

Process Metrics Based Bug Prediction in GitHub Projects: Does Inequality Matter?

V. Milias*, A. Adhikari[†], L. Gasparini[‡] and C. M. Valsamos[§]
 Email: {v.milias*, a.adhikari[†], l.gasparini[‡], c.valsamos[§]}@student.tudelft.nl

Abstract—Defect prediction has been attracting a lot of attention in the recent years, and the choice of the right set of metrics is crucial for the effectiveness of predictors. Previous research has shown that to this regard, process metrics score better than code metrics in most of the cases. In this paper we investigate whether the contributions patterns influence the performance, stability and portability of the predictors based on process metrics. In particular, we compare the predictors performance in two sets of projects: one with high contributions inequality and one with low inequality. Our research suggests that the predictors trained on projects with low contributions inequality, give better results regarding performance.

Keywords—Bug prediction, Process metrics, Github, Gini coefficient.

I. INTRODUCTION

In software development, defect prediction has been shown to be very important. The effectiveness of these predictors is highly influenced by the used metrics. From prior research [1], [2], [3], [4] we can conclude that process metrics are important indicators of defect proneness. Rahman and Devanbu [1] investigated the efficacy of process and code metrics from different angles. They compared the effectiveness of process and code metrics based prediction models, while taking performance, stability, portability and stasis into account. First, they came to the conclusion that process metrics are, in general, more useful. Second, they found out that the values of code metrics change less from one release to another compared to process metrics. This leads to the stagnation of the prediction model. Process metrics look more flexible to incorporate relevant changes for better bug prediction.

We test this from a different angle by researching whether process metrics are equally useful in projects with different characteristics. In this paper we investigate whether the contributions patterns influence the characteristics of bug predictors based on process metrics. For this reason we consider two sets of projects, all written in Java and maintained by the Apache Software Foundation: one with low contributions inequality and one with high inequality. A high inequality means that most of the code is written by a few core developers and other contributions come from developers that only write a few lines of code. C. Bird et al. [2] came to the conclusion that high value of ownership per file is associated with a lower number of defects. We then study whether in projects with a relatively low contribution inequality, which translates to low average ownership of the files, the performance of the process metrics based bug predictors is different.

Regarding performance, our research suggests that the bug prediction on low Gini projects is better than on high Gini projects. Regarding stability and portability, we did not find any significant difference between low and high Gini projects.

Our paper is organized as follows: *background and theory*, *methodology*, *results*, *related work* and *conclusion*.

II. BACKGROUND AND THEORY

Bug prediction in software systems has been attracting a lot of interest in the recent years [5], [6], [7]. Efforts have been put into developing models that can successfully identify files that are defect prone, varying in the type of metrics used and the learning models employed.

The metrics can be generally divided in two categories: *process metrics* and *code metrics*. The first involve the characteristics of the contributions to the software, such as the number of developers or the scattering of the changes made to the files, while the latter measure the intrinsic properties of the code, such as the number of methods or the cyclomatic complexity.

In [1] the two categories were compared, on the basis of the performance, stability, portability and stasis of the predictors built on them, and process metrics emerged as a clear winner.

In this paper we investigate whether the ownership of the files in the project plays a role in the performance of bug predictors based on process metrics. A project with a high level of ownership has most of the contributions made by a small number of core developers, while in a project with low ownership the contributions are more equally distributed.

To measure the contributions' inequality we use the Gini coefficient, defined as:

$$G = \frac{\sum_{i=1}^n \sum_{j=1}^n |x_i - x_j|}{2n \sum_{i=1}^n x_i}$$

where x_i represents, in our case, the number of contributions made by contributor i , and n is the total number of contributions.

It is a popular measure, typically used in economics to assess the income inequality of the population, but it has also been successfully employed in the field of software engineering. In [8] for instance, the authors find a negative correlation between the Gini index computed on history of changes and the bug proneness.

The coefficient ranges from 0 to 1. A Gini coefficient of zero means perfect equality while a Gini coefficient of 1 expresses maximal inequality among values (e.g., for a large number of

people, one person has all of the income while the rest has nothing).

