

Empirical Software Engineering

Deriving a Usage-Independent Software Quality Metric

--Manuscript Draft--

Manuscript Number:	EMSE-D-19-00023	
Full Title:	Deriving a Usage-Independent Software Quality Metric	
Article Type:	SI: PROMISE 2018	
Keywords:	Software Quality; Software Usage; Software Faults; Bayesian Networks; NPM packages	
Corresponding Author:	Tapajit Dey University of Tennessee, Knoxville Knoxville, TN UNITED STATES	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:	University of Tennessee, Knoxville	
Corresponding Author's Secondary Institution:		
First Author:	Tapajit Dey	
First Author Secondary Information:		
Order of Authors:	Tapajit Dey	
	Audris Mockus	
Order of Authors Secondary Information:		
Funding Information:	National Science Foundation (1633437)	Dr. Audris Mockus
Abstract:	<p>Context: The extent of post-release use of software affects the number of faults, thus biasing quality metrics and adversely affecting associated decisions. The proprietary nature of usage data limited deeper exploration of this subject in the past.</p> <p>Objective: To determine how software faults and software use are related and how, based on that, an accurate quality measure can be designed.</p> <p>Method: Via Google Analytics we measure new users, usage intensity, usage frequency, exceptions, and release date and duration for complex proprietary mobile applications for Android and iOS. We utilize Linear Regression, Bayesian Network, and Random Forest models to explain the interrelationships and to derive the usage independent release quality measure. To increase external validity, we also investigate usage (downloads) and quality (number of issues) for 4430 popular NPM packages.</p> <p>Results: We found the number of new users to be the primary factor determining the number of exceptions, and found no direct link between the intensity and frequency of software usage and software faults. Crashes increased with the power of 1.02-1.04 of new user for the Android app and power of 1.6 for the iOS app. Release quality expressed as crashes per user was independent of other usage-related predictors, thus serving as a usage independent measure of software quality. Usage also affected quality in NPM, where downloads were strongly associated with numbers of issues. Unlike in mobile case where exceptions per user decrease over time, for 45.8% of the NPM packages the number of issues per download increase.</p> <p>Conclusions: We expect our result and our proposed quality measure will help gauge release quality of a software more accurately and inspire further research in this area.</p>	

Deriving a Usage-Independent Software Quality Metric

Tapajit Dey
Graduate Student
University of Tennessee, Knoxville
Min H. Kao Building, Room 619
1520 Middle Drive
Knoxville, Tennessee, USA - 37996
Email: tdey2@vols.utk.edu
ORCID: 0000-0002-1379-8539

Audris Mockus
Ericsson-Harlan Mills Chair Professor
University of Tennessee, Knoxville
Min H. Kao Building, Room 613
1520 Middle Drive
Knoxville, Tennessee, USA - 37996
Email: audris@utk.edu
Office Phone: 865-974-2265

Conflict of interest disclosure:

Funding: This work was supported by the National Science Foundation under Grant No. 1633437.

Deriving a Usage-Independent Software Quality Metric

Tapajit Dey · Audris Mockus

Received: date / Accepted: date

Abstract Context: The extent of post-release use of software affects the number of faults, thus biasing quality metrics and adversely affecting associated decisions. The proprietary nature of usage data limited deeper exploration of this subject in the past. Objective: To determine how software faults and software use are related and how, based on that, an accurate quality measure can be designed. Method: Via Google Analytics we measure new users, usage intensity, usage frequency, exceptions, and release date and duration for complex proprietary mobile applications for Android and iOS. We utilize Linear Regression, Bayesian Network, and Random Forest models to explain the interrelationships and to derive the usage independent release quality measure. To increase external validity, we also investigate usage (downloads) and quality (number of issues) for 4430 popular NPM packages. Results: We found the number of new users to be the primary factor determining the number of exceptions, and found no direct link between the intensity and frequency of software usage and software faults. Crashes increased with the power of 1.02-1.04 of new user for the Android app and power of 1.6 for the iOS app. Release quality expressed as crashes per user was independent of other usage-related predictors, thus serving as a usage independent measure of software quality. Usage also affected quality in NPM, where downloads were strongly associated with numbers of issues. Unlike in mobile case where exceptions per user decrease over time, for 45.8% of the NPM packages the number of issues per download increase. Conclusions: We expect our result and our proposed quality measure will help gauge release quality of a software more accurately and inspire further research in this area.

Tapajit Dey · Audris Mockus
Department of Electrical Engineering and Computer Science
University of Tennessee, Knoxville
Knoxville, Tennessee, USA
E-mail: tdey2@vols.utk.edu, E-mail: audris@utk.edu

Keywords Software Quality · Software Usage · Software Faults · Bayesian Networks · NPM packages

1 Introduction

Improving quality of software is one of the objectives of software engineering. “Software Quality” has been defined in various ways, but in this paper we take a narrow focus on the manifestation of software defects as crashes observed from the perspective of the users. Observing a software crash, generally speaking, is a manifestation of low quality of the software to a user. Thus, it seems intuitive to measure the quality of software¹ by counting the number of crashes, with more crashes being associated with lower quality. If, for example, we compare two different softwares or two different releases of a software, we first calculate the number of crashes for each and then compare these numbers. However, software with more users tends to see more crashes [11, 19, 37] as each user may exercise it differently. In an extreme case, a software or a release with no users will have no crashes, regardless of its quality. This interdependence of software usage volume and crashes experienced is typically not considered in quality measurement in industry or in empirical studies (although few studies do note that [14, 16]). Ignoring this relationship, however, would misguide quality improvement efforts (avoid quality improvements for releases/ software with low usage) and/or misguided developer performance metrics (reward developers of low-usage products). This analogy can also be extended for software defects (bugs), and by extension for issues raised against a software, since software crashes are manifestations of underlying defects and [24, 19] observed that the number of discovered software defects increases with the number of users, although the relationship between crashes and defects is not very well understood [16].

One possible reason for this oversight is the scarcity of reliable usage data. While the number of defects and crashes reported by users are carefully tracked by most large scale projects (*e.g.* Mozilla Firefox, Ubuntu etc.), tracking the variables related to usage, *e.g.* the number of users, intensity of usage etc. is almost impossible without a reliable monitoring system. Such a system is rarely used by open-source software and even many traditional software-as-a-product systems do not or can not have such capability. Moreover, even when such a dataset is available, it is almost always proprietary, so obtaining and sharing it, even for the software development teams in these proprietary projects, is difficult since the deployment is typically managed by a different team within the organization. Without such data, however, it becomes exceedingly difficult to interpret the quality of a software from the customer reported crashes/defects alone due to the interdependence of usage and crashes/defects [11, 19, 37].

We were able to obtain the usage data for some mobile applications developed by Avaya, *viz.* Avaya Communicator for Android (currently known

¹ in fact, we mean to measure one aspect of the quality of software

as Avaya Equinox®) and Avaya one-X®Mobile SIP for iOS. The usage data was obtained from Google Analytics. For analyzing the usage data for these applications, we used three usage related variables: number of users, usage intensity (average duration of software use per user), and usage frequency (average number of times the app was used by a user), along with two variables describing attributes of the particular release: release date and effective duration of the release, measured by how long the release continued to have new users, and looked at how these variables affect the number of exceptions *i.e.* application crashes.

After the usual data cleaning and variable construction stages, we first applied a linear regression (LR) model to identify the significant predictors for the number of exceptions. Then we used a Bayesian Network (BN) model to discover the interrelationship between the variables. The BN model was generated by using structure search algorithms, and the search method was chosen based on the result of a simulation study. We have presented the detailed result of the simulation study and hope that other practitioners willing to use BN structure search methods in their work might find it useful. Finally, we ran a random forest (RF) model to identify the importance of the variables for predicting the number of exceptions. All analyses in this study was done in R [50]. We found that the frequency and intensity of usage have little impact on the number of exceptions, but the number of users does have a significant impact. Thus, we establish the interdependence between the number of crashes and usage, and finally propose a quality metric that is independent of usage, which would enable us to compare the qualities of different softwares and/or different releases of a software more accurately.

In our previous work [11], we only analyzed the General Availability releases for Avaya Communicator for Android. Since all of the data was from Google Analytics, we had the same set of variables for all the Avaya softwares. We found similar results from that study as well, with the number of new users being the most important factor affecting the number of exceptions. Furthermore, we proposed a quality metric of average number of exceptions experienced by end users (so lower means better) and found it to be independent of other usage metrics.

In this study, we have added the analysis of the development version of Avaya Communicator for Android and the General Availability releases of Avaya one-X®Mobile SIP iOS Client. Although we collected data for several other apps, those were dropped due to having too few releases/ exceptions/ users to give a reliable result. We employed some new data correction steps for correcting the observed number of users and visits (Section 4). We used the same three modeling techniques and the same quality measure and found the result to be very similar for all cases considered. We also added a timeline showing how the perceived quality of the releases vary with time for different releases. Moreover, to increase the external validity of these findings, we also consider the relationship between downloads, a measure of usage, and number of issues, a measure similar to the number of crashes or bugs, for a rather different scenario of NPM packages. Node Package Manager (NPM) is

the package manager for node.js, an open-source, cross-platform JavaScript run-time environment. Since we do not have the number of crashes for these packages, we looked at the number of issues reported for these packages instead. We chose to look at 4430 NPM packages that had more than 10,000 monthly downloads and a GitHub page with issues enabled. From the analysis of NPM packages we found that for only 36 out of 4430 packages (0.8%) the number of daily downloads is not a significant predictor ($p\text{-value} > 0.5$) of the number of issues on that day.

In summary, the primary contributions include an observed strong relationship between the number of exceptions (crashes) and usage by analyzing two different mobile softwares (three different versions) and a quality metric that gives a more complete and reliable measure of quality. The wider implications of these findings suggest the potentially serious problems in existing quality metrics and predictions. Specifically, the organizational goals for a software project quality should take into account usage, and software defect predictors could be improved substantially if the population of users would be taken into account and could be predicted. Moreover, the analysis of the NPM packages established the extendibility of the the concept, thus opening the possibility of wider application of the approach. We have also presented the detailed result of the simulation study we conducted to choose the best performing BN structure search algorithms, which we believe will of of use to practitioners willing to use BN structure search methods in their work. To increase external validity, we investigated usage-quality relationship in a different context: the number of downloads and issues of 4430 NPM packages. We found the number of downloads to be a significant predictor of number of issues for most (99.2%) of them. Finally, we have presented a timeline showing how the perceived quality changes with time during the different releases for the mobile applications and during the life of the NPM packages, and discussed the trends we found.

The rest of the paper is organized as follows: In Section 2 we give some details about both the commercial software and NPM packages we analyzed. Details of the data collected for the study is discussed in Section 3. Section 4 lists the data preprocessing steps we followed. We discuss the simulation study we performed for choosing the best performing BN structure search method that was used in subsequent analysis in Section 5. In Section 6, we present the different modeling approaches used to explain the number of exceptions(crashes) for the commercial software we studied and the results of that study. A comparison of our result and already published results is presented in Section 7. The quality metric we used is discussed in Section 8. The result of the analysis of the NPM packages is presented in Section 9. We discuss the implications of our result in Section 10. In Section 11, we discuss various works in related topics. Finally, we discuss the limitations of our study in Section 12 and conclude the paper in Section 13.

2 Software Studied

For this study, we looked into two different types of software, which are vastly different in nature. The primary focus of the study is on the commercial software developed by Avaya for mobile applications, which are from the telecommunication domain. For external validation of the theory developed using this data, we looked into the NPM packages, which are open-source JavaScript packages used in web-development. In this section, we provide some details about the two to help contextualize our study.

