

基于规则的评论评分

赵海臣

背景

🌀 为了鼓励用户更多地写下自己的实际使用感受，增加产品的信用度，以及提高顾客对商品的认识，需要设计一种算法能对用户所写的评论进行评分，并且能够动态实时地在用户输入评论时及时显示分数，鼓励用户进一步完成评论。并且在达到一定分值后实行消费券奖励。

基础版

- 🌀 小短句计分算法
- 🌀 总句计分算法
- 🌀 分数非线性化
- 🌀 广告识别与屏蔽
- 🌀 无意义评论过滤

小短句评分

✧ 小短句以用户的特殊字符作为界限，如空格、标点之类的。

✧ 每个小短句上限100分

- ★ 有n : $nScore = 30$
- ★ 有a : $aScore = 30$
- ★ 有v : $vScore = 30$
- ★ 附加分: 每多一种其它词性 $oScore = 10$

✧ 优质标签匹配附加分100分

- ★ n+a组合能匹配到我们定义词库标签的，额外
 - $+MatchScore = 100$

✧ 小短句评分

$$SentenceFrac = \frac{\min(100, nScore + aScore + vScore + n * oScore) + MatchScore}{100}$$

小短句单字评分

🌀 目标：需要达到我们预设的字数才能获得奖励，小短句需要结合字数进一步评分，分数分为字数基础分与内容分

★ 假设目标字数阈值为 $TargetWordThr = 100$

★ 假设字数基准分为20%， $WordValue = 20$

● 每写一个字：

$$WordCountScore = \frac{WordValue}{TargetWordThr}$$

★ 假设内容分为80%， $WordContentValue = 80$

● 短句内容加成分：

$$WordContentScore = \frac{WordContentValue}{TargetWordThr} * SentenceFrac$$

🌀 小短句总分：

★ 为了防止一条句子过长，设定小短句的字数上限

● $WordCountLimit = 20$

$$SentenceScore = (WordCountScore + WordContentScore)$$

$$* \min(WordCountLimit, SentenceWordCount)$$

词库标签匹配奖励规则

- ✎ 由于词库标签无论正面还是负面在定制过程中都是根据筛选的高频评论来出的，因此只要能匹配到标签库标签，说明短句是有规则语义的，也是能符合我们标签提取规则的优质评论。
- ✎ 因此小短句中能合成标签库标签，奖励分可以适当提高。

额外特征加分

有一些特有的可配置专用名词或者符号列表，若匹配上，同样有额外的小短句奖励分。

	手机评价上不了图比较遗憾??	
j30231	之前试过别人的，觉得好看就想买一个，之前在香港买的1号色的太淡了，而且有点干，找了很多地方都没有这个6号色的，所以看到聚美有就赶紧下手，滋润而且不带金粉，不会干，颜色也很正，可能之前试的时候只是擦了一点，所以买了这个自己擦得多一些显得很红，下次要试试4号色。。。这个性价比真的很高	107.96
166831	这个套装特别的水润，不油腻，很清爽，特别适合干性皮肤的人。就是感觉水和乳不太好倒出来，其他还好。价格也算划算，刚好是在小美搞活动的时候买的，还送了大礼包。觉得是一次很愉快的购物，等这一套用完效果不错的话，会再次选择御泥坊的东西，包装也挺好看	107.76
303031	近江兄弟小熊套装，一瓶28毫升，一个大太阳用，一个室内用，不过我一般情况下我都拿红色的擦身上，因为脸太油了，干皮妹子适用哦*^o^*。蓝色的话我一般拿她擦脸，不过薄荷味的擦在脸上又麻又热又凉?，不过刚擦完进空调屋不能更舒服→_→，50+的防晒足够用了，卸的话脸上用卸妆水身上用沐浴乳就够了?	107.12
	不得不说，小美家nuk的乳胶奶嘴价格算是比较便宜的了，比旗舰店便宜，实体店经常买不到，但是我想说的是乳胶奶嘴真的特别容易老化、	

评论附带的图片，有额外加分。

线上评论实时评估

线上用户实时评论包括“文本”和“图片”两部分，我们需要在用户输入过程中展示出对评论的系统评判：



- ★ 实时评论基础打分 = 文本分满分100 + 图片分满分100
 - 文本分：上述基于规则的文本分，取上限值100
 - 图片分：一个图片50分，两个100分，取上限值100

