

基于情感分析的“真假美猴王”存疑研究*

张辰麟¹, 王明文²⁺, 谭亦鸣², 陈志明², 左家莉², 罗远胜²

(1. 江西师范大学文学院, 江西, 南昌, 330022;

2. 江西师范大学计算机信息工程学院, 江西, 南昌, 330022)

摘要:《西游记》是我国四大名著之一,“真假美猴王”事件,作为《西游记》的高潮部分,留下了不少伏笔,也引发了多种解读。本文通过运用情感分析的方法,对“真假美猴王”事件前后孙悟空与其他角色的对话进行分析,比较孙悟空在“真假美猴王”事件前后对其他角色的情感值的变化,得到了“孙悟空并没有被如来打死,‘真假美猴王’事件消灭的‘心魔’是孙悟空的反抗精神,事件之后,孙悟空选择屈服于神权”的结论。初步探索了情感分析技术对文学研究的可行性。

关键词:情感分析;文学情感分析;情感词典;《西游记》;真假美猴王

A Research on *Journey to the West* Based on Sentiment Analysis

Chenlin Zhang¹, Mingwen Wang²⁺, Yiming Tan², Zhiming Chen², Jiali Zuo²,
Yuansheng Luo²

(1. College of Chinese Language and Literature, Jiangxi Normal University, Nanchang 330022, China;

2. School of Computer and Information Engineering, Jiangxi Normal University, Nanchang 330022, China

+ Corresponding author)

Abstract: As one of the Four Great Classical Novels, *Journey to the West* left lots of foreshadowing to interpret. In this paper, we carry out a research on Monkey King by using sentiment analysis. We apply NLP technologies: automatic segmentation and sentiment lexicon collection to calculate the sentiment of Monkey King. By judging the changes of the sentiment of Monkey King before and after the case “Real and Fake Monkey King”, we finally proposed the conclusion: Monkey King was not killed by Rulai, the supreme Buddha, but he started to tend to bend to the religious authority after the case. This paper made a primary exploration for sentiment analysis on literary studies.

Key words: sentiment analysis; sentiment analysis on literature; sentiment lexicon; *Journey to the West*; Real and Fake Monkey King

1 引言

《西游记》^[1]是中国古代第一部长篇章回体神魔小说。全篇描写了唐三藏远赴西天求取真经的故事,深刻地揭露了当时社会的现实。作为四大名著之一,《西游记》对中国文学的

* 基金项目: 国家自然科学基金(61462045, 61462043, 61562031), 江西省自然科学基金青年项目(20151BAB217014)。

意义不言而喻，而其中的“真假美猴王”事件，把整部《西游记》的故事推向了高潮。该事件描写了孙悟空和唐三藏两个主要角色之间的矛盾。作者吴承恩为“真假美猴王”事件埋下了不少伏笔，这些伏笔引发了多种解读，相关的讨论主要包括孙悟空的性格与语言的变化^[2]、艺术形象^[3]、思想变化^[4]、紧箍咒和金箍棒的象征意义^{[3][5]}、六耳猕猴的身份^[6]、孙悟空和唐僧的关系^[7]等方面。加上一度在网络上引发热议的“被如来打死的究竟是谁”的问题，现有对“真假美猴王”的解读总结起来分为三个类型：

1. 认为真的孙悟空已经被如来佛祖打死，《西游记》后半部分取经的是六耳猕猴。
2. 认为真孙悟空还活着，死的是六耳猕猴，孙悟空是唐三藏的精神导师^[7]，“真假美猴王”一事是唐三藏的修行，消灭的是唐三藏的“心魔”，《西游记》后半部分孙悟空的桀骜不驯并没有收敛^[4]。
3. 认为真孙悟空还活着，“心魔”，即六耳猕猴^[6]，是孙悟空的反抗精神^[2]，“真假美猴王”一事是孙悟空“心的修行”^[3]，事件之后孙悟空走向逐渐被“同化”，屈服于神权，再无反意的悲剧结局^[8]。

以往对于文学作品的研究和相关讨论一般是基于文献法的定性研究，近年来，随着自然语言处理技术的迅猛发展，不少语言研究者开始利用自然语言处理的新方法和新手段，从定量的角度解决语言的相关问题，但几乎未涉及到文学领域。本文将尝试使用自然语言处理中情感分析的方法，对“真假美猴王”事件进行解读。情感分析又称情感计算，是对带有情感色彩的主观性文本进行分析、处理、归纳和推理^[9]，对文本的情感倾向做出判断^[10]的过程。情感分析技术被广泛运用于微博^{[11]-[15]}、用户评论^{[16]-[18]}中的情感倾向研究与预测。情感分析中情感的划分一般采用三分法^[9]（褒义、贬义、中性）或细粒度分法^[19]，其本质是一个文本分类问题。情感分析研究通常从文本数据挖掘^{[13][18][20]}、建立知识库^{[10][12][21]}、挖掘句法特征^{[16]-[18][22]}等方法入手，通过筛选种子词^{[12][13]}，建立情感词典^{[12][13][23]-[25]}并验证情感词典的有效性的方式，从而得到较为可靠的情感词典以分析文本的情感倾向。

