

基于位置和词性特征的藏文情感三元组抽取模型

斯曲卓嘎^{1,2,3}, 拥措^{1,2,3*}, 赛鸣宇^{1,2,3}

1. 西藏大学信息科学技术学院, 西藏拉萨 850000;
2. 西藏自治区藏文信息技术人工智能重点实验室, 西藏拉萨 850000;
3. 藏文信息技术教育部工程研究中心, 西藏拉萨 850000)

摘要:藏文情感三元组(方面词、情感词、情感极性)是细粒度情感分析的核心任务,对于深入理解藏族情感表达和趋势至关重要。但藏文的独特语言结构和文化背景导致其情感表达方式与其他语言不同,从而增加了细粒度情感分析的复杂性。为了提高藏文情感三元组的提取能力,该文提出了 OpinionNet-OTE-MTL 模型,该模型融合了词性信息、Word2Vec 词向量和绝对位置向量,并通过双向长短期记忆网络(BiLSTM)进行特征提取。其中,由于藏文词性种类较多,该文分析了大量的情感数据集并从中提取出11种词性辅助模型识别。最后,为了验证 OpinionNet-OTE-MTL 模型的有效性,在自构建的2000句藏文细粒度情感分析数据上进行了对比实验和消融实验。消融实验表明,词性较位置信息对模型的影响更大,其三元组抽取 F1 值提高了3.06个百分点;对比实验结果表明将词性和位置特征融入进模型后,在情感三元组提取(Triple)任务上的精确率、召回率和 F1 值较基线实验提高了4.73个百分点、6个百分点、6.14个百分点,融入词性和绝对位置信息使模型能更精确地理解藏文的语法结构和语法规则,从而提升了情感三元组分类任务的准确度。

关键词:藏文; Word2Vec; 词性; 位置特征; 情感三元组

中图分类号: TP391

文献标识码: A

文章编号: 1673-629X(2025)02-0130-08

doi: 10.20165/j.cnki.ISSN1673-629X.2024.0300

Tibetan Sentiment Triplet Extraction Model Based on Location and Part-of-Speech Features

SI Qu-zhuo-ga^{1,2,3}, YONG Tso^{1,2,3*}, SAI Ming-yu^{1,2,3}

1. School of Information Science and Technology, Tibet University, Lhasa 850000, China;
2. Key Laboratory of Tibetan Information Technology and Artificial Intelligence in Tibet Autonomous Region, Lhasa 850000, China;
3. Engineering Research Center of Tibetan Information Technology, Ministry of Education, Lhasa 850000, China)

Abstract: The extraction of Tibetan emotional triplets (aspect words, emotional words, emotional polarity) is a core task in fine-grained sentiment analysis, which is crucial for a deep understanding of Tibetan emotional expressions and trends. However, the unique language structure and cultural background of Tibetan make its emotional expression different from other languages, thereby increasing the complexity of fine-grained sentiment analysis. To enhance the capability of extracting Tibetan emotional triplets, we propose the OpinionNet-OTE-MTL model, which integrates part-of-speech information, Word2Vec word vectors, and absolute position vectors. Feature extraction is performed through Bidirectional Long Short-Term Memory networks (BiLSTM). Given the diverse types of parts of speech in Tibetan, we analyze a large amount of emotional datasets and extract 11 types of parts of speech to assist in model recognition. Finally, to validate the effectiveness of the OpinionNet-OTE-MTL model, comparative experiments and ablation experiments were conducted on a self-constructed dataset of 2000 Tibetan sentences for fine-grained sentiment analysis. The ablation experiments indicated that part-of-speech information had a greater influence on the model than positional information, resulting in a 3.06 percentage points increase in the F1 score of triplet extraction. Comparative experimental results showed that after integrating part-of-speech and positional features into the model, the precision, recall, and F1 score of the emotional triplet extraction task increased by 4.73 percentage points, 6 percentage points, and 6.14 percentage points respectively compared to the baseline experiment. Integrating part-

收稿日期: 2024-06-14

修回日期: 2024-10-16

基金项目: 西藏自治区科技计划项目(XZ202401JD0010)

作者简介: 斯曲卓嘎(1997-),女(藏),硕士研究生,研究方向为自然语言处理;通信作者: 拥措(1974-),女(藏),教授,博导,研究方向为自然语言处理、藏医药古籍智能信息处理。

of-speech and absolute positional information enables the model to better understand the grammatical structure and semantic rules of Tibetan, thereby improving the accuracy of emotional triplet classification tasks.

