

Effective Bug Triage System

Renu Deorankar¹, Shubhangi Vairagar²

¹ME Student, Computer Engineering, University of Pune, Maharashtra

²Assistant Professor, Computer Engineering, University of Pune, Maharashtra

ABSTRACT

Software bugs are the major areas where most of the IT companies spend its 45% of efforts. The most efficient way of reducing bugs in the software is bug triage, which is appropriately assigning the new but to a developer. Automatic Bug triage is conducted by using text classification methods to minimize the manual work. In this paper, the problem of huge data for bug triage is addressed i.e. reducing the amount of data in assigning the bugs to developers. This technique combines feature selection and instance selection for simultaneously minimizing the bug dimension and word dimension. To identify the application order of instance selection and then feature selection, extraction of various attributes from previous bug track record is carried out to build new bug track or dataset. The proposed system tests its performance on the sample dataset of the open source software i.e. Eclipse. The proposed system works on novel data processing techniques to generate small amount of but highly informative bug report data in software maintenance and development

Keywords: Feature selection, application of data preprocessing, Instance selection, Mining software repositories, data management for bug repositories, bug data minimization, bug triage, prediction for reduction orders.

1. INTRODUCTION

Mining various software repositories in these days has become essential and hence employment of data mining tools for mining data for software maintenance is necessary [2]. In the current software development environment, the output of software development is stored into large software repositories for maintaining large scale bug data, e.g. codes, emails, specifications and bugs. Conventional way of software analysis is not convenient for complex and huge data [5]. Software data has been efficiently handled by the data mining techniques in current digital age. (e.g., [7]). In proposed approach for bug data reduction, we scale out the data by analyzing and categorizing the bug data which will indirectly increase the quality of the data. Proposed approach is a combinational process of feature selection along with instance selection simultaneously. Also the proposed system comprised of a novel module to determine the status of the bug checking the current state of the bug whether it's rectified or not or it's assigned to some developer or not.

Software bug management is carried out with an important entity names as bug repository. An open bug repository is available in many open source software, which not only allow users to report about the issues or defects found in the software, but also suggest future enhancements as per user's needs, and also comment or resolve the existing bug reports. The quantity of bug reports have exponentially expanded which is making bug triaging for prominent open source software [2]. Two noteworthy difficulties being confronted in software bug archives are gigantic measure of bug report information and the nature of bug report. A bug archive keeps up as a bug report in a literary depiction shape and is redesigned by status pecking order of the bug settling [1].

The proposed system makes use of feature selection and instance selection with along with historical data for minimizing the bug sizes in the bug repository in order to obtain high quality and small size informational data. The proposed system also will show the graphical analysis of the proposed system and existing system efficiency.

2. LITERATURE SURVEY

A. Towards Effective Bug Triage with Software Data Reduction Techniques

In this paper, we address the issue of information decrease for bug triage, i.e., how to lessen the scale and enhance the nature of bug information. We consolidate occasion choice with highlight choice to all the while lessen information scale on the bug measurement and the word measurement. To decide the request of applying example determination and highlight choice, we remove characteristics from verifiable bug information sets and assemble a prescient model for another bug information set. We exactly examine the execution of information lessening on absolutely 600,000 bug reports of two huge open source ventures, to be specific Eclipse and Mozilla. Our work gives a way to deal with utilizing systems on information handling to shape lessened and excellent bug information in software advancement and support.

B. Efficient Bug Triaging Using Text Mining

Vast open source software ventures get rich rates of submitted bug reports. Triaging these approaching reports physically is blunder inclined and tedious. The objective of bug triaging is to allocate conceivably experienced engineers to new-coming bug reports. To lessen time and cost of bug triaging, we exhibit a programmed way to deal with anticipate an engineer with significant experience to comprehend the new coming report. In this paper, we research the utilization of five term determination strategies on the precision of bug task. Furthermore, we re-adjust the load between engineers taking into account their experience. We direct analyses on four genuine datasets. The exploratory results demonstrate that by selecting a little number of segregating terms, the F-score can be essentially moved forward.