目前，将情感分析技术应用于文学作品分析的相关研究尚不多见。本文将通过构建孙悟空的情感词典，分析“真假美猴王”事件前后孙悟空对其他角色的情感变化，从定量的角度分析“真假美猴王”事件，从而探索情感分析技术对文学研究的可行性。

2 语料选取与自动分词

2.1 语料的选取

为保证研究的真实性，本文研究对象为《西游记》原版，而非现代汉语版作为研究对象，并利用检索的方式，将冒号加双引号（：“……”）作为特征，从原版《西游记》中抽取所有角色的对话，并根据纸质版的《西游记》，人工对其句子错漏、别字等进行了改正，该方法一共抽取了 322307 字的《西游记》人物对话，总共包含 10664 句台词（不包括台词中的诗句部分），约占整个《西游记》总篇幅的一半。而后，从 32 万余字的人物对话中，手工挑选

出所有孙悟空的台词，合计 107009 字，包括 3358 句台词。

2.2 自动分词及分词结果优化

本文选择了 Jieba、NLPIR、Stanford 三种自动分词系统分别对这 3358 句台词进行分词。并随机抽取了《西游记》中的两章（37、81 章共 81 句台词）来验证自动分词的分词效果。验证阶段由两位文学学院的博士和两位古代文学专业的硕士作为专家，采用带权重的投票的策略，得到了一个更为合理的人工分词结果，并以此人工分词结果作为黄金标准。本文借鉴了机器翻译评测当中的 TER 值^[28]作为评测标准，将人工黄金标准作为参考译文，将自动分词的结果作为机器翻译结果，计算出分词符号被插入、删除、替换和移动等操作的编辑次数，以便衡量分词（Segment）效果的好坏， 本文将 TER 公式运用如式（1）：

$$TER = \frac{\text{编辑次数}}{\text{人工分词的次数}} \quad (1)$$

表 1 展示了三种分词系统及三种分词投票得到的分词结果的 TER 值：

表 1 三种分词系统及投票的 TER 值

	Beam search	Shifts tried	Total TER
Jieba	106	25	0.21024649589173514 (435.0/2069.0)
NLPIR	90	9	0.31706138231029485 (656.0/2069.0)
Stanford	100	19	0.26582890285161914 (550.0/2069.0)
Voting	88	7	0.24117931367810536 (499.0/2069.0)

我们发现，由于明清白话中白话文句子较短结构单一，人名、称呼、地名、专有名词等命名实体出现分词错误时很容易黏连到其前后的词语，造成连锁反应。基于此，我们对孙悟空的所有 3358 句台词中的人名^[29]、地名、专有名词^{[30][31]}等进行了人工筛选，并构建了一个 599 词的用户词表添加到分词系统。首先将待分词语料进行预处理，用数字前后加分词符的形式，替换掉用户词表中的词，以避免用户词表与分词系统内置词表的冲突，并采用了音序排列、大词在上的方法进行预处理替换，避免用户词表词与词之间存在包含关系，将替换之后的结果再次进行分词，通过后处理，将用户词表中的词还原，经此步骤，三种分词系统分词结果如表 2：

表 2 加入用户词表后的三种分词系统及投票的 TER 值

	Beam search	Shifts tried	Total TER
Jieba	104	23	0.19719671338811020 (408.0/2069.0)
NLPIR	87	6	0.28129531174480427 (582.0/2069.0)
Stanford	101	20	0.23827936201063316 (493.0/2069.0)
Voting	86	5	0.21991300144997583 (455.0/2069.0)

可以看出，Jieba 分词系统+用户词表的方案 TER 值最低，性能最好，因此最终我们选用了该方案进行全文分词。

3 情感种子词的获取

3.1 情感种子词训练集

《西游记》角色的台词，往往以“某某道：”作为开始，若作者已经预设了这句台词的情感，则会表示为“某某 XX 道：”，如“悟空骂道：”。其中“骂”即是台词的前缀，也是作者留给我们的明确的情感信息。因此，这些带有明确情感前缀的台词适合用来抽取情感种子词。

