

面向软件安全性缺陷修复的开发者推荐方法

孙小兵^{1,2)} 周澄¹⁾ 杨辉¹⁾ 李斌¹⁾

¹⁾(扬州大学信息工程学院, 江苏扬州 225127)

²⁾(复旦大学, 上海市数据科学重点实验室 上海 201203)

摘 要 软件开发与维护过程中常会出现一些安全性缺陷, 这些安全性缺陷会给软件 and 用户带来很大的风险。安全性缺陷在修复过程中, 其修复级别和质量要求往往高于一般性的缺陷, 因此, 推荐出富有安全性经验的开发者及时有效地修复这些安全性缺陷非常重要。现有的开发者推荐技术在推荐开发者时仅仅考虑了开发者的历史开发内容, 很少考虑到开发人员的安全性缺陷修复经验和修复质量等因素, 所以这些技术不适用于安全性缺陷的开发者推荐。本文针对安全性缺陷的修复提出了一种有效的软件开发者推荐方法 SecDR。SecDR 在推荐开发者时不仅考虑了开发者的历史开发内容(与安全性相关), 还分析了开发者的修复质量和历史修复缺陷的复杂度等因素。此外, SecDR 还实现了开发者的多经验级别推荐: 推荐初级开发者修复简单的安全性缺陷, 高级开发者修复复杂的安全性缺陷。本文在三个开源项目(Mozilla, Libgdx, ElasticSearch)上分别对 SecDR 推荐开发者进行有效性验证。通过对比实验证明, SecDR 针对安全性缺陷推荐开发者相比于其他方法(如: DR_PSF)的推荐精度平均高出 19%~42%。另外, 实验对比了 SecDR 与实际开发人员的分配情况, 结果显示 SecDR 可以更好地规避不合理的软件开发者的推荐。

关键词 安全性缺陷; 开发者推荐; 缺陷库; 缺陷分配; 软件维护

中图法分类号 *****

DOI 号 *投稿时不提供 DOI 号* 分类号

SecDR: Developer Recommendation for Security Bugs

SUN Xiaobing^{1,2)} ZHOU Cheng¹⁾ YANG Hui¹⁾ LI Bin¹⁾

¹⁾(School of Information Engineering, Yangzhou University, Yangzhou 225127, China)

²⁾(Shanghai Key Laboratory of Data Science, Fudan University, Shanghai 201203, China)

Abstract Security bugs are commonly emerged bugs during the software development and maintenance, which cause security risks during software deployment. Security bugs need to be fixed with high quality and patched faster than other types of bugs. Recommending developers to fix security bugs is one of the important tasks during the security bug fixing process. Some developer recommendation techniques have been proposed to fix the bugs, but most of these techniques did not recommend developers considering their security experience and fixing quality. In this paper, we propose a novel approach, SecDR (Security Developer Recommendation), to recommend developers considering their fixing quality and fixing complexity of their historical fixed security bugs. In addition, SecDR recommends junior developers for simple bugs, and recommends senior developers for complex bugs. An empirical study on three open source subjects (Mozilla, Libgdx and ElasticSearch) are conducted to evaluate the effectiveness of SecDR. In our study, SecDR is also compared with the state-of-art

收稿日期: 年-月-日; 最终修改稿收到日期: 年-月-日 *投稿时不填写此项*. 本课题得到国家自然科学基金 (No. 61402396, 61472344)、中国博士后自然科学基金(No. 2015M571489)、南京大学软件新技术国家重点实验室开放课题(No. KFKT2016B21)资助、江苏省青蓝工程项目.孙小兵(通讯作者), 男, 1985 年生, 博士, 副教授, 硕士生导师, 主要研究领域为软件维护与演化、软件数据分析学.E-mail: xbsun@yzu.edu.cn.周澄, 女, 1985 年生, 博士生在读, 主要研究领域为软件智能化分析.杨辉, 男, 1990 年生, 硕士研究生, 主要研究领域为软件Bug分配.李斌(通讯作者), 男, 1965 年生, 博士, 教授, 博士生导师, 主要研究领域为数据分析、智能计算.第 1 作者手机号码: 18252740912. E-mail: xbsun@yzu.edu.cn

developer recommendation technique, DR_PSF, to evaluate the effectiveness of developer recommendation. Results show that the accuracy of SecDR is improved over DR_PSF, and the gain values range from 19% to 42%. Moreover, we also compare the results of SecDR with actual developer allocation, and results show that SecDR can effectively recommend developers, which is even better than the developer allocation in the real bug assignment environment.

Key words security bug fix; developer recommendation; bug repository; bug assignment; software maintenance

0 引言

随着计算机技术的迅速发展,人们越来越依赖软件和互联网,计算机带来了便利革新也带来许多隐患。计算机软件开发与使用的过程中,经常出现一些安全性缺陷^[1],导致软件系统产生漏洞。黑客通过漏洞恶意攻击他人计算机,如“勒索病毒”等的蔓延,造成了大量的数据泄露,给用户和软企双方带来巨大的利益损失^[2,3]。因此,安全性缺陷受到软件开发者的广泛关注。如何及时有效地修复安全性缺陷已经成为了工业界和学术界共同关注的问题^[1,2,4]。

相关研究结果表明,安全性缺陷的修复级别,修复复杂度、质量和时效性要求往往比一般性缺陷要高^[2]。例如,在 firefox 项目中,安全性缺陷的修复速度比性能性缺陷快 2.8 倍,被 reopen 的可能性却是性能性缺陷的 2.5 倍,其它类型缺陷的 4.5 倍^[2]。在缺陷修复分配时,安全性缺陷的分配难度更大,往往由于不能很好地及时分配给合适的开发者而面临重新分配,据统计安全性缺陷的再分配次数是一般性缺陷的 3.5 倍^[2]。同时安全性缺陷也是开发者修复历史中相对较少的缺陷类型,表 1 中给出了 Mozilla 项目中开发者安全性缺陷修复数量分布情况,修复数量少于 10 个的开发者有 195 位,占总人数的 86%,只修复过 1 个的开发者有 88 位,占总人数的 39%,绝大部分开发者修复安全性缺陷的经验匮乏。

表 1 Mozilla 项目中开发者安全性缺陷修复数量分布

缺陷数量	>=100	50~100	10~50	<10	=1
开发者数量	2	3	27	195	88

目前,很多相关工作都在研究如何推荐合适

的开发者解决软件中新出现的缺陷。例如,shokripour 等人提出了 Time-Text-Based 方法^[7],该方法利用主题词语和时间因素预测系统中的开发人员对缺陷的熟悉程度从而推荐出一个最合适的开发者。这些推荐技术都可称为泛推荐技术,主要思想是将新缺陷的描述信息与历史修复记录相匹配,遵循匹配值最高推荐原则,广泛面向软件缺陷管理库中的所有缺陷进行推荐,缺失了特定类型缺陷推荐的针对性。推荐结果不可避免偏向于历史开发数据相对丰富的开发者,而这些开发者却不一定拥有丰富的安全性知识和较高的安全性缺陷修复技能,泛推荐技术并不完全适用于安全性缺陷^[14,15]。

