

CẢI TIẾN VIỆC THỰC THI DÒ TÌM NHỮNG BÁO CÁO LỖI TRÙNG NHAU SỬ DỤNG THÔNG TIN CENTROID CLASS MỞ RỘNG

IMPROVING DETECTION PERFORMANCE OF DUPLICATE BUG REPORTS USING EXTENDED CLASS CENTROID INFORMATION

Nhan Minh Phúc¹

Tóm tắt – Trong việc bảo trì phần mềm, những báo cáo lỗi đóng một vai trò quan trọng đối với sự chính xác của những gói phần mềm. Thật không may, vấn đề báo cáo lỗi trùng nhau lại xảy ra, lí do có quá nhiều báo cáo lỗi được gửi đến trong những dự án phần mềm khác nhau, dẫn đến nhiều báo cáo lỗi bị trùng nhau và việc xử lí thường tốn nhiều thời gian và chi phí trong vấn đề bảo trì phần mềm. Nghiên cứu này sẽ giới thiệu một phương pháp dò tìm dựa vào thông tin centroid lớp mở rộng (Extended Class Centroid Information (ECCI)) để cải tiến việc thực thi dò tìm. Phương pháp này được mở rộng từ phương pháp trước đây chỉ sử dụng centroid mà không xem xét đến những tác động của cả hai lớp bên trong là inner và inter. Ngoài ra, phương pháp này cũng cải tiến việc sử dụng normalized cosine trước đây cho việc xác định sự giống nhau giữa hai báo cáo lỗi bằng denormalized cosine. Hiệu quả của phương pháp ECCI được minh chứng thông qua việc thực nghiệm với ba dự án mã nguồn mở là: SVN, Argo UML và Apache. Kết quả thực nghiệm cho thấy rằng, phương pháp ECCI cho kết quả dò tìm tốt hơn những phương pháp khác khoảng 10% trong tất cả các trường hợp khi được so sánh.

Từ khóa: dò tìm trùng lặp, báo cáo lỗi, thông tin centroid lớp, đặc điểm trọng lượng

Abstract – In software maintenance, bug reports play an important role in the correctness of

software packages. Unfortunately, the duplicate bug report problem arises because there are too many duplicate bug reports in various software projects. Handling with duplicate bug reports is thus time-consuming and has high cost of software maintenance. Therefore, this research introduces a detection scheme based on the extended class centroid information (ECCI) to enhance the detection performance. This method is extended from the previous one, which used only centroid method without considering the effects of both inner and inter class. Besides, this method also improved the previous use of normalized cosine in identifying the similarity between two bug reports by denormalized cosine. The effectiveness of ECCI is proved through the empirical study with three open-source projects: SVN, Argo UML and Apache. The experimental results show that ECCI outperforms other detection schemes by about 10% in all cases.

Keywords: duplication detection, bug reports, class centroid information, weighting feature.

I. GIỚI THIỆU

Trong vấn đề bảo trì phần mềm, việc tìm ra những lỗi cũng như những vấn đề không bình thường là một xử lí quan trọng để tránh những rủi ro. Thông thường, những tình huống này sẽ được miêu tả lại và gửi đến hệ thống quản lí báo cáo lỗi như Bugzilla, Eclipse... Sau khi những báo cáo lỗi được gửi, một hoặc nhiều người sẽ được giao nhiệm vụ phân tích những lỗi này và chuyển đến những lập trình viên phù hợp cho việc xử lí lỗi. Theo những nghiên cứu gần đây, vấn đề dò tìm lỗi trùng nhau đang nhận được nhiều sự quan tâm của các nhà nghiên cứu,

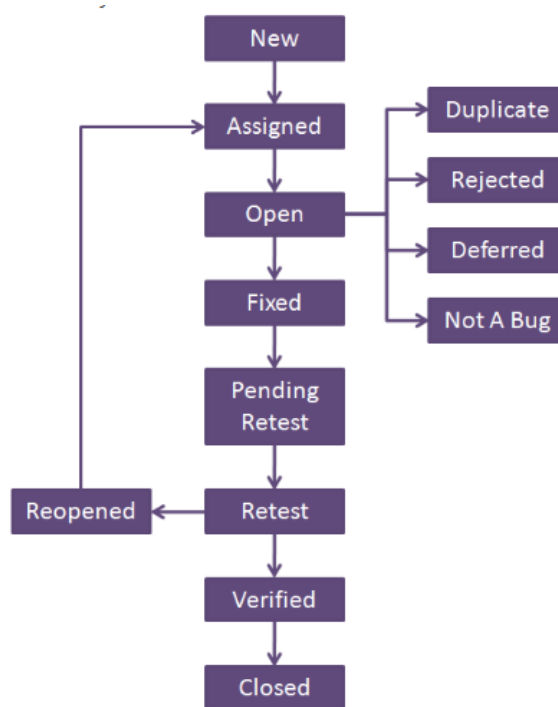
¹Bộ môn Công nghệ Thông tin, Khoa Kỹ thuật và Công nghệ, Trường Đại học Trà Vinh
Email: nhanminhphuc@tvu.edu.vn