我们将所有带有情感前缀的台词挑拣出来，并将一部分例子示于表 3：

表 3 台词前缀情感分类

	示 例
正面情感前缀	笑道（开心、高兴）、唱个大喏道、躬身道、大喜道、跪拜道、甚喜道、即答礼道、赤淋淋跪下、大笑道、满心喜悦、心喜道、你道他叫好作甚……
负面情感前缀	笑道（嘲笑、讥笑等）、骂道、厉声高叫道、滴泪道、咄一声、心中怒发、大惊失色、打了两掌、掣棒高叫、高叫一声、后悔道、慌了……

在明清章回体小说中，“笑道”比较特殊，不仅包含了开心、高兴的笑，也包含了讥笑、讽刺、嘲笑等，正面/负面情感无法从字面上判断，因此对于以“笑道”为前缀的台词，我们单独进行了人工判断。

该阶段一共得到了 400 句台词的正面情感语料和 412 句台词的负面情感语料。

3.2 否定词表与停用词表

本文参考了现代汉语常用的否定词表，包含 35 个否定词，通过检索这 35 个否定词在《西游记》中出现的次数，发现一部分否定词不会在明清白话环境下出现，因此予以删除。并用否定词典中的否定词素，如“不、无、非、没”等作为对象^[32]进行搜索，挖掘了一些明清白话环境下的否定词，并将它们加入了否定词表，最终得到了一个 37 词的否定词表。

结合明清白话句子较短的特点，本文将否定词的支配域设定为否定词所在的分句。如某个分句中出现了否定词，则认为该分句中所有的情感词原本的情感倾向应与实际属于的语料分类相反（未出现分句中有复数否定词的特殊情况）。

通过挑出含有否定词的分句还原到其对应的分类，最终得到正面/负面情感的语料分别为：正面情感语料 2100 个分句，负面情感语料 2073 个分句。

本文尝试使用了四川大学、哈尔滨工业大学等研究机构研究的停用词表（共 1893 个停用词），发现以往的停用词表对明清白话的停用效果不理想，因此除了使用这些停用词表之外，还通过查阅词典的方式，停用了古今常用的数词、量词、代词、介词、连词、助词，以及前文用户词表中的地名、人名（不包括称呼）、否定词表、数量词词组等，加入停用词表。

3.3 种子词情感值计算

该部分的实验方法主要参考了赵妍妍^[22]等学者的情感词典构建方法，并根据《西游记》明清白话的语言性质对方法进行了一定的改动，本文使用式（2）^[33]来计算种子词属于正面/负面的情感倾向。

$$Polar_i = \frac{freq_{i-pos}}{freq_i} \times \ln freq_{i-pos} - \frac{freq_{i-neg}}{freq_i} \times \ln freq_{i-neg} \quad (2)$$

$freq_i$ 为词语 i 出现的频次, $freq_{i-pos}$ 为词 i 在正面情感语料中出现的频次, $freq_{i-neg}$ 为词 i 在负面情感语料中出现的频次, 由于语料规模较小, 词频次取对数使用了自然数 e 为底, 以避免得到的结果值过小, 由于指数为 1 的情况下对数为 0, 该式可以很好地屏蔽偶然事件。如 $Polar_i$ 大于零, 则证明该词属于正面情感词的可能性更高, 如果 $Polar_i$ 小于零, 则证明该词属于负面情感词的可能性更高, 该式得到的结果 $Polar_i$ 的绝对值越大, 证明词语的情感极性越大。

通过去除停用词, 最终得到了一个包括 330 个正面情感词和 332 个负面情感词的种子词词表, 记为 SeedA。通过专家人工对这 662 个词进行了复审, 剔除了其中分词错误词、无情感倾向的一般物质名词。我们发现, 一部分词语正面/负面的情感分类恰巧分反, 这与语料规模较小以及否定词支配域设定有一定的关系。对这些恰好分反极性的词语, 本文采取了两种策略, 一种是将极性分反的词还原到其本该属于的极性词表当中, 生成了包含 228 个正面情感词和 232 个负面情感词的种子词表, 记为 SeedB; 另一种策略是放弃这些词, 得到包含 189 个正面情感词和 198 个负面情感词的种子词表, 记为 SeedC。表 4 展示了情感种子词中正面/负面情感倾向最高的 Top20 词:

表 4 正面/负面情感种子词 Top20

正面情感词	弟子、放心、不敢、陛下、公主、请、兄弟、造化、令郎、指教、不打紧、承、恕、老、莫怕、情、保护、正、照顾、莫怪
负面情感词	呆子、孽畜、夯货、性命、我把你、棒、泼、贼、弄、外公、敢、不好、不得、忒、怪物、无礼、棍、泼怪、怎生、误

