

INTRODUCTION

This project is the second project of the Udacity Data Analyst Nanodegree program and its primary focus is on data wrangling. The data for the project is gotten from twitter WeRateDogs account. Python is the tool used for this project with some other libraries like pandas, matplotlib, request and seaborn. The gathering and cleaning process is documented in a jupyter notebook.

The Data Gathering

Three different data which was later merged into one was used for the project. The first data is the WeRateDogs twitter-archive-enhanced.csv which was provided by Udacity in csv format and contains the following columns: tweet_id, in_reply_to_status_id, in_reply_to_user_id, timestamp, source, text, retweeted_status_id, retweeted_status_user_id, retweeted_status_timestamp, expanded_urls, rating_numerator, rating_denominator, name

Although not all the columns are valid as some have high percentage of null values, some columns need to be tidied, some are not important.

The second data is the image-predictions.tsv, the url link was provided by Udacity was programmatically downloaded using request library, it contains the following columns: tweet_id, jpg_url, img_num, p1, p1_conf, p1_dog, p2, p2_conf, p2_dog, p3, p3_conf, p3_dog

The third data is from the twitter API but since I couldn't get developer access on twitter I use the json file provided by Udacity which I read line by line into pandas dataframe, the lines I read are: id, retweet_count, favorite_count because that's what I require

Data assessing

The data was assessed visually by scrolling through each of it in excel and programmatically using pandas and methods such as .head(), .info() etc from assessing the data, some issues were discovered from the dataset, the issues were separated into tidiness issues and quality issues

Data Cleaning

The first and important step in data cleaning is making a copy of the original dataset, it's the copied dataset that the cleaning is performed on so that the original data remain intact

The datasets were cleaned in three steps which are the

define; here the solution to the problem is defined

code; the code to effect the cleaning is done at this stage

test; the issue is re-assessed to see if it's still there

one of the cleaning challenge I had is combining the pupper, floofer, pupperer and puppo column into one

Conclusion

Data wrangling is an important and useful skill in data analysis because it can be use to extract data from different source and it better for big data compared to excel

In []: