

Information Retrieval Based Nearest Neighbor Classification for Fine-Grained Bug Severity Prediction

Yuan Tian¹, David Lo¹, and Chengnian Sun²

¹*Singapore Management University, Singapore*

²*National University of Singapore, Singapore*

{yuan.tian.2012,davidlo}@smu.edu.sg, suncn@comp.nus.edu.sg

Abstract—Bugs are prevalent in software systems. Some bugs are critical and need to be fixed right away, whereas others are minor and their fixes could be postponed until resources are available. In this work, we propose a new approach leveraging information retrieval, in particular BM25-based document similarity function, to automatically predict the severity of bug reports. Our approach automatically analyzes bug reports reported in the past along with their assigned severity labels, and recommends severity labels to newly reported bug reports. Duplicate bug reports are utilized to determine what bug report features, be it textual, ordinal, or categorical, are important. We focus on predicting fine-grained severity labels, namely the different severity labels of Bugzilla including: `blocker`, `critical`, `major`, `minor`, and `trivial`. Compared to the existing state-of-the-art study on fine-grained severity prediction, namely the work by Menzies and Marcus, our approach brings significant improvement.

I. INTRODUCTION

Software systems usually contain defects that need to be fixed after releases, and in some projects users are allowed to feedback on these defects that they encounter through bug reporting systems such as Bugzilla. With Bugzilla, users can report not only the description of the bug but also estimate the severity of the reported bugs. Unfortunately, although guidelines exist on how severity of bugs need to be assigned, the process is inherently manual that is highly dependent on the expertise of the bug reporters in assigning correct labels. Novice bug reporter might find it difficult to decide the right severity level. Developers (aka. Bugzilla assignee) can later adjust the severity [1] and use this severity information to prioritize which bugs to be fixed first.

As the number of bug reports made is large, a number of past studies have proposed approaches to help users in assigning severity labels, and development team in validating bug report severity [18], [14], [15]. All these approaches combine text processing with machine learning to assign severity labels from the textual description of the reports. Menzies and Marcus develop a machine learning approach to assign the severity labels of bug reports in NASA [18]. More recently, Lamkanfi et al. develop another machine learning approach to assign severity labels of bug reports in several Bugzilla repositories of open source projects [14]. In a later work, Lamkanfi et al. have also tried many different

classification algorithms and investigate their effectiveness in assigning severity labels to bug reports [15]. Menzies and Marcus assign fine-grained labels (5 severity labels used in NASA), while Lamkanfi et al. assign coarse-grained labels (i.e., binary labels: severe and non-severe).

The bug severity prediction tools are not perfect though and there is still room for improvement. Menzies and Marcus reported F measures (i.e., harmonic mean of precision and recall) of 14 to 86% for the different severity labels [18]. Lamkanfi et al. reported F measures of 65% to 75% on Bugzilla reports from different software systems [14]. Thus there is a need to improve the accuracy of the prediction tools further.

In this work, we propose an information retrieval (IR)-based nearest neighbor solution to predict the severity labels of bug reports. We first measure the similarity of different bug reports and based on this similarity we recover past bug reports that are most similar to it. There are various measures that have been proposed in the information retrieval community to measure the similarity between two textual documents [22], [31], [27], [26]. Some of the popular techniques are BM25 and its extensions [26]. BM25 technique and its extensions require some parameters to be learned. We leverage bug reports that have been marked as duplicate to set these parameters. Our hypothesis is that duplicate bug reports would help us to identify what features are important and what are not to measure the similarity between two bug reports. Based on a set of k nearest neighbors, the severity labels of these k similar bug reports are then used to decide the appropriate severity label for a new bug report.

In this work, we focus on predicting fine-grained bug severity labels. We investigate the effectiveness of our proposed approach and compare it with the past study by Menzies and Marcus [18]. Since our approach requires duplicate bug reports, we do not use the NASA data investigated by Menzies and Marcus. Rather, we analyze a large number of bug reports stored in Bugzilla bug tracking systems of Eclipse, OpenOffice, and Mozilla. We focus on predicting five severity labels of Bugzilla namely: `blocker`, `critical`, `major`, `minor`, and `trivial`. Following the work of Lamkanfi et al. [14], [15], we do not consider the severity label `normal` as this is the default option and “many

reports just did not bother to consciously assess the bug severity” [14], [15]. Thus, we treat these reports as unlabeled data.

Our experiments show that we could achieve a precision, recall, and F measure of up to 72%, 76%, and 74% for predicting a particular class of severity labels. Precision quantifies the amount of false positives, while recall quantifies the amount of false negatives. High precision and high recall mean less number of false positives and less number of false negatives respectively. F measure is the harmonic mean of precision and recall. Comparing with the state-of-the-art work on fine-grained severity level prediction by Menzies and Marcus, we show that for most bugs and most severity labels we could improve their approach significantly, especially on hard-to-predict¹ severity labels.

The following lists our contributions:

- 1) We propose an information retrieval based nearest neighbor solution, by leveraging duplicate bug reports, to predict fine-grained severity labels.
- 2) We have experimented our solution and compare it with the state-of-the-art work over a collection of more than 65,000 bug reports from three medium-large software systems: OpenOffice, Mozilla, and Eclipse.
- 3) We show that we can achieve a significant improvement over the state-of-the-art fine-grained bug severity prediction technique, especially on hard-to-predict severity labels.

The structure of this paper is as follows. In Section II, we describe some background material related to bug reporting and text pre-processing. In Section III, we elaborate our approach. We present our experiments and their results in Section IV. We discuss related work in Section V. We conclude and describe future work in Section VI.

II. BACKGROUND

In this section, we describe the bug reporting process, then present standard approaches to pre-process textual documents, and finally highlight $BM25F_{ext}$ to measure the similarity between structured documents.

A. Bug Reporting

To help improve the quality of software systems, software projects often allow users to report bugs. This is true for both open-source and closed-source software developments. Bug tracking systems such as Bugzilla are often used. Users from various locations can log in to Bugzilla and report new bugs. Users can report symptoms of the bugs along with other related information to developers. These include textual descriptions of the bug either in short or detailed form, product and component that are affected by the bug, and the estimated severity of the bug. The format of bug

reports varies from one project to another, but bug reports typically contain the fields described in Table I.

