


```
[39]: df_clean.head(1)
```

```
tweet_id  timestamp                source      text      expanded_urls  rating_name
0  89242064355536193  2017-08-01T16:23:56+00:00  href="http://twitter.com/download/iphon...  This is Phineas. He's a mystical boy. Only eve...  https://twitter.com/dog_rates/status/892420643...

1 rows × 22 columns
```

```
In [40]: df_clean.info()
<class 'pandas.core.frame.DataFrame'>
Int64Index: 2175 entries, 0 to 2355
Data columns (total 22 columns):
#   Column                Non-Null Count  Dtype
---  --
0   tweet_id              2175 non-null    int64
1   timestamp             2175 non-null    object
2   source                2175 non-null    object
3   text                  2175 non-null    object
4   expanded_urls         2117 non-null    object
5   rating_numerator      2175 non-null    int64
6   rating_denominator    2175 non-null    int64
7   name                  2175 non-null    object
8   dog_stage             364 non-null     object
9   retweet_count         2175 non-null    float64
10  favorite_count        2175 non-null    float64
11  jpg_url               1994 non-null    object
12  img_num               1994 non-null    float64
13  p1                    1994 non-null    object
14  p1_conf               1994 non-null    float64
15  p1_dog               1994 non-null    object
16  p2                    1994 non-null    object
17  p2_conf               1994 non-null    float64
18  p2_dog               1994 non-null    object
19  p3                    1994 non-null    object
20  p3_conf               1994 non-null    float64
21  p3_dog               1994 non-null    object
dtypes: float64(6), int64(13), object(13)
memory usage: 390.8+ KB
```

Quality Issue_3

Define

- Correcting data type in tweet_id (from int into string).

Code

```
In [41]: df_clean['tweet_id'] = df_clean['tweet_id'].astype(str)
```

Test

```
In [42]: df_clean.head(1)
```

```
Out[42]:
```

tweet_id	timestamp	source	text	expanded_urls	rating_name
0	89242064355536193	2017-08-01T16:23:56+00:00	href="http://twitter.com/download/iphon... This is Phineas. He's a mystical boy. Only eve...	https://twitter.com/dog_rates/status/892420643...	

```
1 rows × 22 columns
```

```
In [43]: df_clean.info()
<class 'pandas.core.frame.DataFrame'>
Int64Index: 2175 entries, 0 to 2355
Data columns (total 22 columns):
#   Column                Non-Null Count  Dtype
---  --
0   tweet_id              2175 non-null    object
1   timestamp             2175 non-null    object
2   source                2175 non-null    object
3   text                  2175 non-null    object
4   expanded_urls         2117 non-null    object
5   rating_numerator      2175 non-null    int64
6   rating_denominator    2175 non-null    int64
7   name                  2175 non-null    object
8   dog_stage             364 non-null     object
9   retweet_count         2175 non-null    float64
10  favorite_count        2175 non-null    float64
11  jpg_url               1994 non-null    object
12  img_num               1994 non-null    float64
13  p1                    1994 non-null    object
14  p1_conf               1994 non-null    float64
15  p1_dog               1994 non-null    object
16  p2                    1994 non-null    object
17  p2_conf               1994 non-null    float64
18  p2_dog               1994 non-null    object
19  p3                    1994 non-null    object
20  p3_conf               1994 non-null    float64
21  p3_dog               1994 non-null    object
dtypes: float64(6), int64(2), object(14)
memory usage: 390.8+ KB
```

Quality Issue_4

Define

- Correcting data type in timestamp (from string into datetime).

Code

```
In [43]: df_clean['timestamp'] = pd.to_datetime(df_clean['timestamp'])
```

Test

```
In [45]: df_clean.head(1)
```

```
Out[45]:
```

tweet_id	timestamp	source	text	expanded_urls	rating_name
0	89242064355536193	2017-08-01T16:23:56+00:00	href="http://twitter.com/download/iphon... This is Phineas. He's a mystical boy. Only eve...	https://twitter.com/dog_rates/status/892420643...	

```
1 rows × 22 columns
```

```
In [46]: df_clean.info()
<class 'pandas.core.frame.DataFrame'>
Int64Index: 2175 entries, 0 to 2355
Data columns (total 22 columns):
#   Column                Non-Null Count  Dtype
---  --
0   tweet_id              2175 non-null    object
1   timestamp             2175 non-null    datetime64[ns, UTC]
2   source                2175 non-null    object
3   text                  2175 non-null    object
4   expanded_urls         2117 non-null    object
5   rating_numerator      2175 non-null    int64
6   rating_denominator    2175 non-null    int64
7   name                  2175 non-null    object
8   dog_stage             364 non-null     object
9   retweet_count         2175 non-null    float64
10  favorite_count        2175 non-null    float64
11  jpg_url               1994 non-null    object
12  img_num               1994 non-null    float64
13  p1                    1994 non-null    object
14  p1_conf               1994 non-null    float64
15  p1_dog               1994 non-null    object
16  p2                    1994 non-null    object
17  p2_conf               1994 non-null    float64
18  p2_dog               1994 non-null    object
19  p3                    1994 non-null    object
20  p3_conf               1994 non-null    float64
21  p3_dog               1994 non-null    object
dtypes: datetime64[ns, UTC](1), float64(6), int64(2), object(13)
memory usage: 390.8+ KB
```

Quality Issue_5

Define

- Delete rows with jpg_url have missing.