本文针对软件项目中安全性缺陷的修复,提出一种新的开发者推荐方法——SecDR(Security Developer Recommendation)。SecDR 采用分级推荐机制,不仅考虑了开发者与安全性缺陷相关的历史开发内容,还考虑了开发者的修复质量和历史修复缺陷的复杂度。首先,SecDR 预测安全性缺陷的修复复杂度等级,综合评估开发者的历史修复经验预测开发者的修复等级;再对推荐的开发者进行排序,排序时主要考虑开发者的经验等级,修复质量和与新缺陷相关的安全性知识匹配度等因素,并给出一个开发者推荐列表。SecDR 不仅可以有效地推荐出开发经验丰富的高级开发者解决难度较大的软件安全性缺陷,还可以推荐出开发经验较少的初级开发者解决比较简单的安全性缺陷,从而实现开发者的多经验级别推荐。这种分级机制可以很好地协调软件开发者的任务安排,充分利用了项目团队中的有效人力资源,最大满足安全性缺陷时效性原则。

1 面向安全性缺陷的开发者推荐技术

SecDR 包括六个步骤, 具体流程见图 1。首先, SecDR 提取软件缺陷库中的安全性缺陷, 构成独立的安全性缺陷库。当出现新缺陷(下文特指新的安全性缺陷), 将其描述与安全性缺陷库

历史信息进行匹配, 查找相关历史安全性缺陷, 同时对新缺陷的修复复杂度进行预测; 另一方面, SecDR 对库中所有的开发者进行工作经验分析, 并预测开发者的经验等级; 最后, SecDR 根据匹

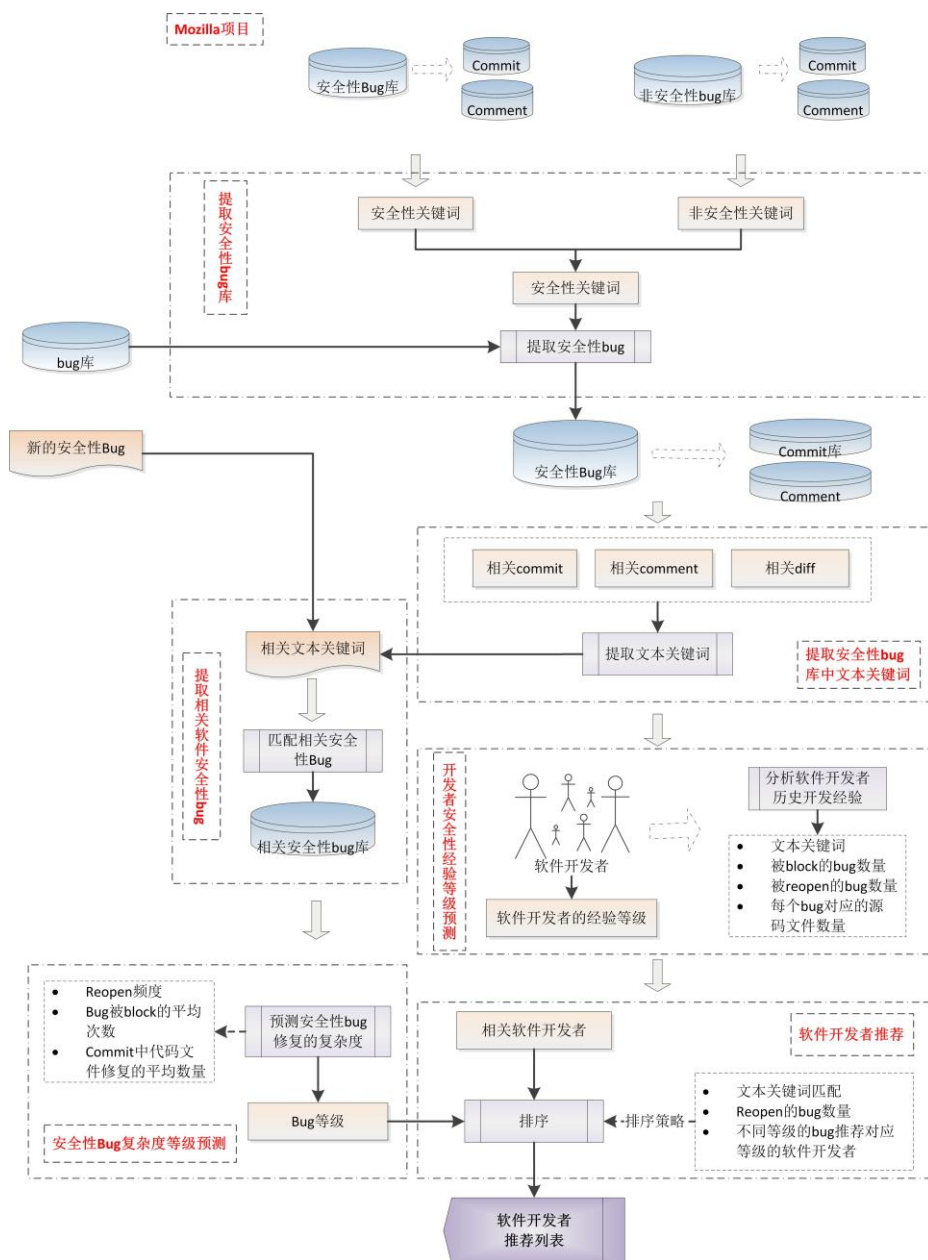


图 1 SecDR 软件开发者推荐流程图

配出的缺陷开发者的经验等级、安全性缺陷修复质量和历史修复安全性缺陷的复杂度等因素综合推荐出一个合适的开发者推荐列表。

1.1 提取安全性缺陷

SecDR 主要针对安全性缺陷推荐富有安全性缺陷相关经验的开发者。然而, 大部分开源系统中并没有将安全性缺陷独立管理, 因此本文首先识别出缺陷库中的安全性缺陷, 创建独立安全性缺陷库作为 SecDR 的训练数据。文中识别安全性

