

Recognizing recurrent development behaviours corresponding to Android OS release life-cycle

Pavel Senin

*CSDL laboratory, ICS Department
University of Hawaii at Manoa
Honolulu, Hawaii
senin@hawaii.edu*

Abstract—Android OS is an open-source Linux-based operating system for mobile devices developed by Open Handset Alliance led by Google. I attempted to apply a novel data mining technique based on SAX approximation and indexing of time-series with TF* IDF weights in order to discover recurrent behaviors within Android development process.

By mining of software process artifact trails corresponding to eight years of OMAP kernel development I was able to train a classifier which successfully recognizes pre- and post-release behaviors.

Keywords—Behaviors Detection, Empirical Study, Source Control System

I. INTRODUCTION

As many other large open-source projects, Android OS was in the development for many years. Android is “an open-source software stack for mobile phones and other devices”, <http://source.android.com/> which is based on the Linux 2.6 monolithic kernel. Android Inc., the small startup company, begun the development. In 2005 company was acquired by Google which formed the Open Handset Alliance - a consortium of 84 companies which announced the availability of the Android Software Development Kit (SDK) in November 2007. The Android OS code is open and released under the Apache License.

There are about two millions of change records registered in the Android SCM by more than eleven thousands of contributors within eight years span. Git is used as a version control system for Android and source code organized into more than two hundreds of sub-projects. These organized by the function (kernel, UI, mailing system, etc.) and underlying hardware (CPU type, bluetooth communication chip, etc.).

By the large body of previous research it is shown that the change metadata is a rich source of software process and developers social characteristics. Relevant to my research subject work showed a possibility to discover recurrent behaviors by use of Fourier Analysis of change events [1], in another work it was shown the relation between the activity time intervals and software product quality [2]. In this paper I extend the previous research by introducing a universal framework for the temporal partitioning of the software

change artifacts their complexity reduction and successive recurrent behaviors discovery.

A. Contribution

To the best of my knowledge, this work is the first attempt to study the applicability of symbolic aggregate approximation and term frequency inverse document frequency weight statistics to the mining of software process artifacts trails. This methodology enables one to significantly reduce the large complexity (dimensionality and noise) of temporal artifacts left by software development and transform real-valued streams into the symbolic representation. These reduced symbolic data streams can be indexed and reused for mining and knowledge discovery.

As an example of a possible data-mining workflow demonstrating the resolving power and correctness of the approach I present a case study of building a classifier for pre- and post-release behaviors.

II. MOTIVATION

Previously it was found that software development as any other human activity could be successfully partitioned by the time of the day reflecting our lifestyle and habits [3] [4]. However external constraints, such as employment and management constraints [5], software release cycle [6] and other found to be able to significantly alter natural activity patterns. Furthermore, within open-source projects with diverse development community scattered over the globe and often following undocumented development process, the natural human activity cycles are often discarded as well as the development and release cycles are significantly altered. Thus, the only feasible way to discover an open-source software process is to analyze its artifacts trails such as SCM logs, bug-and issue tracking systems and mailing lists archives. Essentially these trails are event-series where every time-stamped event has an attached set of metadata. However the complexity of this data (large code-base, number of committers with wide range of experience and commit habits, etc.), the precision of events recall (commit or bug report timestamp is not always a good

indicator of the event, SCM-introduced noise through merge commits etc.) impose a great challenge for the researchers.

These challenges are not new to the data-mining community and an enormous wealth of methods, algorithms and data structures exist for overcoming these issues. One of the solution to facilitate mining of the temporal data is to transformed it into other data representations such as spectral, polynomial, wavelets, piece-wise, symbolic, etc. In this paper, I investigate the application of Symbolic Aggregate Approximation [7] and the TF*IDF statistics [8] to the problem of discovering recurrent behaviors from software process artifacts with application to Android SCM data.

III. RESEARCH QUESTION

Previously a variety of time-series mining algorithms was applied to the problem of finding of trends, periodicity and recurrent behaviors within the software change artifacts trails such as Linear predictive coding and cepstrum coefficients [9], Fourier Transform [1] and coding [10]. These work proved that recurrent behaviors naturally occurring due to the variety of reasons are reflected in the SCM trails and could be recovered with application of data-mining techniques.

