

基于图嵌入的软件项目源代码检索方法*

凌春阳^{1,2}, 邹艳珍^{1,2,3}, 林泽琦^{1,2}, 谢冰^{1,2}, 赵俊峰^{1,2,3}



¹(高可信软件技术教育部重点实验室(北京大学),北京 100871)

²(北京大学 信息科学技术学院,北京 100871)

³(北京大学(天津滨海)新一代信息技术研究院,天津 300450)

通讯作者: 谢冰, E-mail: xiebing@pku.edu.cn

摘要: 源代码检索是软件工程领域的一项重要研究问题,其主要任务是检索和复用软件项目 API(application program interface,应用程序接口).随着软件项目的规模越来越大、越来越复杂,当前,源代码检索一方面需要提高基于自然语言 API 查询的准确性,另一方面需要定位和展示目标 API 及其相关代码之间的关联,以更好地辅助用户理解 API 的实现逻辑和使用场景.为此,提出一种基于图嵌入的软件项目源代码检索方法.该方法能够基于软件项目源代码自动构建其代码结构图,并通过图嵌入对源代码进行信息表示.在此基础上,用户可以输入自然语言问题、检索并返回相关的 API 及其关联信息构成的连通代码子图,从而提高 API 检索和复用的效率.在以开源项目 Apache Lucene 和 POI 为例的检索实验中,该方法检索结果的 *F1* 值比现有基于最短路径的方法提高了 10%,同时显著缩短了平均响应时间.

关键词: API 检索;代码检索;代码图;图嵌入

中图法分类号: TP311

中文引用格式: 凌春阳,邹艳珍,林泽琦,谢冰,赵俊峰.基于图嵌入的软件项目源代码检索方法.软件学报,2019,30(5):1481-1497. <http://www.jos.org.cn/1000-9825/5721.htm>

英文引用格式: Ling CY, Zou YZ, Lin ZQ, Xie B, Zhao JF. Approach to searching software source code with graph embedding. Ruan Jian Xue Bao/Journal of Software, 2019,30(5):1481-1497 (in Chinese). <http://www.jos.org.cn/1000-9825/5721.htm>

Approach to Searching Software Source Code with Graph Embedding

LING Chun-Yang^{1,2}, ZOU Yan-Zhen^{1,2,3}, LIN Ze-Qi^{1,2}, XIE Bing^{1,2}, ZHAO Jun-Feng^{1,2,3}

¹(Key Laboratory of High Confidence Software Technologies (Peking University), Ministry of Education, Beijing 100871, China)

²(School of Electronics Engineering and Computer Science, Peking University, Beijing 100871, China)

³(Peking University Information Technology Institute (Tianjin Binhai), Tianjin 300450, China)

Abstract: Searching software source code and locating software's API (application program interface) are important research issues in software engineering. As software projects are becoming more and more complex, existing search tools mainly face the following two challenges. First, more accurate search results are required in natural language question based search process. Second, the relationships between API are required to illustrate so that these API's underlying logic and usage scenarios are able to be understood more quickly. In this study, an novel approach is proposed to searching a software project's API based on graph embedding. It aims to improve the accuracy of natural language based code graph search. A software project's code graph is built automatically from its source code and they are represented through graph embedding. For a natural language question, a code-connected subgraph, composed by relevant API and their

* 基金项目: 国家重点研发计划(2016YFB1000801); 国家杰出青年科学基金(61525201)

Foundation item: National Key Research and Development Program of China (2016YFB1000801); National Science Foundation for Distinguished Young Scholars (61525201)

本文由智能化软件新技术专刊特约编辑申富饶教授和李戈副教授推荐.

收稿时间: 2018-08-31; 修改时间: 2018-10-31; 采用时间: 2018-12-13

associated relationships, are returned as the best answer. In experiments, Apache Lucene and POI projects are selected as examples to perform some API search tasks. Experimental results show that the proposed approach improves *F1*-score by 10% than existing shortest path based approach, while reduces average response time significantly.

Key words: API search; code search; code graph; graph embedding

在软件开发过程中,开发者最常用的复用方式就是通过检索、调用软件项目提供的 API(application programming interface,应用程序接口)来实现所需的功能.然而随着软件项目源代码的规模变得越来越大,源代码之间的依赖关系越来越复杂,API 检索也变得越来越大困难^[1].比如,开源项目 Lucene 包含了超过 5 377 个类以及 34 042 种方法.而研究人员在研究搜索引擎查询日志时发现,开发者往往花费了大量的时间与精力来查找 API^[2],通用搜索引擎(如 Google)在检索软件代码方面存在许多问题且效率较低^[3].由此也促生了在很多技术性问答网站(如 StackOverflow)上,可以发现大量关于如何使用 API 的提问.

受自然语言处理技术的影响,目前许多软件源代码检索工具,比如 Krugle、Ohloh 和 Sourcerer^[4]等,都是将源代码视为纯文本进行处理.其基本思路是:利用软件源代码中的自然语言特征,基于查询问题与代码的文本相似度进行匹配,返回相似度较高的检索结果.典型的相似度匹配模型包括布尔模型(Boolean model)^[5]、向量空间模型(vector space model)^[5]、BM25 模型^[5]等.然而研究表明,这些工具的准确率不能让人满意.在 Lü 等人^[6]于 2015 年进行的一项研究中表明,Ohloh 返回的排名前十的结果中,只有 25.7%~38.4%是有用的.研究者们指出,导致现有软件源代码检索工具效果不理想的一个主要原因是其缺乏对自然语言查询问题语义的理解.对此,研究者们从深化问题理解的角度提出了许多新的方法,如扩展语义相近的词(semantic similar word)^[7]、引入查询重构策略(query reformulation)^[8]或利用其他在线文档(如 MSDN)来优化问题理解和检索过程^[6]等.虽然这些技术很大程度上提高了 API 定位的准确度,但其需要积累和利用大量的辅助信息,在上述信息缺乏的情况下,往往不能达到预期效果.本文则希望仅从软件项目的源代码出发,帮助开发者提高 API 检索的效率,因为源代码是复用过程中能够直接获得的资源,其他信息有时则难以获得.

此外,现有这些检索工具返回的结果往往只有相关的代码片段信息,缺乏对关联 API 的说明与展示,不便于用户理解目标 API 调用背后的逻辑关联和整体概念,在一定程度上影响了软件复用的效率.为此,研究者们指出,不应单纯地将软件代码视为纯文本.比如在 McMillan 等人和 Chan 等人^[9,10]的工作中,将源代码以有向图的形式组织起来,图中的结点表示代码元素,边表示代码元素之间的关联关系.图结构是对源代码知识的一种自然且有效的表示,被广泛应用于软件文档检索、特征定位等问题.借助代码的结构图,对源代码的检索就转化为图上的搜索过程,这样得到的结果既能包含需要被定位的代码元素,又能充分体现代码之间的关联关系.举例来说,假设在使用开源软件项目 ApacheLucene(<http://Lucene.apache.org/>)的过程中,开发者当前使用 StringField 类对一个字符串的字段进行处理.当需求发生变化,还要对含有浮点数的字段进行处理时,这个开发者就需要使用 FloatPoint 类并考虑如何重构代码.而我们可以从图 1 所示的源代码图中发现,这两个类都继承自 Field 类,实现了 IndexableField 接口,同时具有 tokenStream 方法.因此对开发者来说,一种更好的策略是使用该抽象的父 Field 作为变量类型,并在处理不同查询时赋值为不同的子类,从而能够较好地满足其上述需求.

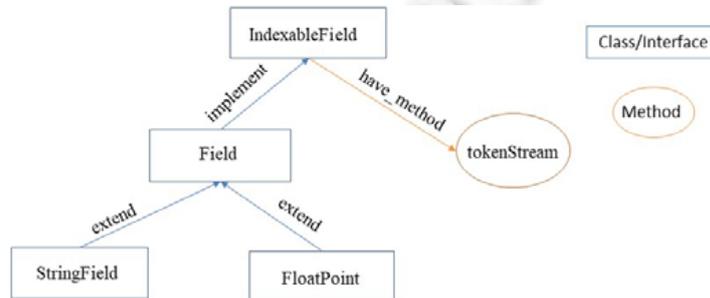


Fig.1 Code graph of Lucene project (Part)

图 1 Lucene 项目的代码图示例(部分)

目前,在基于图结构的 API 检索方面,主要采用了最短路径^[10]、PageRank^[9]等浅层语义分析技术.考虑到在图上进行搜索是一个计算量很大的问题,且要求延迟越短越好,算法的效率至关重要.如果使用在线查询最短路径的方式,则由于搜索过程中经常需要计算图上两个结点之间的距离,将会导致较长的延迟时间.Chan 等人的工作^[10]也提到了这个问题,并利用了一些近似算法进行优化.

