

## 多源社交数据融合的多角度旅游信息感知

郭 彤, 郭 斌, 张佳凡, 於志文, 周兴社

(西北工业大学 计算机学院, 陕西 西安 710129)

**摘 要:** 为了提高用户获取旅游知识的效率, 同时提供有效的旅行辅助, 提出多源社交数据融合的多角度旅游信息感知方法. 对异构的旅游相关数据进行预处理, 提出跨媒体多角度关联方法来连接碎片化旅游信息, 利用“景观-特征刻画”实现景点的多角度刻画. 通过序列模式挖掘算法, 从历史游记数据中得到典型旅游路线并向用户进行推荐. 基于评论文本和图像上、下文的相似性, 将图像和文本结合, 实现跨媒体信息关联. 从大众点评和蚂蜂窝上, 采集国内 8 个热门景点数据进行实验. 结果表明, 采用该方法能够更细粒度地刻画景点, 且提供的旅游路线能够满足不同用户的需求.

**关键词:** 群体智能; 景区感知; 智能推荐; 社交数据融合; 多角度刻画

**中图分类号:** TP 399

**文献标志码:** A

**文章编号:** 1008-973X(2017)04-0663-06

## CrowdTravel: leveraging heterogeneous crowdsourced data for scenic spot profiling

GUO Tong, GUO Bin, ZHANG Jia-fan, YU Zhi-wen, ZHOU Xing-she

(School of Computer Science, Northwestern Polytechnical University, Xi'an 710129, China)

**Abstract:** A multi-aspect tourism information perception method based on multi-source social media data fusion was proposed in order to improve the efficiency of travel knowledge acquisition and provide effective travel assistance for users. A cross-media multi-aspect correlation approach was proposed to connect fragmented travel information after a preprocessing step for the heterogeneous travel related data, which resorted to the “scene-feature characterization” in order to make a multi-aspect characterization. The perception method mined typical travel route from historical travelogues based on the sequential pattern mining, which can be regarded as the recommended route. Connecting cross-media information was based on the similarity between the reviews and the image contexts. Results of experiments over a dataset of eight domestic popular scenic spots, which was collected from Dazhongdianping and Mafengwo, indicate that the approach makes a fine-grained characterization for the scenic spots and the provided travel route can meet different users’ needs.

**Key words:** crowd intelligence; scenic spot profiling; intelligent recommendation; social media data fusion; multi-aspect characterization

群体智能<sup>[1-2]</sup>是指挖掘和关联群体贡献的数据以实现感知对象的多侧面感知与理解. 社交媒体作为一种在线交互平台, 近年来呈现多样化发展的趋势, 大量用户组成虚拟网络社区, 允许用户发布信

息并支持群体用户分享和传播信息.

随着旅游业的快速发展, 来自世界各地的游客喜欢在旅行之后借助社交媒体来表达关于景点的看法, 这种群智贡献的信息可以帮助其他用户进行旅行安

收稿日期: 2016-12-15.

浙江大学学报(工学版)网址: [www.zjujournals.com/eng](http://www.zjujournals.com/eng)

基金项目: 国家“973”重点基础研究发展规划资助项目(2015CB352400); 国家自然科学基金资助项目(61332005, 61373119).

作者简介: 郭彤(1993—), 男, 硕士生, 从事移动群智感知的研究. ORCID: 0000-0003-3910-8337. E-mail: [tongg@mail.nwpu.edu.cn](mailto:tongg@mail.nwpu.edu.cn)

通信联系人: 郭斌, 男, 教授. ORCID: 0000-0001-6097-2467. E-mail: [guob@nwpu.edu.cn](mailto:guob@nwpu.edu.cn)

排. 旅行评论和博客游记是两种主流的社交旅游共享方式, 可以作为旅游信息总结的可靠知识来源. 面对日益增长的评论和游记, 非常需要一种信息感知方法来处理海量旅游信息并为用户提供准确的旅行参考.

