



Technical debt forecasting: An empirical study on open-source repositories



Dimitrios Tsoukalas^{a,b,*}, Dionysios Kehagias^a, Miltiadis Siavvas^a, Alexander Chatzigeorgiou^b

^a Information Technologies Institute, Centre for Research and Technology Hellas, Thessaloniki 57001, Greece

^b Department of Applied Informatics, University of Macedonia, Thessaloniki 54643, Greece

ARTICLE INFO

Article history:

Received 5 November 2019

Received in revised form 3 July 2020

Accepted 3 August 2020

Available online 8 August 2020

Keywords:

Technical debt

Technical debt forecasting

Machine learning

Empirical study

ABSTRACT

Technical debt (TD) is commonly used to indicate additional costs caused by quality compromises that can yield short-term benefits in the software development process, but may negatively affect the long-term quality of software products. Predicting the future value of TD could facilitate decision-making tasks regarding software maintenance and assist developers and project managers in taking proactive actions regarding TD repayment. However, no notable contributions exist in the field of TD forecasting, indicating that it is a scarcely investigated field. To this end, in the present paper, we empirically evaluate the ability of machine learning (ML) methods to model and predict TD evolution. More specifically, an extensive study is conducted, based on a dataset that we constructed by obtaining weekly snapshots of fifteen open source software projects over three years and using two popular static analysis tools to extract software-related metrics that can act as TD predictors. Subsequently, based on the identified TD predictors, a set of TD forecasting models are produced using popular ML algorithms and validated for various forecasting horizons. The results of our analysis indicate that linear Regularization models are able to fit and provide meaningful forecasts of TD evolution for shorter forecasting horizons, while the non-linear Random Forest regression performs better than the linear models for longer forecasting horizons. In most of the cases, the future TD value is captured with a sufficient level of accuracy. These models can be used to facilitate planning for software evolution budget and time allocation. The approach presented in this paper provides a basis for predictive TD analysis, suitable for projects with a relatively long history. To the best of our knowledge, this is the first study that investigates the feasibility of using ML models for forecasting TD.

© 2020 Elsevier Inc. All rights reserved.

1. Introduction

The Technical Debt (TD) notion, a term inspired by the financial debt of economic theory, was introduced in 1992 by Ward Cunningham ([Cunningham, 1993](#)) as a metaphor intended to describe the problem of introducing long-term problems to software products, by not resolving existing quality issues early enough in the overall software development lifecycle (SDLC). The TD metaphor was initially related to software implementation (i.e., at the code level) but was gradually extended to other phases of the SDLC, i.e., software architecture, design, documentation, requirements, and testing ([Brown et al., 2010](#)). In the same manner like financial debt, TD incurs interest payments in the form of increased future software costs, usually caused by poor

design and code quality. To effectively manage the identification, quantification, and repayment of TD during the software development lifecycle, researchers and practitioners have developed and adopted a multitude of theories, methods and tools ([Li et al., 2015](#)).

However, predicting the accumulated TD during the evolution of a software application is an open and challenging research issue, as both the software system and its TD emerge in parallel ([Digkas et al., 2017](#)). The opportunity to predict TD is of paramount importance to software maintainability, which is recognized as one of the most effort-intense activities in the SDLC ([Ampatzoglou et al., 2015a](#)). System engineers and project managers need the right tools and appropriate training support to be able to perform long-term effective software maintenance ([Ampatzoglou et al., 2015a](#)). Therefore, forecasting the evolution of TD could be valuable in assessing the point at which the software product could become unmaintainable and to identify software artifacts, which are prone to accumulate significant levels of TD.

* Corresponding author at: Information Technologies Institute, Centre for Research and Technology Hellas, Thessaloniki 57001, Greece.

E-mail addresses: tsoukj@iti.gr (D. Tsoukalas), diok@iti.gr (D. Kehagias), siavvasm@iti.gr (M. Siavvas), achat@uom.gr (A. Chatzigeorgiou).

Although the topic of predicting the evolution of various aspects directly or indirectly related to the TD concept, such as code smells (Fontana et al., 2016a), fault-proneness (Arisholm and Briand, 2006) and general software evolution trends (Chaikalis and Chatzigeorgiou, 2015), has attracted the attention of both academia and industry, to the best of our knowledge no studies are focusing on the forecasting of TD itself (Tsoukalas et al., 2018). Hence, a method or tool that would provide practical decision-making support by predicting future TD of a software system that is expected to evolve over time can be valuable to software development teams. Consequently, software architects and project managers would be able to gain a better understanding of future TD issues and plan well in advance appropriate refactoring activities for saving maintenance costs.

As a first step towards TD forecasting, in our previous work, we have studied and applied statistical time series models for TD Principal forecasting (Tsoukalas et al., 2019). Statistical time series models are mostly univariate, i.e., they require only the historical data of the variable of interest to forecast its future evolution behavior and have thus been widely used in the literature for predicting software evolution trends, future change requests or software defects (Chaikalis and Chatzigeorgiou, 2015; Yazdi et al., 2014; Raja et al., 2009; Goulão et al., 2012; Kenmei et al., 2008). We used a dataset of 5 real-world open-source Java applications and found that the Autoregressive Integrated Moving Average model ARIMA(0,1,1) can provide accurate TD Principal predictions over a sufficiently long time period for all sampled applications. However, even though the overall ARIMA model performance was satisfactory for short-term TD Principal forecasting (up to 8 weeks ahead), we observed that its predictive performance dropped significantly for long-term predictions. Moreover, we concluded that ARIMA models might prove difficult to tune, as one has to follow the entire Box-Jenkins methodology (Box et al., 2015).

The work presented in this paper is a logical continuation and extension of our previous efforts (Tsoukalas et al., 2019), in order to provide a more complete approach for TD forecasting. More specifically, we believe that we can achieve better scores by trying out more advanced multivariate models able to support feature engineering, i.e., take into account various TD-related features and their combinations to generate better TD predictions. Therefore, while in our previous study (Tsoukalas et al., 2019) the main focus lied on univariate time series forecasting methods, in the present paper we attempt to empirically evaluate the ability of multivariate Machine Learning (ML) methods to adequately forecast future TD trends of software applications and achieve better and more practical results in both short-and long-term predictions. Building multivariate models that, alongside the evolution of the target variable, learn also from the evolution of additional features related to the target variable is a widely-used strategy (Bidarkota, 1998; du Preez and Witt, 2003), since the covariation of time series that follow similar time-based patterns can model interesting interdependencies and therefore improve forecast accuracy. To this end, in this paper, we extend our initial dataset by adding 10 more software applications (15 in total) and investigate whether the combination of software-related metrics acting as TD indicators and already existing ML forecasting methods could lead to the development of novel models that provide predictions about the evolution of TD in a software project. Towards this direction, we have studied and applied various popular ML methods, such as Regression, Regularization, Support Vector Regression, and Regression Trees to forecast the evolution of TD Principal.

The problem that our work attempts to solve can be summarized in the following research question:

RQ: Is the usage of machine learning models on a specific set of Technical Debt indicators a meaningful and accurate approach

to forecasting Technical Debt Principal in a long-lived, open-source software?

The objective of this study is to evaluate the ability of ML methods to model and predict the TD evolution of a software application, based on a set of TD indicators selected as TD predictors. The viewpoint is that of researchers who intend to investigate how different ML approaches can be effectively adopted by project managers and developers to accurately forecast the evolution of TD Principal and thus, support planning and decision-making. The context is an empirical study on TD Principal values of 15 real-world open-source Java applications publicly available in the GitHub repository. The included TD Principal values cover almost 3 years of each application's evolution, which corresponds to nearly 150 snapshots in weekly intervals. As such, the resulting models are expected to be meaningful in the context of the dataset constructed for this study. A positive answer to the formulated research question will suggest that ML models trained on a selected set of Technical Debt indicators can potentially be used as the basis for the construction of a TD Forecasting tool. We will also investigate the extent to which these models can properly capture the evolution of TD Principal values in terms of accuracy and forecasting length.

To shed light on this question, we conducted an empirical study following the roadmap illustrated in Fig. 1. Initially, we studied the relevant literature and identified TD indicators that could act as predictors, such as TD-related features and various Object Oriented (OO) metrics. Afterward, we constructed a relatively large code repository comprising 15 real-world open-source Java applications retrieved from the GitHub¹ online repository. For each application, we collected a subsequent number of snapshots (commits) ranging from 100 to 150 in weekly intervals, spanning up to almost 3 years of each application's evolution. This approach led to a dataset containing 1850 snapshots in total (171M lines of code). In order to extract the identified TD-related features and various Object Oriented (OO) metrics that could act as predictors, we used two popular tools, namely SonarQube² and CKJM Extended³ respectively. This process led to 15 independent application-specific datasets containing TD indicators and TD values for each snapshot. Subsequently, we employed techniques like correlation analysis, univariate and multivariate analysis, which allowed us to select the most statistically significant TD predictors and thus retain as much discriminatory information as possible. Finally, we examined potential ML forecasting models and algorithms that could be applied for TD prediction and compared their accuracy for various forecasting horizons in order to reach safer conclusions regarding the significance of the observed results. To the best of our knowledge, this is the first study in the field of TD that examines the applicability of ML models for TD forecasting.

An overview of the methodology described above is presented in Fig. 1.

The meaningfulness of any forecasting model is also related to its ability of reflecting the developers' perspective on whether the modeled phenomena and predicted evolution are useful. To this end, we have performed a survey to empirically evaluate the meaningfulness of the TD forecasting approach introduced in this study and to investigate the usefulness of the TD Forecasting concept in general, via a questionnaire distributed to representatives of a software company.

The rest of the paper is structured as follows: Section 2 presents the background concepts in the field of forecasting and TD. Section 3 presents the related work in the field of forecasting

¹ <https://github.com/>

² <https://www.sonarqube.org/>

³ http://gromit.iiar.pwr.wroc.pl/p_inf/ckjm/

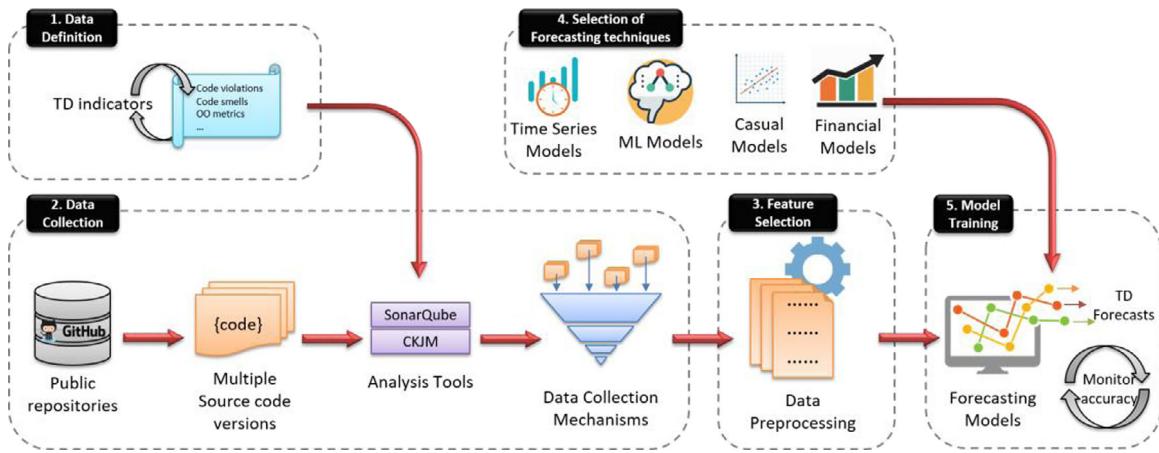


Fig. 1. Paper roadmap.

models and more specifically in their applicability for TD forecasting. Section 4 thoroughly describes data definition, collection and pre-processing steps. Section 5 describes forecasting model training, testing and benchmarking, as well as the current state of technical implementation of the proposed approach. Section 6 presents the results of a survey that has been conducted to empirically evaluate the meaningfulness of the TD forecasting approach introduced in this study. Section 7 reports the limitations and validity threats of this empirical study, while Section 8 discusses significant implications for both research and practice. Finally, Section 9 concludes the paper and discusses ideas for future work.

2. Background

This section provides an overview of relevant background on Forecasting and TD, in order to introduce unfamiliar readers with the main concepts of this paper.

2.1. Forecasting concepts

Forecasting is the process of making predictions of the future based on past and present data, usually by analysis of trends. Being able to predict future values of an observed attribute plays an important role in nearly all fields of science and engineering (Palit and Popovic, 2006). Due to the increasing variety and complexity of forecasting problems over the years, many forecasting techniques have been developed, and continue to be developed until today, each for a special use. The forecasting domain has been influenced, for a long time, by statistical methods that can be classified under two broad categories: causal (or associative) and time series models. Causal models (including the widely used regression analysis) assume that there is a cause-and-effect relationship between the variable of interest and other variables, and therefore try to discover that relationship to forecast future values. Time series models (including the widely used ARIMA model) treat the examined system as a black box and assume that information needed to forecast is contained in a set of time-dependent data that will continue to follow same patterns as in the past (Das, 2012). During the last decades however, Machine Learning (ML) models have drawn attention and have established themselves as serious contenders to classical statistical models in the forecasting community (Bontempi et al., 2013). These models, also called black box or data-driven models, are self-correcting learning algorithms that utilize supervised, unsupervised or reinforcement learning to acquire knowledge of the stochastic dependency between the past and the future, based only on historical data.

The experts' opinions regarding which of the two approaches (i.e., time series and ML) yields more accurate predictions vary. In a recent study by Makridakis et al. (2018), the authors claim that ML methods need to become more accurate, requiring less computer time, and be less of a black box. A major contribution of their paper is in showing that traditional statistical methods are more accurate than ML ones and pointing out the need to discover the reasons involved, as well as devising ways to reverse the situation. However, in their comparisons they made clear that the results might be related to the specific dataset being used. They believe that if the series are much longer in length, ML methods can train their weights more optimally. On the other hand, in related studies Werbos showed that Artificial Neural Networks (ANNs) can achieve better results compared to traditional statistical methods such as linear regression and Box-Jenkins (ARMA, ARIMA) approaches (Werbos, 1974, 1988). A similar study by Lapedes and Farber (1987) concludes that ANNs can be successfully used for modeling and forecasting nonlinear time series. Recently, other models appeared such as regression trees, support vector regression and nearest neighbor regression (Friedman et al., 2001; Alpaydin, 2010).

2.1.1. Causal or associative models

Causal, or associative, models assume that the variable that needs to be forecasted is somehow related to other variables in the environment through a cause-and-effect relationship. In this case, the forecasting challenge is to discover the relationships between the variable of interest and these other variables. These relationships, which can be very complex, take the form of a mathematical model, which is used to forecast future values of the variable of interest. Some of the best-known causal models are regression models, such as Linear and Multivariate regression, or regularization models, such as Ridge and Lasso regression.

Linear and Multivariate regression are the most commonly used techniques for modeling the relationship between two or more independent variables and a dependent variable by fitting a linear equation to observed data. In the simplest case, the Linear regression model allows for a linear relationship between the forecast variable and a single predictor variable. When there are two or more predictor variables, Multivariate regression is used. The main advantages of these techniques are their simplicity and that they are supported by many popular statistical packages. During Linear and Multivariate regression, the coefficients of these variables are estimated using the least squares method. However, quite often simple linear regression models are suffering from over-fitting or under-fitting. Hoerl and Kennard (2000) and Tibshirani (1996) regression are some of the

simple techniques to reduce model complexity, reduce multicollinearity and prevent over-fitting by applying regularization, i.e., add some more constraints to the loss function. In the case of Ridge regression, those constraints are the sum of squares of the coefficients multiplied by the regularization coefficient (lambda). This regularization type is called L2. Lasso regression, works in a similar way but instead of adding squares to the loss function it adds absolute values of the coefficients. As a result, during the optimization process, coefficients of unimportant features may become zero, which allows for automated feature selection. This regularization type is called L1.

2.1.2. Machine learning methods

ML models are self-correcting learning algorithms that utilize various forms of learning, such as supervised, unsupervised or reinforcement, to predict new outcomes based on previously known results. Although most of these methods have existed for a long time, it is only during the last decades that they have drawn attention due to the constantly improving models, data and processing capacities. While traditional statistical forecasting techniques use only strictly formatted historical data, ML forecasting can take advantage of several data sources, since data can be of different source, format, dimensionality, etc. However, if not handled correctly, these methods can suffer from serious drawbacks such as the lack of interpretability (black box), expensive computational requirements or overfitting. Some of the most widely used ML forecasting models are Support Vector regression, K-Nearest Neighbor regression, Decision Trees, Random Forest regression and various ANN variants, such as Multi-Layer Perceptron, Bayesian and Generalized Regression Neural Networks (Makridakis et al., 2018).

Support Vector Machines (SVM) were originally developed for solving classification problems but have been later extended to the domain of regression problems. The goal of Support Vector regression (SVR) (Drucker et al., 1997) is to find a function that approximates the actually obtained target values for all the training data, and has a minimum generalization error. To achieve this, it tries to learn a non-linear function by linearly mapping features into high-dimensional, kernel-induced feature space. K-Nearest Neighbor regression (Altman, 1992) is a nonparametric regression method basing its forecasts on a similarity measure, i.e., the Euclidean distance between the points used for training and testing the method. Thus, given a number of inputs, the method picks the closest training data points and sets the prediction as the average of the target output values for these points.

As in the case of SVM, Decision Trees were originally developed for solving classification problems but were later extended to the domain of regression problems. A Regression Tree (RT) (Breiman, 2017) is a variant of decision trees that is built through an iterative process of splitting the data into partitions, and then splitting it up further on each of the branches. Initially all of the samples in the training set are put together in one node. The algorithm chooses an independent variable with values that minimize the sum of the squared deviations from the mean in the separate parts. An enhanced version of the RT algorithm is Random Forest (RF) method (Ho, 1998). RF is an ensemble of Decision Trees trained with the “bagging” method (Breiman, 1996). Bagging repeatedly selects a random sample with replacement of the training set and fits trees to these samples. After training, predictions for unseen samples can be made by averaging the predictions from all the individual regression trees or by taking the majority vote.

2.2. Technical debt concepts

Nowadays, TD is seen as an important part of software management, as many studies have identified several causes for its creation. Fowler (2003, 2009) states that software development debt is usually a consequence of time pressure. Kruchten et al. (2012) assign TD to YAGNI decisions (You Ain't Gonna Need It) that often result in unjustified and unnecessary investments in new features, architecture, over engineering, etc. Fowler (2009) proposes TD quadrant, a 2×2 matrix (Intentionality x Awareness), to visualize four different pathways that lead to TD. McConnell (2012) suggests a similar categorization, arguing that TD may be unintentional and intentional. Unintentional debt is often a consequence of poor coding practices, while intentional debt is a result of non-optimal decisions that are committed on purpose. Suryanarayana et al. (2014) point out that extreme situation when accumulated TD is enormous and cannot be paid off could lead to technical bankruptcy. Moreover, several recent studies have highlighted the need to analyze TD from SDLC point of view. Li et al. (2015) classify different TD types into ten levels based on their occurrence during the main phases of a software development process (i.e., requirements, design and architecture, implementation, testing, building, documentation, infrastructure, versioning, and defects).

2.2.1. Technical debt main components

The main component of TD is the *Principal*, which refers to the cost that has to be paid in order to eliminate the debt, i.e., the effort required to address the difference between the current and the optimal level of design-time quality. Depending on the type of TD, this can be translated into different kinds of activities, such as code refactoring, documentation updates or improving test coverage (Li et al., 2015). The second main component of TD is *Interest*, which is composed of two parts: (i) the *interest amount*, i.e., the potential penalty in terms of increased effort and decreased productivity that will have to be paid in the future as a result of not completing these tasks in the present (Seaman and Guo, 2011), and (ii) the *interest probability*, i.e., the probability that the artifact that contains the debt will undergo maintenance. When this additional effort (interest) reaches a level that makes maintenance so difficult and expensive that the system is no longer financially viable, the project is declared bankrupted (Ampatzoglou et al., 2018).

2.2.2. Technical debt indexes

In an attempt to provide an empirical TD quantification and assessment, various TD indexes, that is, indexes that offer an evaluation of the overall quality (in terms of TD) of a software application, have been proposed by researchers and subsequently implemented as industrial tools (Fontana et al., 2016b). To quantify their TD indexes, these tools, initially gather their atomic data by calculating several TD indicators, such as OO metrics, software quality metrics, violations or code and architectural smells. Subsequently, to assess the quality of both the architecture and the code of an application they employ well-known models for modeling TD, such as the ISO/IEC 25010 standard (ISO/IEC, 2011), and the Software Quality Assessment based on Lifecycle Expectations (SQALE) (Letouzey and Ikkiewicz, 2012) methodology among others (Curtis et al., 2012; Marinescu, 2012; Letouzey, 2012).

Regarding the various TD indexes and corresponding tools that have been proposed, in their study, Curtis et al. (2012) are based on Software Economics theories and quantify TD as the cost of violating architectural and code rules, giving three levels of severity to violations: high, medium and low. To achieve that, they introduce a cost function that quantifies principal and interest taking as input the number of must-fix violations, the time

required to fix each violation, and the cost for fixing a violation. To further support their findings, they integrate their formula into CAST,⁴ a tool that quantifies TD by identifying architectural and code violations and categorizing them by quality attributes. In another work, Marinescu (2012) introduces a novel framework for assessing TD using a technique for identifying architectural smells (called design disharmonies), detected by evaluating different metric-based rules that cover the majority of the aspects of design, such as complexity, coupling, and encapsulation. The impact of disharmonies is formulated as an index that uses three factors for its calculation, namely influence, granularity, and severity. This framework was integrated into inFusion,⁵ a tool that evaluates software quality by providing a global score known as the Quality Deficit Index (QDI). In a related study, Letouzey (2012) presents the widely used SQALE method for monitoring and assessing the quality and TD of the source code of a software application. One of the most representative tools for assessing the TD of a software product using the SQALE method is SQuORE,⁶ a commercial quality management tool that uses four indicators namely: efficiency, portability, maintainability, and reliability to calculate code TD. For each of these indicators, a set of quality rules is assigned. One of the advantages of this tool is that it takes into account source code, unit tests, documentation quality, available functional requirements, etc. resulting in a more accurate and complete calculation of TD. Finally, SonarQube⁷ is a widely-used open source platform for continuous inspection of code quality that provides analysis functionalities and a wide range of metrics for measuring code quality attributes. During the previous years, SonarQube used the SQALE method to assess the TD of a software product but has now switched to a different method. More specifically, it now checks code compliance against a set of classified coding rules and if the code violates any of these rules, it considers it as a violation or a TD item. Other popular quality assessment tools that worth mentioning are Sigrid,⁸ Structure101,⁹ NDepend,¹⁰ and Teamscale.¹¹

3. Related work

TD, a metaphor inspired by the financial debt of economic theory, indicates quality compromises that can yield short-term benefits in the software development process, but may negatively affect the long-term quality of software products. Numerous techniques, methods, and tools have been proposed over the years for estimating and managing TD, providing a variety of options to the developers and project managers of software applications. However, apart from managing TD, predicting its future value is equally important since this knowledge is expected to facilitate decision-making tasks regarding software implementation and maintenance, such as incurring or paying off TD instances. In this section, we investigate the state-of-the-art and examine the major contributions that have been made until today in the field of TD forecasting.