针对这些问题,本文提出一种基于图嵌入的软件项目源代码检索方法.图嵌入是一种表示学习技术,能够将图上的结点映射为低维空间中的实值向量,通过向量之间的关系来表达图上的结构信息.对于本文的问题而言,图嵌入技术一方面能够有效表达软件代码图中的深层结构信息,与最短路径等方法相比,更充分地考虑了结点的上下文和关联关系;另一方面,图嵌入过程可以离线进行,在用户在线查询阶段可以立即参与计算,时间复杂度仅是常数级别,大大提高了算法效率.

对比现有工作,本文主要贡献包括:

- (1) 提出了一种将软件源代码进行图嵌入表示的方法,该方法能够基于软件项目源代码,自动构建其代码图结构,并通过图嵌入对源代码进行信息表示.
- (2) 提出了一种基于图嵌入的软件项目源代码检索方法,该方法能够基于项目源代码的图嵌入技术表示,将自然语言查询问题语义匹配到代码子图,提高了检索的准确性,同时展示了 API 的关联关系;
- (3) 实现了一个基于图嵌入的软件项目 API 检索工具原型,并通过 2 个具体软件项目中关于源代码检索的实际问题为例,验证了本文方法的有效性.与 Chan 等人提出的方法相比,本文方法在召回率、F1 值上有了显著提升,并有效减低了检索响应时间.

本文第 1 节介绍图嵌入技术.第 2 节详细描述本文提出的基于图嵌入的软件项目源代码检索方法.第 3 节通过实验验证本文方法的有效性.第 4 节介绍相关工作.第 5 节对本文工作进行总结并展望未来研究工作.

1 图嵌入

图嵌入(graph embedding)的主要目的是将一张图上的结点映射到低维空间的向量,通过这些向量来体现原本图上的结构信息.这里的结构信息可以是一阶、二阶甚至更高阶的结构,一阶结构信息即保持原本相邻的结点在嵌入空间中的距离仍然很近,而更高阶的结构信息则可以保持原本具有相似上下文的结点在嵌入空间中的距离仍然很近.目前,图嵌入技术在知识表示领域颇受关注,被广泛应用于可视化、顶点分类、关联预测、推荐等任务.现有的图嵌入技术主要可以分为以下几类^[11].

- 基于因子分解的图嵌入方法.该类方法将图上结点之间的关系表示成矩阵,对该矩阵进行因子分解从而得到嵌入的向量.矩阵的类型包括邻接矩阵、拉普拉斯矩阵(Laplacian matrix)、卡兹相似矩阵(Katz similarity matrix)等,根据不同的矩阵类型有不同的分解方法,其中,Belkin 等人^[12]提出的 Laplacian Eigenmaps 就是对拉普拉斯矩阵做特征值分解得到的;Ahmed 等人^[13]提出的 GF 算法则是基于对邻接矩阵做因子分解.这种方法的时间复杂度是结点数的平方量级,可扩展性较差,对于规模庞大的代码图而言并不适用.
- 基于随机游走的图嵌入方法.该类方法算法能够近似表示图的许多属性,比如结点的中心度和相似度.Perozzi 等人提出的 DeepWalk^[14]就是该方法的典型代表,该算法利用了神经语言模型 SkipGram 的思想,根据周围的邻居结点来预测当前结点的嵌入向量.该类方法通常从某一结点为中心出发进行多次随机游走,最大化观测到的前 k 个结点以及后 k 个结点的概率,能够保持高阶的相似性.当只需要观测图的部分结构或者图太大而无法全部测量时,该类方法是一个不错的选择.但是该类方法通常只能捕捉到一条路径上的局部结构信息,并且难以找到最佳的采样方法,调优困难.
- 基于深度学习的图嵌入方法.该类方法有的利用深度自编码器(deep autoencoder)来实现降维,因为自编码器拥有学习数据中非线性结构的能力.典型地,Wang 等人^[15]提出的 SDNE 包括无监督和有监督两部分:前一部分由自编码器学习能够重新构造邻居的结点嵌入,后一部分则基于 Laplacian Eigenmaps 对那些本来相邻而映射到嵌入空间后却相距很远的结点进行惩罚.还有一些算法利用卷积神经网络来

得到嵌入向量,如 Kipf 等人^[16]提出的一种半监督学习的方法,这些算法的计算开销更小,更适用于相对稀疏的图.这一类方法虽然能更好地捕捉非线性的结构,然而复杂度还是较高,需要大量的计算开销.

本文采用的是 Tang 等人^[17]提出的 LINE(large-scale information network embedding)方法,该方法的可扩展性很好,能够快速地对大规模网络进行嵌入,并且能够同时保留局部和全局的结构信息.该方法严格意义上不属于上述的 3 种类型,而是在模型中显式地定义了一阶和二阶的近似函数以及任意 2 个结点之间的联合概率分布,通过最小化嵌入矩阵与邻接矩阵所代表的两个概率分布之间的 KL(Kullback-Leibler)散度得到图嵌入的向量.定义好目标函数之后,LINE 并没有采用传统的随机梯度下降算法进行优化,而是提出了一种新的边采样的方法(edge-sampling),根据边的权值成比例的概率进行采样,这样能够防止随机梯度下降中由于边的权值过大导致的梯度爆炸问题.为了进一步形象地阐明一阶近似和二阶近似的意义,我们以图 2 为例进行说明.一阶近似能够将直接有边相连的结点映射到更近的距离,比如图 2 中的结点 6 和结点 7,它们之间有一条权重很大的边连接.而二阶近似则能够将具有相同邻居的结点映射到更近的距离,比如图 2 中的结点 5 和结点 6,它们都具有很多相同的邻居结点.因此,一阶近似能够保留更多局部的结构信息,而二阶近似能够保留更多全局的结构信息.同时考虑 2 种近似的结果,就是将结点 5~结点 7 都映射到嵌入空间中更近的距离.

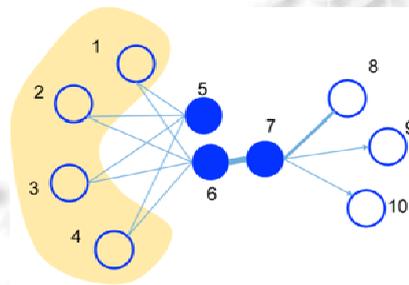


Fig.2 Graph embedding maps nodes with similar structure to closer distance in the embedding space^[17]

图 2 图嵌入将结构相似的结点映射到嵌入空间中更近的距离^[17]

2 本文工作

本文提出了一种基于图嵌入的软件项目源代码检索方法,其基本思想是:利用图嵌入技术有效表达软件代码图中的深层结构信息,在检索过程中充分地考虑了结点的上下文和关联关系.如图 3 所示.

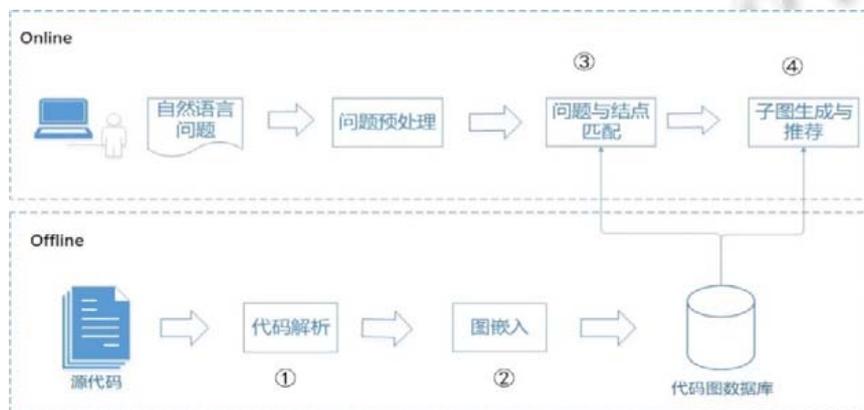


Fig.3 Framework of searching software project's source code based on graph embedding

图 3 基于图嵌入的软件项目源代码检索方法框架

该方法主要包括 4 个部分.

- 1) 代码图的构建:获取项目的源代码并做相应解析,自动构建其代码图.

- 2) 代码图的嵌入:使用图嵌入技术将代码图中的结点表示为向量,以备后续阶段使用.
- 3) 问题与代码图结点的匹配:对用户输入的自然语言问题进行预处理,然后匹配到代码图中结点作为候选结点,并对候选结点进行度量.
- 4) 代码子图的生成与推荐:从第 2 阶段得到的候选结点中挑选合适的结点构成子图,并扩展为连通的子图,将排名最佳的结果返回给用户.

2.1 代码图的构建

在一个软件项目(对于面向对象的函数库而言)的源代码中,主要有两种组成成分:一种是类或接口,另一种是方法.而类与方法之间的关联关系的类型大体上所有面向对象语言都共有,当然,部分关联关系与具体编程语言有关.本文以 Java 语言作为实验对象,考虑了 Java 语言中的如下 6 种关联关系,分别为继承(Inh)、实现(Impl)、成员(Mem)、参数(Par)、返回值(Ret)和调用(Call),用集合 $Rel = \{Inh, Impl, Mem, Par, Ret\}$ 表示(见表 1).在此基础上,我们可以利用现有工具对一个开源软件项目的源码进行解析,构建代码图并存储于图数据库中.

