



BUG REPOSITORY ANALYSIS: A STUDY

Km Rupal Singh¹, Dayashankar Singh²

¹PG Scholar, Deptt. Of Comp. Sc. & Engg.

M.M.M.Engineering College, Gorakhpur (UP), (India)

²Associate Professor, Deptt. Of Comp. Sc. & Engg.

M.M.M.Engineering College, Gorakhpur (UP), (India)

ABSTRACT

In a bug vault at whatever point another bug is accounted for it is critical to consider this new bug report keeping in mind the end goal to recognize it as copy or none copy. For this reason, a lot of work has been done here. In this paper, we experience the beforehand proposed methods and break down the commitment of each work in making the copy recognition in bug store such a vital field of research.

Keywords: “Bug”, “Bug report”, “Bug Triage”, “Life cycle”.

I.INTRODUCTION

The Software programs are vast and complex, in this way keeping in mind the end goal to keep up their quality, support of programming is required. Some extensive programming ventures are being kept up by their bug reports to remedy their mistakes. Bugs are put together by various programming groups i.e. Advancement group, Testing group and end clients, hence bringing about an extensive database of the bug report. Since the bug reports put together by programming groups of engineers, analyzers and end clients so for the same bug many copy bug reports are submitted. For this reason, need to give loads of time to physically recognize whether an approaching report is a copy or not. This has been accounted for in 2005 that for Mozilla "consistently just about 350 bugs create the impression those needs triaging. This is to an extreme degree a lot for just the Mozilla software engineers to deal with" [1], this clarifies the need for a mechanized framework for copy identification in a bug store. In one of the methodologies [2], the copies were viewed as helpful as they add data to the same prior bug report. In different methodologies, the copies were dealt with distinctively [3-4]. We would, along these lines, clarify the bug, its life cycle, and the different fields of a bug report and the different methodologies for copy discovery in the ensuing segments.

II.BACKGROUND

Preparatory

In this area, we clarify some fundamental ideas and portray the life-cycle of bug reports. We likewise give a few insights and also the examination of bug reports.



2.1 Bug-report phrasing

(1) **Bug.** A product bug is a mistake, defect, or blame in a PC program or framework that produces startling outcomes or conduct. There is a qualification amongst "bug" and "issue". An issue could be a bug however not generally a bug. It can likewise be an element request), assignment, missing documentation, et cetera. The way toward finding and evacuating bugs is called "troubleshooting". Bugs might be caused by little coding blunders, however, the consequences of bugs can be not kidding, and making finding and settling bugs a somewhat difficult undertaking.

(2) **Bug report.** A bug report is a product record depicting programming bugs, which is presented by an engineer, an analyzer, or an end-client. Bug reports have numerous different names, for example, "deformity reports", "blame reports", "disappointment reports", "blunder reports", "issue reports", "inconvenience reports", et cetera. Regularly a bug report is made out of its recognizable proof number, its title, when it is accounted for and changed, the seriousness or significance, the software engineer appointed to settle the bug, the determination status (e.g., new, unverified, settled) of the bug, the portrayal of the report (e.g., ventures to imitate the bug, stack follows, and expected conduct), extra remarks (e.g., discourse about the conceivable arrangements), connections (e.g., proposed patches, test cases), and a rundown of reports that should be tended to sometime recently this report is settled.

(3) **Bug storehouse.** Bug archives are regularly utilized as a part of open-source programming activities to permit both engineers and clients to post issues experienced with the product, propose conceivable improvements and remark after existing bug reports. An open bug archive is open and unmistakable for all individuals. As bugs revealed in bug storehouses might be distinguished and explained by all individuals, the nature of the open ventures may enhance.

(4) **Bug-report triage.** Bug-report triage comprises of the accompanying procedure: making sense of whether a bug is a genuine bug, checking whether the revealed bug is a copy of a current bug, organizing bug reports, furthermore, choosing which engineer should take a shot at the bug reports.

(5) **Bug-report duplication.** Copy bug reports allude to the bug gives an account of a similar bug submitted by various columnists, as there are numerous clients communicating with a framework and detailing its bugs. Identifying copy bug reports is a technique of bug triage, which diminishes triaging expense and spares time for designers in settling similar issues.