C. Supporting Bug Investigation using History Analysis

In my exploration, I propose a computerized strategy to bolster bug examination by utilizing a novel investigation of the historical backdrop of the source code. Amid the bug-altering process, engineers spend a high measure of manual exertion exploring the bug with a specific end goal to answer a progression of inquiries regarding it. My exploration will bolster engineers in noting the accompanying inquiries concerning a bug: Who is the most suitable designer to settle the bug?, Where is the bug found?, When was the bug embedded? furthermore, Why was the bug embedded?

D. Automated, Highly-Accurate, Bug Assignment Using Machine Learning and Tossing Graphs

The quantity of reported bugs in expansive open source undertakings is high and triaging these bugs is an essential issue in software support. As a stage in the bug triaging process, relegating another bug to the most suitable engineer to alter it, is not just a period devouring and dull errand. Triager who considers a bug and allocates it to a designer, additionally should know about engineer exercises at various parts of the undertaking. It is clear that just a couple of engineers have this capacity to complete this progression of bug triaging. The primary objective of this paper is to recommend another way to deal with the procedure of performing programmed bug task. The data expected to choose the best designers to settle another bug report is separated from the form control archive of the venture. Not at all like all the past proposed approaches which utilized Machine Learning and Information Retrieval strategies, this examination utilizes the Information Extraction (IE) techniques to remove the data from the product archives. The proposed approach does not utilize the data of the bug storehouse to settle on choices about bugs keeping in mind the end goal to get better results on tasks which don't have numerous altered bugs. The point of this examination is to prescribe the genuine fixers of the bugs. Utilizing this methodology, we accomplished 62%, 43% and 41% correctness on Eclipse, Mozilla and Gnome ventures, individually.

E. Formal models for expert finding in enterprise corpora

Hunting an association's report storehouses down specialists gives a practical answer for the assignment of master finding. We introduce two general techniques to master looking given a report accumulation which are formalized utilizing generative probabilistic models. The first of these straightforwardly models a specialist's learning in light of the reports that they are connected with, whilst the second finds archives on theme, and afterward finds the related master. Shaping solid affiliations is urgent to the execution of master discovering frameworks. Subsequently, in our assessment we look at the changed methodologies, investigating an assortment of relationship alongside other operational parameters, (for example, topicality). Utilizing the TREC Enterprise corpora, we demonstrate that the second technique reliably outflanks the first. A correlation against other unsupervised methods, uncovers that our second model conveys superb execution.

3. ARCHITECTURE

3.1 Bug Triage:

Vital role of bug triage is to appropriately assign a developer to the corresponding bugs. The developer will reopen or rectify the bug only if the bug is assigned to the developer. Contingent on the status upgraded by the [1]. Bug triage plans to allot a proper designer to alter another bug, i.e., to figure out who ought to settle a bug.