Developers (in particular bug triagers) would then verify these symptoms and fix the bugs. They could make adjustment to the severity of the reported bug. There are often many reports that are received and thus developers would need to prioritize which reports are more important than others – the severity field is useful for this purpose. As bug reporting is a distributed process, often the same bug is reported by more than two people in separate bug reports. This is known as duplicate bug report problem. The developer/triager would also need to identify these duplicate bug reports so as not to waste bug fixing efforts.

B. Text Pre-Processing

Tokenization. A token is a string of characters, and includes no delimiters such as spaces, punctuation marks, and so forth. Tokenization is the process of parsing a character stream into a sequence of tokens by splitting the stream at delimiters.

Stop-Word Removal. Stop words are non-descriptive words carrying little useful information for retrieval tasks. These include linking verbs such as “is”, “am” and “are”, pronouns such as “I”, “he” and “it”, etc. Our stop word list contains 30 stop words, and also common abbreviations such as “I’m”, “that’s”, “we’ll”, etc..

Stemming. Stemming is a technique to normalize words to their *ground* forms. For example, a stemmer can reduce both “working” and “worked” to “work”. This better allows a machine learning algorithm to capture the similarity between two or more words. We used Porter stemming algorithm [25] to process our text.

C. $BM25F$ and Its Extension

We present $BM25F$, and $BM25F_{ext}$. The first is a standard document similarity function, the latter is the extended $BM25F$ proposed in [26] to handle longer query documents.

$BM25F$ Similarity Function. $BM25F$ is a function to evaluate the similarity between two structured documents [21], [32]. A document is structured if it has a number of fields. A bug report is a structured document as it has several textual fields, i.e., *summary* and *description*. Each of the fields in the structured document can be assigned a different weight to denote its importance in measuring the similarity between two documents.

Before we proceed further, let’s define a few notations. Consider a document corpus D consisting of N documents. Also, each document d has K fields. Let’s denote the bag of words in the f^{th} field as $d[f]$ for $1 \leq f \leq K$.

$BM25F$ similarity function has two primary components which assign *global* and *local* importance to words. The *global* importance of a word t is based on its inverse document frequency (IDF). This IDF score is inversely

¹F measures of these labels are much lower than the others.

Table I
FIELDS OF INTEREST IN A BUG REPORT

Field	Description
Summ	<i>Summary</i> : Short description of the bug which typically contains only but a few words.
Desc	<i>Description</i> : Detailed description of the bug. This includes information such as how to reproduce the bug, the error log outputted when the bug occurs, etc.
Prod	<i>Product</i> : Product that is affected by the bug.
Comp	<i>Component</i> : Component that is affected by the bug.
Sev	<i>Severity</i> : Estimated impact of the bug to the workings of the software. In Bugzilla, there are several severity levels: <code>blocker</code> , <code>critical</code> , <code>major</code> , <code>normal</code> , <code>minor</code> , and <code>trivial</code> . There is also another severity level, <code>enhancement</code> which we ignore in this work, as we are not interested in feature requests but only defects.

Table II
EXAMPLES OF BUG REPORTS FROM MOZILLA BUGZILLA

	ID	Summary	Product	Component	Severity
1	525359	replying to an HTML message which includes a contenteditable div leaves Thunderbird compose unusable until restart (from incredimail for example)	Thunderbird	Message Compose Window	major
	543032	Impossible to answer a mail from thunderbird 3.01 after viewing an e-mail sent by Incredimail	Thunderbird	Message Compose Window	critical
2	537897	No way to select engines when setting up to use an existing account	Mozilla Services	Firefox Sync, Backend	normal
	543686	Everything is synced when logging in to an existing account	Mozilla Services	Firefox Sync, UI	normal
3	538953	Using Search bar AND a proxy with password authentication ... keeps asking the password at any key entered	Firefox	Search	normal
	544836	Proxy authentication broken while typing in the search field	Firefox	Search	major

proportional to the number of documents containing a word; it is defined in Equation 1.

$$IDF(t) = \log \frac{N}{N_t} \quad (1)$$

In Equation (1), N_t is the number of documents containing the word t .

Another component prescribes the *local* importance of a word t in a document d . This local importance, denoted as $TF_D(d, t)$, is defined in Equation 2. This is the aggregation of the local importance of the word t for each of document d 's field.

$$TF_D(d, t) = \sum_{f=1}^K \frac{w_f \times occurrences(d[f], t)}{1 - b_f + \frac{b_f \times size_f}{avg_size_f}} \quad (2)$$

In Equation (2), w_f is the weight of field f , $occurrences(d[f], t)$ is the number of times the word t occurs in field f , $size_f$ is the number of words in $d[f]$, avg_size_f is the average size of $d[f]$ for all documents in D , and b_f , which takes the value between 0 to 1, is a parameter that controls the contribution of the size of the fields to the overall score.

Based on the global and local term importance weights, given two documents d and q , each of which is a bag of words, the BM25F score of d and q is:

$$BM25F(d, q) = \sum_{t \in d \cap q} IDF(t) \times \frac{TF_D(d, t)}{k + TF_D(d, t)} \quad (3)$$

In Equation (3), the word t is common in d and q , and k , whose value is greater or equal to zero, is a parameter that controls the contribution of $TF_D(d, t)$ to the overall score. We notice that BM25F has a number of free parameters that need to be tuned: w_f and b_f for each document's field, and k . Given a document containing K fields, BM25F requires $(1 + 2K)$ parameters to be tuned. An optimization technique based on stochastic gradient descent has been used to tune these BM25F parameters [30].