Key words: Tibetan; Word2Vec; part of speech; positional features; emotional triplet

0 引言

随着藏文信息的快速增长,众多藏文文本在互联网平台上变得日益丰富,对其进行情感分析有助于深入理解藏族的传统文化。张俊等研究者^[1]使用了基于情感词典的方法来识别情感倾向。与此同时,袁斌等研究者^[2]利用藏文的语法和语义特点,建立了一个特征空间进行情感分析。另外,孙本旺等研究者^[3]运用深度学习技术,通过 Word2Vec 创建词向量模型,并结合 CNN 和 LSTM 模型对藏文微博进行了情感分析,取得了良好的效果。

藏文情感分析虽然取得了初步进展,但目前的研究多集中在粗粒度情感分类的任务上,对于细粒度的情感分析尚未完全实现。而情感三元组抽取作为情感细粒度分析中最重要的任务之一,它通过识别文本中的方面词、情感词和情感极性,为藏文情感分析提供了更为详尽和精确的信息。不仅有助于准确识别藏文文本讨论的具体方面、相关联的观点以及相应的情感倾向,而且对于理解藏族文化中的情感表达模式和细微差别具有重要意义。然而,与中英文情感分析相比,藏文细粒度情感分析的研究起步较晚,面临诸多挑战。首先,目前缺乏足够规模的标注藏文细粒度数据集,严重限制了情感分析研究的深入和模型训练的有效性。其次,藏语文本的复杂句法结构、多样的修辞手法和丰富的词性,为情感分析带来了额外的难题。

因此,针对上述问题,该文做了以下工作。首先,针对目前没有藏文细粒度情感分析语料导致该研究空白的情况,构建了 2 000 句藏文细粒度情感语料,为后续研究提供了基础数据支持;其次,为了提高词嵌入捕捉藏文词义的丰富信息,使用了 430 MB 数据训练了 Word2Vec 词向量模型;再次,针对 OTE-MTL 模型识别藏文情感三元组效果不佳的情况,提出了融合词性、Word2Vec 向量及绝对位置向量的 OpinionNet-OTE-MTL 模型,增加模型对藏文情感三元组的识别和提取能力;最后,在该文构建的语料上使用 OpinionNet-OTE-MTL 模型进行了对比及消融实验,证明了词性以及位置编码对藏文情感细粒度分析具有良好效果。

1 相关研究

随着自然语言处理技术的进步,藏文情感分析得到了显著推动,为低资源语言的文本分析提供了新的视角和技术支撑。李海刚和于洪志^[4]最早开发了基于情感词库的文本分类系统,使用相似度算法进行情感

判断。2018 年拥措教授和史晓东^[5]总结了藏文短文本研究方法,并分析了情感分析的趋势。公保加羊等人^[6]提出了基于集成学习的藏文情感分析算法,结合现有藏文文本情感分析算法的优势,提升藏文情感分析模型的精度和效率。同年,朱宇雷等人^[7]构建了大规模的情感分析数据集,并采用 Graph-SAGE 网络提取特征。孟祥和等人^[8]则利用多核卷积神经网络和 BiLSTM (Bidirectional Long Short-Term Memory network) 模型提取文本特征。这些方法都有效提升了情感分类的准确度,但均在粗粒度情感分类上进行研究,并未涉及细粒度分析,而细粒度分析对于深入理解和挖掘藏文文本中的情感信息具有重要作用。

