# Wrangle Report

This report briefly describes the data wrangling efforts exerted in this project.

In this project we gathered data from a variety of sources and a variety of formats. 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 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 entirety of this project completed inside the Udacity classroom Using Python and its libraries.

## *The Data wrangling, which consists of:*

- Gathering data
- Assessing data
- Cleaning data

We will begin with gathering data

# Gathering data:

A) **_Enhanced Twitter Archive:_** contains data extracted programmatically from tweet data sent by WeRateDogs to Udacity via email exclusively to be used in this project. The data provides the rating, dog name, dog stage and some other related information.

| 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 https://t.co/MgUWQ76dJU | 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 https://t.co | 12 | 10 | Archie | None | None | None | None |
| This is Darla. She commenced a snooze mid meal. 13/10 happens to the best of us https://t.co/tD36da7qLQ | 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 #BarkW | 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_marlo) #Bar | 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 https://t.co/Zr4hWfAs1H https://t.co/tVJBRMnhxl | 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 https://t.co/ | 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 https://t. | 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 sophisticate | 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 ht | 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 https://t.co/u1XPQMl29g | 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 https://t.co/BxvuXk0U0 | 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 https://t.co/f8dEDcrKSR | 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 https://t.co/y70c | 13 | 10 | Stuart | None | None | None | puppo |

The extracted data from each tweet's text.

B) **_Image Predictions File:_** produced by running every image in the WeRateDogs Twitter archive through a neural network that can classify breeds of dogs*. The results: a table full of image predictions (the top three only) alongside each tweet ID, image URL, and the image number that corresponded to the most confident prediction (numbered 1 to 4 since tweets can have up to four images).

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

Tweet image prediction data

C) ***Additional Data via Twitter API :*** obtained by querying Twitter's API then stored in a txt file called tweet-json Gathering this data requires a Twitter developer account.

## *Assessing Data:* after gathering data they were assessed visually and programmatically for quality and tidiness issues

The following findings were concluded:

### A) *Tidiness*

- Dog stage data is separated into 4 columns.
- All data is related but divided into 3 separate data frames.

### B) *Quality*

- There are 181 retweets as indicated by retweeted_status_id

- Names are typed with underscore separators

- Invalid tweet_id data type

- Invalid timestamp data type

- Missing photos for some IDs (2075 rows instead of 2356)

- Some predicted breeds are invalid (not breed of dogs)

- Tweet_id column is named only id

## *Cleaning data*

The previous issues were cleaned in a high quality and tidy master pandas data frame.