

Data Wrangling Report

Introduction

This Wrangle and Analyze Project is part of Udacity's Data Analyst Program. In this project the data wrangling, which consists of:

- Gathering data, download resources from different references
- Assessing data, for quality and tidiness issues
- Cleaning data, that identified in previous step
- Storing, analyzing, and visualizing wrangled data
- Reporting on data analyses and visualizations

The most tools, libraries and programming language used in this project are:

- Python
- Numpy Library
- Pandas Library
- Requests Library
- Tweepy Library
- Json Library
- Matplotlib Notebook
- Twitter's API

Gathering

Data was gathered from 3 different sources:

1. The 'twitter-archive-enhanced' was provided by Udacity. This file includes a huge variety of variables for each tweet as tweet_id, timestamp, text, name, rating, etc.
2. The 'image_prediction' file was downloaded using the Requests library.
3. Additional data, including favorite count and retweet count were gathered using the Twitter API.

Assessing Data

After gathering data were used the following methods:

- .head()
- .sample()
- .info()
- .value_counts()

Quality issue:

- Tweet_id was the incorrect data type
- The columns contained the 'None' instead of 'NaN'
- Timestamp was the incorrect datatype
- Missing values from images dataset (2075 rows instead of 2356)
- Sources are not readable

Tidiness:

- Dogstage was in 4 columns(doggo, floofer, pupper, puppo), there is no reason to keep them
- Merge all dataframes together as they all contained the proper information

Cleaning Data

This part of the data wrangling was performed in three stages: Define, Code and Test. First, copies of DataFrames were created before cleaning. Then used methods for investigation: `merge()`, `reduce()`, `drop()`, `replace()`, `head()`, etc

Storing

The final dataset was stored as 'twitter_archive_final' in a csv file. At this step, the data was successfully wrangled and therefore ready for analysis and visualization

Analysis and Visualization

Visualizations are provided in `act_report.pdf`