

# wrangle\_report

## 简介

本次项目选取推特昵称为 WeRateDogs 的档案为数据集进行数据清洗分析和可视化。因为这份档案只包含基础信息，还需要另外收集相关数据，一起进行评估和清洗，得出结论。

## 一、 收集数据

- twitter\_archive\_enhanced.csv 推特档案
- tweet\_json.txt 包含转发数(retweet\_count)和喜欢数(favorite\_count)
- image\_predictions.tsv 图像预测文件,对推特中的品种(或其他物体)进行预测的结果

### 收集方式

- twitter\_archive\_enhanced.csv 可以直接下载
- tweet\_json.txt 从 API 下载
- image\_predictions.tsv 从 url 下载

## 二、 评估

通过目测和编程发现数据需要清理:

### 1.质量

#### twitter\_archive\_enhanced 表格

- retweeted\_status\_id 含有转发数据
- tweet\_id 是整数，不是字符串
- rating\_numerator 分子中有低于 10 分的数据
- rating\_denominator 分母有不等于 10 分的数据
- name 中有空值，且含 a、an 等错误
- doggo floofer pupper puppo 有全空的
- 与 image\_predictions 合并后 twitter 表中 jpg\_url 包含无图片数据

#### tweet\_json 表格

- tweet\_id 是整数，不是字符串

## image\_predictions 表格

- tweet\_id 是整数，不是字符串

## 2.整洁度

- twitter\_archive\_enhanced 表中狗狗的地位 (stage):doggo,pupper,puppo,floof(er)成为一列。
- 把 twitter 表、tweet 表、image 表合成一张表，三者的联系是 tweet\_id
- 删除 twitter 的无用列，只保留 tweet\_id,text,rating\_numerator,rating\_denominator,name, stage, tweet\_count,favorite\_count

# 三、清洁

## 1.备份:

通过 copy()对三个数据集合进行备份

## 2.处理缺失数据

删除 twitter\_archive\_enhanced 表中 doggo、floofer、pupper、puppo 这些列

## 3.整洁度

- 对 twitter\_archive\_enhanced 表新增一个变量 stage，暂时以 np.nan 代替，后续会从 text 中提取
- 利用函数 mege 把 twitter\_archive\_enhanced 表和 tweet\_json 合成一张表
- 删除 twitter\_archive\_enhanced 表无用列，只保留 tweet\_id,text,retweeted\_status\_id,rating\_numerator,rating\_denominator, name, stage,retweet\_count,favorite\_count

## 4.质量

- 使用 as.type()把 tweet\_id 的数据类型修改为字符串。因为 tweet 表格已经合并到 twitter 表格里，因此只需对 df\_twitter\_clean 和 df\_image\_clean 进行操作即可。
- 利用 drop 删除 df\_twitter\_clean 里 retweeted\_status\_id 的非空值，删除转发数据
- 使用正则表达式和 pandas 的 str.extract 方法，从 text 中提取 rating\_numerator,rating\_denominator。提取出的分子 rating\_numerator 的,去除特别高的，把范围限定在 11~16。分母 rating\_denominator 为 10。
- 使用正则表达式和 pandas 的 str.extract 方法，从 text 中提取 name
- 使用 str.findall 从 text 中提取 stage

- 把 `df_twitter_clean` 和 `df_image_clean` 合并在一起，方式是 `inner`，去除不包含图片的。

## 5.保存数据

把清洗好的数据存在 `twitter_archive_master.csv` 里。