2.1 The Mobile Applications developed by Avaya

One of the software chosen for this study was Avaya Communicator for Android (currently known as Avaya Equinox®). It integrates the mobile devices of the users with their office Avaya Aura®communications environment and delivers mobile voice and video VoIP calling, cellular call integration, rich conferencing, instant messaging, presence, visual voicemail, corporate directory access and enterprise call logs².

Another software we studied was the Avaya one-X®Mobile SIP for iOS, which provides mobile communications for the iPhone, iPod touch, and iPad through a wireless-enabled SIP Avaya Aura®environment combining enterprise features with the convenience of a mobile endpoint for users on the go. The Avaya one-X Mobile®SIP for iOS appears as an end point in the Aura®environment³.

Avaya is developing large, complex, real-time software systems that are embedded and standalone products. Development and testing are spread through 10 to 13 time zones in the North America, USA, Europe and Asia. R&D department employed many virtual collaboration tools such as JIRA, Git, WIKIs and Crucible. Development teams use Scrum-like development methodologies with a typical 4-week sprint. We consider a 15+ year old software component, the so-called Spark engine. As a software platform, Spark provides a consistent set of signaling platform functionalities to a variety of Avaya telephone product applications, including those of third parties. Spark is a client platform that provides signaling manager, session manager, media manager, audio manager, and video manager. The codebase involves more than 200K files and, over all forks, over 4M commits. The Android software chosen for this study is a fork of the Spark codebase. A more in-depth description of the development process is provided in [12].

2.2 NPM Packages

Node Package Manager or NPM is one of the most active and dynamic software ecosystems at present. It currently hosts more than 800,000 packages, and

² <https://support.avaya.com/products/P1574/avaya-equinox-for-android>

³ <https://support.avaya.com/products/P0949/avaya-onex-mobile-sip-for-ios>

have more than doubled in size in past couple of years (in January 2017, NPM reportedly hosted around 350,000 packages [1]). The popularity of NPM packages have, accordingly, skyrocketed as well. According to [61], “JavaScript is getting more popular all the time, and NPM is being adopted by an ever greater percentage of the JavaScript community.” About 75% of all JavaScript developers used NPM, with about 10 million users, in January 2018, according to [61]. Therefore, NPM is an excellent candidate for this study. Moreover, since they track the number of downloads of all packages in the ecosystem, which, in spite of essentially being a mix of downloads by users, bots, and mirror servers, as explained in [60], is the closest measure of usage we could find for open-source projects. We only looked at the NPM packages that have more than 10,000 downloads a month. According to [60], automated downloads are expected to be around 50 per day, or 1500 per month. Packages with over 10K downloads should, therefore, not be noticeably impacted by downloads by automated sources. We collected the issues associated with the projects from their individual GitHub repositories. We found 4430 projects which had more than 10,000 monthly downloads since January 2018 and also had public GitHub repositories with nonzero number of issues. We collected the number of downloads and the total number of issues for all these packages from 2015-03-01 to 2018-08-31. However, we did not conduct a release by release comparison for these packages, since the release durations vary by a lot for most packages. Since the recorded number of downloads is a mix of downloads by human and non-human users, a release by release comparison would not give a reliable picture of the effect of actual usage by human users on the number of issues. However, the number of downloads by bots are relatively stable and vary only with time [60], so controlling for the date would eliminate the spurious effects of downloads by bots. So, we decided to focus on the entire packages instead of releases of the packages, and measured the effect daily downloads have on number of issues of that package on that day after controlling for the calendar date. Although the number of recorded downloads is a mix of downloads by actual users and bots, mirror servers etc. [60] it is still a far better measure than, *e.g.* number of stars of a GitHub repository which was used in studies like [5] as a measure of popularity, which is little different from usage as we measure it.

3 Details of the Data

3.1 Data on mobile applications collected from Google Analytics

The post-deployment data for the mobile applications were obtained from the Google Analytics platform. Google Analytics is a web analytics service offered by Google that tracks and reports website traffic. It is now one of the most widely used web analytics services on the internet. In addition to traditional web applications it also allows tracking of mobile applications. To do that, the producer of a mobile application needs to set up an account and instrument

Table 1 Measures available in the Original Data

Application Release Version	No. of exceptions†
Operating System version in the user's device	Date of record entry
No. of fatal exceptions†	No. of new visits†
No. of visits†	Time on site†
Details on user's mobile device: brand, category(mobile or tablet) and model	No. of new users†
No. of total users†	Sessions per user†

their mobile application to send certain events to Google Analytics. Notably, it works for the mobile applications investigated in this study.

We collected data for a number of mobile applications developed by Avaya from Google Analytics, but some of the datasets turned out to be unusable for this study, for reasons ranging from very low volume of collected data (*e.g.* Avaya Communicator for Android - Experimental Releases) to zero recorded exceptions making an analysis impractical (*e.g.* Avaya One-X®ScsCommander). The following datasets were found usable:

- Avaya Communicator for Android - General Availability and Development versions.
- Avaya one-X®Mobile SIP for iOS - General Availability versions.

The data was collected between December 2013 and May 2016, although the exact time varies across the applications. Although we are primarily focused on the General Availability (GA) versions, since only these versions are available for end-users, we also decided to look into the development version for Avaya Communicator for Android, since we have detailed data available for these versions and we wanted to see if it shows a different characteristics from the GA versions.

The original data obtained from Google Analytics had measures for the variables listed in Table 1, aggregated at a per-day granularity, meaning that each entry in the original data table contained the measures for the numerical variables (marked with a † symbol in the table) for each unique combination of date, application release version, operating system version, mobile device brand, category, and model. We had the same set of variable for all the applications listed above.

It is important to note that Google Analytics releases only aggregate data even to developers of the application and limits the number of REST API calls, so one can not, for example, retrieve usage data for every calendar second or get exact time of the events. The daily counts split by release version of the application, OS version, and type of device, provided sufficiently fine granularity for our analysis.

3.2 Data on NPM Packages

We used the API provided by NPM for collecting daily downloads of the 4430 NPM packages we studied. (The API documentation is available in: <https://github.com/npm/registry/blob/master/docs/download-counts.md>).

To obtain the metadata information for every package in NPM, we wrote a “follower” script, as described in <https://github.com/npm/registry/blob/master/docs/follower.md>. The output contained the metadata information for all releases of all packages in NPM. From this we extracted the URL of GitHub repositories of the packages. Some NPM packages do not have a valid GitHub URL, so those were dropped from subsequent analysis. Using the Rest API provided by GitHub we collected information on the issues for all these NPM packages. Finally, we used the issue creation dates to construct a dataset of the total number of issues per day. We used the total number of issues instead of the number of open issues because we are interested in the number of issues encountered by the users of the packages. Whether an issue is resolved or not depends on a number of factors, *e.g.* the number of developers, the responsiveness of the developers, the number of packages managed by each developer, the complexity of the problem; most of which are unrelated to usage, so we decided using the total number of issues is a much more reasonable option.

Our final dataset for the NPM packages contained the number of recorded downloads per day and the total number of issues reported for that package until that day. The date range was from 2015-03-01 to 2018-08-31.

4 Data Preprocessing

This section contains the data cleaning, transformation, and variable construction steps undertaken prior to the application of the different modeling methods. The NPM data was relatively clean, so we did not employ any pre-processing techniques on the collected data for the NPM packages.

4.1 The Google Analytics data for the Mobile Applications

Removal of variables before aggregation: Upon initial investigation into the data, we found that no. of exceptions and no. of fatal exceptions were exactly the same, as recorded by Google Analytics, so we removed the no. of fatal exceptions from the dataset. Only fatal exceptions were recorded for this application, *i.e.*, crashes that require a complete restart of the mobile application and, potentially, may affect the operating system itself. This is not surprising since the bulk of the functionality for the application was written in C++ and called from the mobile OS via Native Interface.

We did not consider the variables related to mobile device details and operating system versions because the application, as noted above, was primarily written in C++ and the user interface aspects that vary greatest among devices and versions of OS were not likely to have influence. To validate that assumption we investigated and found no correlation of exceptions with either variable.

Data correction: This additional preprocessing step was not a part of our previous study [11]. We found during careful inspection of the data that some

of the releases had non-zero number of users but zero new users in the dataset. This obviously hints at some part of the data being missing. So, as a data correction step, we modified the number of new users so that in the chronological order of the data, the cumulative number of new users is never less than the number of users for a day.

Aggregating data to per-release granularity: We had some missing values in the data, however, most of the missing data was about the mobile devices and since we didn't use them in our analysis, we got rid of that problem by simply dropping the variables. Since our aim is to model the quality of the different releases, we aggregated the data to a per-release granularity, from the the original data that was recorded in per-day granularity. The raw data contained 177 different GA releases and 25 development releases for the Avaya communicator for Android and 11 GA releases for the Avaya mobile SIP for iOS. We dropped 4 GA releases for the Avaya communicator for Android from further consideration because a significant portion of observations were missing. The result of aggregation, however, was two new variables: start date (first day for which we have a record for that release) of a release, and end date (last date for which we have a record for that release) of a release, which in turn helped create another variable: duration of a release. We did not to keep the end date in the final table, since duration and start date can be used to compute the end date.

Verifying the correctness of Release date: The original data involves only the usage aspects and the version information of the software. The project under consideration was relatively new and it was one of the early attempts for the team to deploy mobile software on Android and iOS. As such, not everything was well documented and also was rapidly evolving over time and no record of the exact release dates for most of the releases was available. We did manage to get release dates for some of releases from Google Play Store/ Apple App Store, but not all the release dates were available. For the releases with dates available on Google Play Store/ Apple App Store, the official release dates from Avaya records, and the start dates obtained from the data were either very close or exactly the same, so we do not have a reason to doubt the dates obtained from the data.

Removal of variables post aggregation: The numerical variables were aggregated to give a sum for each variable. Upon further inspection, we found the number of users, new users, visits, and new visits to be highly correlated. In the second iteration, we removed the variable "sessions per user", because aggregating it directly is meaningless, and we were not sure how it was originally calculated by Google Analytics (was it a mean or a median? were new users or total users counted?). We also removed the "total users" and "total visits", because while summing up the new users/visits for each day gives an accurate measurement of the total number of new users/visits for a release, it is not guaranteed that summing up total users/visits does the same due to possible double counting the number of users/visits.

Final list of variables: Keeping the goal of our study in mind, the variables we have after the initial cleaning steps give us necessary information for

Table 2 Measures in the Aggregated Data Table

<i>Release variable</i> - Start Date for the release (Release.Date)	<i>Release variable</i> - Effective Duration of the release (Release.Duration)
<i>Post-Release defects</i> - Total No. of exceptions (Exceptions)	<i>Usage variable</i> - Average time on site per user (Usage.Intensity)
<i>Usage variable</i> - Total number of new users (New.Users)	<i>Usage variable</i> - No. of visits per user (Usage.Frequency)

a model of post-release defects and software usage. In our list of variables, we have the total number of exceptions *i.e.* post-release defects. As for measures related to software usage, we have the total number of new users; the “Time.On.Site” variable, normalized by the number of users of a release, provides a measure for the temporal intensity of usage per user; and the number of visits per user is a measure for the frequency of usage. We also have two variables related to each individual release: the start date *i.e.* the release date gives a measure for the calendar time of each release, and is useful in gaining insight about if the number of post-release defects and software usage vary with time, and the duration of a release, which could have an effect on the number of exceptions and the number of new users, since these variables were not normalized with duration. Since we only have a limited amount of data, we restricted ourselves to use only these six variables. Our final aggregated data table had the measures listed in Table 2, with the corresponding variable names we used in the model enclosed in brackets.