近年来, 针对旅游领域进行知识挖掘的相关研究<sup>[3-5]</sup>, 通过不同的角度展开工作. 例如, Rattenbury 等<sup>[6]</sup>利用每个地点标签的分布来提取地标; Hao 等<sup>[7]</sup>研究从游记中挖掘代表性的地标信息来推荐旅游路线. 由于上述方法都是从单一角度来刻画景点, 没有把文本和图像关联起来, 导致最终效果不好. 多源社交数据融合<sup>[8]</sup>有以下两个优点: 1) 图像可以作为附加信息对文本内容进行补充, 特别是在旅游领域, 文本缺乏足够的表现力; 2) 利用具体的图像可以使用户对感兴趣的景观有更直观的了解. 一项最近的研究<sup>[9]</sup>提出利用评论和游记来挖掘特征, 并生成可视化总结的方法. 该研究只是利用游记中的图像结合挖掘的特征生成可视化总结, 没有进一步从中发现和感知路线信息, 最终应用比较简单. 本文把景点内的一处人文或自然景观看作景点的一个标签, 结合人们对某个景观的描述(美丽的、壮观的等)构成“景观-特征刻画”组合, 这种关联模式可以帮助人们从多角度了解景点信息.

基于以上考虑, 本文提出多源社交数据融合的多角度旅游信息感知方法. 首先利用跨媒体多角度关联方法将碎片化的评论数据关联起来, 进而通过序列模式挖掘从游记中挖掘典型旅游路线并进行推荐, 最后基于评论和图像上下文之间的相似性来实现跨媒体信息关联. 本文通过从大众评论和蚂蜂窝上, 采集国内流行的 8 个景点的数据进行实验. 实验结果表明, 本文提出的方法能够更细粒度地刻画景点, 并且提供的旅游路线能够满足不同用户的需求.

## 1 相关工作

近年来, 群体智能的应用领域越来越广泛, 包括事件发现和刻画、商业智能、活动推荐等. 特别是在旅游领域, 大量用户通过社交媒体发布的信息十分丰富, 利用这些内容进行旅游知识挖掘成为研究人员关注的问题之一<sup>[10]</sup>.

一般来说, 旅游知识挖掘可以分为两个方面: 一方面是从海量社交网络中挖掘丰富旅游知识, 另一方面是旅游推荐. 目前, 通过社交网络挖掘丰富旅游知识的研究大都基于传统 MDS 理论. 比如, Lin 等<sup>[11]</sup>使用位置和主题等信息, 提出目前最流行的旅游总结算法之一. Wan 等<sup>[12]</sup>通过区分内部文档和文

档间连接, 改善了传统的基于图的排序算法. Li 等<sup>[13]</sup>使用结构化支持向量机来选择句子. 上述方法只可以给出整体的一般描述, 不能从多角度进行发现. 对于旅游推荐, 利用特征挖掘是一种常见的思路. 基于特征挖掘的方法<sup>[14-15]</sup>通常对高频名词短语进行一些约束来识别产品特征, 因此经常会得到非特征结果, 还可能漏掉低频的特征<sup>[16-17]</sup>. 基于特征的方法需要对各种参数进行调整, 这使得它们很难被另一个数据集合复用. 除特征挖掘外, 由于概率主题模型(probabilistic latent semantic analysis, PL-SA)在低维空间上发现主题的效果较好, 已经在各种文本挖掘任务中被运用. 现有的概率主题模型都不适用于本文提出的问题, 因为它们没有考虑和解决前面提到的评论和游记数据的内在局限.

为了更好地满足不同用户的需求, 提出的方法能够同时实现上述两个方面, 即在从海量社交网络中挖掘丰富旅游知识的同时, 实现旅游推荐.

## 2 系统框架

系统框架如图 1 所示. 图中, 输入是某一个景点的评论、游记和景观词的集合,  $P = \{R, T, N\}$ ,  $R$  表示评论,  $N$  表示该景点内所有主要景观的名称集合,  $T = \{C, I\}$  代表游记, 其中  $I$  表示图像,  $C$  表示图像的上下文. 输出是一组按照推荐旅游路线的顺序排列的景观-特征刻画, 利用高质量且有代表性的图像对景观-特征刻画进行可视化.