$BM25F_{ext}$ Similarity Function. $BM25F$ is particularly developed to compute similarity of a short document (i.e., query) with a longer document. It is typically used for search engines, where user queries are usually short and consist of only a few words. However, bug reports are longer textual documents – the description field of a bug report can contain a few hundred words. Thus, since we want to have a similarity function that measures the similarity of two bug reports each of which are relatively long textual documents, there is a need to extend $BM25F$. Sun et al. [26] address this need by proposing $BM25F_{ext}$ which considers the term frequencies in queries; it has the following form.

$$BM25F_{ext}(d, q) = \sum_{t \in d \cap q} IDF(t) \times \frac{TF_D(d, t)}{k + TF_D(d, t)} \times W_Q$$

$$\text{where } W_Q = \frac{(l + 1) \times TF_Q(q, t)}{l + TF_Q(q, t)} \quad (4)$$

$$TF_Q(q, t) = \sum_{f=1}^K w_f \times occurrences(q[f], t) \quad (5)$$

In Equation (4), for each common word t appearing in document d and query q , its contribution to the overall $BM25F_{ext}$ score has two components: one is the product of IDF and TF_D inherited from $BM25F$; and the other is the local importance of word t in the document q – denoted as W_Q . W_Q follows the word weighting scheme of Okapi BM25 [16]. Parameter l , whose value is always greater than or equal to 0, controls the contribution of the local importance of word t in q to the overall score – if $l = 0$, then the local importance of t in q is ignored, and $BM25F_{ext}$ becomes $BM25F$.

In Equation (5), the contribution of each word t is the summation of the product of w_f , which is the weight of field f , with the number of occurrences of t in the f^{th} field of q . Different from TF_D , defined in Equation 2, to compute TF_Q , we do not perform any normalization. We do not perform normalization as retrieval is being done with respect to a single fixed query – we want to rank bug reports based on their similarities to a given query bug report.

$BM25F_{ext}$ requires an additional free parameter l in addition to those needed by $BM25F$. This brings the total numbers of parameters for $BM25F$ to $(2 + 2K)$. These parameters can be set by following a gradient descent approach presented in [26].

III. PROPOSED APPROACH

In this section, we describe our proposed approach. We first summarize our approach. We then highlight two major components of our approach.

A. Overall Framework

Our framework assigns a severity label to a bug report BQ in question by investigating prior bug reports with known severity labels in the pool of bug reports $BPool$. The high-level pseudocode of our approach, named IR Based Nearest Neighbour Severity Prediction Algorithm, is shown in Figure 1. The algorithm would first find the top-k nearest neighbors (Line 1) and then predict the label by considering the labels of these nearest neighbors (Lines 2-3).

Our framework thus consists of two major components: similarity computation, which is an integral part of finding nearest neighbors, and label assignment. In the similarity computation component, we measure the similarity between two bug reports. We leverage duplicate bug reports as training data to assign features that are important to measure how similar two reports are. We use an extended BM25 document similarity measure for the purpose. In the label assignment component, given a bug report whose severity is to be predicted, we take the nearest k bug reports based on the similarity measure. These k bug reports are then used to predict the label of the bug report.

Procedure INSPECT

Inputs:

BQ : Bug report in question

$BPool$: Historical bug report pool

Output: Predicted bug report severity label

Methods:

1: Let $NNSet$ = Find top-K nearest neighbors of BQ in $BPool$

2: Let $PredictedLabel$ = Predict label from $NNSet$

3: Output $PredictedLabel$

Figure 1. IR Based Nearest Neighbour Severity Prediction Algorithm

B. Similarity Computation

A bug report contains more than textual features, it also contains other information such as *product*, *component*, etc. We want to make use of all these features, textual and non-textual, to detect the similarity among bug reports. To do this, given two bug reports d and q , our similarity function $REP(d, q)$ is a linear combination of four features, with the following form where w_i is the weight for the i -th feature $feature_i$.

$$REP(d, q) = \sum_{i=1}^4 w_i \times feature_i \quad (6)$$

Each weight determines the relative contribution and the degree of importance of its corresponding feature. Features that are important to measure the similarity between bug reports would have a higher score. Each of the four features along with their definitions are given in Figure 2. There are two types of features: textual and non-textual; we elaborate them in the following paragraphs.

Textual Features. The first feature of Equation (7) is the textual similarity of two bug reports based on the *summary* and *description* fields as measured by $BM25F_{ext}$ similarity function described in Section II. The second feature is similar to the first one, except that *summary* and *description* fields are represented by bags of bigrams (a bigram is two words that appear consecutively one after the other) instead of bags of words (or unigrams).

Non-Textual Features. The other two features have binary values (0 or 1) based on the equality of the *product* and *component* fields of d and q .

The similarity function REP defined in Equation (6) has 16 free parameters in total. For $feature_1$ and $feature_2$, we compute textual similarities of d and q over *two* fields: *summary* and *description*. Computing each of the two features requires $(2 + 2 \times 2) = 6$ free parameters. Also, we need to weigh the contributions of each of the 4 features in Equation (6). Thus overall, REP requires $(2 \times 6 + 4) = 16$ parameters to be set. Table III lists all these parameters.

The above metric is similar to the one proposed by Sun et al. [26] except we remove several features: one is a binary feature that compares the types of the reports: defect,

$$feature_1(d, q) = BM25F_{ext}(d, q) // \text{of unigrams} \quad (7)$$

$$feature_2(d, q) = BM25F_{ext}(d, q) // \text{of bigrams} \quad (8)$$

$$feature_3(d, q) = \begin{cases} 1, & \text{if } d.prod = q.prod \\ 0, & \text{otherwise} \end{cases} \quad (9)$$

$$feature_4(d, q) = \begin{cases} 1, & \text{if } d.comp = q.comp \\ 0, & \text{otherwise} \end{cases} \quad (10)$$