Log-transformation of variables: The release date was converted from the Date format to numeric format, which resulted in the values for the release date variable being represented by the difference in days from Unix time (counted from 1970-01-01). We found that all of the variables under consideration had a long-tailed distribution, so we took logarithm of them. The distribution of the variables of GA releases of Avaya communicator for Android is shown in Figure 1. The distribution of the variables of other applications is available in our GitHub repository: https://github.com/tapjday/release_qual_model.

5 Simulation Study for Selecting the Best BN Structure Search Method

5.1 Basics of Bayesian Network models

Bayesian Network [30,55] is a type of Probabilistic Graphical Model (PGM), which explicitly represents the conditional dependency/independence as a directed acyclic graph where variables represent nodes and dependencies represent links, and thus this representation can be used as a generative model⁴.

⁴ A generative model specifies a joint probability distribution over all observed variables, whereas a discriminative model (like the ones obtained from regression or decision trees) provides a model only for the target variable(s) conditional on the predictor variables. Thus, while a discriminative model allows only sampling of the target variables conditional on the predictors, a generative model can be used, for example, to simulate (*i.e.* generate) values

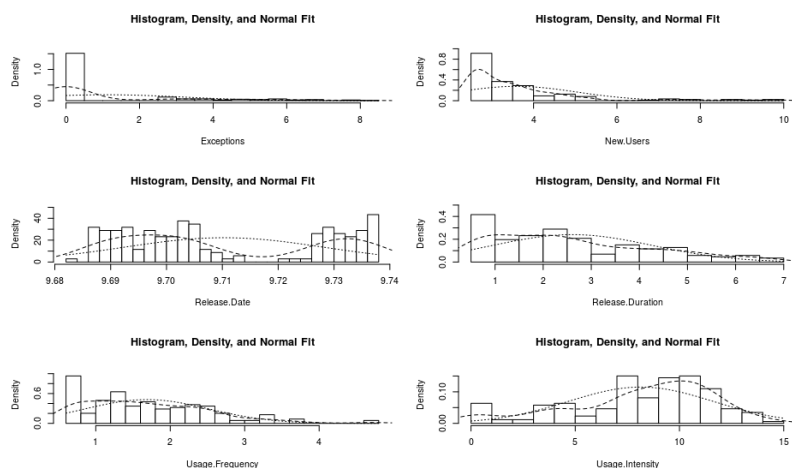


Fig. 1 Distribution of the variables after transformation: GA releases of Avaya communicator for Android

Bayesian Networks models can be useful in the context of Software Engineering research [16] due to having several advantages over regression models. To be precise, regression analysis is a very simple BN where there is one directed link from each independent variable to dependent variable. BNs, therefore, can help with multicollinearity, a common problem with software engineering data [63,59,7,33], that is present in our data as well, by linking independent variables.

Another variety of PGM that we did not use in this paper (details in Section 12) is the Markov random fields that represent the interrelationships between variables as undirected graphs. They differ in the set of independencies they can encode and the factorization of the distribution that they induce [30].

Bayesian Network Model construction: Despite the promises of BNs, they tend to be quite sensitive to data, and operational data, is often problematic [35,66]. Careful preprocessing, therefore, is needed to ensure a reliable and reproducible result. Two primary ways to use BNs exist. With the first approach the graph represents dependencies obtained from domain experts. The graph may include prior distributions about the parameters of the overall model. The data is then used to calculate the posterior distribution and to make inference. The second approach puts minimal a-priori assumptions about the model and focuses on the search for the best graphical representation for a given dataset (structure learning). This is an NP-hard problem [8], but a number of different heuristic structure learning algorithms are available. Due to the lack of any strong theory connecting the variables we are considering, we decided to use the structure search method for BN model construction. Since

of any variable in the model, and consequently, to gain an understanding of the underlying mechanics of a system, generative models are essential.

our goal is to find a Bayesian network model for the data, we didn't examine the methods that do not result in a Directed Acyclic Graph (DAG). We found that the *bnlearn* package in R implements a wide range of BN searching methods for continuous, discrete, or a mixed set of variables and the corresponding families of scoring functions and also has a good number of examples. These methods were also shown to be able to recover the underlying network for a protein-signaling-chain (in Biology) in [54]. We, therefore, use this package for our analysis. In addition to the methods implemented in *bnlearn* package, we investigated some methods from a few other packages which can be interfaced with the *bnlearn* package.

Due to the potential inconsistencies of the BN models, we performed our modeling in two stages. First, we considered all available BN structure methods in the *bnlearn* package and ran a simulation based study to find the methods that are most accurate and then we used those methods on our data to create the final model.

Methods considered:

The different BN structure search methods we considered are listed below:

- *Greedy Hill-Climbing search algorithms*(HC) [43,54]
- *Hybrid algorithms*(Hybrid) [43,54]
- *Posterior maximization* using *deal* package in R [54,6] .
- *Simulated Annealing* using *catnet* package in R [3,54].
- *MAP (maximum a posteriori estimation) Bayesian Model Averaging* (MAP) [43,54]

All structure search algorithms try to maximize some form of a network score. Among the various scores available, BIC score is the suitable one when the goal is to create an explanatory model from non-informative prior models [56,57]. BIC score is used for discrete data while the Gaussian equivalent of BIC (bic-g) score is used for continuous data.

The results, *i.e.* the structure and the parameters resulting from a structure search algorithm, are often noisy, meaning that different settings induce slightly different networks. To mitigate this effect we use non-parametric bootstrap model averaging method described in [17], which provides confidence level for both the existence of edge and its direction. This enables us to select a model based a confidence threshold. Authors of [17] argue that threshold is domain specific and needs to be determined for each domain. For instance, a threshold of 0.95 indicates that only the edges that appeared in more than 95% of the bootstrap optimized models were selected.

Many applications of BNs discretize the data prior to applying the structure learning methods, so we considered it as a possibility as well. Using continuous data works best when the random variables (possibly after a transformation) have Gaussian distribution. While using discrete data does not require such assumptions, obtaining the optimal discretization for a dataset is in itself an NP-hard problem [9]. Choosing a sub-optimal discretization technique may result in spurious or missed relationships, which can in turn result in

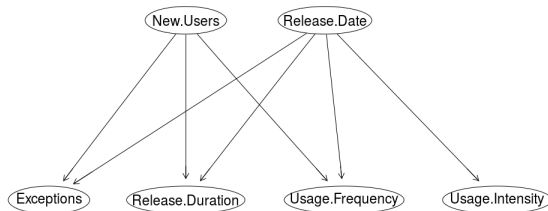


Fig. 2 Custom model used for Simulation Study

incorrect dependencies being reported in the resulting model. Given the pros and cons of both types of methods, we use methods of both types for our simulation study. As we are interested in creating a generative model, we had to use a discretization method that is unsupervised. The basic problem with commonly used supervised methods (*e.g.* Chi-square, or MDLP discretization algorithms) is that they optimize discretization to improve explanatory power for a single response variable. This is not suitable for a BN structure search, because we do not know which variables will be responses (have arrows pointing to them) and which will be independent (have no incoming arrows) *a-priori*. While some research on multidimensional discretization methods exists [49], we are not aware of any that have a robust implementation.

Simulation Study:

We performed the simulation study by first creating a random BN (see Figure 2) with six nodes, since we also have six variables in our final list (Table 2). For demonstration purposes we use the same variable names. We fitted this graph with our data on GA releases of Avaya communicator for Android (log-transformed and scaled) to generate values for the coefficients for each edge. This model was used in our simulation study going forward. We created 1000 different datasets from the BN structure in Figure 2, and applied the different structure search algorithms (both continuous and discrete versions, where available) listed above. Our performance metric is finding how many times the different algorithms can recover the underlying structure from the simulated data.

Other than testing the methods themselves, we also tested whether or not we should discretize the data. We tried different discretization methods, *viz.* equal interval, equal frequency, and k-means clustering based discretization methods from the *arules* package [21], and the Hartemink⁵ discretization methods in the *bnlearn* package.

Except for the *Posterior maximization* using *deal* package, which can't be bootstrapped, all other results were bootstrapped, so we tested different thresholds in our simulation study as well. Finally, for the the Hybrid search algorithm, in which conditional independence tests are performed to restrict the search space for a subsequent greedy search, there are many re-

⁵ Hartemink's pairwise mutual information method[22].

Table 3 Result of Simulation Study

Method	Exact	Off-by-one
HC	0.574	0.264
MAP	0.596	0.214
Hybrid- si.hiton.pc	0.000	0.019
Hybrid- mmpc	0.000	0.016
Hybrid- gs	0.000	0.011
HC-D-F	0.000	0.010
Hybrid- iamb	0.000	0.010
Hybrid- mmpc -D-H	0.000	0.008
Hybrid- si.hiton.pc -D-H	0.000	0.008
HC-D-H	0.000	0.007
Hybrid- mmpc -D-F	0.000	0.007
Hybrid- si.hiton.pc -D-F	0.000	0.006
Hybrid- iamb -D-F	0.000	0.005
Hybrid- gs -D-F	0.000	0.004
Hybrid- gs -D-H	0.000	0.004
Hybrid- iamb -D-H	0.000	0.002

Table 4 Result of Simulation Study: Different Thresholds

Method	Threshold	Exact	Off-by-one
MAP	0.85	0.68	0.25
MAP	0.80	0.67	0.25
MAP	0.90	0.67	0.26
MAP	0.95	0.66	0.27
MAP	1.00	0.66	0.27
MAP	0.75	0.66	0.21
HC	0.65	0.63	0.23
HC	0.70	0.63	0.23
HC	0.75	0.63	0.23
HC	0.80	0.63	0.23
HC	0.85	0.62	0.24
HC	0.55	0.62	0.23
HC	0.60	0.62	0.23
MAP	0.70	0.62	0.21
HC	0.90	0.60	0.26
MAP	0.65	0.58	0.17
HC	0.95	0.57	0.29
MAP	0.60	0.43	0.14
MAP	0.55	0.33	0.11
HC	1.00	0.19	0.47

strict methods available, *viz.* gs” (Grow-Shrink), ”iamb” (IAMB), ”fast.iamb” (Fast-IAMB), ”inter.iamb” (Inter-IAMB), ”mmpc” (Max-Min Parent Children), ”si.hiton.pc” (Semi- Interleaved HITON-PC), ”chow.liu” (Chow-Liu), ”aracne” (ARACNE) [32], and we tested all of these restrict options in our simulation study.

The result of the simulation study is shown in Table 3, which shows the fraction of times exact structures and off-by-one structures⁶ were generated by each method in the simulation. The result varies with the chosen threshold, so in Table 3, we show the overall performance of the different methods which generated an exact or off-by-one structure at least once in the simulation. For the hybrid search methods, we list mention the restrict option that was used, and the ‘-D’ suffix indicates a discretization method was used to discretize the data prior to applying a structure search method. ‘-D-H’ indicates Hartemink discretization method and ‘-D-F’ indicates Equal-Frequency discretization method. It is clear from the table that only HC and MAP methods can effectively reproduce the correct underlying structure around half of the times and they create more off-by-one structures than others, indicating the error rate is the lowest for these methods.

In Table 4, we show the fraction of times exact and off-by-one models were generated by HC and MAP methods, which performed the best among the methods considered, for different thresholds. It can be seen that using a moderately high threshold between 0.75 and 0.9 gives good results for both HC and MAP, while higher thresholds for HC and lower thresholds for MAP give worse results. Using the optimal threshold creates models that have more than one wrong and/or missing edge only 7-14% of the times.

The result of the simulation study had the following findings:

⁶ one extra / missing / reversed edge

- Using structure search algorithms on the continuous data resulted in much more frequent recovery of the original BN structure compared to discretized data.
- Bootstrapping improves the stability of the results considerably.
- The bootstrapped Hill-Climbing search and MAP Bayesian Model Averaging algorithms outperformed all others both in terms of accuracy and runtime, being able to recover the underlying structure more than 63% of the times and making no more than one error 86% of times with optimal thresholds.