Software evolution is a term used in software engineering to refer to the process that starts with the development and then provides incremental updates of the software. According

to Lehman's laws of software evolution, software systems must evolve over time or they will become irrelevant (Lehman, 1980). Gaining a higher level of information about the evolution of large software systems is a key challenge in dealing with increasing complexity and decreasing software quality (Gall and Lanza, 2006). For this reason, the attempts to analyze, understand and predict the evolution of a software system have increased considerably in the last years (Godfrey and German, 2008), and nowadays, the terms software evolution and software maintenance are often used as synonyms (Mens, 2008). In his work, Mens (2008) stresses the need to develop better predictive models for measuring and estimating the cost and effort of software maintenance and evolution activities with higher accuracy. Therefore, the improvement of these models can be proven of great value in software development, since being able to estimate the future evolution of a software product, could provide valuable insight for its quality as well.

According to ISO/IEC 25010 (ISO/IEC, 2011), which is a well-accepted international standard, the notion of software quality is hierarchically decomposed into a set of quality attributes, like maintainability, reliability, and security. A multitude of quality models have been proposed over the years allowing the assessment and/or prediction of these quality attributes individually (Wagner, 2009; Van Koten and Gray, 2006; Zhou and Leung, 2007). For instance, Wagner (2009) implements a model based on Bayesian Belief Networks for assessing and predicting the maintainability of a software application based on a set of software metrics. Similarly, Van Koten and Gray (2006) try to predict object-oriented software maintainability by applying a Bayesian network, while Zhou and Leung (2007) approach the same problem by using multivariate adaptive regression splines.

Since quality attributes are relatively abstract and difficult to be measured directly from the artifacts of software products (e.g., source code), ISO/IEC 25010 (ISO/IEC, 2011) further decomposes them into a set of more concrete quality properties (e.g., complexity), which can be directly quantified through common metrics (e.g., McCabe's Cyclomatic Complexity). Similarly, to the high-level quality attributes, a large number of methods have been proposed to estimate the future evolution of software quality properties and metrics used to calculate them, such as future number of changes (Yazdi et al., 2014; Kenmei et al., 2008; Chug and Malhotra, 2016; Elish and Elish, 2009), software defects (Raja et al., 2009; Nagappan et al., 2006), fault-proneness (Arisholm and Briand, 2006; Goulão et al., 2012; Gondra, 2008; Khoshgoftaar et al., 2002), code smells (Fontana et al., 2016a), and vulnerabilities (Roumani et al., 2015). The majority of these methods try to approach the subject by applying time series or ML prediction models on individual software properties based on the analysis of available information (historical data, trends, source code metrics, etc.).

A commonly used technique to analyze the evolution of software systems is time series analysis. In their study, Yazdi et al. (2014) model the evolution of the design of software systems by applying ARMA time series to several typical projects successfully. Based on the empirical results the authors point out that time series models can predict the future changes of the next revisions of the systems with sufficient accuracies. In another study, Kenmei et al. (2008) use time series models to forecast future change requests evolution and to identify trends based on data collected from three large open source applications. They highlight that time series are capable to model change requests and act as a support tool for project staffing and planning. Likewise, Raja et al. (2009) use the time series approach to predict defects in software evolution. They use defect reports for eight open source projects and build time series models to predict software defects which lead to the conclusion that the model may be

⁴ <https://www.castsoftware.com/>

⁵ inFusion tool is no longer supported and has been evolved into <http://www.aireviewer.com/>.

⁶ <https://www.vector.com/int/en/products/products-a-z/software/square/>

⁷ <https://www.sonarqube.org/>

⁸ <https://www.softwareimprovementgroup.com/solutions/sigrid-software-assurance-platform/>

⁹ <http://structure101.com/products/workspace/>

¹⁰ <https://www.ndepend.com/>

¹¹ <https://www.cqse.eu/en/products/teamscale/landing/>

used to facilitate planning for software evolution budget and time allocation. Finally, [Goulão et al. \(2012\)](#) build a time series model to forecast the change requests evolution based on data collected from Eclipse's change request tracking system. Additionally, they include the identification of seasonal patterns and tendencies, which is important to validate that usage of seasonal information significantly improves the estimation ability of this model, when compared to other ARIMA models.

In addition to time series analysis, multiple studies address the problem of forecasting the evolution of various aspects of software quality by employing ML techniques. In their study, [Chug and Malhotra \(2016\)](#) introduce a benchmarking framework for predicting the number of changes, and therefore the maintainability of a software application, using OO metrics as predictors. Through their framework, they compare the effectiveness of 17 ML techniques (including linear regression, decision trees, SVM and genetic algorithms) over seven open source systems. They conclude that although good predictive performance is achieved by almost all ML techniques, the genetically adaptive learning models perform better than the others do. In a similar study, [Elish and Elish \(2009\)](#) compare various ML techniques, such as multivariate linear regression, SVM, ANN, TreeNet, and regression trees, also for predicting maintainability through the number of line changes. Their results indicate that competitive prediction accuracy is achieved when applying the TreeNet model. In the same way, [Fontana et al. \(2016a\)](#) compare 16 different supervised ML techniques for code smell detection using 74 software systems. They report that the highest performance is obtained by using J48 and Random Forest algorithms, while code smells can be detected with very high accuracy. Regarding fault-proneness prediction, in a study conducted by [Arisholm and Briand \(2006\)](#), the authors propose a multivariate regression model for predicting fault-prone components of object-oriented legacy systems by using history change and fault data from previous releases. Moreover, in ([Gondra, 2008](#)), Gondra et al. propose the use of ML to predict software fault-proneness. Their approach first employs sensitivity analysis to select software metrics that are more likely to indicate the existence of errors, and afterward, trains an ANN to predict future fault-proneness. In a relative study, [Nagappan et al. \(2006\)](#) use principal component analysis on code metrics to build regression models that accurately predict the likelihood of post-release defects. Finally, [Khoshgoftaar et al. \(2002\)](#) use regression and classification trees to identify fault and non-fault prone modules on multiple releases of a large scale legacy telecommunications system, concluding that these algorithms result in predictions with satisfactory accuracy and robustness.

The multitude of models that are available in the literature for predicting the evolution of specific quality attributes and quality properties reveal the importance of quality prediction and forecasting in the software engineering community. However, with the evolution of a software system, accumulated TD is evolving as well. Since TD is an indicator of software quality (with an emphasis on maintainability), predicting its future value is considered equally important. Various studies have focused on analyzing the evolution of TD and its impact on software development, from different perspectives ([Digkas et al., 2017; Ampatzoglou et al., 2018; Digkas et al., 2018; Tan et al., 2018; Ampatzoglou et al., 2015b; Chatzigeorgiou et al., 2015](#)). In their study, [Ampatzoglou et al. \(2015b\)](#) highlight the need for knowing TD evolution, while stressing the need for project managers to be able to preserve a software product maintainable for as long as possible. For that purpose, [Chatzigeorgiou et al. \(2015\)](#) introduce the term "breaking point", which refers to the point in time when the accumulated interest will be equal to the TD principal, i.e., the cost becomes higher than the benefit. Trying to expand this work, [Ampatzoglou et al. \(2018\)](#) instantiate and

validate FITTED, a framework that assesses the breaking point of source code modules to support decision making with respect to investments on improving quality of a software, thus providing managers with an insightful decision-making tool. Hence, forecasting the evolution of TD principal and interest could be valuable for estimating the point in which the software product could become unmaintainable.

To effectively predict how the TD of a software system will progress in the future in order to improve the TD repayment strategy, it is necessary to constantly monitor and analyze its evolution. While the previously mentioned studies indicate that there has been extensive research with respect to forecasting the evolution of quality attributes and properties, directly or indirectly related to TD, only a few contributions exist so far regarding TD forecasting ([Skourletopoulos et al., 2014; Tsoukalas et al., 2019](#)), indicating that it is a scarcely investigated field. The need for forecasting the evolution of TD has been highlighted by a recent study by [Tsoukalas et al. \(2018\)](#), in which the authors raise the awareness of the gap in the field of TD. They claim that an interesting topic would be to investigate different efficient ways to produce TD forecasting models for accurate prediction of TD principal and interest evolution. In addition, they stress that it would be useful to examine if TD forecasting could foster the development of high-quality software products. In a first attempt towards this issue, [Skourletopoulos et al. \(2014\)](#) introduce the concept of predicting TD for Software as a Service (SaaS) systems, by exploiting COCOMO, a software cost model proposed by [Boehm et al. \(1984\)](#). However, their study is limited only to cloud computing systems. In another study ([Tsoukalas et al., 2019](#)), the authors empirically evaluate the ability of time series analysis to model and predict TD evolution in long-lived, open-source software projects. They find that the autoregressive integrated moving average model (ARIMA) can provide accurate predictions over a fairly long period of up to 8 weeks. However, they observe that predictive power decreases considerably for longer forecasting horizons.

Under those circumstances, being able to forecast not only the evolution of software quality but also the evolution of TD principal and interest of a software system in the future is of great significance and value. Through our study, we identified some interesting open issues that should be addressed through further research. In particular, no concrete contributions exist in the related literature regarding TD forecasting, while there is still a large volume of potential metrics and techniques that have not been used and that could potentially enhance the completeness of the software quality forecasting concept. Such a work would enable project managers and developers to support decision-making in uncertainty and plan precise payback strategies, in order to manage TD promptly and avoid unforeseen situations long-term.

4. Data definition, collection and preparation

For the execution of this study, we aimed at combining different TD-related features and metrics into a common dataset (source triangulation) with the purpose of investigating if and to what extent multivariate ML models can be used in order to accurately predict the TD evolution of software applications. This section describes in detail the definition, collection and pre-processing of the dataset that was used later as input by the produced TD forecasting models. As a first step towards creating the TD-related dataset, we studied the literature and selected an initial set of TD indicators. As soon as the appropriate TD indicators were selected, we downloaded multiple consecutive snapshots (commits) of 15 open-source projects and then used

the SonarQube¹² and CKJM Extended¹³ tools to extract these indicators, along with the TD Principal value of each snapshot. Once the data collection step had finished, we performed data pre-processing on the collected data. Techniques such as descriptive statistics, correlation analysis, and feature selection were applied on the dataset to prepare it as an input for forecasting models. Finally, we restructured each application specific dataset to a format that can be used as input to the forecasting models. In what follows, the above procedure is presented in detail.

4.1. TD indicator definition

TD indicators allow to discover TD items by analyzing different artifacts created during the SDLC. Most TD indicators proposed in the literature are related to software metrics (Li et al., 2015; Alves et al., 2016) that allow the assessment of attributes, features, or characteristics of software artifacts. In the context of object-oriented (OO) programming, various sets of metrics, such as the metric suit proposed by Chidamber and Kemerer (1994) or the Quality Model for Object Oriented Design (QMOOD) (Bansya and Davis, 2002), make it possible to characterize the size, complexity, coupling and cohesion of the code among others. These metrics have been widely used in the literature to predict maintenance effort and maintainability (Riaz et al., 2009; Fioravanti and Nesi, 2001), which is the quality attribute that is most closely related to TD. Besides OO metrics, code smells are also a well-known indicator of the presence of code TD (Alves et al., 2016; Palomba et al., 2018). Code smells are warning signs indicating possible deeper problems in the design or code of software, often resulting from the violation of at least one programming principle (Fowler, 1999). These problems may impede the software maintenance process and impose the need for code refactoring (Fontana et al., 2012). In addition, Automatic Static Analysis (ASA) tools, such as FindBugs¹⁴ or Checkstyle¹⁵, are also widely used to indicate TD (Vetro', 2012; Zazwarka et al., 2014; Izurieta et al., 2012). ASA tools allow the analysis of source code in search for bugs or violations of good programming practices that can cause failures or quality decay of the software. Most of these violations can be removed through refactoring to avoid unforeseen problematic situations (Zazwarka et al., 2013).

As discussed above, OO metrics, code smells, issues extracted from ASA tools, and software quality metrics extracted from quality assessment tools have been widely used in the literature as indicators able to monitor and quantify TD and the quality of software maintainability in general. In the approach presented in this paper, we treat these indicators coming from different sources as potential TD predictors (source triangulation) and combine them with already existing forecasting methods to develop novel models that provide predictions about the future evolution of TD in a software application. Two of the most popular and widely used tools for calculating such TD indicators are SonarQube and CKJM Extended. SonarQube is an open source platform for continuous inspection of code quality that provides analysis functionalities and a wide range of metrics for measuring quality attributes of code, tests, and design. As of today, it has been adopted by more than 120K organizations¹⁶ including nearly more than 100K public open-source projects.¹⁷ In this study, SonarQube has been used as proof of concept for research purposes, since according to two recent studies on Technical Debt Management (Li et al., 2015;

Ampatzoglou et al., 2015a), it is the most frequently used tool for estimating TD principal. In addition, another reason for selecting this tool is the fact that it is highly customizable, allowing the users to adjust the standard out-of-the-box set of rules (named “sonar way”) that it provides in order to better meet their needs. In a relevant study (Lenarduzzi et al., 2019), the authors suggest that companies should continuously re-consider the adopted SonarQube rules based on business’s objectives and preferences. In a similar way, developers and users of the TD Forecasting tool described in this study could fine-tune the rule-set of SonarQube prior to obtaining TD-related measurements, so that predictions can be tailored based on the company’s critical needs.

Therefore, in the present work, we opted for the TD-related metrics¹⁸ that are provided by SonarQube, as our primary TD Principal predictors. The version of SonarQube used within the context of this work is 6.7.4. Furthermore, SonarQube was also used to compute the target variable, i.e., to quantify the TD Principal of the selected software applications. To do so, SonarQube checks code compliance against a set of classified coding rules and if the code violates any of these rules, it considers it as a violation or a TD item. For each of the identified TD items, SonarQube computes the remediation time (i.e., estimated effort) needed to refactor it and considers it as TD.

To complement the TD predictor set we decided to account also for the popular Chidamber and Kemerer (C&K) metrics and Quality Model for Object Oriented Design (QMOOD) metrics (Chidamber and Kemerer, 1994; Bansya and Davis, 2002). The reason behind this choice is that C&K metrics, such as DIT, NOC, RFC, LCOM and WMC, and QMOOD metrics, such as DAM, MOA, and CAM have been intensively studied for their ability to predict maintainability and maintenance effort (Riaz et al., 2009; Fioravanti and Nesi, 2001) (see “Studies” column in Table 1). One of the main limitations of SonarQube tool is the lack of OO detection mechanisms. Therefore, to collect OO metrics for our applications, we chose the popular CKJM Extended (Jureczko and Spinellis, 2010), an extended version of the CKJM open-source tool able to calculate a wide range of metrics¹⁹ (including those defined in C&K and QMOOD suites), by processing the bytecode of Java files. Indicators extracted by CKJM Extended can be calculated at the source-code level, and can be used to assess well-known quality properties associated with the architecture of a software application, such as complexity, coupling, cohesion, and inheritance among others. In addition, CKJM Extended calculates C&K metrics strictly according to the original (1994) definition by Chidamber and Kemerer.

In Table 1, the metrics that were selected as TD indicators and therefore used as independent variables for the creation of our dataset are presented along with a short description. Moreover, to strengthen the TD indicators selection, we provide references to studies that relate each metric with TD and the quality of software maintainability in general. The first half of the table describes metrics computed by SonarQube, while the last half describes metrics extracted by CKJM Extended. The target variable, i.e., the variable that we want to forecast, is denoted here as *total_principal*. We define *total_principal* as the effort (in minutes) to fix all issues and we compute it as the sum of code smell, bug, and vulnerability remediation effort.

At this point, it should be noted that an obvious choice of related metrics to be included as independent variables in the multivariate models that we investigate in this work are the constituent components of TD Principal (i.e., *total_principal*) as computed by SonarQube, i.e., the code smells, bugs, and vulnerabilities. However, apart from the TD Principal constituent

¹² <https://www.sonarqube.org/>

¹³ http://gromit.iiar.pwr.wroc.pl/p_inf/ckjm/

¹⁴ <http://findbugs.sourceforge.net/>

¹⁵ <https://checkstyle.sourceforge.io/>

¹⁶ <https://www.sonarqube.org/>

¹⁷ <https://sonarcloud.io/explore/projects>

¹⁸ <https://docs.sonarqube.org/latest/user-guide/metric-definitions/>

¹⁹ http://gromit.iiar.pwr.wroc.pl/p_inf/ckjm/metric.html

Table 1
TD indicators.

Metric	Description	studies
Project-level metrics (computed by SonarQube)		
		<i>Technical debt metrics</i>
sqale_index	Effort to fix all code smell issues. The measure is stored in minutes.	
reliability_remediation_effort	Effort to fix all bug issues. The measure is stored in minutes.	
security_remediation_effort	Effort to fix all vulnerability issues. The measure is stored in minutes.	
total_principal	Effort to fix all issues. The sum of the three metrics mentioned above, i.e., code smell, bug and vulnerability remediation effort. The measure is stored in minutes.	
		<i>Reliability metrics</i>
bugs	Total number of bug issues of a project.	Xuan et al. (2017) Griffith et al. (2014) and Digkas et al. (2017)
		<i>Security metrics</i>
vulnerabilities	Total number of vulnerability issues of a project.	Griffith et al. (2014) Digkas et al. (2017) and Siavvas et al. (2019)
		<i>Maintainability metrics</i>
code_smells	Total number of code smell issues of a project.	Digkas et al. (2017), Fontana et al. (2012) Zazworska et al. (2011) Olbrich et al. (2010) Charalampidou et al. (2017) Sjøberg et al. (2012) and Mamun et al. (2019)
		<i>Size metrics</i>
comment_lines	Number of lines containing either comment or commented-out code of a project.	Griffith et al. (2014) and Karus and Dumas (2012)
ncloc	Number of physical lines of a project that contain at least one character, which is neither a whitespace nor a tabulation nor part of a comment.	Digkas et al. (2017) Mamun et al. (2019) Lucca et al. (2004) Eski and Buzluca (2011) and Elish and Elish (2009)
		<i>Coverage metrics</i>
uncovered_lines	Number of lines of code of a project, which are not covered by unit tests.	Griffith et al. (2014)
		<i>Duplication Metrics</i>
duplicated_blocks	Number of duplicated blocks of lines of a project.	(Griffith et al., 2014) (Nugroho et al., 2011) and (Marinescu, 2012)
		<i>Complexity metrics</i>
complexity	The Cyclomatic Complexity of a project calculated based on the number of paths through the code.	Giger et al. (2012) Bruntink and van Deursen (2006) and Singh and Saha (2012)
Class-level metrics aggregated at project-level (computed by CKJM Extended)		
		<i>Complexity metrics</i>
AMC	Average method complexity: The average method size for each class (number of Java binary codes in the method), averaged for all project classes.	Karus and Dumas (2012) and Zhou and Xu (2008)
WMC	Weighted Methods per Class: The total number of methods that a class contains weighted by their complexity values, averaged for all project classes.	Eski and Buzluca (2011) Giger et al. (2012) Singh and Saha (2012) Van Koten and Gray (2006) Shatnawi and Li (2008) Zhou and Leung (2007) and Elish and Elish (2009)
DIT	Depth of Inheritance tree: The depth of inheritance tree for each class from the object hierarchy top, averaged for all project classes.	Giger et al. (2012) Singh and Saha (2012) Van Koten and Gray (2006) Shatnawi and Li (2008) (Zhou and Leung, 2007) (Zhou et al., 2012) and (Elish and Elish, 2009)

(continued on next page)

components that have an evident effect on the target variable itself, we decided to investigate also other, not directly related metrics (presented in Table 1), which however are known to act as TD indicators, in order to examine whether they have an

equally (or more) significant impact on TD Principal. To this end, the extensive feature selection analysis reported in Section 4.3 will allow us to come up with the best predictors set, tailored to our dataset and experimental setup.

Table 1 (continued).

Metric	Description	studies
NOC	Number of children: The number of immediate descendants (i.e., children) of a class, averaged for all project classes.	Giger et al. (2012) Bruntink and van Deursen (2006) Shatnawi and Li (2008) Zhou and Leung (2007) Zhou et al. (2012) and Elish and Elish (2009)
RFC	Response for a class: The number of local methods plus the number of methods called by class methods, averaged for all project classes.	Eski and Buzluca (2011) Giger et al. (2012) Bruntink and van Deursen (2006) Singh and Saha (2012) Van Koten and Gray (2006) Shatnawi and Li (2008) Zhou and Leung (2007) and Zhou et al. (2012)
<i>Coupling metrics</i>		
CBO	Coupling between objects: The total number of classes coupled to a given class, averaged for all project classes.	Eski and Buzluca (2011) Giger et al. (2012) Singh and Saha (2012) Shatnawi and Li (2008) and Zhou et al. (2012)
Ca	Afferent coupling: The total number of other classes that call methods of the given class, averaged for all project classes.	Karus and Dumas (2012) Zhou et al. (2012) and Elish (2010)
Ce	Efferent coupling: The total number of other classes that their methods are called by the given class, averaged for all project classes.	Karus and Dumas (2012) Zhou et al. (2012) and Elish (2010)
CBM	Coupling between methods: The total number of parent classes to which a given class is coupled, averaged for all project classes.	Eski and Buzluca (2011) Giger et al. (2012) and Singh and Saha (2012)
IC	Inheritance coupling: The total number of new or redefined methods of a class to which all its inherited methods are coupled, averaged for all project classes.	Eski and Buzluca (2011) and Giger et al. (2012)
<i>Cohesion metrics</i>		
LCOM	Lack of Cohesion in Methods: The number of methods pairs in a class that are not interrelated through the sharing of some of the class fields, averaged for all project classes.	(Eski and Buzluca, 2011) Giger et al. (2012) Van Koten and Gray (2006) Shatnawi and Li (2008) Zhou and Leung (2007) Zhou et al. (2012) and Elish and Elish (2009)
LCOM3	Lack of Cohesion in Methods: Similar to LCOM but ranging from 0 to 2.	Singh and Saha (2012) and Zhou et al. (2012)
CAM	Cohesion Among Methods: This metric computes the relatedness among methods of a class based on their parameter lists, averaged for all project classes (Range 0 to 1).	Eski and Buzluca (2011)
<i>Other metrics</i>		
NPM	Number of Public Methods: The total number of methods in a class that are declared as public, averaged for all project classes.	Bruntink and van Deursen (2006) Singh and Saha (2012) Van Koten and Gray (2006) Zhou and Leung (2007) and Elish and Elish (2009)
DAM	Data Access Metric: The ratio of the number of private or protected fields to the total number of fields declared in the class, averaged for all project classes (Range 0 to 1).	(Kaur and Anand, 2013) and Goyal and Joshi (2014)
MOA	Measure of Aggregation: The total number of data declarations (class fields) whose types are user defined classes, averaged for all project classes.	Kaur and Anand (2013) and Goyal and Joshi (2014)

4.2. Collection of data

To start the dataset construction process, we initially selected 15 popular open-source applications from the GitHub²⁰ repository. The selected 15 applications have different sizes and belong to different application domains, which range from Networking Software (e.g., Kafka, Dubbo, OKHttp, Retrofit, Openfire, WebSocket) to Business Software (e.g., OFBiz), and from Scientific Software (e.g., SystemML) to Utilities Software (e.g., Commons IO, Guava, Jenkins, ZXing). The selection criteria were based on the software popularity, activity level, data availability, and the

Java programming language. More specifically, we exploited the advanced search mechanism provided by GitHub by selecting only Java projects in the “Languages” filter, and then sorting the results based on “Most Stars” filter. This resulted to an initial set of Java applications ranked by their popularity. Next, to assess the activity level of each application, we used the “Insights” GitHub functionality to examine the total number of commits and compare it to the lifespan of the application. We selected only applications whose commit activity was frequent (at least once per week) and long-lived (at least 3 years). Moreover, since SonarQube and CKJM Extended both work on compiled classes to compute metrics values, selected projects needed to be compilable, i.e., not producing any errors during compiling.

