tweet's sources are masked
	8- Ratings are not calculated

### iii. Cleaning

First observation was the presence of inconsistent name of *'id'* column in *api\_df*, that is renamed into *'tweet\_id'* to be consistent with that of archive and image-prediction dataframes.

The dtypes of *'tweet\_id'* columns were *'integr'* in all dataframes, as well as the dtype of *'timestamp'* column in the *archive\_df* was *'string'*. These issues were solved using *astype()* function to convert *tweet\_id* from *integr* to *string*, and *to\_datetime()* function to convert the *timestamp* dtype from *string* to *date time*.

Contrarious to the tidiness rules, the dog stage in *archive\_df* was represented in four separated column, [*'doggo'*, *'pupper'*, *'floofer'*, *'puppo'* ], so the four columns were combined into one column called *dog\_stage*.

Similarly, the columns of predictions *p1,p2,p3* and its confidence *p1\_conf, p2\_conf, p3\_conf* in the *image\_prediction\_df* were unified into single prediction column according to the highest confidence and the boolean values in the *p1\_dog, p2\_dog, p3\_dog* columns.

One of the principal key points for this project is to analyze original tweets, not the retweets or replies. Through the assessment step, it was obvious the presence of non-null records in the *in\_reply\_to\_status\_id*, and *retweeted\_status\_id* columns of the *archive\_df*; referring to the presence of retweet and replies. So, these non-null records were dropped, then all retweets and replies columns [*'in\_reply\_to\_status\_id'*, *'in\_reply\_to\_user\_id'*, *'retweeted\_status\_id'*, *'retweeted\_status\_user\_id'*, *'retweeted\_status\_timestamp'*] were dropped totally since they are unnecessary in this analysis.

Also, tweets missing *expanded\_urls* in the *archive\_df*, were dropped as they were actually missing the image itself.

The invalid values in the *name* column of *archive\_df* were replaced with *'None'*

The sources of the tweets were masked at the end of *html* string. The clean types of sources were extracted and displayed clearly.

The rating system of WeRateDogs follows a scale of 1 to 10, but it is invariably may have ratings exceeds maximum, so it was decided to only calculate the ratings by dividing the numerator by 10 in a new column *[ratings]*. Then the *rating\_numerator* and *rating\_denominator* columns were dropped.

At the end of the cleaning steps, the dataframes were merged into a master\_df using *tweet\_id* and storied as *twitter\_archive\_master.csv* for further iteration of the wrangling process and data analysis.