We consider this study one of the contributions of the paper, and hope that it would be useful for researchers using BN structure learning techniques.

6 Explaining Exceptions for the Commercial Software

As mentioned earlier, we conducted our analysis in three stages: first, we used linear regression (LR) on the data with the number of exceptions as the response variable; then, we used Bayesian Network (BN) modeling approach to identify the interrelationship between the variables; and finally, we used a random forest (RF) model to verify the results.

We chose LR for the simplicity, robustness and ease of interpretation. To better understand interrelationships among variable (since LR is not applicable for sets of highly correlated predictors) we used BN models. Finally, to establish the predictive capabilities of our models we used RF, which is known as one of the best Machine Learning classifiers. That way we could both obtain the most insight and also to validate our findings through the use of radically different approaches.

6.1 Linear Regression Model

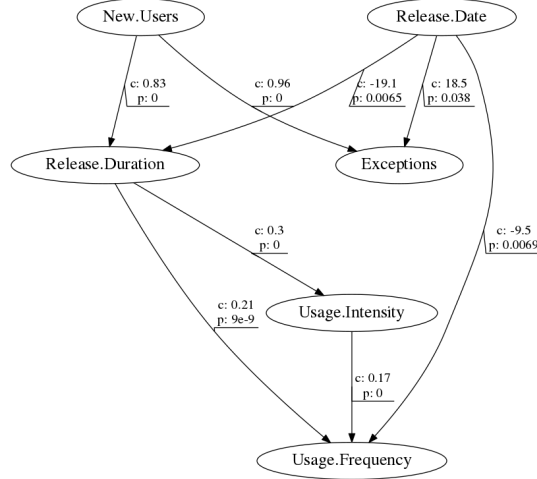
We first used a linear regression model to discover the significant variables affecting the number of exceptions. The output of the fitted model is shown in Table 5. The model resulted in a decent fit for the GA releases of Avaya Communicator for Android, given the sample size of 173, with adjusted R^2 of 0.481 (which is higher than 0.435, the adjusted R^2 reported with our previous approach [11]), and the variable “New.Users” was the most significant predictor while “Release.Date”, “Usage Intensity”, and “Usage Frequency” were also statistically significant. For the development releases of the same, only “New.Users” was significant, and it resulted in a good fit, with adjusted R^2 of 0.781. Finally, for the GA releases of Avaya mobile SIP for iOS, the value of R^2 was 0.874, and “New.Users”, “Release.Date”, and “Release.Duration” were the significant variables.

6.2 Bayesian Network Model

One key assumption for applying the continuous BN structure search algorithms is that the variables have a distribution close to a Gaussian distribu-

Table 5 Summary Result of LR model for “Exceptions”

	GA releases of Android			Dev releases of Android			Mobile SIP for iOS		
	Estimate	Std. Error	p-value	Estimate	Std. Error	p-value	Estimate	Std. Error	p-value
(Intercept)	-145.7827	72.2447	0.0452	-514.753	411.865	0.227	998.2528	284.0527	0.0170
New.Users	0.8172	0.1229	0.0000	1.08755	0.155	1.11e-06	1.2766	0.1969	0.0013
Release.Date	14.7756	7.4449	0.0488	52.884	42.383	0.227	-103.8231	29.4034	0.0167
Release.Duration	0.1361	0.1290	0.2930	-0.031	0.297	0.917	-2.0129	0.3381	0.0019
Usage.Frequency	-0.5388	0.2179	0.0144	-0.321	0.526	0.550	0.8246	0.9057	0.4043
Usage.Intensity	0.1261	0.0615	0.0419	0.025	0.183	0.892	1.0318	0.5286	0.1084

**Fig. 3** Final BN Model for GA releases of Avaya Communicator for Android (with c: coefficients after fitting the transformed, but unscaled data, p: p-value for the link)

tion. To satisfy this modeling assumption, we scaled all the variables to unit scale. The variable “Exceptions” still had a long tailed distribution, but the distributions of the other variables were much closer to normal distribution.

According to the result of the simulation study, we decided to use bootstrapped hill-climbing search and MAP Bayesian model averaging methods for constructing the Final BN models for our datasets and considered the model that resulted from both the methods. The resultant BN model for the GA releases of Avaya Communicator for Android is shown in Figure 3, which shows “New.Users” and “Release.Date” are parent nodes of “Exceptions”. Figure 4 shows the final BN Model for Development releases of Avaya Communicator for Android, in which only “New.Users” is the parent of “Exceptions”, and Figure 5 shows the final BN Model for GA releases of Avaya mobile SIP for iOS, where once again “New.Users” and “Release.Date” are parent nodes of “Exceptions”.

Every bootstrap run was performed over 500 bootstrap samples, and a hill-climbing search with 100 random restarts was applied on each sample to find the best fitting network, so in essence, each resultant network was obtained by averaging 50,000 candidate networks. We used a Threshold of 0.85, as it seemed optimal from our simulation study.

The result form a bootstrap run shows the relative strength of the link and the relative confidence for the direction of the link. In Table 6 we have shown

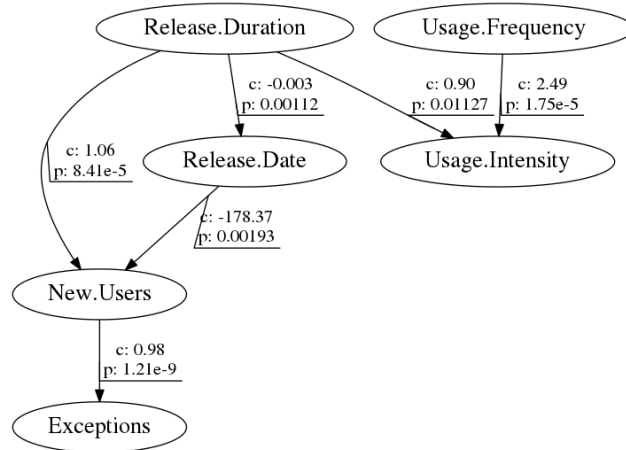


Fig. 4 Final BN Model for Development releases of Avaya Communicator for Android (with c: coefficients after fitting the transformed, but unscaled data, p: p-value for the link)

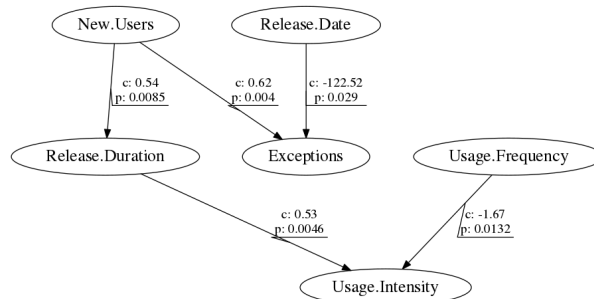


Fig. 5 Final BN Model for GA releases of Avaya mobile SIP for iOS (with c: coefficients after fitting the transformed, but unscaled data, p: p-value for the link)

the result from one bootstrap run of the HC method for all possible edges for the GA release data of Avaya Communicator for Android. If an edge has $< 50\%$ confidence in its direction, then the edge appears in the opposite direction in our model. Although Bayesian Networks are sometimes interpreted as causal relationships [47], there are disagreements on how that should be done. We, therefore, are not interpreting these relationships as causal here. All observed links, therefore, indicate the presence of observed correlation (and are empirical in nature) and the direction is a property of the topological ordering of nodes in a DAG, and affects the total probability distribution of the variables.

The BN models were fitted to the unscaled data, and the resulting coefficient of each link is also shown in the figures. The p-value for each link was calculated from a linear model with the source nodes as predictors and the destination node as the response variable, e.g. the p-value for the link from “New.Users” to “Exceptions” was calculated by looking at the result of: $\text{lm}(\text{Exceptions} \sim \text{New.Users} + \text{Release.Date})$.

Table 6 Example bootstrap result - GA releases of Avaya Communicator for Android

from	to	strength	direction
Exceptions	New.Users	1.00	0.34
Exceptions	Release.Date	0.86	0.47
Exceptions	Release.Duration	0.46	0.50
Exceptions	Usage.Frequency	0.75	0.78
Exceptions	Usage.Intensity	0.35	0.47
New.Users	Exceptions	1.00	0.66
New.Users	Release.Date	0.20	0.62
New.Users	Release.Duration	1.00	0.71
New.Users	Usage.Frequency	0.71	0.85
New.Users	Usage.Intensity	0.34	0.64
Release.Date	Exceptions	0.86	0.53
Release.Date	New.Users	0.20	0.38
Release.Date	Release.Duration	1.00	0.63
Release.Date	Usage.Frequency	0.97	0.82
Release.Date	Usage.Intensity	0.66	0.77
Release.Duration	Exceptions	0.46	0.50
Release.Duration	New.Users	1.00	0.29
Release.Duration	Release.Date	1.00	0.37
Release.Duration	Usage.Frequency	0.90	0.55
Release.Duration	Usage.Intensity	1.00	0.53
Usage.Frequency	Exceptions	0.75	0.22
Usage.Frequency	New.Users	0.71	0.15
Usage.Frequency	Release.Date	0.97	0.18
Usage.Frequency	Release.Duration	0.90	0.45
Usage.Frequency	Usage.Intensity	1.00	0.22
Usage.Intensity	Exceptions	0.35	0.53
Usage.Intensity	New.Users	0.34	0.36
Usage.Intensity	Release.Date	0.66	0.23
Usage.Intensity	Release.Duration	1.00	0.47
Usage.Intensity	Usage.Frequency	1.00	0.78

We fitted the model to the transformed, but unscaled data (for easier interpretation of results).

By looking at the p-values for the links, we can say that all the links in the BN models are statistically significant. Links having a negative coefficient indicate an inverse relationship between the parent and the child node. The performance of explanatory models is evaluated by the fraction of deviance explained by the model. Our model explains 45.5% and 34.7% of the variation in “Exceptions” (adjusted R^2 value of the model) for GA and development releases for Avaya Communicator for Android respectively and 20% for GA releases of Avaya mobile SIP for iOS.

6.3 Random Forest Model

As a verification step to identify the important variables affecting the number of exceptions, we used a Random Forest model to fit the data, with “Exceptions” as the response variable. The variable importance plot for the GA release data of Avaya Communicator for Android, as shown in Figure 6, indicates that “Release.Date” and “New.Users” are the two most important variables. We ran a 10-fold cross-validation exercise with “Exceptions” as the response variable, the resultant R^2 varied between 0.1176 and 0.6455 (mean: 0.4557, standard deviation: 0.1719). The high standard deviation is likely caused by

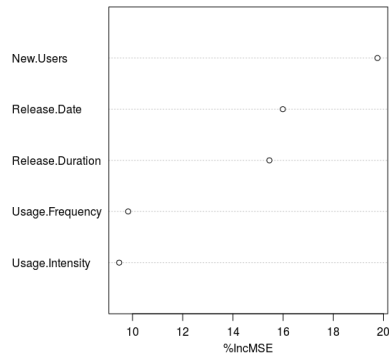


Fig. 6 Variable Importance Plot of RF model for “Exceptions” for GA release data of Avaya Communicator for Android

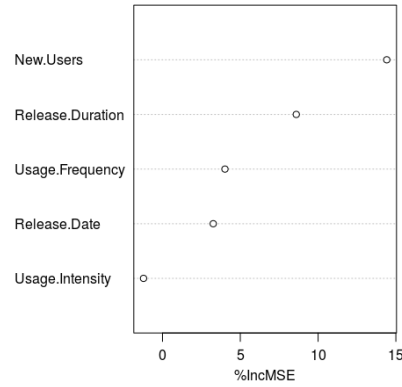


Fig. 7 Variable Importance Plot of RF model for “Exceptions” for development release data of Avaya Communicator for Android

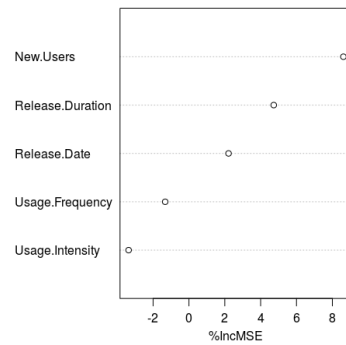


Fig. 8 Variable Importance Plot of RF model for “Exceptions” for GA releases of Avaya mobile SIP for iOS

our relatively small sample size of 173. This result reinforces the result we got from earlier analyses.

