

This reports briefly explains and describes my wrangling effort in the project

The dataset used in this project is the tweet archive of a popular user in twitter known as WeRateDogs which rates different types of people's dogs with a humorous comments about the dog. These ratings almost always have a denominator of 10. The numerators, though? Almost 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 wrangling process done is divided into 3 three steps:

1-Gathering data

2- Assessing data

3-Cleaning data

Each step will be briefly explained below.

### First: Gathering data

The data is collected from three different resources related to the same topic

#### 1- Enhanced Twitter Archive

This contains basic tweet data for all 5000+ of WeRateDogs tweets (the data is filtered to be only 2356 row) this data set contains 17 columns which are **tweet\_id** , **timestamp**, **source**, **text**, **rating\_numerator**, **rating\_denominator**, **name** , **dog type** etc.....

| text  | rating_numerator | rating_denominator | name     | doggo | floofer | pupper | puppo |
|---|------------------|--------------------|----------|-------|---------|--------|-------|
| This is Phineas. He's a mystical boy. Only ever appears in the hole of a donut. 13/10 <a href="https://t.co/MgUWQ76dJU">https://t.co/MgUWQ76dJU</a>                                     | 13               | 10                 | Phineas  | None  | None    | None   | None  |
| This is Tilly. She's just checking pup on you. Hopes you're doing ok. If not, she's available for pats, snugs, boops, the whole bit. 13/10  | 13               | 10                 | Tilly    | None  | None    | None   | None  |
| This is Archie. He is a rare Norwegian Pouncing Corgo. Lives in the tall grass. You never know when one may strike. 12/10 <a href="https://t.co/ID36da7qLQ">https://t.co/ID36da7qLQ</a> | 12               | 10                 | Archie   | None  | None    | None   | None  |
| This is Darla. She commenced a snooze mid meal. 13/10 happens to the best of us <a href="https://t.co/ID36da7qLQ">https://t.co/ID36da7qLQ</a>   | 13               | 10                 | Darla    | None  | None    | None   | None  |
| This is Franklin. He would like you to stop calling him "cute." He is a very fierce shark and should be respected as such. 12/10 #BarkWeek  | 12               | 10                 | Franklin | None  | None    | None   | None  |
| Here we have a majestic great white breaching off South Africa's coast. Absolutely h*ckin breathtaking. 13/10 (IG: tucker_mario) #BarkWeek  | 13               | 10                 | None     | None  | None    | None   | None  |
| Meet Jax. He enjoys ice cream so much he gets nervous around it. 13/10 help Jax enjoy more things by clicking below   |                  |                    |          |       |         |        |       |
| <a href="https://t.co/Zr4hWfAs1H">https://t.co/Zr4hWfAs1H</a> <a href="https://t.co/tVJBRMnhxl">https://t.co/tVJBRMnhxl</a>   | 13               | 10                 | Jax      | None  | None    | None   | None  |
| When you watch your owner call another dog a good boy but then they turn back to you and say you're a great boy. 13/10 <a href="https://t.co/h">https://t.co/h</a>                      | 13               | 10                 | None     | None  | None    | None   | None  |
| This is Zoey. She doesn't want to be one of the scary sharks. Just wants to be a snuggly pettable boatpet. 13/10 #BarkWeek <a href="https://t.co/h">https://t.co/h</a>                  | 13               | 10                 | Zoey     | None  | None    | None   | None  |
| This is Cassie. She is a college pup. Studying international doggo communication and stick theory. 14/10 so elegant much sophisticated  | 14               | 10                 | Cassie   | doggo | None    | None   | None  |
| This is Koda. He is a South Australian deckshark. Deceptively deadly. Frighteningly majestic. 13/10 would risk a petting #BarkWeek <a href="https://t.co/h">https://t.co/h</a>          | 13               | 10                 | Koda     | None  | None    | None   | None  |
| This is Bruno. He is a service shark. Only gets out of the water to assist you. 13/10 terrifyingly good boy <a href="https://t.co/u1XPQM29g">https://t.co/u1XPQM29g</a>                 | 13               | 10                 | Bruno    | None  | None    | None   | None  |
| Here's a puppo that seems to be on the fence about something haha no but seriously someone help her. 13/10 <a href="https://t.co/BxvuXk0Uk">https://t.co/BxvuXk0Uk</a>                  | 13               | 10                 | None     | None  | None    | None   | puppo |
| This is Ted. He does his best. Sometimes that's not enough. But it's ok. 12/10 would assist <a href="https://t.co/f8dEDorKSR">https://t.co/f8dEDorKSR</a>                               | 12               | 10                 | Ted      | None  | None    | None   | None  |
| This is Stuart. He's sporting his favorite fanny pack. Secretly filled with bones only. 13/10 puppared puppo #BarkWeek <a href="https://t.co/y70k">https://t.co/y70k</a>                | 13               | 10                 | Stuart   | None  | None    | None   | puppo |

#### 2- Image Predictions File

This file contains every image in the WeRateDogs Twitter archive through a **neural network** that can classify breeds of dogs and the results of this process is a data set contains a lot of image predictions. This data set consists of 2075 row and 12 columns which are : **tweet\_id** , **jpg\_url** , **img\_num** , **p1** (first prediction) , **p1\_conf** , **p1 p1\_dog** , **P2**(second prediction), etc.....