Figure 2. Features in the Retrieval Function

Table III
PARAMETERS IN *REP*

Parameter	Description
w_1	weight of $feature_1$ (unigram)
w_2	weight of $feature_2$ (bigram)
w_3	weight of $feature_3$ (product)
w_4	weight of $feature_4$ (component)
$w_{summ}^{unigram}$	weight of <i>summary</i> in $feature_1$
$w_{desc}^{unigram}$	weight of <i>description</i> in $feature_1$
$b_{summ}^{unigram}$	<i>b</i> of <i>summary</i> in $feature_1$
$b_{desc}^{unigram}$	<i>b</i> of <i>description</i> in $feature_1$
$k_1^{unigram}$	k_1 in $feature_1$
$k_3^{unigram}$	k_3 in $feature_1$
w_{summ}^{bigram}	weight of <i>summary</i> in $feature_2$
w_{desc}^{bigram}	weight of <i>description</i> in $feature_2$
b_{summ}^{bigram}	<i>b</i> of <i>summary</i> in $feature_2$
b_{desc}^{bigram}	<i>b</i> of <i>description</i> in $feature_2$
k_1^{bigram}	k_1 in $feature_2$
k_3^{bigram}	k_3 in $feature_2$

enhancement, etc., another is a feature that computes the difference between the reported severities, and the other is a feature that computes the difference between the versions. Since we only consider defects, and we assume that severity label is not available, we could not use the first two of the three omitted features to compute similarity between bug reports. We do not use the last feature as we do not have the complete version information for all subject programs which requires manual crawling of the web. *REP* parameters are tuned using gradient descent. We take a training set consisting of duplicate bug reports, and follow the same approach as proposed in the work by Sun et al. [26]. We include the above description to ensure that our paper is self-explanatory.

C. Label Assignment

Leveraging the similarity measure, we locate the top- k nearest neighbors of a bug report in question. We then aggregate the contribution of each bug report to predict the label of the bug report. We compute the weighted mean of the labels of the neighbors as the predicted label. We map the labels into integers and order them from the most severe to

the least severe. The labels `blocker`, `critical`, `major`, `normal`, `minor`, and `trivial` are mapped to 0, 1, 2, 3, 4, and 5 respectively.

Consider a set of nearest neighbors $NNSet$ of a bug report BQ . Also let $NNSet[i]$ be the i th nearest neighbor, $NNSet[i].Label$ be the label of the i th nearest neighbor (expressed in integer), and $NNSet[i].Sim$ be the similarity of BQ with $NNSet[i]$. The predicted label is computed by the following formula:

$$\left\lfloor \frac{\sum_{i=0}^k (NNSet[i].Sim \times NNSet[i].Label)}{\sum_{i=0}^k (NNSet[i].Sim)} + 0.5 \right\rfloor$$

The above formula aggregates the label of each neighbor based on its similarity with the target bug report BQ . The higher is a neighbor similarity with BQ , the more powerful it is in influencing the label of BQ . The formula ensures that the label would fall into the range. We use the floor operation and the “+ 0.5” to round the resultant label to the nearest integer.

As bug reports with normal severity are treated as unlabeled data, we ignore the contributions of these reports. In case the k neighbors of a new bug report whose severity label is to be predicted are *all* assigned normal label, we simply assign label `major` to the new bug report.

Example. To illustrate the above, we present an example. Consider a bug report BQ , with top-3 neighbors N_1 , N_2 , and N_3 with labels 5, 4, and 3 respectively. Let the *REP* similarity scores of BQ with each of the neighbors to be:

$$REP(BQ, N_1) = 0.5$$

$$REP(BQ, N_2) = 0.45$$

$$REP(BQ, N_3) = 0.35$$

The assigned label of BQ would then be:

$$\begin{aligned} &= \left\lfloor \frac{\sum_{i=0}^3 (REP(BQ, N_i) \times N_i.Label)}{\sum_{i=0}^3 (REP(BQ, N_i))} + 0.5 \right\rfloor \\ &= \left\lfloor \frac{(0.5 \times 5 + 0.45 \times 4 + 0.35 \times 3)}{(0.5 + 0.45 + 0.35)} + 0.5 \right\rfloor \\ &= \left\lfloor \frac{(2.5 + 1.8 + 1.05)}{1.3} + 0.5 \right\rfloor \\ &= 4 \end{aligned}$$

IV. EXPERIMENTS

In this section, we highlight the datasets that we use in this study, followed by our experimental settings. We then present the measures used to evaluate the approaches, followed by our results. Finally, we also mention some threats to validity.

A. Datasets

We chose the bug repositories of three large open source projects: OpenOffice, Mozilla and Eclipse, as the three

projects have different backgrounds, implementation languages and users, which can help generalizing the conclusions of our experiments. OpenOffice is a multi-platform and multi-lingual office suite. Mozilla is a not-for-profit community producing open-source software and technologies used by other applications, such as the Firefox browser and Rhino JavaScript interpreter. Eclipse is a large project aiming to build a flexible development platform for all lifecycles of software development.

We extract three datasets from the open source projects by collecting reports submitted within a period of time. Each dataset only contains defect reports, whereas feature requests and maintenance tasks are filtered away. We use the final assigned severity labels in the defect reports as the ground truth. Table IV details the three datasets. We construct a training set by selecting the first M reports of which 200 reports are duplicates, in order to tune the parameters in the retrieval function REP , regardless of the size of the resultant bug report set for training. Those M reports are also used to simulate the initial bug repository for all experimental runs. The number M for the 3 datasets are given in sub-column # All of column *Training Reports* in Table IV. The rest of the reports are used for testing the prediction approach, shown in column *Testing Reports*.

B. Experimental Settings

We propose an online evaluation approach that mimics how severity prediction could be used in practice. At each experimental run, we iterate through the reports in the set of testing reports in chronological order. Once we reach a report R , we apply a severity prediction tool to predict the severity label of R . This would be the recommendation given to the user/developer on the severity of the bug report. As the accuracy of all existing severity prediction techniques are still low, humans/developers/triggers cannot be completely taken out from the picture. At the beginning of the next iteration, we add R and its true label (we assume triggers make the right decision and give a correct feedback) to the pool of bug reports $BPool$ in Figure 1. After the last iteration is done, we measure how good the recommendations are. Similar online evaluation approaches have been used in evaluating studies on the detection of duplicate bug reports [27] and on the recommendation of developers to fix bug reports [29].