线上实时展示分的非线性映射

🌀 用户实时输入评论时，我们能够直接获得“文本”+“图片”的情况，新评论肯定没有点赞数、评论数，所以我们只需要考虑“文本”+“图片”作为实时评论分。

🌀 文本最大分100+图片最大分100=200，并且分数的增长偏线性，容易被顾客猜到逻辑，在展示给顾客实时评分时，应该对分数进行非线性映射：

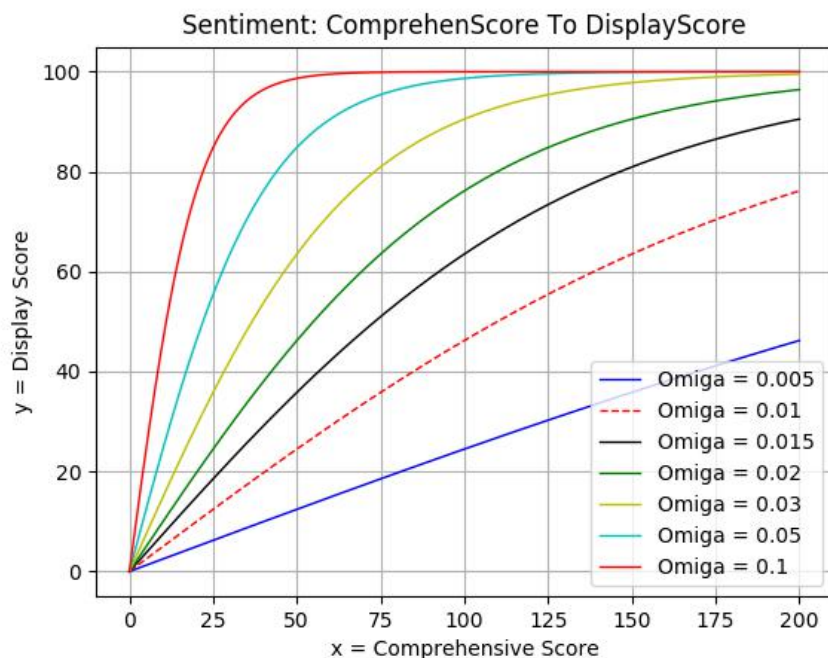
- ★ 将0-200的分数增长，合理地映射到0-100区间
- ★ 分数最开始的时候增长较快，鼓励用户，后续慢慢涨幅减慢，敦促用户完成更多评论。

非线性实时分数映射公式

采用sigmoid进行非线性映射：

$$DisplayScore = 200 * \left(\frac{1}{1 + e^{-\omega * ComprehensionScore}} - 0.5 \right)$$

★ 不同参数映射率，文本分满分100，图片分满分100，横坐标是二者加起来的分，纵坐标是看到的分



参数的选择与实时分数的分类与效果

🌀 sigmoid的 ω 选择0.02，考虑以下几点：

- ★ 文本特别特别好，获得100分，没有一个图片，图片分0分，那显示的是100对应的76分
- ★ 图片特别丰富，100分，但没文字，文字0分，那也是76分
- ★ 用户在输入过程中会慢慢感知分数增加，开始增加很快，慢慢会变慢

🌀 因此

- ★ 假设80分以上是优秀评论，用户只写特别好的评论，或者特别多的图片，也只能达到76分，需要图文结合才能达到好评论的标准
- ★ 90分对应的是150，一般情况下是两张图片 + 至少50分以上较高质量的评论，这个标准可以作为优质评论

🌀 效果：

- ★ 内容丰富 = 必须要有图片，1张图片加中量评论，或者，2张以上图片加一些文字
- ★ 非常完美 = 必须要有图片，并且文字写得足够丰富

广告屏蔽

✎ 广告主体一定要留下通讯方式，一般都是以微信号或QQ或QQ群，这几个联系方式的主体都是数字与字母

✎ 屏蔽规则：

- ★ 去掉非汉字的特殊字符外连续的纯数字或者字母数字拼接，长度大于6，并且，非重复字符串(有66666...,11111...等网络数字连词可能会被误认为广告)

✎ 处理规则：

- ★ 实时输入时，假装没发现，正常计分，但是永远达不到奖励值，发布体查看能看到评论，但别人看不到。

无意义评论过滤

☞ 无意义评论常见以下几种：

- ★ 重复字/字符串阈值，高于该阈值的句子当作垃圾级，永远达不到奖励值，并且仅自己可见
 - 连续重复字阈值：
 - $\text{ConsWordsAsJunk} = 10$
 - 连续重复字符串阈值：
 - $\text{ConsStringAsJunk} = 5$
- ★ 乱按键盘打出的无意义句子
 - 分词器若分不出词性，短句分数达不到要求；
 - 若无充分的标点与空格，永远不可能达到我们的奖励分值；
 - but, 无法充分防范，只能尽量降低误判概率。

关键敏感词过滤

✎ 将一些额外的过滤词设置为一个过滤词词表，可以通过过滤词词表对一些垃圾评论进行过滤。

THE END

THANK YOU!

