

# ChangeChecker: A Tool for Defect Prediction in Source Code Changes based on Incremental Learning Method

Zi Yuan

Department of Computer Science  
Beihang University  
Beijing, China  
yuanzi@sei.buaa.edu.cn

Lili Yu

Software Testing Center  
Second Artillery  
Beijing, China  
xueer-123@263.net

Chao Liu

Department of Computer Science  
Beihang University  
Beijing, China

Linghua Zhang

Software Research and Development Center  
Oracle  
Beijing, China

**Abstract**—In software development process, software developers may introduce defects as they make changes to software projects. Being aware of introduced defects immediately upon the completion of the change would allow software developers or testers to allocate more resources of testing and inspecting on the current risky change timely, which can shorten the process of defect finding and fixing effectively. In this paper, we propose a software tool called ChangeChecker to help software developers predict whether current source code change has any defects or not during the software development process. This tool infers the existence of defect by dynamically mining patterns of the source code changes in the revision history of the software project. It mainly consists of three components: (1) incremental feature collection and transformation, (2) real-time defect prediction for source code changes, and (3) dynamic update of the learning model. The tool has been evaluated in a large famous open source project Eclipse and applied to a real software development scenario.

**Keywords**—software engineering; source code change; defect prediction; incremental learning

## I. INTRODUCTION

In the past few years, more and more attention has been paid to software quality. Software defects, which affect software quality negatively, are considered to be the key risk and major cost driver for both companies that develop software and companies that consume software systems in their daily business [1]. Normally, if software defects are not fixed timely, they will affect the software in three aspects negatively. Firstly, defects that settle in software for a long time may cause new defects or subsequent error modifications [2]. Secondly, if these defects are discovered and reported late, it will be costly for the developer to reacquaint himself with the changed source code, which will increase the cost of defect fix effort. Thirdly, if these defects are reported by customers after software release, it will damage the software reputation significantly. If there is a

tool that can accurately predict whether one or more defects have been introduced immediately after a change is made, it will enable developers to take steps to fix the defects timely.

In this paper, the desired tool called ChangeChecker has been proposed based on machine learning classification method, which takes a source code change as a learning instance, a series of characteristics of it as features, and predicts whether this change introduces defects or not. Given by the data characteristics of the source code changes, there are totally four problems to be handled in building ChangeChecker:

The first problem is which features of a source code change closely related to defects, and how to collect and transform them. The source code changes possess the characteristic of data stream and they constantly arrive with the passing of time. Therefore, the whole data set cannot be obtained by one-time collection and transformation. Incremental way is desirable.

The second problem is that the decision making should be real time. When a new change occurs, the classification result must be given at once. Therefore, it is impractical to rebuild the learning model every time a new change arrives.

The third problem is that due to the variation of underlying defect generation process with time, concept drift is inevitable. Therefore the predictive power of source code change features varies as time: features that in the past have been important, become redundant with the passing of time and new high-predictive features arise that have not been considered before. Therefore, there is a pressing need for a mechanism that can update the learning model dynamically to handle the concept drift.

The fourth problem is that the class distribution of the change data set is imbalanced and unstable. Usually, the number of source code changes that did not introduce defects is much larger than that of the defect introducing changes. But

sometimes the class distribution will vary with the passing of time. Therefore, it is very essential to bring up a mechanism that can mitigate the negative effects of class imbalance on the classification performance.

In order to solve the four problems mentioned above, the ChangeChecker proposed in this paper is mainly made up of three parts: (1) incremental feature collection and discretization, which is used to solve the first problem; (2) real-time defect prediction with incremental Naïve Bayes Classifier, which is used to solve the second problem; (3) dynamic update of the learning model, including incremental feature selection in dynamic feature space and dynamic update of the threshold for classification, which is mainly used to handle the third and fourth problems.