在自然语言处理(NLP)领域,词性和位置特征对于提升情感分析的准确率起着至关重要的作用。在 2017 年,Wang Wenya 等人提出了 CMLA^[9]模型,该模型利用多层注意力机制共同提取文本中的方面项和观点项。2018 年,Li Xin 团队提出了 HAST^[10]模型,通过融合多种技术手段提升了方面项提取的准确性;同一时期,赵富团队^[11]提出了 CWPAT-Bi-LSTM 模型,通过结合词性和注意力机制,进一步提高了情感分析的准确率。2019 年,姚艳秋^[12]采用了 LS-SO 算法,考虑情感词的上下文信息来提高分类的准确性;同时,王行甫^[13]等进一步优化了细粒度情感分析,开发了一种结合词性、位置和单词情感的内存网络模型,并引入了 POSP-CNNAM 机制,显著提升了情感分析的准确性。到了 2020 年,Lim 团队^[14]结合文本和地理位置信息,运用 CNN 和 BiLSTM 网络,成功提升了 Twitter 情感分析的准确度。2021 年 Xu Lu 等人^[15]提出新的位置感知标记方案和 JET 模型,该模型使用特征因式分解来捕获三元组元素间的交互作用;同年,薛芳团队^[16]通过双层词性感知技术,强化了情感分析模型的能力。在 2023 年,Wang Junlang 团队^[17]提出了一个融合词性信息的预训练语言模型框架,有效提升了在缺乏明显情感词汇时的文本情感分析性能;而杜孟洋团队^[18]则运用自注意力机制和句法依赖,进一步提升了情感分类的准确率。同年,周雨婷等人^[19]提出了基于语义增强和指导路由的方法,通过键值对网络动态学习文本特征,优化了情感三元组的抽取。到了 2024 年,李增伟^[20]提出的 SSES-SPAN 模型,通过特征编码和图卷积网络技术,不仅提高了情感分类的准确率,还为理解复杂语言表达提供了新的视角。同时,赵园春^[21]提出了基于细粒度标记的抽取算法,通过卷积注意力机

“T-POS”，表示正面情感；而紧随其后的情感词位置为(6,8)则被标记为“S”。同理，在示例二中，位置为(7,8)方面词被标记为“T-NEG”，表示负面情感，位置为(4,5)情感词标记为“S”。位置为(16,17)方面词被标记为“TT-POS”，表示正面情感，位置为(20,23)情感词被标记为“SS”，通过以上示例可以判断出藏文的方面词一般出现在句中起始或中间位置而紧随其后的是情感词，即位置信息的融入使方面词、情感词关系更明确有助于判断情感三元组。

2.4 模型

OPE-MTL 模型^[24]是由 Zhang Chen 等人在 2020 年提出的情感三元组识别模型，它通过多任务学习框架有效地抽取文本中的方面和观点，并准确解析它们之间的情感依赖关系。该模型采用的多头部架构和先进算法提高了处理的灵活性和准确性。但是，由于模型缺乏对藏文词义的理解，对藏文细粒度分析存在着挑战。为了克服这一难题，该文提出将词性和位置特征与 OTE-MTL 模型相结合，以更深入地理解藏文情感三元组，即 OpinionNet-OTE-MTL 模型。该模型通过两个主要步骤实现情感三元组的提取，首先是预测阶段，其次是解码阶段。在预测阶段，将位置嵌入、词嵌入与词向量嵌入相结合构建句子的编码表示，再输入至双向长短期记忆网络 (BiLSTM) 中进行学习。提取方面词和情感词的特征表示，并通过两个序列标记器进行方面和情感标记，同时使用双向评分器预测词对之间的情感依赖。接着通过联合训练最小化标记损失和依赖解析损失进行多任务学习。在解码阶段，利用启发式规则从提取的方面词、情感词和情感极性生成最终的三元组输出。模型框架如图 1 所示。