Ngày nhận bài: 03/01/2017; Ngày nhận kết quả bình duyệt: 27/03/2017; Ngày chấp nhận đăng: 10/05/2017

lí do chính là số lượng báo cáo lỗi trùng nhau đã tăng đến 36%. Cụ thể với dự án của Eclipse được thống kê từ tháng 10/2001 đến tháng 8/2005, có 18,165 báo cáo lỗi, trong đó những lỗi trùng nhau chiếm tới 20%. Ngoài ra, theo dữ liệu của Firefox được thống kê từ tháng 5/2003 đến tháng 8/2005, có 2,013 báo cáo lỗi trùng nhau, trong đó 30% là những báo cáo lỗi trùng nhau. Gần đây theo Mozilla [1], từ 01/2009 đến 10/2012, mỗi tháng họ phải xử lý gần 2,837 lỗi với sự hỗ trợ gần 2,221 lập trình viên. Từ số liệu thống kê cho thấy, số lượng những báo cáo lỗi trùng nhau là rất lớn, điều này cho thấy tầm quan trọng của việc đưa ra những giải pháp trong việc xử lý lỗi trùng nhau là hết sức cần thiết và cấp bách. Vì vậy, việc nhận biết những báo cáo lỗi tự động đóng vai trò rất quan trọng và mang lại nhiều lợi ích. Thứ nhất, nó tiết kiệm được thời gian và công sức con người cho việc phân tích lỗi. Thứ hai, những thông tin chứa trong những báo cáo lỗi trùng nhau có thể rất hữu ích cho việc tìm ra nguyên nhân và cách xử lý lỗi.

Quy trình báo cáo lỗi được thực hiện như Hình 1. Khi một báo cáo lỗi vừa được gửi đến, nó sẽ được gán trạng thái "New". Sau đó, lỗi sẽ được bộ phận kiểm tra lỗi (tester) kiểm tra, nếu đây là lỗi thật sẽ được giao cho một lập trình viên tương ứng để xử lý, khi đó, trạng thái báo cáo lỗi sẽ là "Assigned". Trạng thái "Open" là khi lập trình viên bắt đầu phân tích và tiến hành xử lý lỗi. Nếu quá trình kiểm tra phát hiện báo cáo lỗi này đã được báo trước đó rồi, khi đó gán trạng thái là "Duplicate". Trạng thái "Rejected" được gán nhãn khi tester phát hiện lỗi này không có thật. Nếu báo cáo lỗi mà khi xử lý lỗi liên quan đến quá nhiều yếu tố có thể ảnh hưởng đến phần mềm, khi đó lỗi này sẽ được sửa trong phiên bản sau và báo cáo lỗi được dán nhãn "Deferred". Trạng thái "Not a bug" được gán khi tester phát hiện lỗi này không phải là một lỗi phần mềm mà thuộc chức năng phần mềm không hỗ trợ. Trạng thái "Fixed" được gán khi lập trình viên đã xử lý xong lỗi và chuyển đến bộ phận kiểm tra lỗi để kiểm tra lại. "Pending retest" là trạng thái mà báo cáo lỗi đang trong quá trình kiểm tra lại. "Retest" là trạng thái báo cáo lỗi được kiểm tra lại để biết lỗi đã sửa xong hay chưa. Nếu tester phát hiện vẫn còn lỗi, khi đó báo cáo lỗi sẽ được

gán "Reopen", và báo cáo lỗi này sẽ được xử lý lại. Nếu tester xác nhận báo cáo này đã được sửa xong, khi đó sẽ được gán nhãn "Closed".



Hình 1: Mô hình báo cáo lỗi

Theo tìm hiểu trong những năm gần đây, tình hình nghiên cứu về báo cáo lỗi trùng nhau trong các kho phần mềm mở tại Việt Nam còn rất hạn chế và hầu như chưa có, hầu hết những nghiên cứu chỉ tập trung ở nước ngoài. Tuy nhiên, về phương pháp phần lớn họ sử dụng mô hình không gian vector (Vector Space Model) kết hợp với việc tính độ giống nhau giữa hai báo cáo lỗi [1]–[8]. Gần đây phương pháp xử lý ngôn ngữ tự nhiên [9] đã được giới thiệu, phương pháp này được thực hiện kết hợp với thông tin thực thi của báo cáo lỗi, mặc dù kết quả cho thấy có sự cải thiện trong việc dò tìm lỗi trùng nhau so với những phương pháp trước, nhưng hiệu quả vẫn còn khá hạn chế. Chính vì điều này, phương pháp ECCI được giới thiệu với việc sử dụng xử lý ngôn ngữ tự nhiên cơ bản kết hợp với centroid class để tăng độ chính xác trong việc dò tìm những báo cáo lỗi trùng nhau, do phương pháp này xem xét đến những tác động của cả hai lớp bên trong là inner và inter. Kết quả thực nghiệm đã cho thấy

phương pháp này có sự cải tiến đáng kể so với những phương pháp trước đây.

II. VẤN ĐỀ ĐÒ TÌM LỖI TRÙNG NHAU

Khi người dùng sử dụng phần mềm mà phát sinh lỗi, thông tin báo cáo lỗi khi đó sẽ được gửi đến hệ thống quản lý phần mềm tương ứng. Một thông tin báo cáo lỗi là một dữ liệu có cấu trúc bao gồm nhiều trường như: tóm tắt lỗi (summary), mô tả lỗi (description), hệ điều hành sử dụng (OS)... như trong Hình 2.