The remainder of this paper is structured as follows. Firstly, we give a series of definitions which closely correlate with the subject of this paper. Secondly, we describe the system framework and working process of ChangeChecker. Then, we discuss the evaluation of ChangeChecker on a large open source project and give a demonstration in a real application scenario. Finally, we draw the conclusion of this paper.

## II. DEFINITIONS

In this section, we present several necessary definitions and formal representations which closely correlate with the subject of this paper.

### A. Indication Labeling

Given a sequence of source code changes  $C = (c_1, c_2, \dots, c_t)$ , the label of a change  $c_i (i = 1, 2, \dots, t)$  is defined as:

$$l(c_i) \in \{Clean, Buggy\} \quad (1)$$

where a change  $c_i$  denotes a source code file commit, *Clean* denotes that there is no defect in change  $c_i$ , and *Buggy* denotes that there is at least one defect in change  $c_i$ .

### B. Features

In this paper, there are mainly five aspects considered as constructing features of source code changes, namely *where* and *when* it happens, *what* happens, *who* manipulates it, and *why* it happens. Each aspect corresponds to a feature group. Formally, we have:

$$F_i = F_i^{where} \cup F_i^{what} \cup F_i^{when} \cup F_i^{who} \cup F_i^{why} \quad (i = 1, 2, \dots, t) \quad (2)$$

where  $F_i$  denotes the available feature set at time step  $i$ , and  $F_i^{where}$ ,  $F_i^{what}$ ,  $F_i^{when}$ ,  $F_i^{who}$ , and  $F_i^{why}$  denote the feature groups from the five aspects.

1) **Feature Group Where:** In this paper, *where* doesn't indicate the exact location of a source code change but the context of it. The characteristics of source code files that have been touched reflect context information of the source code change. It is well understood that if it is difficult for a source code file to be understood thoroughly, the risk of changing it

will be high. Hence, some metrics reflecting the complexity and activity of the touched file are believed to be related to the difficulty of understanding the file and used as features.

a) **Halstead Metrics:** Halstead Metrics were proposed by Maurice Halstead, who argued that the harder the code to read, the more defect prone the modules are [3]. Halstead Metrics are measures of lexical complexity. There are four basic Halstead Metrics (i.e. total number of operators, total number of operands, unique number of operators and unique number of operands) and two derived Halstead Metrics (i.e. volume and programming effort) considered.

b) **McCabe Metric:** Introduced by Thomas McCabe, the idea behind McCabe Metrics is to capture the structural complexity level of a code [4]. The assumption is that it is more likely for the number of defects to increase as the source code gets more complex. Cyclomatic complexity, which is the most famous metric in McCabe Metrics, is considered in this paper.

c) **Graph Metrics:** As indicated by Zimmermann [5], dependencies between software entities are crucial for their successful operation. We assume that the dependency between source code files may correlate with the injection of defects. In this study, we construct a weighted multiple type dependency graph  $G = \langle V, E \rangle$ , where  $V$  denotes the set of nodes (i.e. source code files), and  $E$  denotes the set of edges (i.e. dependencies between two files).  $\forall e_i \in E (i = 1, 2, \dots, M)$ ,  $t(e_i)$  and  $w(e_i)$  denote the type and the weight of the edge respectively. There are four types of dependencies considered, namely method invocation, field access, inheritance, and implementation. The weight indicates the amount of dependencies.  $\forall v_j \in V (j = 1, 2, \dots, N)$ ,  $S_{v_j} = \{e_{s_1}, e_{s_2}, \dots, e_{s_m}\}$  denotes the set of edges that take  $v_j$  as the start point and  $T_{v_j} = \{e_{t_1}, e_{t_2}, \dots, e_{t_n}\}$  denotes the set of edges that take  $v_j$  as the end point. Formally, the two Graph Metrics weighted out-degree and weighted in-degree are represented as follows:

$$D^+(v_j) = \sum_{k=1}^m w(e_{s_k}) \quad D^-(v_j) = \sum_{l=1}^n w(e_{t_l}) \quad (3)$$

d) **History Metrics:** Recently, researchers have shown the usefulness of collecting History Metrics from software repositories for defect prediction models [6]. It is assumed that if a source code file has been changed or fixed many times by many developers in the history, the risk of changing it in the future will be high. Hence, there are totally two History Metrics considered in this paper, namely the number of revisions in the last year and the number of fixes in the last year.

e) **Text Information:** Text data contains plenty of semantic function information about source code, which may relate to the injection of defects. Terms in file text and its directory hierarchy are considered as features in this paper.