系统框架包含以下 4 个部分: 1) 文本预处理: 将景点所有的评论数据作为输入, 挑选至少包含一个景观词且具有高信息熵的句子作为下一步的输入; 2) 多角度旅游信息感知: 利用上一步得到的文本内容, 结合景观词集合, 挖掘每个景观相对应的特征; 3) 旅游路线推荐: 从每篇游记中提取一条旅游路线, 将热度最高的路线呈现给用户; 4) 跨媒体信息关联: 通过比较图像上、下文和评论之间的文本相似性, 投票选择与景观对应的有代表性的图像.

### 2.1 文本预处理

预处理分为 3 步: 分词, 去除停用词, 挑选高信息熵句子.

首先使用 FudanNLP 对评论进行分词处理, 为每个景点建立一个单独词典  $V$ , 然后过滤掉中文停顿词. 结合该应用实际背景, 同时过滤与摄影相关的词, 最终效果会更好. 本文假设关于景点的信息均匀分布在评论数据集中. 此外, 根据信息论知识可知, 当一个随机事件的结果已知时, 信息熵可以反映携

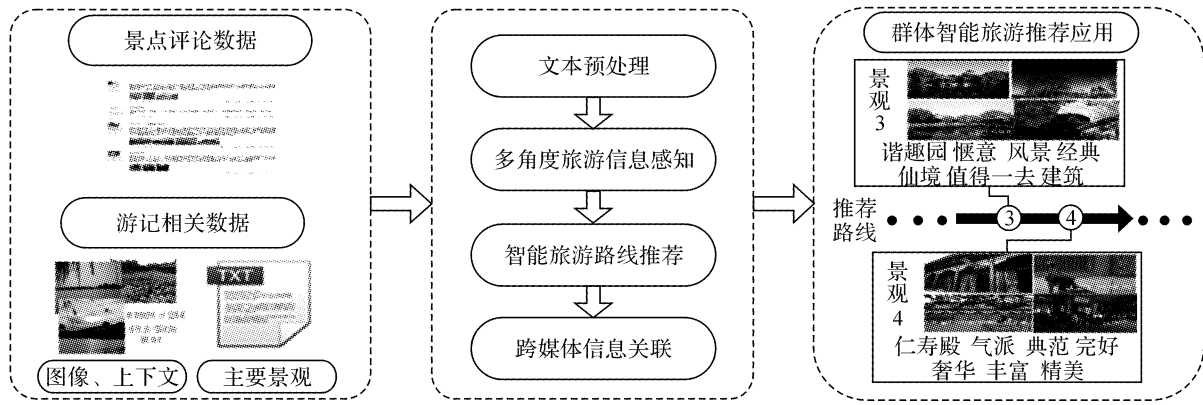


图1 基于多源社交数据融合的多角度旅游信息感知系统框架

Fig. 1 Overview of multi-aspect tourism information perception based on multi-source social media data fusion

带信息的平均数量. 首先计算每个单词  $w$  的信息熵, 如下所示:

$$H(w) = -p(w) \log p(w). \quad (1)$$

式中:  $p(w)$  为评论句子集合中每个单词  $w$  出现的概率,  $p(w) = f(w, R) / \sum_w t_f(w', R)$ , 其中  $t_f$  为单词  $w$  的词频.

本文提出多角度景观-特征刻画组合, 若一个评论句子不包含任何景观词, 则认为它对于下一步信息感知没有帮助. 为了降低计算成本, 不会继续作为下一步的输入.  $R$  中至少包含一个景观词的句子  $s$  的信息熵用句中每个单词  $w_k$  信息熵的总和来计算:

$$H(s) = \sum_k H(w_k). \quad (2)$$

按照下式所示的方法挑选具有高信息熵的重要的句子:

$$S = \{s \mid H(s) > \varepsilon, s \in R\}. \quad (3)$$

式中:  $\varepsilon$  为通过实验调节设置的阈值参数,  $S$  为从评论  $R$  中挑选出来的句子集合.