²⁰ <https://github.com>

For each application that met the above criteria, approximately 150 snapshots (commits) in weekly intervals were fetched, spanning up to 3 years of each project's evolution. More precisely, we opted for the last commit in every analyzed week as the time point of analysis. The rationale behind this option, i.e., to ensure fixed and weekly time intervals between the commits is twofold. First, ensuring fixed time distance between the retrieved samples (i.e., commits) is critical for the reliability of the produced forecasting models. Secondly, collecting snapshots at weekly rather than daily intervals is a more viable solution as rarely do projects keep daily commits. In addition to this, another reason that led us to the decision to take snapshots at weekly rather than daily intervals or even to analyze all consecutive commits, is to avoid as many periods of inactivity as possible, but not eliminate them. Periods of inactivity are consecutive snapshots in which no significant new changes were committed to a codebase. While including many such snapshots could potentially introduce a high level of noise in the dataset, including a reasonable number of 'low or no activity' periods is important for an accurate model: Generating forecasts that (potentially) indicate future periods of inactivity could prove useful in practice for project managers in decision-making activities. Choosing longer intervals (e.g., monthly) would probably reduce the periods of inactivity even further, but it would result in significantly fewer data and thus significantly lower forecasting performance, whereas the produced models would be able to provide forecasts only at a monthly basis and not for shorter periods (e.g., some weeks ahead), which would restrict their practicality in decision-making.

The approach described above led to a codebase containing up to 1850 snapshots in total (171M lines of code). A sufficiently high number of applications is fundamental to reach a conclusion that does not depend on a specific dataset, allowing to generalize the obtained results. For the purpose of constructing the dataset, a dedicated crawler was created. In order to facilitate the reproducibility and the extensibility of the present study, as well as the construction of similar datasets, this crawler has been made available online.²¹ In Table 2, the applications that were selected for constructing the codebase are presented in detail.

After fetching the source code of each snapshot for the 15 selected applications, we proceeded to the next step, i.e., using SonarQube and CKJM Extended (as described in Section 4.1) in order to analyze each snapshot and build 15 application-specific datasets consisting of the TD indicators described in Table 1. We chose the format of each application-specific dataset to be the following: each row contains a specific snapshot of the application in chronological order (time series), whereas the columns contain the values of the TD indicators, plus one column containing the value of total TD principal for that particular snapshot. This format helped us also during the forecasting model construction phase described later.

Since the work presented in this paper aims at modeling TD evolution of the entire software project (system), rather than predicting the TD of individual software artifacts (e.g., classes), we performed data collection for each application at the system level. In other words, each application snapshot (commit) provided a single observation in the application-specific dataset. TD-related metrics extracted by SonarQube analysis are computed at system-level by default, so no further modifications were needed. However, most of the metrics extracted by CKJM Extended, such as *DIT*, *RFC*, and *NOC*, are originally defined at the class level. Therefore, those metrics could not be directly used as independent variables. For this purpose, we aggregated CKJM metrics at system-level, i.e., we used their weighted mean among

classes. More specifically, in our approach the system-level value of each metric is the aggregation of its class-level values weighted by the lines of code of each class, divided by the total lines of code of the system under analysis. This aggregation approach has been used in relevant studies by Baggen et al. (2012) and Wagner et al. (2015), but also in some of our previous studies as well (Siavvas et al., 2017a,b).

4.3. Data preparation

Selection of independent (input) variables is a critical part in the design of a ML algorithm. Each additional input unit adds another dimension and contributes to the "curse of dimensionality" (Bellman, 2003), a phenomenon in which performance degrades as the number of inputs increases. Furthermore, irrelevant or partially relevant features can negatively impact model performance. Thus, after constructing our dataset, the next step is to provide a clear understanding of the statistical attributes of our variables, and then to reduce the number of input variables described in Table 1 by keeping only the most important ones, i.e., the ones that are highly significant for TD Principal forecasting. Techniques like correlation analysis, univariate and multivariate analysis will allow us to retain as much discriminatory information as possible.

In order to study the statistical significance of each indicator over the TD quality and be able to safely perform dimensionality reduction of our dataset, also known as feature selection, we need to maximize diversity and representativeness by considering a comparable number of different heterogeneous applications. The dataset we constructed and described in Section 4.2 for forecasting purposes consists of 15 applications, a number that may not be suitable for generalizing our findings and reach to a generic conclusion regarding feature selection. Therefore, to increase the size of our dataset for feature selection purposes, in addition to those applications, we have also exploited a benchmark repository that consists of the 100 most popular Java libraries (e.g., Junit, Xerces, HyperSQL, etc.) retrieved from the Maven Repository.²² The same dataset was used in the study by Siavvas et al. (2017a) for calibrating a Quality Assessment Model, as well as in a similar study (Siavvas et al., 2017b) for investigating the interrelationship of software metrics and specific vulnerability types. For the purpose of this work, we further extended the repository by adding 110 more Maven applications, based on their popularity. As a result, the final benchmark repository contains 210 open-source software Java applications, comprising approximately 30 million lines of code, which is considered an adequate number for the purpose of identifying most significant TD indicators. To extract all required indicators from the extended dataset, as in the case of the initial 15 applications, we analyzed the source code of each application using SonarQube and CKJM Extended, as described in Section 4.2. Finally, we merged these metrics with the metrics obtained from the last snapshot of each of our analyzed applications presented in Table 2, leading to an extended dataset containing metrics from a total of 225 applications. We chose to add only the last snapshot of the analyzed applications presented in Table 2 in order to ensure equal representativeness, since each software application in the benchmark repository is represented by only one commit. Furthermore, we assumed that the last snapshot of each project would be more mature and bigger in size compared to its previous snapshots. This extended dataset can be found online.²³

At this point, it should be noted that the extended dataset was used only for correlation analysis and feature selection purposes,

²² <https://mvnrepository.com/>

²³ <https://sites.google.com/view/technical-debt-forecasting/main>

²¹ <https://sites.google.com/view/technical-debt-forecasting/main>

Table 2

Applications of the TD dataset.

Application name	Analyzed weekly snapshots	Last snapshot LOC	Total commits	GitHub contributors	GitHub stars	Description
#	Timeframe					
apache/kafka	150	30/10/2015 – 07/09/2018	116.000	7.055	621	14.8k
apache/commons-io	150	11/09/2015 – 27/07/2018	30.000	2.303	53	657
apache/ofbiz	100	04/11/2016 – 28/09/2018	243.000	24.427	20	678
apache/systemml	150	02/10/2015 – 10/08/2018	200.000	6.039	59	798
apache/groovy	150	25/12/2015 – 02/11/2018	210.000	16.806	290	3.6k
apache/nifi	100	04/11/2016 – 28/09/2018	289.000	5.599	280	1.9k
apache/incubator-dubbo	100	28/07/2017 – 01/02/2019	67.000	4.139	279	20.3k
google/guava	150	25/12/2015 – 02/11/2018	114.000	5.190	219	35.8k
square/okhttp	150	18/12/2015 – 26/10/2018	24.000	4.438	198	35.8k
square/retrofit	100	27/01/2017 – 21/12/2018	7.900	1.782	131	34.8k
jenkinsci/jenkins	150	25/03/2016 – 01/02/2019	147.000	29.265	599	14.8k
spring-projects/spring-boot	100	04/11/2016 – 28/09/2018	15.000	25.004	646	28.6k
TooTallNate/java-WebSocket	100	17/03/2017 – 25/01/2019	5.200	954	62	1.9k
zxing/zxing	100	10/03/2017 – 01/02/2019	29.000	3.524	96	24.6k
igniterealtime/Openfire	100	18/11/2016 – 12/10/2018	100.000	9.214	114	2.1k

as the additional 210 applications obtained from the benchmark repository do not contain project history (past commits), and therefore, are not suitable for forecasting model experiments. After feature analysis and selection, we switched back to our original dataset, containing 15 applications with their commit history (1850 commits in total) for forecasting model training.

4.3.1. Descriptive statistics

Descriptive statistics is the term given to the analysis of data that helps describe data in a meaningful way, allowing for simpler interpretation. It provides simple summaries about the sample and about the observations that have been made. Descriptive statistics include measures of central tendency, such as the mean, median, and mode, and measures of variability, such as standard deviation, variance, the minimum and maximum variables, and the kurtosis and skewness. After extracting the metrics of each application (using SonarQube and CKJM Extended) and merging them into a common dataset as described in the previous section, the descriptive statistics of TD indicators calculated based on the extended dataset are presented in Table 3. As a reminder, the extended dataset comprises a superset of the dataset that was constructed for forecasting purposes (i.e., the 15 applications presented in Table 2) and the additional 210 applications that were

added at a later stage for feature selection purposes. Therefore, for the computation of the statistical metrics presented in Table 3 we have included all 225 applications. For the conduction of our experiments, we used the Python programming language and more specifically the Pandas²⁴ data analysis library.

Metrics that vary little are not likely to be useful predictors. In our case, from Table 3 we observed that for all metrics there are significant differences between the lower 25th (lower) percentile, the median, and the 75th (upper) percentile, thus showing strong variations. Therefore, all metrics were selected to be used for subsequent analysis. We also observed that, as it is the case with software engineering data (Chidamber and Kemerer, 1994), most of our metrics are highly skewed, which means that few outlier observations may substantially affect the results, if not treated carefully. To mitigate this risk, in the rest of the analysis that follows we opted for techniques that perform well when the distribution of values in the feature space cannot be assumed.

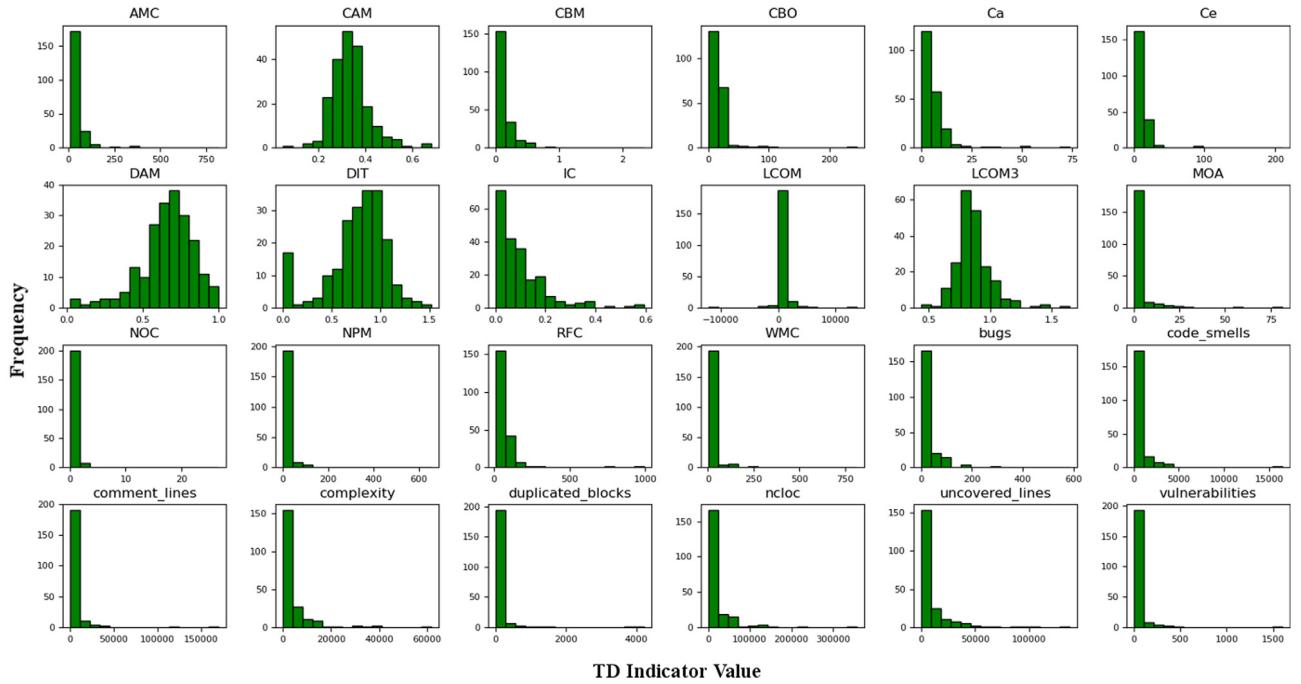
In Fig. 2, histograms of each metric are presented to further complement our initial analysis. We used histograms to further examine the normality and skewness of each metric. A normal

²⁴ <https://pandas.pydata.org/>

Table 3

Descriptive statistics of TD indicators (extended dataset).

Metric	Mean value	Standard deviation	Minimum value	Lower quartile	Median value	Upper quartile	Maximum value	Skewness	Kurtosis
bugs	29.665	61.148	0	3	9	30	585	5.286	37.601
vulnerabilities	38.818	125.935	0	1	7	32	1598	9.675	114.507
code_smells	858.660	2095.092	2	94	235	732	16442	5.551	35.262
comment_lines	5113.780	15587.257	1	387	1381	3889	170698	8.203	77.661
ncloc	19524.278	37253.222	175	2655	7068	20438	357664	5.262	38.091
uncovered_lines	9815.507	17856.102	32	1302	3297	10288	137326	3.905	19.212
duplicated_blocks	110.177	429.327	0	0	12	47	4219	7.771	67.183
complexity	3914.024	7420.959	11	421	1352	4284	61862	4.275	23.544
AMC	53.370	74.382	9.407	25.374	36.565	52.876	819.191	6.697	58.288
WMC	36.399	82.087	2.085	12.638	18.921	31.524	814.007	7.136	57.322
DIT	0.760	0.309	0	0.655	0.813	0.964	1.523	-0.936	0.977
NOC	0.527	1.900	0	0.080	0.261	0.481	26.770	12.828	176.915
RFC	71.453	90.004	11.618	39.508	54.295	78.112	993.299	7.631	69.342
CBO	17.177	21.071	0	9.056	13.554	18.983	245.469	7.139	68.779
Ca	6.243	8.145	0	2.463	3.850	7.399	73.761	4.974	31.903
Ce	12.425	18.100	0	5.922	9.664	13.024	210.225	7.710	73.955
CBM	0.143	0.231	0	0.026	0.085	0.164	2.340	5.366	41.706
IC	0.098	0.106	0	0.025	0.071	0.137	0.595	2.054	5.452
LCOM	442.525	1552.824	-12160.67	61.584	153.151	503.916	13693.182	0.885	46.386
LCOM3	0.863	0.151	0.435	0.776	0.840	0.916	1.648	1.377	5.060
CAM	0.337	0.081	0.050	0.278	0.333	0.377	0.683	0.858	3.011
NPM	23.251	51.913	0.250	8.009	12.287	21.058	655.642	9.550	109.341
DAM	0.661	0.178	0.024	0.574	0.682	0.779	1.000	-0.952	1.526
MOA	3.579	7.849	0	0.761	1.591	3.263	81.424	6.701	56.123

**Fig. 2.** Histograms of TD indicators (extended dataset).

distribution is symmetric and bell-shaped. We observed that most of the metrics are not normally distributed.

To further validate this finding we used boxplots, which can be found online²⁵ as supportive material. Boxplots provide a standardized way of displaying the distribution of data. These graphs divide the dataset into a five-number summary: the minimum, first quartile (Q1), median, third quartile (Q3), and maximum. Outlier observations are often easy to identify by inspecting a boxplot, since they are plotted as individual points lying outside the boundaries set by the minimum and maximum values. We observed that some of the metrics seem to have outlier

observations. We further investigated these findings during the univariate (one variable outlier) and multivariate (two or more variable outlier) analysis described below.

4.3.2. Correlation analysis

As previously stated, performance of a ML algorithm degrades as the number of inputs increases. Irrelevant or partially relevant features can negatively affect model performance. In order to successfully reduce the number of input variables by keeping only the most important ones, we first applied Spearman's rank correlation coefficient (ρ) (Spearman, 1987) analysis between TD Principal as computed by SonarQube and each TD indicator described in Table 3 for the 225 applications. Spearman's rank correlation was selected, as it is a nonparametric test that is

²⁵ <https://sites.google.com/view/technical-debt-forecasting/main>

not sensitive to outliers. Additionally, it does not assume any distribution for the studied data, which is important in our case, as our data did not seem to follow any known distribution. To interpret the strength of the correlations Cohen (1977) suggestion was used. According to Coehen et al. a correlation less than 0.3 is considered weak, between 0.3 and 0.5 is considered moderate, and above 0.5 is considered strong. Finally, to ensure that the observed associations did not occur by chance, the correlations were tested for statistical significance. For this purpose, the *p-value* of each correlation was examined. A *p-value* of 0.05 means that we are 95% confident that the observed association has not occurred by chance. Hence, we examined the statistical significance of the observed correlations at 95% level of confidence.

In Fig. 3, the correlation between each metric is illustrated based on color warmness, i.e., the more red a box, the higher the correlation between the corresponding metrics. We focused on the last row, which represents the correlation between our dependent variable, i.e., TD Principal and each independent variable, i.e., TD indicators.

To further complement Spearman's correlation analysis, the correlations between TD and each TD indicator, as well as the significance (*p-value*) of each correlation are presented in numbers in Table 4. To facilitate the readability of the correlation table, alphabetical symbols are used to denote the strength of each correlation based on the Coehen's et al. suggestion described above. In particular, the values marked with (a), (b) and (c) correspond to weak, medium and strong correlations respectively. In addition, statistically significant *p-values* ($p = < 0.05$) are marked in bold, while not statistically significant *p-values* ($p > 0.05$) are in regular font.

As a first step towards feature selection, metrics that have either a low correlation score (i.e., correlation < 0.3), or a non-significant statistical correlation (i.e., *p-value* > 0.05) with respect to TD were marked as candidates for removal. In particular, five metrics were identified as non-correlated (i.e., CBM, IC, LCOM3, CAM, DAM), while 3 metrics had statistically insignificant correlations (i.e., RFC, CE, LCOM) with respect to TD and consequently were filtered out, leaving 16 out of 24 TD indicators for further analysis.

4.3.3. Univariate analysis

After applying the correlation analysis as the first filter towards feature selection, we considered to apply a univariate regression analysis between each remaining metric (TD indicator) and the TD for the extended dataset (225 applications). The importance of controlling for potential confounders in empirical studies of object-oriented products has been emphasized in the study by El Emam et al. (2001). Univariate regression focuses on determining the relationship between one independent variable (i.e., each metric) and the dependent variable (i.e., TD Principal) and has been widely used in software engineering studies to examine the effect of each metric separately (Arisholm and Briand, 2006; Zhou and Xu, 2008; El Emam et al., 2001). Thus, we used this method as a second filter, to help us with the process of removing metrics whose underlying relationship is not statistically significant to TD. During descriptive statistics however, we observed that most of the metrics were highly skewed. In order to render the data suitable for univariate regression analysis we applied the natural logarithm $\log(\ln)$ transformation to the values of the remaining metrics (Makridakis et al., 2018). Using the natural logarithm reduces the skew of the response and predictors (linear regression assumptions include normal distribution of the residuals).

Table 5 summarizes the results of the univariate linear regression analysis for each metric, applied on the extended dataset

Table 4
Spearman's rank correlation of TD indicators (extended dataset).

Metric	Correlation with total principal	<i>p-value</i>
Project-level metrics (computed by SonarQube)		
bugs	0.784 ^c	8.10023e-45
vulnerabilities	0.722 ^c	5.25492e-35
code_smells	0.962 ^c	1.0044e-118
comment_lines	0.838 ^c	3.77265e-91
ncloc	0.917 ^c	2.20628e-58
uncovered_lines	0.929 ^c	2.45155e-56
duplicated_blocks	0.846 ^c	1.23839e-84
complexity	0.915 ^c	2.35973e-83
Class-level metrics aggregated at project-level (computed by CKJM)		
AMC	0.349 ^b	8.79842e-10
WMC	0.408 ^b	2.83379e-09
DIT	-0.471 ^b	6.05013e-13
NOC	0.396 ^b	8.26175e-09
RFC	0.463 ^b	0.106134
CBO	0.568 ^c	0.0376558
Ca	0.562 ^c	2.88947e-19
Ce	0.492 ^b	0.442171
CBM	0.053 ^a	7.94597e-19
IC	0.003 ^a	3.53632e-14
LCOM	0.385 ^b	0.958871
LCOM3	0.112 ^a	6.03566e-07
CAM	-0.144 ^a	2.4988e-11
NPM	0.440 ^b	0.015035
DAM	-0.168 ^a	2.09667e-07
MOA	0.337 ^b	1.66596e-12

^aWeak correlation

^bMedium correlation

^cStrong correlation

Table 5
Univariate analysis results of TD indicators (extended dataset).