缺陷的技术主要借鉴 Gegick 等人所提出的基于文本挖掘的识别技术^[1]。如图 1“提取安全性缺陷库”模块所示, 本文主要采用开源项目(如 Mozilla¹)中所有缺陷对应的 commit 库和 comment 数据; 抽取关键词, 构建安全性关键词词库; 最后, 根据安全性关键词词库提取项目中与安全性相关的缺陷组成一个新的安全性缺陷库。

¹ <https://www.mozilla.org/en-US/security/advisories/>

应的所有 commit 中, 一共修改过的文件数量。

根据以上信息可以预测新缺陷的修复复杂度。首先统计整体安全性缺陷库中被 reopen 的缺陷比例, 被阻止的缺陷比例, 以及每条缺陷平均修改的源码文件数量; 然后计算相关安全性缺陷库中相对的 3 个数值; 最后预测新缺陷的修复复杂度, 其计算方法如式(3):

$$C_n = \lceil r_r - r_w \rceil + \lceil b_r - b_w \rceil + \left\lceil \frac{f_r - f_w}{f_r + f_w} \right\rceil \quad (3)$$

其中, r_r 和 r_w 分别代表相关安全性缺陷库和整体安全性缺陷库中被 reopen 的缺陷比例, b_r 和 b_w 分别代表相关安全性缺陷库和整体安全性缺陷库中被阻止的缺陷比例, f_r 和 f_w 分别代表相关缺陷库和整体安全性缺陷库中平均每条缺陷修改的文件数量, C_n 代表新缺陷修复复杂度的度量值。从式(3)中, 可以看到每项的取值都在(-1,1)的范围内, 而所有项都取原值的上界, 即最终取值{0,1}两种情况。所以 C_n 的最终取值有四种情况{0,1,2,3}。当 $C_n=0$ 时, 定义新缺陷的修复复杂度为“简单”, 当 $C_n=1$ 或 2 时, 定义新缺陷的修复复杂度为“一般”, 当 $C_n=3$ 时, 定义新缺陷的修复复杂度为“困难”。

1.5 开发者经验等级预测

由于安全性缺陷修复的及时性、高质量要求等特征, 本文主要从两个角度分析软件开发人员的历史经验: (1)软件开发者的历史开发内容; (2)软件开发缺陷修复的专业性(质量和复杂度两个因素)。

SecDR 主要利用一些主题词语反映软件开发者的历史开发内容, 这些主题词语主要来源于: (1)开发者历史修复过的缺陷的描述; (2)修复缺陷时修改过源码文件的相关 diff 内容; (3)相关缺陷对应的 comment 信息。

SecDR 主要从三个方面反映开发者修复安全性缺陷的专业能力: (1)开发者历史修复的被阻止缺陷的比例(复杂度因素); (2)被 reopen 的比例(质量因素); (3)平均每条缺陷的修改文件的数量(复杂度因素)。本文假设开发者能越好地修复复杂的缺陷, 他们修复该缺陷的专业性就越强。

通过对开发者的经验等级进行评估, SecDR 可根据缺陷修复的复杂度差异有针对性地推荐出不同经验级别的开发者。从而实现: 推荐初级软

件开发者修复简单的安全性缺陷, 高级软件开发者推荐复杂的安全性缺陷。这样可避免总是高级软件开发者的偏向推荐, 从而让开发团队更好地协作工作。具体评估方法如式(4):

$$E_p = \left\lceil \frac{k_p - k_a}{k_p + k_a} \right\rceil + \lceil b_p - b_a \rceil + \lceil r_p - r_a \rceil + \left\lceil \frac{f_p - f_a}{f_p + f_a} \right\rceil \quad (4)$$

其中, k_p 和 k_a 分别表示该开发者拥有的关键词量和系统中平均每位开发者拥有的关键词量, b_p 和 b_a 分别表示该软件开发者和系统中平均每位开发者修复过的被阻止缺陷的比例, r_p 和 r_a 分别表示该软件开发者和系统中平均每位开发者修改过被 reopen 的缺陷的比例, f_p 和 f_a 分别代表该开发者和系统中平均每位开发者在修复缺陷时, 平均每个缺陷修复的文件数量。 E_p 代表该软件开发者的等级系数, 可以看出, 式(4)中的每项都是取值上界, 取值范围为{0,1}。因此, E_p 的取值范围在{0,1,2,3,4}五种情况, 当 $E_p=0$ 或 1 时, 假定该开发者为“初级开发者”, 当 $E_p=2$ 时, 该开发者为“中级开发者”, 当 $E_p=3$ 或 4 时, 该开发者为“高级开发者”。因此, 最终开发者的经验等级被分为“初级”、“中级”、“高级”三个等级。

1.6 开发者推荐

开发者推荐主要分为两个阶段: 相关开发者推荐和排序。

SecDR 根据文本相似度推荐出与新缺陷有相关经验的开发人员。首先, 从相关安全性缺陷库中抽取出所有关键词并进行查重进一步去除重复单词, 然后将这些关键词与系统中每一位软件开发者的关键词进行相似度计算, 具体计算如式(5):

$$R_p = \frac{|k_p \cap k_r|}{|k_p|} \quad (5)$$

其中, k_p 为该开发者拥有的关键词数量, $k_p \cap k_r$ 表示该开发者与相关安全性缺陷共同拥有的关键词数量。式(5)的分母只考虑了软件开发者的关键词, 而没有考虑相关安全性缺陷库中关键词。这样可以将初级开发者更有效地推荐出来, 避免了高级软件开发者的偏向推荐, 以备下一步结合新缺陷的复杂度合理推荐出相应经验级别的

开发者。

在推荐出相关开发者后, 对他们进行排序。首先, SecDR 计算每位相关开发者的排序权值, 其计算方法如式(6):

$$W_p = R_p \times \theta^{|C_n - E_p|} \quad (6)$$

其中, R_p 是软件开发者的历史开发经验与新缺陷的相关系数, 其代表了开发者的历史开发内容与新缺陷的匹配程度, 其计算方法见公式(5), θ 表示一个权值调整参数, 其取值范围为(0,1), C_n 表示新缺陷的修改复杂度值, E_p 表示开发者的经验等级值, 他们的差值很好地实现了开发者不同经验等级的推荐。 r 代表开发者修复的缺陷是否被 reopen 过, 这个因素可以说明开发者的历史修复质量。三个参数的具体取值方法如式(7)(8)(9):

$$C_n = \begin{cases} 1(\text{简单}) \\ 2(\text{一般}) \\ 3(\text{复杂}) \end{cases} \quad (7)$$

$$E_p = \begin{cases} 1(\text{初级}) \\ 2(\text{中级}) \\ 3(\text{高级}) \end{cases} \quad (8)$$

$$r = \begin{cases} 0(\text{没有 reopened} \cdot \text{bug}) \\ 1(\text{有 reopened} \cdot \text{bug}) \end{cases} \quad (9)$$