Recently Lin&Keogh [7] discovered and explored the universal applicability and resolving power of Symbolic Aggregate Approximation by application of the technique to a wide range of time-series data mining problems and in particular to motif and discords detection.

In this exploratory work I am investigating the applicability of this technique to the discovery of recurrent behaviors within SCM trails. Thus, the research questions I am resolving are:

- Which kinds of SCM data need to be collected for such analyzes?
- What is the optimal way of data representation and a data storage configuration?
- Which partitioning (slicing) is appropriate and which set of parameters one should use for SAX approximation?
- Which distance metrics serves best for measuring similarity and dissimilarity?
- What is the general mining workflow, and which parameters are crucial for result?

IV. EXPERIMENTAL METHODS

A. Data collection and organization

Two XML files were offered for the MSR challenge. These contain the most of the information obtainable from Google-hosted git source code repository as well as from Google-hosted bug and issue tracking system. While the issues and comments trail contain nearly complete information, the change trail provided for a contains only the high-level change information.

Table I
SLIDING WINDOW, PAA AND ALPHABET SIZE CHOICES.

Sliding window size	PAA size	Alphabet size
one week (7 days)	4	3
two weeks (14 days)	6	5
month (30 days)	10	5

The thirteen data fields of the change trail XML file provide information about the revision tree, author and a committer identification, change message and affected targets. Since I am focusing on the mining of temporal patterns for inferring recurrent behaviors, in addition to the existing data I collected auxiliary information about change. By creating a local mirror and by iterating over existing commits hashes I was able to recover auxiliary data for 68% of existing commits. The rest of changes which is about 32% of total change information belongs to legacy projects and is unrecoverable due to the changes in Android repository.

For every recoverable change record I collected a summary of added, modified or deleted files as well as a summary about LOC changes: added, modified or deleted lines. All this information was stored in the Trajectory database backend. Main tables of the database correspond to change and issue events; these tables accompanied with change targets tables and issue comments as well as tables for contributors. Overall, the database was normalized and optimized for the fast retrieval of change and issue information using SQL language.

Following the previous research I partitioned the change trails by the time of the day and by natural time periods of one week, two weeks and a month. These data was aggregated into the intermediate representation within MySQL database.

B. Symbolic approximation and indexing

SAX indexing depends on the three parameters which are required. First parameter is the sliding window, second is the PAA approximation size and third is the SAX Alphabet size.

In this work I have selected three sizes for sliding window: 7 days, 14 days and 30 days as these represent an intuitive and logical intervals of a week, two weeks and a month. For the PAA reduction I choose 4 steps for 7 days window, 6 steps for a bi-weekly interval, and 10 steps for a monthly window. Finally the alphabet, I choose 3 letters for weekly window, and 5 letters for bi-weekly and monthly windows. All these are summarized in the Table I.

By applying SAX to the pre-aggregated data I have obtained their symbolic temporal representation. Do I need a figure explaining SAX transform here?

C. Token-based distance metrics

For the experiments I have selected three similarity metrics. First one is based on the SAX min distance and

the Euclidean distance which applied to vector of tokens sorted by frequency observed in the aggregated symbolic stream corresponding to selected stream class, aggregation parameters, user, project and time-interval of the interest.

$$D(S, T) = \sqrt{\sum SAXdist(s_i, t_i)^2} \quad (1)$$

where $SAXdist$ is the the SAX distance based on the Normal alphabet.

As an alternative I have tried the Jaccard similarity between two sets S and T which is simply

$$J_\delta(S, T) = \frac{|S \cup T| - |S \cap T|}{|S \cup T|} \quad (2)$$

and the $TF * IDF$ similarity or which defined as a dot product

$$TFIDF(S, T) = \sum_{\omega \in S \cap T} V(\omega, S) \cdot V(\omega, T) \quad (3)$$

where

$$V(\omega, S) = \frac{V'(\omega, S)}{\sqrt{\sum_{\omega'} V'(\omega, S)^2}} \quad (4)$$

is a normalization of $TF * IDF$ (product of token frequency and inverse document frequency):

$$V'(\omega, S) = \log(TF_{\omega, S} + 1) \cdot \log(IDF_{\omega}) \quad (5)$$

where IDF_{ω} is a measure of the general importance of the pattern among all users

$$IDF_{\omega} = \frac{|D|}{DF(\omega)} \quad (6)$$

where $|D|$ is cardinality of D - the total number of users, and $DF(\omega)$ is the number of users having ω pattern in their activity set.

