2356 non-null object pupper 2356 non-null object puppo dtypes: float64(4), int64(3), object(10) memory usage: 313.0+ KB We found three issues in the information, The data type of datetime is a string. · Some IDs are float. · Missing values. Next, take a look in the contents and we found: source href="http://twitter.com/download.iphone href="http://twitter.com/download/iphone" <a href="http://twitter.com/download/iphone" <a href="http://twitter.com/download/iphone" <a href="http://twitter.com/download/iphone" • Redundant contents in the source column which made them unreadable. The information we need are squared in red, and the sources we sorted out from the dataframe are "Twitter for Iphone", "Make a Scene", "Twitter Web Client", and "TweetDeck". Invalid ratings seen in numerator and denominator: rating_numerator rating_denominator 420 10 666 10 182 10 960 Strange names such as: name such quite quite quite The names and rating numbers came from texts and replies, thus some issues might occur in other columns. Most important of all, we have to make the data as tidy as possible, which means some columns shouldn't be dispersed. doggo floofer pupper puppo None The columns above are all telling the same thing. **Image Prediction** In image prediction, we will gather the data from the Udacity server, and name it as df image predictions. Since all the data we gathered were grouped by tweet ids, the issues weren't as much as the previous dataframe. By looking at the data information, minor issues we found are as follows: <class 'pandas.core.frame.DataFrame'> RangeIndex: 2075 entries, 0 to 2074 Data columns (total 12 columns): tweet_id 2075 non-null int64 2075 non-null object jpg_url img_num 2075 non-null int64 2075 non-null object p1 2075 non-null float64 pl_conf pl_dog 2075 non-null bool p2 2075 non-null object p2_conf 2075 non-null float64 2075 non-null bool p2_dog р3 2075 non-null object 2075 non-null float64 p3_conf 2075 non-null bool p3_dog dtypes: bool(3), float64(3), int64(2), object(4) memory usage: 152.1+ KB p1, p2, p3 in the data corresponded to the top three confident prediction of the dogs, in this case, they should be categories instead of strings. The figure we are going to see next is the first count rows of same tweet pictures (url) with the different ids. Some pictures were seen more than one time, but with different ids, we can know that they are retweets. https://pbs.twimg.com/media/ChK1tdBWwAQ1flD.jpg 2 https://pbs.twimg.com/media/CkjMx99UoAM2B1a.jpg https://pbs.twimg.com/media/Ct2qO5PXEAE6eBO.jpg https://pbs.twimg.com/media/CvT6IV6WEAQhhV5.jpg https://pbs.twimg.com/media/CtVAvX-WIAAcGTf.jpg Name: jpg_url, dtype: int64 **Twitter API Query** The last dataframe we will generate is from the twitter api of our developer account. We will need four passwords: "consumer key", "consumer secret key", "access token", and "access secret token". We will then create our last dataframe with four columns which we need it later: "date_time", "tweet_id", "retweets", "favorites". Fortunately, the datafram we got this time is clean. **Cleaning Data** Issue Types **Quality Issues** df_archive The data type of the timestamps should be datetime. IDs should be strings or integers. Missing values found in several columns. The source column should be more readable. Some rating numerators and denominators are not real ratings. Some names are detected as other words. df_image_predictions p1, p2, p3 should be categorical instead of strings. · Some pictures are retweets. **Tidiness Issues** • The dog stage spread across 4 columns - doggo, floofer, pupper, puppo. • Information about tweets is spread across three different dataframes. Our process of data cleaning starts with tidiness issues. We will first make a copy of all three dataframes and merge them into a master dataframe to make our next cleaning steps more convenient. Secondly, we will combine all dog stages into one unique column, and check if every stages match the texts. **Tidiness Issues** 1. Merging the 3 dataframes into one named "df_new" by pd.merge() Result: The new dataframe have 31 columns in total. 1. Combining different dog stages into one column. • Count the number of all stages and thier combination using groupby. Merge the stage columns and create a new column "stage". Inspect the rows not within the four stages and refine them. doggo floofer pupper puppo count None None None None 1740 1 None None None 23 puppo None None 210 None pupper None floofer None None 7 65 doggo None None None 5 1 doggo None None puppo doggo 11 None pupper None floofer 1 None None doggo We can see there were some rows with multiple stages. Now let's combine all columns with the '+' operator, and use comma to seperate columns with multiple stages. Next, we will look through every texts or pictures to confirm the correct dog stage. Finally, we will drop the old columns 'doggo', 'floofer', 'pupper', and 'puppo'. The cleaning result will be: Result: 2058 non-null object dtypes: bool(3), float64(7), int64(6), object(12) memory usage: 424.1+ KB Quality Issues 1. Drop NaNs, duplicates and unmeaningful columns. Result: <class 'pandas.core.frame.DataFrame'> Int64Index: 1986 entries, 0 to 2057 Data columns (total 22 columns): 1986 non-null int64 tweet_id 1986 non-null object times tamp 1986 non-null object source 1986 non-null object expanded_urls 1986 non-null object 1986 non-null int64 rating_numerator 1986 non-null int64 rating_denominator 1986 non-null object name 1986 non-null int64 retweets 1986 non-null int64 favorites 1986 non-null object jpg_url