Issue #4002
Issue list: (6 of 6) First Last Prev Next Show list

Issue #: 4002
Component: argouml
Subcomponent: AndroidMDA module
Status: NEW
Resolution:
Platform: All
OS: Windows XP
Version: 0.21.1
Priority: P3
Issue type: Direct
Target milestone:
Reporter: ashkali (Ashik Ali)
Add CC:
CC: None defined

Assigned to: rastaman (Ludovic Mabre)
URL:
* Summary: Not able to load XMI created by AndroidMDA into ArgouML
Status:
whiteboard:
Attachments: Date/Name: Description: Submitted

Description: Opened: Thu Feb 23 19:42:00 -0800 2006 Sort by: Oldest first | Newest first

Hi all, I've tried the latest ArgouML 0.20 today. The following are the steps performed:-
1) Created a new AndroidMDA J2EE project from Maven
2) Imported the AndroidMDA profile (androida-profile-3.0.xml.zip) into ArgouML. All the androida specific stereotypes and tagged values are loaded correctly.
3) Tried importing the XMI file generated by AndroidMDA into ArgouML 0.20, but the following exception occurred which stopped the import action:-

Hình 2: Ví dụ về các thông tin trong một báo cáo lỗi

Trường tóm tắt lỗi thường là những mô tả ngắn gọn về vấn đề lỗi phát sinh, trong khi đó trường mô tả lỗi thường được xem là quan trọng nhất, lí do trường này mô tả chi tiết về lỗi phát sinh cũng như thao tác người dùng thực hiện gây ra lỗi. Trường hệ điều hành sẽ cho biết thông tin hệ điều hành của người dùng khi sử dụng phần mềm gây ra lỗi, điều này cũng giúp dễ dàng hơn cho lập trình viên trong việc khắc phục lỗi phần mềm. Ngoài ra, nó cũng có phần bình luận cho những người báo cáo lỗi khác bình luận. Nếu một báo cáo lỗi là báo cáo đầu tiên, nó được gọi là báo cáo lỗi chính (master bug report). Ngược lại, nó sẽ được gắn lỗi trùng nhau sau khi được xử lí kiểm tra giống báo cáo lỗi chính. Trong Hình 3, báo cáo lỗi có mã số 983 được thông báo trùng với báo cáo lỗi trước đó có mã số 88. Để dò tìm những báo cáo lỗi trùng nhau, đầu tiên, chúng ta phải rút trích những thông tin văn bản từ những báo cáo lỗi. Thông thường, một báo cáo lỗi bao gồm những thông tin như nội dung tóm tắt lỗi, phần mô tả lỗi, hệ điều hành...

Bug 983 : Implement Import in the File Menu
Last modified: 2008-11-03 05:56

Project: processing Version: unspecified Component: pde

Status: RESOLVED Resolution: DUPLICATE of bug 88 Priority: P2 Severity: normal

Platform: All OS: All

Reporter: orgicus Assigned To: fry

Attachment	Type	Created	Size	Actions
Description		Opened: 2008-10-25 15:18		

An Import feature would be very handy for easier collaboration. Something like:
File > Import > from Directory
It would just copy the contents of the selected Directory to the default Sketches folder (or whatever is set in Preferences)
File > Import > from ZIP
It would unzip the contents of the selected ZIP to the default Sketches folder (or whatever is set in Preferences)
I can't implement this by my own at the moment, but I think it's an accomplishable task and it would help people exchange sketches easier.
Thank you!

Additional Comment #1 From fry 2008-11-03 05:56

*** This bug has been marked as a duplicate of 88 ***

Hình 3: Ví dụ một báo cáo lỗi trùng nhau trên SVN

III. PHƯƠNG PHÁP ĐÒ TÌM LỖI TRÙNG NHAU

A. Tổng quan về xử lí dò tìm lỗi

Để xác định một báo cáo lỗi vừa được người dùng gửi đến có trùng với những báo cáo lỗi đã được gửi trước đây hay không bằng phương pháp ECCI, phương pháp này được kế thừa và cải tiến từ phương pháp sử dụng đặc điểm lớp trong centroid [10], trong đó, chúng tôi xem xét cả hai đặc điểm trọng lượng bên trong lớp để cải thiện cho việc phân loại báo cáo lỗi, cũng như xem xét thông tin lớp liên quan đến trọng lượng từ. Trong nghiên cứu này, một lớp được định nghĩa như một cụm báo cáo lỗi trùng nhau. Trong tập dữ liệu, việc xem xét báo cáo lỗi trùng nhau dựa vào thông tin được đánh dấu trong báo cáo lỗi có dạng "This bug has been market as a duplicate of <bug report ID>" như ví dụ trong Hình 3. Khi đó, thông tin centroid có thể được trích ra từ mỗi cụm để tính sự giống nhau giữa các báo cáo lỗi. Toàn bộ quy trình xử lí báo cáo lỗi trùng nhau theo phương pháp ECCI được thực hiện như sau:

1. Xử lí ngôn ngữ tự nhiên
2. Tính trọng lượng đặc điểm lớp trong báo cáo lỗi

3. Tính ECCI centroid

4. Tính sự giống nhau giữa các báo cáo lỗi sử dụng Denormalized Cosine

5. Sắp xếp các báo cáo lỗi trùng nhau

Hình 4 cho thấy toàn bộ quy trình xử lý báo cáo lỗi trùng nhau theo phương pháp ECCI, bao gồm năm bước, các bước thực hiện sẽ được mô tả chi tiết bên dưới.