## 2.2 多角度旅游信息感知

特征词挖掘会受到噪声和不同表达方式的影响. 为了解决这些困难, 提出一种新的增量学习方法.

作为输入数据的一部分, 已经获取了景点  $P$  中的主要景观集合  $N = \{n_k\}_{k=1}^K$ . 根据与景观  $n$  相关的句子集合  $S_n = \{s_i\}_{i=1}^N \in S$  为该景观提取特征词, 构成集合  $W$ . 特征词包括名词和形容词, 能够很好地描述游客对该景观的评价.

借助文献[6]的假设, 采用一种简单有效的贪心策略来挖掘特征词. 事实上, 对于景点内部的每一个独立景观, 只有一个小的单词集合是有价值的. 利用一个预过滤步骤来获得  $F(w_i, S_n)$  较大的单词  $w_i$ , 得到集合  $W_F$ , 这样可以有效地降低计算成本.  $F(w_i, S_n)$  表示  $p(w_i | S_n)$  和  $p(w_i | S)$  之间的约

束关系:

$$F(w_i, S_n) = f(x), \quad x = p(w_i | S_n) - p(w_i | S) > 0. \quad (4)$$

其中引入了一个递增逻辑函数  $f(x) = 1/(1 + e^{-x})$ .

具体地说, 首先挑选  $F(w, S_n)$  最大的单词  $w \in W_F$ , 然后通过求解  $\arg \max_{w_i} j(w_i)$  从  $W_F \setminus W^*$  中选择下一个单词  $w_i$ . 利用  $W^* = W^* \cup \{w_i\}$  不断对特征词集合进行更新, 直到  $W_F = \emptyset$ . 特征词提取算法总结如算法1所示. 在之后的增量学习过程中, 对于新加入的句子, 没有必要处理所有的单词, 仅仅需要基于已经得到的特征词集合进行更新, 这样可以有效地降低计算复杂度.

### 算法1: 特征词提取

输入: 景观  $n$  以及点评句子集合  $S_n$ , 参数  $\lambda$

输出:  $W^*$

1.  $W_F = \{w \mid F(w, S_n) > \lambda\}$
2. 从  $W_F$  中挑选出  $F(w, S_n)$  最大的单词  $w$
3.  $W^* = W^* \cup \{w\}$
4. 重复
5. 通过求  $\arg \max_{w_i} \varphi(w_i)$ , 从  $W_F \setminus W^*$  中选择下一个单词  $w_i$
6.  $W^* = W^* \cup \{w_i\}$
7. 直到  $W_F = \emptyset$
8. 返回  $W^*$

## 2.3 旅游路线推荐

旅游路线是指游客从进入到离开景点所遵循的路线, 是一个包括若干景观名称的有序序列. 如果很多游客都选择同一条旅游路线, 那么该路线可以作为推荐信息呈现给用户. 游记中图像和上、下文是按照作者的写作顺序组织的, 在大多数情况下, 可以看作用户的旅游路线.

与关联规则挖掘中发现频繁项集的过程不同, 在这一步所处理的输入和输出都是有序的. 序列模

式挖掘与关联规则挖掘不同,可以从离散数据集中发掘频繁出现的有序事件或子序列.本文根据实际需求对传统序列模式挖掘算法进行改进,新的算法如算法 2 所示.

