

## **Report on WeRateDogs data wrangling**

To wrangle WeRateDogs Twitter data, I went through 3 consecutive steps; gather, assess, and clean datasets.

### **Firstly, gathering 3 datasets from 3 different resources:**

- 1- Enhanced\_Twitter\_Archive (csv file), provided by Udacity (downloadable) and contains basic tweet data for 2356 records/tweets including information such as tweet's text, rating, name, stage (i.e. doggo, floofer, pupper, puppo).
- 2- Image\_Predictions (tsv file), obtained by using 'requests' library and the provided url. The resultant table includes image predictions (the top three only) of each tweet id, image URL, and the image\_num that corresponded to the most confident prediction.
- 3- Additional information were gathered from Twitter API as json\_txt file using the authentication keys provided by Twitter API. df\_api was filtered to include only id, retweet\_count and favorite(like)\_count columns.

The three datasets were converted into pandas dataframes for assessing/cleaning/analyzing/visualization.

### **Secondly, assessing the dataframes for quality and tidiness:**

Now, we have 3 dataframes; df\_archive, df\_image\_predictions, df\_api

The dataframes have been explored manually and programmatically. The manual assessment was based on opening the files either in jupyter notebook web-browser or from saved files on my computer, then scrolling up and down to get familiar with the data structure. Regarding the programmatic assessment, that was performed by using different methods in pandas such as; .head(), tail(), .sample() , .shape to get an overview of the data structure. Also, .info() was used to check the data entries for each variable along with their data types, .duplicated() to check duplicated records , .isnull() to check null entries, and .value\_counts() for specific variables to deepen in the data details. Though, I did not focus on .describe() since most numerical data are available in rating variables that have outliers (probably, as erroneous values or exaggeration from the dogs' lovers). So, I was happy with sorting out the rating score as a string.

### **Thirdly, cleaning the dataframes by using (Define-Code-Test) protocol:**

#### **Quality (create a copy of each dataframe to work with, so we still have the original data)**

- 1- Quality\_1 (validity/consistency): dtype of tweet\_id in the three dataframes is integer, so convert into a string to avoid performing any mathematical calculation on that column. The dtype of timestamp column in df\_archive is object, then convert into datetime (might be used in the future) (use as.type()).
- 2- Quality\_2 (consistency): rename id column in df\_api as tweet\_id (parent key), for consistency with other dataframes (use .rename()).

- 3- Quality\_3 (completeness/consistency): some data entries in df\_archive don't have image\_url, so filter records with incomplete information by using tweet\_id in df\_image\_predictions as a guide (make a list of tweet\_id in df\_image, then clean using .isin()).
- 4- Quality\_4 (accuracy): Remove retweets and replies from df\_archive since we are interested in only the original tweets (assign the retweets and replies entries, each to a new variable then drop them from df\_archive).
- 5- Quality\_5 (consistency): filter df\_image\_prediction from retweet and replies records by matching against the cleaned df\_archive using tweet\_id as a guide (use isin()).
- 6- Quality\_6: in df\_archive, drop columns; in\_reply\_to\_id, in\_reply\_to\_user\_id, retweeted\_status\_id, retweeted\_status\_user\_id and retweeted\_status\_timestamp, since they are redundant, and useless (use .drop()).
- 7- Quality\_7: replace 'None' values in the last 4 columns (dogs'\_developmental \_tages) with pandas NaN (first, replace 'None' with empty spaces .replace(), then fillna after merging).
- 8- Quality\_8: replace erroneous values in name column by first creating a pattern from the tweet text, then a for loop to replace 'a', 'an' names with their correct values from the text. Then replace 'None' entries with NaNs.
- 9- Quality\_9: replace erroneous entries in column name with correct values extracted from the text pattern (.str.extract()).
- 10- Quality\_10 (accuracy): extract source names from the HTML tag (.str.extract()).

### **Tidiness**

- 1- Tidiness\_1: In df\_archive, combine 4 columns of dogs' stages in one single column (one variable) by creating a new column, dog\_stage and add the data string from 4 columns into it, then fill in the empty spaces with NaN.
- 2- Tidiness\_2: In df\_image\_predictions, rename columns (p, p\_conf, and p\_dog) with more descriptive names to be more understood (use .rename(columns={'oldname':'newname'})).
- 3- Tidiness\_3: Each type of observational unit forms a separate table, i.e., df\_image\_predictions and df\_api are displaying some extra attributes of the tweets in df\_archive. So, after cleaning, combine the 3 dataframes into one pandas dataframe (use .merge() on tweet\_id).

### **Storing data**

After cleaning and merging the dataframes, save as (twitter\_archive\_master.csv) to act on.