



INTERNATIONAL JOURNAL OF PURE AND APPLIED RESEARCH IN ENGINEERING AND TECHNOLOGY

A PATH FOR HORIZING YOUR INNOVATIVE WORK

IMPROVING SOFTWARE DEVELOPMENT PROCESS THROUGH DATA MINING TECHNIQUES

RADHIKA S. MOREY

Prof Ram Meghe College of Engineering, Sant Gadge Baba University (SGBAU), Amravati, Maharashtra

Accepted Date: 05/03/2015; Published Date: 01/05/2015

Abstract: Research in the fields of software quality, maintainability requires the analysis of large quantity of data, which originate from software projects. It is a challenging task in pre-processing the data and synthesizing the composite results. It is very often an error prone task. Data mining techniques are generally considered as the best for preprocessing data. There may be the chance of not getting the proper result for boosting up the software industry, but upto certain extent it will be surely helpful in research field or to the researchers for further innovative development.

Keywords: Data mining, Software engineering, Feature selection.

Corresponding Author: MS. RADHIKA S. MOREY



PAPER-QR CODE

Access Online On:

www.ijpret.com

How to Cite This Article:

Radhika S. Morey, IJPRET, 2015; Volume 3 (9): 804-811

INTRODUCTION

Software development in earlier times was focused on creating algorithms and implementing through programs. This technique was capable of solving almost all software programs but evolution of sophisticated hardware required continuous project planning resulted in low productivity, heavy maintenance cost which led to failure of user expectation and led to stagnation of software development. This crisis was caused by the fact that was no formal method and methodologies support tools for project management.

Software community realized that to solve large software intensive system, they borrowed ideas from other fields of engineering and incorporated into software development. This was the origin of software engineering. Operation. The study of these approaches; led to the application of engineering tools in software development. Knowledge discovery process in databases (KDD) borrowed from data mining gave rise to data mining algorithm that were capable of solving many complicated software problems. Data mining, a branch of computer science is the process of extracting patterns from large data sets by combining methods from statistics and artificial intelligence with database management. Data mining is seen as an increasingly important tool by modern business to transform data into business intelligence giving an informational advantage.

Software engineering data contains a wealth of information about a project's status, progress, and evolution. Here objective is to refine these data through an algorithm so that it can be processed to get a well-defined software process. This can be done by putting the data into the format for further processing in subsequent stages.

Using data mining technique we will explore the potential of this valuable data in order to better manage in subsequent stages of development and produce high quality software systems, which can be delivered on time and within stipulated budget. Mining algorithms fall into four main categories:

- Frequent pattern mining—finding commonly occurring patterns;
- Pattern matching—finding data instances for given patterns;
- Clustering—grouping data into clusters; and
- Classification—predicting labels of data based on already labelled data. [1]

2. EXISTING LAYER MODEL

Firstly existing model gives rise to the challenges of mining software engineering data. Second, development frontier of data mining practice in software engineering. Third, intended to analyze successful cases of mining SE data. Finally, intended to give an overview on commonly used data mining tools. Our overview will help the participants gain a better understanding of available tools. The participants can use such tools in order to explore their data and integrate data mining techniques in their research and day to day work.

Existing model in Fig 1 has got certain drawbacks for which developers are facing problems in conducting research with software repository.

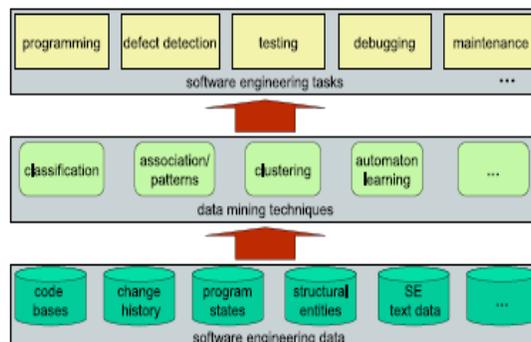


Figure 1. Overview of mining SE data

The main problem developers face working with software repository is that it is difficult to set up experiments that can be extended or replicated as a result of which they invariably face hurdles which the researcher must overcome in order to experiments with huge data chunks from large number of projects.

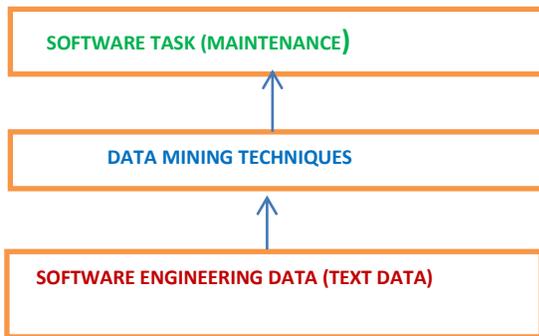


Fig 2: PROPOSED LAYER MODEL.