#### 算法 2:智能旅游路线推荐

输入:游记数据集  $T$ ,景观名称集合  $N$

输出:推荐游览路线  $r$

1. for 每一篇  $T$  中的游记  $a$
2.  $c \leftarrow \text{FM}(N, a)$
3. IF Judge( $c$ ) = True
4. 将  $C$  放入集合  $P$  中
5. end
6.  $L_1 = \{n | n \in N, \text{support}(n) \leq \text{min\_sup}\}$
7. for ( $k=2; L_{k-1} \neq \emptyset, k++$ ) do
8.  $C_k = \text{GC}(L_{k-1})$
9. for 每一条  $P$  中的路线  $c$  do
10.  $C_k$  中所有被包含在  $c$  中的候选者加 1
11.  $L_k = \{c | c \in C_k, \text{support}(c) \geq \text{min\_sup}\}$
12. end
13.  $r = L_k$  中的最长序列

算法 2 中,  $\text{FM}(N, a)$  根据景观集合  $N$  在游记  $a$  中利用模糊匹配算法得到一条旅游路线.  $\text{Judge}(c)$  用来对路线  $c$  进行判断,若  $c$  中不存在重复的景观且其中包含的景观数量大于等于景点内总景观数的  $1/3$ ,则返回 TRUE. 这是为了保证提取的路线足以完整描述一条景点内活动轨迹.  $\text{support}(n)$  用来计算支持度,  $\text{min\_sup}$  为最小支持度阈值.  $\text{GC}(L_{k-1})$  根据  $L_{k-1}$  利用连接和剪枝操作产生新的候选者.

#### 2.4 跨媒体信息关联

根据观察可知,用户通常会为游记中的图像  $I$  在相邻的位置给出简单的文字描述  $C$ . 通过提取景观句以及游记中图像的上、下文可以构建文本到图像的关联,采用多数表决的方式来为景观  $n$  选择对应图像. 具体过程如下.

图像聚类:首先提取游记中包含上下文  $c_I$  的图像  $I$ ,得到关于景点  $P$  的图像集合  $I_P$  和上下文集合  $C_P$ . 在集合  $I_P$  上利用谱聚类,基于视觉内容特征矢量分成视觉上的不同集群  $L_P = \{l_1, l_2, \dots, l_{|I|}\}$ . 若某个上下文  $c_I$  被标注为  $l_i$ ,则表示与其相邻的图像聚类的结果在  $l_i$  中.

投票选择图像聚类:对于  $S_n$  中的每个句子  $s$ ,根据余弦理论从  $C_P$  中查找最相似的上下文句子  $c_I$ ,投票决定可能关联的图像集群  $L_a$ . 该算法总结如算法 3 所示.

#### 算法 3:跨媒体信息关联

输入:景观  $n$  以及相应的  $S_n, C_P, L_P$

输出:  $L_a$

1.  $N(l_{c_I})$  初始化为 0
2. for  $S_n$  中的每一个  $S$  do
3. for  $C_P$  中的每一个  $c_I$  do
4.  $\text{score}(s, c_I) = \cos(s, c_I)$
5. end
6. 选择得分最高的  $c_I, N(l_{c_I}) = N(l_{c_I}) + 1 (l_{c_I} \in L_P)$
7. end
8. 对所有  $N(l_{c_I})$  排序,选择得票数最高的作为  $L_a$
9. 返回  $L_a$

可得与每个景观相关联的图像集群  $L_a$ ,借助亲和传播算法<sup>[18]</sup>从图像集群  $L_a$  为每个景观-特征刻画组合挑选有代表性的图像.

### 3 实验验证

通过大量实验来验证采用该方法能够更细粒度地刻画景点,并且提供的旅游路线能够满足不同用户的需求.

#### 3.1 数据集和实验设置

通过大众点评和蚂蜂窝获取国内 8 个流行景点数据. 使用各景点名作为搜索词,查询得到评论、游记和主要景观集合. 实验数据如表 1 所示.