(6) **Bug following framework.** A bug following framework (additionally called imperfection following framework) oversees bug reports and designers who settle bugs. Bug following frameworks is intended to track detailed programming bugs. Bugs are put away in a bug archive, which is the significant segment of a bug following framework. To submit and resolve bugs naturally, engineers of numerous prominent open-source ventures (e.g., Mozilla, Shroud, and Linux bit) utilize bug following frameworks (e.g., Bugzilla, Jira, Mantis, Trac).

2.2 The life-cycle of a bug

Figure 1 demonstrates the life-cycle of a bug report.

1. **New:** When a tester find a bug and posting it very first time then the status of the defect is "NEW".
2. **Open:** After an analyzer has posted a bug, the lead of the analyzer endorses that the bug is honest to goodness and he changes the state as "OPEN".
3. **Allot:** Once the lead changes the state as "OPEN", he relegates the bug to comparing designer or engineer group. The condition of the bug now is changed to "ALLOT".
4. **Test:** Once the engineer settles the bug, he needs to allocate the bug to the testing group for next round of testing. Before he discharges the product with bug settled, he changes the condition of bug to "TEST". It indicates that the bug has been settled and is discharged to testing group.
5. **Conceded:** The bug, changed to conceded state implies the bug is relied upon to be settled in next discharges. The purposes of changing the bug to this state have many components. Some of them are need of the bug might be low; the absence of time for the discharge or the bug might not have a real impact on the product.
6. **Rejected:** If the designer feels that the bug isn't honest to goodness, he rejects the bug. At that point, the condition of the bug is changed to "REJECTED".
7. **Confirmed:** Once the bug is settled and the status is changed to "TEST", the analyzer tests the bug. On the off chance that the bug is absent in the product, he affirms that the bug is settled and changes the status to "CONFIRMED".

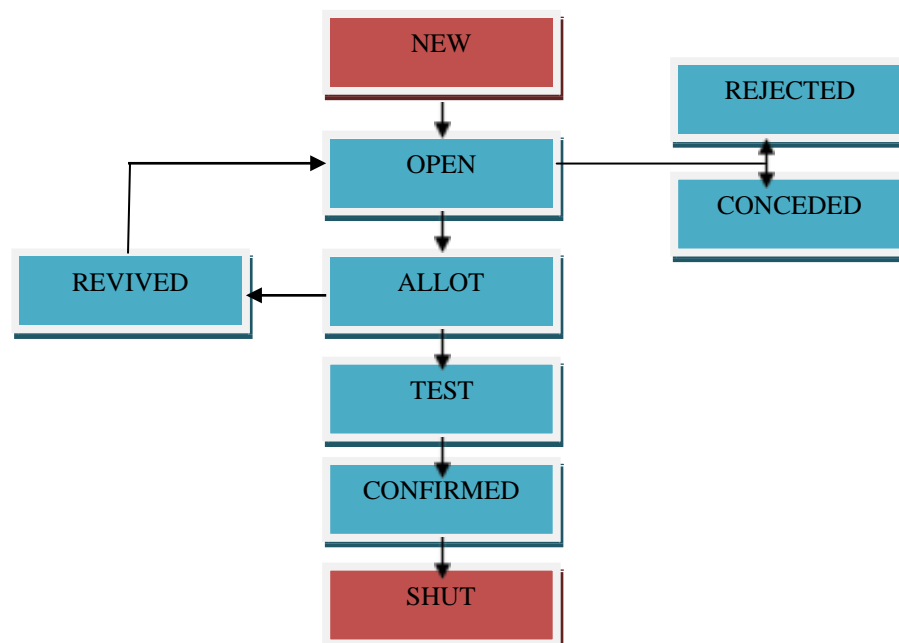


Figure 2: The life-cycle of a bug



8. **Revived:** If the bug still exists even after the bug is settled by the engineer, the analyzer changes the status to "REVIVED". The bug navigates the life cycle by and by.
9. **Shut:** Once the bug is settled, it is tried by the analyzer. On the off chance that the analyzer feels that the bug never again exists in the product, he changes the status of the bug to "SHUT". This state implies that the bug is settled, tried and affirmed.

III.LITERATURE SURVEY

C.Sun, D.Lo, S.C.Khoo and J.Jiang, [13] utilized a bug following framework, diverse analyzers or clients may present different reports on similar bugs, alluded to as copies, which may cost additional support endeavors in triaging and settling bugs. Keeping in mind the end goal to distinguish such copies precisely, in this paper propose a recovery work (REP) to gauge the likeness between two bug reports. It completely uses the data accessible in a bug report including not just the likeness of printed content in synopsis and depiction fields yet additionally the comparability of no textual fields, for example, item, segment, adaptation, and so on. The downsides of that framework are there is no ordering structure of bug report store to accelerate the recovery procedure.