V. CLUSTERING

While the application of the SAX distance and Jaccard similarity yields a single number allowing the use of the convenient clustering libraries such as *hclust()* and *kmeans()* of R for manipulation with sparse vectors produced by the $TF * IDF$ application I used *sparcl()* R package for hierarchical clustering and my own implementation of *k - means* clustering.

VI. RESULTS

A. Kernel-OMAP life cycle patterns discovery

I arbitrary selected the Android kernel-OMAP as one of the large Android OS sub-projects. OMAP is a proprietary system on chips (SoCs) for portable and mobile multimedia applications based on general-purpose ARM architecture processor provided by Texas Instruments.

As a training set for building dictionaries of pre and post-release patterns I choose three Android releases: *Android 1.0*, *Android 1.5* and *Android 2.0*. For each release in the set I selected four weeks before the release as *pre-release* and

Table III
PRE- AND POST-RELEASE DEVELOPMENT PATTERNS CLASSIFICATION RESULTS FOR KERNEL-OMAP.

Release	Classification	Release	Classification
1.6 -pre	-	beta -pre	+
1.6 -post	+	beta -post	+
2.2 -pre	+	2.0.1 -pre	+
2.2 -post	+	2.0.1 -post	-
1.1 -pre	+	2.1 -pre	+
1.1 -post	+	2.1 -post	+
2.3 -pre	+	2.2.1 -pre	+
2.3 -post	+	2.2.1 -post	-

four weeks after release as *post-release* training intervals. As an additional constraint for this experiment I have selected contributors affiliated with *@android.com* e-mail domain and a telemetry stream of *added_lines*. After pre-generated patterns retrieval with SQL query I applied hierarchical clustering based on $TF * IDF$ similarity as a sanity test. The almost perfect clustering picture 1 indicates that there are significant differences in the pre-release and post-release weekly behaviors among selected contributors.

While hierarchical clustering is a good sanity test for the separation of data, the performance of K-means clustering is the most valuable metrics [11]. I performed k-means on the symbolic representation of data using $TF * IDF$ statistics and the Euclidean distance. Algorithm converged after two iterations separating pre- and post-release dictionaries with a single mismatch for the Android 2.0 pre-release.

Android kernel-OMAP hierarchical clustering stream ADDED_LINES, user mask ``*@google.com``

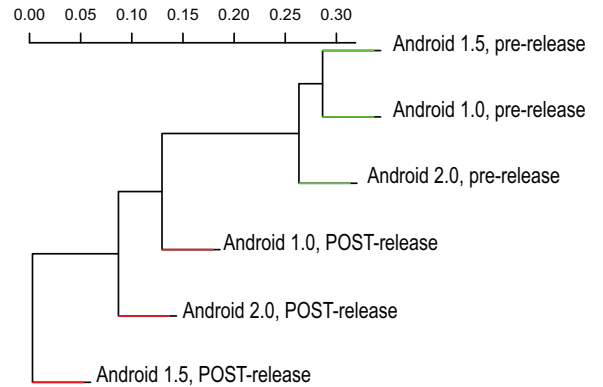


Figure 1. Hierarchical clustering of pre- and post- commit patterns.

By using centroids of resulting clusters as a basis for pre and post-release patterns I tested the classifier on the rest of Android releases for which kernel-OMAP remained an active project (some of the pre- and post- release interval contained non or only trivial patterns making them ineligible for classification). This classifier was able to successfully classify more than 81% of pre- and post-release behaviors (Table III).

Table II
WORDS AND THEIR $TF * IDF$ WEIGHTS.

release	"bbac"	"abca"	"babc"	"bbba"	"bcaa"	"bcbb"	"ccaa"	"cbaa"	"bbcb"	"bbbb"	"bbbc"
2.0-post	0.63	0	0.63	0	0	0	0	0.39	0.24	0.06	0
1.0-post	0	0.93	0	0	0	0	0	0	0	0.09	0.36
1.5-post	0	0	0	0	0	0	0	0	0.79	0.61	0
1.5-pre	0	0	0	0.23	0.23	0.91	0	0.14	0.18	0	0.09
2.0-pre	0	0	0	0	0	0	0	0	0	1	0
1.0-pre	0	0	0	0	0	0	0.79	0	0	0.08	0.61