大众点评的每一条评论都较短,一般只包含若干个简单的句子,缺乏足够的图像信息. 相比之下,蚂蜂窝的游记数据信息量较少、噪声更大,有着丰富的高品质图像及相应的上下文. 平均每一个景点可以得到 2 100 条评论和 3 100 篇游记. 首先对数据集进行预过滤,去除重复的评论、图像以及与景点不相关的游记等,同时限制图像必须要有相应的上下文信息,最后平均每个景点的数据集合包括 1800 多条

表 1 旅游信息数据集  
Tab. 1 Travel information dataset

景点	景观数	评论	游记
颐和园	20	6 415	2 504
故宫	10	5 256	1 775
云南丽江	15	1 543	4 511
张家界	15	1 116	3 024
九寨沟	12	1 102	4 021
大唐芙蓉园	15	1 023	412
庐山	10	588	2 018
黄山	12	443	2 860

评论和2 000+幅图像及上下文. 为了表示图像内容,为每一幅图像提取5种视觉特征构成一个809维向量,其中包含81维色矩、37维边缘直方图、120维小波纹理特征、59维LBP特征和512维GIST特征<sup>[19]</sup>. 根据实验比较可知,算法1中的参数 $\lambda$ 设置为0.6. 对于每个景观,选择排名前6的特征词来表示.

### 3.2 性能评价

如图2所示为蚂蜂窝官方提供的颐和园旅游路线,虽然路线数量很多,但这些路线只是简单地基于地理位置信息得到的,用户很难做出选择. 在现实中,多人选择同一路线表明该路线是最流行且用时更合理,这种推荐信息可以为用户提供可靠参考. 本文利用群智贡献的数据实现旅游路线推荐,从大量用户发布的旅游社交数据中挖掘出热度最高的路线作为推荐.

此外,相比于传统只针对单个数据源的研究,本文提出的方法从多个侧面对景观进行刻画,同时利用

- ①东宫门②仁寿殿③德和园④文昌院⑤玉澜堂、宜芸馆  
⑥乐寿堂⑦长廊⑧排云殿⑨佛香阁⑩石舫⑪耕织⑫如意门  
→①东宫门②仁寿殿③德和园④文昌院⑤玉澜堂、宜芸馆⑥乐寿堂  
⑦长廊⑧排云殿⑨佛香阁⑩石舫⑪十七孔桥⑫铜牛⑬新建宫门  
→①东宫门②仁寿殿③德和园④文昌院⑤玉澜堂、宜芸馆⑥乐寿堂  
⑦长廊⑧排云殿⑨佛香阁⑩苏州街⑪澹宁堂⑫谐趣园⑬东宫门  
→①东宫门②苏州街③澹宁堂④谐趣园⑤仁寿殿⑥德和园⑦玉澜堂  
⑧文昌院⑨乐寿堂⑩长廊⑪排云殿⑫佛香阁⑬石舫⑭耕织图⑮如意门  
→  
.....