4 全文情感词挖掘

4.1 全文点互信息计算

通过计算了全文的词频发现, 除停用词外, 大部分词的词频数为 1。由于语料规模较小, 通过初步试验, 我们发现这些词给情感分析带来的噪声影响尤为严重, 因此本文剔除了词频为 1 的词, 得到全文待计算情感倾向的词语数量为 2442 个。并利用上阶段生成的三个种子词表: SeedA、SeedB、SeedC, 分别放入全文语料中进行点互信息计算, 而后以式 (3) 来确定待计算情感词 n 的极性:

$$Polar_n = \sum_{i \in pos} \log \frac{p(n|i)}{p(n)} - \sum_{j \in neg} \log \frac{p(n|j)}{p(n)} \quad (3)$$

其中 pos 为正面情感的种子词集, neg 为负面情感的种子词集, 若 $Polar_n$ 大于零, 则该词属于正面情感词的可能性较大, 如小于零则属于负面情感词的可能性较大, $Polar_n$ 绝对值越大则证明极性越高。由于训练集高达 28000 字左右而孙悟空所有的台词只有 10 万字, 为了防止过拟合, 本文将情感种子词规模最小的 SeedC 按照情感词的极性排列, 进行了五等

分，分别作 20%SeedC、40%SeedC、60%SeedC、80%SeedC，并分别作为种子词表进行了点互信息运算，得到全文情感词典。

4.2 确定情感词典规模

本文将 7 个种子词集所计算出的情感词典，按照情感值极性大小排列，并根据情感词典的规模向下取整进行十等分，加上种子词集自身，共得到不同种子词规模、不同互信息规模的 77 个情感词典。

为了确定情感词典的规模，我们从全部语料中随机抽取了 90 句台词，包括 377 个分句作为验证集，并让 4 名专家以人工标注的方式，标注每一个分句的情感极性，正面感情的标为 1，负面感情的标为-1，无明确感情的标为 0，如出现分歧则采取投票的方式。

以 77 个情感词典分别对这 377 个分句进行情感打分。打分的规则为：句中出现正面情感词+1 分，分句中出现负面情感词-1 分，每出现一次否定词则整个分句最后的得分结果乘以-1。最后得到的结果如果大于 0，则记为 1，即正面情感句子，如小于 0 则记为-1，如等于 0 则记为 0。77 个情感词典对句子的打分和人工打分的一致性准确率展示如表 5：

表 5 77 个情感词典的情感判断准确率

	PMI+0%	PMI+10%	PMI+20%	PMI+30%	PMI+40%	PMI+50%	PMI+60%	PMI+70%	PMI+80%	PMI+90%	PMI+100%
20%SeedC	0.5756	0.59151	0.56499	0.53316	0.5305	0.52255	0.52785	0.52785	0.50928	0.48806	0.47215
40%SeedC	0.5809	0.56764	0.5756	0.55438	0.55438	0.54642	0.50663	0.48276	0.47745	0.48011	0.4748
60%SeedC	0.59947	0.5809	0.5756	0.55172	0.55703	0.54642	0.53316	0.51724	0.51724	0.49867	0.49072
80%SeedC	<u>0.61273</u>	0.59416	0.57294	0.54907	0.53846	0.55438	0.54907	0.51194	0.50133	0.50133	0.50663
SeedC	<u>0.61804</u>	<u>0.60477</u>	0.59151	0.56764	0.56499	0.5305	0.53846	0.53581	0.52255	0.51459	0.50663
SeedB	<u>0.62069</u>	0.59151	0.57825	0.54377	0.53581	0.5305	0.54111	0.5305	0.49867	0.49867	0.48806
SeedA	0.51194	0.51194	0.49867	0.45358	0.44032	0.42971	0.41645	0.4244	0.4244	0.41114	0.39788

注：下划线为较好结果，下同

我们发现，互信息并不能起到正面的作用，主要原因是验证集中无情感句占到了 198 个，即总数量的 52.5%，因此投入的情感词数量越少，正确率反而显得越高。

4.3 无情感句的处理与错误距离

我们以人工判定无情感的句子数量 198 为基线，对判定无情感句子数量大于 198 的情感词典，赋予其一个惩罚因子，实际情感词典判断正确的句子数量 N 应以式（4）表示：

$$N = \text{判断正确的数量} - (\text{判断为无情感的句子数量} - 198) \quad (4)$$

判定无情感句子数量不足 198 的不做处理，通过这一步处理，各情感词典实际准确率如表 6 所示：

表 6 77 个情感词典的情感判断准确率（处理后）

	PMI+0%	PMI+10%	PMI+20%	PMI+30%	PMI+40%	PMI+50%	PMI+60%	PMI+70%	PMI+80%	PMI+90%	PMI+100%
20%SeedC	0.26525	0.38727	0.41645	0.44032	0.4748	0.4748	0.49867	0.50663	0.50928	0.48806	0.47215
40%SeedC	0.313	0.35809	0.42971	0.44562	0.49602	0.51194	0.50663	0.48276	0.47745	0.48011	0.4748

