

引用格式: 桂志鹏, 胡晓辉, 刘欣婕, 等. 顾及地理语义的地图检索意图形式化表达与识别[J]. 地球信息科学学报, 2023, 25(6): 1186-1201. [Gui Z P, Hu X H, Liu X J, et al. Map retrieval intention formalization and recognition by considering geographic semantics[J]. Journal of Geo-information Science, 2023, 25(6): 1186-1201.] DOI: 10.12082/dqxxkx.2023.230019

顾及地理语义的地图检索意图形式化表达与识别

桂志鹏^{1,2,3,4*}, 胡晓辉^{1,5}, 刘欣婕¹, 凌志鹏¹, 姜屿涵², 吴华意^{2,3,4}

1. 武汉大学遥感信息工程学院, 武汉 430079; 2. 武汉大学 测绘遥感信息工程国家重点实验室, 武汉 430079; 3. 湖北珞珈实验室, 武汉 430079; 4. 地球空间信息技术协同创新中心, 武汉 430079; 5. 重庆市地理信息和遥感应用中心, 重庆 401127

Map Retrieval Intention Formalization and Recognition by Considering Geographic Semantics

GUI Zhipeng^{1,2,3,4*}, HU Xiaohui^{1,5}, LIU Xinjie¹, LING Zhipeng¹, JIANG Yuhan², WU Huayi^{2,3,4}

1. School of Remote Sensing and Information Engineering, Wuhan University, Wuhan 430079, China; 2. State Key Laboratory of Information Engineering in Surveying, Mapping and Remote Sensing, Wuhan University, Wuhan 430079, China; 3. Hubei Luoqia Laboratory, Wuhan 430079, China; 4. Collaborative Innovation Center of Geospatial Technology, Wuhan 430079, China; 5. Chongqing Geomatics and Remote Sensing Center, Chongqing 401127, China

Abstract: Mainstream map retrieval methods for spatial data infrastructures are mainly based on metadata text matching or image similarity calculation, but such approaches lack active perception and understanding of user retrieval intention, and in turn fail to truly meet user requirements. While existing intention recognition methods are incapable to express and recognize map retrieval demands with joint constraints of complex geographic concepts. To address this issue, this paper proposes a map retrieval intention formalization and recognition method by considering geographic semantics, aiming to improve the accuracy of map retrieval in an intention-driven and explainable manner by using relevance feedback samples. More specifically, a formalization model constrained by geographic ontology in the form of "intention-sub-intention-dimension component" is designed for expressing user's map retrieval intention. With the support of the formalization model, a recognition algorithm based on Minimum Description Length (MDL) principle and Random Merging (RM) strategy, named MDL-RM, is proposed by treating intention recognition as a combinational optimization problem. MDL-RM takes the description length of the sample set from relevance feedback as the optimization goal, merges samples randomly with the assistance of geographic ontologies and semantic similarities among geographic terminologies to generate sub-intention candidates, and searches the optimal intention using a greedy search approach. In order to evaluate the accuracy of recognized intention, we proposed a semantic metric, named Best Map Average Semantic Similarity (BMAS), and calculated it along with Jaccard index in five typical map retrieval scenes. Meanwhile, we analyzed the time cost and the influence of parameter settings and validated the

收稿日期 2023-01-13; 修回日期: 2023-02-21.

基金项目: 国家自然科学基金项目(42090011、41971349); 国家重点研发计划项目(2021YFE0117000)。[**Foundation items:** National Natural Science Foundation of China, No.42090011, 41971349; National Key Research and Development Program of China, No.2021YFE0117000.]

作者简介: 桂志鹏(1982—), 男, 宁夏吴忠人, 博士, 副教授, 主要从事高性能地理计算与时空大数据分析相关研究。
E-mail: zhipeng.gui@whu.edu.cn

effectiveness of random merge and sample augmentation strategy. The experimental results on the synthetic data demonstrate that the proposed method has higher accuracy and sample noise tolerance in most retrieval scenes comparing with the method based on Gene Ontology (RuleGO) and the Decision Tree learning method with Hierarchical Features (DTHF). The random merge strategy can reduce average computing time effectively without declining accuracy, and the sample augmentation strategy facilitates retrieval intention recognition even when the sample size is as low as 20. The proposed method is expected to be adapted and applied into geoportals and catalogue services to improve the service quality and user experiences upon the sharing and discovery of geographic information resources.

Key words: geographic information retrieval; intention formal expression; user relevance feedback; geographic ontology; semantic similarity; greedy search; minimum description length principle

***Corresponding author:** GUI Zhipeng, E-mail: zhipeng.gui@whu.edu.cn

摘要:主流地图检索方法多基于元数据文本匹配或图像内容相似度计算,缺乏对用户意图的主动理解,导致检索结果欠佳;而现有意图识别方法无法准确表达与识别复杂地理概念联合约束的地图资源检索需求。为此,本文提出一种顾及地理语义的地图检索意图形式化表达与识别方法,旨在利用相关反馈样本“感知”用户需求,以提升检索精度。该方法通过地理本体约束“意图-子意图-维度分量”模型的构建,实现检索需求的语义化描述;并将意图识别视为组合优化问题,基于最小描述长度准则、顾及地理概念从属关系的样本随机合并策略及贪心搜索实现最优意图识别。实验结果表明,相比基于频繁项集挖掘的RuleGO、决策树的DTHF算法,本文方法具有更高的识别准确度与噪声容忍度;随机合并策略可在不降低识别准确性的情况下有效缩短平均求解耗时;样本增强策略保证算法在样本规模仅为20时仍具有较高识别准确度。该方法可望应用于地理信息门户,提升各类地理信息资源共享与服务品质。

关键词:地理信息检索;意图形式化表达;用户相关反馈;地理本体;语义相似度;贪心搜索;最小描述长度准则

1 引言

地图是描述地理对象与现象时空分布、演化过程及其要素相互作用的有效工具,是地理信息展示与传播的载体,对社会经济发展决策及科学研究具有重要支撑作用^[1-2]。随着对地观测技术的进步、共享服务平台的建设及公众参与度的提高,符合制图规范且包含元数据描述的图幅、图层对象等地图资源无论在数量还是质量上均显著提升^[3]。以业界应用广泛的网络地图服务(Web Map Service, WMS)为例,互联网中存在超过4万条可用的WMS,所提供的30多万张地图图层涵盖气候、海洋、能源、地质等诸多主题^[4]。然而如此海量丰富的地图资源的利用并不充分,亟需一种高效准确的检索方法帮助用户从中定位和发掘兴趣资源。

目前地图资源检索主要包括元数据文本匹配与基于图像内容检索2种方式。前者通过匹配用户输入的检索词与资源元数据中标题、描述、提供者等字段实现检索。然而元数据规范、描述语言及著录习惯存在差异且检索词运用具有不确定性,文本字段匹配无法反映用户在语义层面的需求,导致检索错误或不全^[5]。利用地理本体中的概念与关系扩

展查询式^[5-7]或将文本匹配转化为概念语义相似度计算^[8-10]可实现顾及语义的地理信息检索,提高查全查准率。考虑到元数据字段可能存在缺失或“图文不符”的问题,有研究利用多模态信息增强元数据以提升匹配质量^[11-14]。针对元数据文本无法全面描述资源内容的问题,基于图像内容的检索从视觉特征层面辅助资源匹配^[15-17]。此外,优化资源缩略图质量^[3]、信息可视化方案及人机交互机制^[18]能够改善用户体验,提升检索效率。然而,上述检索方法及优化策略缺乏对用户需求的主动理解^[19],存在需求与检索条件之间的不匹配,即“意图鸿沟”问题。特别是各类地图资源涉及的领域知识庞杂,增加了检索词的构造难度;而地图的类型丰富、表达方式和制图风格多样,地理要素符号化和制图综合使得地图表达高度抽象化,基于图像内容的检索过度关注视觉特征细节却难以显式描述用户需求的类属概念,导致检索结果“只得其形、未得其意”。因此,地图精准检索不能仅停留在元数据文本与视觉相似性匹配层面,需要从意图识别角度理解用户检索需求。

基于意图的信息检索可为地图资源检索需求理解提供借鉴,主要包括意图的形式化表达与识别

2个方面。目前意图形式化表达多采用预定义类别与带权词项集合2种方式。基于预定义类别的方法在建立意图分类体系的基础上,利用决策树^[20]、SVM^[21]或深度学习^[22]分类模型实现意图识别。然而类别仅能表示检索目标的大致范围,无法刻画用户需求的具体内容。对此,结合检索结果与隐式反馈数据提取带权词项集合形式的需求描述,可综合表达意图内容及其强度分布^[23-25],但此方法未考虑词项语义及逻辑关系,难以对用户兴趣资源进行精确定位。Zhang等^[26]设计“概念-属性”语义层次化意图模型,并结合概念分类器与属性分类器构建用户意图;然而该模型仅关注资源描述的内容(如图像中的对象),未考虑其他层面的用户偏好(如提供者)。Fariha等^[27]基于溯因推理,利用检索结果示例集合推断由约束条件及其逻辑关系构成的查询式,以辅助非专业用户表达检索意图,但该方法仅能推断逻辑“与”关系,无法应对一次检索任务中涉及多个子意图(即逻辑“或”)的情况。Zhang等^[28]结合伪相关反馈与编码器-解码器框架生成自然语言形式的意图描述,但自然语言无法直接用于检索且编码器-解码器训练需要大量数据,适用场景受限。综上,由于上述形式化表达模型未显式建模意图维度及子意图,无法体现地图资源检索需求多意图多维属性约束的特点;同时,现有意图识别方法大多将反馈数据中提取的词项视为独立约束项,未利用地学概念间的从属关系捕获反馈数据的语义关联,因而无法确定意图各维度约束的概念层级(例如,用户需要的地图其关键地理要素类型是“河流”,还是包含“河流”、“湖泊”、“冰川”的上位概念“水体”?),导致识别结果难以准确反映用户检索需求。

针对上述传统地图检索方法未显式考虑用户意图及现有意图驱动的检索方法在形式化表达与语义利用方面的不足,本文构建“意图-子意图-维度分量”模型表达意图,并提出一种基于最小描述长度(Minimum Description Length, MDL)准则与随机合并(Random Merging, RM)策略的地图检索意图识别算法MDL-RM。该算法以用户相关反馈样本为输入,使用地理本体中的概念作为意图维度分量取值,并利用概念从属关系合并相关反馈样本,为意图形式化表达与候选子意图生成提供语义支撑;将意图识别视为组合优化问题,采用贪心算法最小化样本集合的编码长度以搜索最优意图,并使用样本随机合并策略加速迭代,实现地图检索意图的显式表达与识别。