2) **Feature Group What:** The source code snippets, which have been added, deleted in a source code change, reflect the content about this change. The content contains rich information about what has been done in the change. We take

terms in these snippets as features in this paper. In addition, the number of lines which have been deleted and added in a change is also considered to describe the size of the change.

3) **Feature Group When:** The time when the change occurs can usually capture a developer's habit and work cycle (e.g. Some developers always introduce more defects after midnight) [7]. In this paper, we take the time of day and the day of week as two primary features in feature group *when*.

4) **Feature Group Who:** Since defects are introduced into the software by developers, the characteristics of developers should be considered. In this paper, we use login name as a feature to distinguish different developers.

5) **Feature Group Why:** Mocus [8] has pointed out that the intent of the software change (e.g. defect fixing) is related to the injection of defects. Since commit message contains rich information about the change intention, terms in it are considered as features in this paper. Meanwhile, the length of the message, which is used to quantify the information contained in the commit message, is also taken as a feature based on the hypothesis that the more complex the intent of a change is, the more defect prone the change is.

These feature groups almost cover all aspects of a source code change, and are considered to be closely related to the injection of defects.

### C. Learning Model

In this paper, the learning model that ChangeChecker based on will be updated with time. At each time step  $t$ , it can be denoted as a five tuple:

$$L_t = \langle ST_t, CP_t, TH_t, FS_t, CL_t \rangle \quad (4)$$

where  $ST_t = \{st_t^1, st_t^2, \dots, st_t^n\}$  denotes a collection of statistics tables,  $n$  indicates the number of available features, and the table  $st_t^i (i = 1, 2, \dots, n)$  holds the number of appearances of the value of feature  $i$  for different classes;  $CP_t = \{cp_t^1, cp_t^2, \dots, cp_t^q\}$  denotes a collection of cut points set,  $q$  indicates the number of quantitative features, and  $cp_t^j (j = 1, 2, \dots, q)$  denotes cut points set used to discretize the quantitative feature  $j$ ;  $TH_t$  denotes the value used to determine the threshold of classification according to calculated probabilities;  $FS_t$  denotes the most predictive features  $(f_{s_1}, f_{s_2}, \dots, f_{s_m})$  selected from  $F_t$  based on  $ST_t$  and an incremental feature selection method; and  $CL_t$  denotes the probabilities that the current instance belongs to each class, which is obtained by an incremental classifier that calculates the probabilities using  $ST_t$  and the selected features  $FS_t$ . In this paper, the Naïve Bayes (NB) is chosen as the incremental classifier because it is simple, effective, efficient and robust. Furthermore, the Information Gain (IG) algorithm is taken as the feature selection method because it is inherently incremental and compatible with NB.

### D. Learning Task

We use an incremental learning framework. At every time step  $t$ , we have historical change sequence with labels, i.e.,

$C = (c_1, c_2, \dots, c_t)$ . When a new code change  $c_{t+1}$  arrives, the task is to predict the label of  $c_{t+1}$ . To this end, we build the learning function  $L_t$  based on the feature set  $F_t = (f_1, f_2, \dots, f_n)$ . We apply  $L_t$  to predict the label of  $c_{t+1}$ . Formally, we have:

$$L_t(c_{t+1} | f_1, f_2, \dots, f_n) \rightarrow l(c_{t+1}) \quad (5)$$

## III. SYSTEM FRAMEWORK

The tool based on an incremental learning framework works as follows. (1) Feature collection and transformation. As source code changes constantly arrive with the passing of time, their features should be collected and transformed incrementally; (2) Prediction. ChangeChecker predicts the newly arrived source code change with the current learning model; (3) Update. ChangeChecker collects the user feedback to update the learning model continually, including statistics tables updating, cut points sets updating, and classification threshold updating. Fig. 1 shows the working process of ChangeChecker.