Metric	R2	<i>p-value</i>	Standard error	Relationship
bugs	0.625	0.000	0.048	+
vulnerabilities	0.528	0.000	0.045	+
code_smells	0.931	0.000	0.019	+
comment_lines	0.685	0.000	0.036	+
ncloc	0.823	0.000	0.033	+
uncovered_lines	0.847	0.000	0.029	+
duplicated_blocks	0.692	0.000	0.031	+
complexity	0.812	0.000	0.032	+
AMC	0.107	0.000	0.164	+
WMC	0.131	0.000	0.131	+
DIT	0.143	0.000	0.518	-
NOC	0.037	0.005	0.311	+
CBO	0.298	0.000	0.136	+
Ca	0.317	0.000	0.132	+
NPM	0.160	0.000	0.127	+
MOA	0.098	0.000	0.141	+

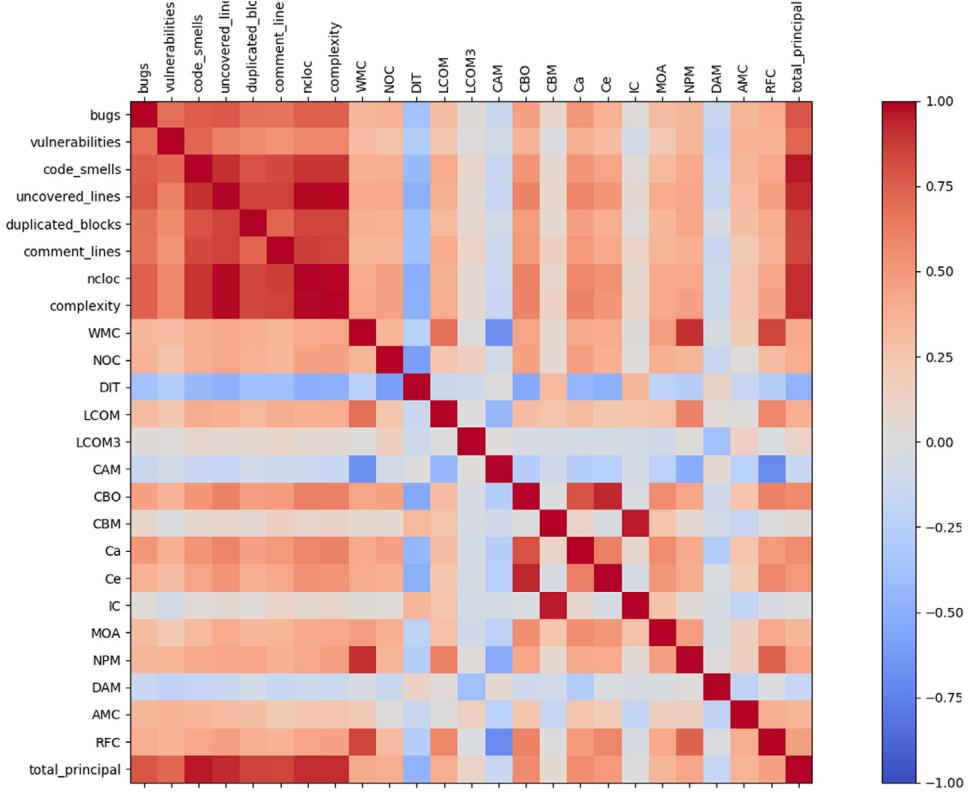


Fig. 3. Spearman's rank correlation of TD indicators (extended dataset).

(225 applications). Column “R2” gives the coefficient of determination, i.e., the proportion of the total variation in the dependent variable that is explained by the model. Columns “p-value”, “Standard error”, and “Relationship” show the statistical significance, the standard error, and the sign of the regression coefficient for the independent variable, respectively.

We set the significance level at $\alpha = 0.05$. Metrics with p -values lower than 0.05 ($p\text{-value} \leq 0.05$) are considered statistically significant to TD, and therefore can be selected as TD indicators for further analysis. On the contrary, metrics with p -values greater than 0.05 can be removed from further analysis, since they are not considered statistically significant. In our case, all 16 remaining metrics had p -values lower than 0.05. Therefore, no further metrics were dropped during this step.

4.3.4. Multivariate analysis

Since univariate analysis did not filter out any metrics, we proceeded with applying multivariate regression analysis (Efroyimson, 1960) as a final filtering step towards feature selection. While univariate analysis is used to examine the effect of each independent variable on the target variable separately, multivariate analysis examines the common effectiveness of a set of independent variables at predicting the dependent variable. Multivariate analysis is usually combined with Stepwise regression (Efroyimson, 1960), a feature selection method in which the choice of predictive variables is carried out by an automatic procedure, thus allowing for removing independent variables based on their significance (p -values). Backward Elimination, a special type of Stepwise regression, involves starting with all candidate variables, testing the deletion of each variable using a multiple linear regression, deleting the variable whose loss gives the most statistically insignificant deterioration of the model fit (i.e., highest p -value), and repeating this process until no further variables can be deleted without a statistically significant loss of fit (i.e., until all remaining variables have p -values less than the user-defined

significance level). This technique has been widely used in empirical software engineering studies to examine the effects of combined metrics on the software quality, defect or code smells prediction (Arisholm and Briand, 2006; Nagappan and Ball, 2005; Challagulla et al., 2005; Munson and Khoshgoftaar, 1990; Khoshgoftaar and Munson, 1990). During this step, all metrics reported in Table 5 were examined since all of them were found to be statistically significant (i.e., $p\text{-value} < 0.05$) during the univariate regression analysis.

While one can argue that another round of feature selection might be unnecessary, we decided to perform it mainly for two reasons. First, this additional filtering layer will capture instances where an independent variable that was found to be significant with respect to the dependent variable while being independently examined (univariate analysis), may not have significant predictive power when combined with other variables. Therefore, including it to the final set may lead to redundant information and increase model complexity. Second, the “sliding window” method described in Section 4.4 will extend each initial sample of the dataset by including past information and future information simultaneously into a single row. If for example, we decide to leave the size of independent variable set as is, i.e., 16 features, and we want to include information up to 2 lags in the past (+1 for the current lag), the final set will comprise $3 * 16 = 48$ independent variables. While ML models generally support complex relationships between variables, this data reframing approach may result in a dramatic increase in the number of features, and therefore increase the complexity of the prediction algorithms to a point where the performance drops significantly. That being said, we considered to keep our independent variables set as small as possible, but without losing much of its explanatory power.

Before starting the Backward Elimination process, a significance level has to be set. This value acts as a significance threshold that determines the stopping point, i.e., the point at which

we no longer need to drop any independent variables (i.e., predictors). In our case, we set the significance level value at 0.05 to examine the statistical significance of each TD indicator to act as a predictor at a 95% level of confidence. Subsequently, the Backward Elimination process involved performing multiple iterations of fitting a multiple linear regression model with all possible predictors, inspecting the *p-values* of each predictor, and then finding and removing the most insignificant predictor, i.e., the predictor with the highest *p-value*. As long as there was a predictor that could be removed (i.e., its *p-value* is greater than 0.05), the process was repeated by fitting a new model excluding the previously removed predictors. The process stopped when all remaining predictors had *p-values* less than the significant level of 0.05.

After multiple iterations of applying Backward Elimination, details of the final multivariate linear regression model are shown in [Table 6](#). As can be seen, the final model has four covariates, meaning that four metrics, namely *bugs*, *code smells*, *duplicated blocks*, and *afferent coupling* (*Ca*), seem to have the most significant impact on TD and act as good TD predictors, at least for the dataset under investigation (225 applications). For each covariate, we provide its coefficient, the standard error, the *t-ratio* and the statistical significance (*p-value*) of the coefficient. The *t-ratio* is the ratio of the coefficient estimate to its standard error. Since our sample is relatively large, a *t-ratio* greater than 1.96 (in absolute value) suggests that our coefficients are statistically significantly different from zero at the 95% confidence level. We observe that all remaining predictors have high *t-ratio* values (>1.96). In addition, the *p-value* for each covariate tests the null hypothesis that the coefficient is equal to zero (no effect). We observe that after performing Backward Elimination, all remaining predictors have low *p-values* (< 0.05), thus they are likely to be meaningful since the null hypothesis is rejected. The intermediate results that we obtained throughout the various iterations of the Backward Elimination process can be found online²⁶ and provide information regarding the metric that was eliminated at each iteration until reaching the final set of TD indicators shown in [Table 6](#). We also provide details regarding the coefficient, the standard error, the *t-ratio*, and the statistical significance (*p-value*) of each metric during every iteration of the process.

Finally, to strengthen the feature selection process followed, we tested the optimal TD predictors selected above for multicollinearity. Multicollinearity is a phenomenon where two or more predictors show high intercorrelations, i.e., they are highly linearly related. While correlation between a predictor and the target variable is an indication of good model performance, correlation among the predictors is usually an issue. If this issue is not taken care of during the feature selection analysis, it can later cause unpredictable variance and lead to overfitting, as the model cannot ascertain how important a feature is to the target variable. One of the most common ways to identify and quantify the severity of multicollinearity in a linear regression analysis is the Variance Inflation Factor (VIF) (Marquardt, 1970). The VIF is calculated by taking each predictor, regressing it against every other predictor in the model and then using the produced coefficient of determination (R^2) into the following formula:

$$VIF = \frac{1}{1 - R^2} \quad (1)$$

VIF values range from 1 upwards. As a rule of thumb, a VIF value between 1 and 5 indicates that a predictor is moderately correlated with the other predictors, while a value between 5 and 10 indicates that multicollinearity is likely present and thus, the predictor should be removed. We computed VIF factors for each of the predictors presented in [Table 6](#). As can be seen, all VIF values are considerably less than 5, indicating that our final TD predictor set does not suffer from multicollinearity.

Table 6
Multivariate analysis model of TD indicators (extended dataset).

Metric	Coefficient	Standard error	t-ratio	p-value	VIF
bugs	0.1075	0.027	4.056	0.000	2.503
code_smells	0.7489	0.029	25.899	0.000	3.574
duplicated_blocks	0.1276	0.020	6.494	0.000	2.752
Ca	0.1816	0.040	4.567	0.000	1.247

All the analysis described above was performed by using the Python programming language and more specifically the scikit-learn²⁷ ML library. To conclude, among the initial 24 metrics (TD indicators) under investigation, four of them were found to have statistically significant effects on TD. Therefore, the optimal TD predictors extracted through this process were *bugs*, *code smells*, *duplicated blocks*, and *afferent coupling* (*Ca*). These metrics will be considered as input to the forecasting models during the model training phase described in Section 5.

4.4. Sliding window method

In general, ML models do not directly support the notion of observations over time. As a result, time series data usually need to be re-framed in a form suitable for supervised learning problems before used for forecasting tasks. To understand this notion, an example of dummy data collected in temporal order is presented in [Table 7](#). Each row represents a sample of data collected at a specific lag (timestamp). Columns 2 to 4 hold the values of independent variables X1 to X3 respectively, while column Y1 holds the value of the target variable. One thing that is apparent in this table is that the structure of the data does not quite fit the supervised learning framework. Two problems arising from this particular data format are the following: First, if the dataset is used in this format during model training, no past information will be included in the samples, due to the fact that each row only includes information about one specific lag. Second, since the target variable of each row points to a current lag value, the model will learn to make estimations only for the current lag (rather than forecasts).

A benefit of using ML models over traditional statistical approaches (e.g., ARIMA) is their ability to support more than one input features. Trying to take advantage of this, we used a method called “sliding window” (Dietterich, 2002) to transform the dataset in a format that integrates into a single sample multiple prior time steps as inputs (X) to predict future time steps as output (Y). In short, this method extends each initial sample of the dataset by including past information and future information simultaneously into a single row. This approach is described in more detail below.

The number of past time steps that we want to include as input into each sample is called the “window width” or size of the lag. As a first step, the width of the sliding window needs to be chosen. Window width, illustrated as a red box in [Table 7](#), corresponds to the number of rows, i.e., the current lag (indicated with a red arrow) plus a number of past lags that will be merged into a new single row. In this example, supposing that *t* is the current lag, the red box in [Table 7](#) indicates that independent variables of the samples at lags *t* and *t*-1 (one step in the past) will be merged into one new row that incorporates not only current but also past information. Additionally, the desired forecasting horizon, illustrated as a blue box in [Table 7](#), needs to be chosen. More specifically, the blue box in this example indicates that we want forecasts for 1 step-ahead, thus the Y value of *t*+1 sample will be selected as the target variable. In case we wanted to prepare the

²⁶ <https://sites.google.com/view/technical-debt-forecasting/main>

²⁷ <https://scikit-learn.org/stable/>

Table 7
Dataset collected in temporal order.

timestamp	X			YI
	X1	X2	X3	
0	10	100	1000	10000
1	20	200	2000	20000
2	30	300	3000	30000
3	40	400	4000	40000
4	50	500	5000	50000
...

Table 8

The reframed dataset after applying the sliding window approach.

Index	X						Y
	X1 (t-1)	X2 (t-1)	X3 (t-1)	X1 (t)	X2 (t)	X3 (t)	
0	10	100	1000	20	200	2000	30 000
1	20	200	2000	30	300	3000	40 000
2	30	300	3000	40	400	4000	50 000
3	40	400	4000

dataset for 2 steps-ahead forecasts, $t+2$ value would be selected as the target variable, and so on. The above process will result in a new row, as depicted in [Table 8](#). The process is repeated by shifting the two boxes simultaneously over the samples, one step at a time, creating new rows until the window reaches the end of the table. Applying the above transformation will result in a reframed dataset that uses one past lag plus the current lag of independent variables to forecast 1 step-ahead. The reframed dataset is presented in [Table 8](#).

There is no standard answer regarding the choice of the window width, i.e., the number of past lags that will be merged per row. This choice usually depends on the number of independent variables, the length of the forecasting horizon and the forecasting model itself. Therefore, during the initial window width selection, a balance needs to be found between the model complexity and the optimal prediction quality. It is often a good idea to test different numbers by training an algorithm and see what values work better for different forecasting horizons, based on an error minimization criterion. For instance, we found out that choosing a window width of 2 lag observations resulted in the minimum Mean Absolute Error (MAE) when trying to forecast for 5 steps ahead, for most of the application-specific datasets across different models. Adding more than 2 lags simply increased models complexity, without profound impact on the model accuracy. Respectively, for longer forecasting horizons, a larger window appeared to be more suitable and resulted in better model performance.

We restructured each application-specific dataset using this method, depending on the forecasting length we wanted to test our models, to make it suitable for supervised ML. Once a time series dataset is prepared this way, any of the standard linear and non-linear ML algorithms can be applied, as long as the order of the rows is preserved. A fragment of the Apache Kafka reframed dataset, after applying the sliding window approach is presented in [Table 9](#).

Another particular challenge that emerges from the concept of TD forecasting is the need to make multi-step forecasts, that is, forecasts for more than one time-step into the future. This need is driven by the fact that we are trying to capture the entire future TD evolution of a software application, rather than the TD value at a particular time step in the future. There are three main approaches that ML methods can use to make multi-step forecasts: (i) the *Direct* approach, where a separate model

is developed to forecast each forecast lead time, (ii) the *Recursive* approach, where a single model is developed to make one-step forecasts, and the model is used recursively where prior forecasts are used as input to forecast the subsequent lead time, and (iii) the *Multiple output* approach, where a single model with multiple outputs is developed, capable of predicting the entire forecast sequence in a one-shot manner.

Most ML-based regression models, with the exception of ANNs, do not directly support more than one outputs. Hence, we excluded *Multiple output* approach. Moreover, the sliding window method described above assumes that the dataset is reformed in a multivariate way, i.e., it includes lag observations from independent variables. As a result, *Recursive* approach is also excluded because it would require also forecasted values of the independent variables to predict further than one step ahead. Therefore, we adopted the *Direct* approach, which means that separate models will be developed to forecast each forecasting horizon. In practice, this means that when trying to forecast for N steps-ahead, the dataset will be reframed N times by following the sliding window approach described above, where each time the dependent variable Y will point to a value from $t+1$ to $t+N$ steps ahead. Subsequently, N separate models will be created, each dedicated to forecast one future point starting from $t+1$ up to $t+N$. Finally, the outputs of the models will be merged into a common vector that depicts the entire forecasted TD evolution up to N steps ahead.

5. Machine learning approach for TD forecasting

In the previous section, we first introduced various software-related metrics that have been widely used in the literature as TD indicators and, then described the data collection process we followed in order to prepare our initial application-specific datasets. Subsequently, during the feature selection process described in [Section 4.3](#), we reduced the initial 24 features (TD indicators) to 4 in order to reduce model complexity. The optimal TD predictors selected were *Code Smells*, *Bugs*, *Duplicated Blocks* and *Afferent Coupling* (*Ca*). Finally, we restructured each application-specific dataset using the sliding window method to make it suitable for supervised ML. In this section, we examine the ability of various ML models to forecast the evolution of TD Principal for each application-specific dataset based on the selected TD predictors. To do so, we train and test the selected models for various forecasting horizons ranging from 1 to 40 steps (weeks) ahead by means of time series validation. Obtained prediction errors of the investigated algorithms are compared among the various forecasting horizons and their benchmarking and evaluation results are documented thoroughly.

5.1. Model training, testing and benchmarking

In this section, we investigate the ability of linear and non-linear ML models to forecast TD evolution of 15 software applications. To do so, we applied a collection of ML models such as

Table 9

The Apache Kafka dataset reframed for 1 step-ahead forecasts using a sliding window with a width of 2.

X	Index	cs (t-2)	bu (t-2)	db (t-2)	ca (t-2)	cs (t-1)	bu (t-1)	db (t-1)	ca (t-1)	cs (t)	bu (t)	db (t)	ca (t)	Y
														TD (t+1)
0	1277.0	89.0	29.0	7.985	1267.0	91.0	31.0	7.961	1102.0	88.0	33.0	7.986	23782.0	
1	1267.0	91.0	31.0	7.961	1102.0	88.0	33.0	7.986	1104.0	89.0	37.0	8.064	23791.0	
2	1102.0	88.0	33.0	7.986	1104.0	89.0	37.0	8.064	1107.0	89.0	37.0	7.995	23852.0	
3

Where cs = Code smells, bu = Bugs, db = Duplicated blocks, ca = Afferent coupling, t = Time step

Multivariate Linear Regression (MLR), Ridge and Lasso regression, Stochastic Gradient Descent (SGD), Support Vector Regression (SVR) with both linear and Gaussian kernel, and Random Forest regression and compared their results for each application-specific dataset. Most of these models have been extensively compared and evaluated in the literature for their ability to predict software quality attributes, such as Maintainability (Riaz et al., 2009; Chug and Malhotra, 2016; Jin and Liu, 2010; Malhotra and Lata, 2018; Elish and Elish, 2009) and Security (Walden et al., 2014; Chowdhury and Zulkernine, 2011; Shin et al., 2011), or lower-level software properties, such as code smells (Fontana et al., 2016a) and defects (Challagulla et al., 2005). However, choosing the most appropriate ML model is often the result of trial and error, as the predictive performance of these algorithms strongly depends on the size and structure of the data. Therefore, within the context of this paper, we considered investigating a broad spectrum of ML models in order to account for highly diverging data relationships that may govern the different application-specific datasets and overcome the limitations of different techniques. The selected models are briefly described below:

- Multivariate linear regression (MLR) is the most commonly used technique for modeling the relationship between two or more independent variables and a dependent variable by fitting a linear equation to observed data. During MLR, the coefficients of the variables are estimated using the least squares method. The main advantages of this technique are its simplicity, interpretability and the fact that it performs well when the relationship to be modeled is not extremely complex. In addition, it is supported by many popular statistical packages. However, quite often simple MLR models are suffering from overfitting.
- Ridge and Lasso regression are simple techniques that aim to reduce model complexity and thus, prevent overfitting by applying regularization, i.e., add some constraints to the loss function. In the case of Ridge regression, those constraints are the sum of squares of the coefficients multiplied by the regularization coefficient (lambda). This regularization type is known as L2. Lasso regression works similarly but instead of adding the squares of the coefficients to the loss function, it adds absolute values. As a result, during the optimization process, coefficients of unimportant features may become zero, which acts as an automated feature selection. This regularization type is known as L1. The main advantages of these regularization techniques, apart from the fact that they prevent overfitting, are the simplicity and computational efficiency of the produced model. However, regularization models are often suffering from high bias error.
- Gradient descent is the process of minimizing a function by following the gradients of the loss function. This involves knowing the form of the loss as well as the derivative so that the function can move towards the minimum value. Stochastic Gradient Descent (SGD) regression is based on gradient descent, but instead of updating coefficients based on the derivative of the data, the algorithm updates the

coefficients based on the derivative of a randomly chosen sample. In that way, SGD allows the function to converge and overcome local minima faster. Because of the randomness involved, the main advantage of SGD is its ability to perform well with noisy data.

- The goal of Support Vector regression (SVR) is to find a function that approximates the target values for all the training data with the minimum generalization error. To achieve this, it tries to learn a non-linear function by linearly mapping features into high-dimensional, kernel-induced feature space. The main advantage of SVR is the efficient non-linear data handling by using the kernel trick. In addition, SVR supports regularization capabilities (L2 Regularization) that prevent overfitting. However, hyper-parameter tuning and choosing an appropriate kernel function can be proven a difficult task.
- Starting with the base case, a Regression Tree (RT) is a variant of decision trees that is built through an iterative process of splitting the data into partitions on each of the decision nodes, known as binary recursive partitioning. An enhanced version of the RT is Random Forest (RF) method. RF is an ensemble of RTs trained with the “bagging” method. Bagging repeatedly selects random samples by replacing the training set and fits trees to these samples. After training, predictions for unseen samples are made by averaging the predictions, or by taking the majority vote of all the individual RTs. The main advantages of RF are its interpretability and the fact that it is great at learning complex, highly non-linear relationships. However, RF models are slower and require more memory compared to the other models presented. In addition, RF models are prone to major overfitting due to the training nature of decision trees.

For the conduction of our experiments, we used the Python programming language and more specifically the *scikit-learn*²⁸ ML library. For reasons of brevity, the ML approaches presented below will focus mainly on the Apache Kafka software system. However, results and model comparisons will include also the rest of the applications.

Once our dataset is ready for supervised learning, the next step is to train and validate the performance of the selected algorithms. Validation methods extensively used in ML, such as k-fold cross-validation, cannot be directly used with time series data due to the temporal order in which values were observed. Hence, observations cannot be randomly split into groups without respecting the temporal order. To better assess prediction accuracy and compare different models we adopted the Walk-forward Train-Test validation method (Stone, 1974), a strategy inspired by k-fold cross-validation. Walk-forward Train-Test validation is a commonly used way to evaluate time series models performance, based on the notion that models are updated when new observations are made available. In brief, during Walk-forward Train-Test validation a subset of n consecutive points extracted from the original time series is used to train an initial model.

²⁸ <https://scikit-learn.org/stable/>

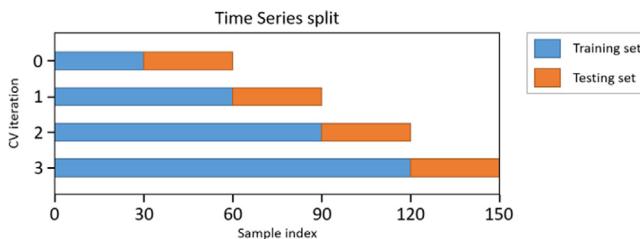


Fig. 4. Walk-forward Train-Test validation.

