

文章编号: 2095-2163(2022)12-0062-08

中图分类号: TP391

文献标志码: A

基于时间序列的方面级网络舆情情感演化模型

董光文, 袁健

(上海理工大学 光电信息与计算机工程学院, 上海 200093)

摘要: 针对网络舆情情感分析主题词抽取不精确和文本静态化分析问题, 论文提出了一种基于时间序列的方面级网络舆情动态情感演化模型 ARMA-ALEE。通过方面级情感分类模型获取方面词和情感极性值, 并对方面词使用过滤算法优化, 再通过困惑度和 JS 散度确定最终方面词个数, 进一步地还基于 ARMA 时间序列模型对方面词、方面词强度和方面词相关性的 ARMA-ALEE 模型动态地进行网络舆情情感演化分析。实验表明, 该模型的情感演化研究取得了较好的结果。

关键词: 情感分析; 主题提取; 情感演化; ARMA 时间序列

Aspect-level network public opinion sentiment evolution model based on time series

DONG Guangwen, YUAN Jian

(School of Optical-Electrical and Computer Engineering, University of Shanghai for Science and Technology, Shanghai 200093, China)

[Abstract] Aiming at the problems of inaccurate subject word extraction and static analysis of texts in sentiment analysis of network public opinion, this paper proposes an aspect-level dynamic sentiment evolution model ARMA-ALEE of network public opinion based on time series. The aspect words and sentiment polarity values are obtained through the aspect-level sentiment classification model, and the filter algorithm is used to optimize the aspect words, after that the final number of aspect words is determined by the perplexity and JS divergence. Thereafter the evolution of network public opinion sentiment is dynamically analyzed for aspect words, the strength of aspect words and ARMA-ALEE model related to aspect words based on ARMA. Experiments show that the emotional evolution research of this model has achieved good results.

[Key words] sentiment analysis; topic extraction; emotional evolution; ARMA time series

0 引言

随着互联网的飞速发展, 网络社交平台已经逐渐成为新兴的舆论载体。当面对一些热点新闻或者突发事件时, 对网络社交平台中的相关言论进行有效分析, 实时了解当前热点或新闻事件的舆情演变发展趋势, 并在必要时采取行动施以重点监控, 保障网络舆情平稳发展, 从而为人们营造一个积极健康的良好网络环境。

常见的情感演化研究大多偏向于对静态文本的分析或以热门话题为基础进行主题词辨识, 同时也要有海量数据作为支持, 而舆情事件在热点初期却较难察觉, 若能对舆论情感进行实时动态的细粒度分析就可以准确把握舆情的动态和趋势, 对舆情的预测和调整具有重要的实用价值。

1 相关研究

1.1 情感分析研究

近年来, 在情感分析领域, 国内外学者已取得了可观成果, 研究上也主要集中在情感极性分析、多分类情感分析和方面级情感分类三个方向。

常用的情感分析技术研发初期就是以情感词典为主的研究方法, 这种方法需要依赖人工去构建词典, 并要不断地扩充词典, 情感分析效率并不高^[1]。后来学者们分别使用有监督和无监督的机器学习方法来进行情感分析研究, 在一些模型上取得了较好的效果, 但在此过程中也需进行特征工程的构建, 往往消耗不少人力^[2-3]。目前, 深度学习技术已经成为该领域主流的研究方法。特别是循环神经网络(RNN)、卷积神经网络(CNN)以及注意力机制在情感分析中的组合使用, 使情感分析技术已然日臻成

基金项目: 国家自然科学基金(61775139)。

作者简介: 董光文(1995-), 男, 硕士研究生, 主要研究方向: 自然语言处理、深度学习; 袁健(1971-), 女, 博士, 副教授, 硕士生导师, 主要研究方向: 自然语言处理、数据挖掘、深度学习等。

通讯作者: 袁健 Email: yuanjianwq@163.com

收稿日期: 2022-03-27

熟。Lv 等人^[4]提出一种上下文和方面记忆网络 (CAMN) 方法来解决方面级情感分析问题,引入了深度记忆网络、双向长短期记忆网络和多重注意力机制,能够更好地捕捉文本中的情感特征,获得文本方面级情感分析结果。

1.2 情感演化研究

情感演变主要是对含有情感的主观信息进行分析,并从情感的态度和角度对情感在时间中的演变进行分析。面对各类突发情况下不断涌现的网络舆论热点事件,国内外的学术界从多个角度对其情感演变进行了全方位的分析 and 探讨。

在网络舆情情感演化分析方面,邢云飞等人^[5]以“江歌案”为例,从情感的极端和情感的强弱入手,探讨了其演变及变化规律。钱进宝^[6]以“穹顶之下”为例,建立以词汇相关性为基础的文字情感矢量模型,在 K-medoids 中加入历史代价函数,可以对网络上的热门事件进行动态的情感演变分析,从而避免了以往仅限于对静止的数据进行分析的不足。戴杏云等人^[7]在统计用户关系、用户影响力等指标的基础上,建立基于网络的动态情感图的分析模型,从而为控制和指导社会网络舆论提供了基础。张柳等人^[8]以“学术不端”为例,从舆论发展的角度来分析情感演变的规律,运用了词云图和情感知识图谱,分别揭示了爆发期、蔓延期和衰退期用户使用高频率词和情感分配的演变规律。