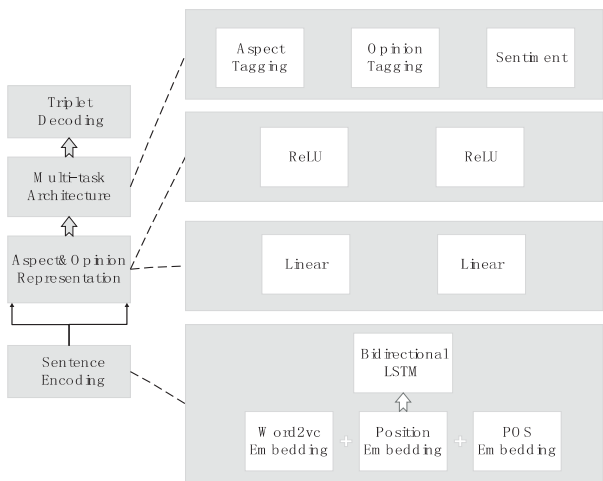


图 1 OpinionNet-OTE-MTL 模型框架

OpinionNet-OTE-MTL 模型主要使用双向长短期记忆网络 (BiLSTM) 来处理输入文本，以捕获每个词的上下文信息及其在句中的位置特征。其中，双向 LSTM 通过以下操作进行特征提取：

$$\text{Embedding} = e_i \oplus p_i \oplus t_i \tag{3}$$

$$h_i = [\text{Lstm}(\text{Embedding})][\text{Lstm}(\text{Embedding})] \tag{4}$$

式中， e_i 是由 Word2Vec 获得的词嵌入， p_i 是词的位置向量， t_i 是词的词性表示，Embedding 是融合了位置、词性和词信息的上下文表示， h_i 是利用 Lstm 对 Embedding 进行特征提取获得的更丰富的语义信息，使得生成的上下文向量序列更为精确。其次，模型通过线性降维和 ReLU 激活函数提取藏文方面和情感词的关键特征。这一步旨在去除对后续计算不必要的信息，减少过拟合风险，并剔除与方面标记和情感词标记无关的特征。图 2 为 OpinionNet-OTE-MTL 模型的输入输出样例图，通过在模型的输入添加词性和位置信息，可以辅助模型对句子的方面词和情感词的词性和位置编码以及情感态度进行判断，从而能够准确抽取出文本的情感三元组。结合 Word2Vec、词性及位置编码输入至 BiLSTM 模型，还可以增强模型对藏文文本情感内容的理解。

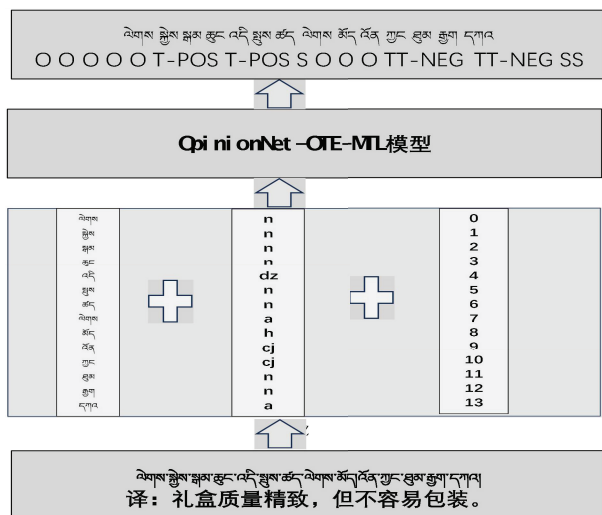


图 2 输入输出样例图

3 实验数据与评价指标

3.1 实验数据

在藏文细粒度情感分析研究领域，由于缺乏公开的语料库导致无法进行细粒度情感分析，因此，为了评估 OpinionNet-OTE-MTL 模型的有效性，该文构建了一个新的藏文细粒度文本语料库，用于情感三元组识别测试。构建步骤如下：首先，从中文社交媒体平台（如微博、电商和淘宝）收集大量产品的用户评价；其次，为了确保翻译的情感准确性，选择了微软翻译工具。基于微软翻译在精确捕捉情感极性和深入理解上下文方面的显著优势，因此使用微软翻译，将从多个中文社交媒体平台收集的数据集转换成了藏文；最后，对构建的藏语文料库进行了详尽的人工标注和校对，确

百分点和 16.44 百分点。通过实验结果分析发现,虽然方面词提取层面 HAST 模型略优于 OTE-MTL,但对整体情感三元组提取和情感词提取的效果在 OTE-MTL 模型上表现最佳。因此后文实验在 OTE-MTL 上融合词性与绝对位置信息。