For the development releases of Avaya Communicator for Android, the variable importance plot is shown in Figure 7. “New.Users” is again the most important variable, followed by “Release.Duration”. For 3 out of the 10 folds, the R^2 value was negative, which basically means the fitted value had a different trend than the data [42]. It varied between -3.15 and 0.89 (median: 0.06, mean: -0.22, standard deviation: 1.19)

For the GA releases of Avaya mobile SIP for iOS, the random forest model resulted in very poor prediction performance, because we had a sample size of only 11. However, the variable importance plot still shows the number of new users is the most important variable, as can be seen from Figure 8.

7 Comparison with published results

In this section we compare our findings with already reported results that studied other commercial applications. The number of users for most of the releases we studied are very small, with a median of 7 users per release, although a few releases have more than 16,000 users. On slide 22 of his presentation [24], Caper Jones reported that the number of defects increase 2 to 3 times for a 10 fold increase in the number of users (from 1 to 10 and 10 to 100) for a software of similar complexity (between 10,000 and 100,000 function points). However, they were looking at the number of defects, and typically the number of exceptions is larger than the number of defects, because one defect could cause crashes for multiple users (or multiple crashes for a single user). The study published in [37] was done for a system with many more users (around 4,000 to 16,000), however, they reported that for a two-fold increase in the number of users the number of Modification Requests (MR tickets) increase around 1.25 times, which is more than what would have been predicted by our model (1.02 – 1.04) for Android apps, but less than what we have (1.6) for the iOS app. Although we were unable to do a direct comparison to another mobile application, these findings add more context to our result, and indicates the necessity of further studies that publish their datasets to understand the usage-fault relationship in a wider range of applications.

8 Application of our model: A derived measure of Quality

In order to arrive at the usage independent quality measure, we follow the framework of establishing laws governing relationships among measures of software development proposed in [34]. Law is an equivalent of invariance, i.e. a function of measures that is constant under certain conditions. In this case we want it to be constant for releases that have the same quality. First, the law requires a plausible mechanism and second, an empirical validation. Each new user may have a different type of phone, operating system, service provider, geographic region, and usage pattern. It is reasonable to assume that some of these configurations lead to software malfunction manifested as an exception. This provides us with a plausible mechanism on how precisely more new users of one release might generate more exceptions even if we have two releases of identical quality. To obtain empirical validation of this postulated mechanistic relationship we rely on our models, all of which show the number of software exceptions to be dependent on the number of users and on software release date. Therefore, we arrive at the following software law that is applicable for the investigated context: the average number of exceptions experienced by each user should, therefore, be independent of usage and depend only on the qualities of a software release.

In this section we test the above evidence-based hypothesis and provide the result of an analysis with the **number of exceptions per user** as a response variable (“Quality”) representing software quality. *This is actually a measure*

for faultiness, so a lower value of “Quality” indicates the actual quality of the software perceived by end users is better.

The value of the “Quality” variable (not log transformed) was seen to be varying between 0 and 10.85 (mean: 0.45, median: 0, standard deviation: 1.48) for the GA release data of Avaya Communicator for Android, between 0 and 22.83 (mean: 1.12, median: 0, standard deviation: 4.55) for development versions of the same, and for the GA releases of Avaya mobile SIP for iOS it varied between 0 and 0.5 (mean: 0.0488, standard deviation: 0.15) .

Similar to the previous analysis, we applied Linear Regression, Bayesian Network search, and Random Forest modeling approaches on the dataset containing this quality measure and the remaining variables, all of which were log-transformed.

The result, as expected, shows that the quality of a software, measured by average number of faults experienced by each user, has no dependence on other usage variables. The LR model of GA releases of Avaya Communicator for Android (Table 7) suggest that other variables have some effect on the quality variable. The BN model(Figure 9), obtained with a threshold of 0.85 from a bootstrapped Hill-Climbing structure search model, indicates the “Quality” variable depends only on the “Release.Date” variable. Finally, the result of 10-fold cross-validation with the RF model (Variable Importance plot in Figure 12) indicates that the “Release.Date” variable is much more important compared to others, and the two usage related variables are of much lower importance. The R^2 value for the LR model was 0.1455 in this case, and for the 10-fold cross-validation (RF model) it varies between -3.368 and 0.643 (mean: -0.292, standard deviation: 1.371). The implication of negative R^2 , as mentioned before is that the fitted model shows a very different trend [42].

For the development versions of Avaya Communicator for Android, all the predictor variables turned out be insignificant in both the LR (Table 7) and the BN (Figure 10) models. The RF model gives a really bad fit in the 10-fold cross-validation as well, with R^2 values ranging from negative infinity to 0.64 (median: -5.74), indicating the predictors are very poor. Still, the two usage related variables have the lowest importance in the variable importance plot as seen in Figure 13.

Finally, for the GA releases of Avaya mobile SIP for iOS, the LR model (Table 7) as well as the BN model (Figure 11) shows that release date and release duration have effect on the “Quality” variable, but the two other usage variables have no effect. The RF model again gives a really bad fit in the 10-fold cross-validation, with R^2 values ranging from negative infinity to -36.43, indicating the predictors are very poor. Still, the two usage related variables have the lowest importance in the variable importance plot as seen in Figure 14.

Table 7 Summary Result of LR Model for “Quality”

	GA releases of Android			Dev releases of Android			Mobile SIP for iOS		
	Estimate	Std. Error	p-value	Estimate	Std. Error	p-value	Estimate	Std. Error	p-value
(Intercept)	-59.6470	19.9155	0.0032	-285.0676	180.0109	0.1290	28.5606	1.1271	0.0000
Release.Date	6.1428	2.0503	0.0031	29.3616	18.5247	0.1287	-2.8833	0.1167	0.0000
Release.Duration	0.0514	0.0271	0.0600	0.0854	0.1152	0.4671	-0.0923	0.0012	0.0000
Usage.Frequency	-0.1593	0.0618	0.0108	-0.3136	0.2372	0.2011	0.0006	0.0037	0.8831
Usage.Intensity	0.0395	0.0175	0.0253	0.0595	0.0820	0.4764	0.0001	0.0022	0.9614

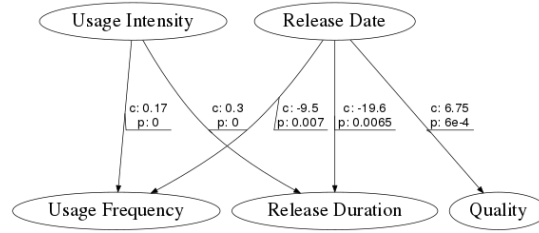


Fig. 9 Bayesian Network Model for “Quality” - GA releases of Avaya Communicator for Android (with c: coefficients after fitting the transformed, but unscaled data, p: p-value for the link)

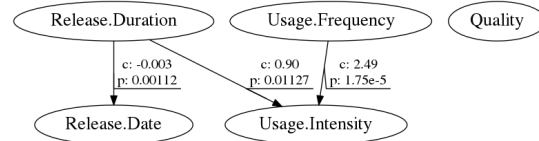


Fig. 10 Bayesian Network Model for “Quality” - Development releases of Avaya Communicator for Android (with c: coefficients after fitting the transformed, but unscaled data, p: p-value for the link)

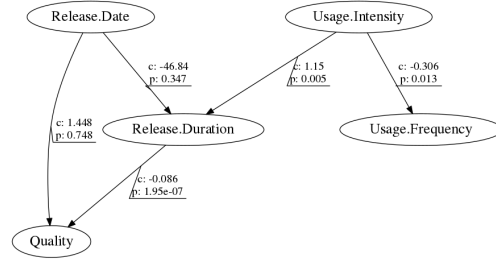


Fig. 11 Bayesian Network Model for “Quality” - GA releases of Avaya mobile SIP for iOS (with c: coefficients after fitting the transformed, but unscaled data, p: p-value for the link)

The results from these analyses clearly indicate that the quality measure defined by the number of exceptions per user is independent of software usage, and, therefore, suitable for comparing the quality of software development process among different releases of a software.

8.1 Timeline of Quality

We wanted to see how the perceived quality of the releases of the different mobile applications described above change with time. As a general trend, we observe that most of the exceptions occur right after the release date. then, as the number of users keep increasing with time, the value of the quality variable drop and come to a stable value. In this paper we show only the timeline for GA releases of Avaya mobile SIP for iOS (Figure 15), since the other two softwares had a lot of releases, making them difficult to identify from the plot.

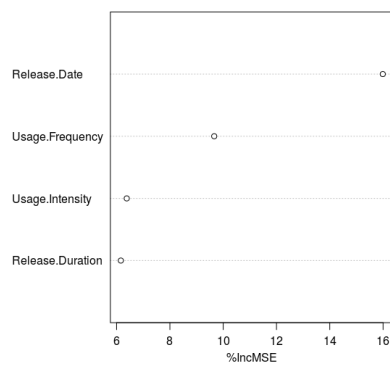


Fig. 12 Variable Importance plot from the Random Forest Model for Quality Variable - GA releases of Avaya Communicator for Android

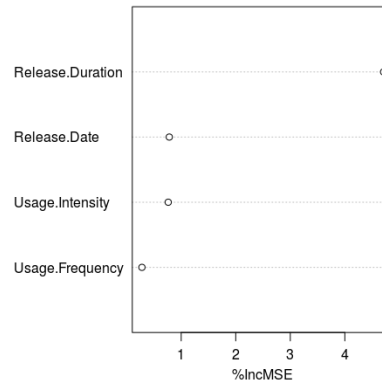


Fig. 13 Variable Importance plot from the Random Forest Model for Quality Variable - Development releases of Avaya Communicator for Android

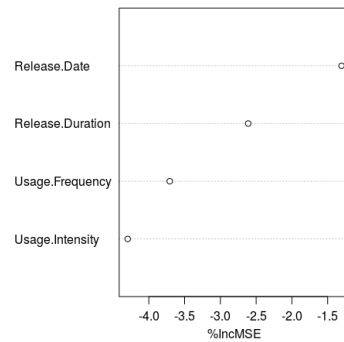


Fig. 14 Variable Importance plot from the Random Forest Model for Quality Variable - GA releases of Avaya mobile SIP for iOS

The other two are available in our GitHub repository:
https://github.com/tapjdey/release_qual_model.

9 Analysis of the NPM Data: Results

As mentioned before, for the analysis of the NPM data we focused on the timeline of the entire packages from 2015-03-01 to 2018-08-31, instead of the individual releases of a package to reduce the effect of the downloads by bots and other automated sources. Since we only had three variables (no. of issues, no. of downloads, calendar date) in our dataset, we decided to go for the simpler liner model. We had the total number of reported issues for the package as the response variable and observed how much the number of daily downloads, before and after controlling for the calendar date, explained it.