1) *Xử lý ngôn ngữ tự nhiên*: Như Hình 2 và Hình 3, nội dung báo cáo lỗi, ngoài những thông tin hữu ích mô tả lỗi, còn chứa những thông tin không thật sự có ích cho việc tự động dò tìm lỗi trùng nhau, ví dụ những từ "and, or, not, but, very..." hay những dấu câu như dấu gạch ngang, dấu ngoặc đơn... Vì vậy, việc loại bỏ những từ không cần thiết này rất quan trọng, ảnh hưởng nhiều đến sự chính xác của các phương pháp dò tìm. Trong bước này, mỗi báo cáo lỗi sẽ được rút trích thông tin từ hai trường chính trong báo cáo lỗi gồm trường tóm tắt lỗi (summary), mô tả lỗi (description), do các thông tin từ hai trường mô tả đầy đủ và có nghĩa để hỗ trợ việc xử lý lỗi. Sau đó, thông tin này sẽ được xử lý thông qua các bước xử lý ngôn ngữ tự nhiên ở mức cơ bản gồm tách từ (tokenization), tiếp theo là loại bỏ những từ không có nghĩa (stop words), ví dụ những từ như "the, and, or,..."; tiếp theo, tiến hành chuyển tất cả các dạng biến thể của một từ trở về từ gốc (stemming). Những thao tác xử lý ngôn ngữ tự nhiên cơ bản này được hỗ trợ bởi công cụ hỗ trợ WTool (Word Vector Tool). Công cụ này giúp việc xử lý các thao tác xử lý ngôn ngữ tự nhiên nhanh và dễ dàng hơn.

2) *Tính trọng lượng đặc điểm lớp trong báo cáo lỗi*: Trong quy trình xử lý báo cáo lỗi, việc tính đặc điểm trọng lượng lớp vô cùng quan trọng, nó ảnh hưởng trực tiếp đến kết quả xác định sự giống nhau giữa các báo cáo lỗi. Mỗi từ trong các báo cáo lỗi sẽ được xác định và chuyển sang mô hình không gian vector tương ứng với một trọng lượng. Phương pháp ECCI được thừa kế và cải tiến từ Class-Feature-Centroid(CFC) [11], [10] và trọng lượng đặc điểm lớp [12]. Trong CFC, trọng lượng của từ w_{ij} được tính như sau:

$$w_{ij} = b \frac{DF_{t_i}^j}{C_j} \times \log\left(\frac{|C|}{CF_{t_i}}\right)$$

Bảng 1: Các công thức tính trọng lượng bên trong lớp inner

Tên công thức	Chức năng
EXP-DF (CFC)	$I_{inner}^i = b \frac{DF_{t_i}^j}{C_j}$
TF	$I_{inner}^i = tf_{ijk}$
EXP-TF	$I_{inner}^i = b^{tf_{ijk}}$
EXP-TF-DF	$I_{inner}^i = b^{tf_{ijk} \times \frac{DF_{t_i}^j}{C_j}}$

Trong đó, t_i là từ (term) trong báo cáo lỗi, $DF_{t_i}^j$ là số báo cáo lỗi chứa t_i của lớp C_j , $|C_j|$ là số báo cáo lỗi trong lớp C_j , $|C|$ là tổng số lớp, $CF_{(t_i)}$ là số lớp chứa t_i , và b là tham số lớn hơn một, dùng để điều chỉnh cho trọng lượng w_{ij}

trong đó CFC, $b \frac{DF_{t_i}^j}{C_j}$ xem xét đến số báo cáo lỗi chứa mức độ xuất hiện thường xuyên của một từ bên trong lớp. Công thức log xem xét mức độ giống như IDF (inverse document frequency) truyền thống. ECCI được cải tiến từ CFC và trên cơ sở dựa vào [11]. Khi đó, mức độ thường xuyên của một từ tf_{ijk} của t_i trong báo cáo lỗi d_k , thuộc lớp C_j được tính như sau:

$$tf_{ijk} = \frac{fre(t_i)}{fre(t_i) + d + h \times \frac{dl}{dl_{avg}}}$$

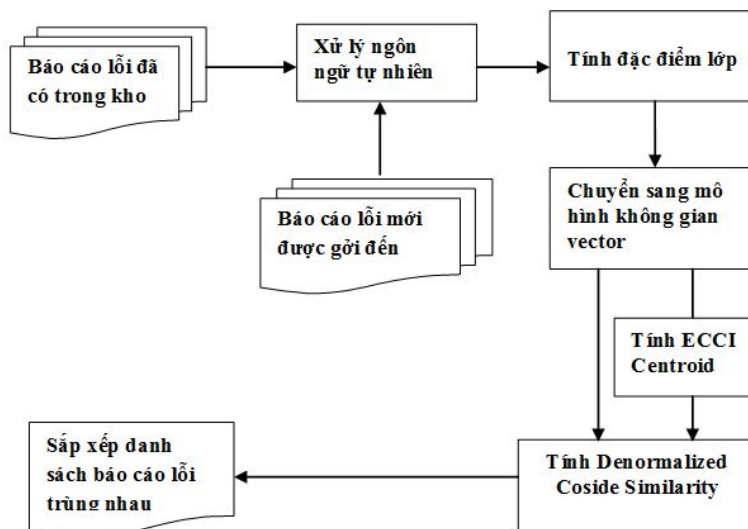
Trong đó, $fre(t_i)$ là số lần xuất hiện của t_i trong báo cáo lỗi d_k hoặc của lớp C_j , d là tham số điều chỉnh tránh cho mẫu số bằng 0, h là tham số ảnh hưởng đến chiều dài của báo cáo lỗi, dl là chiều dài của báo cáo lỗi d_k hoặc tổng chiều dài của báo cáo lỗi trong lớp C_i , dl_{avg} là trung bình của chiều dài các báo cáo lỗi. Nếu $t_i \in d_k$, khi đó dl_{avg} được tính như sau:

$$dl_{avg} = \frac{\sum_{d_m \in C} dl(d_m)}{\sum_{C_n \in C} |C_n|}$$

Trong đó, $|C_n|$ là số báo cáo lỗi trong C_n Nếu $t_i \in C_j$ nhưng $t_i \notin d_k$, khi đó:

$$dl_{avg} = \frac{\sum_{d_m \in C} dl(d_m)}{|C|}$$

Trong đó, $|C|$ là tổng số lớp, d và h là hai tham số và nó có thể nằm trong một khoảng giá trị tùy theo tập dữ liệu. Tuy nhiên, nghiên cứu này chỉ xác định $0.3 \leq d \leq 0.8$ và $1.5 \leq h \leq 20.0$ để tìm ra giá trị tốt nhất cho d và h .