Then, accuracy of the model is tested against future time steps and prediction is evaluated against the known value to compute prediction errors. Finally, the time window is expanded to include the known values into the training set and the process is repeated. Validation results are combined (e.g., averaged) over the rounds to give an estimate of the model's predictive performance. Using Walk-forward Train-Test validation will result in more models being trained, and in turn, a more accurate estimate of the performance of the models on unseen data. Fig. 4 provides a visualization of the Walk-forward Train-Test validation behavior.

The Apache Kafka dataset consists of 150 observations (snapshots). For Walk-forward Train-Test validation we chose the number of splits = 5, meaning that training set will start from 25 samples and will expand up to 125 samples during the last iteration. The test set will constantly contain 25 observations. The number of splits = 5 was chosen such that each train/test group of data samples is large enough to be statistically representative of the broader dataset. A larger number of splits would result in overly small train/test groups, which in turn would suffer from large variability (Kohavi et al., 1995). It is worth mentioning here that number of splits = 5 was chosen also for the other application datasets that contain 150 observations. For those that contain 100 observations, we chose the number of splits = 4 to maintain the train/test group size analogy. To test predictive performance of our models for different future horizons, we repeated the whole validation process five times, where predictions were made for the next n+1 (1 week), n+5 (5 weeks), n+10 (10 weeks), n+20 (20 weeks), and n+40 (40 weeks) future steps respectively.

Before the learning process begins, a hyper-parameter tuning process must take place in order to increase models' predictive performance. A model hyper-parameter is an external attribute of the model. In contrast to typical model parameters, e.g., the coefficients of a Linear Regression model, the value of a hyper-parameter cannot be estimated from data during the training process. Hyper-parameter examples may include the penalty parameter C of the error term in SVM, the number of trees in the Random Forest, etc. In order to tune our models in the best possible way, we used the GridSearchCV²⁹, a python implementation of the Grid-search method (Feurer et al., 2015). Grid-search is commonly used to find the optimal hyper-parameters of a model that result in the most accurate predictions, by performing an exhaustive search over specified parameter values for an estimator. We chose R^2 (coefficient of determination) as the objective function of the estimator to evaluate a parameter setting. R^2 is the proportion of the variance in the dependent variable that is predictable from the independent variable(s). We performed hyper-parameter selection on every application-specific dataset during the 5-fold Walk-forward Train-Test validation described above to avoid overfitting and ensure that the selected models have a good degree of generalization.

We evaluated and compared the forecasting performance of the investigated models using the Mean Absolute Percentage Error (MAPE). The MAPE is a popular measure for forecast accuracy that uses absolute values to measure the size of the error in percentage terms. MAPE has two advantages. First, the absolute values keep the positive and negative errors from canceling out each other. Second, because relative errors do not depend on the scale of the dependent variable, this measure allows for comparing forecast accuracy between differently scaled time-series data (e.g., different software applications). The equation of MAPE is given below:

$$MAPE = \frac{100}{n} \sum_{i=1}^n \frac{|y_i - \hat{y}_i|}{y_i} \quad (2)$$

where n is the number of observations, y_i is the actual value and \hat{y}_i is the forecast value.

To further complement model evaluation, we also computed the Mean Absolute Error (MAE) as well as the Root Mean Squared Error (RMSE). Both of these errors are widely used in forecasting tasks. MAE measures the average magnitude of the errors in a set of predictions, without considering their direction. RMSE is a quadratic scoring rule that also measures the average magnitude of the error. Both MAE and RMSE express average model prediction error in units of the variable of interest. The equations of MAE and RMSE are given below:

$$MAE = \frac{\sum_{i=1}^n |y_i - \hat{y}_i|}{n} \quad (3)$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (4)$$

Again, n is the number of observations, y_i is the actual value and \hat{y}_i is the forecast value.

In Table 10, we report a comparison of prediction errors of the regression models trained on the Apache Kafka dataset for multiple (1, 5, 10, 20 and 40) time steps (weeks) into the future. Prediction errors in each cell of the table are averaged values of the testing errors for all train-test splits that were performed during Walk-forward Train-Test validation. Prediction errors indicated in bold are averaged values of the specific models that were created for each week-ahead prediction category (i.e., 1-week, 5-weeks, 10-weeks ahead models, etc.).

As a reminder, since we adopted the Direct approach described in Section 4.4, each examined model provides a single output, that is, the predicted value of the horizon that it was trained to provide forecasts for. In practice, this means that the values of the errors presented below refer to a forecast for a specific individual point in the future, not the entire forecasted evolution of up to that point. We present aggregated forecasts that illustrate the entire evolution of the examined applications in the Model Execution phase described in Section 5.2.

Fig. 5 illustrates the MAPE of the forecasting models, averaging the five forecasting horizon cases (i.e., 1, 5, 10, 20 and 40 weeks ahead) under investigation. Fig. 6, Fig. 7, Fig. 8, Fig. 9 and Fig. 10 illustrate the MAPE of the forecasting models for 1, 5, 10, 20 and 40 steps (weeks) ahead respectively.

By observing Fig. 5, it is clearly depicted that linear models, such as MLR, Lasso, Ridge and SVR(linear) Regression, have generally lower MAPE values and perform better than non-linear models, such as SVR(rbf) and Random Forest Regression. Moreover, we observe that among linear models, the best accuracy is demonstrated by models that apply Regularization in order to prevent overfitting, i.e., Lasso and Ridge Regression. More specifically, Lasso Regression is the best candidate with an average

²⁹ https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html

Table 10
Apache Kafka TD predictions using Walk-forward Train-Test validation.

Model	Weeks ahead	MAE (min)	RMSE (min)	MAPE (%)
MLR	1	594.799	819.085	1.267
	5	2655.770	3094.180	5.906
	10	3105.604	3613.806	6.595
	20	4163.185	4948.930	8.146
	40	6214.496	6790.829	11.072
	Average	3346.771	3853.366	6.597
Lasso regressor	1	430.106	676.957	0.881
	5	1474.881	1768.435	3.055
	10	2240.932	2615.345	4.455
	20	2894.261	3239.653	5.444
	40	4700.503	5004.199	8.655
	Average	2348.137	2660.918	4.498
Ridge regressor	1	438.024	682.877	0.898
	5	1579.645	1869.698	3.260
	10	2579.991	2922.568	4.979
	20	3239.642	3627.422	5.977
	40	5061.284	5428.619	9.177
	Average	2579.717	2906.237	4.858
SGD regressor	1	789.139	1049.557	1.690
	5	2812.558	3253.977	6.225
	10	3069.528	3606.347	6.419
	20	3940.336	4718.533	7.617
	40	6383.504	7063.643	11.311
	Average	3428.834	3970.349	6.706
SVR regressor (linear)	1	571.416	830.770	1.214
	5	3457.010	4004.375	7.567
	10	2811.954	3344.931	6.153
	20	3760.544	4551.495	7.564
	40	6569.625	7304.984	11.89
	Average	3428.834	3970.349	6.706
SVR regressor (rbf)	1	4758.278	5254.790	9.561
	5	5753.172	6257.845	11.819
	10	4657.056	5348.217	9.798
	20	3687.333	4485.870	7.591
	40	5240.294	5794.354	9.835
	Average	4819.227	5428.215	9.721
Random forest regressor	1	4273.391	4731.381	8.824
	5	5454.023	5936.521	11.259
	10	4430.855	4994.763	9.252
	20	3425.365	4173.618	6.798
	40	3664.173	3947.600	7.007
	Average	4262.788	4767.657	8.660

Apache Kafka - Mean Absolute Percentage
Error - Average of 1-40 weeks ahead

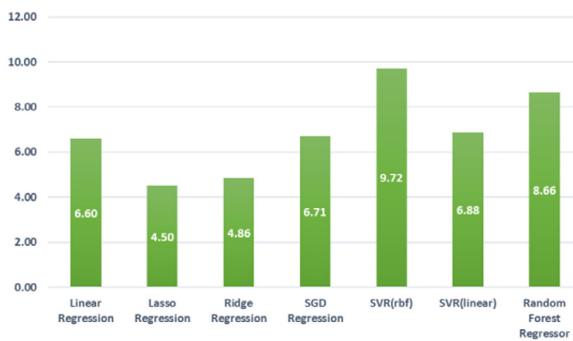


Fig. 5. Apache Kafka TD predictions — MAPE averaged for all steps-ahead using Walk-forward Train-Test validation.

Apache Kafka - Mean Absolute Percentage
Error - 1 week ahead

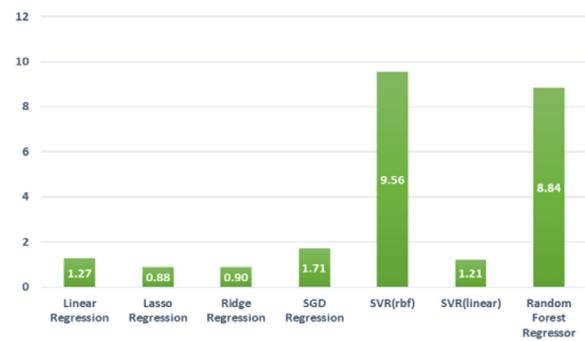


Fig. 6. Apache Kafka TD predictions — MAPE for 1 step-ahead using Walk-forward Train-Test validation.

MAPE value of 4.5%, followed by Ridge Regression with an average MAPE value of 4.86%.

When it comes to shorter forecasting length (i.e., 1–10 weeks ahead), the difference between linear and non-linear model performance becomes even clearer. By having a look at Fig. 6, we notice that forecasting the TD of the Apache Kafka project for

1 step ahead (1 week) using Lasso Regression gives a MAPE of 0.88%, while for the same horizon SVR with a Gaussian kernel gives 9.56% and Random Forest Regression gives 8.84%. Correspondingly, by having a look at Fig. 7 and Fig. 8 we observe that Lasso Regression for 5 steps (5 weeks) and 10 steps (10 weeks) ahead gives a MAPE of 3.06% and 4.46%, while for the

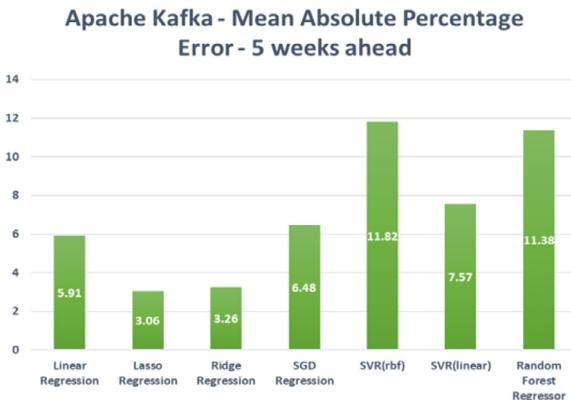


Fig. 7. Apache Kafka TD predictions – MAPE for 5 steps-ahead using Walk-forward Train-Test validation.

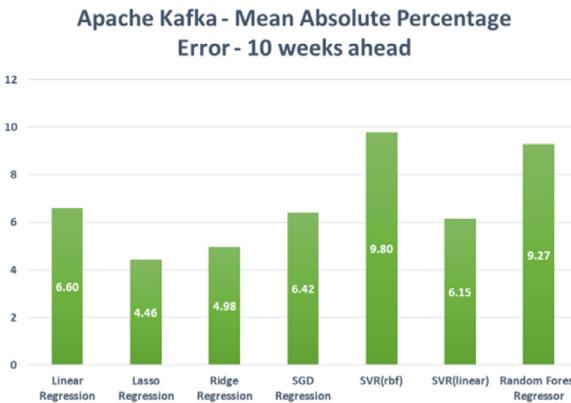


Fig. 8. Apache Kafka TD predictions – MAPE for 10 steps-ahead using Walk-forward Train-Test validation.

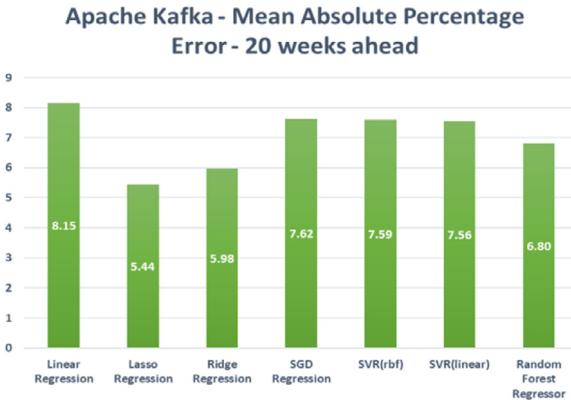


Fig. 9. Apache Kafka TD predictions – MAPE for 20 steps-ahead using Walk-forward Train-Test validation.

same horizons SVR with a Gaussian kernel gives 11.82% and 9.80% respectively.

Linear models are the best candidates even for a forecasting length of 20 steps (20 weeks) ahead, as can be seen by Fig. 9. However, an interesting observation is that while their predictive power drops significantly as we try to forecast longer into the future, non-linear models seem to have an almost stable performance over the holdout sample for all steps ahead. This could be an indicator that for even longer lengths, non-linear models could perform better than the linear ones. Indeed, by having a look in Fig. 10, we observe that for a forecasting horizon of 40 steps (40

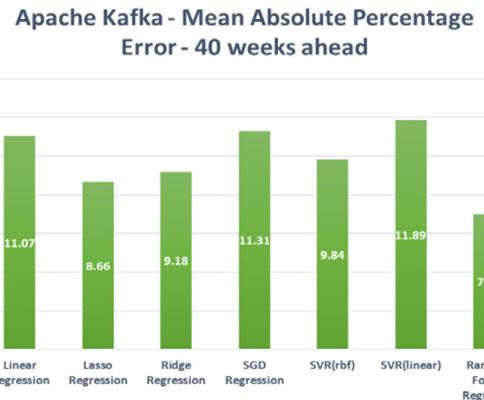


Fig. 10. Apache Kafka TD predictions – MAPE for 40 steps-ahead using Walk-forward Train-Test validation.

15 Projects - Mean Absolute Percentage Error - Average of 1-40 weeks ahead

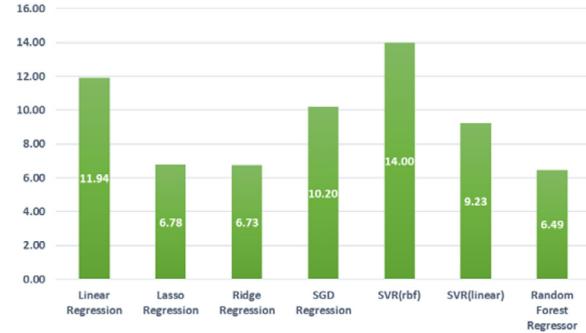


Fig. 11. 15 projects TD predictions – MAPE averaged for all steps-ahead using Walk-forward Train-Test validation.

weeks) ahead, Random Forest Regression is the best candidate and gives the lowest MAPE (7.01%).

To further examine the ability of the investigated algorithms to forecast TD Principal and get an understanding of how the models perform, we repeated the same experiments for each of the 15 applications in our dataset. We will not go through each project one by one, but instead we will provide averaged scores. Detailed results of applying Walk-forward Train-Test validation for the rest of the applications can be found at the Appendix section (Tables 14–27). Fig. 11 below illustrates the MAPE of the forecasting models, averaging the five forecasting horizon cases (1, 5, 10, 20 and 40 weeks ahead) and the 15 software applications under investigation. Fig. 12, Fig. 13, Fig. 14, Fig. 15 and Fig. 16 illustrate the MAPE of the forecasting models for 1, 5, 10, 20 and 40 steps (weeks) ahead respectively, averaging the 15 software applications under investigation.

Similarly to the Apache Kafka case, we observe that for shorter forecasting lengths, linear models that apply Regularization, such as Lasso and Ridge regression, have generally lower MAPE values and higher performance compared to the non-linear models. As depicted in Figs. 12 and 13, Ridge and Lasso Regression models are the best candidates with MAPE values of 1.44%–3.91% and 1.39%–4.11% respectively. We also observe that again, the predictive power of linear models drops significantly as we forecast longer into the future. However, the non-linear Random Forest Regression algorithm seems to have an almost stable performance over the holdout sample for all steps (weeks) ahead. In fact, starting from 20 (Fig. 15) up to 40 steps ahead (Fig. 16), we observe that Random Forest Regression is the best candidate and

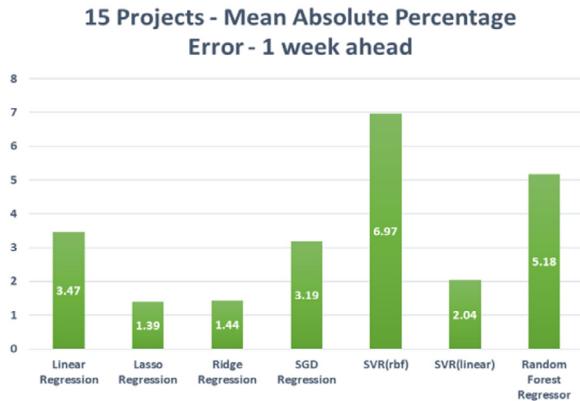


Fig. 12. 15 projects TD predictions – MAPE for 1 step-ahead using Walk-forward Train–Test validation.

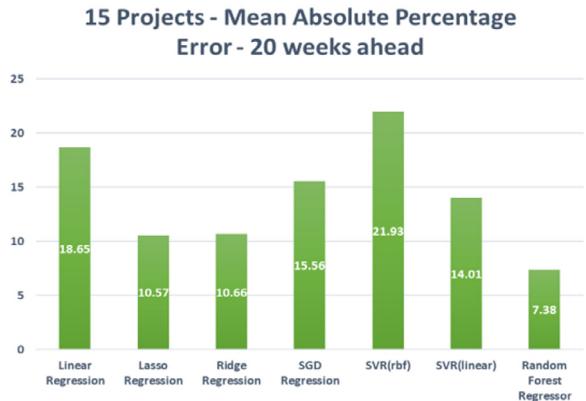


Fig. 15. 15 projects TD predictions – MAPE for 20 steps-ahead using Walk-forward Train–Test validation.

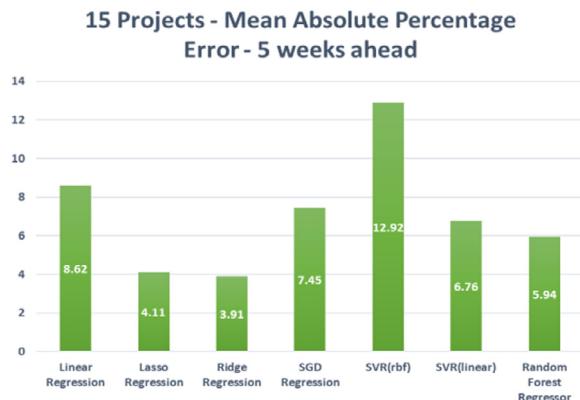


Fig. 13. 15 projects TD predictions – MAPE for 5 steps-ahead using Walk-forward Train–Test validation.

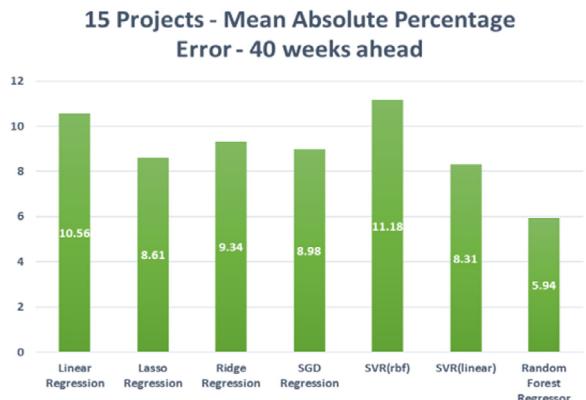


Fig. 16. 15 projects TD predictions – MAPE for 40 steps-ahead using Walk-forward Train–Test validation.

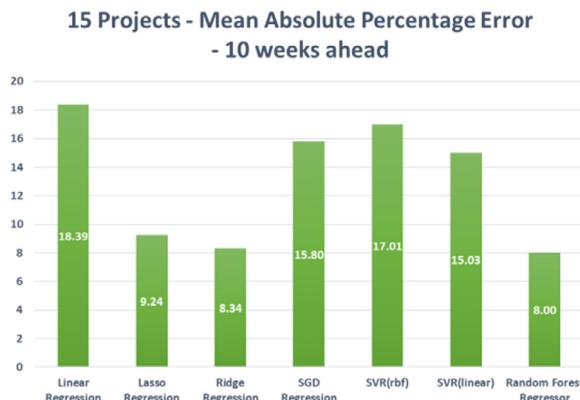


Fig. 14. 15 projects TD predictions – MAPE for 10 steps-ahead using Walk-forward Train–Test validation.

performs better than the other models giving the lowest MAPE value of 7.38% and 5.94% respectively.

To sum up, an interesting finding that we can extract from the analysis of the experiments is the observation that linear models that apply Regularization, i.e., Lasso and Ridge Regression, are capable of achieving high forecasting performance for shorter forecasting lengths (<10 weeks ahead), while the non-linear Random Forest Regression is performing better than the rest of the investigated models for longer forecasting lengths (>10 weeks ahead). The fact that the above results are observed in all of the 15 application-specific datasets is also interesting

and of high significance, whereas it increases our confidence regarding the generalizability of the aforementioned findings. Although we cannot be sure that these results may apply to similar applications, they do make a valuable contribution to the beginning of the TD forecasting landscape composition. Therefore, a TD forecasting tool could leverage the predictive power from both of these algorithms combined to deliver good predictions and adequately forecast future TD Principal trends of software applications.