Unfortunately, the classification based approaches employed in [18] (i.e., Severis) is slow. For around 4,000 bug reports of OpenOffice, employing the online evaluation approach would mean re-training the classification model for around 4,000 times – we re-train the model everytime a new user feedback is received. This took us more than 10 hours. As the number of bug reports increases the runtime increase in a super-linear fashion as at each step in the online evaluation approach more bug reports need to be investigated to train the model. Thus, we also evaluate the

existing approach in an offline manner – we take a set of bug reports that we use to train REP to train Severis. We then use trained Severis to assign labels to the remaining set of bug reports.

We perform both offline and online evaluation for Severis on OpenOffice bugs. We show that the results of these two evaluation approaches do not differ much. We only perform offline evaluation for Severis for the other two bug report datasets: Mozilla, and Eclipse. As our approach is fast and relies on nearest neighbors, we only do the online strategy.

C. Evaluation Measures

We use the standard measures of precision, recall, and F measure for each severity label to evaluate the effectiveness of Severis and INSPECT. F measure is the harmonic mean of precision and recall and it is often used to measure if an increase in precision/recall outweighs a loss in recall/precision. The same measures were used by Menzies and Marcus to evaluate Severis [18]. The definitions of precision, recall, and F measure for a severity label S^L are given below²:

$$precision(S^L) = \frac{\# \text{ reports correctly labeled with } S^L}{\# \text{ reports labeled with } S^L}$$

$$recall(S^L) = \frac{\# \text{ reports correctly labeled with } S^L}{\# \text{ reports that should be labeled with } S^L}$$

$$F \text{ Measure}(S^L) = 2 \times \frac{precision(S^L) \times recall(S^L)}{precision(S^L) + recall(S^L)}$$

D. Comparison Results

We compare INSPECT with parameter k set to 1 (i.e., 1-nearest neighbor) and Severis on the three datasets. We present the results in the following sub-sections.

1) *OpenOffice Results*: The result of INSPECT on bug reports of OpenOffice is shown in Table V. Different from the other three programs in OpenOffice there are only five severity levels [19]. We map them to `critical`, `major`, `normal`, `minor`, and `trivial`. Again we drop `normal` from our analysis. We note that we can predict the `critical`, `major`, `minor`, and `trivial` severity labels by F measures of 36.0%, 74.0%, 39.8%, and 22.2% respectively. The F measure is very good for `major` severity label but is poorest for `trivial` severity label.

The result for Severis (offline) is also shown in Table V. We note that Severis can predict the `critical`, `major`, `minor`, and `trivial` severity labels by F measures of 25.6%, 75.1%, 20.5%, and 1.2% respectively. Comparing these with the result of INSPECT, we note that INSPECT can improve the F measure for `critical`, `minor`, and `trivial` labels by a relative improvement of 41%, 94%,

²# reports refers to number of reports.

Table IV
DETAILS OF DATASETS

Dataset	Period		Training Reports			Testing Reports		
	From	To	#Duplicate	#All	#All - #Normal	#Duplicate	#All	#All - #Normal
OpenOffice	2008-01-02	2010-12-21	200	2,986	617	488	20,438	3,356
Mozilla	2010-01-01	2010-12-31	200	4,379	1,273	1,802	68,049	16,490
Eclipse	2001-10-10	2007-12-14	200	3,312	500	6,203	175,297	43,587

Table V
PRECISION, RECALL, AND F MEASURE FOR INSPect, SEVERIS [OFFLINE], AND SEVERIS [ONLINE], ON OPENOFFICE

Severity	INSPect			Severis [Offline]			Severis [Online]		
	Precision	Recall	F Measure	Precision	Recall	F Measure	Precision	Recall	F Measure
critical	33.2%	39.3%	36.0%	40.7%	18.6%	25.6%	58.5%	15.4%	24.4%
major	72.4%	75.6%	74.0%	63.9%	91.0%	75.1%	63.2%	96.2%	76.3%
minor	44.2%	36.2%	39.8%	39.0%	13.9%	20.5%	42.2%	7.7%	13.0%
trivial	26.8%	18.9%	22.2%	6.7%	0.7%	1.2%	60.0%	1.0%	2.0%

and 1748% respectively. For the `major` label, INSPect lose out to Severis by only 2%. Thus for OpenOffice, in general our proposed approach INSPect performs better than Severis.

Although expensive (in terms of runtime; it takes more than 10 hours to complete), we also run Severis using the online evaluation approach and present the result in Table V. We notice that the result using the online evaluation, although requires much more computation time, generally is not better than using offline evaluation. There is a small increase in F measure for `major` and `trivial`; However, for `critical` and `minor` there is a small reduction in F measure.

2) *Mozilla Results*: The result of INSPect on bug reports of Mozilla is shown in Table VI. We note that we can predict the `blocker`, `critical`, `major`, `minor`, and `trivial` severity labels by F measures of 32.6%, 65.9%, 54.3%, 35.9%, and 35.8% respectively. The F measure is very good for `critical` severity label but is poorest for `blocker` severity label.

The result for Severis is also shown in Table VI. Note that we only run the offline version of Severis as the online version takes much time, and our experiment with OpenOffice shows that employing online or offline evaluation does not affect the performance of Severis. We note that Severis can predict the `blocker`, `critical`, `major`, `minor`, and `trivial` severity labels by F measures of 0.4%, 65.1%, 59.7%, 3.4%, and 2.2% respectively. Comparing these with the result of INSPect, we note that we can improve the F measures for `blocker`, `critical`, `minor`, and `trivial` labels by a relative improvement of 8,038%, 1.2%, 957%, and 1,528% respectively. For the `major` label, INSPect lose out to Severis by 9.0%. Thus for Mozilla, in general our proposed approach INSPect performs better than Severis.

3) *Eclipse Results*: The result of INSPect on bug reports of Eclipse is shown in Table VII. We note that we can predict the `blocker`, `critical`, `major`, `minor`, and `trivial` severity labels by F measures of 26.0%, 29.0%, 57.8%,

40.3%, and 26.5% respectively. The F measure is very good for `major` severity label but is poorest for `blocker` severity label.