Hình 4: Ví dụ một báo cáo lỗi trùng nhau trên SVN

- Chỉ số tác động bên trong lớp inner

Với việc mở rộng thông tin dựa vào lớp, khi đó, bốn công thức để tính chỉ số tác động bên trong lớp inner được giới thiệu, và được tiến hành thực nghiệm để tìm ra một công thức tốt nhất. Bảng 1 cho thấy bốn công thức dùng để tính trọng lượng bên trong lớp inner.

- Chỉ số tác động bên trong lớp inner

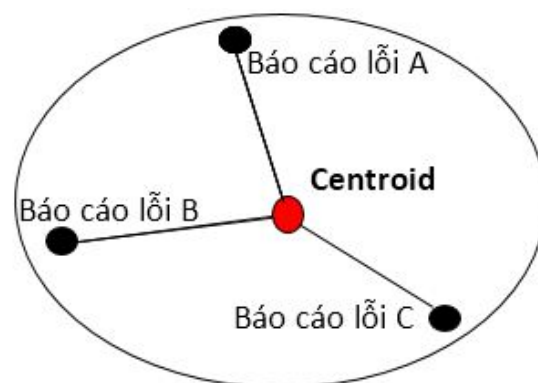
Để tăng cường độ chính xác trong việc phân loại báo cáo lỗi đối với chỉ số bên trong lớp I_{inner} , trong trường hợp này, ta sử dụng theo phương pháp CFC:

$$I_{inner}^i = \log\left(\frac{|C|}{CF_{t_i}}\right)$$

Nếu từ t_i xuất hiện trong tất cả các lớp, khi đó $I_{inner}^i = 0$, do $|C| = CF_{t_i}$, Nếu từ t_i xuất hiện chỉ trong một lớp, khi đó $I_{inner}^i = \log|C|$. Trong trường hợp này, t_i có sự phân biệt tốt nhất trong các lớp báo cáo lỗi trùng nhau.

3) Centroids và ECCI centroids: Phương pháp trong [2] sử dụng mô hình không gian vector cho cụm báo cáo lỗi của centroid. Trong phương pháp này, những báo cáo lỗi trùng nhau của cùng một nhóm thì được xem như một cụm, và vector centroid chính là trung bình cộng của các báo cáo lỗi trong cùng nhóm này như trong Hình 5, khi đó, được xem như là một báo cáo lỗi mới. Điều này có nghĩa là khi một báo cáo mới được gửi đến, nó sẽ được so sánh với vector centroid của

những cụm đã có trong kho lỗi thay cho việc so sánh với từng báo cáo lỗi. Trong khi đó, centroid mở rộng sử dụng trong phương pháp ECCI cũng sử dụng giống centroid này, tuy nhiên, điểm khác biệt là nó sử dụng lớp, trong đó, xem xét đến các lớp inner và inter như đã đề cập phần 2) và 3) bên trong cùng một centroid. Điều này giúp cải thiện được việc so sánh chính xác hơn giữa hai báo cáo lỗi. ECCI centroids (EC) là một trong những thành phần quan trọng hỗ trợ việc tìm ra sự giống nhau giữa các báo cáo lỗi, nó là trung bình cộng của các vector báo cáo lỗi trong cùng một lớp C_j :



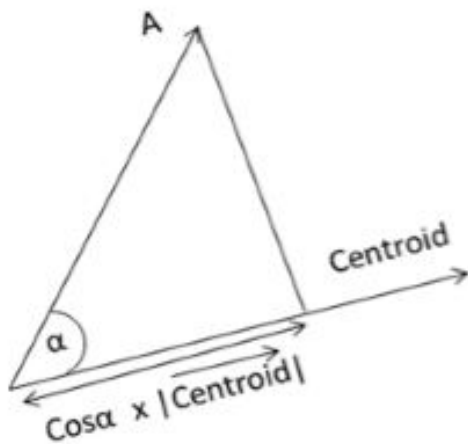
Hình 5: Mô hình centroid

$$\vec{EC}_j = \frac{1}{|C_j|} \sum_{\vec{d}_k \in C_j} \vec{d}_k$$

4) *Tính sự giống nhau giữa các báo cáo lỗi với denormalized cosine*: Bước này sẽ tiến hành xác định sự giống nhau giữa các báo cáo lỗi, khi có một báo cáo lỗi mới được gửi đến, nó sẽ được tính toán để xác định nó có trùng lặp với những báo cáo đã tồn tại trước đó hay chưa. Phương pháp truyền thống sử dụng cosine similarity truyền thống như sau:

$$Sim(\vec{d}_a, \vec{d}_b) = \frac{\vec{d}_a \cdot \vec{d}_b}{|\vec{d}_a| \cdot |\vec{d}_b|}$$

Tuy nhiên, với cosine similarity, nó không xem xét sự tác động của dữ liệu \vec{d}_b , mà chỉ cho thấy sự khác nhau giữa hai báo cáo lỗi \vec{d}_a và \vec{d}_b . Vì vậy, ECCI sử dụng denormalized cosine [5] để xem xét trong những thay đổi của centroid khác nhau trong các lớp, điều này giúp cải thiện việc dò tìm trùng nhau trong các báo cáo lỗi. Hình 6 cho thấy cách tính sử dụng denormalize cosine.