5.2. Model execution

Following the construction and benchmarking of our models described in Section 5.1, this section presents indicative examples of model execution as well as indicative visualizations of the forecasting results. As reported in Section 4.4, we decided to adopt the Direct approach, meaning that separate models were developed to forecast each forecast lead time. In Fig. 17 below, we provide an example of forecasting the TD Principal evolution of Apache Kafka application for 20 steps (weeks) ahead using Random Forest regression, which during the model validation phase was reported to have a stable performance and perform better than the other examine models for longer forecasting lengths. The red line denotes the forecast, while the blue line denotes the ground truth. It should be noted that the samples covered by the red line (i.e., test set) were excluded during the model-training phase. Behind the scenes, according to the adopted Direct approach, 20 models were executed, one for each specific length of interest (starting from 1 step to 20 steps), while their forecasted TD values

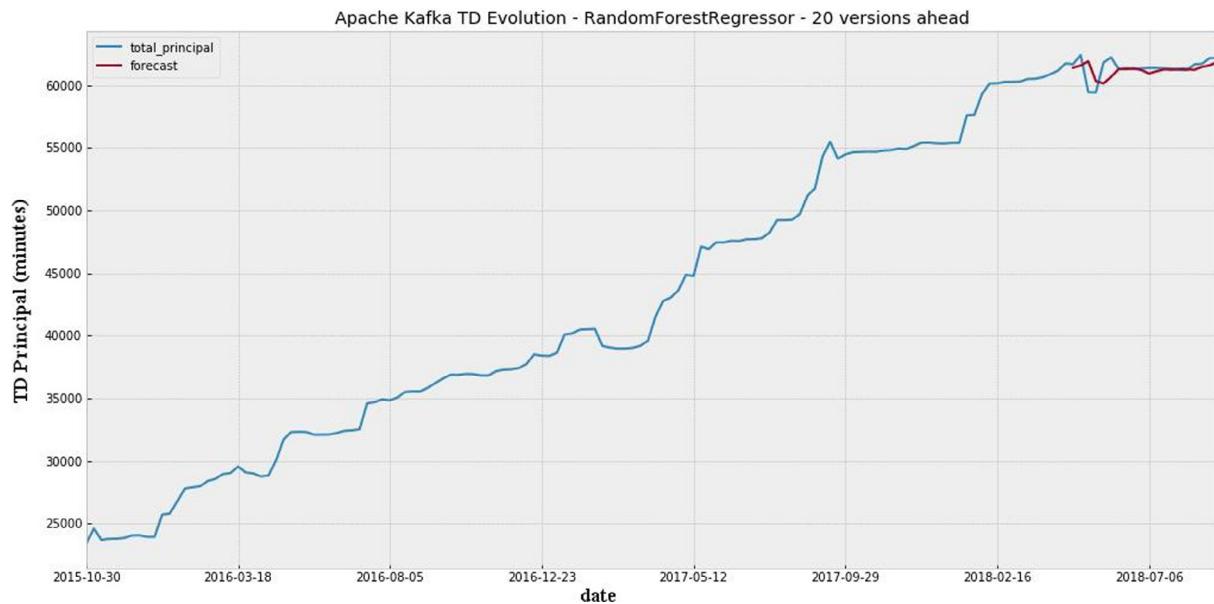


Fig. 17. Apache Kafka TD Principal forecasting for 20 steps ahead using Random Forest and the Direct approach.

where aggregated into a common vector, and then plotted as the projected TD evolution.

Indicative visualizations illustrating the forecasting results of Random Forest regression for 20 steps ahead for the rest of the 14 applications are provided in the online³⁰ supportive material. As can be seen in both Fig. 17 and the online material, similar observations can be made for all 15 applications under investigation. In particular, the Random Forest regression seems to provide meaningful long-term forecasts for each one of the studied cases (i.e., software applications). In fact, the selected algorithm is able to capture the trend of the future evolution of the TD Principal, whereas in most of the cases the future value of the TD Principal is also captured with a sufficient level of accuracy. For reasons of brevity, we do not provide illustrations of TD evolution forecasts for the rest of the algorithms. However, similarly to the case of long-term TD Principal forecasting using Random Forest regression, satisfying forecasts were also obtained by using Regularization models (i.e., Lasso and Ridge regression) for shorter forecasting lengths, as expected by the results that were reported during model benchmarking in Section 5.1.

5.3. Technical implementation

The work presented in this paper introduces an approach aiming to cover the existing gap in the field and set the foundations towards methods and accompanying tools able to deliver TD forecasts and therefore assist developers and project managers in taking proactive actions regarding TD management activities. The SDK4ED³¹ European project aims to address this challenging issue by implementing the proposed approach in the form of a tool, i.e., the TD Forecasting tool, as a part of the integrated TD Management (TDM) framework. To this end, an envisaged TD Forecasting tool has been implemented as individual standalone tool in order to facilitate its applicability in practice. This tool consists of a backend server dedicated to the deployment of a set of forecasting models, a web service that exposes the server, and an interactive Graphic User Interface (GUI) that allows the invocation of forecasting models and displays the results, providing users with insightful information for the future evolution

of TD. Both the backend and frontend of the TD Forecasting tool are components of the overall SDK4ED Dashboard, which forms the final outcome of the SDK4ED project. The TD Forecasting tool, integrated into a preliminary version of the SDK4ED Dashboard, can be found online³² (currently being used for development purposes). The main screen of the tool is provided in Fig. 18.

The main screen of the TD Forecasting tool comprises a dropdown button, two interactive plots and one table. The dropdown button allows the user to select the forecasting horizon for which they would like to see predictions for. Once the forecasting horizon is selected, the backend server invokes the proper forecasting algorithm (depending on the selected horizon) and returns the predictions back to the GUI, which in turn parses the result. Then, the interactive plots showing the ground truth (green) and the predicted (red) TD Principal evolution appear on the screen. The first plot shows the entire evolution followed by the forecasted evolution of the application, whereas the second plot focuses solely on the forecast, giving a more fine-grained view. In addition to the plots, a complementary table comprising the detailed results of the forecasts is presented at the bottom-right part of the screen. This table presents the forecasted TD values for the upcoming weeks, as well as the difference between the current TD value and the forecasted TD values per week, which may serve as an indicator of whether the TD Principal will increase or decrease, and to what extent. This additional information is expected to help the developers take even more informed decisions regarding the prioritization of their TD repayment activities.

6. Case study

In this section, we present the results of an industrial study conducted to empirically evaluate the meaningfulness of the TD forecasting approach introduced in this work and to investigate the extent to which this approach can provide valuable insights and affect developers' decisions regarding the evolution of software, via a questionnaire distributed to representatives of a software company.

³⁰ <https://sites.google.com/view/technical-debt-forecasting/main>

³¹ <https://sdk4ed.eu/>

³² <http://160.40.52.130:3000/tdforecast>

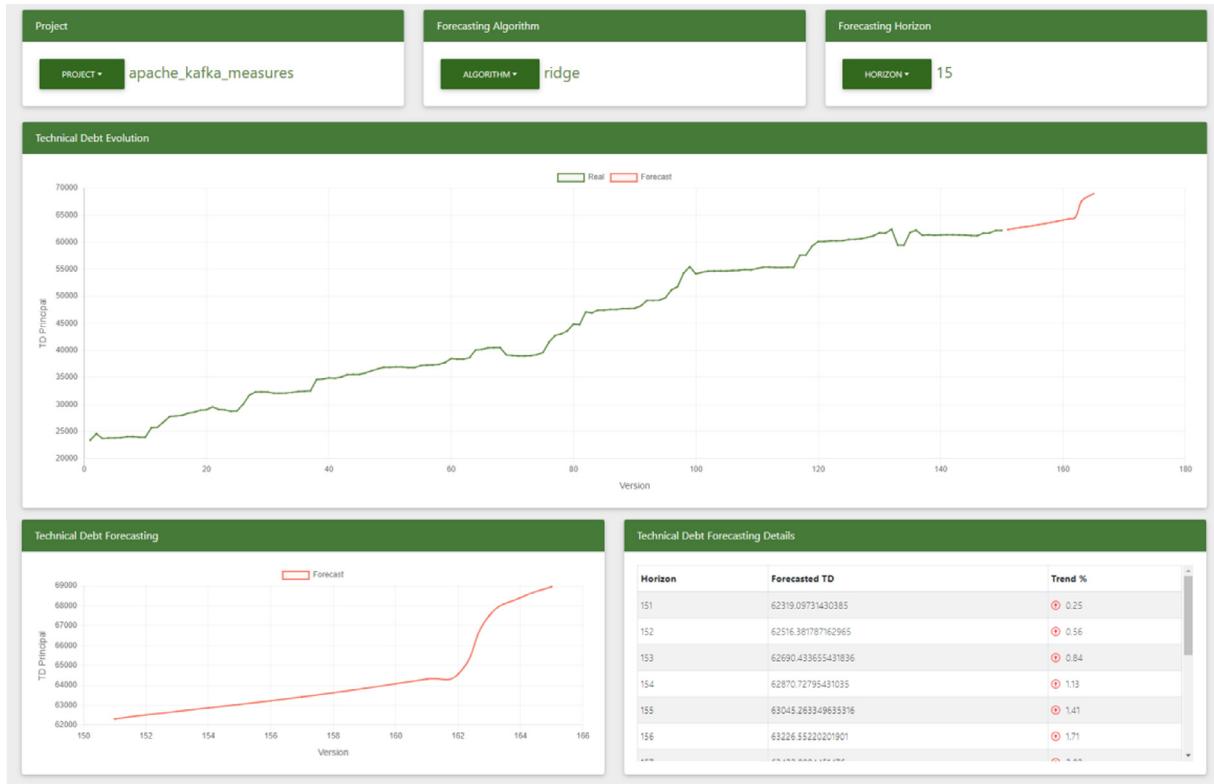


Fig. 18. Main screen of the TD Forecasting tool.

6.1. Survey design

To minimize the possibility that current work will remain a statistical exercise detached from real software development practices, we have designed a survey through which we seek feedback from practitioners. The goal of this survey is twofold: (a) to empirically evaluate the meaningfulness and accuracy of the TD forecasting methodology proposed in this study, and (b) to empirically assess the usefulness and acceptance of the TD Forecasting concept in general, especially for software companies that deal with daily TD management activities.

For the purposes of this survey, we have involved a Greek department (located in Thessaloniki, Greece) of a large European software company (hereafter referred to as *Company*) that provides IT development services in multiple technologies to private and public organizations. The *Company* employs more than 2,200 highly-skilled professionals worldwide. However, the *Company* wants to keep its anonymity, thus all records in the dataset have been anonymized, and no personalized information can be provided, either about the company and its projects, or the case study participants. In the case of this survey, although participants might not be extremely familiar with the TD concepts and terminology, they are all experienced in issues related to quality assessment, since the *Company* uses SonarQube for continuous inspection of code quality during its software development process. A possible lack of experience in TD terminology has been considered during the design of the data collection instrument. Furthermore, the fact that the *Company* uses SonarQube as a quality inspection tool, allows us to partially validate also the SonarQube TD measurement mechanism, in a sense that we can determine, through communication with the developers, if the TD measurements are in line with real events occurring during the software development process.

As a survey instrument, we opted for a questionnaire, which is described in detail below and can be found also online.³³ The most important part of developing a questionnaire is the selection of questions. In our survey, this process was governed by the guidelines provided by Kitchenham and Pfleeger (2002): (a) keep the amount of questions low, (b) questions should be purposeful and concrete, (c) answer categories should be mutually exclusive, and (d) the number, the order and the wording of questions should avoid biasing the respondent. To this end, we constructed a questionnaire with 13 main questions (4 multiple-choice and 9 short-answer), organized into three main parts (see Table 11), and an introductory part (2 questions). The questionnaire begins with the introductory part (i.e., Part-1) where participants are asked to provide some demographic information, such as their role and years of experience in the *Company*. Subsequently, in Part-2, participants are first introduced to some background information on the concept of TD and its main components (such as TD Principal, inefficiency types, etc.) and then asked to rate, on a Likert scale, a group of questions that aim to evaluate the usefulness of TD Forecasting. The last two parts of the questionnaire, i.e., Part-3 and Part-4, refer to some project-specific questions (to be able to provide valid answers in Part-3 and Part-4, the participants should have been actively involved in the development of these projects). The process of designing and formulating the questions of Part-3 and Part-4 are thoroughly described below.

During our visit to the premises of the *Company*, we were given access to a dedicated SonarQube instance hosting the analysis results of a large set of software applications, in order to find the most suitable candidates for the design of our survey. The criteria we relied on to select an application are as follows. First, the application needs to be developed in Java programming language. Second, it needs to be constantly maintained and thus,

³³ <https://forms.gle/Jjg8RoA55m1EwMJ77>

Table 11
Survey instrument.

ID	Question
Part 1 – Demographics	
Q1.1	What is your role in the company?
Q1.2	How many years of experience do you have in this position?
Part 2 – The usefulness of TD Forecasting	
Q2.1	How useful is it to have an estimation of the current TD Principal of a software project?
Q2.2	How useful is it to have a forecast of the future TD Principal of a software project?
Q2.3	To what extent would a forecast of the TD Principal make you consider changing the planned future development of a project?
Q2.4	Supposed that a forecast shows an increasing trend of the TD Principal, what actions would you take to repay TD?
Part 3 – TD Principal Evolution	
Project A	
Q3.1	Case 1: A temporal TD increase (6/1/19–9/1/19). What is the cause of this change?
Q3.2	Case 2: A sharp TD increase (18/1/19–27/2/19). What is the cause of this change?
Q3.3	Case 3: A sharp TD increase (4/4/19–12/5/19). What is the cause of this change?
Q3.4	Case 4: A gradual TD increase (27/10/19–14/2/20). What is the cause of this change?
Project B	
Q3.5	Case 1: A gradual TD increase (23/2/19–28/7/19). What is the cause of this change?
Q3.6	Case 2: A sharp TD decrease (28/7/19–18/8/19). What is the cause of this change?
Q3.7	Case 3: A sharp TD increase (21/11/19–19/2/20). What is the cause of this change?
Part 4 – TD Principal Forecasting	
Project A	
Q4.1	The latest commits of Project A show relatively stable TD evolution. However, the forecast for 10 commits ahead shows a gradual increase in the TD principal. By having a look at this forecast, would you change anything in the planned development process? Would you consider performing code refactoring in order to prevent this increase?
Project B	
Q4.2	The latest commits of Project B show a gradual increase in the TD evolution. However, the forecast for 10 commits ahead shows a slight decrease in TD principal increasing rate. By having a look at this forecast, would you change anything in the planned development process? Would you consider investing in enhancing already existing, or adding new functionalities?

to provide a relative long history of commits, as well as the associated SonarQube analysis measurements available for these commits. Finally, its SonarQube analysis measurements need to contain (at the minimum) *code smells*, *bugs*, *duplicated blocks*, *lines of code* and the TD Principal itself. The first three metrics are required for the construction of the TD forecasting models, since they were selected as the most statistically significant TD predictors based on the analysis conducted within Section 4. Unfortunately, the fourth TD predictor, i.e., *afferent coupling* (*Ca*) was not available for this analysis, since the Company does not use the CKJM Extended tool. Furthermore, we were not given access to the source code of the applications, in order to execute the CKJM Extended tool ourselves. However, we believe that the unavailability of one independent variable will not significantly affect the forecasting performance of our models, especially since *afferent coupling* (*Ca*) was found to have the highest standard error among the final four variables, as presented in Table 6.

Based on the above criteria, we ended up with two software applications belonging to the Business Software domain, namely *Project A* and *Project B*, with a size of 58K LoC and 384K LoC respectively. Specifically, analysis data from the SonarQube database of the Company included 62 commits for *Project A* and 51 commits for *Project B*. At this point, it should be noted that although committing code updates at a weekly basis (i.e., at the end of each week) is considered an integral part of the Company's routine, we were informed that there existed a few cases where SonarQube did not run for a particular week, mainly due to technical issues. As a result, collecting SonarQube analysis data at fixed weekly intervals (as introduced in Section 4.2 to establish the dataset) was no longer a viable approach. To overcome this issue, we decided to replace the concept of weekly snapshots with

that of consecutive commits. Therefore, within the context of this survey and more specifically in Part-4 of the questionnaire where we present TD forecasts for each application, forecasting for 10 steps ahead is referring to 10 commits rather than 10 weeks ahead. The SonarQube measurements of the two anonymized software applications selected during this step of the process can be found online.³⁴

As a first step, we parsed the collected SonarQube analysis data, performed some required pre-processing steps and extracted two plots (one for each application) illustrating the entire TD Principal evolution of the two applications under investigation. Subsequently, we manually inspected the plots and extracted selected periods where we identified abrupt (but interesting) TD Principal trends, i.e., a trend showing a gradual increase|decrease, sharp increase|decrease, temporary increase|decrease, etc. Therefore, in Part-3 of the questionnaire, for each of the identified cases, participants are asked what the root cause of these abrupt trends was. For instance, if a trend shows a sharp TD Principal decrease during some period, maybe that is due to a code deletion or code refactoring that removed TD-ill code. If the participants are aware of specific actions they performed that justify these abrupt trends (e.g., refactoring, code additions or deletions, deadlines, etc.), then it means that SonarQube TD measurement mechanism can capture these changes and is in line with real events occurring during the software development process.

As a second step, we applied forecasting models to predict the future TD Principal evolution of both *Project A* and *Project B* for 10 steps ahead. To do so, we exploited a prototype of the

³⁴ <https://sites.google.com/view/technical-debt-forecasting/main>

TD Forecasting tool described in Section 5.3 in order to extract forecasting plots in an automated way, that is, without repeating the tedious process of model training, testing and benchmarking through the usage of Python scripts. Therefore, in Part-4 of the questionnaire, a TD forecast for each application is presented to the participants and they are asked if they would be willing to change anything in the planned development process based only on the projected TD Principal evolution. For instance, if the TD evolution of a specific application has been constant up to a point but forecasts show a sharp TD increase in the future, we ask them if they would consider performing refactoring to prevent that increase. Respectively, if the TD evolution has been gradually increasing up to a point but forecasts show a sharp TD decrease, we ask them if they would consider investing in enhancing new functionalities instead of performing refactoring. In that way, the practical usefulness and meaningfulness of TD forecasts is, at least partially, evaluated using qualitative feedback from the participants.

The majority of the questions in Part 2 of the questionnaire have been answered on a Likert Scale ranging from 1 to 5, with the exception of the last question (Q2.4) which gives the respondents the following options: (a) "Refactoring", (b) "Writing new code that is TD-free", (c) "No actions", and (d) "Other: ____". However, the last two parts of the questionnaire, i.e., Part-3 and Part-4, refer to some project-specific questions and require a short description, so they have been answered by providing a "short answer" text box.

6.2. Survey analysis and results

In this section, we present the results of the survey study, through presenting some demographics and subsequently analyzing and discussing the answers of the participants. In total, we obtained four (4) complete answers. The reason that we received only four complete answers is due to the fact that Part-3 and Part-4 of the questionnaire require deep knowledge of the applications under analysis. In fact, during our communication with the company, we specifically requested that participants should be actively involved in the development of these projects so that they can provide valid answers. To facilitate the process of distinguishing between the responses of the four participants but at the same time maintain their anonymity, the participants of this survey are hereafter referred to as P1 to P4. Regarding the demographics of the participants extracted from Part-1 of the questionnaire, three out of four participants (P2, P3, and P4) are working in the *Company* as Software Developers, while one participant (P1) is working as a Software Architect. Moreover, participants' years of experience in this position range from 2 to 12 years, with a mean value of 6.25 years.

First, we analyzed the answers from Part-2 to understand the usefulness of TD forecasting in an industrial context. More specifically, regarding Q2.1 "How useful is it to have an estimation of the current TD Principal of a software project?", on a Likert scale ranging from: 1 – "Not Useful" to 5 – "Very Useful", two out of four of the participants (P1 and P4) chose "Useful" (option 4), while the remaining two (P2 and P3) chose "Very Useful" (option 5). These responses suggest that TD Principal monitoring is perceived as highly important for software development companies, as it can provide valuable information regarding the effort and, in turn, the cost that is required for maintaining and extending a software application.

Regarding Q2.2 "How useful is it to have a forecast of the future TD Principal of a software project?", on a Likert scale ranging again from: 1 – "Not Useful" to 5 – "Very Useful", two out of four participants (P1 and P2) chose "Useful" (option 4), while the remaining two (P3 and P4) chose "Very Useful" (option 5).

These responses suggest that TD Principal forecasting is of great significance and value for software development companies, since they would be able to gain a better understanding of future TD issues and plan well in advance appropriate refactoring activities for saving maintenance costs.

Regarding Q2.3 "To what extent would a forecast of the TD Principal make you consider changing the planned future development of a project?", on a Likert scale ranging from: 1 – "Not at all" to 5 – "To a great extent", three participants (P1, P2, and P4) chose "To a moderate extent" (option 4), while the remaining one (P3) chose "To a great extent" (option 5). These responses suggest that all participants would consider changing the planned future development of a project based on a forecast of the TD Principal. In fact, this statement is further evaluated through specific questions presented to the participants in Part-4 of this questionnaire.

Finally, regarding Q2.4 "Supposed that a forecast shows an increasing trend of the TD Principal, what actions would you take to repay TD?", three out of four participants (P1, P2, and P4) responded with "Refactoring", while the remaining one (P3) responded with "*First make sure that the new code will have less TD and then, when time plan allows it, refactor existing code*". While code refactoring is a well-established approach for TD repayment (Fowler, 1999), clean code has recently emerged as a promising TD prevention strategy. By inspecting respondents' answers, we notice that while P1, P2, and P4 would opt for refactoring the already existing code to repay TD, the latter answer indicates that P3, possibly forced by strict deadlines that require the delivery of new functionalities, would prefer to increase the overall quality of the project by writing new TD-free code, i.e., clean code that contributes positively to the overall TD. This question is also further assessed through specific questions presented to the participants in Part-4 of this questionnaire.

Subsequently, we analyzed the answers from Part-3 of the questionnaire. In Part-3, for a series of identified cases where the TD Principal of Project A and Project B showed abrupt trends, participants were asked what was the root cause of these changes. The main goal of this part is to assess whether participants are aware of specific actions the development team had performed that justify these abrupt trends (e.g., refactoring, code additions or deletions, deadlines, etc.), and therefore to validate that SonarQube TD measurement mechanism can capture these changes and is in line with real events occurring during the software development process. Participants' answers are summarized in Tables 12 and 13, which refer to comments regarding observed TD Principal trends of Project A and Project B respectively. By inspecting the tables, it can be seen that three participants (P2, P3, and P4) are working on Project A, while two participants (P1 and P3) are working on Project B (with P3 working on both projects). The figures of the identified cases illustrated in the tables can be also found online.³⁵

By reading participants' comments regarding the identified TD Principal trends in Table 12, we can note that in almost all cases, developers involved in Project A are aware of specific actions they performed that justify these abrupt trends. More specifically, according to the participants, the temporal TD Principal increase of Project A depicted in Q3.1 can be attributed to a large amount of functionality that was added from less experienced engineers in order to meet deadlines, followed by improvements to the code that led TD Principal to decrease again. Similarly, the sharp TD Principal increase depicted in Q3.2 can be attributed to rapid code additions that had to be completed in limited time and without proper testing, in order to deliver the 1st version of the project to a customer. Subsequently, the sharp TD Principal

³⁵ <https://sites.google.com/view/technical-debt-forecasting/main>

Table 12
Participants comments on TD Principal trends of Project A.