综合前文论述可知,目前网络舆情情感演化的研究大多着重于舆情主题的挖掘、传播的特征和生命周期模型等方向展开研究。研究时则需要大量的数据做支撑,也就是只有当舆论成为热点时才能更好地选择准确的主题、抽取出特征或划分生命周期,而当某舆论处于发展阶段的初期时却较难被发现,这将导致舆论分析的效果欠佳。通常来说,人们对舆情的情感往往都是动态的,舆情情感的波动也会和某突发话题的发展趋势密切相关,若不考虑时间发展的维度,对网民们的情感动态演化很难做出有效判断。

基于此,本文引进了时间序列 ARMA 模型,并在方面级情感分析 CAMN 模型^[4]的基础上,提出了基于时间序列的方面级网络舆情情感演化模型 (Aspect - level network public opinion sentiment evolution model based on time series, ARMA-ALEE)。该模型的创新点如下:

(1) 在方面级情感分析基础上获取方面词和情感极性,对方面词使用过滤和优化算法以提高其精

度,并在方面词的基础上进一步提取主题词来做过滤优化后分析,进行更加细粒度的舆情演化分析。

(2) 提出了 ARMA-ALEE 情感动态演化模型,引入 ARMA 时间序列模型,基于 ARMA 对方面词、方面词强度和情感强度动态倾向性训练,实现网络舆情动态情感演化分析。

2 ARMA-ALEE 情感演化模型

ARMA-ALEE 模型的整体结构如图 1 所示,该模型先对实时文本数据集进行预处理,并按时间顺序进行划分,然后利用 CAMN 模型对每个语料集进行处理,对处理后的结果进行方面词优化、情感强度计算和方面词强度的算法实现,接着将基于 ARMA 时间序列模型实现 ARMA-ALEE 情感演化算法,最后进行情感演化分析及可视化。

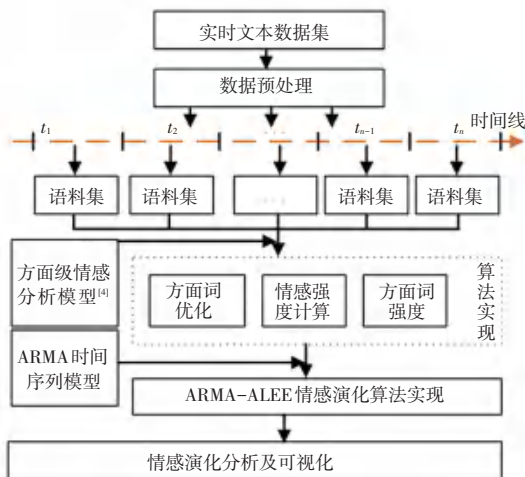


图 1 ARMA-ALEE 模型整体框架图

Fig. 1 Overall framework diagram of ARMA-ALEE model

2.1 方面词的优化

使用主题词对舆情数据进行演化分析时,会有许多与数据集关联不大且没有实用价值的主题词,为了避免对网络舆情情感演化的影响,不少学者对无用主题的过滤也做了一定的研究^[9]。本文将对文本数据中每个句子的方面级进行研究,确定每个句子的方面属性,采用方面属性代替主题词属性对网络舆情情感演化进行分析。

由于文本数据集中的长度参差不齐(尤其是针对微博),这就导致方面分类有时不精确、或者方面分类过多等问题,对网络舆情情感演化造成了一定的影响,本文将对文本数据集中获取的所有方面词进行优化操作。

2.1.1 方面词过滤框架

以时间为演化发展线索,利用 CAMN 模型^[4]获

取到每个时间段内的方面词,对提取的方面词进行过滤处理,提高方面词对网络舆情中情感演化分析的效果。本文对方面词过滤的流程如图2所示。由图2可看到,首先把文本数据集以某个时间段为间隔划分开,将对应的数据集分配到相应的时间段内,基于CAMN模型^[4]获得每个时间段内的方面主题词及其个数。接下来,对方面词进行过滤,剔除一些没有价值的方面词,以防止在相邻时间段内对相关主题的辨识和判断。最后,通过算法对经过筛选后的方面词确定最优方面词个数和邻近时间段内方面词之间的相关性。