4.3 消融实验

为了验证提出的融合词性和位置向量的 OpinionNet-OTE-MTL 模型的合理性,设计了消融实验评估词性和位置对模型的影响效果,实验结果如表 8 所示。与 OTE-MTL 模型相比,仅融入词性的 POS-OTE-MTL 在三元组提取任务(triple)上表现出了一定的提升,其中 F1 值提高了 3.06 百分点。这表明词性信息对于识别文本中的三元组关系具有一定的帮助。在方面词提取任务时,POS-OTE-MTL 模型相比 OTE-MTL 模型,F1 值提升了 3.09 百分点,因此词性对提取方面词有关键作用,有助于句子结构识别、语义角色揭示、情感表达关联及上下文理解。其次,与 OTE-MTL

模型相比,在三元组(triple)提取任务上,融入位置信息的 Position-OTE-MTL 模型 F1 值提升了 2.83 百分点,与仅融入词性相比减少了 0.23 百分点。因此,结果证明了融入词性信息对三元组识别的影响大于位置信息。最后,当 OpinionNet-OTE-MTL 模型融合了词性和位置信息时,其性能提升更为显著。在三元组提取(triple)任务上,与仅融入词性的 POS-OTE-MTL 模型相比,OpinionNet-OTE-MTL 模型的精确率、召回率和 F1 值分别提高了 4.68 百分点、2.25 百分点和 3.08 百分点。而在与仅融入位置信息的 Position-mtl 模型相比时,在方面词(ap)提取任务上,OpinionNet-OTE-MTL 模型的精确率、召回率和 F1 值分别提高了 4.76 百分点、1.5 百分点和 3.33 百分点。消融实验的结果证实了 OpinionNet-OTE-MTL 模型通过整合词性和位置信息,相辅相成增强了模型各子模块的有效性,并在提升情感三元组任务的性能上发挥了重要作用。

表 8 消融实验结果 %

模型	ap(方面词)			op(情感词)			triple(三元组)		
	精确率	召回率	F1 值	精确率	召回率	F1 值	精确率	召回率	F1 值
OTE-MTL	38.08	42.74	40.28	62.38	66.74	64.49	43.80	22.99	30.16
Position-mtl	40.25	46.99	43.36	59.16	63.74	61.37	39.64	28.24	32.99
POS-OTE-MTL	40.16	47.99	43.37	61.35	65.74	63.36	43.85	26.74	33.22
OpinionNet-OTE-MTL	45.01	48.49	46.69	57.00	58.99	57.98	48.53	28.99	36.30

5 结束语

藏文的情感三元组提取任务作为新兴的细粒度情感分析领域,具备巨大的发展潜力,相对于传统的情感分类,它提供了更为深入和详细的情感理解。针对目前藏文情感三元组研究较少,并且藏文细粒度情感抽取困难的原因,该文提出了 OpinionNet-OTE-MTL 模型,该模型通过结合词性和其在文本中的位置信息,更全面地捕获了文本的情感信息,词性的融入可以更好地定位方面词和情感词,位置的融入则增强了模型对情感倾向的精确识别能力。通过对比和消融实验均证明了该模型的有效性。

此外,藏文作为一种语法结构复杂的语言,该文仅考虑了词性和位置信息,且标注为人工标注,因此存在着主观性、成本高和忽略上下文的问题,未来将继续探索并融合更多的藏文语法结构和语义信息,以增强藏文情感分析的细致度和全面性。同时,在数据构建方法上采用机器标注结合人工校正的方法,以及期望在后续研究中,通过深入挖掘藏文文法,并将依存关系分析和注意力机制等先进技术应用于模型,进一步提升

情感三元组提取的准确度。

参考文献:

- [1] 张俊,李应兴.基于情感词典的藏文微博情感分析研究[J].硅谷,2014,7(20):220.
- [2] 袁斌,江涛,于洪志.基于语义空间的藏文微博情感分析方法[J].计算机应用研究,2016,33(3):682-685.
- [3] 孙本旺,田芳.基于深度学习算法的藏文微博情感计算研究[J].计算机技术与发展,2019,29(10):55-58.
- [4] 李海刚,于洪志.藏文文本情感分类系统设计[J].甘肃科技纵横,2011,40(1):106-107.
- [5] 拥措,史晓东,尼玛扎西.短文本情感分析的研究现状——从社交媒体到资源稀缺语言[J].计算机科学,2018,45(s1):46-49.
- [6] 公保加羊,拉玛杰,官却多杰,等.基于深度学习的藏文情感分析研究[J].青海科技,2023,30(1):56-60.
- [7] 朱宇雷,德吉卡卓,群诺,等.基于图神经网络结合预训练模型的藏文短文本情感分析研究[J].中文信息学报,2023,37(2):71-79.
- [8] 孟祥和,于洪志.融合音节和词条特征的藏文文本情感分类研究[J].中文信息学报,2023,37(2):80-86.
- [9] WANG Wenya, PAN S J, DAHLMEIER D, et al. Coupled

- multi-layer attentions for co-extraction of aspect and opinion terms[C]//Proc of AAAI conference on artificial intelligence. Palo Alto; AAAI, 2017: 273-284.
- [10] LI Xin, BING Lidong, LI Piji, et al. Aspect term extraction with history attention and selective transformation[C]//Proc of the 27th international joint conference on artificial intelligence. Palo Alto; AAAI, 2018: 4194-4200.
- [11] 赵富, 杨洋, 蒋瑞, 等. 融合词性的双注意力 Bi-LSTM 情感分析[J]. 计算机应用, 2018, 38(S2): 103-106.
- [12] 姚艳秋, 郑雅雯, 吕妍欣. 基于 LS-SO 算法的情感文本分类方法[J]. 吉林大学学报: 理学版, 2019, 57(2): 375-379.
- [13] 王行甫, 王磊, 苗付友, 等. 结合词性、位置和单词情感的内存网络的方面情感分析[J]. 小型微型计算机系统, 2019, 40(2): 383-389.
- [14] LIM W L, HO C C, TING C Y. Sentiment analysis by fusing text and location features of geo-tagged tweets[J]. IEEE Access, 2020, 8: 181014-181027.
- [15] XU Lu, LI Hao, LU Wei, et al. Position-Aware tagging for aspect sentiment triplet extraction[C]//Proc of conference on empirical methods in natural language processing. Stroudsburg; ACL, 2020: 2339-2349.
- [16] 薛芳, 过弋, 李智强, 等. 基于双层词性感知和多头交互注意机制的方面级情感分析[J]. 计算机应用研究, 2022, 39(3): 704-710.
- [17] WANG J, LI X, HE J, et al. Enhancing implicit sentiment learning via the incorporation of part-of-speech for aspect-based sentiment analysis[C]//China national conference on chinese computational linguistics. Singapore; Springer, 2023: 382-399.
- [18] 杜孟洋, 王红斌, 普祥和. 融入词性自注意力机制的方面级情感分类方法[J]. 吉林大学学报: 理学版, 2023, 61(6): 1375-1386.
- [19] 周雨婷, 代金鞘, 刘嘉勇, 等. 一种基于语义增强和指导路由机制的方面级情感三元组抽取方法[J]. 四川大学学报: 自然科学版, 2023, 60(5): 112-120.
- [20] 李增伟, 刘帅. 语义和句法依赖增强的跨度级方面情感三元组抽取[J]. 计算机系, 2024, 33(6): 201-210.
- [21] 赵园春, 韩虎, 徐学锋. 细粒度标记的结点自适应方面情感三元组抽取[J/OL]. 计算机工程与应用, 1-11 [2024-08-16]. <http://kns.cnki.net/kcms/detail/11.2127.tp.20240626.1435.006.html>.
- [22] 黄梓芑, 曾碧卿, 陈鹏飞, 等. 基于语言特征增强的方面情感三元组抽取[J/OL]. 计算机工程, 1-12 [2024-09-01]. <https://doi.org/10.19678/j.issn.1000-3428.0069260>.
- [23] 郑文超, 徐鹏. 利用 word2vec 对中文词进行聚类研究[J]. 软件, 2013, 34(12): 160-162.
- [24] ZHANG Chen, LI Qiuchi, SONG Dawei, et al. A multi-task learning framework for opinion triplet extraction[C]//Proc of conference on empirical methods in natural language processing. Stroudsburg; ACL, 2020: 819-828.