Table 1 Relationship types in code graph

表 1 代码图中的关系类型

关系名称	关联的结点类型	描述
继承(Inh)	类/接口 A-类/接口 B	表示类/接口 A 是类/接口 B 的父类/接口
实现(Impl)	类 A-接口 B	表示类 A 实现了接口 B
成员(Mem)	方法 A-类/接口 B	表示方法 A 是类/接口 B 的成员
参数(Par)	方法 A-类/接口 B	表示类/接口 B 是方法 A 的参数
返回值(Ret)	方法 A-类/接口 B	表示类/接口 B 是方法 A 的返回值
调用(Call)	方法 A-方法 B	表示方法 A 调用了方法 B

代码图(code graph). 代码图 $G=(V_G, E_G)$ 是一个有向无权图,顶点集 V_G 是所有的类和方法的结点,边集 E_G 中任一条边 $e(u, v), \exists R \in Rel$, 使得 uRv . 其中, uRv 表示 u, v 之间存在 R 类型的关联关系.

现在主流的 Java 源码解析工具有很多,例如 CheckStyle, JAVAssist, Yasca, JDT 等.这些工具都现成可用,可以解析 Java 项目并提取出代码元素以及它们之间的关系.本文使用的是 Eclipse 的 JDT,作为一个轻量级的代码解析工具,它可以迅速将 Java 代码构建成一个 DOM 结构的抽象语法树(AST).代码中的每个元素都对应 AST 上的一个结点,通过遍历 AST 就可以得到所有代码元素和它们之间的关系,用来构建软件代码结构图.

本文使用图数据库 Neo4j(<https://neo4j.com/>)来存储代码图.Neo4j 是一个用 Java 实现,完全兼容 ACID 的图数据库,其本身数据组织的结构和我们需求的代码结构图非常相似,也是通过结点与结点之间的关系构成整张图.除此之外,Neo4j 还提供了诸如索引(index)、遍历(traversal)等机制以及一种类 SQL 的查询语言 Cypher,以方便在图上进行更高级的操作.

2.2 代码图的嵌入

本文使用图嵌入的原因是在第 3 阶段,即代码子图的生成与推荐时,需要计算图上任意两个顶点之间的距离.图嵌入技术可以使得这一步骤变得十分方便与快速.如果采用在线查询两个顶点的最短路径的方式则将非常耗时,而图嵌入过程却是离线的,在构建完代码图后即可完成.

本文使用了 Tang 等人^[17]提出的名为 LINE 的方法,该方法的可扩展性很好,时间复杂度较低,能够快速地对大规模网络进行嵌入,并且同时保留局部和全局的结构信息,能够适用于有向图、无向图、带权图等多种类型.LINE 的源代码可以在 Github 上找到,本文用 Java 语言重新实现了该算法.

结点之间的距离. 代码图嵌入以后,代码图中的任一顶点 $v \in V_G$, 都映射到一个向量 $r_v \in \mathbb{R}^d$, 其中, d 表示该向量的维度;任意 2 个顶点 u, v 之间的距离定义为相应的向量之间的欧氏距离,即 $dist(u, v) = \|r_u - r_v\|_2$.

2.3 问题与结点的匹配

从这一阶段开始,便进入在线查询的过程.当用户输入一个自然语言的问题后,我们先对问题进行预处理,然后匹配到代码图中一些候选结点.

2.3.1 问题与结点的预处理

本文采用简单的词袋模型(bag of words),将问题视为一些无序的词组成的集合,忽略其句式和位置等信息.由于自然语言文本还存在单词之外的一些符号,如标点符号等,本文将去除非数字和非字母的其他符号,并以此作为分隔符进行切词.切词以后,进行停用词处理.自然语言包含许多无实际意义的功能词,比如冠词和介词等.这些词存在十分普遍而包含信息的却极少,人们一般称之为停用词(stop words).本文去除的停用词包括常见的英文停用词、Java 保留字以及项目领域特有的功能词(比如 Lucene),等.

代码图中结点(类或方法)的预处理是将该类或方法的名字通过切词得到一个词袋,用来描述该结点的语义信息.本文使用驼峰切词法进行处理,比如 QueryParser 这个类名,切词后的结果就是{query,parser}.

问题(query). 问题 Q 是由一些去除停用词之后的单词的集合组成,将该单词集合记为 $BOW_Q=\{q_1,\dots,q_n\}$.

结点单词集合. 代码图 G 中的任一结点 v ,都有一个与之关联的单词集合,记为 $BOW_v=\{s_1,\dots,s_m\}$.

2.3.2 产生候选结点

候选结点集合. 问题 Q 所匹配到的候选结点集合是一个集族,记为 $Candidate_Q=\{C_1,\dots,C_n\}$,其中, $C_i(1\leq i\leq n)$ 是单词 q_i 匹配到的代码图中结点的集合,其大小记为 $|C_i|$.

对于 Query 词袋中的每个词,我们都为之在图上找到一系列的候选结点,这样做是为了充分保留原问题中的语义信息.接下来需要明确对于每个单词 q ,如何匹配到它所对应的候选结点.本文使用文本匹配的方法,主要考虑了如下 5 种匹配情形.

- (1) 全名匹配.单词 q 本身就是就是一个类或方法的全名,比如 $q=QueryParser$.
- (2) 部分匹配.单词 q 被结点 v 所对应的词袋所包含,即 $q\in BOW_v$,比如 $q=query$.
- (3) 词根化匹配.单词 q 与结点 v 所对应的词袋中的词经词根化后相同.词根化(stemming)能够将语义上相同而由于语法等原因导致词形不同的词抽取出同样的词根,本文采用 Snowball(<http://snowball.tartarus.org/>)作为词根化工具.比如 $q=parse$,由于 parse 和 parser 词根化后都是 pars,因此 q 会匹配到 QueryParser 这个结点.
- (4) 缩略规则匹配.单词 q 经过手工构造的一些规则转换后,被结点 v 所对应的词袋包含.这是为了解决自然语言的词汇与代码元素中的命名不一致的问题,代码元素的命名为了方便,很可能只采取了自然语言单词的缩写,比如 Document 在代码中的命名可能是 Doc,Number 在代码中的命名可能是 Num,等.
- (5) 同义词匹配.单词 q 与结点 v 所对应的词袋中的词,经过相似度计算被判定为同义词.同义词是自然语言处理领域的一个重要研究问题,在自然语言与代码元素匹配的过程中也存在这个问题,比如问题中的词是 delete,但代码中某个方法的名字可能是 remove,这种情况也应该被考虑.为了解决这个问题,本文使用了在 Wikipedia 上预训练的 glove(<https://nlp.stanford.edu/projects/glove/>)词向量.通过计算两个词向量的余弦相似度来表示语义的相近性,当大于特定阈值时,就认为这两个词是同义词.

$$Similarity(u,v) = \frac{\mathbf{g}_u \cdot \mathbf{g}_v}{\|\mathbf{g}_u\| \cdot \|\mathbf{g}_v\|} \quad (1)$$

2.3.3 候选结点的度量

上一步骤为了保证对原问题中词汇的覆盖程度,产生了大量的候选结点.但这些结点不应该都是同等地位的,有的候选结点与问题语义更贴近,有的却徒增了许多噪音,因此需要对结点与问题的文本相似程度进行度量.本文考虑两种评价指标:一是该结点词袋与问题中相关的词越多越好;二是该结点引入的不相关的词越少越好,即当结点与问题相关的词数量相同时, BOW_v 越小的结点质量越好.可类比于信息检索中的准确率与召回率,具体定义如下.

$$score_{relevant}(v) = \frac{|BOW_v \cap BOW_Q|}{|BOW_Q|} \quad (2)$$

$$score_{irrelevant}(v) = \frac{|BOW_v \cap BOW_Q|}{|BOW_v|} \quad (3)$$

为了考虑同义词的影响,在计算两个词袋的交集时,需要进行特殊处理.如果结点词袋中的某个单词不在 Query 词袋中出现过,那么利用词向量计算该单词与 Query 中所有单词的余弦相似度,取其中的最大值作为该单词对交集的贡献.因此,该交集的值不一定是整数.

结点权重. 以 $score_{relevant}$ 和 $score_{irrelevant}$ 的 $F1$ 值作为该结点与问题的相关度的度量,称为该结点的权重,记为 $w(v)$,其计算公式如下所示.

$$w(v) = \frac{2 \times score_{relevant}(v) \times score_{irrelevant}(v)}{score_{relevant}(v) + score_{irrelevant}(v)} \quad (4)$$

2.4 子图的生成与推荐

得到候选结点集合之后,需要从中挑选合适的结点来构成能够反映问题语义的子图.从大量的候选结点中构成子图的搜索空间非常巨大,如何高效地生成子图并度量其优质程度是问题的重中之重.最后,将生成的子图扩展成为连通的子图返回给用户.因此,本阶段可分为两个部分:子图的生成与度量,子图的扩展与推荐.