Since we want to analyze the differences from different perspectives, we compare the *performance*, *stability* and *portability* of the defect predictors on two sets of projects, one with high contribution inequality and one with low inequality.

Portability, as noted in [9], is particularly important for projects missing historical data to train the predictors on. For instance, this is the case for newly born projects. In this case it is useful to be able to train the predictors on a different existing project for which the data is available.

This leads to our research questions:

- **RQ1:** Does the inequality of contributions affect the *performance* of bug predictors based on process metrics?
- **RQ2:** Does the inequality of contributions affect the *stability* of bug predictors based on process metrics?
- **RQ3:** Does the inequality of contributions affect the *portability* of bug predictors based on process metrics?

III. METHODOLOGY

A. Projects studied

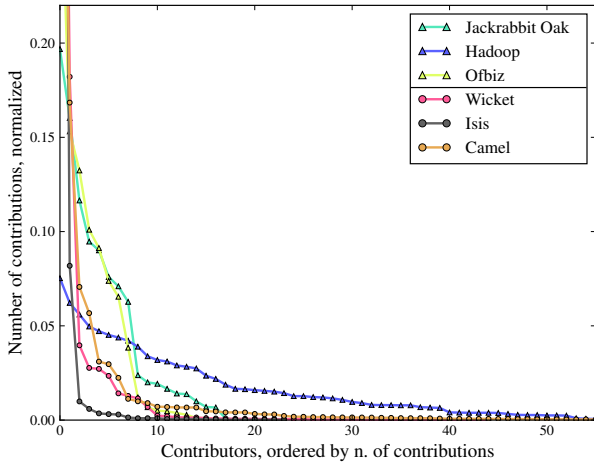


Fig. 1. Contributions distribution

In our study we considered the 6 projects that are listed in table I. To choose them, we considered all the Java projects maintained by the Apache Software Foundation. We computed the Gini index for each of them using the contributions data obtained from the Github API. We filtered all the projects to consider only those with more than 10 contributors and 1000 commits, to have enough data to train the predictors effectively. Then we chose the 3 projects with the highest Gini index and the 3 projects with the lowest, that satisfied our requirements. The high Gini projects consist of 21 releases and the low

Gini ones consist of 20 releases. Subsequently, the number of releases is almost equal among the two categories of projects.

In figure 1 the contributions patterns of the two sets of projects are plotted, where the x axis represents the individual contributors (sorted by the number of contributions) and the y axis is the number of contributions, normalized by the area under the line (sum of all the contributions).

All the projects, while developed using the SVN versioning system, are mirrored to their respective GIT repositories on Github. For each release, we extract all the process metrics using the GIT history.

The projects use JIRA as the issue tracking system, from which we extract the defect information. From the bugs listed as fixed in JIRA we extract the commits fixing them and then we use GIT BLAME to find in which release the bug was introduced.

B. Predicting defects

The choice of the classifier has been shown not to be significant in terms of performance advantage for models built on process metrics [1], [4]. We decided to use logistic regression to compare different models.

For the performance, we train the model on a release i and we test it on release $i + 1$. For what concerns stability, we train the model on release i and we test it on the releases $i + 1, i + 2, \dots$. Regarding portability, we train the predictor on a release of a project and we test it on a release of a different project.

The model that we use is a binary classifier, meaning that it only predicts a file as containing a bug (*defective*) or not (*clean*).

The resulting logistic regression models give also information about the usefulness of the process metrics. A low p-value of a coefficient means that this metric is significant (lower than 0.05). The order of usefulness of different metrics of projects with high and low Gini coefficients is compared.

