

Sentiment Analysis Method combining Word Embedding with Personalized Text Structure

张兴华, 沙瀛

中国科学院信息工程研究所, 北京
中国科学院大学网络空间安全学院, 北京
xing_hua_zhang@126.com shaying@iie.ac.cn

问题

句子/篇章级的情感分析是一项非常重要的自然语言处理任务, 目前针对该任务已存在大量出色的研究工作。

由于大部分的文本来源于社交媒体, 社交网络中的个性化信息实则蕴含着文本之间的情感关联。我们的主要任务就是获取社交文本之间的情感关系, 以进一步提升情感分析的性能。

我们的方法

句子/篇章级的情感分析目前已有的研究工作主要分为两类: ①基于文本内容的情感分析[1]; ②显式融合社交网络个性化特征的情感分析[2]。

基于对上述工作以及社交个性化信息(话题依赖性、用户依赖性和社交一致性, 如图1所示)的探究, 我们提出一种结合文本语义特征与文本之间隐式情感关联特征的方法——融合个性化文本结构的词嵌入式情感分析。

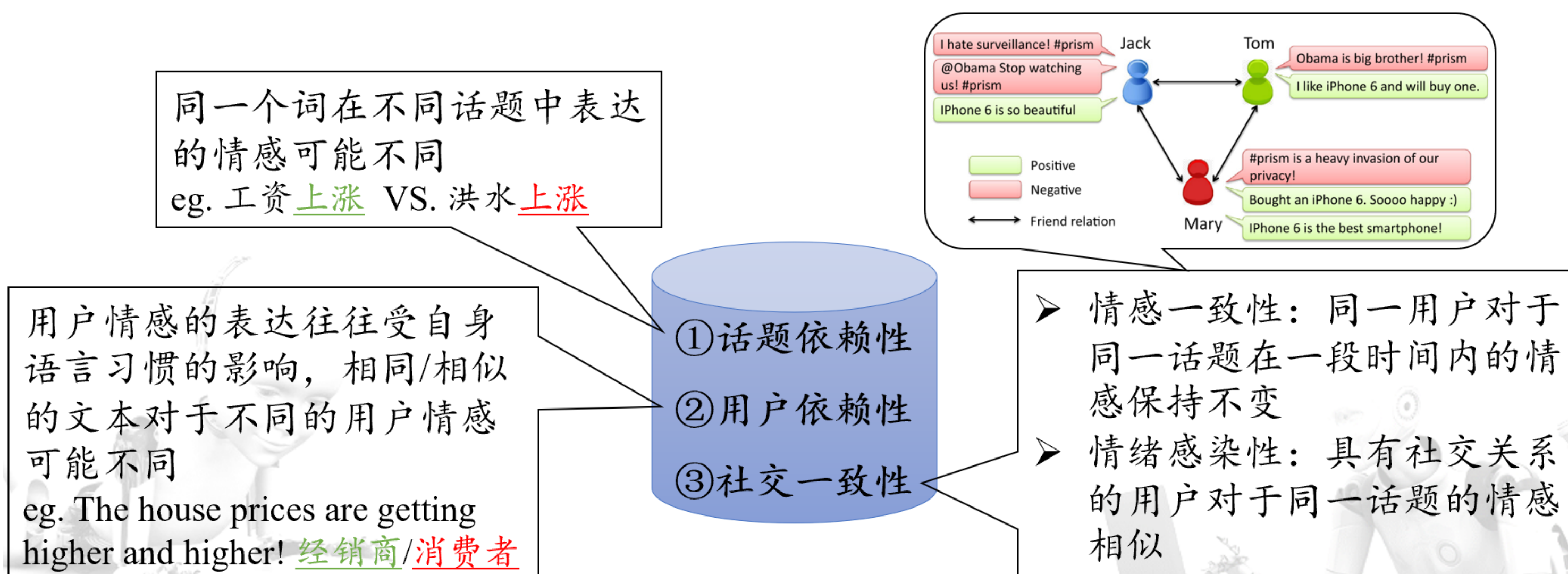


图1 社交个性化信息

融合个性化文本结构的词嵌入式情感分析方法