A. Software Engineering Data

SE text data include bug reports, e-mails, code comments, and documentation for API methods. Common types of text mining algorithms include text clustering, classification, and matching. Example text clustering applications include clustering bug reports to detect duplicate bug reports and thereby reduce inspection efforts, and assigning reports to specific developers to fix the bugs. Example text classification applications include recommending assignment of a new bug report to a specific developer based on the past assignment of old bug reports. Example text matching applications include searching keywords in code comments, API documentation, or bug reports, and detecting duplicates of a given bug report among old reports. [2]

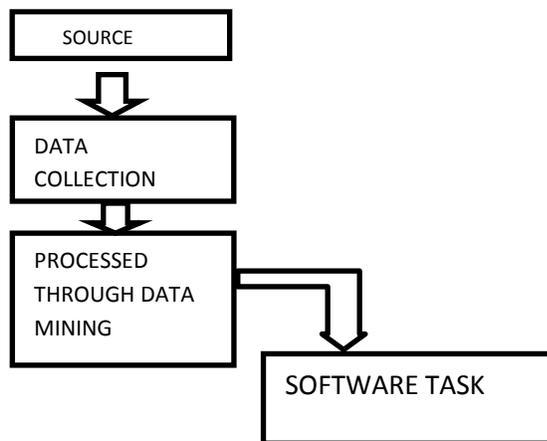


Fig 3: WORKFLOW OF PROPOSED MODEL

Software engineering data may have missing data. Missing data analysis is a wide research area over the past 30-35 years.

The missing data ignoring techniques simply delete the cases that contain missing data. Because of their simplicity, they are widely used (Roth, 1994) this approach has two forms:

1) *List wise deletion (LD)*: is also referred to as case deletion, case wise deletion or complete case analysis. This method omits the cases containing missing values. It is easy, fast, does not 'invent data', commonly accepted and is the default of most statistical packages. The drawback is that its application may lead to a large loss of observations, which may result in too small data sets if the fraction of missing values is high and particularly if the original data set is itself small, as is often the situation for software project estimation (Myrtveit et al., 2001).

2) *Pair wise deletion (PD)*: is also referred to as the available case method. This method considers each feature separately. For each feature, all recorded values in each observation are considered (Strike et al., 2001) and missing data are ignored. This means that different calculations will utilize different cases and will have different sample sizes, an undesirable effect. The advantage is that the sample size for each individual analysis is generally higher than with complete case analysis. It is necessary when the overall sample size is small or the number of cases with missing data is large.

The missing data toleration techniques use a probabilistic approach to handle missing data. They do not predict missing data but assign a probability to each of the possible values. Thus they are internal missing data treatment strategies, which perform analysis directly using the data set with missing values. The missing data imputation techniques estimate missing values for the missing cases and insert estimates obtained from other reported values to produce an estimated complete case. The common forms are as follows:

1) *Mean imputation (MI)*: is also referred to as unconditional mean imputation. This method imputes each missing value with the mean of reported values. The disadvantage is that it leads to underestimation of the population variance.

2) *Regression imputation (RI)*: is also referred to as conditional mean imputation. The regression model is built using the complete observations. It tends to perform better than MI, but still underestimates variance.

3) *Hot-deck imputation (HDI)*: methods fill in missing data by taking values from other observations in the same data set... The choice of which value to take depends on the observation containing the missing value. Randomly choosing observed values from donor cases is the simplest hot-deck method.

4) *Multiple imputations*: means that the missing data are imputed $m > 1$ times, with a different randomly chosen error term added in each imputation. In this method, each missing value is replaced by a set of m plausible values drawn from their predictive distribution. After performing multiple imputations, there are m complete, imputed data sets [3].

B) Feature Selection Method

Feature selection is a pre-processing technique commonly used on high dimensional data. Feature selection studies how to select a subset or list of attributes or variables that are used to construct models describing data. Its purposes include reducing dimensionality, removing irrelevant and redundant features, reducing the amount of data needed for learning, improving algorithms' predictive accuracy, and increasing the constructed models' comprehensibility.

Feature selection extracts the most important set of attributes for model training. a feature-selection algorithm is part of the classification rule. This is why feature selection must be included when using cross-validation error estimation.

Feature Selection for Unlabeled Data: Unsupervised learning, or cluster analysis, aims to group similar objects. Many clustering algorithms assume that domain experts have determined relevant features. But not all features are important. And the presence of many irrelevant features can even misguide clustering results.