C. Process metrics

The process metrics that we use can be found in table II. All process metrics, except for *DDEV*, are relative to the duration of the release and at a file level. We calculate the number of commits made to a file *COMM*, the total number of contributors who modified the file *ADEV* and the total number of developers who contributed to this file up to this release *DDEV*. Regarding the changes in a file we take into consideration the *ADD*, *DEL*, *OWN* and *MINOR* process metrics. *ADD* and *DEL* metrics measure the amount of lines added and deleted respectively, normalized by the total number of added and deleted lines in the file. With the *OWN* metric we get the information about the number of lines added by the owner of the file. The owner of a file is the developer that authored the highest amount of lines. *MINOR* metric is the sum of commits made by the minor contributors [2]. These metrics have proven to work well and are widely used in literature and research [1], [2], [10], [4], [11].

TABLE I. PROJECTS ANALYZED. LOW GINI PROJECTS ARE SHOWN AT THE TOP, HIGH GINI PROJECT AT THE BOTTOM.

Project name	Releases	Gini index	Commits	Contributors
Jackrabbit Oak	1.5.1, 1.5.2, 1.5.3, 1.5.4, 1.5.5, 1.5.6, 1.5.7	0.561391	11679	19
Hadoop	2.5.0, 2.5.1, 2.6.0, 2.6.1, 2.6.2, 2.6.3, 2.7.0, 2.7.1	0.662045	14673	82
Ofbiz	11.04.02, 11.04.03, 11.04.04, 11.04.05, 11.04.06	0.683986	17848	19
Wicket	6.18.0, 6.19.0, 6.20.0, 7.0.0, 7.1.0, 7.2.0	0.919441	5898	44
Isis	1.8.0, 1.9.0, 1.10.0, 1.11.1, 1.12.0, 1.12.1, 1.12.2, 1.13.0	0.934642	7301	23
Camel	2.12.0, 2.13.0, 2.14.0, 2.15.0, 2.16.0, 2.17.0, 2.18.0	0.953537	25619	244

TABLE II. PROCESS METRICS

Name	Description
COMM	Commit Count
ADEV	Active Dev Count
DDEV	Distinct Dev Count
ADD	Normalized Lines Added
DEL	Normalized Lines Deleted
OWN	Owner's Contributed Lines
MINOR	Minor Contributor Count

D. Evaluation

The two evaluation metrics used to assess the accuracy of the predictions are F_1 score and AUC . AUC is a metric designed to describe the accuracy of the predictor in a threshold invariant manner. Different thresholds lead in fact to different numbers of misclassifications (FN, FP) and correct classifications (TP, TN).

F_1 score considers both *precision* and *recall*, where

- *Precision* is defined as $\frac{TP}{TP+FP}$ and measures the relevance of the predictions. A high precision (> 0.5) means that most of the files classified as containing a bug did indeed contain one.
- *Recall* is defined as $\frac{TP}{TP+FN}$. In our case it measures the ratio of the files predicted as containing a bug over the total number of files that were discovered as doing so later on.

and is given by the harmonic mean of *precision* and *recall*, as follows

$$F_1 = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

The F_1 score ranges from 0, when precision or recall are zero, to 1 when they are both 1.

The AUC metric is derived from the *ROC curve* plot. The curve represents the relation between the recall (true positive rate) and the false positive rate ($\frac{FP}{FP+TN}$) at different thresholds, and starts always in the point (0,0) ending in the point (1,1). The best possible classifier, meaning the one having 0 false positives and 0 false negatives, is represented by the point in the upper left (0,1), called *perfect classification* while a *random classifier* results in a point in diagonal.

The AUC (*area under the curve*) is the area under the *ROC curve* and can be interpreted as the probability that the classifier would correctly classify an element better than the random classifier.

These evaluation metrics are broadly used in literature. [1], [2], [12], [4], [11].

IV. RESULTS

In table III the average number of files, defective files and *post*-defective files are listed. By *post*-defective file we mean a file that introduced a bug that was discovered only after the next release. For this files the defect information can only be used for testing and not for training.

TABLE III. DEFECT STATISTICS PER PROJECT

Project	Average number of files	Average number of defective files	Average number of <i>post</i> -defective files
Jackrabbit	2818	16	9
Hadoop	5968	134	82
Ofbiz	1408	7	2
Wicket	3167	16	8
Isis	5031	18	8
Camel	11949	155	54

A. Performance

RQ1: Does the inequality of contributions affect the *performance* of bug predictors based on process metrics?