60%SeedC	0.34483	0.40053	0.44562	0.45623	0.49602	0.51724	0.53316	0.51724	0.51724	0.49867	0.49072
80%SeedC	0.37401	0.42706	0.44297	0.49867	0.51194	<u>0.55438</u>	<u>0.54907</u>	0.51194	0.50133	0.50133	0.50663
SeedC	0.4695	0.49337	0.51989	0.5305	<u>0.55968</u>	0.5305	<u>0.53846</u>	0.53581	0.52255	0.51459	0.50663
SeedB	0.50928	0.5252	0.5305	0.53581	0.53581	0.5305	<u>0.54111</u>	0.5305	0.49867	0.49867	0.48806
SeedA	0.51194	0.51194	0.49867	0.45358	0.44032	0.42971	0.41645	0.4244	0.4244	0.41114	0.39788

77 个情感词典的正确率变化趋势如图 1:

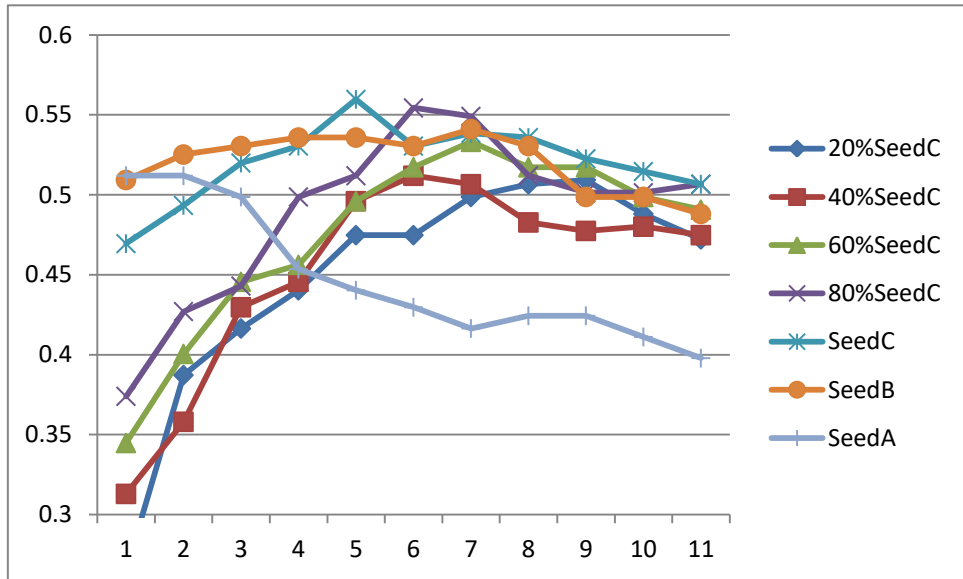


图 1 77 个情感词典的情感判断正确率变化趋势

我们选取了准确率最高的五个情感词典，引入了错误距离的概念，认为情感判断错误的类型不同，其代价也不同，从一定意义上可以反映情感词典的好坏。因此本文对这五个情感词典判定错误的性质进行了分析，把无情感判定为有情感，或将有情感判定为无情感，记为 -1 分，将正面情感判断为负面情感或者负面情感判断为正面情感，记为 -2 分。五个性能最优的情感词典的错误距离合计如表 7:

表 7 情感词典 Top5 的错误距离

	SeedC+40%	SeedC+60%	SeedB+60%	80%SC+50%	80%SC+60%
0→±1	68	87	83	70	73
±1→0	70	53	56	64	62
Contary	26	34	34	34	35
0-1rate	0.18037135	0.23076923	0.22015915	0.18567639	0.19363395
1-0rate	0.18567639	0.14058355	0.14854111	0.16976127	0.16445623
C-rate	0.06896552	0.09018568	0.09018568	0.09018568	0.0928382
Score	<u>-190</u>	-208	-207	-202	-205
Accuracy	<u>55.97%</u>	53.85%	54.11%	55.44%	54.91%

最终选择了 SeedC+40%点互信息的情感词作为《西游记》情感分析的情感词典，该词典错误距离合计-190，情感三分法下准确率为 55.97%。

5 角色情感计分

5.1 角色分类

“真假美猴王”这一事件出现在《西游记》的第五十七到五十八回，由于“真假美猴王”主要描述孙悟空和唐三藏之间的矛盾，孙悟空从第十三回开始保唐三藏取经，至五十六回诛草寇为止，为事件发生之前；第五十九回至第一百回取到真经，《西游记》全篇完结，为事件发生之后。事件发生之前涵盖了 44 章，1720 句台词，事件发生之后涵盖了 42 章，1377 句台词，事件前后语料规模相差较小，具有可比性。