*分数非线性化(未使用)

✧若分数过于线性化，容易造成刻板并且可能会被猜到规则，因此需要对分数的增长进行非线性化处理。

✧非线性化公式：将0-1区域内映射出一条xy弧线。

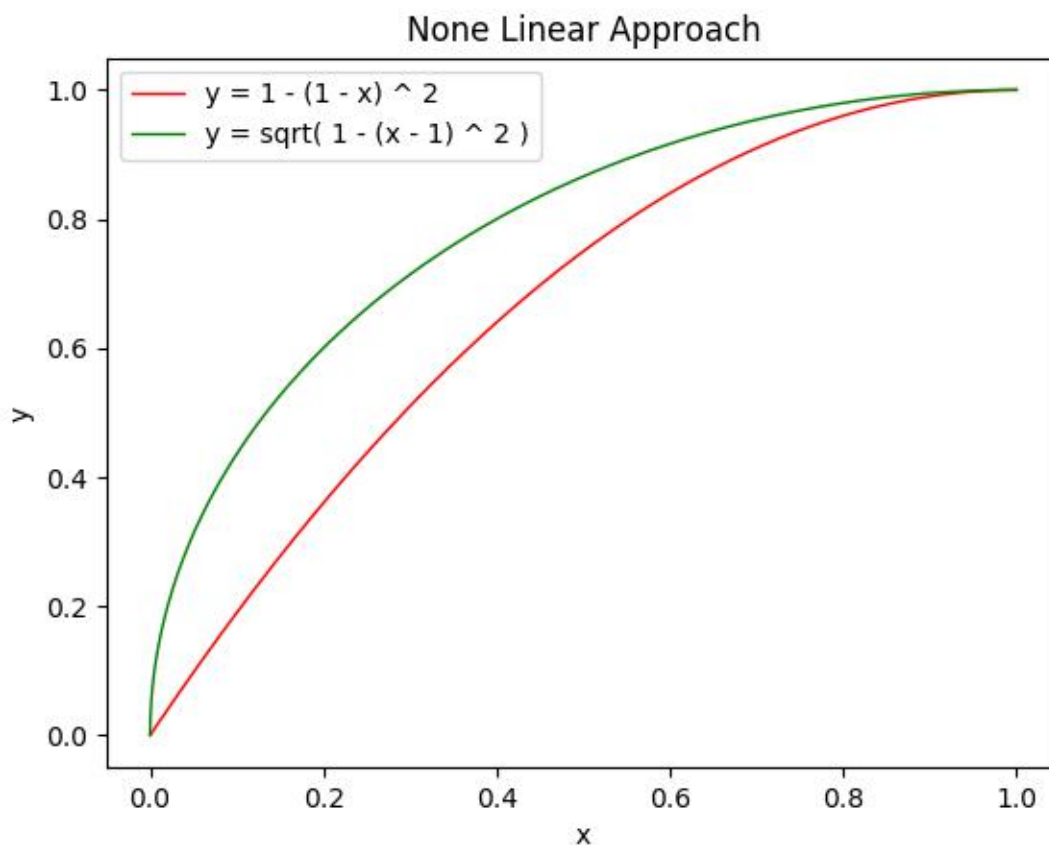
★ 二次方非线性逼近：

$$y = 1 - (1 - x)^2$$

★ 圆周非线性逼近：

$$y = \sqrt{1 - (x - 1)^2}$$

*分数非线性化(未使用)



*最终分数(未使用)

✎假设目标奖励分ScoreWithGift = 80

✎使用圆周非线性逼近的情况下，我们的动态展示分数为：

$$DynamicScore = ScoreWithGift * \sqrt{1 - ((\frac{\sum SentenceScore}{100}) - 1)^2}$$

✎当 $\frac{\sum SentenceScore}{100} > 1$ 之后，圆周非线性方程会发生异常，因此，在圆周方程下：

$$DynamicScore = ScoreWithGift * \frac{\sum SentenceScore}{100}$$