在式(7)中, 当缺陷复杂度等级为简单时, 取值为 1, 一般时取值 2, 复杂时取值 3。式(8)中, 软件开发者为初级开发者时取值 1, 中级开发者取值 2, 高级取值 3。式(9)中, 当该软件开发者修复过的缺陷没有被 reopen 时取值 0, 否则取值 1。

最终, 每位开发者都对应一个排序权值 W_p , SecDR 根据 W_p 对所有相关软件开发者排序, 得到一个排序列表, 排序在前的开发者更适合修复该安全性缺陷。

2 技术实现

SecDR 的工具界面如图 3 所示, 主要分为四个部分。

界面 1 是 SecDR 的初始界面, 界面中主要展示了一个搜索框和一个搜索按钮。用户可以在搜索框中输入一条缺陷描述, 点击“SEARCH”按钮即可获得开发者推荐的相关结果。例如, 搜索: Overflow nsTSubstring::Replace Prep causes memory-safety bugs in string library。用户点击“SEARCH”按钮之后, 系统为用户推荐一系列与缺陷描述相关的开发者, 所有的推荐结果经过排序后按照适合度由高到低依次展示。

界面 2 中给出了推荐结果, SecDR 一共推荐 170 位软件开发者, 分别由 17 个分页页面展示, 每个页面展示 10 位。其中, eddy 为推荐结果中最

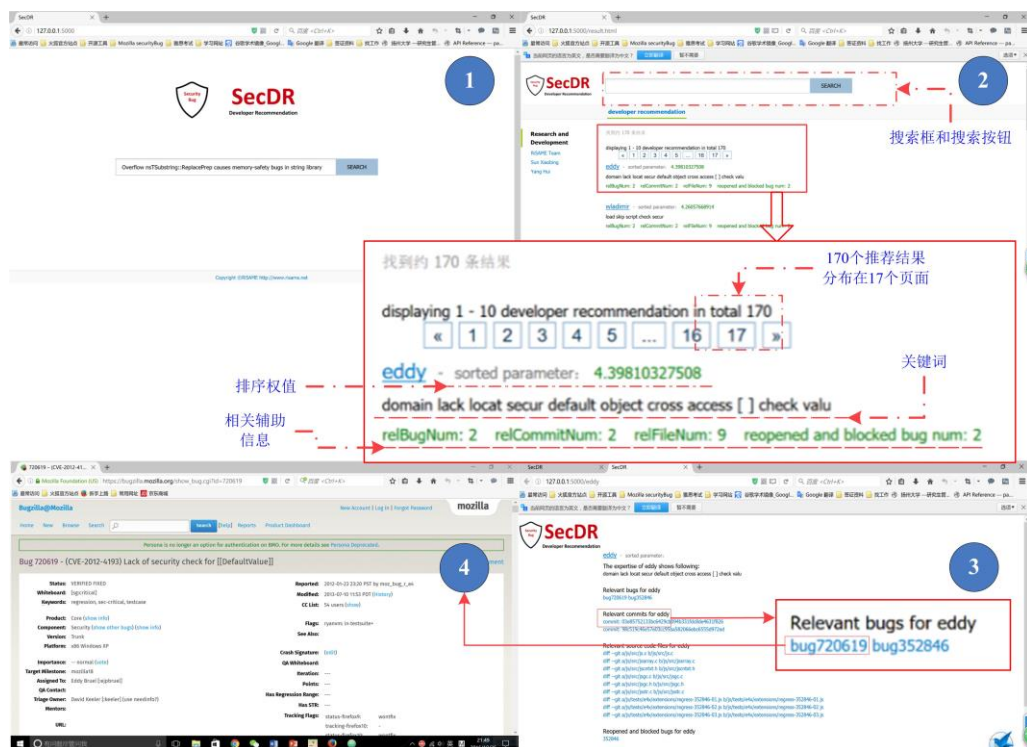


图 3 SecDR 工具界面

合适的开发人员。此外，为了方便交互，推荐结果界面中也给出了搜索框和对应的搜索按钮，用户可方便地在该界面上进一步搜索其他缺陷。SecDR 对所有推荐的开发人员的相关工作经验进行了分析。由界面 2 可看到，每一位开发者名字下面显示了一系列关键词，这些关键词反映该开发者与此缺陷相关的工作经验，可供搜索人员查看。此外，开发人员名字后面显示了排序的权值，由图可知，eddy 的排序权值为 4.40 最高，所以 he 是最合适的开发者。人名下面第二行绿色信息显示了每一位开发者在修复该缺陷时可以参考的相关辅助信息，其中，relbugNum 代表相关安全性缺陷的数量，relCommitNum 代表相关 commit 的数量，relFileNum 代表与该缺陷相关的源码文件数量，以及 reopened and blocked bug num 代表相关复杂缺陷的数量。这些信息都可以供开发人员参考，他们可根据这些信息进一步判断最合适的开发者从而做出最佳选择。

为了进一步辅助开发者修复该安全性缺陷，SecDR 在推荐合适的开发者之外，还给每一位开发者推荐一系列参考信息。这些信息主要帮助该开发者进一步理解和修复该缺陷。界面 3 中主要

给出了详细的相关缺陷、commit 和源码文件的信息。这些信息都是以链接的方式给出，开发者可以点击相关链接进一步查看详细内容。例如，用户可以在界面 3 中点击 bug 720619 链接，系统界面即可跳转到界面 4，并且进一步查看 bug 720619 了解相关内容。

3 实验验证

3.1 实验对象

为了验证 SecDR 推荐开发者的有效性，分别在 Mozilla、Libgdx² 和 ElasticSearch³ 三个开源项目上进行实验验证。经过 1.1 中安全性缺陷库的提取，这三个项目的安全性相关信息显示在表 2 中。其中，“项目名称”展示了实验对象的项目名称。对于 Mozilla 项目，“安全性缺陷数量”和“安全性开发人员数量”分别表示了本实验中应用到安全性缺陷的数量和对应的修复这些缺陷的开发者数量。在 Libgdx 和 ElasticSearch 两个项目中，“安全性缺陷数量”表示本实验中安全性缺陷数量，

² <http://libgdx.badlogicgames.com/nightlies/docs/api/>

³ <https://www.elastic.co/>

“安全性开发人员数量”代表了修复这些提取后的缺陷的软件开发者数量。“时间阶段”表示本实验应用到所有安全性缺陷报告的时间跨度。

表2 实验对象及其相关信息

项目名称	安全性 缺陷数量	安全性 开发人员数量	时间阶段
Mozilla	3258	234	2005年9月 -2016年7月
Libgdx	5281	485	2010年2月 -2016年5月
ElasticSearch	3261	192	2010年3月 -2016年5月

在表2中, Mozilla项目以Mozilla Firefox和Mozilla Thunderbird等为主,项目中的所有缺陷发布在bugzilla平台上,由人工分类安全性软件缺陷并有针对性地管理和修复。Libgdx是一个跨平台的2D/3D的游戏开发框架,由Java/C/C++语言编写而成,提供独立的接口支持各种游戏开发。ElasticSearch是一个基于Lucene的搜索服务器,提供了一个分布式多用户能力的全文搜索引擎。

3.2 实验设计

为了验证SecDR推荐的有效性,本文提出了以下三个实验研究问题:

问题1: SecDR推荐开发者的精度如何?