2.4.1 子图的生成与度量

本文从两个方面来评价子图的质量:一是子图中的结点与问题之间的文本相似程度尽可能地高,每个结点的文本相似程度即结点的权重 $w(v)$;二是子图中结点之间的距离越近越好,这将起到消除歧义的作用.比如子图中的一个结点是 *Document* 类,现在有两种同名的方法 *add*,它们的文本相似程度是相同的,那么距离 *Document* 类更近的那个方法更可能是相关的,而距离远的结点可能是一个属于其他类的方法,与本问题无关.

首先,我们定义如何从候选结点的集合中选择子图的顶点集.

顶点集. $Candidate_Q = \{C_1, \dots, C_n\}$ 构成的子图 G' 的顶点集 $V' = \{C_{i,j} | x_{i,j} = 1\}$, 其中, $x_{i,j}$ 是二元变量,表示 C_i 中的第 j 个结点是否被选中,满足约束条件 $\sum_{j=1}^{|C_i|} x_{i,j} \geq 1, x_{i,j} = x_{i',j'} \text{ iff } C_{i,j} = C_{i',j'}$.

这里的第 1 个约束条件是为了使问题中的每个词所对应的候选结点集合中至少有 1 个结点被选中,以此保证对原问题语义的覆盖度.另外,同一个结点可能在不同的集合中重复出现,因此只要该结点在某个集合中被选中,则其他集合中也必须选中它,这就由第 2 个约束条件保证.

接下来,我们定义一个子图 G' 的度量函数:

$$score(G') = \sum_{u,v \in V', u} \frac{dist(u,v)}{w(u) \times w(v)} \quad (5)$$

其中的距离 $dist(u,v)$ 即第 2.2 节中所述的使用图嵌入向量计算的欧氏距离.子图的生成过程即求解这个优化问题 $minscore(G')$, 这可以看成是一个二次优化问题.为了方便求解,本文提出一种基于柱状搜索(*beam search*)的方法.柱状搜索可以看作对贪心法的一种优化,每次将保留前 k 个最佳的候选结果.具体如算法 1 所示.

算法 1. 柱状搜索子图.

输入:候选结点集合 $Candidate_Q = \{C_1, \dots, C_n\}$.

输出:排名最高的一个子图.

1. $C \leftarrow \bigcup_{i=1}^n C_i$
2. Sort C by weight,按权值从高到低排序
3. $Beam \leftarrow Top(k, C)$,取前 k 个权值最高的
4. for $i \leftarrow 2$ to n do
5. for each $V' \in Beam$ do
6. for each $c \in C_i$ do
7. $delta \leftarrow \sum_{v \in V'} \frac{dist(c,v)}{w(c) * w(v)}$
8. $newV' \leftarrow V' \cup \{c\}$

9. $cost \leftarrow cost + delta$, 增加该结点后子图的总的代价
10. $Add(newV', newBeam)$, 将结果加入到一个新的 $beam$ 中
11. **end**
12. **end**
13. $sort\ NewBeam\ by\ cost$, 按代价从低到高排序
14. $Beam \leftarrow Top(k, NewBeam)$, 取前 k 个代价最小的
15. **end**
16. **return** $Top(1, Beam)$

对于输入的候选结点集合,按权值排序后取最高的 k 个结点作为起始结点,其中,参数 k 就代表着柱状的大小($beam\ size$).第7行~第10行是先根据公式(5)中的度量函数计算加入当前结点后增加的代价,然后生成一个新的子图加入候选集合中,并计算累积的代价.第13行、第14行仍是只保留列表中排序靠前的 k 个结果.最后,我们返回排名最高的一个子图.

2.4.2 子图的扩展与推荐

利用图嵌入的向量来计算距离,带来了计算的便利和速度的提升,但由于在选择结点的过程中我们并没有考虑实际存在的边,因此到目前为止,我们只得到了子图的顶点集.为了给用户提供更形象地理解,我们应该将这些结点是如何连接起来一同展示给用户,即从顶点集扩展成一个连通的子图.

本文利用将这个问题定义成给顶点集 V' 构造一棵最小生成树($minimum\ spanning\ tree$, 简称 MST),将 V' 中任意两个顶点之间的最小跳数定义为它们之间的边的权重,这样做就意味着用尽可能少的边将所有顶点连接起来.具体算法如下表的算法2所示,其中,第5行的 $FindShortestPath$ 函数即寻找结点集合 X 和 Y 之间最短路径,而每两个结点之间的最短路径可以通过 $Cypher$ 语句在代码图中进行查询;第8行、第9行是将这条路径上所经过的结点和边都加入到扩展后的子图中.直到所有结点都被包含进来,该算法便得到了一个连通的子图.

算法2. 扩展为连通子图.

输入:顶点集 V' .

输出:扩展后的子图 G^e .

1. $V^e \leftarrow V', E^e \leftarrow \emptyset$
2. $X \leftarrow randomly\ select\ a\ node\ from\ V'$
3. $Y \leftarrow V' - X$
4. **while** $Y \neq \emptyset$ **do**
5. $v, path \leftarrow FindShortestPath(X, Y)$
6. $X \leftarrow X \cup \{v\}$
7. $Y \leftarrow Y - \{v\}$
8. $V^e \leftarrow V^e \cup path.V$
9. $E^e \leftarrow E^e \cup path.E$
10. **end**
11. $G^e \leftarrow (V^e, E^e)$
12. **return** G^e

3 实验与实例

基于上述方法,本文设计并实现了基于图嵌入的软件 API 检索工具原型.下面我们通过一些实验来验证本文所提出方法和工具的有效性.这些实验主要用于回答下列研究问题.

- 研究问题 1: 本文所提出的方法是否能够有效地定位用户所需要的 API?

本文方法将源代码组织成代码结构图的形式而非仅视为纯文本来处理,其中利用了图嵌入技术,定义了子

图的度量函数,并使用基于柱状搜索的算法进行代码子图的检索.我们希望本文所提出的方法在回答开发者提出的实际自然语言问题时,能够有效提升软件项目 API 检索的效果.为了验证这一点,我们选择了两个著名的开源项目以及与它们相关的自然语言问题进行了验证,并与传统基于文本匹配的方法进行了对比.

- 研究问题 2:本文中所提出的图嵌入方法与现有其他检索方法相比效果如何?

在本文所提出的方法中,我们使用了图嵌入方法来挖掘代码图的结构信息,在度量结点之间的距离时,使用图嵌入向量之间的距离,并由此定义了子图的度量函数.我们关心图嵌入方法是否能够有效地表达代码结构图中蕴含的语义信息,从而改进代码检索的效果.因此,我们设计了相应的实验,将本文所使用的图嵌入方法与现有的基于最短路径的方法进行了对比,以验证图嵌入方法的有效性.

3.1 实验设计

我们选取了两个著名的开源软件项目——Apache Lucene 和 Apache POI 作为实验对象,自动构造了这两个项目的代码图,并分别对其进行了基于自然语言的代码检索实验.实验设计的细节列举如下.

3.1.1 代码图构造

对于每个软件项目,我们基于它的源代码中构建一个代码结构图.Lucene 和 POI 项目的源代码都可以通过其官网的连接或 Github 直接下载.一个软件项目往往会有很多个不同的版本,我们选取了这两个项目其中被广泛使用版本进行代码结构图的提取.与代码结构图相关的统计信息在表 2 中进行了展示,包括源代码版本、类或接口的结点数量、方法的结点数量、关联关系的数量以及构造时间.我们在一台 3.40GHz 双核处理器、8GB 内存的服务器上,解析生成两个项目的代码结构图的时间开销分别为 39min 和 31min.

Table 2 Information of code graphs for experimental projects

表 2 实验所用项目的代码图信息

项目名称	源代码版本	类/接口结点数	方法结点数	总结点数	关联关系数	构造时间(min)
Lucene	6.3.0	5 377	34 042	39 419	255 880	39
POI	3.14	3 678	32 001	35 679	211 692	31

在 Chan 等人的工作中,他们实验所使用的源代码数据是 JSE,它是一个非常通用的底层函数库.本文所关注的问题是在软件项目复用时如何有效进行 API 检索,我们实验中使用的 Lucene 和 POI 是领域特定的函数库/软件项目.基于功能实现的需要,其 API 之间的关联往往非常紧密,使得这类领域特定的软件项目在实际中的进行代码检索和复用的难度更大.从表 2 中可以看出, Lucene 项目代码图中,结点(类和方法等)数量是 39 000 左右,但它们之间的关联关系数量达到了 255 000.同时,在 POI 项目的代码图中结点数和关联关系的数量关系也与此类似.

3.1.2 API 检索问题

进行 API 检索的自然语言提问多种多样,为了保证问题的真实有效性,本文使用了 2 种客观的方式来获得实验所需的问题.