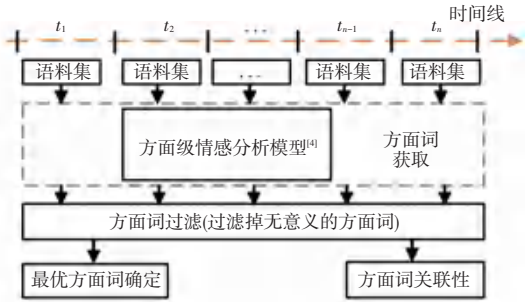


图2 方面词过滤图

Fig. 2 Aspect words filtering diagram

2.1.2 时间段内方面词过滤

对于每个时间段内的方面词数,一些方面词在文本中出现的概率极低或毫无相关,不但加大运算的难度,还将导致与无关话题之间的不必要联系,从而对方面主题进行演化的正确分析产生不利的作用。如果一个方面词在数据集中出现的比例越高,那么该方面词在某个时间片段内的重要程度越大。相反,如果某个方面词在数据集中出现的比例越低,通常就可把该方面词视为不重要的词语,这种出现次数较少的方面词也不会形成网络舆论。本文对方面词进行过滤筛选分为2个阶段,具体过程如下:

(1)基于方面词分布的边缘方面词辨识和筛选。利用CAMN模型^[4]获取到的方面词在每个时间段内分布概率差异较大,在同一时间段内发生频率较高的方面词,则是该时间段内较为核心的热点方面词,也是研究情感演化的关键因素。相反,出现概率较小的方面词,往往被边缘化或者说是毫无意义的,同时也会对情感演化的分析结果产生影响。因此,针对使用CAMN模型^[4]所获得的某个时间段内的方面词数,依据其在数据集中的分布情况来设定临界点,设定方式如下:在数据集中,计算每个方面词A累积的概率P,将一个时间段内的数据集总量N进行标准化处理,得出其在数据集中的权重值

W,将权重W从大到小依次排序,并选取其平均值为筛选阈值,此处需用到的数学公式为:

$$p(w_i | z) = \frac{1}{N} \quad (1)$$

$$P_i = \sum_{i=1}^n p(w_i | z) \quad (2)$$

$$W_i = \frac{P_i}{\sum_{j=1}^n P_j} \quad (3)$$

其中,N表示数据集中所有方面词A的总数; $p(w_i | z)$ 表示单独一个方面词在N中所占的比例; P_i 表示单个方面词的累加概率之和; W_i 表示一个方面词在数据集中的权重。

(2)基于方面词分布的无用方面词的辨识和筛选。经过上一步筛选后,把方面词汇聚在一起进行概率的分布,这些方面词之间的关联意义通常是用来描绘某一话题相关意义或者发展趋势的。假如某方面词和大多数方面词毫无关系且也不具备发展联系,就会被视为无意义并筛选掉。

利用信息熵法对表达对象的方面词倾向性进行衡量。信息熵是一种信息不稳定的度量方法,一个方面词可以看作一系列随机的方面词,当其在数据集中出现的可能性越大时,其信息熵值越低,也就越能突出所要表达的内涵。对经过上一步筛选出的每个方面词进行信息熵计算,具体可由如下公式计算求得:

$$Entropy(A) = -K \sum_{j=1}^m P_j \ln(P_j) \quad (4)$$

其中,Entropy(A)表示方面词A的信息熵; P_j 表示在方面词A中第j个词语出现的概率;K表示一般的常数;m表示方面词A中所包含的词语的个数。

2.1.3 相邻时间段内方面词相关性

在邻近时间段内的方面词中,仅有相互关联的方面词之间才可能会存在相互演化的关系。方面词的相似性是用来衡量方面词之间的相似程度,使用“方面词—单词”概率分布来计算方面词之间的相似性。

在相似度计算过程中,余弦相似度是用2个矢量夹角的余弦值作为衡量矢量相似性的指标。经过滤后得到的每个时间段内的方面词是由一系列的词语组成,而不是以传统的字词矢量来表达,所以相对于余弦相似性,概率分布的距离公式在衡量方面词之间的相似度时略有优势。KL的离散,即KL距离或

者相关熵是对同一时间点上的 2 种概率分布的重要度量,能够表示出 2 个方面词之间的差异情况。当 2 个方面词随机分布相同时, KL 距离为 0, 而随着 2 个方面词随机分配差异的加大, KL 距离也随之增大。推导的计算公式可写为:

$$KL(p \parallel q) = \sum_{i=1}^n p(x_i) \log \frac{p(x_i)}{q(x_i)} \quad (5)$$

其中, $KL(p \parallel q)$ 表示 2 个方面词概率分布为 p 和 q 的 KL 距离; x_i 表示概率分布为 p 和 q 的第 i 个方面词; n 表示 2 个概率分布为 p 和 q 的方面词的总个数。

由于 KL 散度是非对称化的, 故根据其理论给出另一种变种 JS 散度, 从而通过 JS 散度将 KL 散度转化为真实的距离度量, 如式(6)所示:

$$JS(p \parallel q) = \frac{1}{2} KL(p \parallel \frac{p+q}{2}) + \frac{1}{2} KL(q \parallel \frac{p+q}{2}) \quad (6)$$

