

## Summary of wrangling process

The process to do data wrangling can be separated into 3 main steps starting from gathering data, then assessing quality of it, and cleaning it using programmatic approach.

### First gathering data

The data has been received from three different sources: twitter, provided csv, and download direct from Udacity server.

I have started the process by import necessary python library consisted of tweepy, pandas, json, requests, numpy, and % matplotlib inline

Then

- Gather data from twitter using code written. Save it in 'twitter\_json.txt', read tweet\_json.txt file line by line into dataframe called it df.
- Read provided csv ('twitter-archive-enhanced.csv') store it in dataframe called df1.
- Gather dog breed classification data from udacity server stored it in dataframe called img\_pred

After that I take a quick overview of what information provided inside using both python and excel.

### Second Accessing data

In this process, the data has been checked step by step. Few main things I am looking for was whether it has missing value or not, how is the data quality, and how is data tidiness. In case of issues detected, I will document and classify clearly into either quality issues or tidiness issues.

Twitter data (df): 6 quality issues and 2 tidiness issues have been identified.

#### Quality

- 3 columns (contributors, coordinates, geo) have no data at all.
- Missing values columns are extended\_entities, in\_reply\_to\_screen\_name, in\_reply\_to\_status\_id, in\_reply\_to\_status\_id\_str, in\_reply\_to\_user\_id, in\_reply\_to\_user\_id\_str, place, possibly\_sensitive, possibly\_sensitive\_appealable, quoted\_status, quoted\_status\_id, quoted\_status\_id\_str, retweeted\_status)
- 15 columns have data completed for all rows (created\_at, display\_text\_range, entities, favorite\_count, favorited, full\_text, id, id\_str, is\_quote\_status, lang, retweet\_count, retweeted, source, truncated, user)
- Column 'created\_at' has incorrect datatype. It should be date time
- Column 'Source' has information of source (twitter for iphone, twitter web client) together with non related information
- Column 'display\_text\_range' has stored as object

#### Tidiness

- Not need:  
'contributors','coordinates','entities','extended\_entities','favorited','full\_text','user','id\_str',  
'truncated','entities','extended\_entities','in\_reply\_to\_status\_id','in\_reply\_to\_user\_id',  
'in\_reply\_to\_status\_id\_str','in\_reply\_to\_user\_id\_str','in\_reply\_to\_screen\_name',  
'geo','coordinates','place','contributors','is\_quote\_status',  
'possibly\_sensitive','possibly\_sensitive\_appealable','quoted\_status','quote\_status\_id','quote  
d\_status\_id\_str','retweeted','retweeted\_status'
- Need to combine data from 3 table using 'id' as key to create final table

Twitter-archive-enhanced.csv (stored as df1). Identified 2 qualities issue and 2 tidiness issues.

### Quality

- Column 'rating\_numerator' has wrong information (not equal to 10)
- Column 'timestamp' has +0000 inside

### Tidiness

- Not need column: 'in\_reply\_to\_status\_id', 'in\_reply\_to\_user\_id',  
'timestamp','retweeted\_status\_timestamp','expanded\_urls'

Dog breed classification (img\_pred): 2 tidiness issues have been identified.

### Tidiness

- Not need column: jpg\_url, img\_num
- Interested only p1 prediction. Remove p2 and p3

## Final steps clean the data

The cleaning process starting from define what to clean, write a code, and then test the code to make sure everything work as expected.

Twitter file (df): remove not necessary columns, change column 'created\_at' to date time, remove unnecessary text to collect twitter source from column 'source', and change text lenght to integer from column 'display\_text\_range'

Twitter-archive-enhanced.csv (stored as df1): remove not necessary column, change column 'rating\_denominator' to 10 when figure is not equal to 10

Dog breed classification (img\_pred): remove not necessary column

## Combine table is the last step

Finally, I have combined 3 files together to create final table in order to further use for analysis. First combine df with df1 in order to get rating, numerator, denominator, name, and type of dogs. After that, combine it with img\_pred to get dog breed information.

After finished combining data. The file 'twitter\_archive\_master.csv' has been exported for further reference.