

## 清洗过程

获取数据后，有twitter、twitter\_extra、image\_predictions3个dataframe。首先是需要解决什么问题，哪些数据是需要的。然后，3个dataframe以什么为基准合并。检查发现有相同的tweet\_id，以此为基准合并dataframe。然后开始系统清洗。去找质量问题和整洁度问题。

**根据思路整理出需要验证的问题：**

- 1、在twitter、twitter\_extra、image\_predictions中tweet\_id是不是唯一的。验证3个dataframe中个tweet\_id是否唯一，3个dataframe是否有交集。
- 2、检查twitter、twitter\_extra、image\_predictions是否都含有同tweet\_id。
- 3、检查twitter中丢失数据量。通过检查非空的数量。（in\_reply\_to\_status\_id、in\_reply\_to\_user\_id、retweeted\_status\_id、retweeted\_status\_user\_id、retweeted\_status\_timestamp）
- 4、检查是否所有tweet\_id都能对应url。
- 5、检查狗名是否正确
- 6、检查每列中格式是否正确
- 7、检查twitter dataframe中rating\_numerator、rating\_denominator是否正确
- 8、清除twitter\_extra中retweet\_count为0的项
- 9、image\_predictions中p1、p2、p3中明显错误的项目

## 验证问题后发现问题如下：

### 数据质量问题

- (1) twitter的数据frame中不需要处理的数据in\_reply\_to\_status\_id、in\_reply\_to\_user\_id、r (1) twitter的数据frame中不需要处理的数据in\_reply\_to\_status\_id、in\_reply\_to\_user\_id、retweeted\_status\_id、retweeted\_status。
- (2) twitter的数据frame中狗的名字字段含有None,a,an,a,such,quite,the，需要去掉或填充None类型。
- (3) twitter的数据frame中狗的名字字段None为字符串格式，需要修改为python中None。
- (4) twitter的数据frame中timestamp、retweeted\_status\_timestamp应该为时间格式。
- (5) twitter的数据frame中doggo、floofer、pupper、puppo中None不是str应该为None格式
- (6) twitter的数据frame中，有的推送的内容中含有2条狗，rating\_numerator、rating\_denominator有2个评分，这种情况就取2个多平均值。
- (7) twitter的数据frame中，有的推送的内容中含有多个“数字/数字”格式内容，提取内容了错误的rating\_numerator、rating\_denominator。
- (8) twitter的数据frame中，有的rating\_denominator>10的部分数据中，有的有用数据，需要修改。例如为99/90，可以理解为9条狗总分，修改为11/10
- (9) twitter\_extra的数据frame中retweet\_count为0的异常值
- (10) image\_predictions的数据frame中p1、p2、p3中明显错误的异常值

### 数据整洁问题

- (1) 所有的数据在多个dataframe中，需要合并到一个dataframe进行分析。以twitter、twitter\_extra、image\_prediction中对应id为基准合并。
- (2) twitter的数据frame中expanded\_urls含有多个url，需要清理掉

(3) twitter的dataframe中oggo、floofer、pupper、puppo需要合并成1列