JS 散度的扩散系数一般为 0 至 1, JS 散度的数值越低, 则表示两者的相似度越高。

2.1.4 最优方面词的确定

经过方面词过滤算法处理后, 每个时间段内方面词数量是不确定的, 方面词的个数会影响对情感演化分析的效果。如果同一时间段内方面词个数太多, 将会导致方面主题过于分散、且舆情方向过多, 不能突出核心的演化方向。相反, 如果同一时间段内方面词个数过少, 舆情分析则容易向一个方向发展, 就可能会忽略掉一些潜在的方面主题方向。

困惑度是衡量一个语言模型好坏的指标, 困惑度越低, 说明该模型具有较好的泛化能力^[10]。给出的数学定义可表示为:

$$Perplexity(D) = \exp \left(\frac{\sum_{d=1}^M \log p(w_d)}{\sum_{d=1}^M N_d} \right) \quad (7)$$

其中, $Perplexity(D)$ 表示该模型困惑度的大小; D 表示数据集中的测试集; M 表示测试集中含有的时间段数; N_d 表示第 d 个时间段中包含的方面词数量; $p(w_d)$ 表示第 d 个时间段中所含方面词分布的概率。

当潜在方面词的数量增多时, 该模型的困惑度也就越低, 但是往往会会有一个拐点, 表明该模型的泛化能力得到了显著的改善, 从而可以通过这个拐过来估算方面词最佳数量。然而, 仅靠困惑度来判断方面词数量通常不准确, 还需要综合考虑其它的因

素。

主题平均相似度是一种度量各个主题词之间相似度的平均差异程度的指标^[11], 通常使用 JS 散度对其进行较好的衡量, 使用主题相似度来计算方面词的平均相似度, 计算方法可由式(8)表示为:

$$avgsim(T_i, T_j) = \frac{\sum_{i=1}^{k-1} \sum_{j=i+1}^k JS(T_i \parallel T_j)}{K \times (K-1) \div 2} \quad (8)$$

其中, $avgsim(T_i, T_j)$ 表示所有方面词之间的方面平均相似度; T_i 和 T_j 分别表示不同的 2 个方面词; $JS(T_i \parallel T_j)$ 表示 2 个方面词之间的 JS 散度。

JS 离散指当 2 个方面词的数值较大时, 则表示相似性越大。当方面词数目增多时, 方面词之间的相似程度总体上呈现上升的态势, 但同时也会出现一个拐点。

当方面词数量增加时, 方面词间的平均相似性会变大, 而困惑度将会呈现变小的趋势, 但也都会有显著的转折点, 将二者结合起来确定最优的主题词个数, 从而使模型的泛化能力得到显著的提高。

2.2 网络舆情情感演化实现方法

2.2.1 ARMA 时间序列模型

ARMA 时间序列模型也称为自回归移动平均模型, 包括 2 个方面: 自回归模型 (AR) 和移动平均模型 (MA)^[12]。定义时间序列 $t = (t_1, t_2, \dots, t_n)$, 假设在一定的时间内, 一个特定时间点的数值与前面的 p 个序列的数值和前面 q 个输入的随机干扰相关, 从而可以对接下来的时间点进行预测。假定 t_i 被前面 p 个时间顺序的数字所影响的自回归过程的计算方法具体见如下:

$$t_i = \eta_1 t_{i-1} + \eta_2 t_{i-2} + \dots + \eta_p t_{i-p} + e_i \quad (9)$$

其中, $\eta_1, \eta_2, \dots, \eta_p$ 表示自回归系数, e_i 表示误差项。

在不同的时序上, 误差项 e_i 之间存在着相关性, 其移动的平均值的计算方法如式(10)所示:

$$e_i = \mu_1 \varepsilon_{i-1} + \mu_2 \varepsilon_{i-2} + \dots + \mu_q \varepsilon_{i-q} + \varepsilon_i \quad (10)$$

其中, $\mu_1, \mu_2, \dots, \mu_q$ 表示移动的平均系数, ε_i 表示白噪声序列。

联立式(9)~(10)得到 ARMA 模型的计算公式, 即:

$$t_i = \sum_{m=1}^p \eta_m t_{i-m} + \sum_{n=1}^q \mu_n \varepsilon_{i-n} + \varepsilon_i \quad (11)$$

进一步地, 对 ARMA 动态预测模型的流程步骤可做阐释分述如下:

Step 1 首先对时间序列中的每个数值 t_i 进行

均值化处理,然后对数值 t_i 进行稳定性检测。如果不稳定,就进行差分计算,直至差分后的数据平滑为止。

Step 2 对稳定后的数据进行白噪声测试,当检测到平滑的白噪声数据时,利用自相关函数(ACF)和偏相关函数(PACF)求出 ARMA 的阶 p 、 q ,并利用 StatsModels 包来拟合 ARMA(p, q),接着对不同组合(p, q)来计算最小信息准则 AIC 的值,接下来选择 AIC(p, q) 值中的最小阶数作为值(p, q)的估计。