2 地图检索意图形式化表达模型设计

地图检索意图识别需要设计简洁且具有一定表达能力的模型描述用户检索需求及约束条件,包括子意图、意图维度分量及逻辑关系。用户在一次检索任务中往往具有多种需求,如将“海水温度”与“海水盐分”2种地理要素同时作为检索目标,即子意图;并且每种需求可以从多个维度综合表达,如“美国水体”涉及空间范围与地理要素2个维度(图1)。因此,用户意图表达模型应包含子意图与维度分量2种元素。在逻辑关系方面,各子意图之间为逻辑“或”关系,子意图各维度分量之间为逻辑“与”关系,即检索结果应至少符合一个子意图上所有维度分量的约束条件。此外,用户可能明确表示不需要某类地图(即逻辑“非”关系),该情况本文暂未考虑。

基于上述分析,本文构建一种“意图-子意图-维度分量”三层嵌套的地图检索意图模型,其树结构表示如图2所示。用户地图检索意图 I 由若干个子意图组合而成,如式(1)所示。各子意图包含若干维度分量,如式(2)所示。

$$I = \bigvee_{k=1}^m I^k \quad (1)$$

$$I^k = \bigwedge_{i=1}^d v_i^k \quad (2)$$

式中: m 为子意图个数; I^k 为第 k 个子意图; \vee 表示逻辑“或”关系; d 为维度分量个数; v_i^k 为 I^k 中第 i 个维度的取值; \wedge 表示逻辑“与”关系。

各维度分量选用地理本体中的概念作为取值来源,以便利用概念从属关系表达不同抽象层次的需求,实现检索范围的精准约束。例如,本文使用地球与环境术语语义网(Semantic Web of Earth and Environmental Terminology, SWEET)中的概念“河流”、“植被”、“土地利用”、“气候”、“土壤”等作为地理要素维度取值。若某个维度取值 v_i^k 为对应本体的根节点,如地理要素维度取值为SWEET根节点“事物”,则认为用户在此维度无偏好;若所有维度取值均为根节点,则此意图不具备检索约束能力,本文将其视为无意图。

意图维度分量的选取影响检索约束条件表达的范围。本文参考空间数据元数据^[29-31]与专题地图编制^[32]规范,并结合常见的地图检索需求,构建层

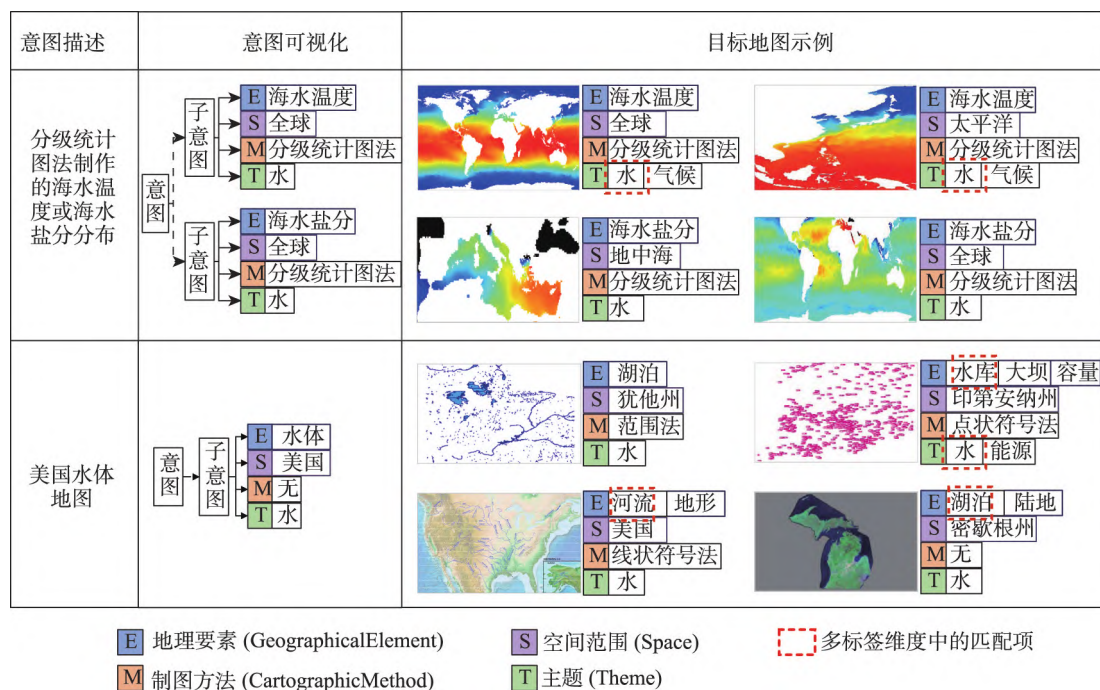


图1 地图检索意图与目标地图示例

Fig. 1 Examples of map retrieval intentions and corresponding target maps

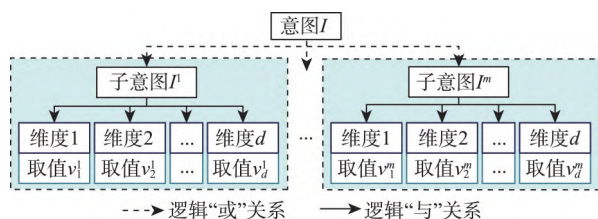


图2 地图检索意图的树结构表示

Fig. 2 Tree-based map retrieval intention representation

次化的地图检索意图维度树(图3)。该维度树包括“内容”、“空间”、“时间”、“制图”和“其他”5个顶层维度,涉及地图的数据采集、制图与发布使用等阶段,以支持复杂多样的检索需求表达,同时为地图内容的描述提供参考框架。在实际应用中,可结合具体检索场景从该维度树中选择合适的子树作为意图维度分量的集合。

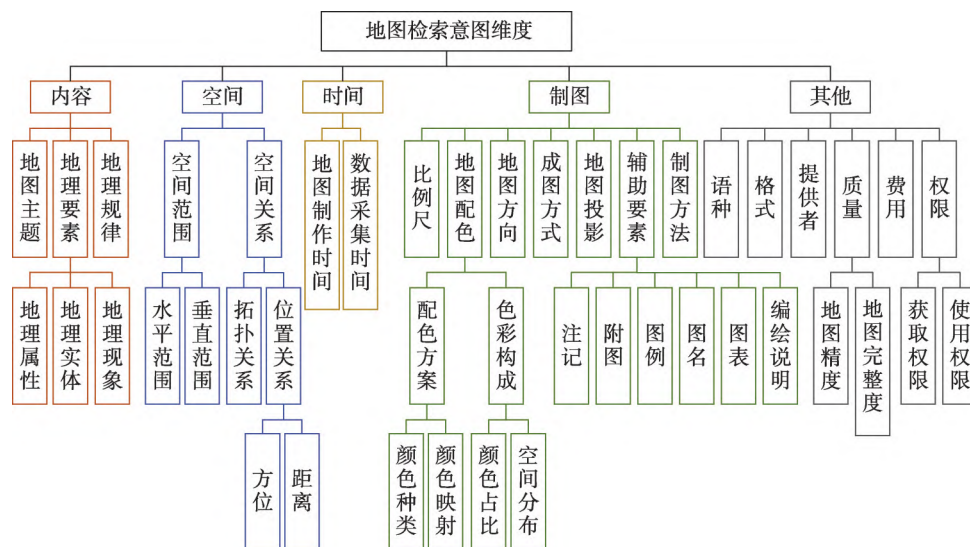


图3 地图检索的潜在意图维度

Fig. 3 Potential intention dimensions for map retrieval

3 基于MDL准则与随机合并策略的意图识别方法

意图识别可视为在意图空间中搜索最符合用户检索需求的意图方案的过程,故本文将其作为组合优化问题进行求解。在意图识别中,检索语句、检索结果与相关反馈样本集合均可作为用户需求载体。其中,相关反馈样本集合基于用户对初始检索结果与检索目标的相关性判断传递需求(相关的为正样本,不相关的为负样本),蕴含更为丰富的偏好信息,可隐式表达检索意图^[33]。MDL准则是一种模型筛选依据,常用于贝叶斯网络结构学习^[34]与规则归纳^[35]。它能够平衡模型复杂度与描述能力^[36],抑制语义等价但形式过于复杂的意图识别结果。因此,本文在意图形式化表达的基础上,基于MDL准则将反馈样本集合的编码长度作为目标函数,采用顾及地理语义样本随机合并策略生成候选子意图以缩小搜索范围,并使用贪心算法迭代求解最优意图。在检索流程中,本文算法以现有检索方式为前置捕获用户需求,并提供识别结果用于检索精化,如图4所示。

3.1 基于MDL准则的意图识别目标函数定义

在满足准确性的前提下,简洁的意图表达方式通常更符合用户思维习惯^[37]。MDL准则认为对于一组描述数据的模型,能产生最多数据压缩效果的模型最优,故而兼顾简洁性与准确性。因此,本文将相关反馈样本集合视为待压缩数据,将意图视为描述数据的模型,通过设计反馈样本集合编码方案及编码长度计算方法(采用信息传输与编码理论中的香农熵),搜索最短编码方案作为意图识别结果。

3.1.1 函数定义

由于本文采用统一的意图形式化表达模型描

述各种意图,基于MDL准则设计目标函数时无需考虑不同意图间的表达差异,故可使用二分编码^[34]定义样本集的编码长度,如式(3)所示。

$$L(S, I) = L(I) + L(S|I) \quad (3)$$

式中: $L(S, I)$ 为给定意图 I 时反馈样本集合 S 的编码长度(即总编码长度); $L(I)$ 与 $L(S|I)$ 分别为意图编码长度与给定意图后的样本集合编码长度。其中意图编码长度由子意图数量的编码长度与各子意图编码长度构成^[36,38],如式(4)与式(5)所示。

$$L(I) = L_N(m+1) + \sum_{k=1}^m L(I^k) \quad (4)$$

$$L(I^k) = \log\left(\prod_{i=1}^d |C_i|\right) \quad (5)$$

式中: m 为子意图数量; $L_N(m+1)$ 为子意图数量编码长度; $L(I^k)$ 为第 k 个子意图 I^k 的编码长度; d 为意图维度数量; C_i 为第 i 个维度对应本体的概念集合; $|*|$ 表示集合中元素的数量。由于子意图数量的上界未知,故计算 m 的编码长度时采用不需要上界先验的正整数通用编码^[39]。在正整数通用编码中,某正整数 x 的编码长度 $L_N(x) = \log(x_0) + \log(x) + \log(\log(x)) + \dots$, 其中 $x_0 \approx 2.865\ 064$, 省略项指从 $\log(x)$ 开始依次对前一项取对数得到的所有非负项。由于本文需要计算无意图(即 $m=0$)时的总编码长度,故使用 $m+1$ 作为实际编码的正整数。为避免子意图出现概率差异所导致的子意图识别率失衡,本文采用等长编码。