The result for Severis is also shown in Table VII. We note that Severis can predict the `blocker`, `critical`, `major`, `minor`, and `trivial` severity labels by F measures of 0.0%, 28.5%, 56.0%, 0.2%, and 0.0% respectively. The F measures of Severis are zeros for `blocker` and `trivial` as it does not assign any bug report to those severity labels. Comparing these with the result of INSPect, we note that we can improve the F measure for `blocker`, `critical`, `major`, `minor`, and `trivial` labels by a relative improvement of infinity, 1.7%, 3.2%, 20,055%, and infinity, respectively. INSPect does not lose out to Severis for any label. Thus for Eclipse, clearly INSPect performs better than Severis.

E. Varying parameter k

Our proposed approach INSPect takes in one user defined parameter k . In the previous experiments we set k to 1. We want to investigate the effect of changing the parameter k on the overall effectiveness of our solution. We plot the effect of varying k ($k = 1, 5, 10, 20$) on F measure for OpenOffice, Mozilla, and Eclipse datasets in Figures 3, 4, & 5 respectively. When we increase k , we consider more nearest neighbors. This might increase accuracy as in effect we are tapping more to the “wisdom of the masses”. However, this might also reduce accuracy as the additional neighbors might not be that similar anymore to the target bug report.

From the figures, for OpenOffice, the F measure of `critical` increases as we increase k . However, the F measures of `minor` and `trivial` decrease as we increase k . For Mozilla, the F measure of `major` slightly increases as we increase k , however, for three severity labels, `blocker`, `minor`, and `trivial`, their F measures decrease as we increase k . For Eclipse, the F measure of `critical` increases as we increase k . However, the F measures of three severity labels, `blocker`, `minor` and `trivial` decrease

Table VI
PRECISION, RECALL, AND F MEASURE FOR INSPECT AND SEVERIS [OFFLINE] ON MOZILLA

Severity	INSpect			Severis [Offline]		
	Precision	Recall	F Measure	Precision	Recall	F Measure
blocker	33.9%	31.3%	32.6%	100%	0.2%	0.4%
critical	64.0%	67.8%	65.9%	82.6%	53.7%	65.1%
major	53.5%	55.2%	54.3%	43.9%	93.1%	59.7%
minor	38.9%	33.4%	35.9%	50.5%	1.8%	3.4%
trivial	38.4%	33.6%	35.8%	19.7%	1.1%	2.2%

Table VII
PRECISION, RECALL, AND F MEASURE FOR INSPECT AND SEVERIS [OFFLINE] ON ECLIPSE

Severity	INSpect			Severis [Offline]		
	Precision	Recall	F Measure	Precision	Recall	F Measure
blocker	25.2%	27.0%	26.0%	0.0%	0.0%	0.0%
critical	28.2%	29.8%	29.0%	22.3%	39.7%	28.5%
major	58.0%	57.5%	57.8%	48.2%	66.8%	56.0%
minor	42.4%	38.4%	40.3%	7.6%	0.1%	0.2%
trivial	28.2%	25.0%	26.5%	0.0%	0.0%	0.0%

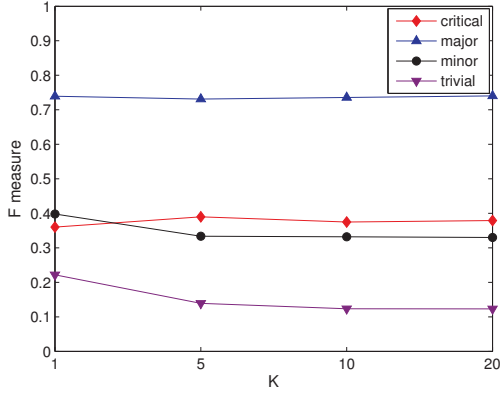


Figure 3. OpenOffice: Varying k and Its Effect on F Measure

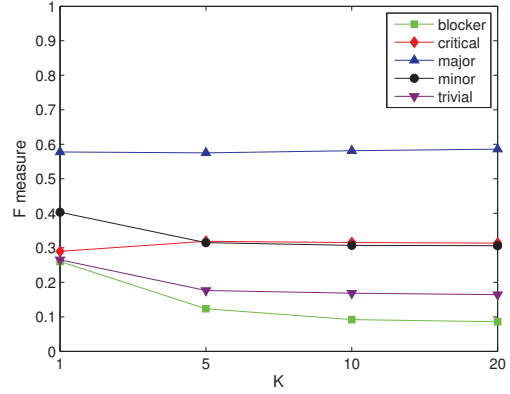


Figure 5. Eclipse: Varying k and Its Effect on F Measure

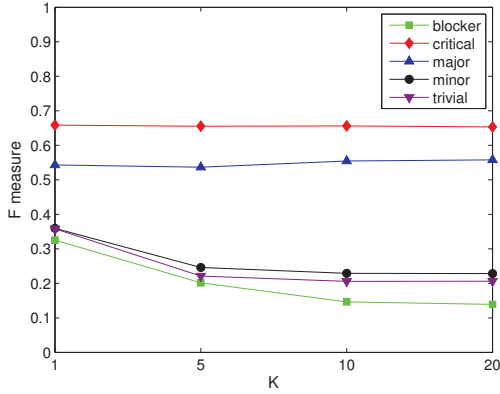


Figure 4. Mozilla: Varying k and Its Effect on F Measure

as we increase k .

F. Threats to Validity & Discussion

We consider three threats of validity: threats to construct validity, threats of internal validity, and threats of external validity.

Threats to construct validity relates to the suitability of

our evaluation metrics. We use standard metrics used in classification and prediction namely: precision, recall, and F measure. These measures have been used before by Menzies and Marcus to evaluate Severis [18].

Threats of internal validity refers to errors in our experiments. We extract the severity labels from the various Bugzilla bug tracking systems. We assume that except the normal label, the severity labels recorded in Bugzilla are the final severity labels that are deemed correct. We use these ground truth labels to measure how good our predictions are. A similar assumption and experimental setting were also made in prior studies [14], [15].

