

WeRateDogs 的推特档案包括了5000 多条推特的基本信息。档案中有一列包含每个推特的文本，用这一列数据提取了评分、狗的名字和“地位”。

推特json部分，我们根据tweet_id提取了转发数和喜爱数。

图像预测部分根据了图片进行了狗狗种类的预测，并给出了预测的可信度。

我们分别从本地读取推特档案csv、API收集转发数和喜爱数、网上下载图像预测tsv文件。然后将上述三个表格根据tweet_id进行合并，成一个完整的数据框。

数据评估流程

1. 查看数据格式是否合适，比如tweet_id应为int64
2. 查看数据是否重复
3. 挨个查看每一列的数据有否异常，并记录问题供后续清洗
 1. 查看twitter发送时间是否连续
 2. 查看text是否重复
 3. 查看评分分子是否过大或者过小，是否单位一致
 4. 查看评分分母是否不为10
 5. 查看姓名列是否提取正确
 6. 查看转发数和喜爱数有没有异常值
 7. 查看狗狗的doggo、floofer、pupper、puppo四列是否完整，和满足独热编码的要求
 8. 查看一下预测文件的性质，比如可信度的特征、有多少不是狗

整理出以下问题：

质量问题

1. tweet_id存在重复
2. 质量问题2：分子出现极端值 1776
3. 质量问题3：分子分母成比例，导致同列的分子单位不一致，即不为 1/10
4. 质量问题4：分母不全都为10
5. 质量问题4：部分名字提取不正确，有 765个"None" 和 67个"a"，还有一个叫'space'，不是大写，可能有问题
6. 质量问题5：狗狗的stage不全，还有"blep", "snoot", "BlepiPen"
7. 质量问题6：有一部分用户上传的图片不是狗狗照片，不应该被记录在这个项目中。
8. 有些情况，存在狗狗有两个stage，查看后，发现这种情况对应的情况包括：图片中有2只狗狗、作者进行比较的时候提及了另一stage等，因为情况复杂且数量相对很少（11个），可以删掉。
9. 漏将floofl归类为floofer。

整洁问题

1. 将3个表描述的是同一个事，应合并成一个表
2. 不应该将狗的stage分成doggo、floofer、pupper、puppo四列，以独热编码方式记录，应该记为一列'stage'，而为真的原列名为内容。