J.Xuan, H.Jiang, Z.Ren, J.Yan, and Z.Luo [14] propose a semi-regulated content arrangement approach for bug triage to keep away from the inadequacy of named bug reports in existing administered approaches. This new approach joins Naive Bayes classifier and desire augmentation to exploit both marked and unlabeled bug reports. This approach prepares a classifier with a few marked bug reports. At that point, the approach iteratively names various unlabeled bug reports and prepares another classifier with names of all the bug reports. In that weighted proposal list for the semi-directed approach give a weighted suggestion rundown to increasing the semi-regulated approach utilizing probabilistic marks of unlabeled bug reports. In light of this weighted suggestion list, enhance the order exactness for the semi-directed approach. The disadvantages of these approach for programmed bug triage with a bug archive there is no bug triage module part consolidating with the bug store in true applications.

T. M. Khoshgoftaar, K.Gao, and N. Seliya [15] reason characteristic choice and imbalanced information: Problems in programming deformity forecast. To deal with imbalanced deformity information.

P. S. Bishnu and V. Bhattacharjee [16] reason programming flaw expectation utilizing quad tree-based k-implies bunching calculation. In that paper procedure the deformity information with quad tree based K-implies bunching to help surrender forecast. In that product, measurements anticipate an incentive for singular programming curio (e.g. Source code record, a class or module). The product ancient rarity contains blame as per the separated highlights of the curio.

S.Shivaji, E.J.Whitehead, Jr. R. Akella, and S.Kim [17] reason decreasing highlights to enhance code change based bug forecast. In that paper a structure to inspect numerous component determination calculations and expel clamor highlight in grouping based imperfection expectation. It doesn't contain how to quantify the commotion protection in imperfection expectation and how to surrender clamor information.

Dang et al. [18] made a model that spots more weight on stack outlines nearer to the highest point of the stack and supports stacks whose coordinated capacities are also dispersed from each other. This system experiences a

proposed $O(n^3)$ grouping calculation.

In 2005, Brodie et al. [19] exhibited an approach that standardizes the call stack to evacuate non-discriminative capacities and additionally smoothing recursive capacities and looks at stacks utilizing weighted alter remove.

Schröter et al. [20] observationally contemplated engineers' utilization of stack follows in investigating and found that bugs will probably be settled in the best 10 edges of their particular crash stack follow, additionally affirming the amazing hugeness of the best k stack outlines in crash report bucketing, which is likewise validated all the more as of late by Wu et al. [21].

Bartz et al. [22] additionally utilized alter remove on the stack follow, yet a weighted variation with weights gained from preparing information. Subsequently, they could consider other information in the crash report beside the stack follow. The weights learned proposed some intriguing discoveries: substituting a module in a call stack brought about a significantly higher separation; also, the call stack alter remove was observed to be the most noteworthy weighted factor, notwithstanding the thought of other crash report information, affirming the instinct in the writing of the stack follow's significance.

IV. APPROACHES

There are three fundamental strides in the framework preprocessing, preparing a discriminative model and recovering copy bug reports in Figure 2.

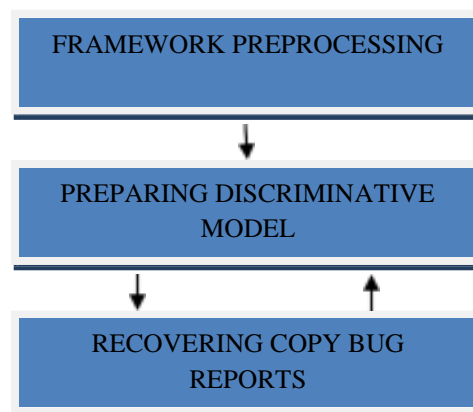


Figure 2: Steps for classifying the Bug Reports

4.1 Step 1: Framework Preprocessing

Information pre-preparing is the underlying advance towards copy identification utilizing Natural Language handling strategies. In this stage, the information is prepared before its utilization. For this reason, following advances are taken:

1) Tokenization: in tokenization, the record is parsed and isolated into tokens by expelling the delimiters (spaces, accentuation marks).

2) **Stemming:** In stemming the words are decreased to their ground shapes, for instance, a stemmer will lessen „runs“ , „running“ to its ground frame „run“ .