To evaluate the importance of the inequality of contributions on the performance of the bug predictors we used the logistic regression model per release.

More specifically, we used the following tactic. For each project, we trained our model on a release and we tested it on the next release. By doing that for all the (release, next-release) pairs we calculated the mean of the AUC and the F_1 score for all the pairs of all projects. We present these results in Figure 2 where we also group them by their Gini index.

In general, we can not report any significant difference for the AUC between the two groups but we can notice larger F_1 scores on the low Gini projects. The high values of the AUC can be explained by unbalanced data (number of defective

files \ll number of non-defective files). Additionally, we can observe that for both groups the values of both AUC and F_1 are fluctuating. Moreover, we calculated the means of the AUC s and F_1 scores for each group.

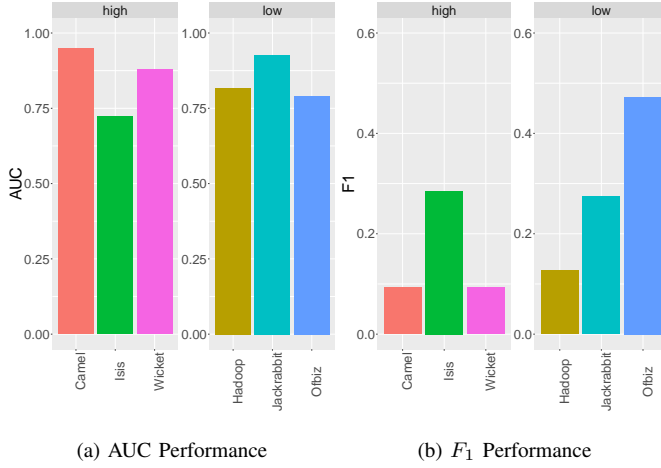


Fig. 2. Performance per project

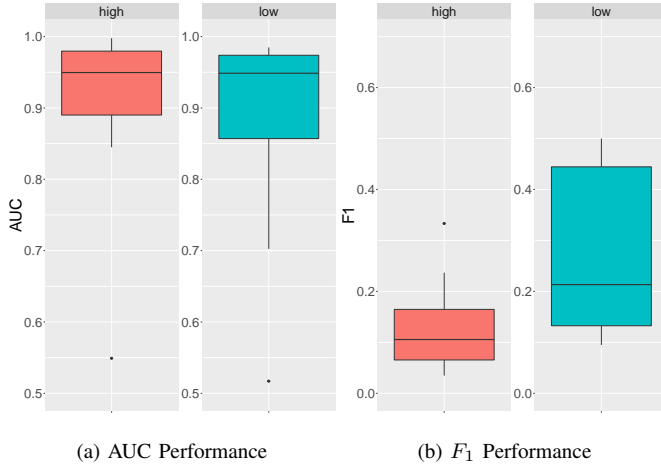


Fig. 3. Performance of low and high Gini projects

We present these results in Figure 3. It is clear from the figure, that the AUC is similar for both groups, around 0.95, however the F_1 score for low Gini projects is higher than the one for high Gini projects, with a mean of 0.13 and 0.23 respectively, which is an interesting outcome.

We think that this is due to the variability of the feature vectors. We suspect that the values of the process metrics of a low Gini project vary more than high Gini projects, hence a better performance for low Gini projects. For example, when the commits are more equally distributed among the contributors the active developers count is more informative. In the high Gini projects the number of the contributors who modify each file is rarely changes (small number of the active

contributors) in comparison with the low Gini, where the number of the active contributors is higher.