Such a model is usually trained using software measurement and defect data from previous development experiences, and the model is then applied to predict the quality of target (under-development) program modules. As a result, practitioners can strategically allocate project resources and focus more on those potential problematic modules. Typically, feature selection techniques are divided into two categories: wrapper-based approach and filter-based approach. The wrapper-based approach involves training a learner during the feature selection process, while the filter based approach uses the intrinsic characteristics of the data based on a given metric, for feature selection and does not depend on training a learner. The primary advantage of the filter-based approach over the wrapper-based approach is that it is computationally faster. Saeys et al. [7] investigated the use of ensemble feature selection techniques, where multiple feature selection methods were combined to yield results. Liu et al. [8] introduced the concept of active feature selection, and investigated a selective sampling approach to active feature selection in a filter model setting.

C) Software Maintenance

Software maintenance is one of the major concerns of software development and maintenance organizations. "Software maintenance is the process of modifying a software system or component after delivery to correct faults, improve performances or other attributes, or adapt to a changed environment." [12] Although the software maintenance phase starts after the delivery of product to the client, it covers a major portion of the cost, effort and time involved in the project

1) Issues in Software Maintenance:

Impact Analysis, Complex code and Architecture, Lack of Understanding, Undefined Process and Procedure for Maintenance, Involvement of Senior Staff, Maintenance Cost Estimation, Improper User Training, and Dependency on Outside Supplier, Lack of Documentation, Measuring Maintenance and Support Service. Motivation of Support personal. [13]

2) Managing Software Maintenance Costs by Developmental Techniques and Management Decision during Envelopment

1. Strive for Commonality
2. Apply Industrial Engineering Practices to Software
3. Engage
4. Adopt a Holistic Approach to Sustainment
5. Develop Highly Maintainable Systems and Software
6. Manage the Off-the-Shelf Software
7. Plan for the Unexpected
8. Analyze and Refine

The Software Sustainment Business Case (use Parametric software sustainment cost estimates) [14]

Software maintenance deals with the management of such changes, ensuring that the software remains correct while features are added or removed.

3. FUTURE WORK

More emphasis can be given in data collection, data preprocessing, and feature selection. There can be some advanced software engineering technique for pre-processing data. We can also combine data mining technique and software engineering technique for better pre-processing.

We should also concentrate on software maintenance cost. Future plan work on this platform includes development of plug –ins for support of other sources like web services and full-fledged API for accessing project meta data in simple way.

4. CONCLUSION

Intensive software project requires handling huge number of data. Data can be text data. It is necessary to pre-process the data through data mining technique. It may or may not be useful for software industry but it will be surely Beneficial for research field

REFERENCES

1. J. Han and M. Kamber, *Data Mining: Concepts and Techniques*, Morgan Kaufmann, 2000
2. Tao Xie and Suresh Thummalapenta, North Carolina State University David Lo, Singapore Management University Chao Liu, Microsoft Research "Data mining for software engineering"
3. A new imputation method for small software project data sets Qinbao Song a*, Martin Shepperd b a Xi'an Jiaotong University, Xi'an, Shaanxi 710049, China b Brunel University, Uxbridge, UB8 3PH, UK Received 17 March 2005; received in revised form 30 April 2006; accepted 3 May 2006 Available online 16 June 2006
4. Y. Jiang, J. Lin, B. Cukie, and T. Menzies. *Variance analysis in software fault prediction models*. In Proceedings of the 20th IEEE International Symposium on Software Reliability Engineering, pages 99–108, Bangalore-Mysore, India, Nov. 16-19 2009.
5. T. M. Khoshgoftaar, P. Rebours, and N. Seliya. *Software quality analysis by combining multiple projects and learners*. Software Quality Journal, 17(1):25–49, March 2009.
6. S. Lessmann, B. Baesens, C. Mues, and S. Pietsch. *Benchmarking classification models for software defect prediction: A proposed framework and novel findings*. IEEE Transactions on Software Engineering, 34(4):485– 496, July-August 2008
7. Y. Saeys, T. Abeel, and Y. Peer. *Robust feature selection using ensemble feature selection techniques*. In Proceedings of the European conference on Machine Learning and Knowledge Discovery in Databases - Part II (2008), pages 313–325, 2008.
8. H. Liu, H. Motoda, and L. Yu. *A selective sampling approach to active feature selection*. Artificial Intelligence, 159(1-2):49–74, November 2004
9. Evolving feature selection". Huan Liu, Arizona State University.
10. A Platform for Software Engineering Research Georgios Gousios, Diomidis Spinellis Department of Management Science and Technology Athens University of Economics and Business Athens, Greece.