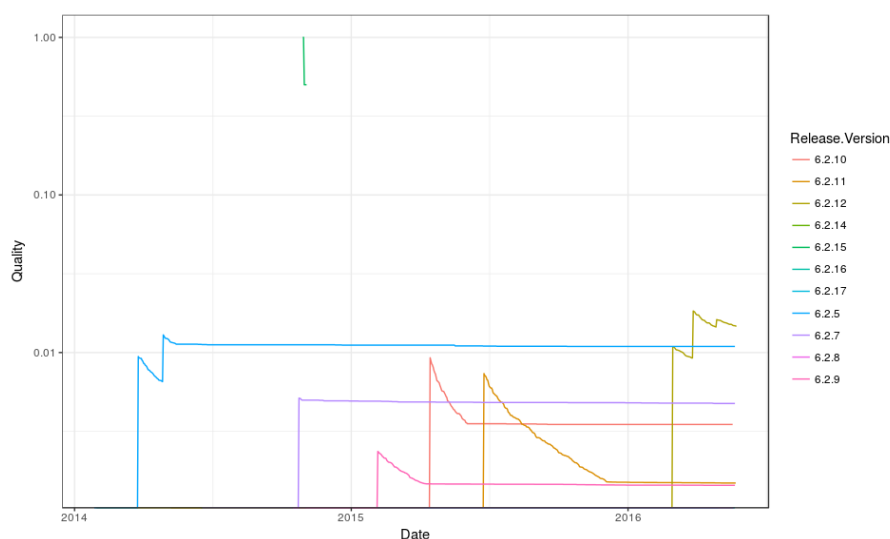


Fig. 15 Timeline for Quality Variable - GA releases of Avaya mobile SIP for iOS

We found that out of the 4430 packages, only for 36 packages the p-value of the predictor variable was more than 0.05 before adjusting for the calendar date, and after adjusting for the calendar date, p-value was less than 0.5 (*i.e.* the predictor was deemed significant) for all 4430 packages.

The R^2 values of the fitted models varied between 0 and 0.8637 (median: 0.4982, standard deviation: 0.0618) before adjusting for the calendar date, and between 0.019 and 0.999 (median: 0.9195, standard deviation: 0.0618) after adjusting for the calendar date. So, we can say that the number of daily downloads is a significant and important predictor for the number of issues encountered by the users for most of the packages and the effect is more pronounced when the effect of automated downloads is controlled by calendar date.

Since we just established that the number of issues of an NPM package depends on the number of daily downloads, a similar quality metric of number of issues per download should be applicable in this situation as well.

To check the quality of different packages, we looked at the minimum and median value of the quality metric for the 4430 packages. We didn't consider the absolute maximum value, since some packages had zero downloads for a few days, driving the value of the quality metric to infinity. So, we used the 90th quantile value as a proxy for the maximum value. We looked at the packages for which the value of the quality metric was more than 1. The threshold was chosen because we were looking at the packages of really high values of the quality metric, and thus were of poor quality. We found that for 340 packages the 90th quantile value of the quality metric was more than 1. The number was 100 and 3 when we looked at the median and minimum values respectively. The three packages for which the total number of issues over the number of daily downloads was more than 1

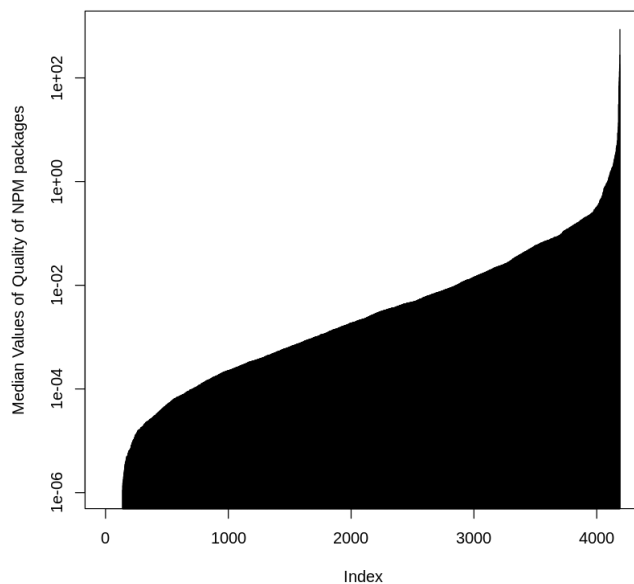


Fig. 16 Histogram of Median Values of Quality of NPM packages

were '@ngrx/store', '@protobufjs/fetch', and '@protobufjs/inquire'. Overall, we found that the packages for which the value of the quality metric was more than 1 were mostly packages from a big project that were relatively less downloaded, *e.g.* 'babel-plugin-transform-es2015-block-scoped-functions' from babel project, 'react-scripts' from facebook-react project etc. There were a few other packages that had very few downloads during most of its life-cycle since 2015, but had an increase in popularity later on, and thus were selected in our list of packages. However, since they had very few downloads for a long time, the median or maximum value of the quality metric was more than 1 (*e.g.* 'bubbleify'). For illustration, in Figure 16 we are showing the histogram of the median value of the quality metric for the 4430 NPM packages, which gives some idea about the overall quality distribution of the packages in the NPM ecosystem. We can see that around 75% of the packages in NPM have a median value of the quality metric less than 0.01, which mean, overall, for around 75% of the NPM packages, less than 1 in 100 regular users ever (since we are looking at the total number of issues) file an issue.

Further inspection showed that the value of the quality variable increases with time for almost half of the packages (2030 out of 4430 packages, 45.8%), unlike what we observed for the mobile applications, where for almost all of the releases of the three softwares, the value of the quality variable decreased with time.

Here we also show the timelines of the quality variable we defined (*i.e.* in this case, number of issues per download) for a few selected well-known

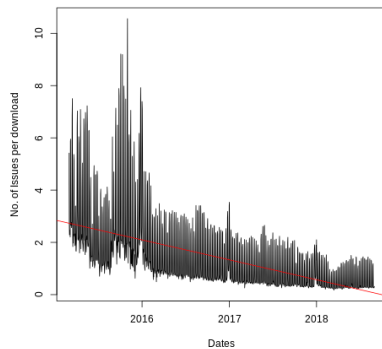


Fig. 17 Timeline for Quality Variable for NPM package: angular

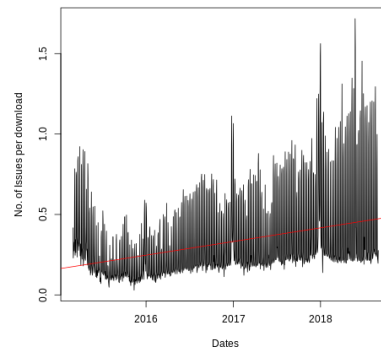


Fig. 18 Timeline for Quality Variable for NPM package: babel

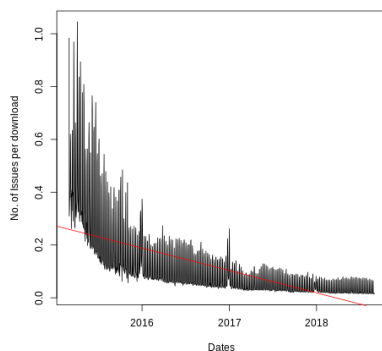


Fig. 19 Timeline for Quality Variable for NPM package: eslint

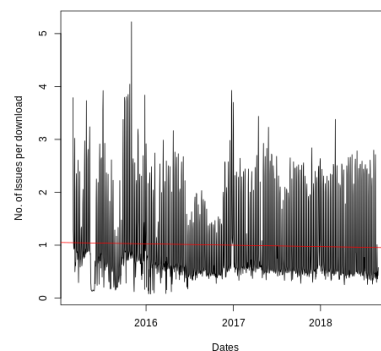


Fig. 20 Timeline for Quality Variable for NPM package: ember-cli

NPM packages for illustration. We see that “angular”(Figure 17) and “eslint”(Figure 20) have a trend similar to what we saw for the mobile apps, with the value of the quality variable decreasing with time, but “babel”(Figure 18) is showing an increase in the value, while for “ember-cli”, the trend is overall constant over time.

10 Implication of our Findings

Our analysis makes it evident that the number of new users is the most important variable in explaining various post-release variables, as seen in all three of the mobile applications as well as for the NPM packages, where the number of downloads was found to be an important predictor for the number of issues. The analysis also indicates that more new users for a release indicate more exceptions being found for the software and, for the GA releases of both the apps analyzed, longer activity for the release (the duration of a release mea-

1 sures how long a release is actively used by users, not the time between two
2 releases, since the releases overlap). This suggests that users may be reluctant
3 to upgrade (or are encouraged to stay) on better-quality releases. Our findings
4 are in agreement with findings of [37, 19, 41] that consider post-release defects
5 for a completely different server software system.
6

7 The release date also affects no. of exceptions, as can be observed by looking
8 at the coefficients. It provides some insight on how this software has evolved.
9 Even after compensating for the effect the number of users have on the num-
10 ber of exceptions, the number of exceptions are increasing with time for the
11 Android app, whereas is decreases for the iOS app. This may indicate that
12 as the software, the OS, as well as the hardware are becoming more complex
13 with time, which is consistent with a rapid growth of functionality and the
14 size of associated code base. The Android app is seeing more crashes due to
15 the variations in the devices and the OS, whereas for iOS, since the devices as
16 well as the OS versions are tightly controlled, the users are seeing less issues,
17 although we have no explicit evidence to support our speculation.
18

19 We found from the timeline that, as a general trend, most of exceptions
20 occur very early after the release, then, as the number of users increase, the
21 perceived quality of the release improves as well.
22

23 An interesting observation from the model is the lack of any direct re-
24 lationship between exceptions and the intensity or frequency of usage. One
25 possibility is that exceptions happen for specific Android/ iOS version/ Phone
26 combination and the way each user is exercising application's functionality.
27 Users for whom the application crashes must wait for the next release. This
28 would lead to the observed phenomena where only the new users increase
29 the number of crashes, which was observed more clearly from the timeline of
30 crashes as well. The duration an application is used by individual users was
31 found to have a much smaller effect on reported defects than the number of
32 new users in prior work [19, 37, 40] as well. In particular, it was observed that
33 most of the issues happen soon after deploying the release and the chances
34 of reporting a defect for a new release drops very rapidly with time after
35 installation.
36

37 From the analysis of the NPM packages we observed that the number of
38 downloads is a significant predictor for the number of issues for most of the
39 NPM packages, and when controlled for the calendar date, which compensated
40 for the variations in the downloads by automated sources, it was a significant
41 predictor for all the NPM packages. So, a similar quality measure was used
42 for this case as well. We found by looking into this metric that, overall, for
43 around 75% of the NPM packages, less than 1 in 100 regular users ever (since
44 we are looking at the total number of issues) file an issue. However, unlike the
45 three mobile apps, where the value of our quality metric decreases with time
46 for all releases, for the NPM packages the quality metric sometimes increases
47 or remain relatively constant over time (around 45.8% of the time).
48

49 Our data, scripts, and more detailed results are available in our GitHub
50 repository: https://github.com/tapjdey/release_qual_model.
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

We found that the exceptions are a result of more new users and the extent of usage does not appear to have a direct effect on the number of new users.

Overall, none of the three models indicate that “Usage.Frequency” or “Usage.Intensity” have any effect on the “Quality” variable. We, therefore, suggest that the exceptions per new user, or a metric similar to that, can be used as a software development quality metric to compare quality of different releases and to quantifying the impact of variables representing the development process of a release.

11 Related Work

Although software quality has always been a common topic in software engineering [4, 29], most of the studies have focused on pre-release data, primarily due to the developers’ concern about finding the appropriate balance between the amount of testing required and the quality of software (e.g. [52, 10]). There have been a number of works on predicting and improving the software quality as well (e.g. [36, 65, 25, 39]). Comparatively, studies about post-deployment quality and dynamics have been less frequent [31, 27]. However, a number of studies have looked at the aspects of software quality metrics, especially the quality perceived by the customers, e.g., [41, 37, 19, 51, 35]. A notable non-academic work involves a study of mobile app monitoring company’s (Criticism) data [18]. The author of the news article found it necessary to normalize crash data by the number of launches. Finally, an empirical investigation between release frequency and quality on Mozilla Firefox has been investigated in [28].