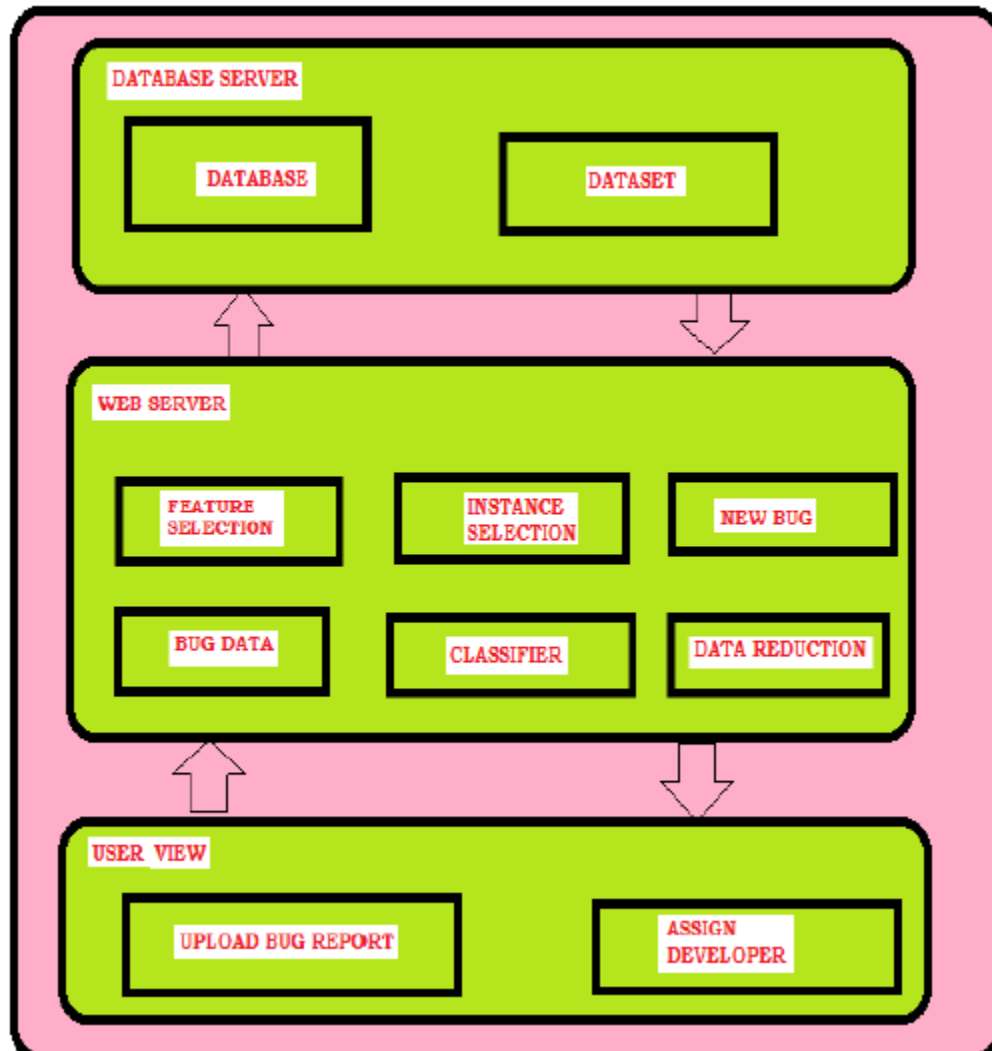


Fig -1: Architecture of proposed system

3.2 Data Reduction:

Here the feature selection and instance selection is applied to reduce the bug report data size so as to get low scale data but a higher quality data. In the proposed work, to relieve the work expense of designers, bug triage information diminishment has two basic points, 1) information scale decrease and 2) precision change of bug information. In contradict to printed information displaying like as in existing framework, we point in building the subset of the first dataset for pre-processing, material before a current bug triage system.



Fig -2: Data Reduction

3.3 Instance Selection

Instance selection technique is corresponding with mining of data tasks such as clustering and classification :

- It's an uncommon process of identifying valid, novel, and efficiently important, and finally understandable patterns in data mining. Selecting the small subset of the data from the entire data as if the whole data is being processed.
- The standard outcome of instance selection is independent of a model:

$$P(M_s) == P(M_w).$$

Evaluation measures:

Direct Measure: Maintaining the similarity between subset of data and original data. Ex) Entropy, histograms, moments.

Indirect Measure: For an example, Instance selection is better or worse in predictive accuracy is done on basis of a classifier. Traditional evaluation technique in sampling, clustering and classification, can be used in performance assessment of instance selection. Ex) Recall, Precision.

3.4 Feature Selection

- It select a minimum set of features such that the probability distribution of different classes given the values for those features is as close as possible to the original distribution given the values of all features [1].
- Avoid Reduce # of patterns in the patterns, easier to understand.
- Do create new attributes that can capture the important information in a data set much more efficiently than the original attributes.
- Use the smallest representation which is enough to solve the task.

Heuristic methods:

- Step-wise forward selection
- Step-wise backward elimination
- Create new attributes that can capture the important information in a data set much more efficiently than the original attributes:
- Three general methodologies:
 - 1.Feature extraction.
 - 2.Domain-specific.
 - 3.Mapping data to new space (see: data reduction)

3.5 Historical Data Module:

Historical data module is used to reduce the bug report data. This module track wise analyses the history of the bug and thereby accordingly reduce the size [3].