B. Stability

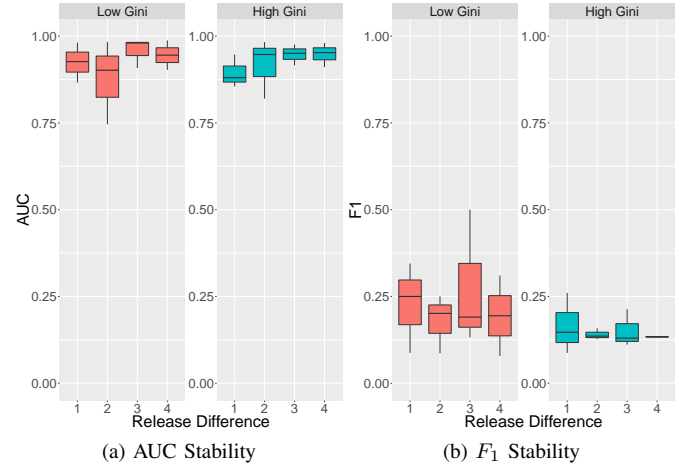


Fig. 4. Stability of low and high Gini projects

RQ2: Does the inequality of contributions affect the *stability* of bug predictors based on process metrics?

Stability is evaluated by training on each release and testing the resulting predictor on the following future releases. Figure 4 shows the mean of the AUC and the F_1 measure when predicting 1 to 4 releases in the future, using each release. The figure shows that for both low and high Gini projects the performance does not change significantly on future releases e.g. there is not a significant difference between projects with high and low Gini regarding stability.

C. Portability

RQ3: Does the inequality of contributions affect the *portability* of bug predictors based on process metrics?

TABLE IV. TRAINING PROCEDURE OF PORTABILITY

Training set	Test set
High Gini project	High Gini projects
High Gini project	Low Gini projects
Low Gini project	High Gini projects
Low Gini project	Low Gini projects

We were interested in seeing how the inequality of contributions affects the portability of process metrics. All projects have several releases so we take the average of the metrics to compare the two groups. We train a predictor on one

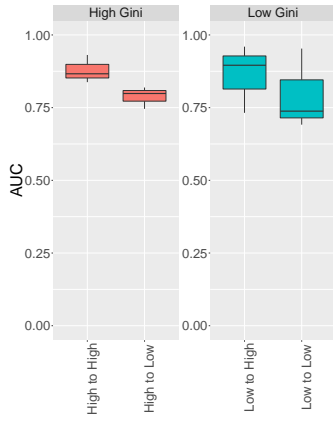


Fig. 5. AUC Portability

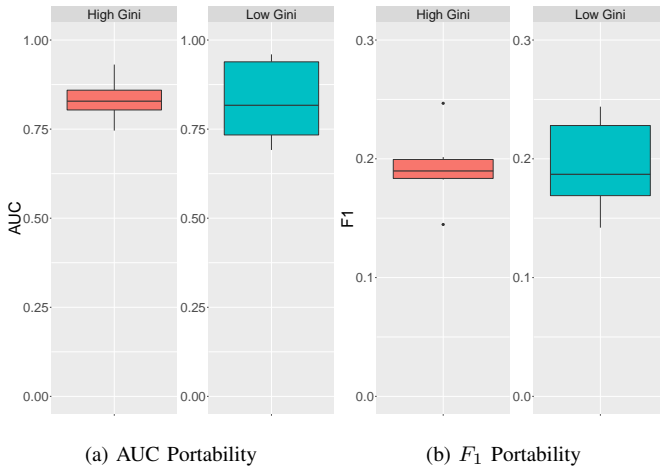
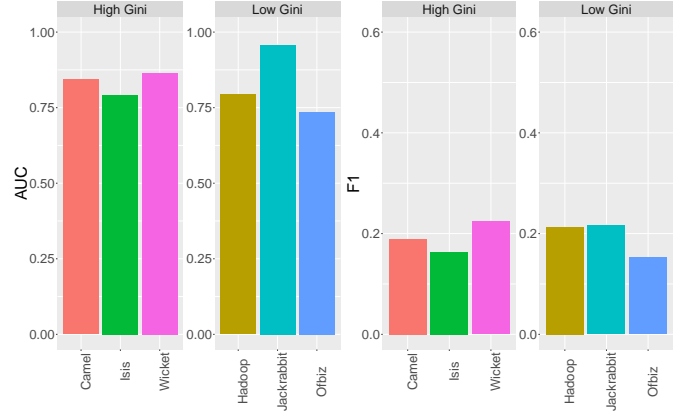