图2 官方提供的颐和园旅游路线

Fig. 2 Official travel routes of Summer Palace

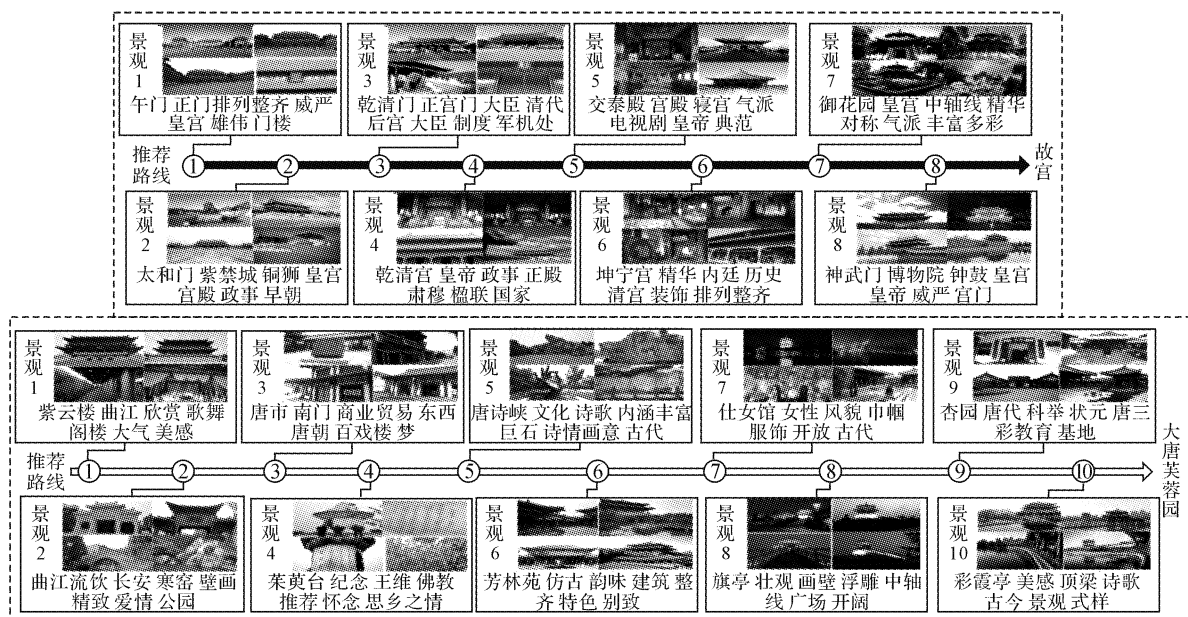


图3 景区刻画和旅游信息推荐

Fig. 3 Scenic spot profiling and travel information recommendation

表2 景观-图像关联准确率

Tab. 2 Accuracy rate of scenes and images correlation

景点	准确率/%	景点	准确率/%
颐和园	85.0	九寨沟	91.7
故宫	84.4	大唐芙蓉园	82.5
云南丽江	90.0	庐山	81.3
张家界	92.5	黄山	91.7

多源数据融合对景观-特征刻画组合进行可视化. 通过对实验结果进行人为鉴定,景观-特征刻画与图像关联的准确率统计如表2所示,平均准确率达到87.5%.

图3以北京故宫、西安大唐芙蓉园为例,展示最终结果. 故宫是国内最流行的景点之一,群智贡献的数据规模非常大,用以验证本文方法的正确性. 大唐芙蓉园更小众,官方提供的参考信息不足,可以通过群体智能丰富信息,进一步体现本文方法的实际效果. 为了评价该方法的效果,开展一个小范围用户调查,了解用户对这种呈现形式的体验感受. 要求参与者考虑3条准则:1)图像与景观-特征刻画之间的一致性(0表示不一致,5表示非常一致);2)挖掘得到的特征与景观的相关程度(0表示不相关,5表示非常相关);3)对这种可视化总结方式的满意度(0表示不满意,5表示非常满意). 邀请25名参与者进行调查,对结果进行平均之后得到的结果如图4所示. 本文的工作得到了积极的反馈,即使对于大唐芙蓉园这样小众的景点也可以得到令人满意的效果,进一步验证了该方法在旅游领域的潜力.

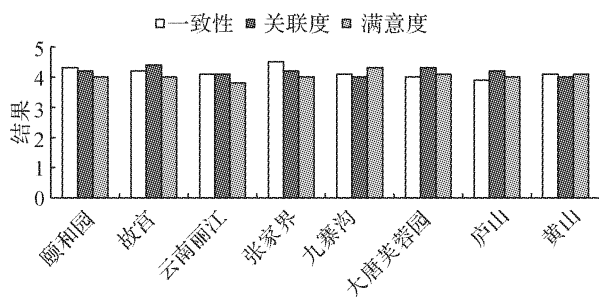


图4 用户调查结果

Fig. 4 User study results

## 4 结 语

本文提出的多源社交数据融合的多角度旅游信息感知方法可以更细粒度地刻画景点,最终结果能够满足不同用户的需求.在多角度旅游信息感知过程中,提出跨媒体多角度关联方法来连接碎片化旅游信息,然后利用序列模式挖掘从游记中挖掘典型旅游路线并进行智能推荐,最后采用投票方式实现跨媒体信息关联.实验结果表明,采用本文方法得到了令人满意的效果.

对于今后的工作,可以充分利用图像内容来改进文本到图像关联模块的性能,使得最终结果更准确.此外,在已推荐旅游路线的基础上挖掘更多有价值的信息是一个有趣的方向,可以更好地满足用户需求.