开发者推荐的精确度直接影响到软件维护的效率,因此,问题一主要验证SecDR能否有效地推荐合适的开发者修复安全性缺陷。

问题2: 与一般性缺陷的泛推荐技术相比,如:DR_PSF, SecDR推荐软件开发者的精度提高了多少?

泛推荐技术在推荐开发者时很少针对安全性缺陷的特性进行推荐。而SecDR在推荐开发者时,综合考虑了开发者的安全性缺陷修复的经验,而且实现了不同经验级别开发者的推荐。所以,问题2主要用于对比泛推荐方法DR_PSF^[15]和本方法SecDR在推荐开发者时的精确性,验证SecDR能否提高开发者的推荐精度。

问题3: 与已有项目开发者的实际分配相比,SecDR推荐开发者的合理性如何?

在缺陷实际修复过程中,很多缺陷的修复需要reopen,或者修复一段时间后需要重新修复,此时的缺陷分配可能不合理。而SecDR在推荐软件开发者时,采用分级推荐机制,主要实现了“初

级开发者修复简单缺陷,高级开发者修复复杂缺陷”的思想。所以,问题3主要验证SecDR推荐不同经验等级的开发者能否改善缺陷实际情况分配不合理的现象。

3.3 实验方法

本实验选择了安全性缺陷库中最近修复的200个缺陷作为测试集,其他缺陷作为SecDR的训练集。

问题1: 为了验证SecDR的推荐精度,使用recall值对本方法进行度量^[5]。其中recall的计算方法如式(10):

$$recall @ k = \frac{1}{r} \times \sum_{i=1}^r \frac{|RD(r_i) \cap AD(r_i)|}{|AD(r_i)|} \quad (10)$$

其中, r 代表训练缺陷的数量。 k 代表针对每条缺陷, SecDR推荐开发者人数的数量。 $RD(r_i)$ 代表SecDR针对 bug_i 推荐的软件开发者, $AD(r_i)$ 代表实际修复 bug_i 的软件开发者。在计算recall值时,针对每一条缺陷,分别推荐不同数量的开发人员,例如:令 $k=1, 5, 10$ 。

本实验不用另一个普遍的度量标准precision,是因为在实际修复过程中,只有一个软件开发者修复软件缺陷,所以推荐结果最多只有一个是正确的,而其他的推荐结果都是错误的,所以precision不适合本实验的度量。

问题2: SecDR主要针对安全性软件缺陷推荐合适的开发人员,而现有的开发者推荐方法并没有针对安全性缺陷推荐。因此,本问题主要对比泛推荐方法DR_PSF和SecDR在推荐安全性缺陷开发者时的效果。在对比时,主要应用gain值进行推荐结果的比较,其计算方法如式(11):

$$gain @ k_{SecDR-DR_PSF} = \frac{recall @ k_{SecDR} - recall @ k_{DR_PSF}}{recall @ k_{DR_PSF}} \times 100\% \quad (11)$$

其中, $recall @ k_{SecDR}$ 和 $recall @ k_{DR_PSF}$ 分别代表SecDR和DR_PSF推荐 k 个软件开发者的recall值,其计算方法如式(10)。

问题3: 本问题主要对比SecDR推荐与实际分配修复人员的合理性。在SecDR推荐开发者时,很多推荐结果与实际修复缺陷的人员不一致,此时,需要对比SecDR的推荐结果与实际修复人员哪一个更合理。主要从两方面进行对比:

(1) 首先统计了推荐十位开发人员时,实际开发者不在推荐列表中的所有缺陷。然后,根据SecDR预测缺陷的复杂度等级分别比较实际开发者的经验等级和推荐列表中第一位开发者的经验

等级, 从而评估两者之间的合理性。具体评估步骤如下: 1. 统计缺陷复杂度与实际开发者经验等级一致的缺陷比例; 2. 统计缺陷复杂度与推荐列表中排名第一开发者的经验等级一致的缺陷比例; 3. 比较两者比例数据, 如果实际开发者的比例较大, 说明实际分配较为合理; 否则, 推荐结果较为合理。

(2) 另一方面, 找出测试集中所有被 reopen 或被重新修复的缺陷, 统计出实际开发者在推荐列表中的比例和不在推荐列表中的缺陷比例。如果实际开发者不在推荐列表的比例较大, 说明 SecDR 可以更好地规避不合理的软件开发者的推荐。

这样, 通过以上步骤可以对 SecDR 的推荐结果与实际分配的情况进行对比, 从而评估出 SecDR 推荐开发者的有效性。

3.4 实验结果

问题 1: 本实验中, SecDR 分别推荐 1 位、5 位和 10 位软件开发者作为推荐结果。为了验证结果的有效性, 本文主要利用 recall 作为度量标准分别计算推荐 1、5、10 位开发人员时的 recall 值。表 3 的“SecDR”列主要显示了本方法的 recall 结果值。从表中可看到, 针对实验中的三个开源项目(Mozilla, Libgdx, ElasticSearch), SecDR 在推荐 1 位开发者时, 有效值在 8.0%到 11.0%范围内, 平均推荐结果是 9.5%; 推荐 5 位开发者时, 有效值范围在 17.5%到 23.5%, 平均结果为 20.5%; 当推荐 10 位开发者时, 范围在 36.0%到 51.5%, 其平均结果是 45.3%。实验结果表明, SecDR 在推荐 1 位、5 位和 10 位开发者时平均精度分别为 9.5%、20.5%和 45.3%。因此, SecDR 针对安全性修复的开发者推荐效果是可接受的。