Threats of external validity refers to the generalizability of our findings. We consider repositories of three medium-large software systems: Eclipse, OpenOffice, and Mozilla. We consider a total of more than 65,000 bug reports. This is larger than the number of bug reports considered in prior studies [18], [14], [15]. Furthermore, the three projects are written in different programming languages, and have different background and user groups. Also, all our studies make use of open source repositories where data is publicly

available. We do not use the datasets from NASA used in [18] and made available in the Promise repository as they do not have information on duplicate bug reports. Note that duplicate bug reports are common phenomenon, and many of them are found in many open source Bugzilla tracking systems, c.f., [22], [31], [13], [27], [26].

V. RELATED WORK

In this section, we highlight related studies on bug severity prediction, bug report analysis, and text mining for software engineering.

A. Past Studies on Bug Severity Prediction

There are a number of studies that predict the severity of bug reports [15], [14], [18]. We highlight these studies in the following paragraphs.

Menzies and Marcus predict the severity of bug reports in NASA [18]. They first extract word tokens from bug reports, and then perform stop word removal and stemming. Important tokens are then identified using the concept of term frequency-inverse document frequency, and information gain. These tokens are then used as features for a classification approach named Ripper rule learner [6]. Their approach is able to identify fine grained bug report labels, which are the the 5 severity levels used in NASA.

More recently, Lamkanfi et al. predict the severity of bug reports from various projects' Bugzilla bug tracking systems [14]. They first extract word tokens and pre-process them. These tokens are then fed to a Naive Bayes classifier to predict the severity of the corresponding bug. Different from the work by Menzies and Marcus, they predict coarse grained bug severity labels: severe, and non-severe. Three of the six classes of severity in Bugzilla (`blocker`, `critical`, and `major`) are grouped as severe, two of the six classes (`minor`, and `trivial`) are grouped as non-severe, and `normal` severity bugs are omitted from their analysis.

Extending the above work, Lamkanfi et al. also try out various classification algorithms to predict the severity of bug reports [15]. They show that Naive Bayes performs better than other mining approaches on a dataset of 29,204 bug reports.

Our approach extends the above research studies. Similar to Menzies and Marcus's work, we detect fine grained bug report labels. Similar to the work by Lamkanfi et al. we consider bug reports on Bugzilla repositories of various open source projects. We compare our approach with that of Menzies and Marcus on a dataset containing more than 65,000 bug reports and show that we could gain significant F measure improvements.

B. Other Studies Analyzing Bug Reports

In a related research area, recently a number of techniques are proposed for duplicate bug report retrieval [22], [31],

[13], [27], [26]. Many of these approaches propose various ways to measure the similarity of bug reports to help developers in assigning bug reports as either duplicate or not. Runeson et al. propose a formula that considers the frequency of common words appearing in both documents as a similarity measure [22]. Wang et al. use both term frequency and inverse document frequency as a similarity measure [31]. They also consider a special situation where runtime traces are available and could be used to compute the similarity between bug reports. In practice, however, only a small minority of bug reports come with runtime traces. Jalbert and Weimer propose yet another term frequency based similarity measure [13]. Sun et al. propose a technique that leverages SVM for duplicate bug report detection [27]. In their later work, Sun et al. propose an approach to measuring the similarity of bug reports using an enhanced BM25F document similarity measure [26]. These recent advances in bug report similarity measurement could potentially be leveraged to categorize bug reports into various severity classes. Our work shows that they are indeed useful for this purpose.

Another line of research is categorization of bug reports to reduce maintenance effort. Anvik et al. [2], Cubranic and Murphy [7], Tamrawi et al. [28] propose various techniques to automatically assign the right developer for a new report. Huang et al. categorize bug reports into those related to capability, security, performance, reliability, requirement, and usability [12]. Pordguski et al. [20] and Francis et al. [8] propose approaches to group reported software failures, by analyzing the corresponding execution traces. Gegick et al. identify security bug reports using text mining [9]. The approach to some extent is similar to the work of Lamkanfi et al. [14], in that it categorizes bugs into two categories. However rather than categorizing bug into: severe and non-severe, it categorizes bug into: security-related and non-security-related.

Previous work also conducts empirical studies on bug repositories. Anvik et al. study the characteristics of bug repositories and show interesting findings on the number of reports that a person submit and the proportion of various resolutions [3]. Sandusky et al. study the nature, impact and extent of a bug report network in one large open source development community[23]. Hooimeijer and Weimer predict the quality of bug reports by a novel descriptive model built based on surface features of over 27,000 bug reports from several open source projects [11]. Bettenburg et al. describe characteristics of good bug reports by surveying Eclipse, Mozilla and Apache developers [4].

C. Text Mining for Software Engineering

There are many studies that utilize various forms of text analysis and mining for software engineering purposes. Haiduc et al. propose a method to summarize source code to support program comprehension [10]. The work proposes

an approach to extract informative yet succinct text to characterize source code entities so that developers can better understand a large piece of code. Sridhara et al. propose an approach to detect code fragments implementing high level abstractions and describe them in succinct textual descriptions [24].

Marcus and Maletic propose an approach to link documentation to source code using Latent Semantic Indexing [17]. Chen et al. proposed an approach to link textual documents to source code by combining several techniques including regular expression, key phrases, clustering and vector space model [5].

Similar to the above studies, we also extend a text mining approach to solve problem in software engineering. Different from the above studies, we investigate a different problem namely the prediction of fine-grained bug report severity label from its text. Our approach combines nearest neighbor classification and an extension of BM25 document similarity function.

VI. CONCLUSION AND FUTURE WORK

Severity labels are important for developers to prioritize bugs. A number of existing approaches have been proposed to infer these labels from textual fields of bug reports. In this work, we propose a new approach to infer severity labels from various information available from bug reports: textual, and non-textual. We make use of duplicate bug reports to weigh the relative importance of each piece of information or features to determine the similarity between bug reports. This similarity measure is then used in a nearest-neighbor fashion to assign a severity label to a bug report. We have compared our approach to the state-of-the-art approach on fine-grained severity label prediction, namely Severis. Extensive experiments on tens of thousands of bug reports taken from three medium-large software systems: Eclipse, OpenOffice, and Mozilla, have been performed. The result shows that we can improve the F measure of the state-of-the-art approach significantly, especially on hard-to-predict severity labels.