4. MATHEMATICAL MODEL

The mathematical formulation of the proposed system can be done on the basis of Deterministic Finite Automata (DFA) and Context Free Grammar (CFG) to represent the entities and the various transitions taking place between the entities.

The DFA for the proposed system comprised of 5 major components as mentioned below:

$$S = \{\Sigma, \lambda, \delta, \psi, F\}$$

Where,

Σ represents the entity is dataset.

λ represents the sequence of operations.

δ represents the operation or transitions being performed for bringing the project from one state to another.

Ψ is the collection of all states of proposed system.

F is the final state of the proposed system which depicts the output if the proposed system.

The proposed system states are as mentioned below:

q0- input accepting state

q1- dataset collection

q2- dataset loading

q3- add data into dataset

q4- get employer details

q5- data reduction

q6- bug report analysis

q7- bug rectification

F- trace bug history

The proposed system goes from one state to another state as follows:

q0→q1 where → dataset collection.

$$\delta_{q0 \rightarrow q1}(\text{dataset collection})$$

q1→q2 where → dataset loading.

$$\delta_{q1 \rightarrow q2}(\text{dataset loading})$$

q2→q3 where → add data into dataset.

$$\delta_{q2 \rightarrow q3}(\text{add data into dataset})$$

q3→q4 where → get employer details.

$$\delta_{q3 \rightarrow q4}(\text{get employer details})$$

q4→q5 where → data reduction.

$$\delta_{q4 \rightarrow q5}(\text{data reduction})$$

q5→q6 where → bug report analysis.

$\delta_{q5 \rightarrow q6}$ (bug report analysis)

q6→q7 where → bug rectification.

$\delta_{q6 \rightarrow q7}$ (bug rectification)

q7→F where → trace bug history.

$\delta_{q7 \rightarrow F}$ (trace bug history)

5. WORKING OF PROPOSED SYSTEM

The working of the proposed system is described as below:

- Dataset Collection and dataset loading into application of bug assignment and analysis system.
- User authentication on privilege level for bug assignment operation and bug resolution operation.
- Data Preprocessing and loaded data analyzing based on differential symbols.
- New bug or old bug reassigning facility for Bug triaging.
- Data Reduction and instance selection module for differentiating instances of bug status.
- Determining the Feature selection from the instance selection process.
- View various Bug reports as per bug status.
- System evaluation module for graphical analysis.

6. EXPERIMENTAL RESULTS

The experimental results of the proposed system shows that the input random large data for bug reports is analyses under combinational implementation of instance and feature selection thereby making the bug data smaller and highly efficient for bug triaging and assigning the new tasks as per the developer history of bug fixing. The instances of the bugs are differentiated based on the bug status after assigning the bugs to the developers. Unless the bugs are assigned to any developer any data regarding bug triaging cannot be analyses for future bug status tracking. Thus the proposed system aids in bug triaging based on historical data of bug reports and thereby make software bug triaging easier with reduction in size of bug repository and bug reports.