问题 2: 为对比 SecDR 与已有的开发者推荐方法(如 DR_PSF)推荐安全性软件开发者的精确性, 本文计算 SecDR 对比于 DR_PSF 的增益(gain 值)。其计算结果在表 3 的“gain 值”列中, 可以看到, SecDR 相比于 DR_PSF 的 gain 值范围在 9.3%到 69.2%。在推荐 1 位、5 位、10 位软件开发者时, gain 值的平均结果为 41.8%, 19.2%, 38.5%。这说明了当分别推荐 1 位、5 位、10 位软件开发者时, SecDR 推荐开发者的精度比 DR_PSF 平均提高了 41.8%, 19.2%, 38.5%。因此, SecDR 相比于泛推荐技术 DR_PSF 可以更有效地针对安全性缺陷推荐开发人员。

表 3 SecDR 和 DR_PSF 开发者推荐结果

(recall@1, 5, 10) 以及对比结果 (gain@1, 5, 10)

项目名称	k	SecDR	DR_PSF	gain 值
Mozilla	1	8.0%	5.5%	45.5%
	5	21.0%	16.5%	27.3%
	10	51.5%	32.0%	60.9%
Libgdx	1	11.0%	6.5%	69.2%
	5	17.5%	13.5%	29.6%
	10	36.0%	32.0%	12.5%
ElasticSearch	1	9.5%	8.0%	18.8%
	5	23.5%	21.5%	9.3%
	10	48.5%	34.0%	42.6%
平均情况	1	9.5%	6.7%	41.8%
	5	20.5%	17.2%	19.2%
	10	45.3%	32.7%	38.5%

问题 3: 在对比实际开发者的分配效果与 SecDR 的推荐结果时, 主要从两个方面进行评估。首先, 将缺陷的复杂度预测值与实际开发者的经验值进行匹配, 当实际开发者的经验等级与缺陷的复杂度等级一致时, 认为该缺陷分配有效, 否则无效。经统计, 图 4 中蓝色柱条给出了实际开发者与缺陷复杂度一致的缺陷比例, 其中 ElasticSearch 项目的值为 32.8%, Libgdx 项目的值为 36.3%, Mozilla 的值为 29.2%, 三个项目的平均值为 32.8%。这说明了开发者实际分配的经验等级与缺陷复杂度等级一致的概率平均为 32.8%。

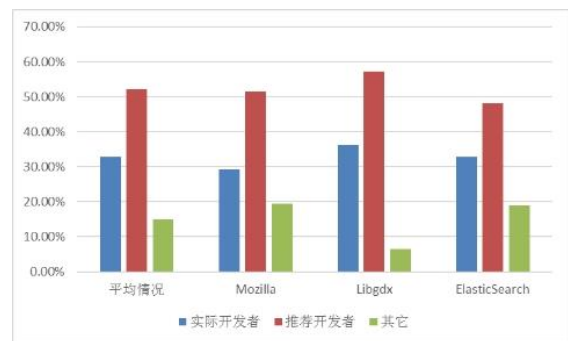


图 4 实际开发者和推荐开发者的经验与缺陷复杂度匹配结果柱状图

然后, 再匹配推荐列表中排名第一位软件开发者的经验等级和缺陷的复杂度, 同样计算匹配一致的缺陷比例。图 4 中的红色柱条代表了推荐结果与缺陷等级一致的缺陷比例, 其中 ElasticSearch 项目的值为 48.1%, Libgdx 项目的值

为 57.2%，Mozilla 的值为 51.4%，三个项目的平均值为 52.2%。这说明了 SecDR 推荐出的开发者等级与缺陷复杂度等级一致的概率平均为 52.2%。因此，SecDR 的推荐结果更好地匹配了缺陷修复的复杂度，这说明 SecDR 的推荐结果在软件开发者经验等级上推荐的更合理。

另一方面，找出测试集中所有被 reopen 或者重新修复的缺陷，分别统计出这些缺陷的实际开发者在推荐列表中的比例和不在推荐列表中的比例。

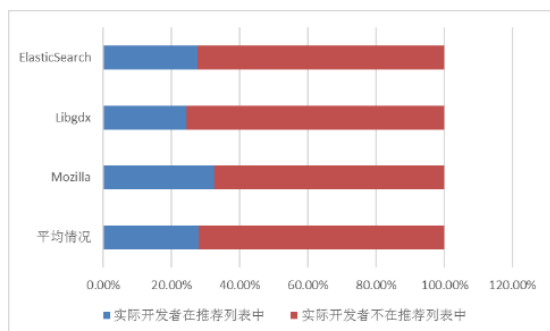


图5 面向测试集中被 reopen 或者重新修复缺陷的开发者推荐结果柱状图

从图 5 中可以看到，蓝色柱条统计了实际开发者在推荐列表中的缺陷比例。其中 ElasticSearch, Libgdx 和 Mozilla 三个项目的比例值分别是：27.6%，24.3%和 32.4%；而被 reopen 或者重新修复的缺陷的实际开发者不在推荐列表中的缺陷比例的平均值为 71.9%。因此，面对这些被多次修复的软件缺陷，SecDR 可以更好地规避不合理的软件开发者的推荐。

4 有效性威胁

在实验过程中存在一些有效性威胁，这些威胁可能会影响实验结果的有效性。

首先，在实验中，所有数据都抓取于开源网站。但是，在抓取过程中，由于网络不稳定的原因，仍然有少量数据会被爬虫程序跳过，从而引起数据不全的现象。虽然遗失的这些数据可能会影响 SecDR 的结果，但这些数据往往是少量的，其影响不会很大。所以，本实验的可靠性也不会受到很大干扰。

第二，本实验主要建立在三个不同应用场景下的开源软件对象上，但对于一些企业内部的软件数据，暂时还无法获取，所以本文不能确保 SecDR 对于企业数据推荐的有效性。另外，本文的实验主要建立在一些相对比较成熟的软件项目

上，这些项目里有很多经验丰富的高级软件开发者，同时也兼备一些经验不太丰富的初级开发者，此类数据的特征比较适合于 SecDR 的个性化推荐。而对于一些刚刚成立的软件项目，由于软件库中的历史数据不够丰富，很难分析出开发者的经验等级以及预测缺陷修复的复杂度，所以 SecDR 不太适合刚建立的软件项目推荐。另一方面，对于当今特别流行的一些软件项目，例如 Spark 等，由于这些项目非常流行，从而会吸引很多开发经验特别丰富的软件开发者新加入项目团队，SecDR 可能会将这些软件开发者分析为初级开发者，因为他们对该项目的开发情况可能还不太熟练。所以，SecDR 的推荐效果更适合一些成立时间较长但相对不太流行的软件项目。