As future work, we plan to improve the accuracy of the proposed approach further. We also plan to embed our solution into Bugzilla to let it be used by many people.

REFERENCES

- [1] http://wiki.eclipse.org/WTP/Conventions_of_bug_priority_and_severity#How_to_set_Severity_and_Priority.
- [2] J. Anvik, L. Hiew, and G. Murphy, "Who should fix this bug?" in *Proceedings of the International Conference on Software Engineering*, 2006.
- [3] J. Anvik, L. Hiew, and G. C. Murphy, "Coping with an open bug repository," in *ETX*, 2005, pp. 35–39.
- [4] N. Bettenburg, S. Just, A. Schröter, C. Weiss, R. Premraj, and T. Zimmermann, "What makes a good bug report?" in *SIGSOFT FSE*, 2008, pp. 308–318.
- [5] X. Chen and J. C. Grundy, "Improving automated documentation to code traceability by combining retrieval techniques," in *ASE*, 2011, pp. 223–232.
- [6] W. Cohen, "Fast effective rule induction," in *ICML*, 1995.
- [7] D. Cubranic and G. C. Murphy, "Automatic Bug Triage Using Text Categorization," in *SEKE*, 2004, pp. 92–97.
- [8] P. Francis, D. Leon, and M. Minch, "Tree-based methods for classifying software failures," in *ISSRE*, 2004.
- [9] M. Gegick, P. Rotella, and T. Xie, "Identifying security bug reports via text mining: An industrial case study," in *MSR*, 2010, pp. 11–20.
- [10] S. Haiduc, J. Aponte, and A. Marcus, "Supporting program comprehension with source code summarization," in *ICSE (2)*, 2010, pp. 223–226.
- [11] P. Hooimeijer and W. Weimer, "Modeling bug report quality," in *ASE*, 2007, pp. 34–43.
- [12] L. Huang, V. Ng, I. Persing, R. Geng, X. Bai, and J. Tian, "AutoODC: Automated generation of orthogonal defect classifications," in *ASE*, 2011.
- [13] N. Jalbert and W. Weimer, "Automated duplicate detection for bug tracking systems," in *DSN*, 2008.
- [14] A. Lamkanfi, S. Demeyer, E. Giger, and B. Goethals, "Predicting the severity of a reported bug," in *MSR*, 2010.
- [15] A. Lamkanfi, S. Demeyer, Q. Soetens, and T. Verdonck, "Comparing mining algorithms for predicting the severity of a reported bug," in *CSMR*, 2011.
- [16] C. D. Manning, P. Raghavan, and H. Schütze, *Introduction to Information Retrieval*. New York, NY, USA: Cambridge University Press, 2008, pp. 232–233.
- [17] A. Marcus and J. I. Maletic, "Recovering documentation-to-source-code traceability links using latent semantic indexing," in *ICSE*, 2003, pp. 125–137.
- [18] T. Menzies and A. Marcus, "Automated severity assessment of software defect reports," in *ICSM*, 2008.
- [19] www.openoffice.org/qa/ooQAReloaded/Docs/QA-Reloaded-ITguide.html#priorities.
- [20] A. Podgurski, D. Leon, P. Francis, W. Masri, M. Minch, J. Sun, and B. Wang, "Automated support for classifying software failure reports," in *Proceedings of the 25th International Conference on Software Engineering*, 2003, pp. 465–475.
- [21] S. Robertson, H. Zaragoza, and M. Taylor, "Simple BM25 Extension to Multiple Weighted Fields," in *Proceedings of the thirteenth ACM international conference on Information and knowledge management*, 2004, pp. 42–49.
- [22] P. Runeson, M. Alexandersson, and O. Nyholm, "Detection of duplicate defect reports using natural language processing," in *ICSE*, 2007, pp. 499–510.
- [23] R. J. Sandusky, L. Gasser, and G. Ripoché, "Bug report networks: Varieties, strategies, and impacts in a f/oss development community," in *International Workshop on Mining Software Repositories*, 2004, pp. 80–84.
- [24] G. Sridhara, L. L. Pollock, and K. Vijay-Shanker, "Automatically detecting and describing high level actions within methods," in *ICSE*, 2011, pp. 101–110.
- [25] www.ils.unc.edu/~keyeg/java/porter/PorterStemmer.java.
- [26] C. Sun, D. Lo, S.-C. Khoo, and J. Jiang, "Towards more accurate retrieval of duplicate bug reports," in *ASE*, 2011.
- [27] C. Sun, D. Lo, X. Wang, J. Jiang, and S.-C. Khoo, "A discriminative model approach for accurate duplicate bug report retrieval," in *ICSE*, 2010.
- [28] A. Tamrawi, T. T. Nguyen, J. Al-Kofahi, and T. N. Nguyen, "Fuzzy set-based automatic bug triaging," in *ICSE*, 2011, pp. 884–887.
- [29] A. Tamrawi, T. T. Nguyen, J. M. Al-Kofahi, and T. N. Nguyen, "Fuzzy set-based automatic bug triaging," in *ICSE*, 2011.
- [30] M. Taylor, H. Zaragoza, N. Craswell, S. Robertson, and C. Burges, "Optimisation methods for ranking functions with multiple parameters," in *Int. Conf. on Information and Knowledge Mgmt (CIKM)*, 2006.
- [31] X. Wang, L. Zhang, T. Xie, J. Anvik, and J. Sun, "An approach to detecting duplicate bug reports using natural language and execution information," in *ICSE*, 2008, pp. 461–470.
- [32] H. Zaragoza, N. Craswell, M. J. Taylor, S. Saria, and S. E. Robertson, "Microsoft cambridge at trec 13: Web and hard tracks," in *TREC*, 2004.