Step 3 利用最小二乘方法对所建立的模型进行求解,得到未知参数 η 和 μ ,对于 $i+1$ 时刻的动态预测计算方法见式(12):

$$t'_{i+1} = \eta_1 t'_i + \dots + \eta_p t'_{i+1} + e_i - \mu_1 \varepsilon_i - \dots - \mu_q \varepsilon_{i+1-q} \quad (12)$$

其中, t'_i 表示零均值时间序列。

2.2.2 方面词强度计算

研究方面词强度在不同时间窗口内的发展趋势,能够反映出一个方面词的稳定性,能够把握一个方面主题的发展方向。用当前时间段内该方面词在所有方面词中所占的比例来表示,计算方法见式(13):

$$AS(A_i) = \frac{P(A_i)}{\sum_{j=1}^m P(A_j)} \quad (13)$$

其中, $AS(A_i)$ 表示时间段内方面词 A_i 的强度; $P(A_i)$ 表示一个方面词在时间段内出现的概率; m 表示一个时间段内方面词优化后的总数量。

2.2.3 情感强度计算

在进行方面级情感分类时,会根据每个方面词的情感极性值分成不同的类别,本文在进行情感强度计算时,选取时间段内的方面词并根据方面词的极性值进行累加求和得到该方面词的情感强度,计算方法的数学公式可表示为:

$$EI(A_i) = \sum_{j=1}^m PV(A_i)_j \quad (14)$$

其中, $EI(A_i)$ 表示方面词 A_i 的情感强度; $PV(A_i)$ 表示一个方面词的情感极性值; m 表示一个时间段内该方面词出现的次数。

2.2.4 情感演化算法实现

针对网上舆论活动中的文本数据进行动态方面级情感演化分析,本文给出了一种动态方面级情感演化分析模型 ARMA-ALEE。ARMA-ALEE 情感动态演化模型的具体工作流程见如下。

输入 网络舆情文本数据集

输出 不同时间段内情感动态演化分析结果

Step 1 对文本进行预处理。

Step 2 舆情演化时间段划分。对数据集根据时间序列上的排序归类进行时间段的划分,本文以时间为单位把对应的数据集划分到一个时间段(根据具体情况以不同单位划分时间)。

Step 3 使用 CAMN 模型^[4]对划分的每个时间段内的数据进行模型训练,获取每个时间段内的数据集所对应的方面词、情感极性值,并用标签进行标记。

Step 4 方面词过滤和确定。根据式(1)~(4)方法筛选掉无用的方面词。根据式(7)~(8)确定最终的方面词个数。

Step 5 方面词强度计算。根据式(13)求出每个方面词对应的方面词强度。

Step 6 方面词相似度计算。根据式(5)~(6)求出方面词之间的相似度。

Step 7 情感强度计算。根据式(14)求出方面词的情感强度。

Step 8 以时间为线索统计数据集特征。将上述步骤中计算得出的方面词强度、情感极性值和方面词相似度按时序分段并合并成文本时间序列集合,数学表示形式如下:

$$DT = (((AS_1, EI_1, JS_1), t_1), ((AS_2, EI_2, JS_2), t_2), \dots, ((AS_n, EI_n, JS_n), t_n)) \quad (15)$$

Step 9 网络舆情情感演化动态倾向性训练。把 DT 作为训练集输入到改进的 ARMA-ALEE 模型中进行迭代训练,把损失函数降低到最小时得到最佳鲁棒性模型。将 15% 的训练集分割成验证集进行校验,然后在校验集上重复校验,获得最优化的超参量组合。接着将验证集和测试集结合,利用 5 折交叉验证方法选取最佳模式,对 $i+1$ 时段下的情感趋势进行动态获取,以 t'_{i+1} 的值作为该时间段内的舆情情感分析结果。

Step 10 网络舆论情感演化分析。在划分的时间段内,根据每个时间段内情感倾向性结果得到每个时间段内方面词和情感极性的变化趋势,进而分析网络舆情情感演化的趋势。

3 实验

3.1 实验数据集

本文以微博“北京冬奥会”为例,根据“北京冬奥会”关键字百度指数数据显示,这一舆情热点事件集

中在 2022 年 2 月 1 日至 2022 年 2 月 25 日, 本文爬取这 25 天内的数据内容进行舆情情感分析。首先对爬取的数据文本进行预处理, 对数据进行清洗, 筛选掉一些无用的文本数据, 最终获取到 159 332 条博文数