第三，由于大部分软件项目中没有现成的安全性缺陷库，SecDR 主要根据 Mozilla 项目中的安全性缺陷库提取新项目中的一些与软件安全性相关的缺陷，并组成一个新的安全性缺陷库。由于项目的应用场景不同，可能会有很多缺陷被错误提取或者遗漏，从而导致 SecDR 处理的源数据的可靠性不足。另外，在提取安全性缺陷时，SecDR 主要根据 Mozilla 项目中的安全性关键词进行提取，例如 buffer overflow, memory leakage 等。但根据相关论文研究，利用关键词预测安全性缺陷时，很可能会遗漏很多安全性缺陷，但预测结果相对比较正确^[1]。这可以说明 SecDR 提取的安全性缺陷基本都是与软件安全性相关的，因此利用提取出的缺陷数据分析软件开发者的安全性相关经验是可靠的，所以 SecDR 推荐的软件开发者在一定程度上也是可靠的。

第四，在预测安全性缺陷修复的复杂度时，首先找到了与该缺陷相关的其他安全性缺陷，然后根据这些相关缺陷修复的复杂度对新缺陷进行预测。本文没有做缺陷预测结果准确性相关的实验，可能会有部分缺陷的预测结果会有偏差，从而影响 SecDR 对不同等级的软件开发者推荐的结果。另外，在研究实验问题三时，在第一方面验证还应用了缺陷的预测等级对推荐结果进行评估，这些数据可能都会影响实验结果。但本文主要研究的是软件开发者的个性化推荐，并不是预测缺陷修复的复杂度问题。因此，在预测复杂度时可能考虑的因素还不够全面，对于此方面的问题我们后面也会进一步研究。

第五，在回答问题一时，主要利用了 recall

值评估 SecDR 推荐的准确性, 也许还有其他的度量标准, 但 recall 度量是软件开发者推荐技术领域里使用最普遍的一种度量, 具有很好的代表性^[5]。在比较 SecDR 与 DR_PSF 推荐安全性开发者时, DR_PSF 也实现了软件开发者的个性化推荐, 但 DR_PSF 更多是区分拥有不同开发经验内容的软件开发者, 从而推荐相关软件开发者完成相应的软件缺陷。而本文主要是针对安全性缺陷进行推荐, 安全性缺陷的修复除了考虑开发者的开发内容, 可能更多地需要考虑软件开发者的安全性经验和修复质量。所以 DR_PSF 与 SecDR 的应用场景略有不同, 其表现的结果差距也较大。但这也说明 SecDR 的推荐更适合于解决安全性缺陷。在评估问题三时, 本文主要通过两方面角度对比推荐结果与实际分配, 也许实际分配时的情况更加复杂, 比如: 最合适的软件开发者暂时请假或有别的任务等, 但这些因素在评估中均被忽略。然而, 对于一些修复情况不太理想的软件缺陷, SecDR 推荐的结果与实际修复的软件开发者不一致的情况会更大, 而对于修复情况较好的缺陷推荐精度较高, 表明了本方法推荐结果的合理性。

最后, 在研究问题 3 对比实际开发者的分配效果与 SecDR 的推荐结果时, 我们采用了 reopen 等指标来评价实际分配修复人员的合理性。但是实际环境中, 由于缺陷修复过程较为复杂, 安全相关的缺陷修复尤为复杂, 因此依据 reopen 等指标来评价实际分配修复人员的合理性并不一定总是合理。后期我们将寻找更合适的度量来进一步验证该推荐的合理性。

5 相关工作

近年来缺陷修复者推荐技术得到广泛关注。Hossen 等人提出了 iMacPro 方法, 主要通过提取源码中的开发人员和维护人员, 基于他们的维护历史推荐一个最合适的软件开发者^[5]。Zhang 等人提出了一种 KSAP 的推荐方法, KSAP 通过构建复杂的缺陷网络推荐系统中的软件开发者^[8]。Yan 等人通过分析文档语义构建了一个 DPLSA 模型, DPLSA 主要建立了文档主题分布来分析软件开发者的历史开发经验从而推荐出合适的开发者^[9]。Zhang 等人开发了一个工具 BUTTER, 应用社交网络分析软件开发者的相关特征来推荐一系列的相关软件开发者^[10]。Wang 提出了一个方法

FixerCache, 在推荐软件开发者时分析了他们的历史开发行为进行推荐^[11]。Zhang 等人在推荐的过程中主要结合了主题模型和软件开发者之间的相互关系来推荐一个最合适的软件开发者解决软件缺陷^[12]。YANG 等人实现了个性化推荐, 在推荐时主要考虑了软件开发者的历史开发内容和开发习惯^[15]。Xia 等人提出了一个工具 DevRec, 这个工具在推荐软件开发者时不仅分析了软件缺陷库, 而且分析了软件开发者的历史开发经验^[16]。

相比以上这些泛推荐技术, 本文在针对安全性缺陷推荐开发者时, 不仅考虑了开发者的历史开发内容和软件安全性知识, 还分析了开发者的历史修复质量以及历史修复缺陷的复杂度。此外, SecDR 预测了新缺陷的修复复杂性, 实现了简单缺陷推荐初级软件开发者, 复杂缺陷推荐高级软件开发者, 从而实现了开发者的多经验级别推荐, 进一步提高了开发者推荐的准确性。

6 总结与展望

本文针对安全性缺陷的修复提出一种新的开发者推荐方法 SecDR, SecDR 在推荐开发者时不仅考虑了开发者的历史开发内容和软件安全性知识, 还分析了软件开发者的历史修复质量以及历史修复安全性缺陷的复杂度。此外, SecDR 预测了新缺陷的修复复杂性, 实现了简单缺陷推荐初级软件开发者, 复杂缺陷推荐高级软件开发者, 从而实现了开发者的多经验级别推荐。为了验证 SecDR 的有效性, 本文在三个开源项目上(Mozilla, Libgdx, ElasticSearch)做了相应的实验验证。在实验验证时, 本文对比了 SecDR 和传统的非安全性缺陷推荐方法 DR_PSF 推荐软件开发者的精度。通过对比实验可发现, 当分别推荐 1,5,10 位软件开发者时, SecDR 的推荐精度比 DR_PSF 平均高出 19%~42%。另一方面, 本文还对比了 SecDR 与实际开发人员的分配情况, 通过相关统计发现 SecDR 推荐软件开发者更合理。然而, 本方法在推荐软件开发者时首先识别了系统中已出现过的安全性缺陷, 在识别安全性缺陷时仅仅从主题词的角度进行了初步安全性预测, 后期仍然需要进一步加强安全性缺陷的识别指标, 从而更精确地构建出安全性缺陷库。另外, 在评估安全性缺陷开发者推荐的实验中, 将补充其他实验对象, 选取其他实验度量指标来评估本文开发者推荐技术, 进一步验证本文技术的有效性。目前的开发