## 参考文献 (References):

- [1] ZHANG D, GUO B, YU Z. The emergence of social and community intelligence [J]. **Computer**, 2011, 44(7): 21-28.
- [2] GUO B, WANG Z, YU Z, et al. Mobile crowd sensing and computing: the review of an emerging human-powered sensing paradigm [J]. **ACM Computing Surveys**, 2015, 48(1): 1-31.
- [3] BIAN J, YANG Y, CHUA T S. Multimedia summarization for trending topics in microblogs [C] // **ACM International Conference on Conference on Information and Knowledge Management**. San Francisco: ACM, 2013: 1807-1812.
- [4] LIU B, HU M, CHENG J. Opinion observer: analyzing and comparing opinions on the Web [C] // **International Conference on World Wide Web**. Chiba: ACM, 2005: 342-351.
- [5] GONG Y, LIU X. Generic text summarization using relevance measure and latent semantic analysis [C] // **SIGIR '01: Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval**. New Orleans: ACM, 2001: 19-25.
- [6] RATTENBURY T, NAAMAN M. Methods for extrac-

- ting place semantics from Flickr tags [J]. **ACM Transactions on the Web**, 2009, 3(1): 1139-1141.
- [7] HAO Q, CAI R, WANG C, et al. Equip tourists with knowledge mined from travelogues [C] // **International Conference on World Wide Web**. Raleigh: DBLP, 2010: 401-410.
- [8] SHEN D, SUN J T, LI H, et al. Document summarization using conditional random fields [C] // **Proceedings of the International Joint Conference on Artificial Intelligence**. Hyderabad: DBLP, 2007: 2862-2867.
- [9] WANG T, BAI C. Understand the city better: multi-modal aspect-opinion summarization for travel [C] // **Web Information Systems Engineering-WISE**. Thessaloniki, Greece: Springer, 2014.
- [10] GUO T, GUO B, ZHANG J, et al. CrowdTravel: Leveraging Heterogeneous Crowdsourced Data for Scenic Spot Profiling and Recommendation [M] // **Advances in Multimedia Information Processing-PCM 2016**. Xi'an, China: Springer, 2016: 617-628.
- [11] LIN C Y, HOVY E. From single to multi-document summarization: a prototype system and its evaluation [C] // **Meeting on Association for Computational Linguistics**. Philadelphia: ACM, 2002: 457-464.
- [12] WAN X, YANG J. Multi-document summarization using cluster-based link analysis [C] // **International ACM SIGIR Conference on Research and Development in Information Retrieval**. Singapore: ACM, 2008: 299-306.
- [13] LI L, ZHOU K, XUE G R, et al. Enhancing diversity, coverage and balance for summarization through structure learning [C] // **International Conference on World Wide Web**. Madrid: DBLP, 2009: 71-80.
- [14] RADEV D R, JING H, STYS M, et al. Centroid-based summarization of multiple documents [J]. **Information Processing and Management**, 2004, 40(6): 919-938.
- [15] MOGHADDAM S, ESTER M. Opinion digger: an unsupervised opinion miner from unstructured product reviews [C] // **ACM Conference on Information and Knowledge Management**. Toronto: ACM, 2010: 1825-1828.
- [16] GUO H, ZHU H, GUO Z, et al. Product feature categorization with multilevel latent semantic association [C] // **ACM Conference on Information and Knowledge Management**. Hong Kong, China: ACM, 2009: 1087-1096.
- [17] HAGHIGHI A, VANDERWENDE L. Exploring content models for multi-document summarization [C] // **Human Language Technologies: Conference of the North American Chapter of the Association of Computational Linguistics**. Boulder: DBLP, 2009: 362-370.
- [18] FREY B J, DUECK D. Clustering by passing messages between data points [J]. **Science**, 2007, 315(5814): 972-976.
- [19] TORRALBA A, MURPHY K P, FREEMAN W T, et al. Context-based vision system for place and object recognition [C] // **IEEE International Conference on Computer Vision**. Beijing, China: IEEE, 2003: 273.