- 对于 POI 项目,我们通过其官方网站上提供的用户指南(<https://poi.apache.org/spreadsheet/quick-guide.html>),抽取了 20 个问题作为第 1 组 Query.每一个问题在用户指南中都有对应的示例代码,我们将这些示例代码中涉及的 API 标注为该问题的 groundtruth.
- 对于 Lucene 项目,由于其官方网站上没有相关教程和示例代码,我们使用 StackOverflow 上关于 Lucene 项目的问答帖子作为来源,采取随机抽样的方法,直到挑选出 20 个与源代码检索相关的问题作为第 2 组 Query.具体抽取方法是:首先,从 StackOverflow 上 2008 年~2016 年的数据中找出带有 Lucene 标签的问题,并按照问题的投票为正,问题有被接受的答案,答案中含有代码片段的条件进行筛选,得到了 923 个问题,再从这些候选问题中进行无放回的随机抽取,人工阅读是否属于源代码检索相关的问题,直到挑选出 20 个相关问题;其次,由于这些问题在 StackOverflow 上分为标题和内容两部分,且标题基本上都可以作为其内容概括,本文就以问题的标题作为代码检索工具的查询输入,即 Query.针对这 20 个 Query,我们还需要确定其对应的目标 API.为此,我们组织了 3 名熟悉 Lucene 项目的研究生,分别阅

读了这些 Query 在 StackOverflow 上的问题和答案,结合帖子的答案和源代码的相关知识,在代码图上进行标注.具体的标注过程是:开发人员以该问题答案中出现的代码片段为基准,并结合源代码的相关文档,利用代码图数据库的查询结果,最终确定该 Query 对应的 API 集合.

本文实验所使用的问题是开发者在实际开发过程中遇到的、真实存在的、用自然语言表述的问题,其中部分问题及其标注见表 3,问题之中可能含有代码元素的名称,比如 lengthNorm,也可能没有.因此我们认为,对这些问题的解答,可以很好地反映各种方法的实际效果.

Table 3 Examples of query and annotated API

表 3 问题与标注 API 示例

编号	问题	标注的 API	所属项目
1	How to set document boost attribute?	Document, Attribute, BoostAttribute, AttributeImpl, BoostAttributeImpl. setBoost	Lucene
2	How to merge several taxonomy indexes?	TaxonomyMergeUtils, IndexWriter, TaxonomyMergeUtils. merge	
3	How is Lucene query boost affected by lengthNorm similarity?	Similarity, TFIDFSimilarity, BoostQuery, TFIDFSimilarity. lengthNorm	
4	Sorting on a numeric field in Lucene.	Field, SortField, SortedNumericSortField	
5	How to create a query using wildcards?	Query, MultiTermQuery, AutomatonQuery, WildcardQuery	
6	How to escape some characters utilizing QueryParser?	QueryParser, QueryParserBase, QueryParserBase. escape	
7	Get string representation of Terms using HighFreqTerms.	HighFreqTerms, Terms, Terms. iterator, HighFreqTerms. getHighFreqTerms, BytesRef. toString	
8	How to get the score of top document in Lucene?	ScoreDoc, TopDocs, Document, TopDocs. getMaxScore	
9	How does regexp query work on Lucene?	RegexpQuery, automaton. RegExp, AutomatonProvider	
10	How to implement a fuzzy search query in Lucene?	FuzzyQuery, Query, MultiTermQuery	
31	How to create workbook and sheet in POI?	Workbook, Workbook. createSheet	POI
32	How to create and set cell style in a workbook?	Workbook, Workbook. createCellStyle, CellStyle, Cell. setCellStyle, CellStyle. setDataFormat	
33	How to set bottom border of a cell?	Cell, Cell. setCellStyle, CellStyle, CellStyle. setBottomBorder	
34	Get document summary information in a workbook.	Workbook, HSSFWorkbook, Workbook. getDocumentSummaryInformation	
35	How to set header center in a sheet?	Sheet, XSSFSheet, Header, Sheet. getHeader, Header. setCenter	
36	Set page number right on the footer.	Footer, HSSFFooter, HeaderFooter, HeaderFooter. page, HeaderFooter. numPages, Footer. setRight	
37	Set italic font in a workbook.	Workbook, Workbook. setFont, Font, Font. setItalic	
38	How to set hyperlink in a cell?	Cell, Cell. setHyperlink, Hyperlink	
39	Get all pictures data from a workbook.	Workbook, Workbook. getAllPictures, PictureData, PictureData. getData	
40	How to set repeating row in a sheet?	Sheet, Row, Sheet. setRepeatingRows	

3.1.3 比较对象

为了考察本文工作的实际效果,我们首先分析比较了本文方法在检索过程中选择第 1 个结果和选择前 3 个结果的情况;其次,我们将这两种情形与其他两种典型的代码检索方法进行了对比.实验方法包括以下几种.

- (1) 本文方法(Top 1).这一检索方法只采用了本文基于图嵌入的代码检索方法的第 1 个结果.
- (2) 本文方法(Top 3).这一检索方法采用了本文基于图嵌入的代码检索方法的前 3 个结果.
- (3) 基于文本匹配的方法.这一检索方法只考虑问题与结点的文本相似度进行匹配,而不考虑结点之间的距离因素,即第 2.4.1 节中子图的度量函数公式(5)中只含权重部分而没有距离部分.
- (4) 基于最短路径的方法.这一检索方法来源于 Chan 等人的工作,其采用最短路径来计算两个结点之间的距离,而非本文方法中使用的图嵌入向量之间的距离.

3.1.4 评价指标

为了评价检索结果的好坏,本文使用了准确率(precision)、召回率(recall)和 F1-值(F1-score)作为评价指标.准确率和召回率的计算与真阳性(true positive)等概念有关,其定义见表 4.

Table 4 Metrics for search results

表 4 检索结果的度量指标

	正确的	不正确的
预测到的	true positives (tp)	false positive (fp)
未预测到的	false negatives (fn)	true negatives (tn)

在本文的实验中,准确率和召回率的定义如公式(6)、公式(7)所示.

$$P(G) = \frac{tp}{tp + fp} = \frac{|V_G \cap V_H|}{|V_G|} \tag{6}$$

$$R(G) = \frac{tp}{tp + fn} = \frac{|V_G \cap V_H|}{|V_H|} \tag{7}$$

其中, V_H 标注结果的结点集合, V_G 本文工具返回结果的结点集合. 这里的 True positives 即两个集合的交集, F1 值即准确率和召回率的调和平均数.

3.2 检索实例

针对表 3 中的实验问题 1 “how to set document boost attribute in Lucene?”, 本文在 Lucene 软件项目的代码图上进行了检索实验. 首先, 这个问题被预处理后的单词集合为 {set, document, boost, attribute}, 对于集合中每一个单词, 根据第 2.3.2 节中提到的匹配方法, 分别得到了 {1623, 1101, 63, 122} 个对应的候选结点, 可以看出, 候选结点的规模相当庞大. 接着, 我们针对每个候选结点计算与问题相关程度作为该结点的权重, 比如 *setBoost* 结点的权重为 0.75, *Attribute* 结点的权重为 0.4, 在这些候选结点的基础上, 按照第 2.4 节中的算法构造问题检索对应的子图. 这里分为两个部分.

- 1) 子图的生成: 基于离线完成的图嵌入向量的结果可以计算两个结点之间的距离, 比如 *Document* 结点和 *setBoost* 结点之间距离为 13.47, 根据本文定义的目标函数, 通过柱状搜索算法可以得到子图中应该包含的候选结点为 {*Document*, *Attribute*, *setBoost*}.
- 2) 子图的扩展: 采用基于最小生成树的方法将子图中的结点扩展为连通的代码子图, 得到的结果如图 4(a)所示, 最终的连通子图共包含了 8 个结点.