The Wrangle Act of WeRateDogs

twitter-archive-enhanced, and name it as df archive.

tweet_id

timestamp

source

name

doggo

floofer

in_reply_to_status_id

in_reply_to_user_id

retweeted_status_id

expanded_urls

rating_numerator

rating_denominator

retweeted_status_user_id

retweeted_status_timestamp

In this project of data wrangling, we will go through three essential processes: gather, access and cleaning. There are three main datasets we are going to deal with, after gathering and accessing, we will then merge them into a master dataframe, and continue

The first data we will see is the twitter archive data, also the data we are going to work with the most. We will read it from a file called

Take a brief look, we can see some pieces that we might need to handle later. To begin with, let's start from the data information:

2356 non-null int64 78 non-null float64

78 non-null float64

2356 non-null object

2356 non-null object

2356 non-null object 181 non-null float64

181 non-null float64

181 non-null object 2297 non-null object

2356 non-null int64

2356 non-null int64

2356 non-null object

2356 non-null object

2356 non-null object

<class 'pandas.core.frame.DataFrame'> RangeIndex: 2356 entries, 0 to 2355 Data columns (total 17 columns):

Table of Contents

finishing the cleaning section.

Twitter Archive

 Twitter Archive Image Prediction Twitter API Query Cleaning Data

iPhone Twitter for iPhone Twitter

img_num

pl_conf

pl_dog

p2_conf

p3_conf

memory usage: 316.1+ KB

p3_dog

stage

p2_dog

p1

p2

р3

1. Make the source content more readable using <code>apply()</code>.

Result:

source

Twitter for iPhone

Twitter for

for iPhone

Twitter for iPhone

Result:

Result:

inferred from text.

1. Refine columns in image prediction.

· Apply our loop results into new columns.

Apply the results into a new column "gender".

in 'text' column, then set the correct float ratings.

Drop the not required columns of image prediction information.

prediction confidence

1986 non-null int64

1986 non-null object 1986 non-null float64

1986 non-null object

1986 non-null object

1986 non-null float64

1986 non-null object

1986 non-null float64

1986 non-null bool

1986 non-null bool

1986 non-null bool

dtypes: bool(3), float64(3), int64(6), object(10)

For example: WeRateDogs® @ @dog_rates · Jul 8, 2017 This is Bella. She hopes her smile made you smile. If not, she is also offering you her favorite monkey. 13.5/10

Create prediction and confidence list and store the true algo with its level of confidence.

memory usage: 232.7+ KB

None male

female

· Loop through the boolean columns and append the prediction and confidence columns into our lists.

dtypes: float64(1), int64(5), object(8)

1. Get the correct dog gender from text column by filtering common phrases and assigning a gender. Using 'None' if no gender can be

1129

633

224 Name: gender, dtype: int64

1. Fix rating numerator and denominators that are not actually ratings by seeking all occurences where there are more than one #/# in 'text' column. Using re.findall() to scrap through the numerator column and str.contains() to view tweets with decimals in rating

• Create a male and female list and loop through all texts to check if it has one of pronouns of male or female.

1986 non-null object

1986 non-null float64

1 8.8K 41.2K 181 expanded_urls rating_numerator rating_denominator text This is Bella. She hopes her https://twitter.com/dog_rates/status/883482846. 10 smile made you sm... Result: text expanded_urls rating_numerator rating_denominator This is Bella. She hopes her https://twitter.com/dog_rates/status/883482846... 10.0 13.50 smile made you sm... 1. Remove unmeaningful names with lowercase by finding with str.isupper() and fill the nulls, then check the results with str.islower() . Result: There are no rows with strange names. 1. Correct all datatypes and check them with .dtypes. Result: object tweet_id datetime64[ns] timestamp

rating_numerator float64 float64 rating_denominator object retweets int64 favorites int64 object jpg_url category stage prediction object confidence float64 gender category dtype: object Finally, we will store all our progress above in a file called twitter archive master.csv, and use them in our data visualization.

source

expanded_urls

text

category

object

object