给定意图后的反馈样本编码长度 $L(S|I)$ 为各子意图覆盖的样本集合 S_k 、 S_k 中正样本数量及剩余样本集合的编码长度之和,如式(6)所示。

$$L(S|I) = \sum_{k=1}^m (|S_k| L_{avg}(S_k|I^k) + \log(|S_k|)) + |S_r| L_{avg}(S_r|I) \quad (6)$$

式中: S_k 为子意图 I^k 覆盖的样本子集; $L_{avg}(S_k|I^k)$ 为 S_k 中样本的平均编码长度; $\log(|S_k|)$ 为 S_k 中正样

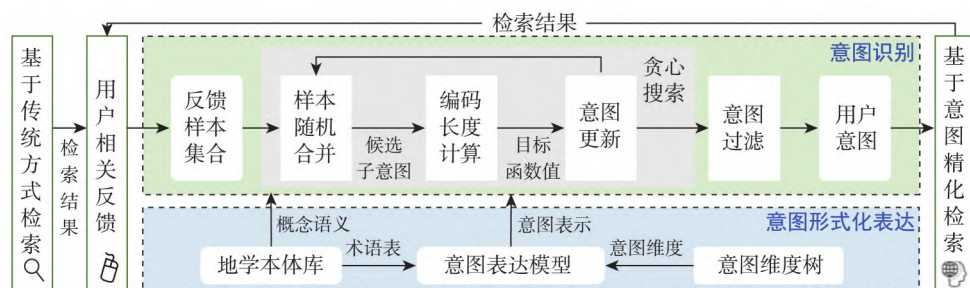


图4 基于MDL准则与随机合并策略的意图识别算法在地图检索中的工作流程

Fig. 4 Overall map retrieval workflow enhanced by an intention recognition method based on MDL principle and random merge strategy

本数量的编码长度(该项用于确保样本编码无损), $L_{avg}(S_r|I)$ 为未被覆盖的剩余样本集合 S_r 中样本的平均编码长度。本文假设反馈样本在各意图维度已标注有若干标签,则样本覆盖的判定规则为若样本某维度存在语义上等价或从属于某子意图对应维度分量取值的标签,则认为样本在此维度满足该子意图;若所有维度均满足该子意图,则该样本被该子意图覆盖。由于用户与检索系统已知各样本标签取值,仅需编码各样本的正负性即可描述反馈样本集合。本文对反馈样本集合采用算术编码,依据香农无噪声编码定理^[40], S_k 与 S_r 中样本正负性的平均编码长度如式(7)。

$$L_{avg}(S_k|I^*) = - \sum_{y \in \{+, -\}} \frac{|S_k^y|}{|S_k|} \log \left(\frac{|S_k^y|}{|S_k|} \right) \quad (7)$$

式中: S_k 与 I^* 分别指代 S_k 与 I^k 或 S_r 与 I ; y 为样本的正负性,即 S_k^+ 与 S_k^- 分别表示 S_k 中正、负样本

集合。本文设定各式中采用以2为底的对数,故编码长度单位为bit。

样本集合编码方案基本工作原理如图5所示。图6展示了针对表1所示相关反馈样本集,给定用

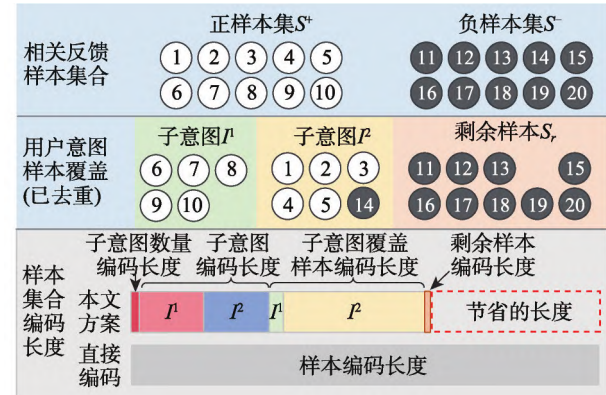


图5 样本集合编码方案工作原理

Fig. 5 Coding mechanism of relevance feedback sample set

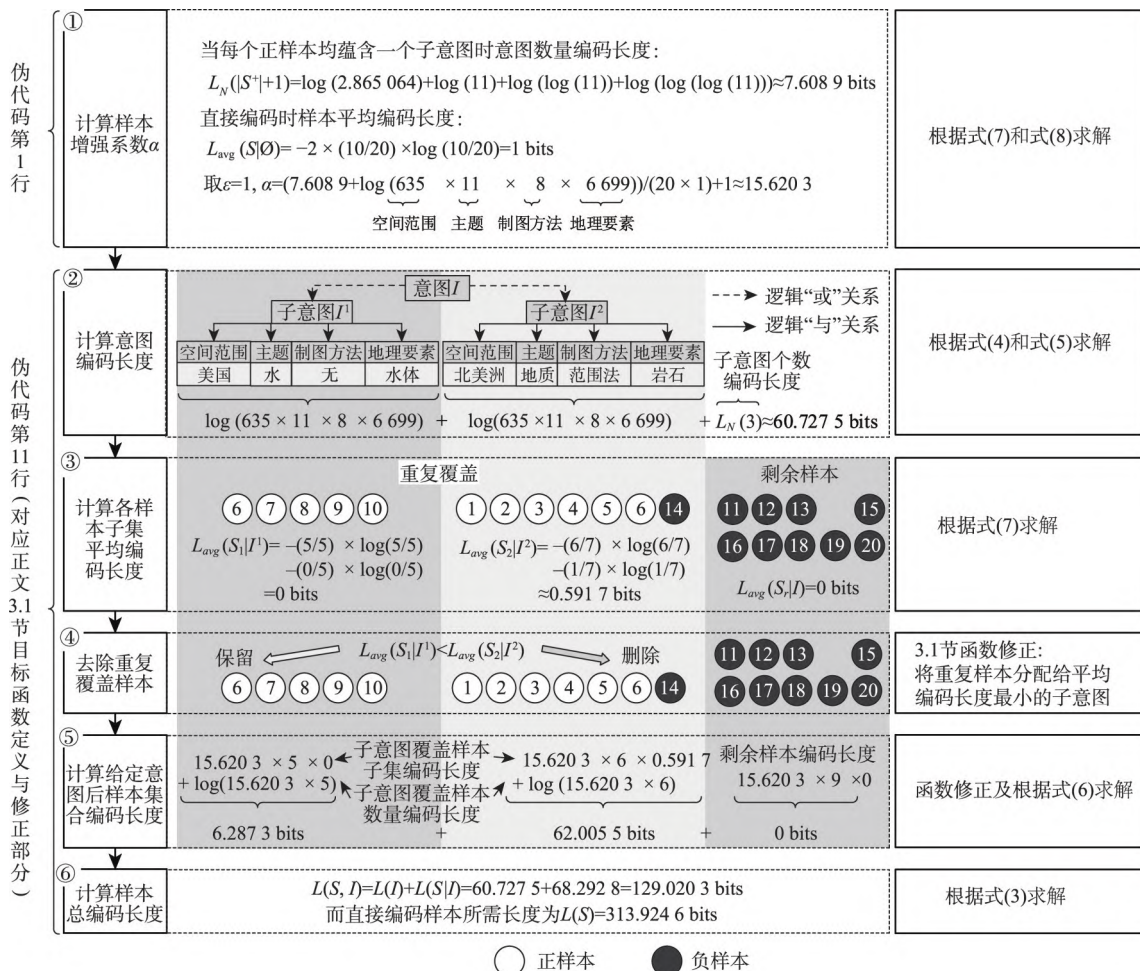


图6 针对给定意图的相关反馈样本集合总编码长度计算示例

Fig. 6 Illustration of coding length calculation of relevance feedback sample set under given retrieval intention

表1 相关反馈样本集合示例
Tab.1 Example of relevance feedback sample set

正样本					负样本				
编号	空间范围	主题	制图方法	地理要素	编号	空间范围	主题	制图方法	地理要素
1	加拿大	地质	范围法	岩石	11	巴西	地质	范围法	沉积岩
2	美国	地质	范围法	变质岩	12	智利	地质	范围法	混合岩
3	北美洲	地质	范围法	沉积岩	13	缅甸州	水	分级统计图	隔水层
4	内华达州	地质	范围法	斜长石	14	美国	地质	范围法	沉积岩
5	墨西哥	地质	范围法	橄榄岩	15	佛罗里达州	水	范围法	防洪堤
6	加利福尼亚州	水、地质	点状符号法、范围法	泉水、沉积岩	16	印第安纳州	生物多样性	分级统计图	鸟类
7	犹他州	水	质底法	湖泊	17	南美洲	水	分级统计图	降水
8	内华达州	水	范围法	湿地	18	加拿大	水	线状符号法	河流
9	佛罗里达州	水	范围法	水库	19	北美洲	地质	范围法	矿物
10	美国	水	线状符号法	河流	20	巴西	地质	分级统计图	土壤

户意图I的样本总编码长度 $L(S,I)$ 的计算过程。其中空间范围、主题、制图方法与地理要素4个意图维度的本体概念数量依次为635、11、8与6699,完整的意图识别计算过程见GitHub^①。由于意图I中2个子意图可较好地覆盖正样本并排除负样本(如图6第②、③步骤对应行所示),故基于意图的编码方案虽然引入了意图编码长度 $L(I)$,但相较直接编码方案能有效压缩样本编码长度 $L(S|I)$,从而缩短总编码长度。此外,由式(4)与式(7)可知本文方案中 $L(I)$ 与子意图数量(意图复杂程度)正相关, $L(S|I)$ 与各样本子集中正负样本比例差异(意图准确程度)呈负相关关系,故理论上最小化 $L(S,I)$ 即可得到兼顾简洁性与准确性的意图识别结果。