据和 67 213 672 条评论文本数据。

3.2 网络舆情情感演化分析

把数据集输入到模型中进行训练, 得到“北京冬奥会”情感演化过程图, 如图 3 所示。

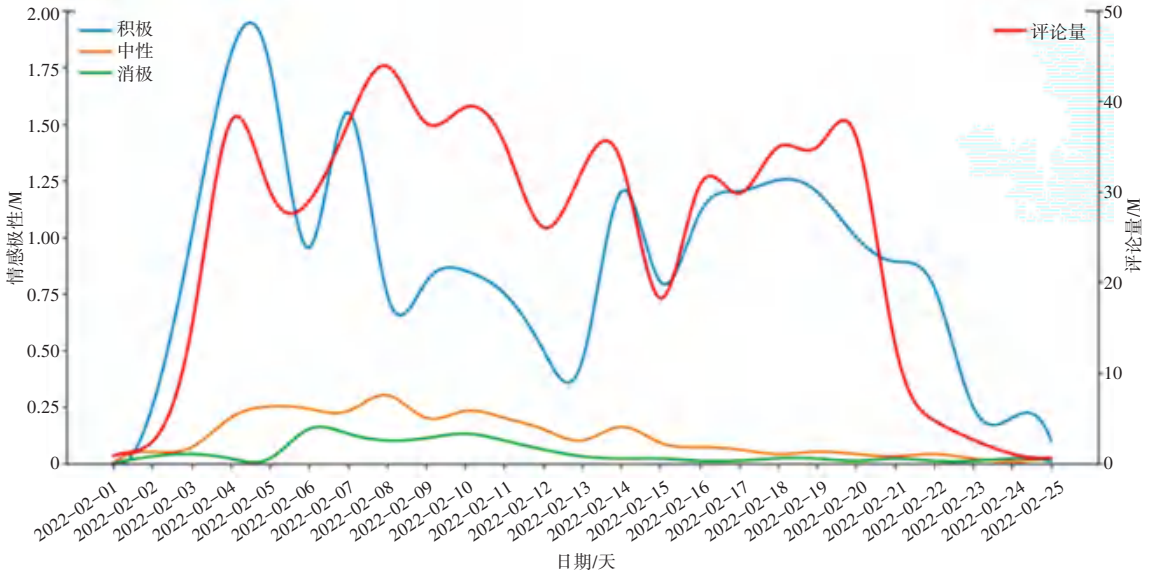


图 3 网络舆情情感演化过程图

Fig. 3 The evolution process of network public opinion

由于数据集过于庞大, 本文选取 2 月 4 日至 2 月 8 日爆发期的一段时间进行详细分析。根据情感演化方面词强度计算结果绘制出词云图, 如图 4 所示, 本文根据图 4 选取舆情热度较高的方面词“冰墩墩”进行分析, 并获取该方面词的相关事件分析表, 详见表 1。

表 1 “冰墩墩”情感演化事件分析表

Tab. 1 “Bingdundun” emotional evolution event analysis table

时间	事件
2022.2.4-	冰墩墩, 冬奥会吉祥物, 深受大众欢迎, 社交网站关于冰墩墩的话题热情不减
...	...

3.2.1 基于方面词的词频和主题的可视化分析

在数据集中, 根据标签标注的方面词找到对应的文本句子, 并将这些句子进行主题词提取优化处理(主题词优化方法同方面词优化方法), 这样就可以清晰地看到该方面词所对应的主题词, 进而便于进行细粒度的舆情分析。

词云图是文本数据集经过分词和去停用词等操作后, 再进行词汇频率的统计, 并对频率高的词汇在大小和颜色方面进行视觉上的对比, 直观表现出文本数据集中所要表达的大致核心意思。本文对热度较高的方面词“冰墩墩”绘制出词云图, 如图 5 所示, 由图 5 便可直观得出该时间段内引起网友们对方面词“冰墩墩”高度关注的高频词汇。采用主题提取模型对该方面词内的主题词进行提取并对其优化处理, 对应的主题提取表见表 2。



图 4 舆情演化方面词词云图

Fig. 4 Word cloud map of public opinion evolution



图5 方面词“冰墩墩”词云图

Fig. 5 Word cloud map of the aspect word “Bingdundun”

表2 方面词“冰墩墩”主题词提取表
Tab. 2 The subject word extraction table of the aspect word “Bingdundun”

序号	主题词
Topic 1	冰墩墩 冬奥会 奥林匹克 中国队 开幕式 喜迎 圆满 祝贺 成功
Topic 2	冬奥会 中国式 北京 底蕴 智慧 科技 战绩 盛名 圆满
Topic 3	冰墩墩 运动员 志愿者 健儿 中国队 金牌 冰壶 奖牌 体魄
Topic 4	冰墩墩 雪容融 吉祥物 冬奥会 憨态可掬 纪念品 宝贝 强健 疫情
Topic 5	冰墩墩 冬奥会 盛会 挥舞 圆梦 迎接 实力 实现 温暖
Topic 6	冰墩墩 可爱 表情 好运 纪念 展览 徽章 头像 陪伴
Topic 7	冰墩墩 发货 抢购 钥匙扣 抽奖 想要 缺货 难买 呜呜
...	...