VII. CONCLUSIONS

The presented approach and workflow is able to recognize recurrent behaviors from software process artifact trails. Its application to Android OS kernel-OMAP subproject SCM artifact trail yielded the pre- and post-release behavior models. In particular, it was shown that error rate of recognition of releases within the same as training project is less than 20%, f

VIII. ACKNOWLEDGEMENT

IX. CONCLUSION

omap-hclust.eps The conclusion goes here. this is more of the conclusion

ACKNOWLEDGMENT

The authors would like to thank... more thanks here

REFERENCES

- [1] A. Hindle, M. W. Godfrey, and R. C. Holt, "Mining recurrent activities: Fourier analysis of change events," in *Software Engineering - Companion Volume, 2009. ICSE-Companion 2009. 31st International Conference on*. IEEE, May 2009, pp. 295–298. [Online]. Available: <http://dx.doi.org/10.1109/ICSE-COMPANION.2009.5071005>
- [2] J. Eyolfson, L. Tan, and P. Lam, "Do time of day and developer experience affect commit bugginess?" in *Proceedings of the 8th Working Conference on Mining Software Repositories*, ser. MSR '11. New York, NY, USA: ACM, 2011, pp. 153–162. [Online]. Available: <http://dx.doi.org/10.1145/1985441.1985464>
- [3] N. E. Klepeis, W. C. Nelson, W. R. Ott, J. P. Robinson, A. M. Tsang, P. Switzer, J. V. Behar, S. C. Hern, and W. H. Engelmann, "The national human activity pattern survey (NHAPS): a resource for assessing exposure to environmental pollutants." *Journal of exposure analysis and environmental epidemiology*, vol. 11, no. 3, pp. 231–252, 2001. [Online]. Available: <http://dx.doi.org/10.1038/sj.jea.7500165>
- [4] I. Gorton and S. Motwani, "Issues in co-operative software engineering using globally distributed teams," *Information and Software Technology*, vol. 38, no. 10, pp. 647–655, Jan. 1996. [Online]. Available: [http://dx.doi.org/10.1016/0950-5849\(96\)01099-3](http://dx.doi.org/10.1016/0950-5849(96)01099-3)
- [5] E. Yourdon, *Death March (2nd Edition)*, 2nd ed. Prentice Hall, Nov. 2003. [Online]. Available: <http://www.worldcat.org/isbn/013143635X>
- [6] M. M. Lehman, "Programs, life cycles, and laws of software evolution," *Proceedings of the IEEE*, vol. 68, no. 9, pp. 1060–1076, 1980. [Online]. Available: <http://dx.doi.org/10.1109/PROC.1980.11805>
- [7] J. Lin, E. Keogh, L. Wei, and S. Lonardi, "Experiencing SAX: a novel symbolic representation of time series," *Data Mining and Knowledge Discovery*, vol. 15, no. 2, pp. 107–144, Oct. 2007. [Online]. Available: <http://dx.doi.org/10.1007/s10618-007-0064-z>
- [8] T. Roelleke and J. Wang, "TF-IDF uncovered: a study of theories and probabilities," in *Proceedings of the 31st annual international ACM SIGIR conference on Research and development in information retrieval*, ser. SIGIR '08. New York, NY, USA: ACM, 2008, pp. 435–442. [Online]. Available: <http://dx.doi.org/10.1145/1390334.1390409>
- [9] G. Antoniol, V. F. Rollo, and G. Venturi, "Linear predictive coding and cepstrum coefficients for mining time variant information from software repositories," in *Proceedings of the 2005 international workshop on Mining software repositories*, ser. MSR '05, vol. 30, no. 4. New York, NY, USA: ACM, 2005, pp. 1–5. [Online]. Available: <http://dx.doi.org/10.1145/1082983.1083156>
- [10] A. Hindle, M. W. Godfrey, and R. C. Holt, "Release Pattern Discovery via Partitioning: Methodology and Case Study," in *Proceedings of the 29th International Conference on Software Engineering Workshops*. Washington, DC, USA: IEEE Computer Society, 2007. [Online]. Available: <http://dx.doi.org/10.1109/ICSEW.2007.181>
- [11] *Initialization of Iterative Refinement Clustering Algorithms*, 1998. [Online]. Available: <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.54.3469>