| tweet_id           | jpg_url   | img_num | p1                       | p1_conf  | p1_dog | p2                 | p2_conf    | p2_dog | p3                          | p3_conf    | p3_dog |
|--------------------|---|---------|--------------------------|----------|--------|--------------------|------------|--------|-----------------------------|------------|--------|
| 892177421306343426 | <a href="https://pbs.twimg.com">https://pbs.twimg.com</a> | 1       | Chihuahua                | 0.323581 | TRUE   | Pekinese           | 0.0906465  | TRUE   | papillon                    | 0.0689569  | TRUE   |
| 891815181378084864 | <a href="https://pbs.twimg.com">https://pbs.twimg.com</a> | 1       | Chihuahua                | 0.716012 | TRUE   | malamute           | 0.078253   | TRUE   | kelpie                      | 0.0313789  | TRUE   |
| 891689557279858688 | <a href="https://pbs.twimg.com">https://pbs.twimg.com</a> | 1       | paper_towel              | 0.170278 | FALSE  | Labrador_retriever | 0.168086   | TRUE   | spatula                     | 0.0408359  | FALSE  |
| 891327558926688256 | <a href="https://pbs.twimg.com">https://pbs.twimg.com</a> | 2       | basset                   | 0.555712 | TRUE   | English_springer   | 0.22577    | TRUE   | German_short-haired_pointer | 0.175219   | TRUE   |
| 891087950875897856 | <a href="https://pbs.twimg.com">https://pbs.twimg.com</a> | 1       | Chesapeake_Bay_retriever | 0.425595 | TRUE   | Irish_terrier      | 0.116317   | TRUE   | Indian_elephant             | 0.0769022  | FALSE  |
| 890971913173991426 | <a href="https://pbs.twimg.com">https://pbs.twimg.com</a> | 1       | Appenzeller              | 0.341703 | TRUE   | Border_collie      | 0.199287   | TRUE   | ice_lolly                   | 0.193548   | FALSE  |
| 890729181411237888 | <a href="https://pbs.twimg.com">https://pbs.twimg.com</a> | 2       | Pomeranian               | 0.566142 | TRUE   | Eskimo_dog         | 0.178406   | TRUE   | Pembroke                    | 0.0765069  | TRUE   |
| 890609185150312448 | <a href="https://pbs.twimg.com">https://pbs.twimg.com</a> | 1       | Irish_terrier            | 0.487574 | TRUE   | Irish_setter       | 0.193054   | TRUE   | Chesapeake_Bay_retriever    | 0.118184   | TRUE   |
| 890240255349198849 | <a href="https://pbs.twimg.com">https://pbs.twimg.com</a> | 1       | Pembroke                 | 0.511319 | TRUE   | Cardigan           | 0.451038   | TRUE   | Chihuahua                   | 0.0292482  | TRUE   |
| 890006608113172480 | <a href="https://pbs.twimg.com">https://pbs.twimg.com</a> | 1       | Samoyed                  | 0.957979 | TRUE   | Pomeranian         | 0.0138835  | TRUE   | chow                        | 0.00816748 | TRUE   |
| 889880896479866881 | <a href="https://pbs.twimg.com">https://pbs.twimg.com</a> | 1       | French_bulldog           | 0.377417 | TRUE   | Labrador_retriever | 0.151317   | TRUE   | muzzle                      | 0.0829811  | FALSE  |
| 889665388333682689 | <a href="https://pbs.twimg.com">https://pbs.twimg.com</a> | 1       | Pembroke                 | 0.966327 | TRUE   | Cardigan           | 0.0273557  | TRUE   | basenji                     | 0.00463323 | TRUE   |
| 889638837579907072 | <a href="https://pbs.twimg.com">https://pbs.twimg.com</a> | 1       | French_bulldog           | 0.99165  | TRUE   | boxer              | 0.00212864 | TRUE   | Staffordshire_bulldog       | 0.00149818 | TRUE   |
| 889531135344209921 | <a href="https://pbs.twimg.com">https://pbs.twimg.com</a> | 1       | golden_retriever         | 0.953442 | TRUE   | Labrador_retriever | 0.0138341  | TRUE   | redbone                     | 0.00795775 | TRUE   |

### 3- Additional Data via the Twitter API

Here we needed to query additional data from Twitter API as this data will help us in the wrangling process and the data set consists of 2354 rows and 3 columns which are **tweet\_id** , **retweet\_count** , **favorite\_count**

## Second:Assessing data

After the gathering process of the data the data need to be assessed and assessing is done visually and programmatically for the quality and tidiness issues

### 1- Quality issues

- There are about 181 retweet in the data frame
- Some ratings denominators not equal 10
- Some ratings numerators are too big
- Timestamp column datatype is an object and it has to be datetime
- Change tweed\_id in all tables into strings
- There are invalid names in name column such as (a , an , The and None)
- missing photos for some id's
- all P names should starts with capital letters
- A lot of of columns contains a big number of null values such as retweeted\_status\_id, retweeted\_status\_user\_id, retweeted\_status\_timestamp so they have to be removed

## 2- Tidness issues

- Merge the three tables to be only one table
- Create one column for all dog types: doggo, floofer, pupper, puppo

### **Third: Cleaning**

The last step in wrangling the data is data cleaning and in this step we have to solve every issue we found in the assessing process to result in high quality and tidy data