3.2.2 方面词情感演化过程分析

方面词“冰墩墩”情感演化过程图如图6所示。

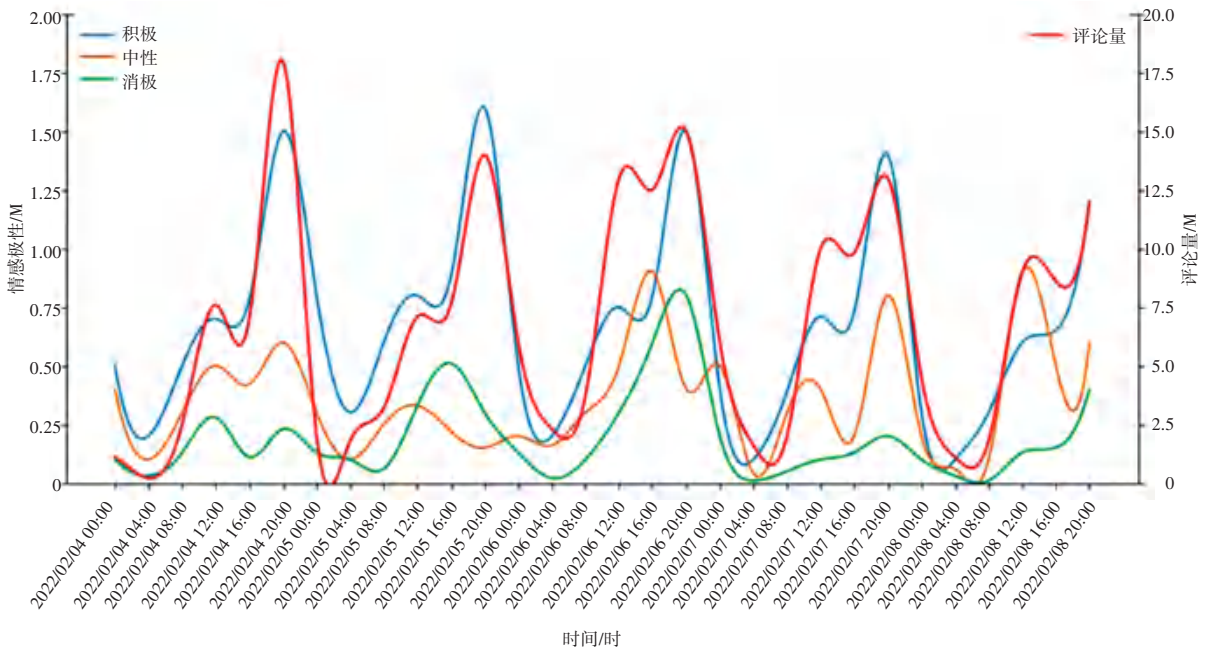


图6 方面词“冰墩墩”情感演化过程图

Fig. 6 The emotional evolution process diagram of the aspect word “Bingdundun”

从图6中红线评论量曲线可以看出,网上关于冰墩墩的言论在2月4日之前也有一定的数量,随着冬奥会开幕式的来临,2月4日人气暴涨,达到了顶峰,在之后几天内的连续传播,并连续出现了高峰,仍有大量网民对此表达自己情感想法。对比情感极性中的积极、中性和消极三条曲线,人们对冰墩墩的评论始终是以积极的态度为主,只有少部分会产生中性和消极的情绪,这也说明了网民们对冰墩墩吉祥物的喜爱之情。

由图6分析可知,蓝色积极情感极性曲线要远

远高于橙色(中性)和绿色(消极)的曲线,分析其中的原因,在图5中,可以看到“可爱”、“吉祥物”、“喜欢”等一些高频的词,体现出人们对“冰墩墩”所表达的积极情感,结合表2,在Topic 1~6这些主题词中,从“表达冬奥会开幕式的举办圆满成功、到表达对冬奥会上运动健儿的骄傲赞扬、到举办冬奥会中体现着中国科技的伟大、再到冰墩墩吉祥物和可爱等”言论中,大都体现着人们言论的积极情感。在图6中,某些时间点人们也表达出了消极的情绪。在图5中一些高频词“难买”、“抢购”等,这与网民

出现情感消极的原因相关。在表 2 的 *Topic 7* 主题词中,进一步表达出网民们的情感极性,表现出网民们对一墩难求的消极情感。

3.3 性能评价

为了验证 ARMA-ALEE 模型的有效性,本文在准确率、召回率和 F_1 值方面对模型的方面主题词的抽取和情感分类极性的判断性能进行评估。对数据集按时间顺序划分,选取 3 个时间段内的数据作为验证数据集,并对这 3 个时间段内的数据进行人工标注标记出主题词,选用 TF-IDF、TF-IDF-Means 主题提取算法和本文的模型算法进行比较,实验结果见表 3。实验结果表明,本文提出的模型在主题词提取优化方面取得了较好的效果,其在准确率、召回率和 F_1 值方面都取得了较好的结果。