Fig. 6. Portability of low and high Gini projects

release of a project and we test it on all releases of the other projects. During the training process the high and low Gini were swapped in the training and the test set accordingly. For one iteration of the training process, the training set included a high Gini project, the test set included all the low Gini projects and we calculate the average for all the evaluation metrics. In the figures the Gini label corresponds to the Gini of the training set. For instance, “high Gini” shows that the training set contained a high Gini project.

In Figure 5 we plot the performance of the model as indicated in table IV. Figure 6 compares the performance of the model per Gini. The AUC of both high and low Gini is close to 0.83. We also don’t observe a significant difference in F_1 metric between high and low Gini projects which are close to 0.19. In Figure 7 we compare the performance of the model per project. We do see that when we have Jackrabbit as the training set the AUC value of the model increases to 0.95. It is worth mentioning that each project achieves an AUC value above 0.75 except the Ofbiz project. The value of the Gini index of Ofbiz project is close to the mean of Gini indexes so



(a) AUC Portability

(b) F_1 Portability

Fig. 7. Portability per project

we cannot conclude why it has the lowest AUC value among the projects.

Our findings, portability-wise, do not indicate that averagely there is a significant difference between high and low Gini projects. However, we can point out that Jackrabbit which has the lowest Gini index achieves the best performance AUC -wise. From all the figures about portability it is obvious that the AUC and F_1 values among low Gini projects vary more in comparison with the values of the high Gini projects.

D. Significance of Metrics

We were interested in discovering the significance of the metrics that we used in our model. More precisely, we focused on exploring if the order of the metrics by their significance is different among the high and the low Gini projects. For that purpose, we ordered the metrics by their significance for each project (by taking the mean of the p-value for the different (release, next-release) pairs of each metric) and grouping these results for the high and the low Gini projects. In table V we present the ordering of the metrics’ significance.

TABLE V. ORDERING OF METRICS BY THEIR SIGNIFICANCE

High Gini projects	Low Gini projects
COMM	ADEV
ADEV	COMM
OWN	DEL
ADD	ADD
MINOR	OWN
DEL	MINOR

The order of significance of the used metrics is different for the two sets of projects. For example, commit count and active developers count are respectively the most important metrics for high Gini and low Gini projects for bug prediction.

V. THREATS TO VALIDITY

A. Generalizability

We analysed three low and three high Gini projects which resulted in a total of 41 releases. Even though the number of

analysed releases is high, more projects have to be considered to generalize the results of our research. We considered only open-source projects, hence our findings are also less generalizable for commercial projects.

Moreover, while there is a significant difference between the Gini indexes of the chosen high Gini and low Gini projects (average difference of 0.3), it would be better if we could increase this difference. This proved to be hard in practice, since when we filtered the projects which were maintained by the Apache Software Foundation, to consider those with more than 10 contributors and 1000 commits, the projects with the lowest Gini indexes didn't contain any information about the defective files. Having more complete bug information would have enabled us to consider projects with lower Gini index for the low Gini group.

Furthermore, we computed the Gini index for each project by taking all of its releases into account, while in our experiments we train on release basis. It could be possible that the Gini indexes of the initial releases are high while in the rest of the releases they are low, because more contributors join the project later. More accurate results could be obtained by labeling the releases as high and low instead of projects.

B. Completeness of Process Metrics

We considered popular process metrics presented in literature [1]. In future work more metrics could also be considered. Change in entropy could be taken into consideration, which measures the scattering of changes to a file [13]. Kim et al. [12] presented a range of different process metrics based on co-commit-neighbors.