3.1.2 函数修正

为保证MDL准则的适用性,须对样本被多个子意图覆盖、反馈样本不足、用户误选等特殊情况进行处理。当样本被多个子意图覆盖时其编码长度将被重复计算,可能导致无法通过编码长度准确比较候选意图优劣。为此,本文将每个样本唯一分配到某一子意图以消除重复覆盖。分配规则首先将样本重复放入覆盖它的所有子意图样本集,计算出各子意图覆盖样本集的平均编码长度,然后将样本分配到平均编码长度最小的子意图样本集,即将式(6)中 S_k 修正为 $S_k - S_{k_exclude}$,其中 $S_{k_exclude}$ 为未在 S_k 中取得最小平均编码长度的重复样本集合。

反馈样本数量不足将使得样本直接编码长度小于基于意图的编码长度,导致意图无法识别。为

此,本文设计样本增强系数 α 对原始样本进行复制。 α 设置的依据为即使每个正样本均蕴含一个子意图,最小化编码长度依然可识别所有子意图,如式(8)所示。

$$\begin{aligned} \alpha|S|L_{avg}(S|\phi) &> L_N(|S^+|+1) + |S^+|\log\left(\prod_{i=1}^d|C_i|\right) \\ L_N(|S^+|+1) + |S^+|\log\left(\prod_{i=1}^d|C_i|\right) & \\ \Leftrightarrow \alpha = \frac{L_N(|S^+|+1) + |S^+|\log\left(\prod_{i=1}^d|C_i|\right)}{|S|L_{avg}(S|\phi)} + \varepsilon \end{aligned} \quad (8)$$

式中: $L_{avg}(S|\phi)$ 为直接编码样本的平均编码长度; ε 为任一大于0的实数。样本增强后,式(6)中给定意图后的样本集合编码长度 $L(S|I)$ 将被修正为 $\sum_{k=1}^m(\alpha|S_k|L_{avg}(S_k|I^k) + \log(\alpha|S_k|)) + \alpha|S_r|L_{avg}(S_r|I)$ 。然而,若用户在反馈过程中存在误选,上述样本增强策略将得到错误的子意图。由于误选的样本占有所有正样本的比例通常较小,本文规定仅保留正样本覆盖比例大于一定阈值 β 的子意图。该阈值将在实验部分讨论。

3.2 基于样本语义合并的候选子意图生成

本文基于正反馈样本提取候选子意图以缩小最优意图搜索范围。为顾及正样本之间的语义关联,本文利用地理本体中的概念从属关系对正样本的标签进行概念泛化生成各维度分量取值,并用逻辑“与”关系组合维度分量得到候选子意图,具体步骤如图7所示。若2个样本在某维度均为单标签且取值相同,则使用此标签作为子意图在该维度的取值;若取值不同则取2个标签在本体中的最近公共

① https://raw.githubusercontent.com/ZPGuiGroup/Whu/Map-Retrieval-Intention-Recognition/master/MDL_RM/Demo%20of%20Iterative%20calculation%20for%20intention%20recognition.pdf。

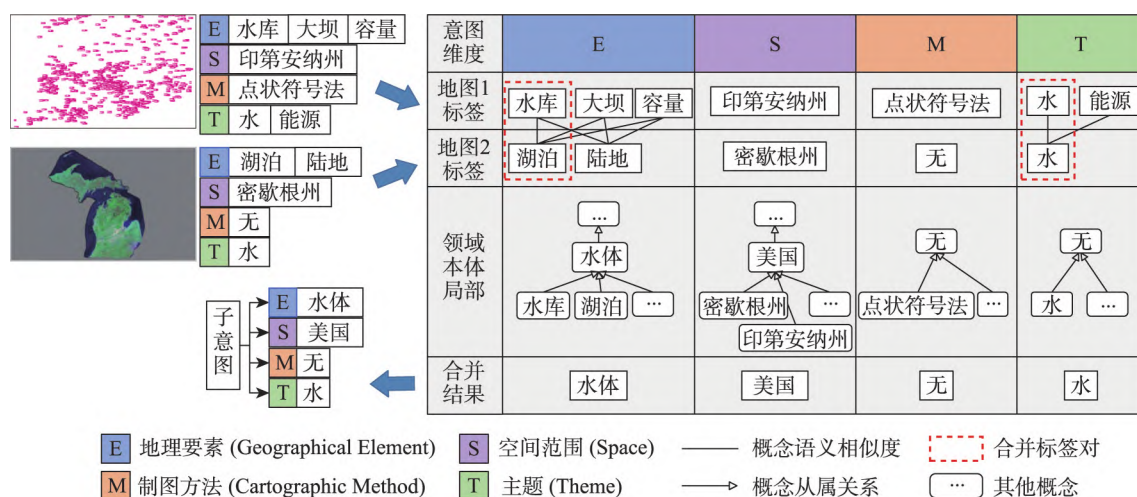


图7 基于反馈样本语义合并的子意图生成

Fig. 7 Sub-intention generation based on semantic merging of feedback samples

祖先 (Lowest Common Ancestor, LCA) 作为取值。如图7中2个样本的空间范围标签“印第安纳州”与“密歇根州”在地名本体 GeoNames 中的 LCA 为“美国”，故将“美国”作为子意图在该维度的取值。若样本在某维度为多标签且取值不同，为避免无关标签的影响，则取具有语义相似度最大的任一标签对的 LCA 作为子意图取值。语义相似度的计算采用改进的 Lin 相似度^[41]，如式(9)所示。

$$Sim(c_1, c_2) = \begin{cases} 1 & c_1 = c_2 \\ \frac{2IC(LCA(c_1, c_2))}{IC(c_1) + IC(c_2)} & \text{其他} \end{cases} \quad (9)$$

式中： c_1 与 c_2 为本体概念； $LCA(c_1, c_2)$ 为其最近公共祖先； $IC(*)$ 为概念的信息量，计算方法参考文献^[42]。

由于地理概念之间的从属关系具有层次性，样本合并生成的子意图可能不具有足够的样本覆盖能力，为此需要进一步将候选子意图视为所有维度均为单标签的样本，与样本或其他候选子意图合并。如某3个正样本“地理要素”维度的标签依次为“时令河”、“常年河”与“湖泊”，若前2个样本已合并得到子意图“河流”，只有进一步与第三个样本合并才能得到真实意图“水体”。

3.3 基于随机合并与贪心搜索的意图识别

在上述候选子意图生成方法的基础上，本文通过随机合并策略缩小搜索空间，并采用贪心算法求解意图以提升算法效率。搜索能力与耗时是影响组合优化算法选取的重要因素。针对正样本数量较多时遍历所有候选子意图生成方案难以满足实

时性要求的问题，本文采用随机合并策略生成一定数量的候选子意图，并保留使总编码长度最小的候选子意图。同时，由于地图检索场景中的子意图数量通常较少，无需全局遍历子意图组合方案即可得到最优解，本文使用贪心算法搜索候选组合方案作为意图识别结果。算法具体流程如算法1所示。首先，初始化子意图集合 R 、总编码长度 L 与剩余样本 S_r ，并计算样本规模增强系数 α 。然后，迭代以

算法1：基于最短描述长度准则与随机合并策略的地图检索意图识别算法(MDL-RM)

输入：反馈样本集合 S ，各意图维度本体 $O=[O_1, O_2, \dots, O_d]$ ，随机合并数量 t_{merge} ，正样本覆盖比例阈值 β
 输出：地图检索子意图集合 R

- 1 calculate α ；/*计算样本增强系数*/
- 2 $R \leftarrow \phi$, $L \leftarrow L(S, \phi)$, $S_r \leftarrow S$ ；/*初始化子意图集合、总编码长度与剩余样本集合*/
- 3 repeat
- 4 $C \leftarrow []$ ；/*初始化候选子意图集合数组*/
- 5 for i from 0 to $t_{merge} - 1$ do
- 6 $a, b \leftarrow \text{random_select}(S_r \cup R)$ ；
- 7 $subI \leftarrow \text{merge}(a, b, O)$ ；/*生成候选子意图*/
- 8 $candR \leftarrow R \cup subI$ ；/*生成候选子意图集合*/
- 9 $C[i] \leftarrow candR$ ；
- 10 end for
- 11 $bestR \leftarrow \arg \min_{candR \in C} L(S, candR)$ ；
/*获取使得总编码长度最小的候选子意图集合*/
- 12 $L_{min} \leftarrow L(S, bestR)$ ；
- 13 if $L_{min} < L$ then
- 14 $R \leftarrow bestR$, $L \leftarrow L_{min}$, $S_r \leftarrow S_r - S_{len(R)}$ ；
- 15 end if
- 16 until $L_{min} \geq L$ ；
- 17 $R \leftarrow \text{filter}(R, \beta)$ ；/*过滤 R 中正样本覆盖比例低于 β 的子意图*/
- 18 return R

下步骤直至无法生成使 L 更小的子意图:① 初始化候选子意图集合数组 C ;② 从子意图集合 R 与剩余正样本集合 S_r^+ 的并集中随机选取2个元素 a 与 b , 合并生成候选子意图 $subI$, 将 $subI$ 加入 R 得到候选子意图集合 $candR$, 此步骤重复 t_{merge} 次;③ 将 R 更新为使 L 减小最多的候选子意图集合, 并更新 L 与 S_r 。最后, 去除 R 中覆盖的正样本占比小于 β 的子意图并输出。

该算法可结合迭代反馈逐步精化意图识别结果。具体而言, 以当前意图识别结果为查询条件获取新的检索结果与用户反馈, 并将此次反馈样本与已有反馈样本合并去重作为新的反馈样本集合, 输入算法得到精化的意图识别结果。

4 实验结果与分析

本文从算法可行性、与基准算法的性能对比、随机合并与样本增强策略有效性及参数敏感性5个方面对 MDL-RM 进行综合验证。具体而言, 实验将:① 验证基于 MDL 准则识别意图的可行性;② 通过对比实验分析本文算法在准确性、鲁棒性与时间效率3个方面的优势与不足;③ 分析不同合并策略对耗时与意图识别准确性的影响;④ 验证样本增强策略在小样本意图识别任务中的有效性;⑤ 探讨参数取值对算法性能的影响。除参数影响分析实验外, 其余实验均设定随机合并数量为50次, 子意图覆盖正样本比例阈值为0.3。实验环境为配置 Intel i7 十二核处理器(主频 3.20 GHz)和 16 GB 内存的台式机, 操作系统为 Ubuntu 16.04, 算法基于 Python 3 实现。为保证结果的稳定性与可靠性, 各实验在样本集合上重复50次并统计平均值。