While Bayesian Networks have been used for software defect prediction for decades, the use of BNs for explanatory modeling in empirical software engineering is still not common despite the promise. A case for use of BNs was made by Fenton et.al. [16, 13], while the earliest publications utilizing BNs we could find [23] constructed search of the structure based on the statistical significance of partial correlations in the context of modeling delays in globally distributed development. [58, 48] considered the application of Bayesian networks to prediction of effort, [15, 44, 45] used Bayesian networks to predict defects, and [46] used BN approach for an empirical analysis of faultiness of a software. On the other hand, Bayesian structure learning is a big domain in itself with a wide range of algorithms, but its use in software engineering context is not very common.

Post-release defect density calculated as a proportion of users who experience an issue within a certain period after installing or upgrading to a new release has been proposed by [38, 41] as a measure of software quality. Hackbarth et al. [20] also found the need to adjust defect counts in their proposed measure of software quality as perceived by customers. We propose a somewhat different measure of quality based on the number of exceptions per user.

In general, software quality is a widely researched topic [26, 29, 53] etc., but in our knowledge, this is the first model-based attempt to obtain a usage independent measure of software quality and the first attempt to model exceptions in mobile applications.

The NPM ecosystem is one of the most active and dynamic JavaScript ecosystems and [62] presents its dependency structure and package popularity. [64] studies the dependency, specifically the lag in updating dependencies in various NPM packages while [2] looked into the use of trivial packages as part of package dependencies for different NPM packages.

The advancements proposed in this paper over the published work are focused on two primary areas: (1) study of the relationship between software faults (issues for NPM packages) and usage using **post-release** data in the context of two proprietary mobile applications and 4430 popular NPM packages, and (2) proposing a usage independent exception-based software quality metric based on our models.

12 Limitations

In terms of modeling aspects, there are some limitations related to the different approaches. The RF model was used for 10 fold cross-validation, and exhibited a rather high value of standard deviation in the R^2 value, and it even went negative a few times indicating a very bad fit [42], likely due to the small sample size and not having good predictors in some models.

While creating the BN model we did not cover all possible ways BNs can be applied to gain insight into the system. For example, we did not investigate the possible existence of any hidden node, or make an effort to formally establish the causal relationship between the nodes. We also did not investigate how the properties of one release affect the subsequent releases, nor did we investigate the presence of any feedback loops.

In the simulation study, although we covered an extensive set of options, we did not try every possible combination of options for the BN structure search exercise.

We also did not use Markov Random Field analysis, which is another probabilistic graphical modeling approach. The primary reason behind choosing the BN approach was that we found an example where this method was used to successfully recover the underlying network [54]. Moreover, it is possible to interpret a BN model as a causal model, and although we did not use that interpretation in this study, our goal is to eventually establish a causal mechanism of how usage affects the number of exceptions/defects experienced by users, so we wanted to use BN from the start.

The accuracy of our result is very much dependent on the Google Analytics data. While we do not have reasons to doubt the accuracy of the counts in Google Analytics data, we would have liked to have better definitions of how it determines “New User”, “Visit”, and, especially, nontrivial to aggregate quantities such as “Visits per User.” Also, it is not clear if Google Analytics

1 distorts data in any way (e.g., by applying differential privacy transformations)
2 for low counts in order to protect the privacy of the users. We do not believe
3 it does, but we have not conducted an experiment to validate that.
4

5 Furthermore, the projects under consideration were relatively new and it
6 was the first attempt for the team to deploy mobile software. As such, much
7 was not well documented and was rapidly evolving over time. As mentioned
8 earlier, we did not have the official release dates for all releases, so we put
9 the start date of the release as the date on which the first usage was reported.
10 However, we did verify the official dates with this reported date for the releases
11 for which we found the release date, and they were very close, but not always
12 exactly the same. This should not affect the overall result, given the total time
13 scale of more than two years. The release end dates, by their nature, have to
14 be estimated based on user activity, since there is no way to force end user
15 to upgrade Android app. For recent releases, therefore, the end date may be
16 censored by our data collection date, hence the duration for these releases
17 might be underestimated.
18

19 It may be possible to collect numerous additional variables that may have
20 an impact on exceptions, for example, the number of changes to the source
21 code made for a release as was done in [37]. Unfortunately, due to the nature of
22 parallel development for multiple releases and products noted in subsection ??,
23 it was virtually impossible to separate the changes that would only affect the
24 specific release on the Android platform. There might be other unobserved
25 variables driving some relationships, but not explored in this study.
26

27 Our model is obtained on a single set of mobile applications from a specific
28 domain, implemented via a rather complex codebase and is certainly not rep-
29 resentative of most mobile applications that tend to be much simpler. Further-
30 more, mobile applications may not represent other types of software further
31 limiting external validity of the results. However, some aspects that we see in
32 the specific application, such as increasing number of faults with the number
33 of users, has been observed in rather different contexts of large-scale server
34 software. This suggests that the model derived in the study may generalize to
35 other domains as well.
36

37 The study of the NPM packages was rather limited in scope, since we only
38 looked at the popular (having more than 10,000 monthly downloads since Jan-
39 uary 2018) packages to reduce the effect of downloads by automated sources.
40 We didn't look at the effects of other factors that affect the number of down-
41 loads, *e.g.* the number of dependents a package has. Although some of the
42 issues could come from users of a dependent package, we didn't actively check
43 the origins of the issues to verify that. We also didn't look at the releases of the
44 packages, because of reasons mentioned before. We didn't differentiate between
45 the types of the issues, because we just wanted to see how many times a user
46 decided to file an issue. Overall, this study was not a direct extension of the
47 previous work, rather, it was an extension of the concept and its application
48 in a different domain.
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

13 Conclusion

From the practical perspective we have established that the extent of use has very strong relationship with the number of exceptions for three large mobile applications from the telecommunication domain. Counting exceptions, therefore, will not accurately measure the quality of software development process but, instead, it would strongly depend on the extent of use. In order to produce a measure that the development team can use to understand and improve quality of their software development process, we proposed to normalize exceptions by usage based on the Bayesian Network models. Notably, a similar normalization was previously proposed in the context of post-release defects that also exhibited strong positive correlation with the number of users. As a larger proportion of applications are mobile and/or delivered as a service, the amount of usage can be relatively easily collected. Consequently, not adjusting software development measures for usage should not be considered as an excusable practice.

From the theoretical perspective we provided the explanation of the relationships among post-deployment quantities using Linear Regression and Bayesian Networks. Linear Regression can be thought as a special Bayesian Network with the response node being potentially connected to each predictor node. Bayesian Networks allow for exploration of relationships among all variables and empirical determination of the relationships exhibited in a particular dataset. For all three mobile softwares analyzed, the number of users was found to be the most significant predictor with both the models. It would be preferable to have each release as a separate categorical predictor, but because for simplicity we chose to use only one observation per release, it could not be done. If, instead, we considered exceptions during different periods of a release, that would have allowed us to introduce such categorical variable and interpret the estimated coefficient for each release as release quality (higher number meaning lower quality).

We also established that it is possible to predict exceptions using Random Forest modeling techniques and that usage plays a key role for the accuracy of these predictions. However, the performance of the predictive model was not consistently good, since, as noted above, prediction is a different task than explanation and, even though it often yields more accurate results, the prediction results may be harder to explain to developers or managers and, therefore, harder to act upon. We believe the findings do have a message for the voluminous research in defect prediction. While defects are not exceptions, usage was also found to affect post-release defects in a similar manner [24, 19, 41]. It would, therefore, be advisable to incorporate forecasts of usage into defect prediction models to increase their accuracy.

Our analysis of the NPM packages established that our approach is extendable to other domains as well. The study revealed that even a less accurate measure of usage like downloads, which, for NPM packages, is a mix of downloads by human users and automated sources, is an important predictor for the number of issues reported, which again is a weakly similar measure to the

number of crashes or bugs. So, our approach can be applied to any situation similar to the ones we studied, even when only proxy measures for usage and crashes/ bugs are available.

We hope that this work will spur more research on software engineering aspects in post-deployment stage because, like mobile applications, modern web applications are even more reliant on usage monitoring not simply from the perspective of crash counting but also because the usability or even revenue stream from the software applications critically depends on how users behave.

From the practical perspective, we hope that any mobile or web software project can easily apply and refine the presented approach of using Google Analytics data to improve the quality of their software. Any Android OS or Apple iOS mobile application can use Google Analytics to monitor application usage and crashes, so the approach should be widely applicable. Despite that, we are not aware of any prior empirical study that would leverages Google Analytics or similar data for software quality modeling.

The result of our simulation study should also be useful for practitioners using Bayesian Network structure search techniques for choosing the best performing methods.

Finally, much more work is needed to gather additional empirical evidence of how software behaves post-deployment. It is important to note that Google Analytics data is available only for application developers, so while each project has the ability to see their app's performance, they can not see data for software created by other organizations. This can be addressed by a) projects sharing their post-deployment data (we have not seen examples of that); or b) publishing findings based on such data in cases such as ours, where the data itself would be impossible to release publicly since it involves numerous, often enterprise, customers who may not agree.

Acknowledgements This work was supported by the National Science Foundation (U.S.) under Grant No. 1633437.

References

1. (2014). URL <https://developers.slashdot.org/story/17/01/14/0222245/nodejss-npm-is-now-the-largest-package-registry-in-the-world>
2. Abdalkareem, R., Nourry, O., Wehaibi, S., Mujahid, S., Shihab, E.: Why do developers use trivial packages? an empirical case study on npm. In: Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering, pp. 385–395. ACM (2017)
3. Balov, N., Salzman, P.: catnet: Categorical Bayesian Network Inference (2016). URL <https://CRAN.R-project.org/package=catnet>. R package version 1.15.0
4. Boehm, B.W., Brown, J.R., Lipow, M.: Quantitative evaluation of software quality. In: Proceedings of the 2nd international conference on Software engineering, pp. 592–605. IEEE Computer Society Press (1976)
5. Borges, H., Hora, A., Valente, M.T.: Understanding the factors that impact the popularity of github repositories. In: Software Maintenance and Evolution (ICSME), 2016 IEEE International Conference on, pp. 334–344. IEEE (2016)
6. Bottcher, S.G., Dethlefsen, C.: deal: Learning Bayesian Networks with Mixed Variables (2013). URL <https://CRAN.R-project.org/package=deal>. R package version 1.2-37