Code

```
In [47]: df_clean = df_clean(df_clean['jpg_url'].notnull())
```

Test

```
In [48]: df_clean.info()
<class 'pandas.core.frame.DataFrame'>
Int64Index: 1994 entries, 0 to 2355
Data columns (total 22 columns):
#   Column                Non-Null Count  Dtype
---  --
0   tweet_id              1994 non-null    object
1   timestamp             1994 non-null    datetime64[ns, UTC]
2   source                1994 non-null    object
3   text                  1994 non-null    object
4   expanded_urls         1994 non-null    object
5   rating_numerator      1994 non-null    int64
6   rating_denominator    1994 non-null    int64
7   name                  1994 non-null    object
8   dog_stage             326 non-null     object
9   retweet_count         1994 non-null    float64
10  favorite_count        1994 non-null    float64
11  jpg_url               1994 non-null    object
12  img_num               1994 non-null    float64
13  p1                    1994 non-null    object
14  p1_conf               1994 non-null    float64
15  p1_dog               1994 non-null    object
16  p2                    1994 non-null    object
17  p2_conf               1994 non-null    float64
18  p2_dog               1994 non-null    object
19  p3                    1994 non-null    object
20  p3_conf               1994 non-null    float64
21  p3_dog               1994 non-null    object
dtypes: datetime64[ns, UTC](1), float64(6), int64(2), object(13)
memory usage: 358.3+ KB
```

Quality Issue_6

Define

- Create 1 column for dog image prediction and 1 column for dog image prediction confidence.

Code

```
In [49]: prediction_dog = []
prediction_dog_confidence = []

def dog_image(df_clean):
    if df_clean['p1_dog'] == True:
        prediction_dog.append(df_clean['p1'])
        prediction_dog_confidence.append(df_clean['p1_conf'])
    elif df_clean['p2_dog'] == True:
        prediction_dog.append(df_clean['p2'])
        prediction_dog_confidence.append(df_clean['p2_conf'])
    elif df_clean['p3_dog'] == True:
        prediction_dog.append(df_clean['p3'])
        prediction_dog_confidence.append(df_clean['p3_conf'])
    else:
        prediction_dog.append('Error')
        prediction_dog_confidence.append('Error')

df_clean.apply(dog_image, axis=1)

df_clean['prediction_dog'] = prediction_dog
df_clean['prediction_dog_confidence'] = prediction_dog_confidence
df_clean = df_clean[df_clean['prediction_dog'] != 'Error']
```

Test

```
In [50]: df_clean.head(4)
```

```
Out[50]:
```

tweet_id	timestamp	source	text	expanded_urls	rating_name
1	89217421306343426	2017-08-01T00:17:27+00:00	href="http://twitter.com/download/iphon... This is Tilly. She's just checking pup on you...	https://twitter.com/dog_rates/status/892174213...	
2	891815181378084864	2017-07-31T00:18:03+00:00	href="http://twitter.com/download/iphon... This is Archie. He is a rare Norwegian Pomerian...	https://twitter.com/dog_rates/status/891815181...	
3	89168959727859688	2017-07-30T15:58:51+00:00	href="http://twitter.com/download/iphon... Darla. She commenced a snack mid meal...	https://twitter.com/dog_rates/status/891689597...	
4	89132755892668256	2017-07-29T16:00:24+00:00	href="http://twitter.com/download/iphon... This is Frankie. He would like you to stop ca...	https://twitter.com/dog_rates/status/891327558...	

```
4 rows × 24 columns
```

```
In [51]: df_clean.info()
<class 'pandas.core.frame.DataFrame'>
Int64Index: 1686 entries, 1 to 2355
Data columns (total 24 columns):
#   Column                Non-Null Count  Dtype
---  --
0   tweet_id              1686 non-null    object
1   timestamp             1686 non-null    datetime64[ns, UTC]
2   source                1686 non-null    object
3   text                  1686 non-null    object
4   expanded_urls         1686 non-null    object
5   rating_numerator      1686 non-null    int64
6   rating_denominator    1686 non-null    int64
7   name                  1686 non-null    object
8   dog_stage             274 non-null     object
9   retweet_count         1686 non-null    float64
10  favorite_count        1686 non-null    float64
11  jpg_url               1686 non-null    object
12  img_num               1686 non-null    float64
13  p1                    1686 non-null    object
14  p1_conf               1686 non-null    float64
15  p1_dog               1686 non-null    object
16  p2                    1686 non-null    object
17  p2_conf               1686 non-null    float64
18  p2_dog               1686 non-null    object
19  p3                    1686 non-null    object
20  p3_conf               1686 non-null    float64
21  p3_dog               1686 non-null    object
22  prediction_dog         1686 non-null    object
23  prediction_dog_confidence 1686 non-null    object
dtypes: datetime64[ns, UTC](1), float64(6), int64(2), object(15)

```