

Data Wrangling Report

The three main steps in the data wrangling process are gathering, assessing, and cleaning (as shown in Fig.1).

WeRateDogs is a project where users tweet about their pets. Text (comments), ratings (such as 11/10, 12/10, etc.), the source of the tweet (such as iPhone), and the URLs are all wrangled. Sadly, this project requires data from a variety of sources which may come in different formats via an algorithm, site scraping, and other data documentation. Let's gather data! first, before anything else!

First: The following were the data sources used in this analysis:

1. Twitter archive enhanced: Udacity contributed this information for the purpose of data wrangling after receiving it from the WeRateDogs community. As of August 1, 2017, this data, which is in comma separated format (.csv), had more than 5000 tweets.
2. Image prediction: This information was generated by a neural network that anticipated dog photographs and classified them according to the different dog breeds. Its contents contain the tweet ID, picture URL, and image number that matched the most accurate forecast, and each tweet's JSON data is stored in a separate file called tweet json.txt.
3. Web scraping: This additional data includes a list of the tweet ids as found in the Twitter archive enhanced but with the number of favourites and retweets for each tweet. Using Python's Tweepy package, one must gather this data using the Twitter API to obtain the JSON data for each tweet, which must then be saved in a tweet json.txt file.

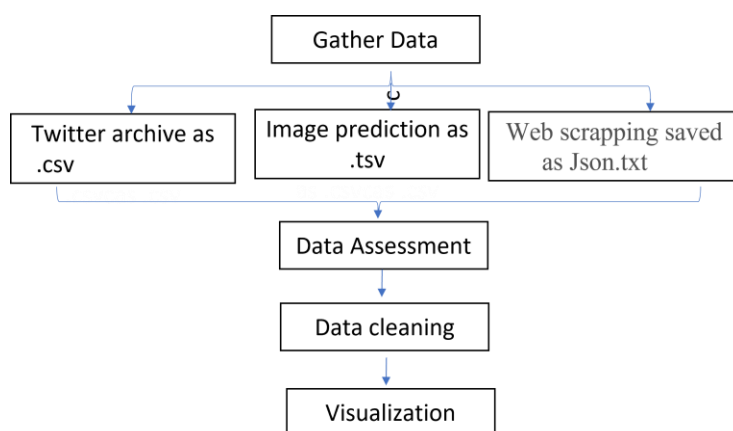


Fig 1: Process flow chart

Assessment and Cleaning

Data type	Type of assessment	Data Assessment	Data cleaning
Twitter archive enhanced as .csv	Visual	'NaN' and 'None' are the same thing +0000 appeared in every single cell The column 'source' has identical 'href' https://twitter.com/dog_rates/status is repeated in the column 'expanded URL' inconsistency issues, a mixture of capital and small letters e.g., row 42	Convert all to NaN Change datatype into datetime that removes +0000

		The column expanded URL has both text and figures	Be consistent with small letters
	Programmatic	In the text column with the index number 870 and tweet_id '761672994376806400' made no sense The names of dog such as 'none', 'a', 'an' occurred 745, 55, 7, and 8 times, respectively. The text column has URLs and hashtag	This may be converted to NaN Separate URLs from text column
Image prediction as. tsv	Visual	False images. Columns p1, p2, and p3 have other animals and other items that are not God.	Filter off the false image the columns
	Programmatic	There are about 281 missing tweet_id because TAE has 2356 while IP has 2075	
Web scrapping saved as Json.txt.csv	Visual	Column label for the tweet_id is represented with 'id' which is not consistent with TAE and IP	
	Programmatic	there are 2 missing rows when compared to the TAE with 2356 rows, RF_tweet has 2354 rows.	
All data frames		Data type of tweet_id	

Tidy issues

1. The timestamp in the twitter archive enhanced dataframe has date and time (plus +0000). The solution is to change the datatype from object into datetime.
2. Tweet_id data type is integer, and should change to string (object)

Result

Storing Data

Save gathered, assessed, and cleaned master dataset to a CSV file named "twitter_archive_master.csv".

```
In [242]: master_data.to_csv("twitter_archive_master.csv", index=False)
```

```
In [243]: #load read_csv("twitter_archive_master.csv")
```

```
In [247]: #F
```

```
Out[247]:
```

	tweet_id	source	text	rating_numerator	rating_denominator	name	year	month	day
0	882422320410584	http://twitter.com/obowabagbaf...	This is Cactus. She is a copper dog. She's...	14	10	Cactus	2017	7	28
1	888888888888888	http://twitter.com/obowabagbaf...	Here's a puppie that adores me and my family :)	12	10	Nah	2017	7	28
2	888888888888888	http://twitter.com/obowabagbaf...	This is Sweet. She's adoring the family dog.	12	10	Sweet	2017	7	24
3	888888888888888	http://twitter.com/obowabagbaf...	This is Russian. Another doggy from another...	12	10	Russian	2017	7	18
4	888888888888888	http://twitter.com/obowabagbaf...	Meet Trup. We don't have any pictures of her.	12	10	Trup	2017	7	8

```
In [94]: master_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 365 entries, 0 to 364
Data columns (total 21 columns):
tweet_id      365 non-null object
source        365 non-null object
text          365 non-null object
rating_numerator  365 non-null int64
rating_denominator  365 non-null int64
name          234 non-null object
year          365 non-null int64
month         365 non-null int64
day           365 non-null int64
hour          365 non-null int64
minute        365 non-null int64
dog_stage     365 non-null object
img_num       307 non-null float64
p1            307 non-null object
p1_dog        307 non-null object
p2            307 non-null object
p2_dog        307 non-null object
p3            307 non-null object
p3_dog        307 non-null object
retweet_count  365 non-null int64
favorite_count 365 non-null int64
dtypes: float64(1), int64(9), object(11)
memory usage: 62.7+ KB
```