

# Report

On

## Data Wrangling

The Waterdogs Twitter project goals included:

- 1) Wrangling the twitter data through the following processes:
  - Gathering Data
  - Assessing Data
  - Cleaning Data
- 2) Storing, analysing and visualizing your wrangled data
- 3) Reporting on the data wrangling efforts and data analyse and visualization

### Gathering Data

The wrangling methods used data from different sources that includes

twitter\_archive\_enhanced.csv

tweet image predictions

Twitter API that uses Tweepy library to gather ] tweet

### Assessing Data

#### Quality Issue:

**twitter-archive-enhanced-2.csv**

- Completeness:  
missing data in the following columns:
  - in\_reply\_to\_status\_id,
  - in\_reply\_to\_user\_id,
  - retweeted\_status\_id,
  - retweeted\_status\_user\_id,
  - retweeted\_status\_timestamp,
  - expanded\_urls.

tweet\_id is an int (applies to all tables)

- Validity:

dog names: some dogs have 'None' as a name, or 'a', or 'an.'

This data-set includes retweets, which means there is duplicated data

Accuracy:

retweeted\_status\_timestamp inconstance data type

ime-stamp (wrong data type)

Consistency:

The Source column still has the HTML tags

rating\_denominator should be a standard 10, but there are a multitude of other values

### **image\_predictions.tsv**

Validity:

p1, p2 and p3 columns have invalid data

Consistency:

the dog breed listed is all lowercase

some columns there is an underscore for multi-word dog breeds.

### **tweet\_json**

Completeness:

Missing Some Data

### **Tidiness Issue:**

#### **twitter-archive-enhanced-2.csv**

The last four columns all relate to the same variable (dogoo, floofer, pupper, puppo).

### **image\_predictions.tsv**

This data set is part of the same observational unit as the data in the 'twitter-archive-enhanced-2.csv' - one table with all basic information about the dog ratings.

### **tweet\_json**

This data set is also part of the same observational unit - one table with all basic information about the dog ratings.

## **Cleaning Data**

After the assessment, cleaning the data through the following means:

### **Define, Code and Test:**

Merge the clean versions of archive, images, and twitter\_counts\_df data frames

Correct the dog types

Remove columns no longer needed.

Change tweet\_id from an integer to a string.

Change the timestamp to correct datetime format.  
Correct naming issues and Standardize dog ratings.