Cases	TD Principal Trend	Participants Comments (what is the cause of these TD changes?)
Q3.1	<p>Temporal increase (6/1/19 – 9/1/19)</p>	<p>P3: "The project was new and a large amount of functionality was added for the first time from less experienced engineers. Then improvements were made to the initial code and this is why there is this decrease after 08/01."</p> <p>P2: "Added 5 new forms to the project along with all the front end and back end code, entities, domains, repositories, services. (06/01-08/01). The decrease from (08/01-09/01) is due to work made on fixing sonar issues regarding TD."</p> <p>P4: "Rapid code development (code additions) in order to meet deadlines without testing by junior engineers"</p>
Q3.2	<p>Sharp increase (18/1/19 – 27/2/19)</p>	<p>P3: "1st iteration to the customer was coming up. Features had to be completed in limited time."</p> <p>P2: "Added 2 new main forms to the project along with all the front end and back end code, entities, domains, repositories, services with code duplication"</p> <p>P4: Same as Q3.1</p>
Q3.3	<p>Sharp increase (4/4/19 – 12/5/19)</p>	<p>P3: "New functionality added after the comments of iteration"</p> <p>P2: "New forms, new module, with code duplication"</p> <p>P4: Same as Q3.1</p>
Q3.4	<p>Gradual increase (27/10/19 – 14/2/20)</p>	<p>P3: "Beta testing of the module was creating requests for bug fixing and some new functionality."</p> <p>P2: "Code expansion. New features with code duplication"</p> <p>P4: "Bug fixing, code additions and deletions and new features implementations all at the same time without testing if previous implementations are broken and without testing if the new ones are stable"</p>

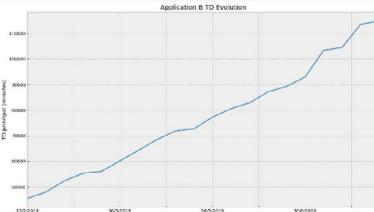
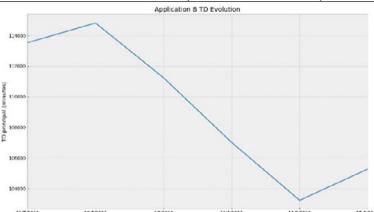
increase depicted in Q3.3 can be attributed to some new additional functionalities that were requested by the customer and thus were quickly added after the 1st delivery of the project. Finally, the gradual TD Principal increase depicted in Q3.4 can be attributed to rapid bug fixing and code expansion, without properly testing previous and new implementations. Based on the above, we could state that the SonarQube TD measurement mechanism is indeed able to capture real events occurring during the software development process.

Similarly to the case of Project A, by reading participants' comments regarding the identified TD Principal trends in Table 13, we can observe that developers involved in Project B are also aware of specific actions they performed that justify these abrupt trends. More specifically, according to the participants, the gradual TD Principal increase of Project A depicted in Q3.5 can be attributed to specific strict deadlines, such as the delivery of the 1st and 2nd version of the product to the customer and a demo for a new contest, that forced the Company to involve more developers into this project. Similarly, the sharp TD Principal decrease depicted in Q3.6 can be attributed to code refactoring and defects fixing that the developers performed during the summer period. Finally, the sharp TD Principal increase depicted in Q3.7 can be attributed to new features that were added for a new customer and affected specific parts of code, thus making developers constantly having

to switch between two customers in a limited time. The comment from the developer stating that SonarQube did not run for a long time does not affect the magnitude of the TD increase but implies that in that case, the increase could be gradual instead of sharp. Similarly to the observations made regarding project A, in the case of Project B we could also state that SonarQube TD measurement mechanism can capture real events occurring during the software development process.

Summarizing the answers of the respondents regarding Q3.1 to Q3.7, that is, Part-3 of the questionnaire, we observe that most of the TD pattern types observed during the TD evolution of Projects A and B can be attributed to similar events that occurred during the software development cycle. More specifically, TD growth (either sharp or gradual) is mainly attributed to rapid code additions, usually without proper testing, in order to implement new features and functionalities that were requested by clients of the Company under strict time constraints. Similarly, temporal TD growth is related to quick and "dirty" code expansions in order to meet deadlines, which however were followed by prompt refactoring actions that improved the TD quality of the recently-added code. On the other hand, TD drop is attributed to heavy refactoring cycles that were performed during relatively relaxed periods, to repay the large amount of accumulated TD and thus, improve the quality and maintainability of the suboptimal

Table 13
Participants comments on TD Principal trends of Project B.

Cases	TD Principal Trend	Participants Comments (what is the cause of these TD changes?)
Q3.5	 <p>Gradual increase (23/2/19 – 28/7/19)</p>	<p>P1: "Many developers writing code with strict deadlines."</p> <p>P3: "1st and 2nd iteration to the customer plus a demo for a new contest in a new country at the same time."</p>
Q3.6	 <p>Sharp decrease (28/7/19 – 18/8/19)</p>	<p>P1: "The cause is refactor of code written quickly."</p> <p>P3: "Free time in summer to spend on refactoring and solving bugs vulnerabilities and code smells"</p>
Q3.7	 <p>Sharp increase (21/11/19 – 19/2/20)</p>	<p>P1: "Sonargube did not run for a long time."</p> <p>P3: "New customer came up asking new features that caused refactor to specific parts of code, plus simultaneously iterations to both customers in limited time."</p>

code introduced in the two projects during a long period of rapid code expansions mentioned above. The above events are in line with the definition of the TD metaphor and verify the necessity for which it was inspired in the first place. The quality compromises made by the Company during the studied period may have yielded the desired short-term benefits, such as the quick delivery of code to the clients, but have resulted in quality decay of the Company's software products, which made developers aware of the need to spend additional time on refactoring actions in order to bring the software back to a maintainable state.

Finally, we analyzed the answers from Part-4 of the questionnaire. In Part-4, TD forecasts for Project A and Project B were presented to the participants and they were asked if they would be willing to change anything in the planned development process based only on the projected TD Principal. In that way, the main goal of this part is to evaluate the practical usefulness and meaningfulness of TD forecasts using qualitative feedback from the developers. Figs. 19 and 20 illustrate 10 steps-ahead forecasts (red line) for Project A and Project B respectively, using Ridge Regression that performed better for short-term predictions. In this point, it should be noted that in the figures included in the questionnaire, the ground truth (blue line) from the starting point of the forecasts and onwards was hidden from the participants. However, we include this information here for reasons of completeness and validation of the forecasting approach presented in this work. In addition, the time point from which we decided to start our forecasts was carefully selected to signal a significant change in the current trend up to that point. The reason behind this choice is that we want to assess the willingness of developers to change their planned development processes, based on a change they cannot foresee just by looking at the past trend. Participants' answers to Q4.1 and Q4.2 are summarized below.

Q4.1 "The latest commits of Application A show relatively stable TD evolution. However, the forecast for 10 commits ahead shows a gradual increase in the TD principal. By having a look at

this forecast, would you change anything in the planned development process? Would you consider performing code refactoring in order to prevent this increase?"

P3: "This project is almost completed and minor changes and additions to functionality are expected. I would focus to solve any blocking, critical and major bugs to prevent software misbehavior during the use from the end user."

P2: "Yes, code refactoring is of utmost importance for not only preventing TD increase but also to decrease TD principal future rate."

P4: "Code refactoring is mandatory at this point but the manager does not approve that. That means that even if the development team wants to implement better architectures and refactor the code in order to be maintainable and re-usable this must be approved by the management. If the management has low priority on producing quality software the TD increase will be continuous. Also, if the code is mainly developed by junior engineers without guidance by senior engineers the TD will be incremental."

By revisiting Q3.4 in Part-3 of the questionnaire, we observe that this gradual TD Principal increase was attributed by the participants to rapid code expansion without properly testing previous and new implementations. This means that the developers invested more in new functionalities rather than refactoring the already existing code. However, in this question (Q4.1), the developers are presented with a forecast that predicts a gradual increase in the TD principal of the application during that period. Based on their answers, we notice that by having a look at this forecast they would reconsider this decision and proceed with refactoring instead of adding new code. In fact, by reading participants' answers to Q4.1, we observe that all respondents are actually willing to take action in order to reduce the TD Principal that is expected to increase based on the forecast shown

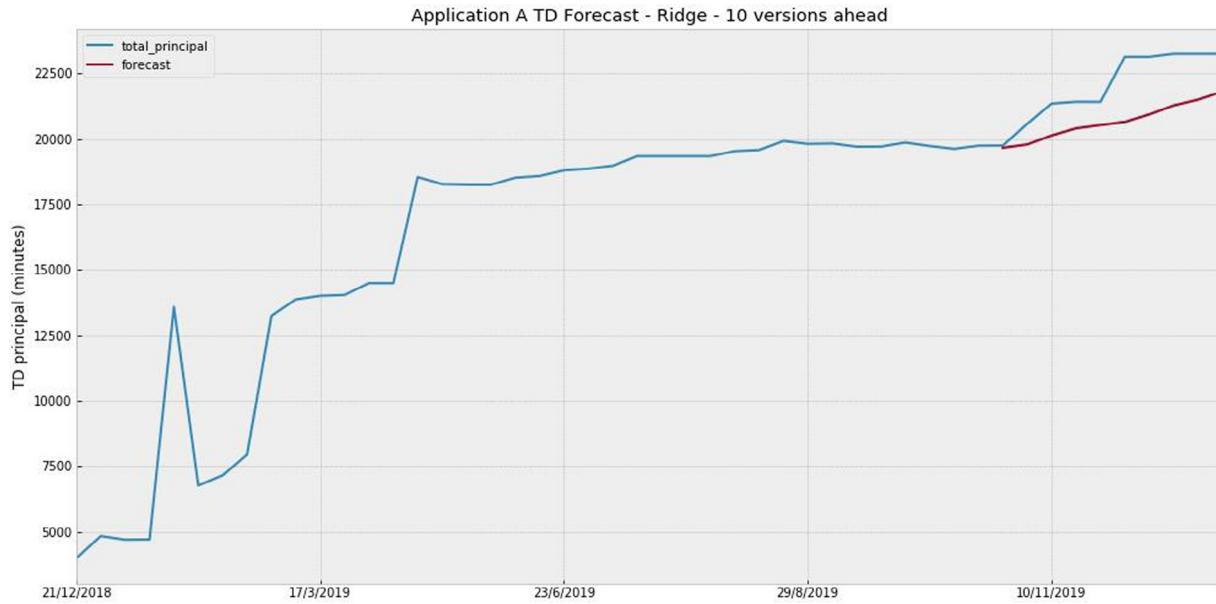


Fig. 19. Project A TD Principal forecasting for 10 steps ahead using Ridge (the ground truth from the starting point of the forecasts and onwards was hidden from the participants).

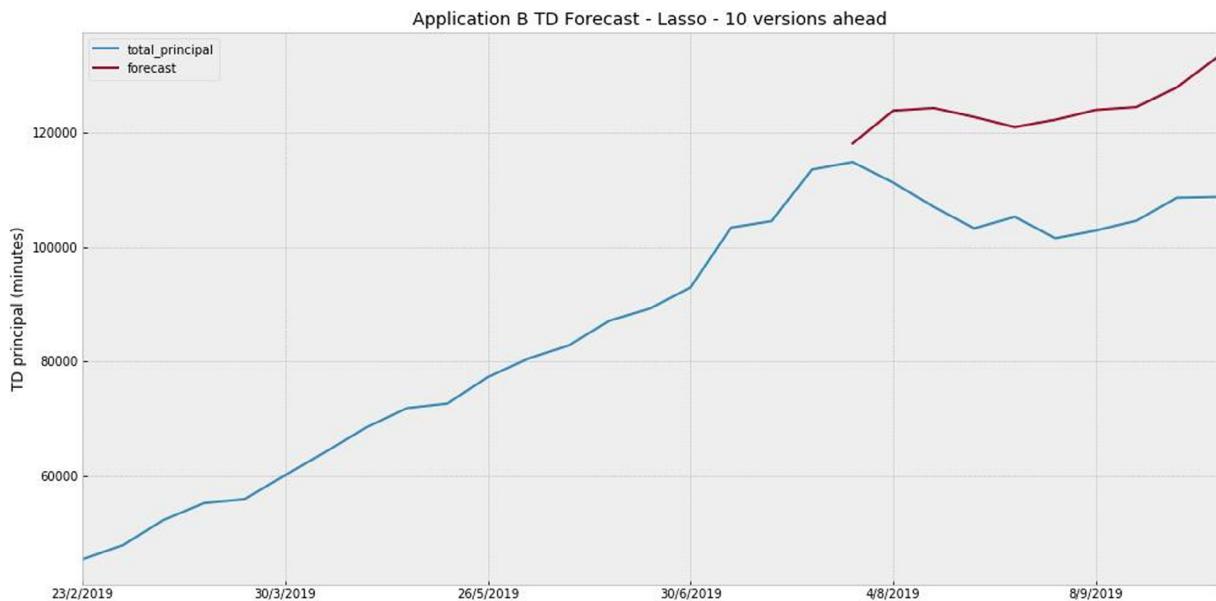


Fig. 20. Project B TD Principal forecasting for 10 steps ahead using Ridge (the ground truth from the starting point of the forecasts and onwards was hidden from the participants).

in Fig. 19. More specifically, the first two respondents state that they would perform refactoring to prevent future TD principal increasing rate and solve any defects that might arise to hinder the expected software behavior. The third respondent also agrees that code refactoring is mandatory at that point in order to prevent this increase. However, he/she states that any deviation from the planned development process must first be approved by the project manager. As a matter of fact, the perfect balance between repaying TD (and therefore increasing software quality) and reducing time to market of a software project is usually hard to achieve and lies at the decision-making abilities of the manager. This confirms the fact that future evolution of software

quality depends heavily on business-related parameters such as planned features, release deadlines etc.

Q4.2 “The latest commits of Application B show a gradual increase in the TD evolution. However, the forecast for 10 commits ahead shows a slight decrease in TD principal increasing rate. By having a look at this forecast, would you change anything in the planned development process? Would you consider investing in enhancing already existing, or adding new functionalities?”

P1: “For sure, we already have done some critical refactorings of code to decrease TD and from now on, that we have better handling of the project, we reassure that each line of code written, will not increase TD, but if it does, we apply refactors at the end of a sprint.”

P3: “This is a still developing project with lots of functionality to be added. Ideally, I would pause the developing process and refactor the existing code. However since future TD seems to decrease I want to try to decrease the TD in the new code, then when time plan allows it go back and refactor problematic areas.”

By revisiting Q3.6 in Part-3 of the questionnaire, we observe that this slight decrease in TD principal was attributed by the participants to code refactoring and defects fixing that the developers performed during the summer period. This means that during a relatively relaxed period, they invested time in cleaning up the already existing code and therefore decreasing the continuously growing TD Principal rate. While there is nothing wrong with this strategy, an alternative solution could be to invest in writing new but TD-free code instead of refactoring the existing one. In fact, writing new ‘clean’ code could be (in the long term) as efficient as code refactoring, especially when considering the difficulty of introducing heavy refactoring cycles in the industry, due to time limitations. In this question (Q4.1), the developers are presented with a forecast that predicts a slight decrease in TD principal increasing rate of the application during that period. Based on their answers, we notice that by having a look at this forecast they would reconsider this decision and proceed with adding new TD-free code instead of refactoring the existing one. In fact, by reading participants’ answers to Q4.2, we observe that both respondents state that they are willing to take advantage of the forecasted slight decrease in TD principal in order to pause refactoring activities and start adding new functionalities to the software application, by focusing more on adding new TD-free code. This would result in both preventing TD from increasing and at the same time delivering new functionalities. More specifically, the first respondent states that she/he would reassure that each new line of code written will not increase TD, but in case it does, they will apply refactoring at the end of the sprint. Similarly, the second respondent states that ideally she/he would pause the developing process and apply refactoring. However, since there is a lot of pending functionality to be added and the forecast shows a slight TD decrease, she/he would focus more on writing new TD-free code.

Through the industrial study reported in this section, we have a first level of validation that (a) actual TD trends reflect the circumstances and/or decisions taken during past development, and (b) forecasts derived through well-studied ML-models can provide valuable insights and affect developers’ decisions regarding the evolution of software. Of course, further research would be needed to solidify any claims on the usefulness of TD forecasting approaches, taking into account the numerous human- and business-related factors that drive the evolution of any software project.

7. Limitations and threats to validity

In this section, we discuss the limitations and validity threats of this empirical study. The accuracy of any forecasting model is by definition constrained, especially in the software domain, where future evolution of software quality depends heavily on numerous business-related factors such as planned features, release deadlines and fluctuations in the size of the development team. Therefore, anticipating such scheduled or unanticipated events would be a challenging endeavor beyond the scope of

our study. We believe that over longer time horizons, repeating phenomena are captured by a project’s history and building a prediction model based on historical data can provide some knowledge on future evolution. As described in Section 6, we have performed a study in an industrial setting to investigate the value of predictions regarding TD evolution. Nevertheless, we acknowledge the inability of the proposed approach to take into account planned or unforeseen business-related events. Apart from the aforementioned limitations, the methodology proposed in this paper suffers from the usual threats to external and internal validity.

External validity refers to the ability to generalize results. The results of the study are unavoidably subject to external validity threats, since the applicability of ML models to forecast TD is examined on a sample set of 15 applications. It is always possible that another set of applications might exhibit different phenomena. Nevertheless, the fact that the selected applications are quite diverse with respect to application domains, size, etc. partially mitigates threats to generalization. In addition, a large part of the proposed methodology consists of constructing forecasting models that learn from past versions and therefore can be easily adapted to any software application, as long as sufficient and reliable historic data are available. A similar threat stems from the fact that our dataset consists of open source Java applications, thus limiting the ability to generalize the conclusions to applications of a different domain or programming language. However, the process of building TD forecasting models described in this paper primarily builds upon the output of the tools used to compute software-related metrics that can act as indicators of the quality attribute of TD. This means that the proposed models can be easily adapted to forecast the TD of applications that are coded in a different programming language, as long as there are tools that support the extraction of software-related metrics that can act as TD indicators for the respective language. This also contributes to mitigating threats to generalization. However, since the dataset does not include industry applications, we cannot make any speculation on closed-source applications. Commercial systems as well as other object-oriented programming languages can be the subjects of further research. Finally, another possible threat to external validity is the small sample size of the survey performed within the context of the case study of this work, as reported in Section 6. In particular, the low number of participants and the small number of investigated software applications used for validation may have insufficient power to provide valuable insights regarding the meaningfulness of the proposed TD forecasting methodology in practice.

Concerning the *internal validity*, i.e., the possibility of having unwanted or unanticipated relationships between the parameters that might affect the variable that we are trying to predict, it is reasonable to assume that numerous other metrics that affect TD might have not been taken into consideration. However, the fact that we constructed our initial set of TD predictors based on software-related metrics that have been widely used in the literature as indicators of the presence of TD, such as OO software metrics, code smells and code issues extracted from ASA tools, limits this threat. Regarding the final selection of TD predictors, if we had limited our feature selection analysis to only correlations between the TD estimates and software-related metrics acting as TD predictors, then there would have been a threat to internal validity. However, we attempted to mitigate this threat through the use of univariate and multivariate regression analysis to further explore the relationships between the dependent and independent variables. Furthermore, in order to study the statistical significance of each indicator over the TD quality and be able to safely perform feature selection, we maximized diversity and representativeness by extending our dataset with additional 210 different heterogeneous applications.

Construct validity refers to the meaningfulness of measurements and that the independent and dependent variables are represented correctly. In this study, the main threats related to construct validity are due to possible inaccuracies in the identification of software-related metrics acting as TD indicators, as well as the identification and quantification of TD itself. In order to mitigate this risk, we decided to use two well-known and widely used tools, namely SonarQube and CKJM Extended. It should be noted that both of these tools were used as a proof of concept of the proposed forecasting methodology. The forecasting approach described in the present study is not dependent on the selected tools, as it could be applied to the measurements produced by other tools, based on user preference. However, the results presented in this study depend on the measurements obtained by these tools and, consequently, on the tools themselves. Therefore, more experimentation is required to assess the correctness of results obtained via other tools. As for the experimented prediction models, we exploited the ML algorithms implementation provided by the scikit-learn library, which is widely considered as a reliable tool. System-level forecasting also poses a threat to the validity of the findings as a tool for guiding TD repayment. In order for the refactoring activities to be more effective, recommendations for TD repayment need to be made at lower levels of granularity (e.g., class-level). However, the main goal of the proposed forecasting approach is to help developers and project managers make high-level decisions on whether there is a need to perform TD repayment in the next period of the project in general, and not to provide fine-grained recommendations on which software components the repayment activities should be focused.

Finally, *reliability* threats concern the possibility of replicating this study. To facilitate such replication studies, we provide an experimental package containing both the dataset and the scripts that were used for our analysis and forecasting model construction. This material can be found online.³⁶ Moreover, the source code repositories of the 15 selected projects are available on GitHub to obtain the same data.

8. Implications to researchers and practitioners

To the best of our knowledge, this is the first study in the field of TD that examines the applicability of ML models for TD forecasting. Across the 15 independently developed open source Java projects, our analysis indicates that linear Regularization models and the non-linear Random Forest regression are able to provide meaningful forecasts of TD evolution, and in most of the cases, with a sufficient level of accuracy. This work has significant implications for both research and practice, despite the limitations noted in the previous section.

8.1. Implications for research

Through our study, we identified some interesting open issues that should be addressed through further research. In particular, although there has been extensive research with respect to predicting the evolution of individual software features, quality attributes, and quality properties that are directly or indirectly related to the TD of a software project, no concrete contributions exist in the related literature regarding TD forecasting. Therefore, we believe that this study has a high impact on the scientific community and therefore, we suggest and encourage researchers to further explore this direction. An interesting topic of future work would be to extensively evaluate TD forecasting techniques

on a broader spectrum of real-world software applications covering different domains or programming languages. In addition, it would be useful to investigate different efficient ways to produce forecasting models for accurate prediction of TD principal and interest evolution, by bringing into the equation other types of software repositories that could be a potential source of TD related data, such as project management and issue-tracking systems. More specifically, mining TD related data from project- and issue-tracking systems, such as the reported effort of fixing bugs on Bugzilla or closing issues on Jira, could provide valuable information towards enhancing the TD forecasting approach. Ultimately, an approach that would pair all the above information with specialized techniques for forecasting, code analysis, software evolution analysis, and natural language processing could pave the way for the advance in the state of the art in this domain.

Predicting the future value of TD interest would be also critical for decision making, as it can be used to timely determine the point at which a software product would become unmaintainable, and therefore to respond promptly through appropriate refactoring activities in order to prevent this situation. More specifically, an interesting way of approaching the problem of TD interest forecasting would be to examine whether forecasting techniques could contribute towards enhancing the process of identifying the “breaking point” of an application, a term introduced by Chatzigeorgiou et al. (2015) that refers to the point in time where the accumulated interest is equal to the TD principal and, thus, the cost becomes higher than the benefit.