表 3 各种算法对主题词提取对比表

Tab. 3 Comparison table of various algorithms for subject word extraction

模型	准确率	召回率	F_1 值
TF-IDF	0.781	0.763	0.769
TF-IDF-Means	0.823	0.806	0.811
ARMA-ALEE	0.918	0.898	0.905

为了验证 ARMA-ALEE 模型的效果是否可行,仍以上述选取的验证集作为实验数据,选取“JST 模型^[13]”、“ASUM 模型^[14]”和“主题-情感联合模型^[15]”进行对比实验,实验结果见表 4。

表 4 各种模型情感演化性能评价表

Tab. 4 Emotional evolution performance evaluation table of various models

模型	准确率	召回率	F_1 值
JST 模型	0.762	0.741	0.756
ASUM 模型	0.798	0.779	0.781
主题-情感联合模型	0.805	0.774	0.789
ARMA-ALEE	0.936	0.918	0.903

从表 4 中可以看出,在验证数据集中模型 ARMA-ALEE 在准确率、召回率和 F_1 值三个指标上都有明显的提升,表明模型 ARMA-ALEE 的性能是远远优于其它对比模型的。从实验结果分析可知,ARMA-ALEE 模型首先在方面级情感分析模型的基础上获取到方面词和情感极性值,又在 ARMA 时间序列模型基础上对优化后的方面词、情感极性值和相似度进行训练,提高了舆情演化主题的准确率,并最终使用 ARMA-ALEE 情感演化算法动态得到网络舆情情感演化结果。

4 结束语

为了提高主题词提取的精确度和实现动态网络

舆情情感演化分析,本文提出了一种基于时间序列的方面级网络舆情情感演化 ARMA-ALEE 模型。经实验验证,本文提出的 ARMA-ALEE 模型在准确率、召回率和 F_1 值方面都优于其它参考模型,证明了 ARMA-ALEE 模型在对网络舆情动态情感演化分析上的优越性。由于新提出的模型要依赖于分类效果较好的方面级情感分析模型,这也是今后需要进一步深入研究的地方。

参考文献

- [1] XU Guixian, YU Ziheng, YAO Haishen, et al. Chinese text sentiment analysis based on extended sentiment dictionary [J]. IEEE Access, 2019, 7: 43749-43762.
- [2] VIJAYARAGHAVAN S, BASU D. Sentiment Analysis in Drug Reviews using Supervised Machine Learning Algorithms [J]. arXiv preprint arXiv:2003.11643, 2020.
- [3] RUZ G A, HENRÍQUEZ P A, MASCAREÑO A. Sentiment analysis of Twitter data during critical events through Bayesian networks classifiers [J]. Future Generation Computer Systems, 2020, 106: 92-104.
- [4] LV Yanxia, WEI Fangna, CAO Lihong, et al. Aspect-level sentiment analysis using context and aspect memory network [J]. Neurocomputing, 2021, 428: 195-205.
- [5] 邢云菲,王晰巍,韦雅楠,等. 新媒体环境下网络舆情用户情感演化模型研究—基于情感极性情感强度理论 [J]. 情报科学, 2018, 36(08): 142-148.
- [6] 钱进宝. 基于演化 K-medoids 方法的微博情感动态分析—以《穹顶之下》为例 [J]. 情报杂志, 2019, 38(03): 155-159, 165.
- [7] 戴杏云,张柳,戴伟辉,等. 社交网络的情感图谱研究 [J]. 管理评论, 2016, 28(08): 79-86.
- [8] 张柳,王晰巍,王铎,等. 微博环境下高校舆情情感演化图谱研究—以新浪微博“高校学术不端”话题为例 [J]. 现代情报, 2019, 39(10): 119-126, 135.
- [9] 杨嘉韵,张慧明. 基于主题-情感融合分析的突发公共卫生事件网络舆情演化研究 [J]. 情报探索, 2021(08): 18-28.
- [10] GROSSMAN D A, FRIEDER O. Information retrieval: Algorithms and heuristics [M]. 2nd ed. USA: Springer Science & Business Media, 2004.
- [11] 关鹏,王曰芬. 科技情报分析中 LDA 主题模型最优主题数确定方法研究 [J]. 现代图书情报技术, 2016(09): 42-50.
- [12] KABÁN A, GIROLAMI M A. A dynamic probabilistic model to visualise topic evolution in text streams [J]. Journal of Intelligent Information Systems, 2002, 18(2): 107-125.
- [13] LIN Chenghua, HE Yulan. Joint sentiment/topic model for sentiment analysis [C]//Proceedings of the 18th ACM Conference on Information and Knowledge Management. Hong Kong: ACM, 2009: 375-384.
- [14] JO Y, OH A H. Aspect and sentiment unification model for online review analysis [C]//Proceedings of the fourth ACM International Conference on Web Search and Data Mining. Hong Kong.: ACM, 2011: 815-824.
- [15] 宋嘉欣. 基于主题-情感联合模型的网络舆情情感演化分析研究 [D]. 秦皇岛:燕山大学, 2020.