

## Wrangle Project Summary

### Wrangling steps:

#### 1- Gathering Data:

- a. Gather data from 'twitter-archive-enhanced.csv': file uploaded directly as csv file (using the read\_csv function).
- b. Gather data from "image\_predictions.tsv": the file is collected from the provided url (using the requests library), then the data, folder created from data (using the OS library), Read the file, using the tab separation (using read\_csv function).
- c. Gather data from "Twitter API": Twitter API json file imported line by line (using json library), convert data to dataframe,
- d. All the previous dataframes are saved as csv files (using pandas function to\_csv).

#### 2- Assessing data: All the dataframes are inspected for quality and tidiness issues (using functions as info, head, tail, columns value\_counts, column sort\_values, and isnull().sum()), and the summary of the findings as follow:

- a. image\_predictions:
  - i. Quality issues:
    - Q1- Column Header is undecryptive: img\_num p1, p1\_conf, p1\_dog, p2, p2\_conf, p2\_dog, p3, p3\_conf, p3\_dog.
    - Q2- some of the dogs types are incorrect, such in p1 column: car\_mirror in row 371, snorkel in row 655, killer\_whale in column 337, mousetrap in row 889
  - ii. Tidiness issues:
    - T1- img\_num column contain the data as the same for the jpg\_url column.
    - T2- the first column to be removed as it the same as index.
- b. twitter\_archive:
  - i. Quality issues:
    - Q4- For the data in column in\_reply\_to\_status\_id: wrong data type (to be object insteade of float).
    - Q5- For the data in column in\_reply\_to\_user\_id: wrong data type (to be object insteade of float).
    - Q6- For the data in column retweeted\_status\_timestamp: wrong data type (to be date insteade of object).
    - Q7- For the data in column in\_timestamp: wrong data type (to be date insteade of object).
    - Q8- For the data in column retweeted\_status\_id: wrong data type (to be object insteade of float).
    - Q9- For the data in column retweeted\_status\_user\_id: wrong data type (to be object insteade of float).
    - Q10- For the data in column retweeted\_status\_timestamp: wrong data type (to be date insteade of object).

Q11- There are many empty arrays in the data frame.

ii. Tidiness issues:

T3- the first column to be removed as it the same as index.

T4- In the columns doggo, floofer, pupper, puppo: need to be melted in one column.

c. twettr api:

i. Quality issues:

Q12- Column Header is undescriptive: extended\_entities, sourc, lang

Q13- There are many missing data.

Q14- For the data in column created\_at: wrong data type (to be date instead of object).

Q15- For the data in column in\_reply\_to\_status\_id: wrong data type (to be object instead of float64).

Q16- For the data in column in\_reply\_to\_status\_id\_str: wrong data type (to be object instead of float64).

Q17- For the data in column in\_reply\_to\_status\_id\_str: wrong data type (to be object instead of float64).

ii. Tidiness issues:

T5- Columns id, id\_str contain the same data.

T6- the first column to be removed as it the same as index.

d. All the 3 data frame:

T7- All the 3 data frame can be merged in one dataframe, using the tweet id as index, and the following columns are replicated between them: in\_reply\_to\_status\_id, in\_reply\_to\_user\_id, source (using duplicated function for all the columns).

3- Cleaning Data:

a. Clean clean twettr api:

i. Remove the unnecessary columns (using drop function), rename the un-descriptive columns (using the rename function).

b. Clean clean twitter archive:

i. Remove the unnecessary columns (using drop function), rename the un-descriptive columns (using the rename function).

ii. Melt the columns doggo, floofer, pupper, puppo to dog\_type new column (using the melt function).

iii. Remove the duplicates of the data frame after melting (using drop\_duplicates function).

c. Clean clean image predictions:

i. Remove the unnecessary columns (using drop function), rename the un-descriptive columns (using the rename function).

d. Merge the new 3 dataframes into one:

i. Merge the dataframes (using merge function).

- ii. Change the None values into NaN values (using replace function).
- iii. Save the new dataframe into csv file (using to\_csv function).

#### References:

- 1- Drop function: <https://cmdlinetips.com/2018/04/how-to-drop-one-or-more-columns-in-pandas-dataframe/>
- 2- Merge function: [https://www.datacamp.com/community/tutorials/python-rename-column?utm\\_source=adwords\\_ppc&utm\\_campaignid=1455363063&utm\\_adgroupid=65083631748&utm\\_device=c&utm\\_keyword=&utm\\_matchtype=b&utm\\_network=g&utm\\_adposition=&utm\\_creative=332602034361&utm\\_targetid=dsa-429603003980&utm\\_loc\\_interest\\_ms=&utm\\_loc\\_physical\\_ms=1005386&gclid=CjwKCAiA17P9BRB2EiwAMvwNyH3bB3luhVIX\\_22iClfPZwCoJjtqPUT2w-3lNkwFKdfitlMfwL1ABoCCRgQAvD\\_BwE](https://www.datacamp.com/community/tutorials/python-rename-column?utm_source=adwords_ppc&utm_campaignid=1455363063&utm_adgroupid=65083631748&utm_device=c&utm_keyword=&utm_matchtype=b&utm_network=g&utm_adposition=&utm_creative=332602034361&utm_targetid=dsa-429603003980&utm_loc_interest_ms=&utm_loc_physical_ms=1005386&gclid=CjwKCAiA17P9BRB2EiwAMvwNyH3bB3luhVIX_22iClfPZwCoJjtqPUT2w-3lNkwFKdfitlMfwL1ABoCCRgQAvD_BwE)
- 3- Rename function: [https://www.datacamp.com/community/tutorials/python-rename-column?utm\\_source=adwords\\_ppc&utm\\_campaignid=1455363063&utm\\_adgroupid=65083631748&utm\\_device=c&utm\\_keyword=&utm\\_matchtype=b&utm\\_network=g&utm\\_adposition=&utm\\_creative=332602034361&utm\\_targetid=dsa-429603003980&utm\\_loc\\_interest\\_ms=&utm\\_loc\\_physical\\_ms=1005386&gclid=CjwKCAiA17P9BRB2EiwAMvwNyH3bB3luhVIX\\_22iClfPZwCoJjtqPUT2w-3lNkwFKdfitlMfwL1ABoCCRgQAvD\\_BwE](https://www.datacamp.com/community/tutorials/python-rename-column?utm_source=adwords_ppc&utm_campaignid=1455363063&utm_adgroupid=65083631748&utm_device=c&utm_keyword=&utm_matchtype=b&utm_network=g&utm_adposition=&utm_creative=332602034361&utm_targetid=dsa-429603003980&utm_loc_interest_ms=&utm_loc_physical_ms=1005386&gclid=CjwKCAiA17P9BRB2EiwAMvwNyH3bB3luhVIX_22iClfPZwCoJjtqPUT2w-3lNkwFKdfitlMfwL1ABoCCRgQAvD_BwE)
- 4- Replace value in columns: <https://www.codegrepper.com/code-examples/delphi/how+to+replace+values+with+nan+in+pandas>