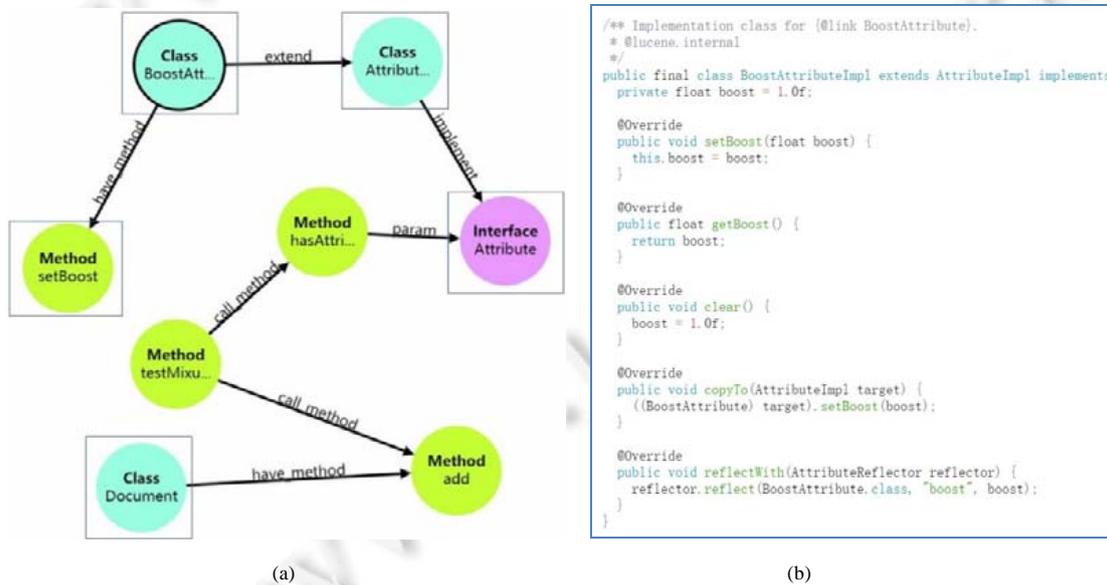


Fig.4 Search result for query “how to set document boost attribute”

图 4 问题“how to set document boost attribute”的检索结果

根据问题的检索结果,我们可以计算对应的准确率、召回率和 $F1$ 值.在图 4 中,用方框标出的 5 个结点是该问题被人工标注的 API,而其余 3 个结点则不属于被标注的.所以该检索结果的准确率为 5/8.而由于方框的 5 个结点刚好覆盖了全部标注的 API,因此该结果的召回率为 1,最终得到的 $F1$ 值为 10/13.

本文的检索结果是一个子图,不仅展示了用户需要的目标 API,还展示了其与相关代码元素的关联.在此基础上,用户可以通过点击图中的 API 结点显示其对应的代码片段.如图 4(b)所示,就是在上述结果中点击结点 *BoostAttributeImpl* 后,能够看到对应的代码体.由此,该检索结果能为用户理解和复用 API 提供有效支持.

3.3 实验结果

针对上述得到的两组各 20 个问题,本文对比了上述 4 种方法检索结果的平均准确率、召回率和 $F1$ 值等评价指标,实验结果见表 5.实验环境是在 Windows 10 系统,2.30GHz 的双核处理器,8GB 内存,Java 版本是 1.8.实验中,图嵌入向量的维度设置为 200,柱状搜索中的柱大小(*beamsize*)设置为 8.

Table 5 Experimental results analysis

表 5 实验结果分析

方法名称	Lucene 项目检索			POI 项目检索			效率分析	
	准确率	召回率	$F1$ 值	准确率	召回率	$F1$ 值	平均响应时间(s)	最长响应时间(s)
本文方法 Top 1	0.53	0.85	0.63	0.70	0.81	0.74	1.39	3.42
本文方法 Top 3	0.55	0.86	0.65	0.72	0.83	0.76	1.39	3.42
基于文本匹配的方法	0.41	0.57	0.43	0.43	0.62	0.49	0.13	0.80
最短路径方法	0.47	0.62	0.52	0.65	0.66	0.64	86.68	257.18

3.3.1 方法效果对比

从表 5 中可以看出,在 Lucene 和 POI 这两个项目中,本文方法的 Top3 都同时获得了最高的平均准确率、召回率和 $F1$ 值.而本文方法 Top 1 与 Top 3 的结果相差很小,说明对于大部分问题本文方法的 Top 1 就是最佳的结果.

基于文本匹配的方法在两个项目上都获得了最低的平均准确率、召回率和 $F1$ 值,表明单纯利用问题与结点的文本相似度而忽视代码结构图中的语义信息,难以取得很好的效果.而本文引入图嵌入技术挖掘代码图结构的方法之后,检索结果得到了显著的提升.

基于最短路径的方法在检索过程中也利用到了代码结构图的信息,只是在计算结点之间的距离时使用的最短路径,而非本文方法中使用的图嵌入向量之间的距离.该方法的效果与纯文本匹配方法相比有了较大提升,但仍然显著低于本文方法的结果.本文方法检索结果的 $F1$ 值比基于最短路径的方法提高了 10%,表明了图嵌入方法的有效性,能够比最短路径方法挖掘到更全面的图结构信息.

此外,我们还对比了不同方法的平均响应时间和最长响应时间.从表 5 中可以看出,纯文本匹配方法的效率最高,平均响应时间和最长响应时间都非常短,因为其计算十分简单,然而其检索效果也是最不令人满意的.而最短路径方法的平均响应时间则非常长,超过了 1min,这对于实际的用户交互而言难以接受.这是因为该方法每次计算距离都需要在线查询两个结点之间的最短路径,使用广度优先搜索的算法最坏情况下也需要 $O(N)$ 的复杂度,这里的 N 表示所有结点的数量.图嵌入过程由于离线已经完成在线计算两个向量的距离时只需要常数时间,所以本文方法在改进检索效果的同时也大大提升了运行的效率.

Chan 等人提出的方法是基于最短路径方法,他们使用的子图度量函数与本文有所差别,他们在论文中报告的实验结果为准确率 0.53、召回率 0.56、 $F1$ 值 0.54,与本文实验中最短路径方法的结果十分接近.我们使用的数据集并不相同,他们的实验中所使用的问题是从 KodeJava 网站上挑选的关于使用 JSE 编程的 20 个示例,并且将问题抽取为多个词组的形式作为输入.而本文实验中使用的 Lucene 和 POI 则是领域特定的函数库,输入也直接是自然语言的问题.

3.3.2 准确率与召回率分析

从实验结果来看,本文方法检索的召回率比较高,但准确率比较低.也就是说,代码子图中除了与 Query 相关

的结点之外,还包含了较多无关的结点.特别是对于 Lucene 项目而言,准确率和召回率的差别比较大,因此我们对 Lucene 项目的 20 个问题做了进一步的分析.

首先,我们对每个问题的准确率进行了具体分析,结果如图 5 所示.其中,斜线底色的柱形表示检索结果中所有的结点(API)数量,纯色柱形表示检索结果中正确的结点数量.可以看到,在我们的检索结果中,通常更多的 API 被返回和推荐出来,大部分子图的规模(结点数)都在 5~8 之间,只有两个问题的子图大小达到了 10.经分析发现,导致这个问题的主要原因是基于最小生成树的子图扩展算法.由于很难断定哪种类型边更加有益,我们将所有边的权值都设为相同,从而产生了大量相同长度的路径,因而在扩展路径上很可能引入没有被标注的结点.虽然我们不能确定这些更多的 API 是必要的,但从召回率较高的事实出发,也就是在保证能够返回用户所需 API 的情况下,我们相信,提供更多的信息对用户理解和使用这些 API 仍然是有益的.而对于子图扩展算法的优化,可以作为本文的进一步工作.

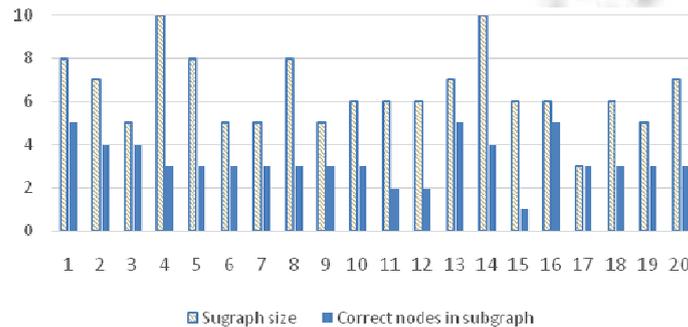


Fig.5 Precision of each question

图 5 每个问题的准确率

我们对每个问题的检索召回率也进行了具体分析,结果如图 6 所示.其中,斜线底色的柱形表示人工标注时的 API 结点数量,纯色柱形表示检索结果中正确的结点数量.从图中可以看出,两种柱形高度十分接近,表明整体的召回率很高.对于其中的 10 个实验示例而言,召回率都达到了 100%,而召回率比较低的只有编号为 13 和 15 的两个问题,柱形高度之差为 2.我们对这些召回率偏低的问题进行了分析,比如编号为 15 的问题“*How to query document have any term in a set?*”,这里影响其实验效果的因素是自然语言的歧义性,问题中的 *set* 是指集合的意思,而源代码中则存在大量将 *set* 作为动词使用的方法名.由于本文在解析问题时只采用了简单的词袋模型,目前尚无法很好地解决这类歧义问题.本文工具最终匹配到的是 *setTerm* 方法,因为该方法含有两个与问题相同的词,结点的权重就很高.在未来工作中,我们将考虑利用代码结构图和图嵌入技术来尽可能地消除自然语言歧义,进一步提升代码检索的效果.