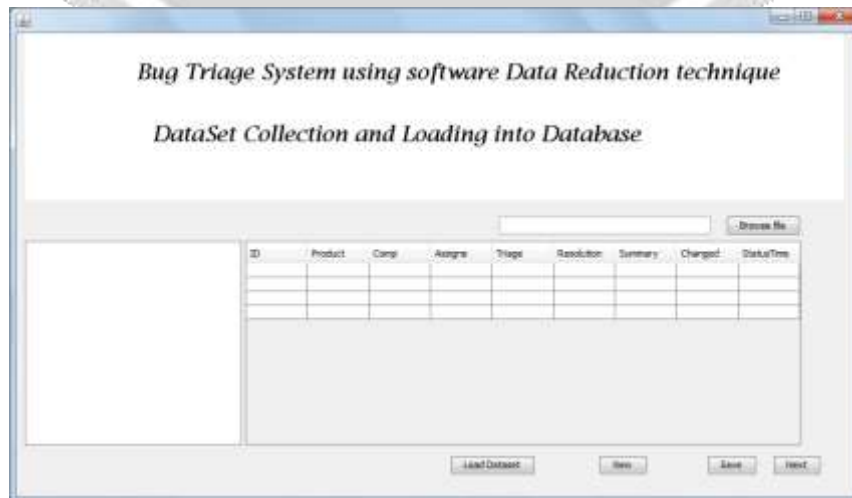


Fig. 6.1 Data loading Form

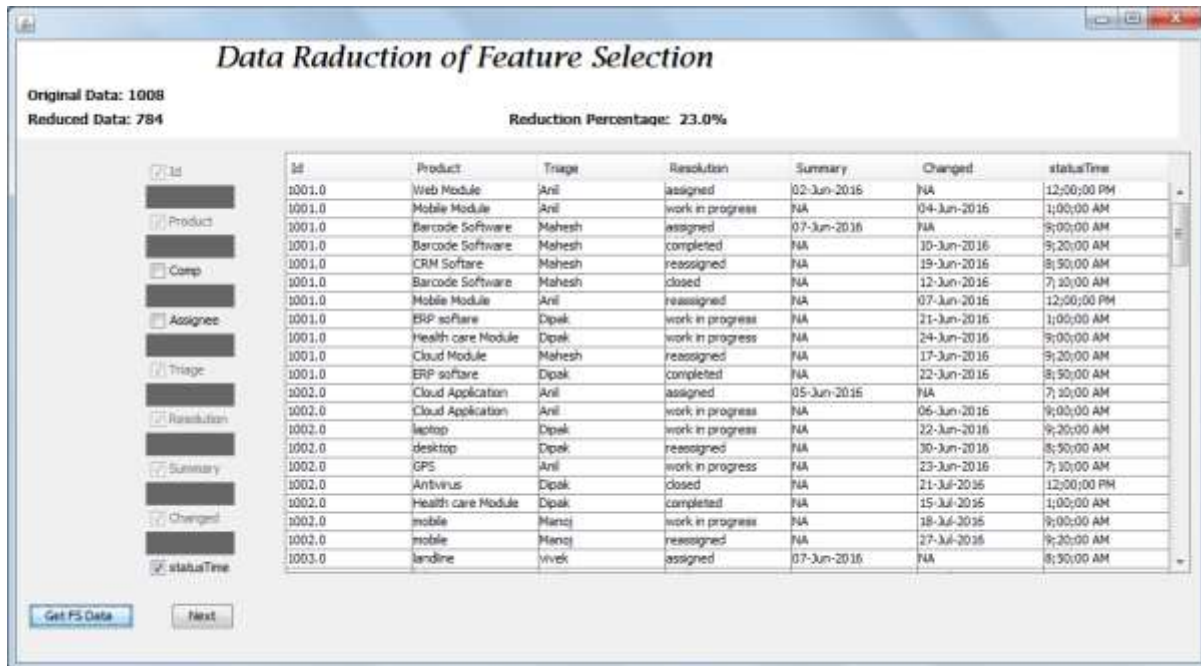


Fig. 6.2 Feature Selection

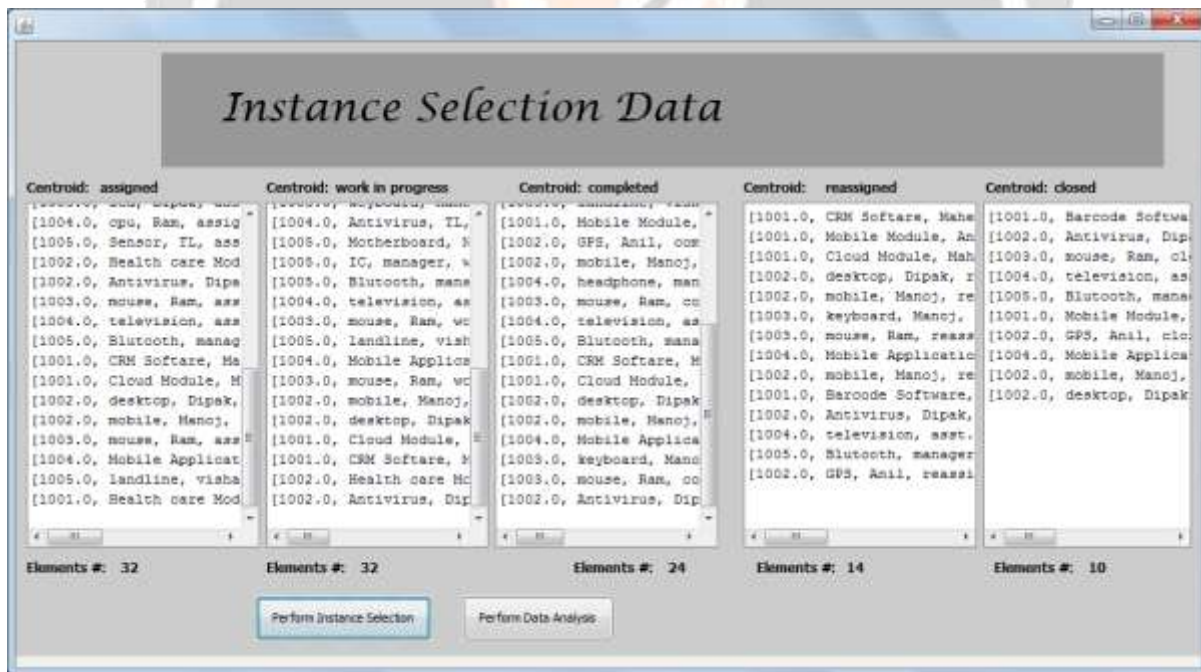


Fig 6.3 Instance Selection

Fig 6.4 Data Analysis Form

7. ACKNOWLEDGMENT

I would like to take this opportunity to express my profound gratitude and deep regard to my guide Prof. Shubhangi Vairagar for her exemplary guidance, valuable feedback and constant encouragement throughout the duration of the project. Her valuable suggestions were of immense help throughout my project work. Her perceptive criticism kept me working to make this project in a much better way. Working under her was an extremely knowledgeable experience for me.

8. REFERENCES

- [1] Jifeng Xuan, He Jiang, Yan Hu, Zhilei Ren, Weiqin Zou, Zhongxuan Luo, and Xindong Wu, "Towards Effective Bug Triage with Software Data Reduction Techniques" *IEEE Transactions on Knowledge and Data Engineering*, vol. 27, no. 1, January 2015.
- [2] Mamdouh Alenezi and Kenneth Magel, Shadi Banitaan "Efficient Bug Triaging Using Text Mining" 2013 academy publisher.
- [3] Francisco Servant "Supporting Bug Investigation using History Analysis" 978-1-4799-0215-6/13 c 2013 IEEE.
