

# Data Wrangle Project:

## weratedogs twitter archive

### Data wrangling, which consists of:

1. Gathering data
2. Assessing data
3. Cleaning data

## 1. Gathering Data

**we gather data from 3 different resources as follow**

- Download file manually
- Download file from the internet
- Access twitter API

### - Download file manually

- Download file (twitter\_archive\_enhanced.csv) and upload it in jupyter enviroment folder.
- Load it in pandas dataframe using pandas read\_csv function then test it using head function

### - Download file from the internet

- Use requests package to get the file from the url provided thensave it.
- Load it in pandas dataframe using pandas read\_csv function then test it using head function

### - Query Data from twitter API

- Create Twitter developer account and get the tokens and credentials needed to access the API.
- Connect to twitter API using tweepy package.
- Get the tweets that has tweet\_id in the twitter\_archive\_enhanced file and save it in tweet\_json.txt each tweet in new line.
- read the tweet\_json.txt and load tweets in dataframe, we grap only the interesting data which are 'tweet\_id', 'retweet\_count', 'favorite\_count', 'followers\_count', 'source'.
- we save each tweet data in dictionary then append the dictionary to list then convert the list to the dataframe

## 2. Assess Data:

### Assess Data visually and programmatically

1. Assess Data visually by open files using excel program and show them also in the jupyter notebook.

## 2. Assess Data programmatically using functions like:

info() to get main informations about datatype and number of null values in each column, head() to get a look at the first rows, sample() to get random rows, shape to get dimension of the dataframe how many rows and columns and other known functions that can be applied to get closer look for the data, you can find all the functions used in wrangle\_act.ipynb

From the visual and programmatic assessment, I found the following issues:

### **Quality Issues:**

1. should remove tweets in twitter\_archive\_df and not in twitter\_api\_df
  - twitter\_archive\_df
2. tweet\_id of type int64 to be easily used
3. timestamp of type object
4. Many None value in the columns name, doggo, floofer pupper and puppo column
5. Inaccurate names in name column like 'a' and 'an'
6. Inaccurate names 'O' in name column
7. 59 null values for the expanded\_urls column
8. Data must contain original tweets no retweets and reply (no tweets has retweeted\_status\_id, in\_reply\_to\_status\_id, in\_reply\_to\_user\_id)
9. Some Columns not needed in the analysis like "in\_reply\_to\_user\_id", "retweeted\_status\_id", "retweeted\_status\_user\_id", "retweeted\_status\_timestamp"
10. numerator column of type int64 and has inaccurate data
11. dominator column has inaccurate data
12. Source column has hard to read values
  - image\_predictions\_df
13. tweet\_id of type int64
14. it has 2075 records which need to be reflected in the other df
15. Inconsistent lower case for some of the predicted bread in the predicted column (p1,p2,p3)
  - twitter\_api\_df
16. tweet\_id of type int64

### **Tidiness Issues:**

1. In twitter\_archive\_df there are four columns (doggo, floofer pupper and puppo) which are values related to one variable
2. In image\_predictions\_df the column names (p1, p1\_conf, p1\_dog, p2, p2\_conf, p2\_dog, p3,

p3\_conf, p3\_dog) not descriptive and need to be merged to three column because each 3 is values to one variable

3. The two dataframes twitter\_archive\_df and twitter\_api\_df should merged in one dataframe because it is related to one observation

### 3. Clean Data

**Clean Data by using the points found in assess data phase using define, code, test procedure.**

**You can find the detailed code and description in wrangle\_act.ipynb file**

→ Before start clean, make clean copy for the three dataframe

#### Quality Clean:

##### First for twitter\_archive\_df

##### 1. should remove tweets in twitter\_archive\_df and not in twitter\_api\_df

- Drop records in twitter\_archive\_df and not in twitter\_api\_df because it is missing tweets which can be deleted.
- use isin and drop functions

##### 2. tweet\_id of type int64

- convert tweet\_id column type from int64 to object -string-
- use .astype(str) function

##### 3. timestamp of type object

- convert timestamp column datatype from object to datetime
- use pd.to\_datetime() function