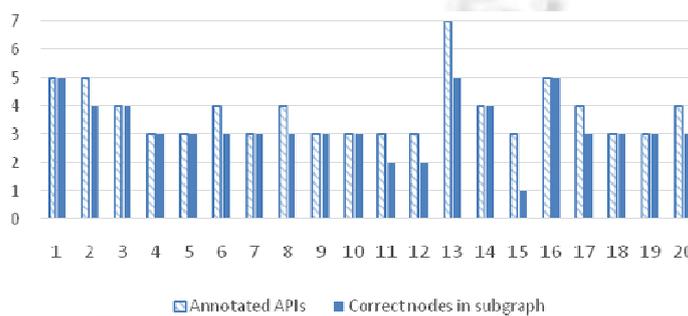


Fig.6 Recall of each question

图 6 每个问题的召回率

3.4 有效性讨论

- 实验评估的合理性.

本文工作的目标是提高软件项目 API 检索与复用的效率,检索结果是一个代码子图.我们关注该子图中的 API 及其关联节点是否能够解决开发者的问题.本文参考 Chan 等人的实验指标,使用准确率、召回率和 F 值作为主要评价,而非信息检索中的 MRR 等排名相关指标.而在进行问题答案的标注过程中,我们从该问题对应的示例代码中提取 API,并通过交叉验证等方式保证标注结果的正确性和客观性.由于实验的问题都对应有被接受的代码片段,因此,如果检索的子图中包含了其中的 API,那么我们认为它的确能够帮助开发者解决问题.在未来工作中,我们也将考虑进行用户研究(UserStudy),以进一步定量探讨给开发者带来的提升效果.

- 影响实验结果的内部因素.

关于检索实例:本文针对 Lucene 和 POI 项目,分别通过两组代码检索实例进行了实验,共 40 个问题.尽管问题数量有限,但我们相信这些问题具有很好的代表性,不会影响实验结论:首先,我们选取的实验问题源于实际,一组问题来源于项目官网的用户指南,另一组问题来源于 StackOverflow 的帖子,这些问题都是开发过程中真实存在的,并且都有相应的代码片段,能够保证问题的较高质量,而对于帖子的筛选和随机抽取,也能够保证其代表性;其次,问题的数量对实验效果没有明显影响,在第 3.3.2 节的实验中可以看到,单个检索实例的评价指标结果与所有实例的结果差异不大.在相关工作中,检索实例的数量基本与本文相似.比如,CodeHow^[6]使用了 34 个问题,Chan^[10]使用了 40 个问题.

- 影响实验结果的外部因素.

关于实验对比方法:由于本文不仅需要定位目标 API,还需要展示其与相关代码元素的关联,为此,检索结果是一个子图,返回代码片段或 API 列表的相关工作不适合与我们直接进行比较.Chan 等人的工作是目前与本文场景最接近的工作,其方法的本质是基于最短路径的.因此,本文实验中以最短路径方法作为比较,能够验证本文方法的有效性.

4 相关工作

本文的相关工作主要包括对源代码 API 检索和图嵌入技术两部分,图嵌入技术已在第 1 节进行了介绍.本节主要介绍源代码检索与 API 推荐的相关技术,重点关注的场景是针对源代码的基于自然语言问题检索.

早期的源代码检索工作主要是将源代码视为纯文本,考虑查询问题与代码之间的文本相似度进行匹配,利用了 TF-IDF, BM25 等信息检索技术. Krugle, Ohloh 和 Sourcerer^[4]就是其中的典型代表. Sourcerer 将解析后的源代码存储在本地的关系型数据库中,并利用 Lucene 进行索引.这些方法的查询方式通常是基于关键词,返回结果是包含查询中关键词或正则表达式的代码片段.然而,由于其缺乏对问题语义的理解,这些工具的实际效果往往不够令人满意.

针对上述问题,近来有许多研究工作结合了自然语言处理技术,从深化问题语义理解的角度提出了一系列的查询精炼(query refinement)技术应用于代码检索,包括查询扩展(query expansion)、查询重构(query reformulation)等.典型地, COCAB^[18]从 StackOverflow 帖子含有的代码片段中提取结构化信息来增强用户的查询,以此解决自然语言与代码元素的词汇不匹配问题. SNIFF^[19]先为代码片段中的 API 寻找到相应文档作为其注解,继而在这些被注解过的代码上进行自然语言的查询. CodeHow^[6]则是通过搜索在线文档识别出一些与问题相关的潜在 API,然后使用扩展布尔模型将这些 API 的信息结合到代码检索的过程中去,从而提高准确度.还有一些方法通过加入语义上相似的词来扩展用户的自然语言查询,这些相似的词可以来源于对网页^[20]或者软件项目^[9,21]的挖掘. DERECS^[22]就预先构造了一个包含大量代码-描述对的语料库,利用它对缺少注释的代码进行描述增强.然而有研究表明:如果在扩展查询时引入了一些不合适的词,将会导致更为糟糕的结果^[23].另一些方法则扩展了用户的查询方式, Gu 等人^[24]的工作中,用户除了可以使用关键词进行查询,还可以加入代码的语法和语义约束条件,从而实现更精确地检索. Wang 等人^[25]提出了一种动态的检索方法,能够结合用户的反馈信息来优化查询. Haiduc 等人^[26]提出了一个名为 Refoqus 的工具,能够预测用户查询的质量,并自动推荐查询重构

的策略.BIKER^[27]从 StackOverflow 上抽取了相似的问题和 API,并利用词嵌入(word embedding)来计算文本描述的相似度.基于机器学习的代码检索和 API 推荐工作也是目前的一个研究热点,比如,ROSF^[28]先通过信息检索的方法生成候选集,再利用有监督学习为候选结果进行重排序.DEEPAPI^[29]使用 RNN Encoder-Decoder 模型,将自然语言问题翻译到 API 调用序列.Function Assistant^[30]则利用语义解析的技术,从代码-文本对中提取特征学习翻译模型.APIRec^[31]从细粒度的代码修改库中训练统计学习的模型,基于上下文上进行 API 推荐.RecRank^[32]则引入路径相关的特征对 APIRec 的结果进行重排序,从而提升了 Top 1 推荐的准确率.虽然这些技术很大程度上提高了 API 定位的准确度,但其需要积累和利用大量的辅助信息,在上述信息缺乏的情况下,往往不能达到预期效果.在本文工作中,我们尤其关注了在没有辅助文档信息和大量查询历史的前提下,如何借助代码自身蕴含的语义信息来辅助提升代码检索的效果.

由于现有代码检索工具往往返回问题相关的代码片段,缺乏对关联 API 的说明与展示,不便于用户理解目标 API 调用背后的逻辑关联,如调用链、类图等,来理解整体的概念^[33,34].为此,一些工作开始关注如何在定位与问题相关 API 的同时,整合、利用和展示代码元素之间的关联信息.典型地,McMillan 和 Chan 等人的工作就将代码组织成有向图的结构,将代码检索问题转化为图上的搜索问题.McMillan 等人提出的工具 Portfolio^[9]利用了 PageRank 和 SAN(spreading activation network)算法,返回图上与问题相关的前 k 个结点作为答案.而 Chan 等人的工作^[10]则更进一步地保证了返回结果是一个连通的子图,从而能够清晰地展现出各个结点之间的连接关系.RACS^[35]针对 JavaScript 框架,构建方法之间的调用关系图,并将自然语言问题也抽取为图形式,寻找与问题的图结构相似的代码片段.许多其他工作也将同样地思想用于文档检索、特征定位等诸多场景^[36,37].受这些工作的启发,本文也将源代码组织为有向图的结构,利用了图嵌入技术来挖掘图结构的深层信息,能够更充分地考虑代码图的上下文信息;同时,由于图嵌入可以离线完成,提高了在线查询算法的效率.

5 总结与展望

为了帮助用户更好地检索和复用软件项目 API,本文提出了一种基于图嵌入的软件项目源代码检索方法.图结构能够很好地反映软件项目源代码中 API 之间的关联,帮助用户理解软件项目;图嵌入技术能够更充分地挖掘图结构的深层信息,并且可以离线完成,有效地提高检索过程中生成子图的效率.基于上述方法,我们设计并实现了相应的软件项目 API 检索工具,并以开源项目 ApacheLucene 和 POI 为例,通过相关实验验证了本文方法的有效性.

在未来的工作中,我们将在图嵌入技术的基础上进一步改进问题与结点的匹配方法,利用词性、句法等自然语言处理技术更好地挖掘问题的语义;同时,在匹配过程中增加对边的类型的分析,在扩展子图时考虑路径上其余结点与问题的相关度,以减少不相关结点的引入,从而提升检索的准确率.此外,我们将进一步探讨在实际过程中用户与搜索结果之间的交互会面临的问题,增强搜索结果的可读性.