Hình 6: Denormalize cosine

5) *Sắp xếp các báo cáo lỗi trùng nhau*: Khi có những báo cáo mới được gửi đến, sẽ được thực hiện kiểm tra và so sánh xem có trùng với những báo cáo đã được gửi trước đó không? Phương pháp ECCI sử dụng thông tin centroid lớp mở rộng, khi đó, những báo cáo lỗi được gửi đến sẽ được tính và sắp xếp theo giá trị giống nhau nhất từ cao xuống thấp theo danh sách top 20.

Bảng 2: Thông tin về datasets của 3 dự án nguồn mở

Mô tả	ArgoUML	Apache	SVN
Ngôn ngữ	Java	C	C
Loại phần mềm	UML Tool	HTTP Server	SCM tool
Kho chứa lỗi	Tigris	Bugzilla	Tigris
Thời gian thu thập	02/2000-05/2007	01/2001-02/2007	03/2001-05/2007
Số báo cáo lỗi	4,613	2,771	2,296
Số báo cáo lỗi trùng	755	614	313

IV. KẾT QUẢ THỰC NGHIỆM

A. Môi trường thực nghiệm

Phương pháp ECCI được tiến hành thực nghiệm với ba kho báo cáo lỗi của những dự án phần mềm mở là Argo UML, Apache, và SVN. Thống kê chi tiết về ba kho phần mềm này được mô tả trong Bảng 2. Để đánh giá phương pháp dò tìm ECCI, đơn vị đo lường gọi là recall rate được sử dụng, nó được tính dựa trên bao nhiêu báo cáo lỗi có thể được dò tìm đúng trong danh sách những báo cáo lỗi trùng nhau, phương pháp này được sử dụng phổ biến trong việc đánh giá kết quả qua tìm báo cáo lỗi trùng nhau và nó được định nghĩa như sau:

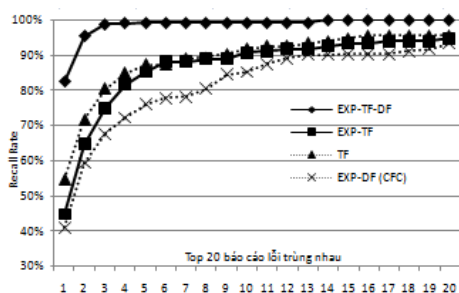
$$\text{Recall rate} = \frac{\text{Số những dự đoán đúng}}{\text{Tổng số những báo cáo lỗi trùng nhau}}$$

B. Những nhân tố tác động đến phương pháp ECCI

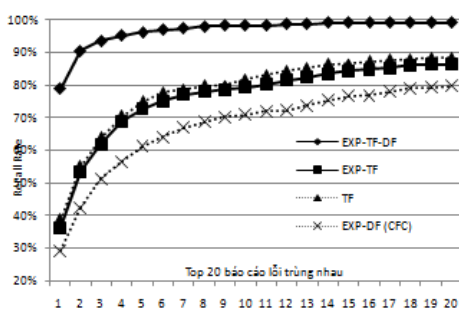
Trong phần này, chúng tôi thảo luận những nhân tố ảnh hưởng đến phương pháp ECCI. Thứ nhất là công thức được chọn cho việc tính trọng lượng đặc điểm lớp. Thứ hai và thứ ba liên quan đến việc xác định giá trị tốt nhất cho các tham số b, d, và h. Cuối cùng là công thức tính sự giống nhau giữa hai báo cáo lỗi sử dụng denormalized cosine.

1) *Trọng lượng đặc điểm lớp*: Trong phần trước, chúng tôi đã giới thiệu bốn công thức khác nhau cho việc tính trọng lượng từng từ trong những báo cáo lỗi dựa vào lớp bao gồm: EXP-DF, TF, EXP-TF, và EXP-TF-DF, khi tiến hành thực nghiệm để tìm ra công thức tốt nhất cho việc tính trọng lượng đặc điểm lớp, qua quan sát kết quả thực nghiệm cho thấy rằng EXP-TF-DF cho kết quả tốt nhất trong bốn công thức. Lí do EXP-TF-DF cho kết quả tốt nhất có thể giải thích là do hỗ trợ nhiều cho thông tin từ dựa vào lớp, điều này cũng giải thích lí do CFC cho kết quả không tốt bởi nó đã không xem xét sự tác động

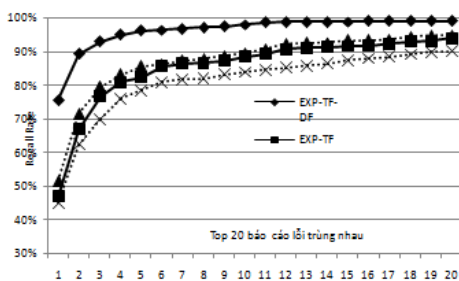
những từ thường xuyên xuất hiện trong báo cáo lỗi dựa vào lớp. Hình 7 cho thấy sự vượt trội của phương pháp EXP-TF-DF.