项目中, 绝大部分开发者修复安全性缺陷的经验比较匮乏, 也就是比较难找到富有安全性经验的开发者, 未来将进一步完善推荐技术, 利用我们所提出的个性化开发者推荐技术^[14,27], 除了推荐开发者, 并且推荐安全性缺陷相关的方案辅助开发人员理解缺陷, 提高开发者的安全性缺陷的修复能力。

参考文献

- [1] M. Gegick, P. Rotella, and T. Xie. Identifying security bug reports via text mining: An industrial case study. In Proceedings of the 7th International Working Conference on Mining Software Repositories. 2010. 11–20.
- [2] S. Zaman, B. Adams, and A. E. Hassan. Security versus performance bugs: A case study on firefox. In Proceedings of the 8th Working Conference on Mining Software Repositories. 2011. 93–102.
- [3] J. Witschey, O. Zielinska, A. Welk, E. Murphy-Hill, C. Mayhorn, and T. Zimmermann. Quantifying developers' adoption of security tools. In Proceedings of the 2015 10th Joint Meeting on Foundations of Software Engineering, ESEC/FSE 2015. 2015. 260–271.
- [4] D. Mitropoulos, G. Gousios, and D. Spinellis. Measuring the occurrence of security-related bugs through software evolution. In Proceeding of the 16th Panhellenic Conference on Informatics. 2012. 117–122.
- [5] K. Hossen, H. H. Kagdi, D. Poshyvanyk. Amalgamating source code authors, maintainers, and change proneness to triage change requests. In: 22nd International Conference on Program Comprehension, ICPC 2014. 130–141.
- [6] H. H. Kagdi, M. Gethers, D. Poshyvanyk, M. Hammad. Assigning change requests to software developers, Journal of Software Maintenance, 2012, 24 (1) : 3–33.
- [7] R. Shokripour, J. Anvik, Z. M. Kasirun, S. Zamani. A time-based approach to automatic bug report assignment, Journal of Systems and Software, 2015, 102: 109–122.
- [8] W. Zhang, S. Wang, Q. Wang, Ksap: An approach to bug report assignment using knn search and heterogeneous proximity. Information and Software Technology, 2016, 70: 68–84.
- [9] M. Yan, X. Zhang, D. Yang, L. Xu, J. D. Kymer. A component recommender for bug reports using discriminative probability latent semantic analysis. Information and Software Technology, 2016, 73: 37–51.
- [10] W. Zhang, G. Han, Q. Wang, Butter: An approach to bug triage with topic modeling and heterogeneous network analysis. In Proceedings of the 2014 International Conference on Cloud Computing and Big Data, CCBDD 2014. 2014. 62–69.
- [11] S. Wang, W. Zhang, Q. Wang, Fixercache. Unsupervised caching active developers for diverse bug triage. In Proceedings of the 8th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement, ESEM 2014. 2014. 25:1–25:10.
- [12] T. Zhang, G. Yang, B. Lee, E. K. Lua. A novel developer ranking algorithm for automatic bug triage using topic model and developer relations. In Proceedings of the 2014 21st Asia-Pacific Software Engineering Conference, APSEC 2014. 2014. 223–230.
- [13] H. Yang, X. Sun, Y. Duan, B. Li, On the effects of exploring historical commit messages for developer recommendation. Chinese Journal of Electronics, 2016, 25(4): 658–664.
- [14] H. Yang, X. Sun, B. Li, J. Hu. Recommending developers with supplementary information for issue request resolution. In Proceedings of the 38th International Conference on Software Engineering Companion. 2016. 707–709.
- [15] H. Yang, X. Sun, Y. D. B. Li. DR_PSF: Enhancing developer recommendation by leveraging personalized source-code files. In The 40th IEEE Computer Society International Conference on Computers, Software and Applications. 2016. 239–244.
- [16] X. Xia, D. Lo, X. Wang, B. Zhou. Accurate developer recommendation for bug resolution. In 20th Working Conference on Reverse Engineering, WCRE 2013. 2013. 72–81.
- [17] Chen Ke-Han, Han Pan-Pan, Wu Jian. User clustering based social network recommendation. Chinese Journal of Computers, 2013, 36(2): 349–359. (in Chinese).
(陈克寒, 韩盼盼, 吴健. 基于用户聚类的异构社交网络推荐算法. 计算机学报, 2013, 36 (2): 349–359)
- [18] Yang Xi-Hui, Li Bin, He peng. A method of project recommendation based on the group software development. Journal of Chinese Computer Systems, 2015, 04(4): 671–676. (in Chinese).
(杨习辉, 李兵, 何鹏, 等. 一种群体软件开发中的项目推荐方法. 小型微型计算机系统, 2015, 04(4): 671–676.)
- [19] D. Mitropoulos, G. Gousios, and D. Spinellis. Measuring the occurrence of security-related bugs through software evolution. In Proceeding of the 16th Panhellenic Conference on Informatics. 2012. 117–122.
- [20] Z. Yin, D. Yuan, Y. Zhou, S. Pasupathy, and L. Bairavasundaram. How do fixes become bugs? In Proceedings of the 19th ACM SIGSOFT Symposium and the 13th European Conference on Foundations of Software Engineering, ESEC/FSE 2011. 2011. 26–36.

- [21] S. B. Saleem, Y. Yu, B. Nuseibeh. An empirical study of security requirements in planning bug fixes for an open source software project. *Month*. 2012.
- [22] M. Felderer, M. Büchler, M. Johns, A. D. Brucker, R. Breu, and A. Pretschner. Chapter one - security testing: A survey. *Advances in Computers*. 2016.101:1–51.
- [23] M. Felderer, P. Zech, R. Breu, M. Büchler, and A. Pretschner. Model-based security testing: a taxonomy and systematic classification. *Journal of Software Testing, Verification, and Reliability*, 2016, 26(2):119–148.
- [24] X. Sun, X. Liu, J. Hu, J. Zhu, Empirical studies on the nlp techniques for source code data preprocessing. In *Proceedings of the 2014 3rd International Workshop on Evidential Assessment of Software Technologies, EAST 2014*. 2014. 32–39.
- [25] M. F. Porter. An algorithm for suffix stripping. *Readings in information retrieval*. Morgan Kaufmann Publishers Inc. 1997. 130 - 137.
- [26] P. Rodeghero, D. Huo, T. Ding, C. McMillan, M. Gethers. An empirical study on how expert knowledge affects bug reports, *Journal of Software: Evolution and Process*, 2016, 28(7): 542-564.
- [27] X. Sun, H. Yang, X. Xia, B. Li. Enhancing Developer Recommendation with Supplementary Information via Mining Historical Commits. To appear in *Journal of Systems and Software*, 2017, <https://doi.org/10.1016/j.jss.2017.09.021>