References:

- [1] Scaffidi C. Why are API difficult to learn and use? ACM Crossroads, 2006,12(4):Article 4.
- [2] Hoffmann R, Fogarty J, Weld DS. Assieme: Finding and leveraging implicit references in a Web search interface for programmers. In: Proc. of the 20th Annual ACM Symp. on User Interface Software and Technology (UIST 2007). 2007. 13–22.
- [3] Stylos J, Myers BA. Mica: A Web-search tool for finding API components and examples. In: Proc. of the Visual Languages and Human-centric Computing (VLHCC 2006). 2006. 195–202.
- [4] Linstead E, Bajracharya S, Ngo T, Rigor P, Lopes C, Baldi P. Sourcerer: Mining and searching Internet-scale software repositories. Data Mining and Knowledge Discovery, 2009,18(2):300–336.
- [5] Baeza-Yates R, Ribeiro-Neto B. Modern Information Retrieval: The Concepts and Technology behind Search. Addison-Wesley, 2011.
- [6] Lü F, Zhang HY, Lou JG, Wang SW, Zhang DM, Zhao JJ. Codehow: Effective code search based on API understanding and extended Boolean model. In: Proc. of the 30th IEEE/ACM Int'l Conf. on Automated Software Engineering (ASE 2015). 2015. 260–270.

- [7] Hill E, Pollock LL, Vijay-Shanker K. Improving source code search with natural language phrasal representations of method signatures. In: Proc. of the 26th IEEE/ACM Int'l Conf. on Automated Software Engineering (ASE 2011). 2011. 524–527.
- [8] Rahman MM, Roy CK. QUICKAR: Automatic query reformulation for concept location using crowdsourced knowledge. In: Proc. of the 31st IEEE/ACM Int'l Conf. on Automated Software Engineering. 2016. 220–225.
- [9] McMillan C, Grechanik M, Poshyvanyk D, Xie Q, Fu C. Portfolio: Finding relevant functions and their usage. In: Proc. of the ICSE. 2011. 111–120.
- [10] Chan WK, Cheng H, Lo D. Searching connected API subgraph via text phrases. In: Proc. of the SIGSOFT FSE. 2012. Article 10.
- [11] Goyal P, Ferrara E. Graph embedding techniques, applications, and performance: A survey. arXiv:1705.02801 [cs.SI], 2017.
- [12] Belkin M, Niyogi P. Laplacian eigenmaps and spectral techniques for embedding and clustering. In: Proc. of the 14th Int'l Conf. on Neural Information Processing Systems. 2001. 585–591.
- [13] Ahmed A, Shervashidze N, Narayanamurthy S, Josifovski V, Smola AJ. Distributed large-scale natural graph factorization. In: Proc. of the 22nd Int'l Conf. on World Wide Web (WWW 2013). 2013. 37–48.
- [14] Perozzi B, Al-Rfou R, Skiena S. Deepwalk: Online learning of social representations. In: Proc. of the 20th Int'l Conf. on Knowledge Discovery and Data Mining. 2014. 701–710.
- [15] Wang D, Cui P, Zhu W. Structural deep network embedding. In: Proc. of the 22nd Int'l Conf. on Knowledge Discovery and Data Mining. ACM Press, 2016. 1225–1234.
- [16] Kipf TN, Welling M. Semi-supervised classification with graph convolutional networks. arXiv preprint arXiv:1609.02907, 2017.
- [17] Tang J, Qu M, Wang M, Zhang M, Yan J, Mei Q. LINE: Large-scale information network embedding. In: Proc. of the 24th Int'l Conf. on World Wide Web (WWW 2015). 2015. 1067–1077.
- [18] Sirres R, Bissyand'e TF, Kim DS, Lo D, Klein J, Kim K, Le Traon Y. Augmenting and structuring user queries to support efficient free-form code search. Empirical Software Engineering, 2018,23(5):2622–2654.
- [19] Chatterjee S, Juvekar S, Sen K. Sniff: A search engine for Java using free-form queries. In: Proc. of the 12th Int'l Conf. on Fundamental Approaches to Software Engineering (FASE 2009). 2009. 385–400.
- [20] Tian Y, Lo D, Lawall JL. Automated construction of a software specific word similarity database. In: Proc. of the CSMR-WCRE. 2014. 44–53.
- [21] Yang J, Tan L. Inferring semantically related words from software context. In: Proc. of the 9th IEEE Working Conf. of Mining Software Repositories (MSR 2012). 2012. 161–170.
- [22] Li X, Wang QX, Jin Z. Description reinforcement based code search. Ruan Jian Xue Bao/Journal of Software, 2017,28(6): 1405–1417 (in Chinese with English abstract). <http://www.jos.org.cn/1000-9825/5226.htm> [doi: 10.13328/j.cnki.jos.005226]
- [23] Sridhara G, Hill E, Pollock LL, Vijay-Shanker K. Identifying word relations in software: A comparative study of semantic similarity tools. In: Proc. of the 16th IEEE Int'l Conf. on Program Comprehension (ICPC 2008). 2008. 123–132.
- [24] Gu YS, Zeng GS. Accurate search method for source code by combining syntactic and semantic query. Journal of Computer Applications, 2017,37(10):2958–2963 (in Chinese with English abstract).
- [25] Wang S, Lo D, Jiang L. Active code search: Incorporating user feedback to improve code search relevance. In: Proc. of the 29th IEEE/ACM Int'l Conf. on Automated Software Engineering (ASE 2014). 2014.
- [26] Haiduc S, Bavota G, Marcus A, Oliveto R, De Lucia A, Menzies T. Automatic query reformulations for text retrieval in software engineering. In: Proc. of the 2013 Int'l Conf. on Software Engineering (ICSE 2013). 2013. 842–851.
- [27] Huang Q, Xia X, Xing ZC, Lo D, Wang XY. API method recommendation without worrying about the task-API knowledge gap. In: Proc. of the 2018 33rd Int'l Conf. on Automated Software Engineering (ASE 2018). 2018. 292–303.
- [28] Jiang H, Nie L, Sun Z, *et al.* ROSF: Leveraging information retrieval and supervised learning for recommending code snippets. IEEE Trans. on Services Computing, 2019,12(1):34–46.
- [29] Gu XD, Zhang HY, Zhang DM, Kim SH. Deep API learning. In: Proc. of the 24th ACM SIGSOFT Int'l Symp. on the Foundations of Software Engineering. ACM Press, 2016. 631–642.
- [30] Richardson K, Kuhn J. Function assistant: A tool for NL querying of API. In: Proc. of the EMNLP. 2017. 67–72.

- [31] Nguyen AT, Hilton M, Codoban M, Nguyen HA, Mast L, Rademacher E. API code recommendation using statistical learning from fine-grained changes. In: Proc. of the 24th ACM SIGSOFT Int'l Symp. on the Foundations of Software Engineering. ACM Press, 2016. 511–522.
- [32] Liu XY, Huang LG, Ng V. Effective API recommendation without historical software repositories. In: Proc. of the 2018 33rd ACM/IEEE Int'l Conf. on Automated Software Engineering (ASE 2018). 2018. 282–292.
- [33] Sillito J, Murphy GC, De Volder K. Asking and answering questions during a programming change task. IEEE Trans. on Software Engineering, 2008,34(4):434–451.
- [34] Sillito J, Murphy GC, De Volder K. Questions programmers ask during software evolution tasks. In: Proc. of the SIGSOFT FSE. 2006. 23–34.
- [35] Li X, Wang ZR, Wang QX, Yan SM, Xie T, Mei H. Relationship-aware code search for JavaScript frameworks. In: Proc. of the 24th ACM SIGSOFT Symp. on the Foundations of Software Engineering. 2016. 690–701.
- [36] Fu K, Wu YJ, Peng X, Zhao WY. A feature location method based on call chain analysis. Computer Science, 2017,44(4):56–59 (in Chinese with English abstract).
- [37] Li Z, Niu J, Wang K, Xin YY. Optimization of source code search based on multi-feature weight assignment. Journal of Computer Applications, 2018,38(3):812–817 (in Chinese with English abstract).

附中文参考文献:

- [23] 黎宣,王千祥,金芝.基于增强描述的代码搜索方法.软件学报,2017,28(6):1405–1417. <http://www.jos.org.cn/1000-9825/5226.htm> [doi: 10.13328/j.cnki.jos.005226]
- [24] 顾逸圣,曾国荪.基于语法和语义结合的源代码精确搜索方法.计算机应用,2017,37(10):2958–2963.
- [36] 付焜,吴毅坚,彭鑫,赵文耘.一种基于子图搜索的特征定位方法.计算机科学,2017,44(4):56–59.
- [37] 李阵,钮俊,王奎,辛园园.基于多特征权重分配的源代码搜索优化.计算机应用,2018,38(3):812–817.



凌春阳(1995—),男,重庆人,硕士生,CCF 学生会员,主要研究领域为软件工程,软件复用,数据挖掘.



谢冰(1970—),男,博士,教授,博士生导师,CCF 高级会员,主要研究领域为软件工程,形式化方法,软件复用.



邹艳珍(1976—),女,博士,副教授,CCF 专业会员,主要研究领域为软件工程,软件复用,信息检索.



赵俊峰(1974—),女,博士,副教授,CCF 高级会员,主要研究领域为软件工程,软件复用,Web 服务,普适计算,大数据分析技术.



林泽琦(1992—),男,博士生,主要研究领域为软件工程,软件复用,知识工程,数据挖掘.