我们根据《西游记》原文，将孙悟空在事件前后所有 3097 句台词的会话对象进行了人工标注，对事件前后孙悟空对话过的所有角色进行分类，由于《西游记》的故事性质，大部分角色仅登场一次，在事件前后均有登场的角色寥寥无几，因此对于取经团队核心，唐三藏、猪八戒、沙悟净三人，每个人单独作为一类，其他角色根据他们的阵营、地位、善恶等，一共分为 10 类，具体的分类如表 8：

表 8 孙悟空在《西游记》中说话对象的分类

	前半部分	后半部分
唐三藏	唐三藏	
猪八戒	猪八戒	
沙悟净	沙悟净	
主要神仙	太白金星、玉皇大帝等	太白金星、哪吒、李天王等
次要神仙	广目天王、雷公电母等	四木禽星、二十八星宿等
佛、菩萨等	如来、观音、文殊菩萨等	宝栋光王、灵吉菩萨、弥勒等
主要妖怪	白骨精、黑风怪、金银角等	黄眉老祖、黄狮精、赛太岁等
次要妖怪	精细鬼、伶俐虫、巴山虎等	小钻风、奔波儿灞灞波儿奔等
土地、水神	通天河老鼋、五庄观土地等	黑松林土地、柳林坡土地等
丁甲护法等	六丁六甲、五方揭谛、四值功曹、护法伽蓝等	
龙王阵营	四海龙王及子孙手下	
平民	刘伯钦、高翠兰、陈澄、陈清等各地百姓	
皇族、官员	宝象国、车迟国、比丘国、朱紫国等各国皇族、官员	

在这些孙悟空对话的对象中，我们舍弃了语料过少、不具有统计意义的与白龙马、后半部分未出现过的与花果山群猴、以及孙悟空的自言自语以及与佛像等物品的对话。

根据上述的分类，本文将语料分为了 26 个子语料。如果孙悟空的对话对象是群体的，如“猪八戒和沙悟净”，则将该台词同时复制到猪八戒、沙悟净两个人的子语料库中。敌对神仙归入妖怪的分类，山贼、草寇、恶兽等归入次要妖怪的分类。一些角色随着故事的发展，所属阵营有所变化，如奎木狼曾经是妖怪，后来归顺成为了四木禽星中的一员。对此我们进行具体分析，将妖怪时期的奎木狼分到主要妖怪的一类，而归顺后则分到次要神仙一类。一些角色精通变化，如打黄风怪时的护法伽蓝，孙悟空以为他是普通的山野老翁，这里就把孙悟空与他对话的台词分到平民一类。

我们把“真假美猴王”事件的三类解读形式化：

假设 1：如果死的是孙悟空，一切都是如来佛祖设局，那么冒名顶替的六耳猕猴对神佛的态度在事件之后相较之前应该具有更高的正面情感倾向，但由于换了人，假孙悟空对取经

团队的另外三个角色的情感几乎不可能与之前的情感有较高相似性。

假设 2: 如果孙悟空没有死,“真假美猴王”是唐僧的修行,是为了缓和唐僧师徒之间的关系,那么孙悟空在事件之后对唐三藏的正面情感倾向应该有明显上升,对神佛等权利阶级应该不会有明显变化。对取经团队的其他主要人员——猪八戒和沙悟净应该基本没有变化。

假设 3: 如果孙悟空没有死,如来安排“真假美猴王”一难是为了消灭孙悟空的“心魔”,让他抛弃主观性和反抗精神,屈服于神权。那么孙悟空在事件之后对唐三藏的正面情感倾向应该有明显上升,对神佛等神权阶级的正面情感倾向应该也有明显的上升,对取经团队的其他主要人员应该基本没有变化。

5.2 情感计算与结果分析

为了验证这三个假设,根据前文所述的句子打分规则,我们对 13 类角色及事件前后 26 个子语料进行了情感打分,结果如下:

表 9 孙悟空对其他角色在事件前后的情感值打分

	总分句数		正面情感		负面情感		正面比率		负面比率		变化率	变化趋势
	前	后	前	后	前	后	前	后	前	后		
唐三藏	1698	1135	493	367	382	182	0.5634	0.6685	0.4366	0.3315	0.1051	↑
猪八戒	1554	1309	399	306	394	294	0.5032	0.51	0.4968	0.49	0.0068	≈
沙悟净	424	321	135	100	82	60	0.6221	0.625	0.3779	0.375	0.0029	≈
主要神仙	274	293	92	134	59	33	0.6093	0.8024	0.3907	0.1976	0.1931	↑
次要神仙	313	162	96	64	38	21	0.7164	0.7529	0.2836	0.2471	0.0365	↑
佛菩萨等	587	325	171	116	101	55	0.6287	0.6784	0.3713	0.3216	0.0497	↑
主要妖怪	851	738	186	159	247	206	0.4296	0.4356	0.5704	0.5644	0.0061	≈
次要妖怪	510	271	125	47	134	70	0.4826	0.4017	0.5174	0.5983	-0.0809	↓
土地水神	160	132	26	37	40	27	0.3939	0.5781	0.6061	0.4219	0.1842	↑
丁甲护卫	79	142	25	61	12	14	0.6757	0.8133	0.3243	0.1867	0.1377	↑
龙王阵营	128	30	39	14	20	1	0.661	0.9333	0.339	0.0667	0.2723	↑
平民	613	470	198	134	87	81	0.6947	0.6233	0.3053	0.3767	-0.0715	↓
皇族官员	378	729	110	222	65	139	0.6286	0.615	0.3714	0.385	-0.0136	≈/↓
土地三类	367	304	90	112	72	42	0.5556	0.7273	0.4444	0.2727	0.1717	↑

注:判定标准以 0.01 为界

实验结果表明,孙悟空在事件前后对猪八戒、沙悟净两位师弟的态度几乎没有变化,基本可以判定取经的还是孙悟空本人而非六耳猕猴,而对主要妖怪和皇族、官员等的态度也几乎没有变化可以佐证这一点,因此拒绝假设 1。事件之后,孙悟空对唐三藏的正面情感有明显上升,但单靠对唐三藏的情感变化尚不能区别假设 2 与假设 3。孙悟空在“真假美猴王”之后,对于神佛群体的态度均明显转好,对主要神仙、次要神仙、佛等正面情感上升幅度很高,尤其是对主要神仙和佛的上升幅度非常大,对土地神、丁甲伽蓝、四海龙王等这些原本孙悟空看不起的神仙,正面情感也有了明显的提升,尤其是对土地神,从“一生好吃没钱酒,偏打老年人”变为事件后的正面情感超过了负面情感,其情感发生了质的变化。由于后半部

分孙悟空对四海龙王等人的台词较少,不具有统计意义。我们尝试将土地神、丁甲伽蓝、四海龙王这三个相似群体的语料合并,其结果依然是正面情感有明显上升。

我们可以明显发现孙悟空对妖怪的负面情感句子数量较多,对师父、神佛等自己人的正面情感句子较多,对八戒的负面情感句子比例明显大于其他群体,这也基本符合对《西游记》中人物关系的认知,证明了本文研究的情感词典在小说人物分析上的有效性。对平民这个群体正面情感倾向有下降的趋势,这可能是由于《西游记》全文中平民角色过多,且往往是故事中的次要人物,因而该语料混杂度过高、噪声过多所造成的,在此不对该群体进行细致讨论。

从上述分析基本可以验证,“真假美猴王”这一事件如来佛祖虽然没有杀死孙悟空,但却消灭了孙悟空的反抗精神,诛“心魔”是一个同化过程^[8]。事件之后,孙悟空对神佛群体的正面情感倾向明显上升,性格逐渐趋向扁平单一,从向往自由渐渐走向了屈服体制的归化之路^[2]。通过情感分析,我们认为,假设3更适合作为“真假美猴王”事件的正确解读。

6 总结与展望

本文从情感分析的视角对传统名著《西游记》中“真假美猴王”事件进行了解读,通过评测现有分词系统,提出了明清白话的分词方案。借助作者吴承恩在《西游记》中遗留下的情感信息确定了情感种子词,并以点互信息的方式生成并验证了适合分析孙悟空这个角色的情感词典。通过对事件前后孙悟空对其他角色情感变化的分析,得出了真的孙悟空并没有死,而是象征着反抗精神的“心魔”被消灭的结论。本文作为一种新的尝试,验证了情感分析技术对文学研究和文学作品角色分析的可行性。

明清小说是中国小说历史上的巅峰时期,四大名著均在其列,明清白话的分词问题值得进一步的探讨,借助自然语言处理的方法对明清白话小说中人物进行情感分析也值得进一步研究。文学作品的分析的语料量级往往比较小,因此本文主要使用的是基于规则的方法,情感分析的准确率较现代汉语语料而言,尚有空间可以改进,因此在日后的研究中,我们将尝试使用更大规模的语料,运用机器学习、神经网络等方法,对明清小说及明清白话进行更加深入进行讨论,进一步提高模型的准确率。