3) **Stop word evacuation:** In this progression, the prevent words are expelled from the information by dispensing with words that are of minimum significance, such words are „the“ , „and“ , „is“ and significantly more.

4.2 Step 2: Preparing Discriminative Model

The rundown and depiction of a bug report are changed over into their equal weight vectors with the assistance of TF-IDF; it is the most well known capacity for measuring the words. TF remains for Text Frequency. It compares to the quantifier of times a term happens in the record, as in (1).

$$TF f(a, b) = 0.5 + \frac{0.5 * f(a, b)}{\sum_{b \in D} f(a, b)} \quad (1)$$

$$Max \{f(w, b): w \in b\}$$

IDF remains for Inverse Document Frequency. It is the logical measure of the significance of a word. It chips away at the suspicion that critical words will happen every now and again in some archive and occasionally over the whole corpus, this can be found in (2), which is the standard equation for IDF.

$$IDF(a, b) = \frac{1}{\log |b|} \quad (2)$$

$$| \{b: D: a \in b\} |$$

And after that the TF-IDF weighting factor is computed as (3), this is a standout amongst the most mainstream words measuring factor which shapes word-vector out of the content.

$$TF-IDF(a, b, D) = TF(a, b) * IDF(a, D) \quad (3)$$

4.2.1 Clustering Algorithm

A. K-Means Algorithm

K-means clustering is a kind of unsupervised learning, which is utilized when you have unlabeled information (i.e., information without characterized classifications or gatherings) in Figure 3. The objective of this algorithm is to discover bunches in the information, with the quantity gatherings, spoke to by the variable K. The algorithm works iteratively to allocate every datum point to one of K cluster in view of the highlights that are given. Information focuses are the cluster in light of highlight closeness. The aftereffects of the K-means clustering algorithm are: The centroids of the K groups, which can be utilized to mark new information and Names for the preparation information (every datum point is doled out to a solitary group).

Instead of characterizing bunches before taking a gander at the information, cluster enables you to discover and break down the gatherings that have framed naturally. . Every centroid of a cluster is a gathering of highlight

esteems which characterize the subsequent gatherings. Looking at the centroid include weights can be utilized to subjectively decipher what sort of gathering each cluster speaks to.

B. Hierarchical Clustering Algorithm

Hierarchical clustering includes making the cluster that has a foreordained requesting start to finish. For instance, all documents and envelopes on the hard circle are sorted out in a progressive system. There are two sorts of various leveled grouping, Divisive and Agglomerative

Divisive Method

In this technique, we relegate the majority of the perceptions to a solitary group and after that segment the cluster into two slightest comparative clusters. At long last, we continue recursively on each group until there is one cluster for every perception.

Agglomerative Method

In this strategy, we relegate every perception to its own particular cluster. At that point, process the similitude (e.g., remove) between each of the clusters and join the two most comparable groups.

4.3 Step 3: Recovering Copy Bug Reports

Order to recover material bug reports from the storehouse.

V.CONCLUSION

The reason for this study is to give a diagram of the procedure that is taken after from the scratch towards the recovery of copies in bug reports. There can be approaches, so this paper principally centered on the essential advances that are followed all in all for copy recognition and gave a brief of how real things are being worked out.

REFERENCES

- [1] Anvik, J., Hiew, L., & Murphy, G. C. (2005, October). Coping with an open bug repository. In *Proceedings of the 2005 OOPSLA workshop on Eclipse technology exchange* (pp. 35-39). ACM.
- [2] Bettenburg, N., Premraj, R., Zimmermann, T., & Kim, S. (2008, September). Duplicate bug reports considered harmful... really? In *Software maintenance, 2008. ICSM 2008. IEEE international conference on* (pp. 337-345). IEEE.
- [3] Sun, C., Lo, D., Wang, X., Jiang, J., & Khoo, S. C. (2010, May). A discriminative model approach for accurate duplicate bug report retrieval. In *Proceedings of the 32nd ACM/IEEE International Conference on Software Engineering-Volume 1* (pp. 45-54). ACM.
- [4] Tian, Y., Sun, C., & Lo, D. (2012, March). Improved duplicate bug report identification. In *Software Maintenance and Reengineering (CSMR), 2012 16th European Conference on* (pp. 385-390). IEEE.
- [5] Tamrawi, A., Nguyen, T. T., Al-Kofahi, J., & Nguyen, T. N. (2011, May). Fuzzy set-based automatic bug