### A. Feature Collection and Transformation

This step can be further divided into three parts, namely raw data collection, feature extraction, and transformation of quantitative features.

1) **Raw Data Collection:** Raw data of source code changes, which is generated during the software development process, should be collected from software history repositories. Here raw data indicates revisions of source code files, commit logs and source code snippets which have been added and deleted in source code changes. This task is mainly implemented by invoking several commands of the version control system.

2) **Feature Extraction:** After collecting the raw data, ChangeChecker further extracts the features mentioned in above section. The techniques used to extract these features are different since there are various data formats in software history repositories.

Some features that describe the complexity of the internal structure of source files (i.e. Halstead Metrics and McCabe Metric) should be extracted using static analysis technique.

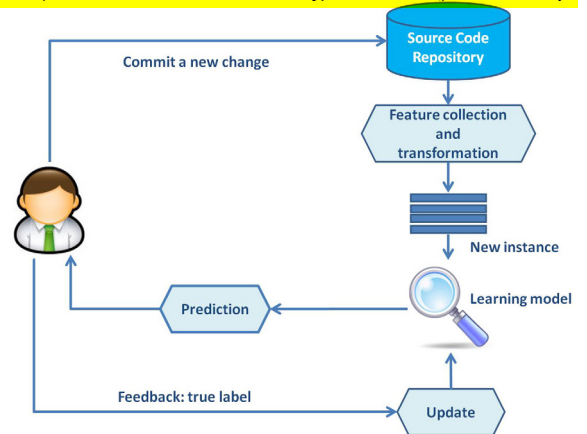


Fig. 1. Incremental learning framework

Features that used to describe the relationships between the current source code file and other files in the project (i.e. weighted out-degree and weighted in-degree) are extracted with incremental building method of software dependency graph proposed in this paper. In order to obtain graph metrics for each revision of the touched file, the evolutionary series of the software dependency graph should be extracted incrementally. As each change arrives, ChangeChecker doesn't rebuild the whole graph, but edits it based on the current change. This editing operation is non-trivial, which should be implemented by converting the source code file change into dependency graph change (i.e. the added nodes and edges, the removed nodes and edges, and the changed nodes and edges). In the converting process, current and previous version of the touched file are compared with each other based on abstract syntax trees to identify modified software entities and their relationships defined in this file. Because the qualified names of some program entities declared in other source files cannot be acquired by parsing each source file alone, the symbol table should be scanned for two times.

Some features can be acquired directly from the version control system (e.g. commit author) or just by simple processing (e.g. time of day and day of week) and computing (e.g. History Metrics).

Text data involved in this paper contains plenty of information about source code changes. Most of the features are extracted from it. ChangeChecker use the BOW and its metamorphosis to generate features from text data, which draws lessons from the work of Kim [9]. For the commit message, which can be considered as common text, features are extracted by taking BOW approach directly. Whereas, for source code file and added/deleted code snippets, in addition to common terms, several important operators (e.g. =, &&, !) should not be ignored. Hence a metamorphosis of BOW called BOW+ [9] is used. Furthermore, the directory and file name are transformed into terms with another metamorphosis of BOW called BOW++ [9], which splits text using slashes and capitals.

3) **Feature Transformation:** In order to conform to the data format required by the statistic table, some quantitative features (e.g. length of commit message, the complexity metrics, etc) are transformed into qualitative features. Given by the stream data characteristics of the source code changes, an incremental flexible frequency discretization method IFFD has been adopted [10]. Every time a new change arrives, ChangeChecker will transform its quantitative features into qualitative features based on the current sets of cut points. According to the value distributions of features, the sets of cut points are constantly updated with time.