7. Briand, L.C., Wüst, J., Daly, J.W., Porter, D.V.: Exploring the relationships between design measures and software quality in object-oriented systems. *Journal of systems and software* **51**(3), 245–273 (2000)
8. Chickering, D.M.: Learning bayesian networks is np-complete. *Learning from data: Artificial intelligence and statistics V* **112**, 121–130 (1996)
9. Chlebus, B.S., Nguyen, S.H.: On finding optimal discretizations for two attributes. In: *International Conference on Rough Sets and Current Trends in Computing*, pp. 537–544. Springer (1998)
10. Dalal, S.R., Mallows, C.L.: When should one stop testing software? *Journal of the American Statistical Association* **83**(403), 872–879 (1988)
11. Dey, T., Mockus, A.: Modeling relationship between post-release faults and usage in mobile software. In: *Proceedings of the 14th International Conference on Predictive Models and Data Analytics in Software Engineering*, pp. 56–65. ACM (2018)
12. Duc, A.N., Mockus, A., Hackbarth, R., Palframan, J.: Forking and coordination in multi-platform development: a case study. In: *ESEM*, pp. 59:1–59:10. Torino, Italy (2014). URL <http://dl.acm.org/authorize?N14215>
13. Fenton, N., Krause, P., Neil, M.: Software measurement: Uncertainty and causal modeling. *IEEE software* **19**(4), 116–122 (2002)
14. Fenton, N., Neil, M., Marquez, D.: Using bayesian networks to predict software defects and reliability. *Proceedings of the Institution of Mechanical Engineers, Part O: Journal of Risk and Reliability* **222**(4), 701–712 (2008)
15. Fenton, N., Neil, M., Marsh, W., Hearty, P., Marquez, D., Krause, P., Mishra, R.: Predicting software defects in varying development lifecycles using bayesian nets. *Information and Software Technology* **49**(1), 32–43 (2007)
16. Fenton, N.E., Neil, M.: A critique of software defect prediction models. *IEEE Transactions on software engineering* **25**(5), 675–689 (1999)
17. Friedman, N., Goldszmidt, M., Wyner, A.: Data analysis with bayesian networks: A bootstrap approach. In: *Proceedings of the Fifteenth conference on Uncertainty in artificial intelligence*, pp. 196–205. Morgan Kaufmann Publishers Inc. (1999)
18. Geron, T.: Do ios apps crash more than android apps? a data dive (2012). <https://www.forbes.com/sites/tomiogeron/2012/02/02/does-ios-crash-more-than-android-a-data-dive>
19. Hackbarth, R., Mockus, A., Palframan, J., Sethi, R.: Customer quality improvement of software systems. *Software, IEEE* **33**(4), 40–45 (2016). URL [papers/cqm2.pdf](#)
20. Hackbarth, R., Mockus, A., Palframan, J., Sethi, R.: Improving software quality as customers perceive it. *IEEE Software* **33**(4), 40–45 (2016)
21. Hahsler, M., Chelluboina, S., Hornik, K., Buchta, C.: The arules r-package ecosystem: Analyzing interesting patterns from large transaction datasets. *Journal of Machine Learning Research* **12**, 1977–1981 (2011). URL <http://jmlr.csail.mit.edu/papers/v12/hahsler11a.html>
22. Hartemink, A.J.: Principled computational methods for the validation and discovery of genetic regulatory networks. Ph.D. thesis, Massachusetts Institute of Technology (2001)
23. Herbsleb, J.D., Mockus, A.: An empirical study of speed and communication in globally-distributed software development. *IEEE Transactions on Software Engineering* **29**(6), 481–494 (2003). URL [papers/delay.pdf](#)
24. Jones, C.: Software quality in 2011: A survey of the state of the art. <http://sqgne.org/presentations/2011-12/Jones-Sep-2011.pdf> (2011). President, Namcook Analytics LLC, www.Namcook.com Email: Capers.Jones3@GMAIL.com
25. Kamei, Y., Shihab, E., Adams, B., Hassan, A.E., Mockus, A., Sinha, A., Ubayashi, N.: A large-scale empirical study of just-in-time quality assurance. *IEEE Transactions on Software Engineering* **39**(6), 757–773 (2013). URL <http://doi.ieeecomputersociety.org/10.1109/TSE.2012.70>
26. Kan, S.H.: *Metrics and models in software quality engineering*. Addison-Wesley Longman Publishing Co., Inc. (2002)
27. Kenny, G.Q.: Estimating defects in commercial software during operational use. *IEEE Transactions on Reliability* **42**(1), 107–115 (1993)
28. Khomh, F., Dhaliwal, T., Zou, Y., Adams, B.: Do faster releases improve software quality?: an empirical case study of mozilla firefox. In: *Proceedings of the 9th IEEE Working Conference on Mining Software Repositories*, pp. 179–188. IEEE Press (2012)

29. Kitchenham, B., Pfleeger, S.L.: Software quality: the elusive target [special issues section]. *IEEE software* **13**(1), 12–21 (1996)
30. Koller, D., Friedman, N.: Probabilistic graphical models: principles and techniques. MIT press (2009)
31. Li, P.L., Kivett, R., Zhan, Z., Jeon, S.e., Nagappan, N., Murphy, B., Ko, A.J.: Characterizing the differences between pre-and post-release versions of software. In: Proceedings of the 33rd International Conference on Software Engineering, pp. 716–725. ACM (2011)
32. Marco Scutari: Learning bayesian networks with the bnlearn r package. *Journal of Statistical Software* **35**(3), 1–22 (2010). URL <http://www.jstatsoft.org/v35/i03/>
33. Mockus, A.: Software support tools and experimental work. In: V. Basili, et al (eds.) *Empirical Software Engineering Issues: Critical Assessments and Future Directions*, vol. LNCS 4336, pp. 91–99. Springer (2007). URL <papers/SSTaEW.pdf>
34. Mockus, A.: Law of minor release: More bugs implies better software quality. <http://mockus.org/papers/IWPSE13.pdf> (2013). International Workshop on Principles of Software Evolution, St Petersburg, Russia, Aug 18-19 2013. Keynote
35. Mockus, A.: Engineering big data solutions. In: ICSE'14 FOSE, pp. 85–99 (2014). URL <http://dl.acm.org/authorize?N14216>
36. Mockus, A., Hackbarth, R., Palframan, J.: Risky files: An approach to focus quality improvement effort. In: 9th Joint Meeting of the European Software Engineering Conference and the ACM SIGSOFT Symposium on the Foundations of Software Engineering, pp. 691–694 (2013). URL <http://dl.acm.org/authorize?6845890>
37. Mockus, A., Weiss, D.: Interval quality: Relating customer-perceived quality to process quality. In: 2008 International Conference on Software Engineering, pp. 733–740. ACM Press, Leipzig, Germany (2008). URL <http://dl.acm.org/authorize?063910>
38. Mockus, A., Weiss, D.: Interval quality: Relating customer-perceived quality to process quality. In: Proceedings of the 30th international conference on Software engineering, pp. 723–732. ACM (2008)
39. Mockus, A., Weiss, D.M.: Predicting risk of software changes. *Bell Labs Technical Journal* **5**(2), 169–180 (2000). URL <papers/bltj13.pdf>
40. Mockus, A., Zhang, P., Li, P.: Drivers for customer perceived software quality. In: ICSE 2005, pp. 225–233. ACM Press, St Louis, Missouri (2005). URL <http://dl.acm.org/authorize?860140>
41. Mockus, A., Zhang, P., Li, P.L.: Predictors of customer perceived software quality. In: Software Engineering, 2005. ICSE 2005. Proceedings. 27th International Conference on, pp. 225–233. IEEE (2005)
42. (<https://stats.stackexchange.com/users/25/harvey-motulsky>), H.M.: When is r squared negative? Cross Validated. URL <https://stats.stackexchange.com/q/12991>. URL:<https://stats.stackexchange.com/q/12991> (version: 2014-05-06)
43. Nagarajan, R., Scutari, M., Lèbre, S.: Bayesian networks in r. Springer **122**, 125–127 (2013)
44. Neil, M., Fenton, N.: Predicting software quality using bayesian belief networks. In: Proceedings of the 21st Annual Software Engineering Workshop, pp. 217–230. NASA Goddard Space Flight Centre (1996)
45. Okutan, A., Yıldız, O.T.: Software defect prediction using bayesian networks. *Empirical Software Engineering* **19**(1), 154–181 (2014)
46. Pai, G.J., Dugan, J.B.: Empirical analysis of software fault content and fault proneness using bayesian methods. *IEEE Transactions on software Engineering* **33**(10), 675–686 (2007)
47. Pearl, J.: Bayesian networks. Department of Statistics, UCLA (2011)
48. Pendharkar, P.C., Subramanian, G.H., Rodger, J.A.: A probabilistic model for predicting software development effort. *IEEE Transactions on software engineering* **31**(7), 615–624 (2005)
49. Perez, A., Larranaga, P., Inza, I.: Supervised classification with conditional gaussian networks: Increasing the structure complexity from naive bayes. *International Journal of Approximate Reasoning* **43**(1), 1–25 (2006)
50. R Core Team: R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria (2017). URL <https://www.R-project.org/>

- 1 51. Rotella, P., Chulani, S.: Implementing quality metrics and goals at the corporate level.
2 In: Proceedings of the 8th Working Conference on Mining Software Repositories, pp.
3 113–122. ACM (2011)
- 4 52. Rubin, J., Rinard, M.: The challenges of staying together while moving fast: An ex-
5 ploratory study. In: Proceedings of the 38th International Conference on Software
6 Engineering, pp. 982–993. ACM (2016)
- 7 53. Schulmeyer, G.G., McManus, J.I.: Handbook of software quality assurance. Van Nos-
8 trand Reinhold Co. (1992)
- 9 54. Scutari, M.: Learning bayesian networks in r, an example in systems biology (2013).
10 <http://www.bnlearn.com/about/slides/slides-useRconf13.pdf>
- 11 55. Scutari, M., Strimmer, K.: Introduction to graphical modelling. arXiv preprint
12 arXiv:1005.1036 (2010)
- 13 56. Shmueli, G.: To explain or to predict? Statistical science pp. 289–310 (2010)
- 14 57. Sober, E.: Instrumentalism, parsimony, and the akaike framework. Philosophy of Science
15 **69**(S3), S112–S123 (2002)
- 16 58. Stamelos, I., Angelis, L., Dimou, P., Sakellaris, E.: On the use of bayesian belief networks
17 for the prediction of software productivity. Information and Software Technology **45**(1),
18 51–60 (2003)
- 19 59. Subramanyam, R., Krishnan, M.S.: Empirical analysis of ck metrics for object-oriented
20 design complexity: Implications for software defects. IEEE Transactions on software
21 engineering **29**(4), 297–310 (2003)
- 22 60. Voss, L.: numeric precision matters: how npm download counts
23 work (2014). URL [https://blog.npmjs.org/post/92574016600/](https://blog.npmjs.org/post/92574016600/numeric-precision-matters-how-npm-download-counts)
24 [numeric-precision-matters-how-npm-download-counts](https://blog.npmjs.org/post/92574016600/numeric-precision-matters-how-npm-download-counts)
- 25 61. Voss, L.: The state of javascript frameworks, 2017 (2018). URL [https://www.npmjs.](https://www.npmjs.com/npm/state-of-javascript-frameworks-2017-part-1)
26 [com/npm/state-of-javascript-frameworks-2017-part-1](https://www.npmjs.com/npm/state-of-javascript-frameworks-2017-part-1)
- 27 62. Wittern, E., Suter, P., Rajagopalan, S.: A look at the dynamics of the javascript package
28 ecosystem. In: Mining Software Repositories (MSR), 2016 IEEE/ACM 13th Working
29 Conference on, pp. 351–361. IEEE (2016)
- 30 63. Yu, P., Systa, T., Muller, H.: Predicting fault-proneness using oo metrics. an industrial
31 case study. In: Software Maintenance and Reengineering, 2002. Proceedings. Sixth
32 European Conference on, pp. 99–107. IEEE (2002)
- 33 64. Zerouali, A., Constantinou, E., Mens, T., Robles, G., González-Barahona, J.: An empir-
34 ical analysis of technical lag in npm package dependencies. In: International Conference
35 on Software Reuse, pp. 95–110. Springer (2018)
- 36 65. Zhang, F., Mockus, A., Keivanloo, I., Zou, Y.: Towards building a universal defect
37 prediction model with rank transformed predictors. Empirical Software Engineering
38 pp. 1–39 (2015)
- 39 66. Zheng, Q., Mockus, A., Zhou, M.: A method to identify and correct problematic software
40 activity data: Exploiting capacity constraints and data redundancies. In: ESEC/FSE’15,
41 pp. 637–648. ACM, Bergamo, Italy (2015). URL <http://dl.acm.org/authorize?N14200>



Click here to access/download
Supplementary Material
p56-Dey.pdf