C. Evaluation

We used several metrics to evaluate the performance of the model and to have a complete picture of the evaluation. More specifically, we might get deceived by seeing the *ROC curve* because we can get very good results for predicting non-defective files correctly. This is why we considered the *F measure* as well. It is worth mentioning that none of the evaluation metrics take cost-effectiveness into account. With cost-effectiveness metrics we value the prediction of the file more if the file is defect dense. Predicting a defect dense file as defective is more crucial than predicting a file as defective with a few defects.

VI. CONCLUSION

In this paper we investigated whether the contributions patterns influence the performance of the predictors, based on process metrics in a releases-oriented setting. In particular, we compared the predictors performance, stability and portability: one with high contributions inequality and one with low inequality. These projects have respectively a high and a low Gini value. We analysed three low and three high Gini projects which resulted in a total of 41 releases. Even though the number of analysed releases is high, more projects have to be considered to generalize the results of our research.

Regarding performance, our research suggests that the bug prediction of low Gini projects are better than high Gini projects. We think that this is due to the variability of the feature vectors. We suspect that the values of the process metrics of a low Gini project vary more than high Gini projects, hence making the metrics more informative. This can be investigated in future work. Regarding stability and portability, we did not find any significant difference between low and high Gini projects.

VII. RELATED WORK

Rahman and Devanbu [1] strongly advise everyone to use process-metric-based models instead of code-metric-based models for bug prediction. Their experiments, though, involve only projects maintained by Apache Software Foundation, which have a high level of ownership. The process metrics that we used are drawn from their research.

Zhimin *et al.* [9] investigated the usability of defect predictions in cross-projects contexts finding that the performance of predictors learned on different projects is comparable to that of predictors that are trained on data from the same project; in the best cases they find that the performance is even higher.

C. Bird *et al.* [2] examine the relationship between ownership and software quality. With their research they give emphasis to the fact that process metrics like *OWN* define the quality of the software being produced and influences the occurrence of defects. Consequently, process metrics are very useful indicators of defect proneness. They came to the conclusion that high value of average ownership per file is associated with less defects. We examine whether inequality in ownership per project influences the prediction model. Regarding the process metrics that we are going to use in our research, we include several metrics like the *MINOR* metric and the *OWN* metric which is used in their research.

Foucault *et al.* [3] explored the relationship between ownership metrics and fault proneness on open source software projects. As a result process metrics that measure ownership of a project play a huge role in defect prediction. Their conclusion is that there is not a strong correlation between ownership metrics and module faults in the seven well known FLOSS systems. Apparently, there is an interesting relationship between contributors, bugs and bug prediction. Inspired by those papers we examine the effectiveness of process-metrics-based predictors in open source projects.

Erik Arisholm *et al.* [4] build predictive models to identify parts of a Java system with a high fault probability. During their research they found out that process metrics improve the performance of the prediction model compared with other models that do not contain process metrics. The process metrics that we are considering are drawn from their research.

Giger *et al.* [8] use the Gini coefficient to investigate the distribution of the changes made to the source code, applied to bug prediction. They find that if the largest amount of changes is carried out by a few developers, a lower amount of bugs can be expected. In our research we apply the Gini coefficient to the overall contributions to the project to measure the general contribution inequality.

Branch *et al.* [14] use the Gini coefficient to analyze software metrics. Their approach focuses on assessing a software product by calculating several metrics. By calculating the software metrics, they use a high order statistic, the Gini coefficient to observe any differences that occur in software systems and to assess the risks and monitor the development process itself. With their work they come to the conclusion that the Gini values are consistent as a project evolves over time. Gini coefficient is traditionally used as an economic metric to measure the distribution of wealth, but there are also research papers in software development domain that use it for many different purposes.