### B. Prediction

The primary purpose of ChangeChecker is to help developers predict whether the incoming source code changes have defects or not. On the arrival of a new change, the tool will make predictions based on the most powerful features and classification threshold at present. When the developer receives

the result, he can check the new change to determine if the prediction is true or not, and then give feedbacks to ChangeChecker.

### C. Update

When the feedback from developers has been received, the new change is considered as a new training instance used to update the current state of the learning model. The update operation can be divided into three parts, statistics tables updating, cut points sets updating, and classification threshold updating. The statistics tables are updated according to the values of the new change's features and its class label. If the new change has any new features which have not appeared in previous changes before, new statistics tables corresponding these features will be created and added into the collection of statistics tables. The incremental flexible frequency discretization method (IFFD) used in ChangeChecker sets the interval frequency of each feature to be a range instead of a constant. The sets of cut points are updated if the addition of new source code change lead to the frequencies of some features' intervals reach upper limit. The threshold is updated according to the prediction result and the true label of the new change. At the end of update, the feature evaluation metrics are then re-calculated and ranked with the Information Gain method. The whole update process is elaborated in Fig. 2 and Algorithm 1.

## IV. EVALUATION

To evaluate ChangeChecker, we have conducted a set of experiments on 5000 sequential code changes from Eclipse JDT data. The result shows the effectiveness of the proposed method. Compared to non-incremental learning model, who has undergone a significant decline as time and got recall lower than 20 percent after 2400 changes, our incremental learning framework with dynamic feature space significantly performs better in dealing with concept drift, and its precision and recall both keep stable at around 70 percent.