##### 4. Many None value in the columns name, doggo, floofer pupper and puppo column

- convert "None value for the column name to np.nan
- convert "None" value for doggo, floofer pupper and puppo to "" to be able to merge in other point in tidiness
- use .replace function

##### 5. Inaccurate names in name column like 'a' and 'an'

- Convert the in accurate value name 'a' and 'an' by an accurate name extracted from the text column
- apply reg expression to text column to extract the right names

##### 6. Inaccurate names 'O' in name column

- Correct the value 'O' in name column
- use .replace function

### **7. 59 null values for the expanded\_urls column**

- Drop rows that has null value on the column expanded\_url
- use the drop and notna functions

### **8. Data must contain original tweets no retweets and reply (no tweets has retweeted\_status\_id, in\_reply\_to\_status\_id, in\_reply\_to\_user\_id)**

- Drop rows that has value in retweeted\_status\_id , in\_reply\_to\_user\_id and in\_reply\_to\_status\_id
- use notnull function and invert sign ~

### **9. Some Columns not needed in the analysis like "in\_reply\_to\_user\_id", "retweeted\_status\_id", "retweeted\_status\_user\_id", "retweeted\_status\_timestamp"**

- Drop column that don't have any value no after deleted retweet and reply which are ["in\_reply\_to\_status\_id", "in\_reply\_to\_user\_id", "retweeted\_status\_id", "retweeted\_status\_user\_id", "retweeted\_status\_timestamp"]
- use drop function

### **10. Numerator column of type int64**

- convert numerator type to float and re-extract it from the text value as there is some inaccurate data
- use astype('float') and extract() functions , regular expression

### **11. Denominator column has inaccurate data**

- Correct rating\_denominator values (denominator should equal 10)
- Re-calculate numerator column for rows that hav Denominator value not equal 10

### **12. Source column has hard to read values**

- Drop source column as we already have same column in twitter\_api\_df which will merge later to this df
- use drop function

## **Second for image\_predictions\_df**

### **13. tweet\_id of type int64**

- convert tweet\_id column type from int64 to object -string-
- use astype(str)

### **14. it has 2075 records which need to be reflected in the other df**

- Drop rows from twitter\_archive\_df with tweet\_id that not in image\_predictions\_df and rows from image\_predictions\_df that not in twitter\_archive\_df
- use drop and isin functions and invert sign ~

### **15. Inconsistent lower case for some of the predicted bread in the predicted column (p1,p2,p3)**

- Change values for p1 , p2, p3 column to be capitalized
- use .str.capitalize() function

## Third For twitter\_api\_df

### 16. tweet\_id of type int64

- convert tweet\_id column type from int64 to object -string-
- use astype(str) function

## Tidiness Clean:

### 1. In twitter\_archive\_df there are four columns (doggo, floofer pupper and puppo) which are values related to one variable

- Combine the four columns (doggo, floofer pupper and puppo) in one column and named it dog\_bread
- use + sign to combine and then .replace to replace empty value with np.nan then drop the four columns
- replace value in dog\_bread column that have multiple bread with mixed bread value

### 2. In image\_predictions\_df the column names (p1, p1\_conf, p1\_dog, p2, p2\_conf, p2\_dog, p3, p3\_conf, p3\_dog) not descriptive and need to be merged to three column because each 3 is values to one variable

- Rename the 9 column ('tweet\_id', 'jpg\_url', 'img\_num', 'propability\_1', 'confidence\_1', 'dog\_1', 'propability\_2', 'confidence\_2', 'dog\_2', 'propability\_3', 'confidence\_3', 'dog\_3') to be more descriptive and merge them into 3 columns 'propability', 'confidence', 'dog'
- use wide\_to\_long function

### 3. The two dataframes twitter\_archive\_df and twitter\_api\_df should merged in one data frame because it is related to one observation

- Merge twitter\_archive\_df\_clean and twitter\_api\_df\_clean into
- use merge function with inner and on tweet\_id

## Storing wrangled data

1. Store twitter\_archive\_master\_df\_clean into twitter\_archive\_master.csv
2. Store image\_predictions\_df\_clean into image\_predictions\_master.csv