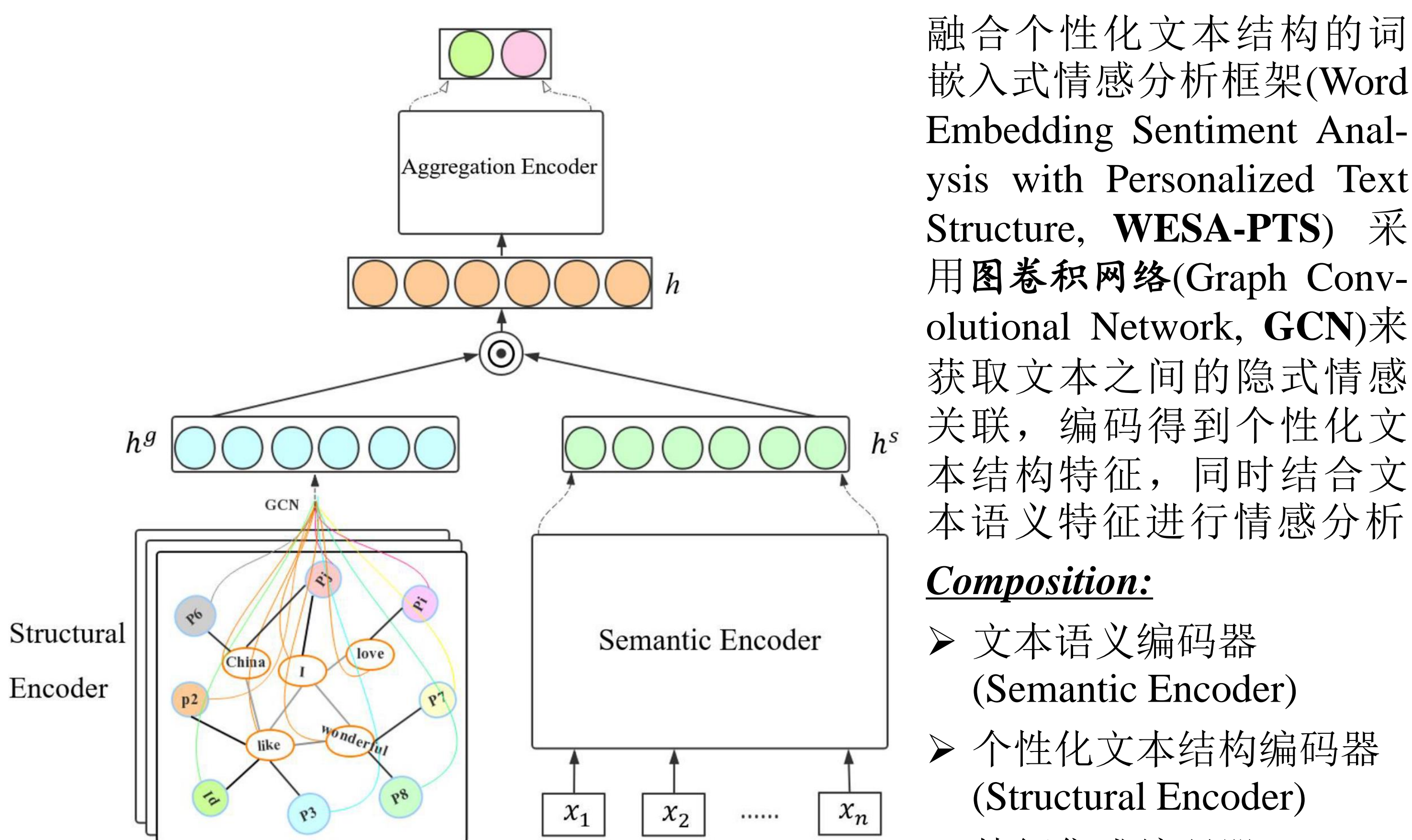


图2 融合个性化文本结构的词嵌入式情感分析框架

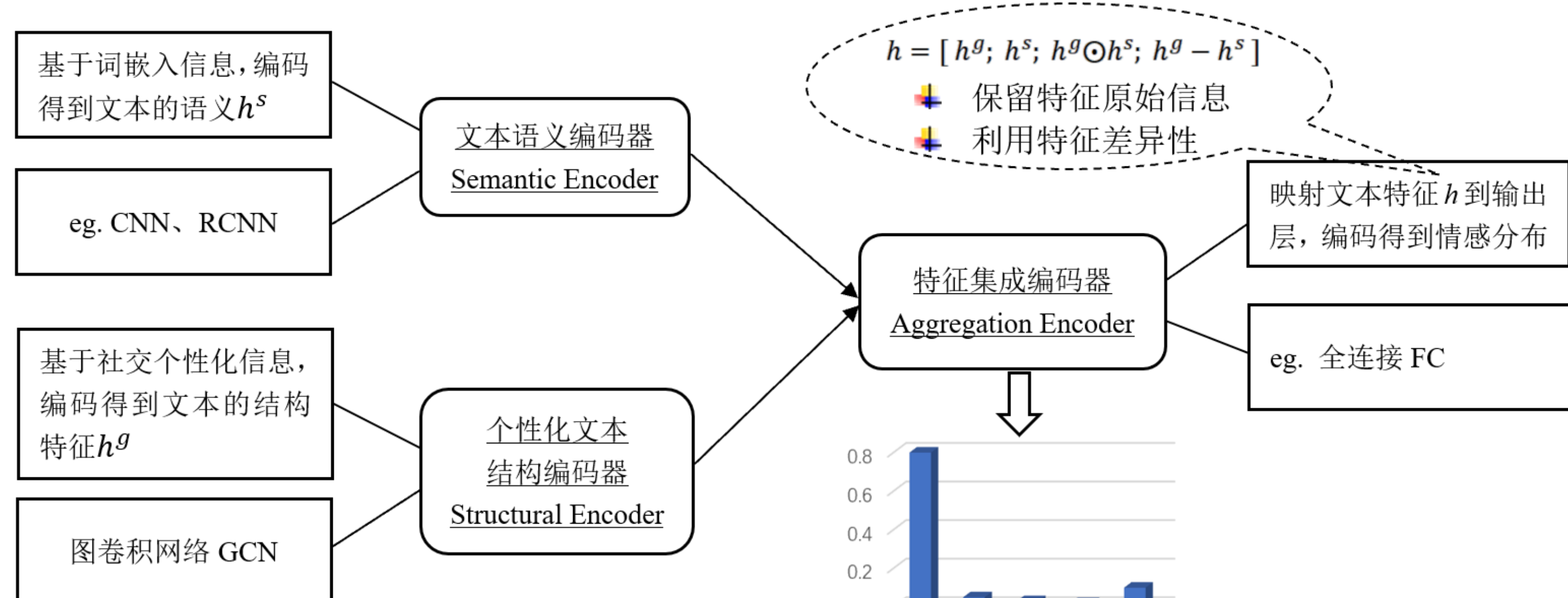


图3 WESA-PTS编码器简介

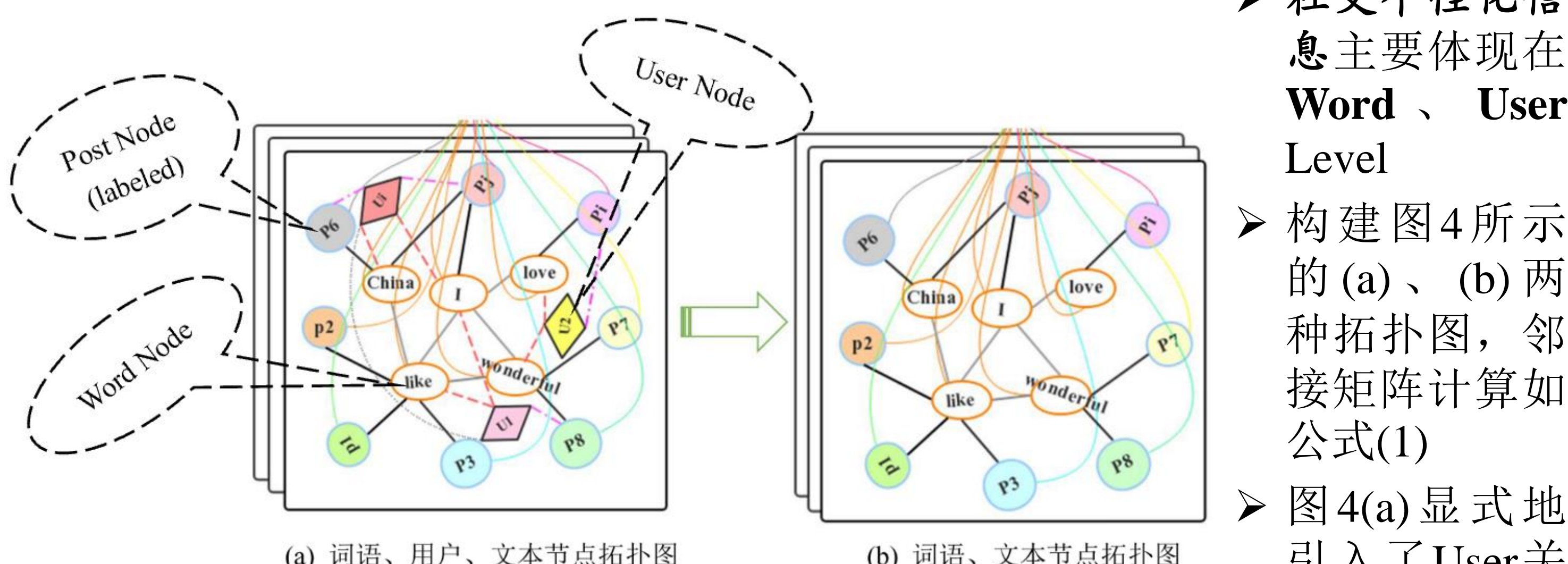


图4 情感分析拓扑图A

GCN Input Layer:

- 社交个性化信息主要体现在 Word、User Level
- 构建图4所示的(a)、(b)两种拓扑图, 邻接矩阵计算如公式(1)
- 图4(a)显式地引入了User关系

$$A_{i,j} = \begin{cases} \text{TF-IDF}(i,j) & i,j \text{ 分别为文本(用户)节点和词语节点} \\ \text{PMI}(i,j) & i,j \text{ 均为词语节点, } \text{PMI}(i,j) > 0, \text{ 否则等于 } 0 \\ P(i,j) & i,j \text{ 分别为用户节点和文本节点} \\ F(i,j) & i,j \text{ 均为用户节点} \\ 1 & i=j \\ 0 & \text{其它} \end{cases} \quad (1)$$

注:

1. TF-IDF(i, j): 计算节点 i 和 j 之间的词频-逆向文档频率。
2. PMI(i, j): 计算节点 i 和 j 之间点互信息量。
3. $P(i, j)$: 文本发表的所属关系, 值为0或1。
4. $F(i, j)$: 社交好友关系, 值为0或1。

Experiments Setting:

Encoder	Semantic Encoder	Structural Encoder	Aggregation Encoder
WESA-PTS			
G2CNN	CNN	GCN	FC
GRCNN	RCNN		

实验结果

实验结果准确率 Accuracy \pm Standard Deviation (%)

实验方法	Yelp-P	Amazon-P	Yelp-F	Amazon-F	
LSTM	89.47 \pm 0.23	86.46 \pm 0.26	48.43 \pm 0.45	62.92 \pm 0.14	
LSTM+attention	87.68 \pm 0.17	85.54 \pm 0.13	52.32 \pm 0.33	64.90 \pm 0.17	
Bi-LSTM	89.72 \pm 0.31	86.63 \pm 0.26	49.01 \pm 0.48	63.62 \pm 0.29	
Bi-LSTM+attention	90.92 \pm 0.06	88.05 \pm 0.07	56.22\pm0.17	66.27 \pm 0.09	
CNN static	88.67 \pm 0.07	86.56 \pm 0.07	52.33 \pm 0.07	65.02 \pm 0.05	
CNN non static	89.70 \pm 0.08	88.29 \pm 0.05	53.20 \pm 0.14	65.70 \pm 0.07	
RCNN	90.91 \pm 0.08	88.79 \pm 0.11	55.70 \pm 0.29	66.36 \pm 0.07	
GCN+user	88.03 \pm 0.00	88.18 \pm 0.00	53.92 \pm 0.02	66.01 \pm 0.00	
GCN	90.28 \pm 0.00	88.61 \pm 0.00	54.54 \pm 0.00	66.26 \pm 0.00	
WESA-PTS(Ours)	G2CNN	90.50 \pm 0.18	89.31 \pm 0.08	55.56 \pm 0.34	66.32 \pm 0.04
	GRCNN	91.27\pm0.08	89.58\pm0.15	55.88 \pm 0.21	66.64\pm0.11

