The dataset that you will be wrangling (and analyzing and visualizing) is the tweet archive of Twitter user @dog_rates, also known as WeRateDogs. WeRateDogs is a Twitter account that rates people's dogs with a humorous comment about the dog. These ratings always have a denominator of 10. The numerators, though? Always greater than 10. 11/10, 12/10, 13/10, etc. Why? Because "they're good dogs Brent." WeRateDogs has over 4 million followers and has received international media coverage.

the data wrangling steps

1: - gathering

First, I gathered the twitter archive data as normal way form pandas.read_csv

Then I gathered the image prediction programmatically and finally

The last one I gathered it by reading the text file line by line and extract the tree columns we want their information

2-

Assessing

The first data frame...

Out[7]:

| | tweet_id | in_reply_to_status_id | in_reply_to_user_id | timestamp | source | text | retweeted_status_id | retwe |
|---|---|---|---|---|---|---|---|---|
| 0 | 892420643555336193 | NaN | NaN | 2017-08-01 16:23:56 +0000 | <a href="http://twitter.com/download/iphone" r... | This is Phineas. He's a mystical boy. Only eve... | NaN | |
| 1 | 892177421306343426 | NaN | NaN | 2017-08-01 00:17:27 +0000 | <a href="http://twitter.com/download/iphone" r... | This is Tilly. She's just checking pup on you.... | NaN | |
| 2 | 891815181378084864 | NaN | NaN | 2017-07-31 00:18:03 +0000 | <a href="http://twitter.com/download/iphone" r... | This is Archie. He is a rare Norwegian Pouncin... | NaN | |
| 3 | 891689557279858688 | NaN | NaN | 2017-07-30 15:58:51 +0000 | <a href="http://twitter.com/download/iphone" r... | This is Darla. She commenced a snooze mid meal... | NaN | |
| 4 | 891327558926688256 | NaN | NaN | 2017-07-29 16:00:24 +0000 | <a href="http://twitter.com/download/iphone" r... | This is Franklin. He would like you to stop ca... | NaN | |

Its quality issues were

-rating_numerator sometimes is lower than 10

-tweet_id must be string not int

-timestamp must be time not string

-unnecessary columnsmissing values in expanded_urls

-unnecessary rating_denominator column

-making name is lower for all

-'None' values in name"None" values in dog_stage


And its tidiness issue was

-dogs kind 4 variables in 4 columns and it should be in one column


The second data frame (image prediction)

Out[14]:

| | tweet_id | jpg_url | img_num | p1 | p1_conf | p1_dog | p2 | p2_co |
|---|---|---|---|---|---|---|---|---|
| 0 | 666020888022790149 | https://pbs.twimg.com/media/CT4udn0WwAA0aMy.jpg | 1 | Welsh_springer_spaniel | 0.465074 | True | collie | 0.15666 |
| 1 | 666029285002620928 | https://pbs.twimg.com/media/CT42GRgUYAA5iDo.jpg | 1 | redbone | 0.506826 | True | miniature_pinscher | 0.07419 |
| 2 | 666033412701032449 | https://pbs.twimg.com/media/CT4521TWwAEvMyu.jpg | 1 | German_shepherd | 0.596461 | True | malinois | 0.13858 |
| 3 | 666044226329800704 | https://pbs.twimg.com/media/CT5Dr8HUEAA-lEu.jpg | 1 | Rhodesian_ridgeback | 0.408143 | True | redbone | 0.36068 |
| 4 | 666049248165822465 | https://pbs.twimg.com/media/CT5IQmsXIAAKY4A.jpg | 1 | miniature_pinscher | 0.560311 | True | Rottweiler | 0.24368 |
| ... | ... | ... | ... | ... | ... | ... | ... | |
| 2070 | 891327558926688256 | https://pbs.twimg.com/media/DF6hr6BUMAAzZgT.jpg | 2 | basset | 0.555712 | True | English_springer | 0.22577 |
| 2071 | 891689557279858688 | https://pbs.twimg.com/media/DF_q7IAWsAEuuN8.jpg | 1 | paper_towel | 0.170278 | False | Labrador_retriever | 0.16808 |
| 2072 | 891815181378084864 | https://pbs.twimg.com/media/DGBdLU1WsAANxJ9.jpg | 1 | Chihuahua | 0.716012 | True | malamute | 0.07825 |
| 2073 | 892177421306343426 | https://pbs.twimg.com/media/DGGmoV4XsAAUL6n.jpg | 1 | Chihuahua | 0.323581 | True | Pekinese | 0.09064 |
| 2074 | 892420643555336193 | https://pbs.twimg.com/media/DGKD1-bXoAAIAUK.jpg | 1 | orange | 0.097049 | False | bagel | 0.08585 |

2075 rows × 12 columns


Its quality issue was

-tweet_id must be string not int


The third data frame (tweet data)

Out[60]:

| | id | retweet_count | favorite_count |
|---|---|---|---|
| **0** | 892420643555336193 | 8853 | 39467 |
| **1** | 892177421306343426 | 6514 | 33819 |
| **2** | 891815181378084864 | 4328 | 25461 |
| **3** | 891689557279858688 | 8964 | 42908 |
| **4** | 891327558926688256 | 9774 | 41048 |
| **...** | ... | ... | ... |
| **2349** | 666049248165822465 | 41 | 111 |
| **2350** | 666044226329800704 | 147 | 311 |
| **2351** | 666033412701032449 | 47 | 128 |
| **2352** | 666029285002620928 | 48 | 132 |
| **2353** | 666020888022790149 | 532 | 2535 |

2354 rows × 3 columns

Its quality issue was

-id must be string not int

And the all-data tidiness issue:

-all data is related but separated to 3 tables

Cleaning

I solved the previous issues using pandas and NumPy

And I did the whole thing on a copy data frames so if I did something wrong does not affect the original file

And then I merged the three data frame in one clean data frame

And finally removed the missing values caused by unmatched ids