4.1 实验数据

本文选取“空间范围”、“地理要素”、“制图方法”与“主题”4个维度开展实验。“空间范围”与“地理要素”对应本体分别选取 GeoNames^[43]局部及 SWEET^[44], “制图方法”与“主题”采用基于文献^[12,32]构建的2层本体。由于现有地图样本标签质量欠佳且缺乏用户相关反馈数据, 而人工标注成本高、耗时长, 本文采用合成数据开展实验。合成数据基于上述4个维度对应本体与预定义意图生成。具体步骤为:① 随机组合各维度取值生成规模为

10 000 000 的单标签样本库;② 设计“无意图”、“单意图单维度”、“单意图多维度”、“多意图单维度”与“多意图多维度”5种意图场景, 并为各场景预定义30种检索意图;③ 基于各预定义意图, 从样本库中分别选取意图覆盖与未覆盖的100个样本作为正、负反馈样本;④ 为各样本集添加4个级别的反馈噪声与6个级别的标签噪声, 最终得到 $5 \times 30 \times 4 \times 6 = 3600$ 个样本集。其中反馈噪声用于模拟相关反馈中的样本误选和漏选现象, 通过交换无噪声样本集合中指定比例的正负样本生成; 标签噪声指正样本中的无关标签, 通过为指定比例的正样本在各维度添加1至4个无关标签生成。

4.2 评价指标

本文从准确性与效率2个方面对算法进行评价。在准确性方面, 使用 Jaccard 系数描述意图识别结果与预定义意图的样本覆盖一致性程度; 算法效率通过意图求解耗时度量。由于 Jaccard 系数无法区分样本覆盖相同但子意图数量与维度分量取值不同的意图识别结果, 本文定义最佳映射平均语义相似度指标 (Best Map Average Semantic Similarity, BMASS), 从语义角度评价意图准确性。Jaccard 系数与 BMASS 取值范围均为 $[0, 1]$ 且值越大准确度越高, 各指标定义如式(10)一式(12)所示。

$$Jaccard(A, B) = \frac{|S_A \cap S_B|}{|S_A \cup S_B|} \quad (10)$$

$$BMASS(A, B) = \max_{f \in F} \frac{1}{|B|} \sum_{k=1}^{|A|} Sim(I^{k_1}, f(I^{k_2})) \quad |A| < |B| \quad (11)$$

$$Sim(I^{k_1}, I^{k_2}) = \frac{1}{d} \sum_{i=1}^d Sim(v_{Ai}^{k_1}, v_{Bi}^{k_2}) \quad (12)$$

式中: S_A 、 S_B 分别为子意图集合 A 、 B 覆盖的样本集合; f 为 A 到 B 的一个单射; F 为所有单射的集合; I^{k_1} 、 I^{k_2} 分别为 A 、 B 中的第 k_1 、 k_2 个子意图; d 为意图维度数量; A 、 B 中第 k_1 、 k_2 个子意图在第 i 个维度上的取值 $v_{Ai}^{k_1}$ 、 $v_{Bi}^{k_2}$ 的相似度 $Sim(v_{Ai}^{k_1}, v_{Bi}^{k_2})$ 由式(9)计算。针对未识别出子意图的情况, 为了计算 Jaccard 系数与 BMASS, 需将意图各维度取值设置为本体根节点(即无意图)。

4.3 算法可行性分析

本实验通过分析子意图搜索迭代过程中编码长度、BMASS 及 Jaccard 系数的变化, 验证基于

MDL准则识别意图的可行性。在除“无意图”场景外的2 880个样本集上重复50次实验统计迭代频次分布,可知最大迭代次数未超过6次且主要集中于3次以内,90%以上的样本集重复实验迭代后BMASS均有所增加(图8(a))。因此,本文分别计算迭代次数为1至3的样本集在迭代过程中各指标的平均值,如图8(b)—图8(d)所示。

由实验结果可知,迭代过程中随着给定意图下的样本编码长度、总编码长度的减少,意图编码长度、BMASS及Jaccard系数逐渐增加。意图编码长度增加说明算法识别出新的子意图或原有子意图合并后覆盖的正样本增加;给定意图下的样本编码长度减少说明各子意图覆盖的样本集合及剩余样本集合中正、负样本比例的差异增大,子意图对其覆盖样本的描述能力增强;BMASS与Jaccard系数增加表明随着迭代次数增加,意图识别结果与预定义意图的样本覆盖一致性与语义相似度均得以提升。因此,基于MDL准则,通过最小化反馈样本编码长度识别用户检索意图是可行的。

4.4 意图识别性能对比

本实验以频繁项集挖掘算法RuleGO^[45]、决策树算法DTHF^[46]为基准,分析MDL-RM在意图识别准确性、算法效率及鲁棒性上的优势与不足。RuleGO利用频繁项集挖掘中的最小支持度过滤生成的项集,能够有效剔除不显著的候选子意图;DTHF利用数据的层级特征定义决策树分裂规则,可处理具有层次语义特征和缺失值的数据。上述对比算法均具有处理反馈噪声、标签噪声与概念间语义从属关

系的能力,可保证算法可比性。为验证算法鲁棒性,实验统计了不同意图场景、反馈噪声与标签噪声级别的3 600个样本集上各指标的平均值。本文结合BMASS取值与求解耗时进行调参,设置随机合并数量为50次,正样本覆盖比例阈值 β 为0.3;RuleGO显著性阈值为0.05,最小支持度阈值为40,单个意图包含的最大词项数为4。图9为评价指标统计结果,图10为意图识别结果示例。

从图8中可以看出,与RuleGO、DTHF相比,本文方法总体更为准确与鲁棒,但求解耗时波动更大。图9(a)中3种算法的BMASS指标均随意图复杂度与噪声比例的提高而降低,但MDL-RM识别的意图相对较为简洁且子意图数量与预定义意图更为接近(图10),从而在大多数场景下具有更高的BMASS值(图9(b))。在标签噪声与反馈噪声容忍度方面,本文算法最优,RuleGO次之,DTHF最差,且反馈噪声对算法的影响大于标签噪声。可以发现RuleGO在部分无意图场景下的准确性优于本文算法(图9(b)中虚线框标识部分),这是因为标签噪声可能增加无意图场景中样本标签的规律性,进而提高低样本占比意图出现的概率,导致算法误将噪声识别为意图。而实验中RuleGO最小支持度阈值设置为40,即仅保留正样本占比大于0.4的子意图,因此对错误子意图的过滤比本文算法更为有效。就反馈噪声而言,其对基准算法RuleGO与DTHF识别准确性的影响大于本文算法。当反馈噪声非零且标签噪声大于0.4时,RuleGO的准确度较差,而MDL-RM在标签噪声小于0.6时仍能保证较高的意图识别准确度。这是

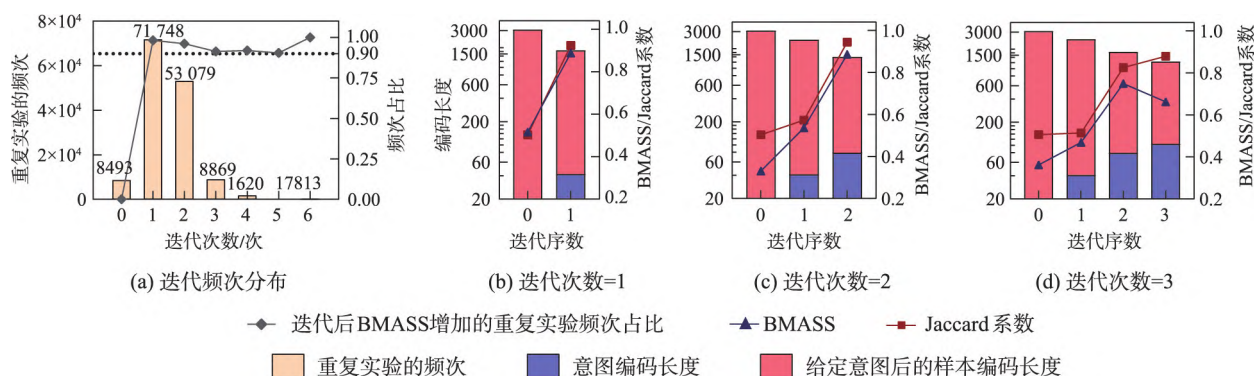


图8 样本集重复实验的迭代频次分布及编码长度、BMASS与Jaccard系数均值的变化

Fig. 8 Frequency distribution of iterations on sample sets and changes of the averages of encoding length, BMASS and Jaccard index