REFERENCES

- [1] F. Rahman and P. Devanbu, "How, and why, process metrics are better," in *Proceedings of the 2013 International Conference on Software Engineering*, ser. ICSE '13. IEEE Press, 2013, pp. 432–441. [Online]. Available: <http://dl.acm.org/citation.cfm?id=2486788.2486846>
- [2] C. Bird, N. Nagappan, B. Murphy, H. Gall, and P. Devanbu, "Don't touch my code!: Examining the effects of ownership on software quality," in *Proceedings of the 19th ACM SIGSOFT Symposium and the 13th European Conference on Foundations of Software Engineering*, ser. ESEC/FSE '11. ACM, 2011, pp. 4–14. [Online]. Available: <http://doi.acm.org/10.1145/2025113.2025119>
- [3] M. Foucault, J.-R. Falleri, and X. Blanc, "Code ownership in open-source software," in *Proceedings of the 18th International Conference on Evaluation and Assessment in Software Engineering*, ser. EASE '14. ACM, 2014, pp. 39:1–39:9. [Online]. Available: <http://doi.acm.org/10.1145/2601248.2601283>
- [4] E. Arisholm, L. C. Briand, and E. B. Johannessen, "A systematic and comprehensive investigation of methods to build and evaluate fault prediction models," *J. Syst. Softw.*, vol. 83, no. 1, pp. 2–17, Jan. 2010. [Online]. Available: <http://dx.doi.org/10.1016/j.jss.2009.06.055>
- [5] M. Shepperd, D. Bowes, and T. Hall, "Researcher bias: The use of machine learning in software defect prediction," *IEEE Transactions on Software Engineering*, vol. 40, no. 6, pp. 603–616, 2014.
- [6] W. Han, C.-H. Lung, and S. A. Ajila, "Empirical investigation of code and process metrics for defect prediction," in *Multimedia Big Data (BigMM), 2016 IEEE Second International Conference on*. IEEE, 2016, pp. 436–439.
- [7] A. Okutan and O. T. Yildiz, "Software defect prediction using bayesian networks," *Empirical Softw. Engg.*, vol. 19, no. 1, pp. 154–181, Feb. 2014. [Online]. Available: <http://dx.doi.org/10.1007/s10664-012-9218-8>
- [8] E. Giger, M. Pinzger, and H. Gall, "Using the gini coefficient for bug prediction in eclipse," in *Proceedings of the 12th International Workshop on Principles of Software Evolution and the 7th Annual ERCIM Workshop on Software Evolution*, ser. IWPSE-EVOL '11. ACM, 2011, pp. 51–55. [Online]. Available: <http://doi.acm.org/10.1145/2024445.2024455>
- [9] Z. He, F. Shu, Y. Yang, M. Li, and Q. Wang, "An investigation on the feasibility of cross-project defect prediction," *Automated Software Engineering*, vol. 19, no. 2, pp. 167–199, 2012.
- [10] A. E. Hassan, "Predicting faults using the complexity of code changes," in *Proceedings of the 31st International Conference on Software Engineering*, ser. ICSE '09. IEEE Computer Society, 2009, pp. 78–88. [Online]. Available: <http://dx.doi.org/10.1109/ICSE.2009.5070510>
- [11] D. Posnett, V. Filkov, and P. Devanbu, "Ecological inference in empirical software engineering," in *Proceedings of the 2011 26th IEEE/ACM International Conference on Automated Software Engineering*, ser. ASE '11. IEEE Computer Society, 2011, pp. 362–371. [Online]. Available: <http://dx.doi.org/10.1109/ASE.2011.6100074>
- [12] S. Kim, T. Zimmermann, E. J. Whitehead, Jr., and A. Zeller, "Predicting faults from cached history," in *Proceedings of the 1st India Software Engineering Conference*, ser. ISEC '08. ACM, 2008, pp. 15–16. [Online]. Available: <http://doi.acm.org/10.1145/1342211.1342216>
- [13] A. E. Hassan, "Predicting faults using the complexity of code changes," in *Proceedings of the 31st International Conference on Software Engineering*. IEEE Computer Society, 2009, pp. 78–88.
- [14] P. Branch, O. Nierstrasz, M. Lumpe, and R. Vasa, "Comparative analysis of evolving software systems using the gini coefficient," *2013 IEEE International Conference on Software Maintenance*, vol. 00, no. undefined, pp. 179–188, 2009.