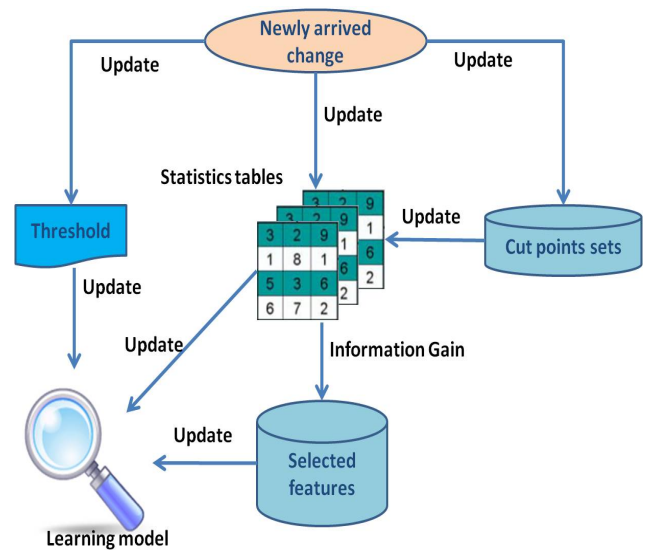


Fig. 2. The update process



Algorithm 1 Update	Algorithm 2 IFFD
<b>Input:</b> $l$ : the label of newly arrived change $c_t$ $l_p$ : the prediction label of change $c_t$ $pro$ : probability of $c_t$ belonging to <i>Buggy</i> $F_c$ : features of change $c_t$ $Labelset$ : the set of labels $QF$ : the set of quantitative features $F_i$ : the set of code changes' features $L_i$ : learning model $\langle ST_i, CP_i, TH_i, FS_i, CL_i \rangle$ <b>Output:</b> $F_{t+1}$ : the set of code changes' features $L_{t+1} \leftarrow \langle ST_{t+1}, CP_{t+1}, TH_{t+1}, FS_{t+1}, CL_{t+1} \rangle$	<b>Input:</b> $F_c$ : features of change $c_t$ $CP_i$ : collection of cut points sets $QF$ : the set of quantitative features <b>Output:</b> $CP_{t+1}$
1: $CP_{t+1} \leftarrow IFFD(F_c, CP_i, QF)$ 2: $TH_{t+1} \leftarrow DynamicThreshold(l, l_p, pro, TH_i)$ 3: <b>for</b> each feature $f \in F_c$ <b>do</b> 4: <b>if</b> feature $f \notin F_i$ <b>then</b> 5:     Add $f$ to set $F_i$ 6: <b>for</b> each label $l_e \in Labelset$ <b>do</b> 7: <b>for</b> each value of feature $f$ <b>do</b> 8: $ST_i[f][value][l_e] \leftarrow 0$ 9: <b>end for</b> 10: <b>end for</b> 11: <b>end if</b> 12: <b>end for</b> 13: <b>for</b> each feature $f \in F_i$ <b>do</b> 14: <b>if</b> feature $f \in F_c$ <b>then</b> 15: $value \leftarrow$ the value of change $c_t$ about feature $f$ 16: $ST_i[f][value][l] \leftarrow ST_i[f][value][l] + 1$ 17: <b>else</b> 18: $ST_i[f][0][l] \leftarrow ST_i[f][0][l] + 1$ 19: <b>end if</b> 20: <b>end for</b> 21: $F_{t+1} \leftarrow F_i$ 22: $ST_{t+1} \leftarrow ST_i$ 23: $F_{sort} \leftarrow$ sort $F_{t+1}$ based on Information Gain and $ST_{t+1}$ 24: $FS_{t+1} \leftarrow$ select subset of $F_{t+1}$ based on $F_{sort}$ 25: $CL_{t+1} \leftarrow$ calculate the classification 26: $L_{t+1} \leftarrow \langle ST_{t+1}, CP_{t+1}, TH_{t+1}, FS_{t+1}, CL_{t+1} \rangle$	1: <b>for</b> each feature $f \in F_c$ <b>do</b> 2: <b>if</b> feature $f \in QF$ <b>then</b> 3: $value \leftarrow$ the value of change $c_t$ about feature $f$ 4: $interval \leftarrow$ get the interval of $value$ based on $CP_i[f]$ 5:     insert $value$ into $interval$ 6: $freq \leftarrow$ get the number of values in the $interval$ 7: <b>if</b> $freq >$ upper limit of the interval frequency <b>then</b> 8: $newpoint \leftarrow$ split $interval$ 9:       add $newpoint$ to $CP_i[f]$ 10: <b>end if</b> 11: <b>end if</b> 12: <b>end for</b> 13: $CP_{t+1} \leftarrow CP_i$
	<b>Algorithm 3 DynamicThreshold</b> <b>Input:</b> $l$ : the label of newly arrived change $c_t$ $l_p$ : the prediction label of change $c_t$ $pro$ : probability of $c_t$ belonging to <i>Buggy</i> $TH_i$ : threshold of classification <b>Output:</b> $TH_{t+1}$
	1: <b>if</b> $l \neq l_p$ <b>then</b> 2: <b>if</b> $l = Clean$ and $pro < 1$ <b>then</b> 3: $TH_t \leftarrow pro$ 4: <b>end if</b> 5: <b>if</b> $l = Buggy$ and $pro > 0$ <b>then</b> 6: $TH_t \leftarrow pro$ 7: <b>end if</b> 8: <b>end if</b> 9: $TH_{t+1} \leftarrow TH_t$