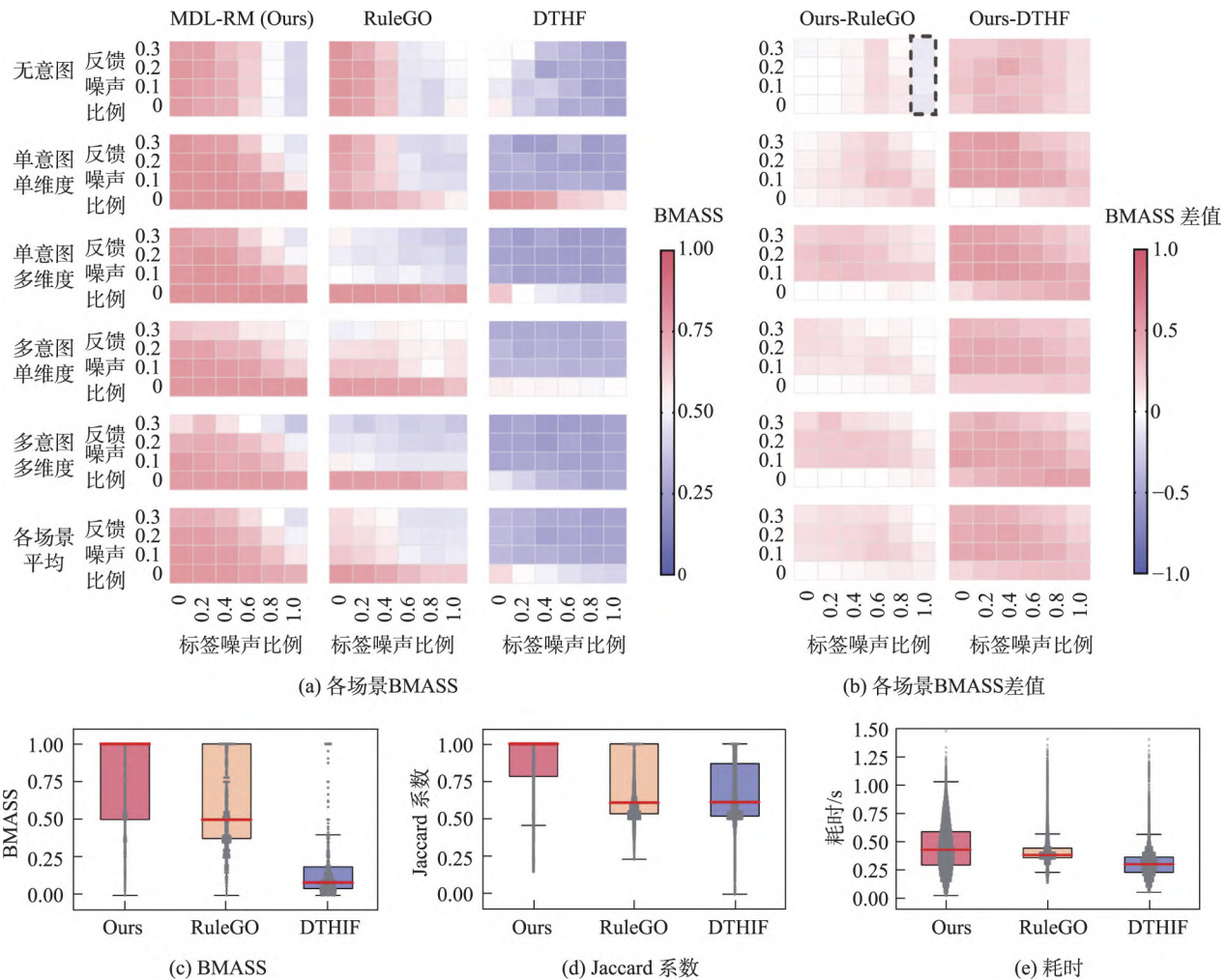


图9 本文算法MDL-RM与RuleGO、DTHF的BMASS、BMASS差值、Jaccard系数及耗时对比
Fig. 9 BMASS, BMASS difference, Jaccard index and computing time of MDL-RM (ours), RuleGO and DTHF

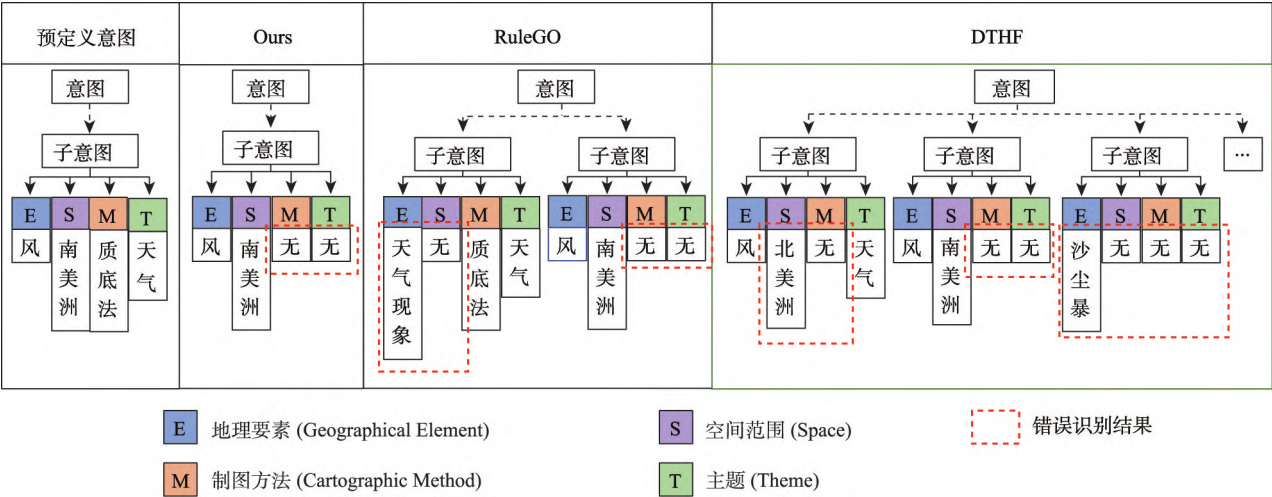


图10 预定义意图与错误识别结果示例
Fig. 10 Examples of predefined and recognized intentions

由于RuleGO为保证子意图的多样性,保留了部分覆盖正样本比例较小的错误子意图,而DTHF为使意图覆盖所有正样本,未剔除错误子意图(图10)。由图9(d)可知3种算法的Jaccard系数波动程度相差不大,但本文算法的中位数更高,说明本文算法的意图识别结果与预定义意图在样本覆盖方面更为一致。由于随机合并策略使得求解过程具有随机性,本文算法的平均求解耗时比RuleGO、DTHF更长且稳定性较差,但多数场景下小于1 s。

4.5 随机合并策略有效性验证

为验证随机合并策略在缩减算法耗时方面的有效性,本文统计使用与不使用随机合并策略时各意图场景下的BMASS、Jaccard系数与耗时,实验结果如图11所示。

从实验结果来看,样本随机合并策略可在不明显降低意图识别准确度的同时显著缩短耗时,使本文算法能够支持实时意图识别。由于用户子意图个数较少且本文各子意图覆盖的样本比例较为均衡,一定数量的随机合并即可使得有效合并出现的概率接近于1,因此在多数场景下是否使用随机合并的BMASS和Jaccard系数无明显差异。然而,不使用随机合并时耗时均值大于3.1 s;随机合并策略仅检查指定个数的合并方案,可显著减少样本合并、目标函数计算等高耗时操作的执行次数,从而将耗时均值降至0.52 s以下。

4.6 样本增强策略有效性验证

本实验通过比较不同规模样本集合在样本增强前后的意图识别准确性,验证样本增强策略应对样本不足问题的能力。为排除样本噪声的影响,对120个有意图无噪声样本集合分别进行随机

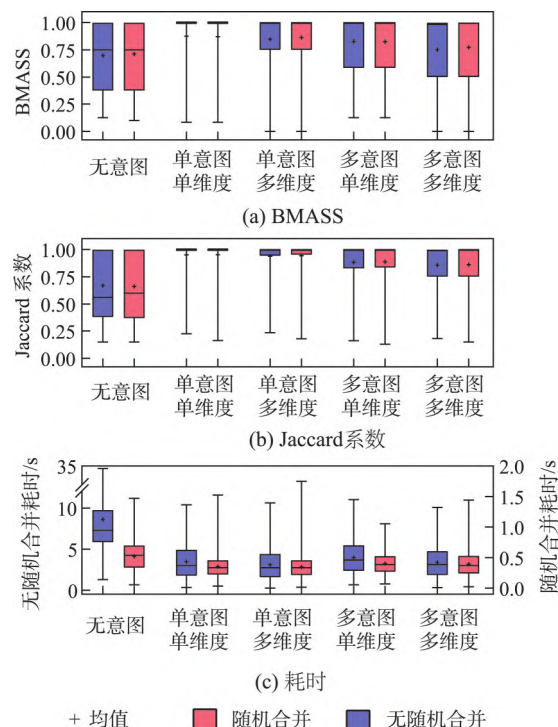


图11 使用随机合并策略前后的BMASS、Jaccard系数与耗时

Fig. 11 Comparison of BMASS, Jaccard index and computing time with and without random merging strategy

采样,生成8个样本数量规模(10、20、40、60、80、100、150、200)共960个样本集合开展实验,各样本集合的正负样本数量相同。样本增强前后的BMASS、Jaccard系数、总编码长度的平均值及标准差如图12所示。

实验结果表明样本增强策略能有效提高样本规模不足情况下意图识别的准确性。当样本规模小于150时,样本增强前的BMASS与Jaccard系数显著低于增强后,而总编码长度大于增强后,且样本数量越少差异越明显。这是因为样本增强前,基

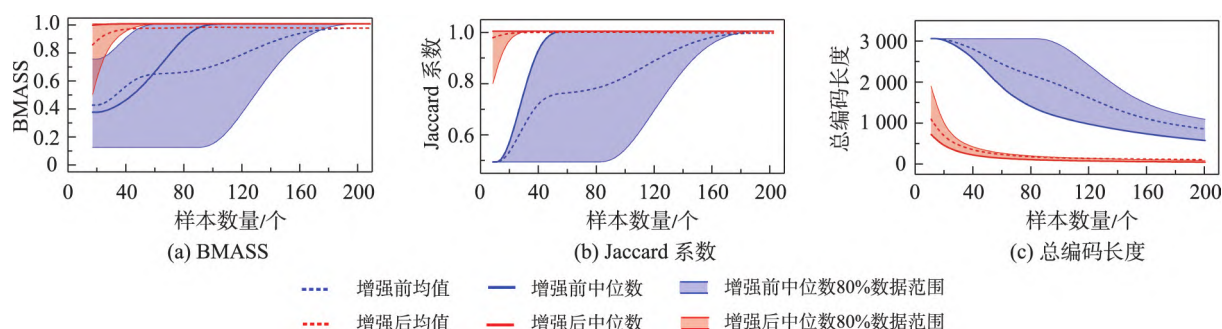


图12 样本增强前后的BMASS、Jaccard系数与总编码长度均值、中位数与分布情况对比

Fig. 12 Averages, medians and ranges of BMASS, Jaccard index and total coding length before and after sample augmentation

于意图编码反馈样本节省的长度不足以抵消编码意图所需长度,导致无法识别出意图;而样本增强后,在样本规模为20时依然获得较高的准确度(Jaccard系数为0.99,BMASS为0.95)。因此,本文方法可通过少量标注样本实现意图识别。

4.7 参数影响分析

本文算法包含“随机合并数量 t_{merge} ”与“正样本覆盖比例阈值 β ”2个参数。为探究上述参数对意

图识别结果及效率的影响,实验统计不同参数值下BMASS、Jaccard系数与耗时的平均值,结果如图13—图14所示。为兼顾BMASS及耗时因素,使得算法性能整体最优,本文探究 t_{merge} 时,将 β 固定为0.3;探究 β 时,将 t_{merge} 固定为50。