a) SVN



b) ArgoUML

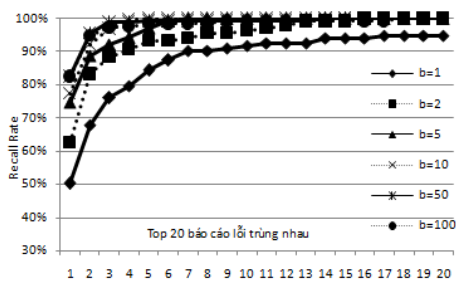


c) SVN

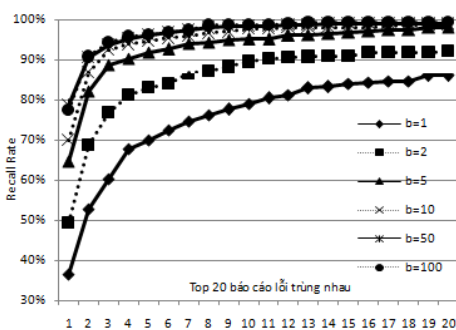
Hình 7: So sánh bốn dạng công thức tính trọng lượng lớp inner

2) *Tham số b*: Tham số b đóng vai trò quan trọng trong EXP-TF-DF. Do đó, việc thực nghiệm để tìm ra giá trị tốt nhất cho tham số b là cần thiết để tìm ra giá trị b tốt nhất trong ECCI. Kết quả thực nghiệm cho thấy rằng giá trị b không có thay đổi nhiều khi b lớn hơn 50, trừ dự án SVVN có tác động nhỏ nhưng không đáng kể. Do đó, ECCI đã sử dụng b=50 cho các thực nghiệm còn lại. Tuy nhiên, giá trị b có thể sẽ có thay đổi

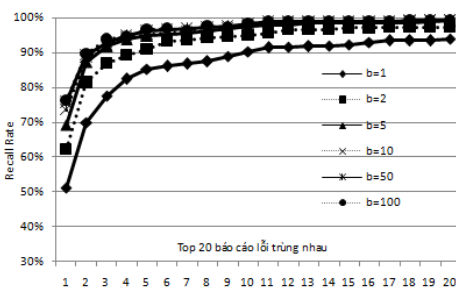
từ dữ liệu thực nghiệm, Hình 8 cho thấy kết quả thực nghiệm với tham số b.



a) SVN



b) ArgoUML

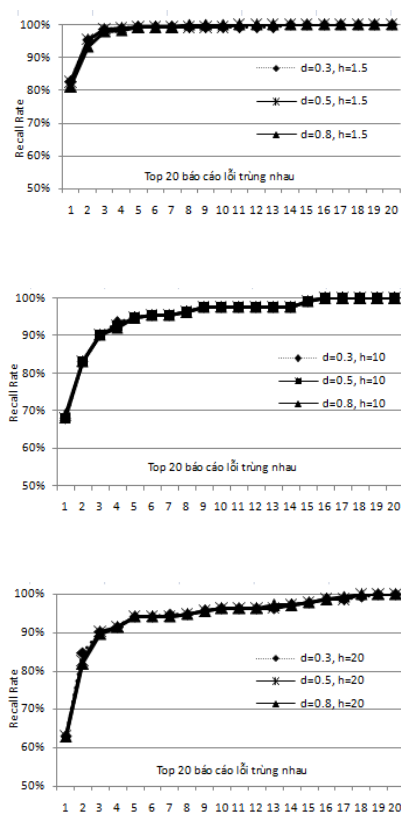


c) Apache

Hình 8: Tìm giá trị tốt nhất cho tham số b

3) *Tham số d và h*: Tham số d và h cũng ảnh hưởng trực tiếp đến EXP-TF-DF. Do đó, việc xác định giá trị tốt nhất cho d và h cũng góp phần quan trọng trong phương pháp ECCI. Và vì vậy, việc thực nghiệm cũng được tiến hành để thấy sự ảnh hưởng của d và h trong EXP-TF-DF, với $0.3 \leq d \leq 0.8$ và $1.5 \leq h \leq 20$. Do có nhiều sự kết hợp giá trị giữa d và h, nên bài báo này chỉ trình bày một vài trường hợp để minh họa chính cho việc tìm ra giá trị tốt nhất cho d và h. Kết quả thực nghiệm như trong Hình 9 đã xác định được

giá trị tốt nhất cho d và h với $d=0.3, h=1.5$. Tuy nhiên, giá trị này có thể đổi tùy theo những tập dữ liệu khác nhau.



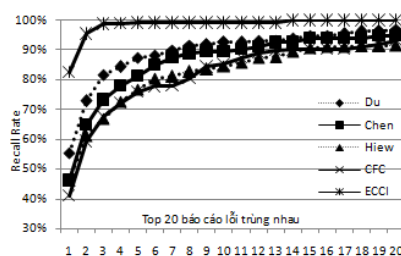
Hình 9: Tìm giá trị tốt nhất cho tham số d và h trên SVN

4) *Denormalized cosine measure*: Trong CFC, việc tính độ giống nhau giữa các báo cáo lỗi sử dụng denormalized cosine measure. Để thấy rõ hiệu quả của nó so với cách truyền thống sử dụng normalized cosine, việc thực nghiệm được tiến hành để so sánh hai phương pháp này. Quan sát kết quả thực nghiệm cho thấy rằng, phương pháp denormalized cosine cho kết quả tốt hơn phương pháp normalized cosine trong cả ba dự án.