参考文献

- [1] (明)吴承恩.西游记[M].北京:人民文学出版社,2011
- [2] 李家宝,王利军.文化视阈中孙悟空反叛与归化再阐释[J].江汉论坛,2018,3:82-85
- [3] 杨扬.象征何盛,悲喜何情,心棒何重?——孙悟空形象再续议[J].东南大学学报(哲学社会科学版),2018,20(1):130-138
- [4] 宁稼雨.孙悟空叛逆性格的神话原型与文化解读[J].文艺研究,2008,10:59-66
- [5] 关四平.从紧箍咒管窥孙悟空的内心世界[J].明清小说研究,2009,2:77-88
- [6] 张振国.《西游记》“六耳猕猴”意象的文化与心理学阐释[J].东南大学学报(哲学社会科学版),2016,18(1):119-123
- [7] 吴光正.《西游记》的宗教叙事与孙悟空的三种身份[J].学术交流,2007,11:140-145

- [8] 郭明友.论孙悟空形象悲剧意蕴的广度与深度[J].名作欣赏, 2010, 17: 44-46
- [9] 赵妍妍, 秦兵, 刘挺.文本情感分析[J].软件学报, 2010, 08: 1834-1848
- [10] 杨立公, 朱俭, 汤世平.文本情感分析综述[J].计算机应用, 2013, 06: 1574-1578
- [11] 李泽魁, 赵妍妍, 秦兵等.中文微博情感倾向性分析特征工程[J].山西大学学报(自然科学版), 2014, 04: 570-578
- [12] 任巨伟, 杨亮, 林鸿飞等.基于情感常识的微博事件公众情感趋势预测[J].中文信息学报, 2017, 02: 169-178
- [13] 赵妍妍, 秦兵, 刘挺等.大规模情感词典的构建及其在情感分类中的应用[J].中文信息学报, 2017, 02: 187-193
- [14] 刘楠.面向微博短文本的情感分析研究[D].武汉: 武汉大学, 2013
- [15] 李勇敢, 周学广, 孙艳等.中文微博情感分析研究与实现[J].软件学报, 2017, 12: 3183-3205
- [16] 叶强, 张紫琼, 罗振雄.面向互联网评论情感分析的中文主观性自动判别方法研究[J].信息系统学报, 2007, 01: 79-91
- [17] 刘鸿宇, 赵妍妍, 秦兵等.评价对象抽取及其倾向性分析[J].中文信息学报, 2010, 01: 84-88
- [18] 王祖辉, 姜维, 李一军.在线评论情感分析中固定搭配特征提取方法研究[J].管理工程学报, 2014, 04: 180-186
- [19] 徐琳宏, 林鸿飞.认知视角下的文本情感计算[J].计算机科学, 2010, 12: 182-185
- [20] 林斌.基于语义技术的中文信息情感分析方法研究[D].哈尔滨: 哈尔滨工业大学, 2006
- [21] 任巨伟, 杨亮, 林鸿飞.情感图式构造及其在文本情感计算中的应用[J].江西师范大学学报(自然科学版), 2013, 02: 130-135
- [22] 赵妍妍, 秦兵, 车万翔等.基于句法路径的情感评价单元识别[J].软件学报, 2011, 05: 887-898
- [23] 王科, 夏睿.情感词典自动构建方法综述[J].自动化学报, 2016, 04: 495-511
- [24] 阳爱民, 林江豪, 周咏梅.中文文本情感词典构建方法[J].计算机科学与探索, 2013, 11: 1033-1039
- [25] 郗亚辉.产品评论中领域情感词典的构建[J].中文信息学报, 2016, 06: 136-144
- [26] 徐琳宏, 林鸿飞, 赵晶.情感语料库的构建和分析[J].中文信息学报, 2008, 01: 116-122
- [27] 陈建美, 林鸿飞, 杨志豪.基于语法的情感词汇自动获取[J].智能系统学报, 2009, 02: 100-106
- [28] Snover M, Dorr B, Schwartz R, et al. A study of translation edit rate with targeted human annotation[J]. Machine Translation Workshop North Bethesda Md, 2006(1):223-231.
- [29] 杨娟.《西游记》神魔描写词语研究[D].重庆: 西南大学, 2016
- [30] 杜贵晨.孙悟空对妖精习称“外公”说之辨误与新解[J].河北学刊, 2015, 02: 95-99
- [31] 周福雄, 李美林.《西游记》中孙悟空使用的骂詈性称谓语研究[J].现代语文(语言研究版), 2016, 11: 87-91
- [32] 张继平.论汉语中的否定概念[J].逻辑与语言学习, 1989, 01: 8-10
- [33] Li F, Pan S, Jin O, et al. Cross-Domain Co-Extraction of Sentiment and Topic Lexicons[C]. Proceedings of the 50th ACL, 2012:410-419.