由图13可知,对于所有样本整体而言,随机合并数量小于50次时BMASS与Jaccard系数变化幅度较大;随着合并数量增加,BMASS先增加后降

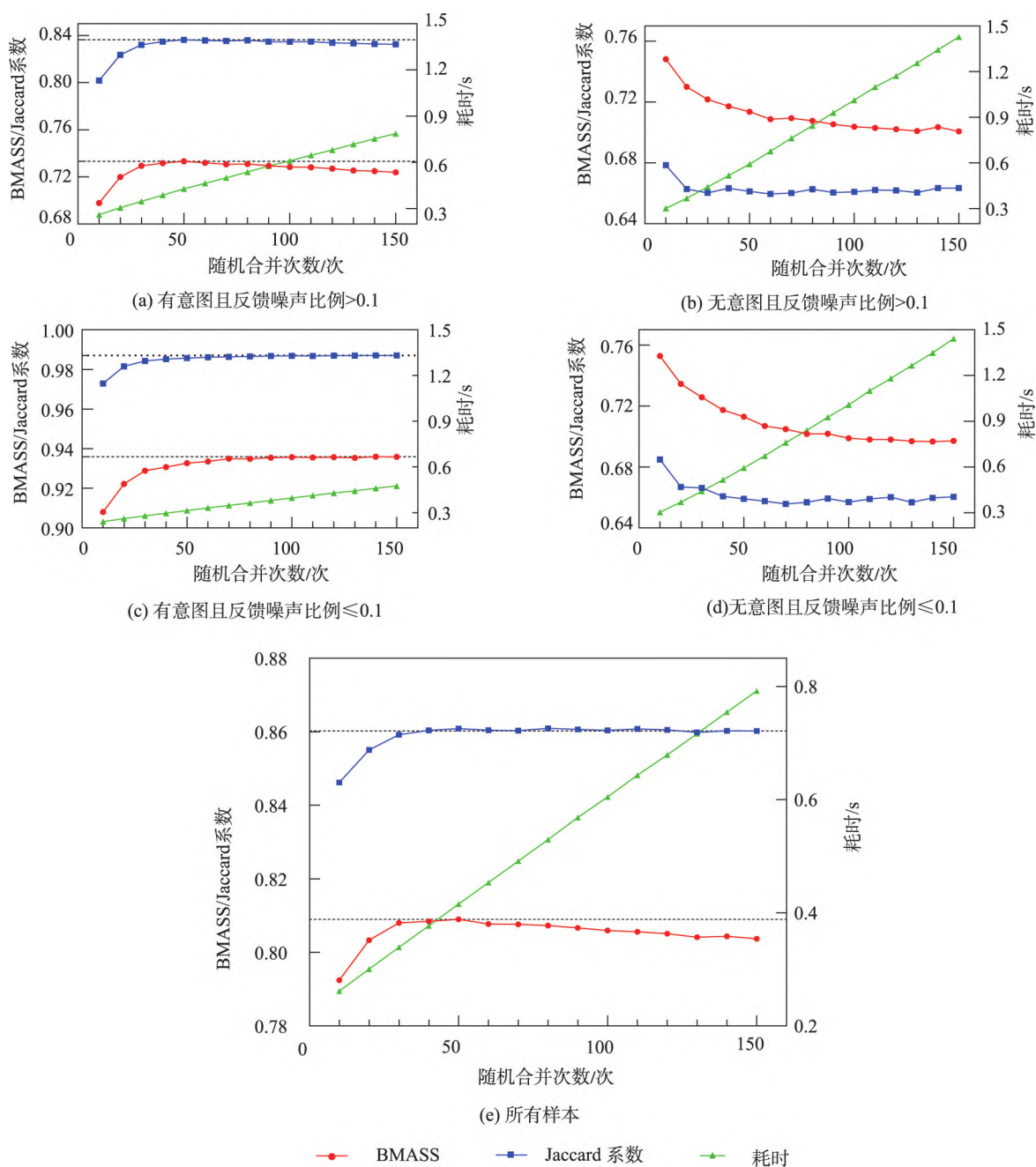


图13 随机合并数量对BMASS、Jaccard系数与耗时的影响

Fig. 13 Influence of numbers of random merge on BMASS, Jaccard index and computing time

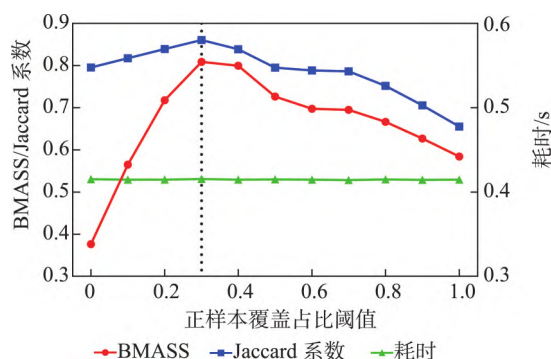


图 14 正样本占比阈值对 BMASS、Jaccard 系数与耗时的影响

Fig. 14 Influence of threshold of positive sample coverage rate on BMASS, Jaccard index and computing time

低, Jaccard 系数先增加后稳定在 0.86 左右, 而耗时始终呈线性增长。本文进一步按照是否有意意图及反馈噪声比例将样本集分为 4 组(图 13(a)—图 13(d)), 探究 BMASS 降低的原因。由于难以保证“无意意图”场景的实验样本合并后不产生任何有效子意图(即某一维度取值非本体根节点), 因此随着随机合并数量的增加, 可能产生错误的意图提取; 同时, 过多的反馈噪声将掩盖真实意图, 使得识别出的意图与预定义意图之间存在偏差, 导致 BMASS 降低。尽管耗时呈线性增长, 但组间存在一定差异。无意意图场景耗时整体高于有意意图场景; 且有意意图场景下, 反馈噪声多则耗时高。这是因为无意意图场景较有意意图场景更易产生低样本占比的错误子意图, 而反馈噪声比例增大会增加出现无效合并的概率, 导致迭代次数与总随机合并次数增加。

图 14 表明, 随着正样本覆盖比例阈值增加, BMASS 与 Jaccard 系数均呈先增加后减少趋势, 并在占比阈值为 0.3 时取得最大值; 求解耗时对该参数不敏感。当比例阈值小于 0.3 时, MDL-RM 未能有效过滤因样本噪声产生的错误子意图, 识别结果与预定义意图存在较大差异; 比例阈值大于 0.3 时, 由于可能存在一定的反馈噪声或多个子意图, 导致某些真实子意图因正样本覆盖占比偏低而被过滤, 意图识别不全。求解耗时稳定在 0.41 s 左右是由于子意图过滤操作对算法整体耗时的影响可忽略。综上, 在用户意图数量及噪声比例未知的情况下, 随机合并数量取值在 40~80 内且正样本覆盖比例阈值为 0.3 时, 能较好地兼顾算法时间效率与准确性。

5 结论与展望

本文提出一种顾及地理语义的地图检索意图形式化表达与识别方法, 旨在克服传统基于元数据文本匹配和图像内容检索的“意图鸿沟”问题, 辅助提高地图资源发现效率。实验结果表明, 与 RuleGO、DTHF 规则归纳算法相比, 本文算法能够准确识别多种场景下的地图检索意图, 并有效应对样本噪声及样本不足问题, 且求解时间可满足实时识别需求。本文主要贡献包括:

(1) 在构建地图检索意图维度树基础上, 设计“意图-子意图-维度分量”嵌套的意图表达模型显式描述子意图、维度分量及其逻辑关系, 可望提升检索过程的可解释性。

(2) 以相关反馈样本为用户需求载体, 将意图识别转化为样本编码长度最小化问题, 通过样本随机合并与贪心算法实现最优意图搜索。

(3) 引入地理本体为子意图构建、意图识别准确度评价指标定义提供语义支撑。

目前本文识别算法仅使用相关反馈中的正样本, 还可考虑利用负样本识别伪用户意图, 并拓展意图形式化表达模型以更全面地描述用户检索需求。此外, 该算法处理的相关反馈样本仅有正负两种标签, 可基于用户鼠标轨迹等隐式反馈将二元“硬标签”扩展到包含用户意图强度信息的“软标签”, 从而提升意图识别准确度。同时, 用户检索意图不仅与当前检索任务有关, 还受长期检索兴趣影响, 可结合用户画像优化算法。

参考文献(References):

- [1] 王家耀, 成毅. 论地图学的属性和地图的价值[J]. 测绘学报, 2015, 44(3): 237-241. [Wang J Y, Cheng Y. Discussions on the attributes of cartography and the value of map[J]. Acta Geodaetica et Cartographica Sinica, 2015, 44(3): 237-241.] DOI:10.11947/j.AGCS.2015.20140406
- [2] 闻国年, 袁林旺, 俞肇元. 地理学视角下测绘地理信息再透视[J]. 测绘学报, 2017, 46(10): 1549-1556. [Lü G N, Yuan L W, Yu Z Y. Surveying and mapping geographical information from the perspective of geography[J]. Acta Geodaetica et Cartographica Sinica, 2017, 46(10): 1549-1556.] DOI:10.11947/j.AGCS.2017.20170338
- [3] 成晓强, 杨敏, 桂志鹏, 等. 信息量与相似度约束下的网络地图服务缩略图自动生成算法[J]. 测绘学报, 2017, 46(11): 1891-1898. [Cheng X Q, Yang M, Gui Z P, et al. An algorithm creating thumbnail for web map services based