- triaging: NIER track. In *Software Engineering (ICSE), 2011 33rd International Conference on* (pp. 884-887). IEEE.
- [6] Nguyen, A. T., Nguyen, T. T., Al-Kofahi, J., Nguyen, H. V., & Nguyen, T. N. (2011, November). A topic-based approach for narrowing the search space of buggy files from a bug report. In *Automated Software Engineering (ASE), 2011 26th IEEE/ACM International Conference on* (pp. 263-272). IEEE.
- [7] Menzies, T., & Marcus, A. (2008, September). Automated severity assessment of software defect reports. In *Software Maintenance, 2008. ICSM 2008. IEEE International Conference on* (pp. 346-355). IEEE.
- [8] Podgurski, A., Leon, D., Francis, P., Masri, W., Minch, M., Sun, J., & Wang, B. (2003, May). Automated support for classifying software failure reports. In *Software Engineering, 2003. Proceedings. 25th International Conference on* (pp. 465-475). IEEE.
- [9] Anvik, J. (2006, May). Automating bug report assignment. In *Proceedings of the 28th international conference on Software engineering* (pp. 937-940). ACM.
- [10] Raymond, E. (1999). The cathedral and the bazaar. *Philosophy & Technology*, 12(3), 23.
- [11] Murphy, G., & Cubranic, D. (2004). Automatic bug triage using text categorization. In *Proceedings of the Sixteenth International Conference on Software Engineering & Knowledge Engineering*.
- [12] Nguyen, A. T., Nguyen, T. T., Nguyen, T. N., Lo, D., & Sun, C. (2012, September). Duplicate bug report detection with a combination of information retrieval and topic modeling. In *Proceedings of the 27th IEEE/ACM International Conference on Automated Software Engineering* (pp. 70-79). ACM.
- [13] Sun, C., Lo, D., Khoo, S. C., & Jiang, J. (2011, November). Towards more accurate retrieval of duplicate bug reports. In *Automated Software Engineering (ASE), 2011 26th IEEE/ACM International Conference on* (pp. 253- 262). IEEE.
- [14] Xuan, J., Jiang, H., Ren, Z., & Zou, W. (2012, June). Developer prioritization in bug repositories. In *Software Engineering (ICSE), 2012 34th International Conference on* (pp. 25-35). IEEE.
- [15] Khoshgoftaar, T. M., Gao, K., & Seliya, N. (2010, October). Attribute selection and imbalanced data: Problems in software defect prediction. In *Tools with Artificial Intelligence (ICTAI), 2010 22nd IEEE International Conference on* (Vol. 1, pp. 137-144). IEEE.
- [16] Bishnu, P. S., & Bhattacharjee, V. (2012). Software fault prediction using quad tree-based k-means clustering algorithm. *IEEE Transactions on knowledge and data engineering*, 24(6), 1146-1150.
- [17] Shivaji, S., Whitehead, E. J., Akella, R., & Kim, S. (2013). Reducing features to improve code change-based bug prediction. *IEEE Transactions on Software Engineering*, 39(4), 552-569.
- [18] Dang, Y., Wu, R., Zhang, H., Zhang, D., & Nobel, P. (2012, June). ReBucket: a method for clustering duplicate crash reports based on call stack similarity. In *Proceedings of the 34th International Conference on Software Engineering* (pp. 1084-1093). IEEE Press.
- [19] Brodie, M., Ma, S., Lohman, G., Mignet, L., Wilding, M., Champlin, J., & Sohn, P. (2005, June). Quickly finding known software problems via automated symptom matching. In *Autonomic Computing, 2005. ICAC 2005. Proceedings. Second International Conference on* (pp. 101-110). IEEE.
- [20] Schroter, A., Schröter, A., Bettenburg, N., & Premraj, R. (2010, May). Do stack traces help developers fix



- bugs? In *Mining Software Repositories (MSR), 2010 7th IEEE Working Conference on* (pp. 118-121). IEEE.
- [21] Wu, R., Zhang, H., Cheung, S. C., & Kim, S. (2014, July). CrashLocator: locating crashing faults based on crash stacks. In *Proceedings of the 2014 International Symposium on Software Testing and Analysis* (pp. 204-214). ACM.
- [22] Bartz, K., Stokes, J. W., Platt, J. C., Kivett, R., Grant, D., Calinoiu, S., & Loihle, G. (2008, December). Finding Similar Failures Using Callstack Similarity. In *SysML*.