- [4] Pamela Bhattacharya, Iulian Neamtiu, Christian R. Shelton, "Automated, Highly- Accurate, Bug Assignment Using Machine Learning and Tossing Graphs", May 2, 2012.
- [5] K. Balog, L. Azzopardi, and M. de Rijke, "Formal models for expert finding in enterprise corpora," in Proc. 29th Annu. Int. ACM SIGIR Conf. Res. Develop. Inform. Retrieval, Aug. 2006, pp. 43–50.
- [6] P. S. Bishnu and V. Bhattacharjee, "Software fault prediction using quad tree-based k-means clustering algorithm," *IEEE Trans. Knowl. Data Eng.*, vol. 24, no. 6, pp. 1146–1150, Jun. 2012.
- [7] H. Brighton and C. Mellish, "Advances in instance selection for instance-based learning algorithms," *Data Mining Knowl. Discovery*, vol. 6, no. 2, pp. 153–172, Apr. 2002. A. K. Uysal and S. Gunal, "A novel probabilistic feature selection method for text classification," *Knowledge-Based Systems*, vol. 36, no. 0, pp. 226–235, 2012.
- [8] S. Kim, H. Zhang, R. Wu, and L. Gong, "Dealing with noise in defect prediction," in Proc. 32nd ACM/IEEE Int. Conf. Softw. Eng., May 2010, pp. 481–490.
- [9] A. Lamkanfi, S. Demeyer, E. Giger, and B. Goethals, "Predicting the severity of a reported bug," in Proc. 7th IEEE Working Conf. Mining
- [10] *Softw. Repositories*, May 2010, pp. 1–10.