- on information entropy and trans-scale similarity[J]. *Acta Geodaetica et Cartographica Sinica*, 2017,46(11):1891-1898.] DOI:10.11947/j.AGCS.2017.20170080
- [4] Gui Z, Cao J, Liu X, et al. Global-scale resource survey and performance monitoring of public OGC web map services[J]. *ISPRS International Journal of Geo-Information*, 2016,5(6):88. DOI:10.3390/ijgi5060088
- [5] 孟婵媛,李勤超,李宏伟,等.基于本体的地理信息查询检索方法研究[J].*测绘科学*,2008,33(S1):251-253,188. [Meng C Y, Li Q C, Li H W, et al. Research on query method of geographic information based on ontology[J]. *Science of Surveying and Mapping*, 2008,33(S1):251-253.] DOI:10.3771/j.issn.1009-2307.2008.07.103
- [6] 苗立志,胥婕,周亚,等.应用描述词汇约简的OGC地理信息服务演绎推理[J].*测绘学报*,2015,44(9):1029-1035,1062. [Miao L Z, Xu J, Zhou Y, et al. OGC geographic information service deductive reasoning based on description vocabularies reduction[J]. *Acta Geodaetica et Cartographica Sinica*, 2015,44(9):1029-1035.] DOI:10.11947/j.AGCS.2015.20140021
- [7] 王东旭,诸云强,潘鹏,等.地理数据空间本体构建及其在数据检索中的应用[J].*地球信息科学学报*,2016,18(4):443-452. [Wang D X, Zhu Y Q, Pan P, et al. Construction of geodata spatial ontology and its application in data retrieval[J]. *Journal of Geo-information Science*, 2016,18(4):443-452.] DOI:10.3724/SP.J.1047.2016.00443
- [8] 邬群勇,郑孝苗,康凌骏.语义地理信息服务的三级匹配发现算法[J].*厦门大学学报(自然科学版)*,2012,51(2):195-199. [Wu Q Y, Zheng X M, Kang L J. The three-level matching algorithm for semantic geospatial information service discovery[J]. *Journal of Xiamen University (Natural Science)*, 2012,51(2):195-199.]
- [9] Gui Z, Yang C, Xia J, et al. A performance, semantic and service quality-enhanced distributed search engine for improving geospatial resource discovery[J]. *International Journal of Geographical Information Science*, 2013,27(6):1109-1132. DOI:10.1080/13658816.2012.739692
- [10] 高勇,姜丹,刘磊,等.一种地理信息检索的定性模型[J].*北京大学学报(自然科学版)*,2016,52(2):265-273. [Gao Y, Jiang D, Liu L, et al. A qualitative method for geographic information retrieval[J]. *Acta Scientiarum Naturalium Universitatis Pekinensis*, 2016,52(2):265-273.] DOI:10.13209/j.0479-8023.2015.113
- [11] 张敏,桂志鹏,成晓强,等.一种WMS领域主题文本提取及元数据扩展方法[J].*武汉大学学报·信息科学版*,2019,44(11):1730-1738. [Zhang M, Gui Z P, Cheng X Q, et al. A text-based WMS domain themes extraction and metadata extension method[J]. *Geomatics and Information Science of Wuhan University*, 2019,44(11):1730-1738.] DOI:10.13203/j.whugis20180083
- [12] Wei Z, Gui Z, Zhang M, et al. Text GCN-SW-KNN: A novel collaborative training multi-label classification method for WMS application themes by considering geographic semantics[J]. *Big Earth Data*, 2021,5(1):66-89. DOI:10.1080/20964471.2021.1877434
- [13] Yang Z, Gui Z, Wu H, et al. A latent feature-based multimodality fusion method for theme classification on web map service[J]. *IEEE Access*, 2019,8:25299-25309. DOI:10.1109/access.2019.2954851
- [14] Hu Y, Gui Z, Wang J, et al. Enriching the metadata of map images: a deep learning approach with GIS-based data augmentation[J]. *International Journal of Geographical Information Science*, 2022,36(4):799-821. DOI:10.1080/13658816.2021.1968407
- [15] 李牧闲,桂志鹏,成晓强,等.多核学习与用户反馈结合的WMS图层检索方法[J].*测绘学报*,2019,48(10):1320-1330. [Li M X, Gui Z P, Cheng X Q, et al. A content-based WMS layer retrieval method combining multiple kernel learning and user feedback[J]. *Acta Geodaetica et Cartographica Sinica*, 2019,48(10):1320-1330.] DOI:10.11947/j.AGCS.2019.20180410
- [16] Hu K, Gui Z, Cheng X, et al. Content-based discovery for web map service using support vector machine and user relevance feedback[J]. *PLoS One*, 2016,11(11):e0166098. DOI:10.1371/journal.pone.0166098
- [17] 葛芸,江顺亮,叶发茂,等.基于ImageNet预训练卷积神经网络的遥感图像检索[J].*武汉大学学报·信息科学版*, 2018,43(1):67-73. [Ge Y, Jiang S L, Ye F M, et al. Remote sensing image retrieval using pre-trained convolutional neural networks based on ImageNet[J]. *Geomatics and Information Science of Wuhan University*, 2018,43(1):67-73.] DOI:10.13203/j.whugis20150498
- [18] Gui Z, Yang C, Xia J, et al. A visualization-enhanced graphical user interface for geospatial resource discovery[J]. *Annals of GIS*, 2013,19(2):109-121. DOI:10.1080/19475683.2013.782467
- [19] 解虹.数字化环境下交互式信息检索[M].北京:科学出版社,2010. [Xie H. Interactive information retrieval in digital environments[M]. Beijing: Science Press, 2010.]
- [20] 陈翀,刘晓兵,徐谷子,等.一种搜索引擎的查询意图发现的新方法[J].*情报学报*,2012,31(3):242-249. [Chen C, Liu X B, Xu G Z, et al. A new method of detecting query intent for search engines[J]. *Journal of the China Society for Scientific and Technical Information*, 2012,31(3):242-249.] DOI:10.3772/j.issn.1000-0135.2012.03.002
- [21] Hanjalic A, Kofler C, Larson M. Intent and its discontents: the user at the wheel of the online video search engine[C]. *Proceedings of the 20th ACM international conference on Multimedia*. New York:ACM, 2012:1239-1248. DOI:10.1145/2393347.2396424
- [22] Zhao J, Chen H, Yin D. A dynamic product-aware learning model for E-commerce query intent understanding

- [C]. Proceedings of the 28th ACM International Conference on Information and Knowledge Management. New York: ACM, 2019:1843-1852. DOI:10.1145/3357384.3358055
- [23] 严锐,李石君.基于查询意图识别与主题建模的文档检索算法[J].计算机工程,2018,44(3):189-194. [Yan R, Li S J. Document retrieval algorithm based on query intent identification and topic modeling[J]. Computer Engineering, 2018, 44(3):189-194.] DOI:10.3969/j.issn.1000-3428.2018.03.032
- [24] 张瑞芳,郭克华.面向个性化站点的用户检索意图建模方法[J].计算机工程与应用,2018,54(6):37-43. [Zhang R F, Guo K H. Novel retrieval intention modeling method for personalized website[J]. Computer Engineering and Applications, 2018,54(6):37-43.] DOI:10.3778/j.issn.1002-8331.1611-0108
- [25] Umemoto K, Yamamoto T, Nakamura S, et al. Search intent estimation from user's eye movements for supporting information seeking[C]. Proceedings of the International Working Conference on Advanced Visual Interfaces. New York: ACM, 2012:349-356. DOI:10.1145/2254556.2254624
- [26] Zhang H, Zha Z J, Yang Y, et al. Attribute-augmented semantic hierarchy: towards bridging semantic gap and intention gap in image retrieval[C]. Proceedings of the 21st ACM international conference on Multimedia. New York: ACM, 2013:33-42. DOI:10.1145/2502081.2502093
- [27] Fariha A, Meliou A. Example-driven query intent discovery: abductive reasoning using semantic similarity[J]. Proceedings of the VLDB Endowment, 2019,12(11):1262-1275. DOI:10.14778/3342263.3342266
- [28] Zhang R, Guo J, Fan Y, et al. Query understanding via intent description generation[C]. Proceedings of the 29th ACM International Conference on Information & Knowledge Management. New York: ACM, 2020:1823-1832. DOI:10.1145/3340531.3411999
- [29] Kuźma M, Bauer H. Map metadata: the basis of the retrieval system of digital collections[J]. ISPRS International Journal of Geo- Information, 2020,9(7):444. DOI: 10.3390/ijgi9070444
- [30] 肖珑,苏品红,姚伯岳.国家图书馆图元数据规范和著录规则[M].北京:国家图书馆出版社,2014. [Xiao L, Su P H, Yao B Y. National library metadata specifications and rules for public maps [M]. Beijing: National Library of China Publishing House, 2014.]
- [31] Content Standard for Digital Geospatial Metadata(FGDC-STD-001-1998)[S].
- [32] 黄仁涛.专题地图编制[M].武汉:武汉大学出版社,2003. [Huang R T. Thematic map compilation[M]. Wuhan: Wuhan University Press, 2003.]
- [33] 查正军,郑晓菊.多媒体信息检索中的查询与反馈技术[J].计算机研究与发展,2017,54(6):1267-1280. [Zha Z J, Zheng X J. Query and feedback technologies in multimedia information retrieval[J]. Journal of Computer Research and Development, 2017,54(6):1267-1280.] DOI: 10.7544/issn1000-1239.2017.20170004
- [34] 曾安,李晓兵,杨海东,等.基于最小描述长度和K2的贝叶斯网络结构学习算法[J].东北师大学报(自然科学版), 2014,46(3):53-58. [Zeng A, Li X B, Yang H D, et al. Bayesian network structure learning algorithm based on MDL and K2[J]. Journal of Northeast Normal University (Natural Science Edition), 2014,46(3):53-58.] DOI:10.11672/dbsdzk2014-03-011
- [35] Proença H M, van Leeuwen M. Interpretable multiclass classification by MDL-based rule lists[J]. Information Sciences, 2020,512(C):1372-1393. DOI:10.1016/j.ins.2019.10.050
- [36] Gr nwald P. A tutorial introduction to the minimum description length principle[M]. Cambridge, MA: MIT Press, 2007.
- [37] 李星洁.形式学习理论与相关归纳逻辑问题研究[D].昆明:云南师范大学,2020. [Li X J. Formal learning theory and related inductive logic problems[D]. Kunming: Yunnan Normal University, 2020.]
- [38] Aoga J O R, Guns T, Nijssen S, et al. Finding probabilistic rule lists using the minimum description length principle[C]. Discovery Science. Cham: Springer International Publishing, 2018:66-82. DOI:10.1007/978-3-030-01771-2_5
- [39] Rissanen J. A universal prior for integers and estimation by minimum description length[J]. The Annals of Statistics, 1983,11(2):416-431. DOI:10.1214/aos/1176346150
- [40] Fano R M. The transmission of information[M]. Massachusetts Institute of Technology, Research Laboratory of Electronics Cambridge, Mass, USA. 1949.
- [41] Lin D. An information-theoretic definition of similarity [C]. ICML. Morgan Kaufmann Publishers Inc, 1998:296-304. <https://dl.acm.org/doi/10.5555/645527.657297>
- [42] Yuan Q, Yu Z, Wang K. A new model of information content for measuring the semantic similarity between concepts[C]. International Conference on Cloud Computing and Big Data. IEEE, 2013:141-146. DOI:10.1109/cloud-com-asia.2013.25
- [43] GeoNames. GeoNames[EB/OL]. <http://geonames.org/>, 2005.
- [44] Raskin R G, Pan M J. Knowledge representation in the semantic web for Earth and environmental terminology (SWEET) [J]. Computers & Geosciences, 2005,31(9): 1119-1125. DOI:10.1016/j.cageo.2004.12.004
- [45] Gruca A, Sikora M. Data-and expert-driven rule induction and filtering framework for functional interpretation and description of gene sets[J]. Journal of Biomedical Semantics, 2017,8(1):1-14. DOI:10.1186/s13326-017-0129-x
- [46] Kinnunen N. Decision tree learning with hierarchical features[D]. Tampere: Tampere University of Technology, 2018.