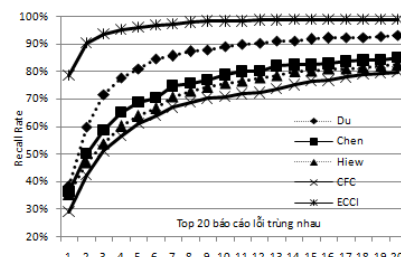
V. SO SÁNH PHƯƠNG PHÁP ECCI VỚI CÁC PHƯƠNG PHÁP KHÁC

Để thấy sự hiệu quả của ECCI với một số phương pháp dò tìm trùng nhau đã được công bố trước đây, cụ thể là phương pháp của Hiew [2], Chen [3], Du [4], CFC [11], thực nghiệm đã tiến hành để so sánh, kết quả so sánh được thấy

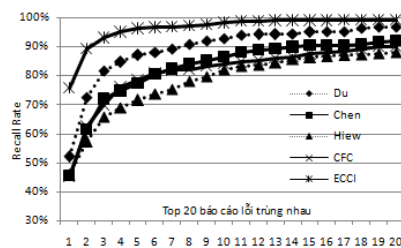
như Hình 10a đến 10c. Từ kết quả thực nghiệm, chúng ta thấy rằng phương pháp ECCI cho kết quả dò tìm tốt hơn so với các phương pháp khác. Điều này cho thấy sự hiệu quả của ECCI, khi xem xét các yếu tố liên quan đến thông tin lớp dựa vào centroid, trong việc tính trọng số của từ trong các báo cáo lỗi, cũng như hiệu quả của cách dùng phương pháp denormalized cosine.



a) SVN



b) ArgoUML



c) Apache

Hình 10: So sánh phương pháp ECCI với các phương pháp khác

VI. KẾT LUẬN

Việc dò tìm trùng nhau của những báo cáo lỗi là một trong những vấn đề quan trọng trong việc bảo trì phần mềm trong những năm gần đây. Trong bài báo này, chúng tôi giới thiệu một phương pháp dò tìm dựa vào thông tin centroid

lớp mở rộng (ECCI) để cải tiến việc thực thi dò tìm những báo cáo lỗi trùng nhau. Kết quả thực nghiệm từ ba dự án mã nguồn mở cho thấy phương pháp này mang lại hiệu quả cao trong việc dò tìm các báo cáo lỗi trùng nhau, đặc biệt là khi so sánh với các phương pháp được giới thiệu trước đây, phương pháp ECCI đã cho kết quả tốt hơn và hiệu quả hơn trong việc dò tìm những báo cáo lỗi trùng nhau khoảng 10% so với các phương pháp trước đó.

(**) Nghiên cứu này được tài trợ từ nguồn kinh phí nghiên cứu khoa học của Trường Đại học Trà Vinh.

TÀI LIỆU THAM KHẢO

- [1] Vincent, Bram Adams MCIS, Polytechnique Montreal, Québec. The Impact of Cross-Distribution Bug Duplicates, Empirical Study on Debian and Ubuntu. *IEEE 15th International Working Conference on Source Code Analysis and Manipulation (SCAM)*. 2015;p. 131–140.
- [2] Lyndon Hiew. Assisted Detection of Duplicate Bug Reports [Master Thesis]; May 2006. The University of British Columbia.
- [3] Zhi-Hao Chen. Duplicate Detection on Bug Report using N-Gram Features and Cluster Shrinkage [Master Thesis]; Jul 2011. YuanZe University.
- [4] Hung-Hsueh Du. A study of Duplication Detection Methods for Bug Reports based on BM25 Feature Weighting [Master Thesis]; Nov 2011. YuanZe University.
- [5] Stephen E Robertson, Steve Walker, Susan Jones, Micheline Hancock-Beaulieu, Mike Gatford. Okapi at TREC-3. *in Proceeding of the Third Text Retrieval Conference (TREC-3)*. 1994;p. 109–126.
- [6] Akihiro Tsuruda, Yuki Manabe, Masayoshi Aritsugi. Can We Detect Bug Report Duplication with Unfinished Bug Reports? *Software Engineering Conference (APSEC) 2015 Asia-Pacific*. 2015;p. 151–158. ISSN 1530-1362.
- [7] Chao-Yuan Lee, Dan-Dan Hu, Zhong-Yi Feng, Cheng-Zen Yang. Mining Temporal Information to Improve Duplication Detection on Bug Reports. *Advanced Applied Informatics (IIAI-AAI) 2015 IIAI 4th International Congress on*. 2015;p. 551–555. ISSN 1530-1362.
- [8] Chengnian Sun, David Lo, Xiaoyin Wang, Jing Jiang, Siau-Cheng Khoo. Discriminative model approach towards accurate duplicate bug report retrieval. *In ICSE 2010: Proceedings of the 32nd international conference on Software Engineering, Cape Town, South Africa*. 2010;IEEE Computer Society.
- [9] Xiaoyin Wang, Lu Zhang, Tao Xie, John Anvik, Jiasu Sun. An Approach to Detecting Duplicate Bug Reports using Natural Language and Execution Information. *in Proceedings of the 30th International Conference on Software Engineering (ICSE '08)*. 2008;p. 461–470.
- [10] Eui-Hong Han and George Karypis. Centroid-Based Document Classification: Analysis and Experimental Results. *in Proceeding of the Fourth European Conference on Principles of Data Mining and Knowledge Discovery (PKDD'00)*. 2000;p. 424–431.
- [11] Hu Guan, Jingyu Zhou, Minyi Guo. A Class-Feature-Centroid Classifier for Text Categorization. *in Proceeding of the 18th International Conference on World Wide Web*. 2009;p. 201–210.
- [12] Xiaoyan Zhang, Ting Wang, Xiaobo Liang, FengAo, YanLi. A Class-based Feature Weighting Method for Text Classification. *Journal of Computational Information System*. 2012;3:965–972.