数据集

- Yelp-P: 2分类
Train 7216
Test 2264
- Amazon-P: 2分类
Train 7003
Test 2335
- Yelp-F: 5分类
Train 10000
Test 5000
- Amazon-F: 5分类
Train 16624
Test 5542

注: 训练集的10%作为验证集

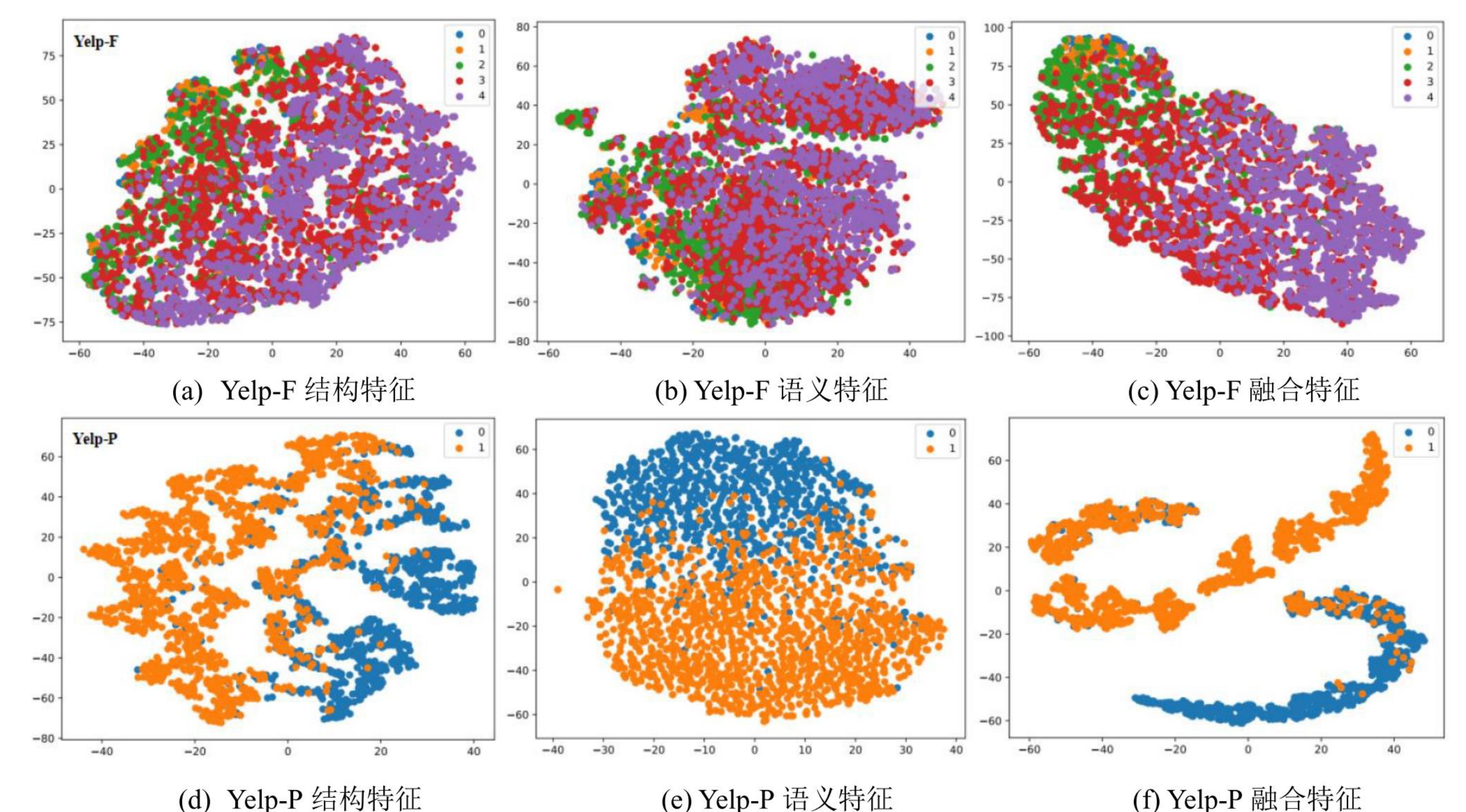


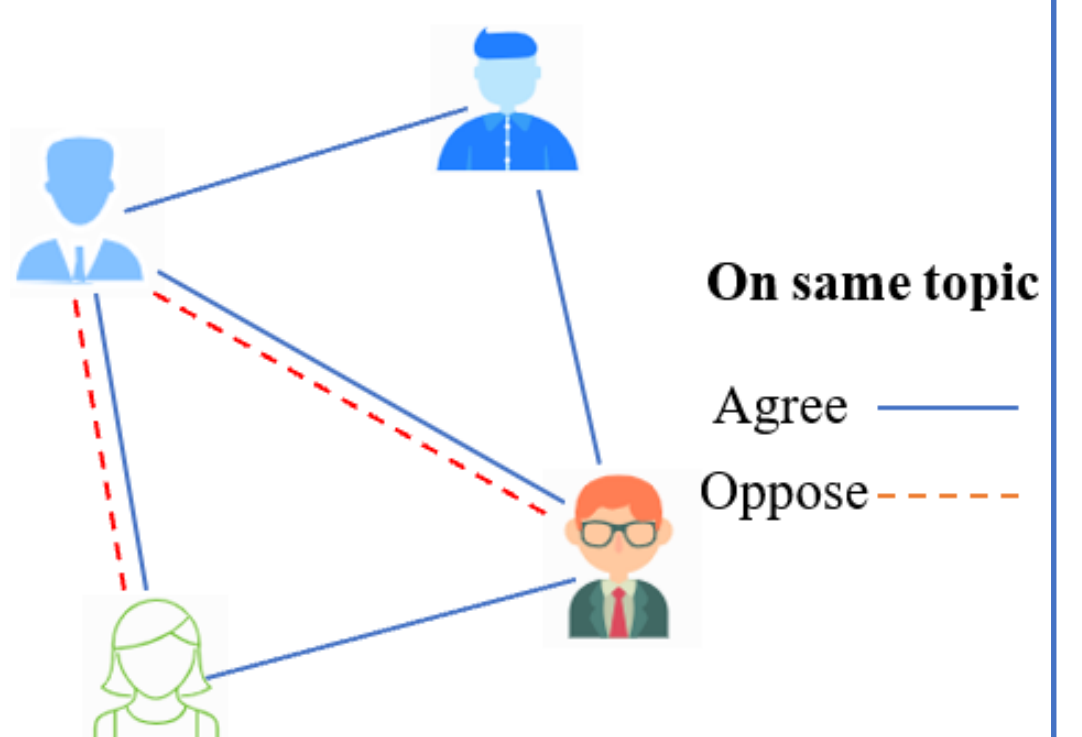
图5 文本结构特征与语义特征及其融合特征的 t-SNE 可视化

Conclusion:

- 社交网络中的个性化信息蕴含着文本之间的情感关联, 隐式地融入社交网络个性化特征的方法要明显优于显式的方式。
- 结合文本语义特征与结构特征进行情感分析的方法明显优于单纯利用语义特征或结构特征的方式。

Discussion:

- ❑ 由于社交环境的复杂多变性, 基于固定社交规则的方法性能往往受到限制(具有社交关系的用户并不完全满足社交一致性, 如右图)。
- ❑ 因文本语义信息的复杂多样性, 同时结合文本关联关系的方法往往能够取得更优性能。



参考文献

- [1] Zhang L, Wang S, Liu B. Deep learning for sentiment analysis: A survey[J]. Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, 2018: e1253.
- [2] Fangzhao Wu, Yongfeng Huang, and Yangqiu Song. Structured microblog sentiment classification via social context regularization[J]. Neurocomputing 175 (2016): 599-609.
- [3] Yao, L., Mao, C., and Luo, Y. Graph convolutional networks for text classification[C]. In Proceedings of the 33rd AAAI Conference on Artificial Intelligence (AAAI), 2019.
- [4] Wu F, Huang Y. Personalized microblog sentiment classification via multi-task learning[C]. Thirtieth AAAI Conference on Artificial Intelligence. 2016.