

Wrangle Report with WeRateDog Data

July 29, 2019

1 Wrangle Report with WeRateDog Data

1.1 Introduction of Data Set :

There are three data sets in total for this project : twitter-archive-enhanced.csv,image-predictions.tsv , tweet-json.txt

1.2 The Goal :

The goal of the data wrangle process is to generate a clean data set called twitter_archive_master.csv for data visualization and analysis later.

2 Data Issues: I found 12 quality issues and 4 tidiness issues

**** Tidiness Issues (Content Issues) ****

- 1) numerator_rating and denominator should me merge in one coloumn instead of two col-
umn.
- 2) Create dog stage variable and remove individual dog stage columns.
- 3) There is a incomplete data in some file so i perform merge operation to combine all the file
for analysis.
- 4) There is inappropriate column name for json file.

**** Quality Issues ****

- 5) In name column there are multiple expression which are wrongly name convert these
wrongly name to none.
- 6) Split date and time in two seperate column instead of one column.
- 7) change 'tweet_id','in_reply_to_status_id', 'in_reply_to_user_id', 'retweeted_status_id',
'retweeted_status_user_id' from float to integer
- 8) Drop the rows when the values of retweet_ids are not none.

- 9) There are missing values in 'in_reply_to_status_id', 'in_reply_to_user_id', 'retweeted_status_id', 'retweeted_status_user_id', 'retweeted_status_timestamp'. The columns can be dropped since they are not relevant for analysis purpose.
- 10) Change missing values in 'name' from 'None' to NaN
- 11) Remove tweet without rating
- 12) Remove extra character after '&' in text column of text file.
- 13) Remove "-" in P1, P2, P3 Using replace function
- 14) Change all string in p1, p2, p3 in lower case.

Clean Data: Tidiness issues first , then quality issues

- 1) numerator_rating and denominator should me merge in one coloumn instead of two col-umn.

** First I convert the numerator and denominator column in float and then i change the the value with actual decimal number by using loc method to access particular row. I simply add the numerator and denominator rating to one column called rating by using astype function and passing str as parameter.

- 2) Create dog stage variable and remove individual dog stage columns.

** I extract (puppo | pupper | floofer | doggo) in one column by using replace method and for multiple value i used multiple variable name. Because few variable is present in two many col-umn.

- 3) There is a incomplete data in some file so i perform merge operation to combine all the file for analysis.

** I simply performe merge operation by calling merge function on the given three file and perform join on 'twitter_id'

- 4) There is inappropriate column name for json file.

** I use column renaming method and rename column twitter_id in place of id.

- 5) In name column there are multiple expression which are wrongly name, convert these wrongly name to none

** I use replace function to convert wrong name to none. First I make list of wrong name after that i perform replace optertaion on that.

- 6) Split date and time in two separte column instead of one column.

** I use str.split function on timestamp column to split in two different column called date and time column.

- 7) change 'tweet_id', 'in_reply_to_status_id', 'in_reply_to_user_id', 'retweeted_status_id', 'retweeted_status_user_id' from float to integer

** I use astype function on the above column which have inappropriate data types.

8) Drop the rows when the values of retweet_ids are not none.

** By using drop function i drop retweet_ids whose value is equal to none.

9) There are missing values in 'in_reply_to_status_id', 'in_reply_to_user_id', 'retweeted_status_id', 'retweeted_status_user_id', 'retweeted_status_timestamp'. The columns can be dropped since they are not relevant for analysis purpose.

** The above column are not useful for further analysis because it have plenty of missing value. So i drop those column by using drop function.

10) Change missing values in 'name' from 'None' to NaN

** Again using replace method replace None to NaN.

11) Remove tweet without rating

** I remove the tweet without rating by using not equal to operator this operator matches the tweet without rating.

12) Remove extra character after '&' in text column of text file.

** By using replace I replace '&' with '&'

13) Remove "-" in P1, P2, P3 Using replace function

** By using replace function. I replace "-" with blank space.

14) Change all string in p1, p2, p3 in lower case.

** By using lower function i changes all the text in p1,p2,p3 column in to the lower text.

3 Conclusion:

In summary, this project was my biggest challenge to date, specifically using the Twitter API to gather the JSON data. Overall, this project was completed successfully and I'm extremely pleased with the new skills I acquired.