## V. DEMONSTRATION

We showcase the ChangeChecker tool using a code evolution parser called "javabear" as the application example, which is being developed by 12 PHD and master students in our lab. There have been 5000 file commits in the version control systems. Fig. 3, 4 and 5 show the user interface of the ChangeChecker. When a developer wants to use the tool ChangeChecker, he should firstly select the "initialize" option and start the initializing process. When he has completed several file modifications and commit them into the version control system, he could verify his commit by selecting the "verify" option. When the verifying process has been finished, the system will list several risk files and their risk area, which might have defects introduced in the last

commits. The developer will check and test these risk files and then give feedback to the system. Finally, updating process is started.

```

ChangeChecker in Helper
ChangeChecker is a cvs commit checking system, which is used to help developers check whether the new commit has any defects or not.

Its work process has been divided into three steps. At first, the developer initialize a ChangeChecker system based on history data in cvs. Then the checker is started just as the cvs commit has been finished and the risk files and their risk area, which might have defects, are listed to warn the developers to check and test them. At last, the system is updated after the developer feedback has been inputted into it.

1) Initialize the ChangeChecker
2) Verify my checkin

Select the option you would like to perform [1-2] 1
Input the name of your project: javabear
[=====]100.0%

Mission Completed, goes to the next step.

```

Fig. 3. The initialization of ChangeChecker

```

1) Initialize the ChangeChecker
2) Verify my checkin
Select the option you would like to perform [1-2] 2
Input your log in name: yuanzi
risk files has been listed.
ID:1
  name:PowerDatabase.java
  path:javabear/src/buaa/sei/javabear/bdb/PowerDatabase.java
  location(line number):1;3-6;9-22
ID:2
  name:DataCenter.java
  path:javabear/src/buaa/sei/javabear/bdb/DataCenter.java
  location(line number):1;5-7;9-43
[=====]100.0%

```

Fig. 4. The Prediction of Risk Files

```

Please input your feedback for ID 1.
1) Clean
2) Buggy
Select your feedback [1-2] 2
Please input your feedback for ID 2.
1) Clean
2) Buggy
Select your feedback [1-2] 1
System begin to update:
[=====]100.0%
Updating completed.

```

Fig. 5. User Feedback and Update

## VI. CONCLUSION

This paper presents a tool called ChangeChecker that is capable of predicting whether the current source code change has any defects during the software development process. The tool concerns about features of source code changes from five different dimensions and is designed based on an incremental learning framework (incremental feature collection and transformation, real-time prediction, and dynamic update of learning model). After evaluation and trial in the real application scenario, it shows that the tool overcomes the concept drift problem, and ensures the high efficiency and good stability of defect prediction.

## REFERENCES

- [1] E. Giger, M. Pinzger, and H. Gall, "Comparing fine-grained source code changes and code churn for bug prediction", in Proc. MSR, 2011, pp.83-92.
- [2] Q. Song, M.J. Shepperd, M. Cartwright, and C. Mair, "Software Defect Association Mining and Defect Correction Effort Prediction", IEEE Transactions on Software Engineering, vol. 32, pp.69-82, 2006.
- [3] M. H. Halstead, "Elements of Software Science. Elsevier", 1977.
- [4] T.J. McCabe, "A Complexity Measure", IEEE Transactions on Software Engineering, vol. 4, pp.308-320, 1976.
- [5] T. Zimmermann and N. Nagappan, "Predicting Subsystem Failures using Dependency Graph Complexities", in Proc. ISSRE, 2007, pp.227-236.
- [6] T.L. Graves, A.F. Karr, J.S. Marron, and H.P. Siy, "Predicting Fault Incidence Using Software Change History", IEEE Transactions on Software Engineering, vol. 26, pp.653-661, 2000.
- [7] J. Eyolfson, L. Tan, and P. Lam, "Do time of day and developer experience affect commit bugginess", in Proc. MSR, 2011, pp.153-162.
- [8] A. Mockus and D. M. Weiss, "Predicting risk of software changes", Bell Labs Technical Journal, pp. 169-180, 2000.
- [9] S. Kim, E.J.W. Jr., and Y. Zhang, "Classifying Software Changes: Clean or Buggy?", IEEE Transactions on Software Engineering, vol. 34, pp.181-196, 2008.
- [10] J. Lu, Y. Yang, and G. I. Webb, "Incremental discretization for naive-bayes classifier", in ADMA'06, 2006, pp. 223-238.