8.2. Implications for practice

Monitoring and forecasting the evolution of TD is highly important for software development companies, as it can provide valuable information regarding the effort and, in turn, the cost that is required for maintaining and extending a software application. Therefore, a TD forecasting methodology integrated into a relative tool, such as the outcome of the present research work, could be crucial for companies that want to remain competitive, while taking planned decisions regarding their TD management activities. In a hypothetical scenario where a software company has to make an investment to a specific application, our TD forecasting tool could provide an effective method to facilitate planning for budget and time allocation. More specifically, when the model predicts a declining number of TD, a project manager can then proactively allocate resources to software enhancements, or to other projects in more need. When the model predicts an increase in TD, an organization a priori can allocate the resources needed to quickly repay it by taking actions such as post development refactoring activities. This research has therefore the potential to make a great economic impact by helping software companies save budget by foreseeing TD accumulation and therefore avoid a potential bankruptcy in the future.

This empirical study has focused exclusively on modeling software evolution at the system level, thus allowing project managers to efficiently prioritize TD activities when dealing with different software applications. When it comes to a specific application however, the developers are often overwhelmed with a large volume of TD liabilities (e.g., code smells, bugs, vulnerabilities, etc.) that they need to fix. This renders the TD repayment procedure tedious, time consuming and effort demanding. In such cases, the significance of prioritizing which software components to refactor is highlighted even further, since fixing TD items in dormant parts of the code does not effectively affect maintenance costs (Schmid, 2013). As a future work, we will investigate the possibility of extending TD forecasting techniques to lower levels of granularity of a software project (e.g., package, class or function level). This would enable for a more granular prioritization of TD

³⁶ <https://sites.google.com/view/technical-debt-forecasting/main>

liabilities by incorporating information retrieved from TD forecasting techniques, allowing for a ranking of a software project's artifacts based on predictions of their long-term accumulated TD values.

9. Conclusions and future work

Technical Debt (TD) refers to inefficiencies during all phases of software development lifecycle that lead to extra maintenance effort. In recent years, TD has attracted the attention of both academia and industry. As a result, there has been a considerable increase in the number and provided functionality of methods and tools that support TD management. TD repayment, a high-level activity of TD management aims to resolve or mitigate TD in a software system by techniques such as reengineering and refactoring. However, a decision regarding whether to repay or not a TD item has different consequences depending on when it is made. This stresses the need for methods and accompanying tools that would enable system engineers and project managers to perform long-term effective software maintenance, by providing insights regarding where and when to apply refactoring. Therefore, what the stakeholders require is a decision-support system to help them make such choices and support decision-making under uncertainty. Under those circumstances, a method or tool able to track and forecast the evolution of TD of a software system could lead to the development of practical decision-making mechanisms aiming to improve the TD repayment strategy and estimate the point in which a software product could become unmaintainable.

The purpose of this paper is to examine whether and to what extend is the usage of ML models a meaningful and accurate approach to forecasting TD Principal in long-lived, open-source software applications (RQ1). Across the 15 independently developed, maintained, and managed open source projects, we have shown that TD Principal patterns can be modeled adequately by ML techniques. More specifically, for forecasting horizons between 1 and 20 weeks ahead, Regularization models (i.e., Lasso and Ridge regression) are able to fit and provide meaningful forecasts of TD Principal evolution. Trying to forecast longer into the future however, we noticed that their predictive power drops significantly. On the contrary, the non-linear Random Forest regression seems to have an almost stable performance over the holdout sample for all examined steps ahead. In fact, for forecasting horizons longer than 20 weeks ahead, Random Forest regression was able to capture the trend of the future evolution of the TD Principal with higher predictive power compared to the linear models. This indicates that a more complex model performs better when the forecast horizon is longer. From the above analysis, we can conclude that ML models constitute a suitable and effective approach for TD Principal forecasting, as they are able to fit and provide meaningful estimates of TD evolution over a relatively long period, while in most of the cases, the future TD value is captured with a sufficient level of accuracy.

To the best of our knowledge, this is the first study in the field of TD that examines the applicability of ML models for TD forecasting and therefore, it constitutes a good basis for future research and experimentation. Future work includes the extensive evaluation of TD forecasting techniques on a broader spectrum of real-world software applications, as well as to lower levels of granularity of a software project (e.g., package, class or function level). To solidify any claims on the usefulness of TD forecasting approaches, we plan to conduct an extended case study where we will provide practitioners with future predictions, track the actual development, and observe the impact of the forecasting results adoption. We also plan to investigate the ability of already examined or new forecasting models to provide more accurate

Table 14
Apache Ofbiz TD predictions using Walk-forward Train–Test validation.

Model	MAE (min)	RMSE (min)	MAPE (%)
MLR	11800.06	14939.92	5.70
Lasso regressor	2586.26	2881.59	1.25
Ridge regressor	3580.03	4043.14	1.73
SGD regressor	3998.45	4749.17	1.93
SVR regressor (linear)	4041.08	4945.64	1.96
SVR regressor (rbf)	5371.22	6492.32	2.60
Random forest regressor	1842.31	2100.87	0.89

Table 15
Apache SystemML TD predictions using Walk-forward Train–Test validation.

Model	MAE (min)	RMSE (min)	MAPE (%)
MLR	4489.38	5151.21	3.78
Lasso regressor	3949.70	4864.14	3.32
Ridge regressor	4125.15	5070.44	3.50
SGD regressor	4182.12	4712.35	3.53
SVR regressor (linear)	8771.71	10464.37	7.49
SVR regressor (rbf)	4201.63	4850.27	3.52
Random forest regressor	4096.68	4661.39	3.46

Table 16
Apache Groovy TD predictions using Walk-forward Train–Test validation.

Model	MAE (min)	RMSE (min)	MAPE (%)
MLR	44141.98	57592.53	20.41
Lasso regressor	24910.16	32278.33	11.56
Ridge regressor	26379.48	33607.69	12.16
SGD regressor	29316.00	38734.54	13.51
SVR regressor (linear)	22512.98	29346.43	10.56
SVR regressor (rbf)	15716.78	20060.46	7.38
Random forest regressor	10384.53	13494.23	4.92

Table 17
Apache Nifi TD predictions using Walk-forward Train–Test validation.

Model	MAE (min)	RMSE (min)	MAPE (%)
MLR	5293.48	6251.46	5.65
Lasso regressor	1920.77	2197.68	1.95
Ridge regressor	2041.70	2320.67	2.09
SGD regressor	2386.59	2678.10	2.47
SVR regressor (linear)	6359.22	7132.60	6.66
SVR regressor (rbf)	3852.61	4252.67	4.03
Random forest regressor	3233.18	3461.34	3.34

Table 18
Google Guava TD predictions using Walk-forward Train–Test validation.

Model	MAE (min)	RMSE (min)	MAPE (%)
MLR	11800.06	14939.92	5.70
Lasso regressor	2586.26	2881.59	1.25
Ridge regressor	3580.03	4043.14	1.73
SGD regressor	3998.45	4749.17	1.93
SVR regressor (linear)	4041.08	4945.64	1.96
SVR regressor (rbf)	5371.22	6492.32	2.60
Random forest regressor	1842.31	2100.87	0.89

predictions for even longer forecasting horizons. Last but not least, we plan to investigate other types of software repositories that could be a potential source of TD related data, such as project management and issue-tracking systems, as well as archived communication between project personnel. More specifically, the analysis of the communication between project personnel could reveal indications of high-TD artifacts that concentrate a large part of maintenance effort. These indications could then be factored in TD forecasting techniques to target these critical – from a maintenance point of view – artifacts. In fact, we believe that there is great potential in mining this information to achieve source triangulation and thus, yield more accurate TD forecasting estimates.

Table 19

Square Okhttp TD predictions using Walk-forward Train–Test validation.

Model	MAE (min)	RMSE (min)	MAPE (%)
MLR	990.90	1116.48	10.24
Lasso regressor	620.66	725.74	6.24
Ridge regressor	434.47	524.53	4.38
SGD regressor	902.72	1018.00	9.36
SVR regressor (linear)	3096.27	3414.23	32.97
SVR regressor (rbf)	749.76	862.88	7.66
Random forest regressor	316.21	383.99	3.20

Table 20

Square Retrofit TD predictions using Walk-forward Train–Test validation.

Model	MAE (min)	RMSE (min)	MAPE (%)
MLR	603.55	678.93	13.25
Lasso regressor	247.41	283.32	5.65
Ridge regressor	239.37	274.13	5.48
SGD regressor	470.46	517.63	10.42
SVR regressor (linear)	902.16	1061.68	19.68
SVR regressor (rbf)	423.43	469.03	9.54
Random forest regressor	198.28	233.01	4.67

Table 21

Spring Boot TD predictions using Walk-forward Train–Test validation.

Model	MAE (min)	RMSE (min)	MAPE (%)
MLR	128.56	151.70	5.50
Lasso regressor	122.08	141.07	5.23
Ridge regressor	124.33	142.50	5.30
SGD regressor	119.47	138.43	5.11
SVR regressor (linear)	642.80	786.51	28.05
SVR regressor (rbf)	141.77	163.29	6.11
Random forest regressor	101.38	116.44	4.30

Table 22

Apache CommonsIO TD predictions using Walk-forward Train–Test validation.

Model	MAE (min)	RMSE (min)	MAPE (%)
MLR	240.97	311.95	6.80
Lasso regressor	160.64	208.43	4.35
Ridge regressor	160.98	210.00	4.36
SGD regressor	233.35	294.95	6.52
SVR regressor (linear)	230.17	311.15	6.38
SVR regressor (rbf)	220.05	289.50	6.10
Random forest regressor	314.67	369.86	8.44

Table 23

Apache Incubator TD predictions using Walk-forward Train–Test validation.

Model	MAE (min)	RMSE (min)	MAPE (%)
MLR	8666.60	10957.26	22.13
Lasso regressor	3221.57	3770.47	8.33
Ridge regressor	3085.79	3576.34	7.97
SGD regressor	6060.73	7022.08	15.40
SVR regressor (linear)	2932.74	3404.16	7.56
SVR regressor (rbf)	2913.77	3432.83	7.51
Random forest regressor	2731.86	3110.55	7.06

Table 24

Java WebSocket TD predictions using Walk-forward Train–Test validation.

Model	MAE (min)	RMSE (min)	MAPE (%)
MLR	652.78	713.54	43.29
Lasso regressor	368.33	412.58	25.98
Ridge regressor	409.87	449.94	29.94
SGD regressor	645.63	703.79	43.82
SVR regressor (linear)	586.90	654.24	39.97
SVR regressor (rbf)	682.55	742.21	50.80
Random forest regressor	444.89	494.17	32.52

Table 25

Zxing TD predictions using Walk-forward Train–Test validation.

Model	MAE (min)	RMSE (min)	MAPE (%)
MLR	205.99	238.87	1.57
Lasso regressor	94.77	105.33	0.72
Ridge regressor	77.02	85.17	0.58
SGD regressor	105.05	119.87	0.79
SVR regressor (linear)	114.15	130.16	0.86
SVR regressor (rbf)	112.55	127.93	0.85
Random forest regressor	79.54	89.46	0.60

Table 26

Jenkins TD predictions using Walk-forward Train–Test validation.

Model	MAE (min)	RMSE (min)	MAPE (%)
MLR	2689.58	3244.52	3.35
Lasso regressor	2694.02	3216.50	3.35
Ridge regressor	3009.70	3519.39	3.73
SGD regressor	2286.04	2774.47	2.83
SVR regressor (linear)	4443.72	5393.16	5.54
SVR regressor (rbf)	2469.89	3029.77	3.05
Random forest regressor	2309.38	2737.34	2.87

Table 27

Openfire TD predictions using Walk-forward Train–Test validation.

Model	MAE (min)	RMSE (min)	MAPE (%)
MLR	9764.29	11254.28	18.30
Lasso regressor	7943.13	9464.07	14.94
Ridge regressor	9844.82	11965.04	17.84
SGD regressor	6584.61	7622.49	11.53
SVR regressor (linear)	12931.75	14325.71	21.48
SVR regressor (rbf)	8136.43	12571.07	14.42
Random forest regressor	10991.14	12090.32	17.55

CRediT authorship contribution statement

Dimitrios Tsoukalas: Conceptualization, Data curation, Formal analysis, Investigation, Methodology, Software, Visualization, Writing - original draft, Writing - review & editing. **Dionyssios Kehagias:** Formal analysis, Funding acquisition, Methodology, Project administration, Resources, Validation, Visualization, Writing - review & editing. **Miltiadis Siavvas:** Data curation, Formal analysis, Software, Visualization, Writing - review & editing. **Alexander Chatzigeorgiou:** Conceptualization, Investigation, Resources, Supervision, Validation, Visualization, Writing - review & editing.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgment

This work is funded by the European Union's Horizon 2020 Research and Innovation Programme through SDK4ED project under Grant Agreement No. 780572.

Appendix

See Tables 14–27.

References

- Alpaydin, E., 2010. *Introduction to Machine Learning*, second ed. The MIT Press.
 Altman, N.S., 1992. An introduction to kernel and nearest-neighbor nonparametric regression. Amer. Statist. 46 (3), 175–185. <http://dx.doi.org/10.2307/2685209>

- Alves, N.S.R., Mendes, T.S., de Mendonça, M.G., Spínola, R.O., Shull, F., Seaman, C., 2016. Identification and management of technical debt: A systematic mapping study. *Inf. Softw. Technol.* 70, 100–121. <http://dx.doi.org/10.1016/j.infsof.2015.10.008>.
- Ampatzoglou, A., Ampatzoglou, A., Avgeriou, P., Chatzigeorgiou, A., 2015b. Establishing a framework for managing interest in technical debt. <http://dx.doi.org/10.5220/0005885700750085>.
- Ampatzoglou, A., Ampatzoglou, A., Chatzigeorgiou, A., Avgeriou, P., 2015a. The financial aspect of managing technical debt: A systematic literature review. *Inf. Softw. Technol.* 64, 52–73. <http://dx.doi.org/10.1016/j.infsof.2015.04.001>.
- Ampatzoglou, A., Michailidis, A., Sarikyriakidis, C., Ampatzoglou, A., Chatzigeorgiou, A., Avgeriou, P., 2018. A framework for managing interest in technical debt: An industrial validation. <http://dx.doi.org/10.1145/3194164.3194175>.
- Arisholm, E., Briand, L.C., 2006. Predicting fault-prone components in a java legacy system.
- Baggen, R., Correia, J.P., Schill, K., Visser, J., 2012. Standardized code quality benchmarking for improving software maintainability. *Softw. Qual. J.* 20 (2), 287–307. <http://dx.doi.org/10.1007/s11219-011-9144-9>.
- Bansya, J., Davis, C.G., 2002. A hierarchical model for object-oriented design quality assessment. *IEEE Trans. Softw. Eng.* 28 (1), 4–17. <http://dx.doi.org/10.1109/32.979986>.
- Bellman, R.E., 2003. *Dynamic Programming*. Dover Publications.
- Bidarkota, P.V., 1998. The comparative forecast performance of univariate and multivariate models: an application to real interest rate forecasting. *Int. J. Forecast.* 14 (4), 457–468. [http://dx.doi.org/10.1016/S0169-2070\(98\)00036-3](http://dx.doi.org/10.1016/S0169-2070(98)00036-3).
- Boehm, B.W., et al., 1984. Software engineering economics. *IEEE Trans. Softw. Eng. SE-10* (1), 4–21. <http://dx.doi.org/10.1109/TSE.1984.5010193>.
- Bontempi, G., Taeib, S.B., Le Borgne, Y.-A., 2013. Machine learning strategies for time series forecasting. In: *Machine Learning Strategies for Time Series Forecasting*. Springer Berlin Heidelberg.
- Box, G.E., Jenkins, G.M., Reinsel, G.C., Ljung, G.M., 2015. *Time Series Analysis: Forecasting and Control*, fifth ed. John Wiley & Sons.
- Breiman, L., 1996. Bagging predictors. *Mach. Learn.* 24 (2), 123–140. <http://dx.doi.org/10.1007/BF00058655>.
- Breiman, L., 2017. *Classification and Regression Trees*. Routledge.
- Brown, N., et al., 2010. Managing technical debt in software-reliant systems. In: Proceedings of the Workshop on Future of Software Engineering Research. FSE/SDP, pp. 47–52. <http://dx.doi.org/10.1145/1882362.1882373>.
- Bruntink, M., van Deursen, A., 2006. An empirical study into class testability. *J. Syst. Softw.* 79 (9), 1219–1232. <http://dx.doi.org/10.1016/j.jss.2006.02.036>.
- Chaikalis, T., Chatzigeorgiou, A., 2015. Forecasting java software evolution trends employing network models. *IEEE Trans. Softw. Eng.* 41 (6), 582–602. <http://dx.doi.org/10.1109/TSE.2014.2381249>.
- Challagulla, V.U.B., Bastani, F.B., Paul, R.A., 2005. Empirical assessment of machine learning based software defect prediction techniques. In: 10th IEEE International Workshop on Object-Oriented Real-Time Dependable Systems. WORDS, pp. 263–270. <http://dx.doi.org/10.1109/WORDS.2005.32>.
- Charalampidou, S., Ampatzoglou, A., Chatzigeorgiou, A., Avgeriou, P., 2017. Assessing code smell interest probability: a case study. In: Proceedings of the XP2017 Scientific Workshops. p. 5. <http://dx.doi.org/10.1145/3120459.3120465>.
- Chatzigeorgiou, A., Ampatzoglou, A., Ampatzoglou, A., Amanatidis, T., 2015. Estimating the breaking point for technical debt. In: IEEE 7th International Workshop on Managing Technical Debt. MTD, pp. 53–56. <http://dx.doi.org/10.1109/MTD.2015.7332625>.
- Chidamber, S.R., Kemerer, C.F., 1994. A metrics suite for object oriented design. *IEEE Trans. Softw. Eng.* 20 (6), 476–493. <http://dx.doi.org/10.1109/32.295895>.
- Chowdhury, I., Zulkernine, M., 2011. Using complexity coupling and cohesion metrics as early indicators of vulnerabilities. *J. Syst. Archit.* 57, 294–313. <http://dx.doi.org/10.1016/j.syarc.2010.06.003>.
- Chug, A., Malhotra, R., 2016. Benchmarking framework for maintainability prediction of open source software using object oriented metrics. *Int. J. Innovative Comput. Inf. Control* 12 (2), 615–634.
- Cohen, J., 1977. *Statistical Power Analysis for the Behavioral Sciences*. Academic Press.
- Cunningham, W., 1993. The WyCash portfolio management system. *ACM SIGPLAN OOPS Messenger* 4 (2), 29–30. <http://dx.doi.org/10.1145/157710.157715>.
- Curtis, B., Sappidi, J., Szynkarski, A., 2012. Estimating the size, cost, and types of technical debt. In: Proceedings of the Third International Workshop on Managing Technical Debt, MTD, 2012, pp. 49–53.
- Das, J., 2012. *Statistics for Business Decisions*. Academic Publishers.
- Dieterich, T.G., 2002. Machine learning for sequential data: A review. In: Joint IAPR International Workshops on Statistical Techniques in Pattern Recognition (SPR) and Structural and Syntactic Pattern Recognition. SSPR, pp. 15–30. http://dx.doi.org/10.1007/3-540-70659-3_2.
- Digkas, G., Lungu, M., Avgeriou, P., Chatzigeorgiou, A., Ampatzoglou, A., 2018. How do developers fix issues and pay back technical debt in the apache ecosystem? In: 2018 IEEE 25th International Conference on Software Analysis, Evolution and Reengineering. SANER, pp. 153–163. <http://dx.doi.org/10.1109/SANER.2018.8330205>.
- Digkas, G., Lungu, M., Chatzigeorgiou, A., Avgeriou, P., 2017. The evolution of technical debt in the apache ecosystem. In: European Conference on Software Architecture. ECSA, pp. 51–66. http://dx.doi.org/10.1007/978-3-319-65831-5_4.
- Drucker, H., Burges, C.J., Kaufman, L., Smola, A.J., Vapnik, V., 1997. Support vector regression machines. In: Proceedings of the 9th International Conference on Neural Information Processing Systems, NIPS, pp. 155–161.
- du Preez, J., Witt, S.F., 2003. Univariate versus multivariate time series forecasting: an application to international tourism demand. *Int. J. Forecast.* 19 (3), 435–451. [http://dx.doi.org/10.1016/S0169-2070\(02\)00057-2](http://dx.doi.org/10.1016/S0169-2070(02)00057-2).
- Efronymson, M., 1960. Multiple regression analysis. In: *Mathematical Methods for Digital Computers*. pp. 191–203.
- El Emam, K., Benlarbi, S., Goel, N., Rai, S.N., 2001. The confounding effect of class size on the validity of object-oriented metrics. *IEEE Trans. Softw. Eng.* 27 (7), 630–650. <http://dx.doi.org/10.1109/32.935855>.
- Elish, M.O., 2010. Exploring the relationships between design metrics and package understandability: A case study. In: 2010 IEEE 18th International Conference on Program Comprehension. ICPC, pp. 144–147. <http://dx.doi.org/10.1109/ICPC.2010.43>.
- Elish, M.O., Elish, K.O., 2009. Application of TreeNet in predicting object-oriented software maintainability: A comparative study. In: 2009 13th European Conference on Software Maintenance and Reengineering. CSMR, pp. 69–78. <http://dx.doi.org/10.1109/CSMR.2009.57>, Mar.
- Eski, S., Buzluka, F., 2011. An empirical study on object-oriented metrics and software evolution in order to reduce testing costs by predicting change-prone classes. In: 2011 IEEE Fourth International Conference on Software Testing, Verification and Validation. ICST, pp. 566–571. <http://dx.doi.org/10.1109/ICSTW.2011.43>.
- Feurer, M., Klein, A., Eggensperger, K., Springenberg, J., Blum, M., Hutter, F., 2015. Efficient and robust automated machine learning. In: Proceedings of the 28th International Conference on Neural Information Processing Systems, Vol. 2. NIPS, pp. 2962–2970. <http://dx.doi.org/10.5555/2969442.2969547>.
- Fioravanti, F., Nesi, P., 2001. Estimation and prediction metrics for adaptive maintenance effort of object-oriented systems. *IEEE Trans. Softw. Eng.* 27 (12), 1062–1084. <http://dx.doi.org/10.1109/32.988708>.
- Fontana, F.A., Ferme, V., Spinelli, S., 2012. Investigating the impact of code smells debt on quality code evaluation. In: Proceedings of the Third International Workshop on Managing Technical Debt. MTD, pp. 15–22. <http://dx.doi.org/10.1109/MTD.2012.6225993>.
- Fontana, F.A., Mäntylä, M.V., Zanoni, M., Marino, A., 2016a. Comparing and experimenting machine learning techniques for code smell detection. *Empir. Softw. Eng.* 21 (3), 1143–1191. <http://dx.doi.org/10.1007/s10664-015-9378-4>.
- Fontana, F.A., Roveda, R., Zanoni, M., 2016b. Technical debt indexes provided by tools: a preliminary discussion. In: 2016 IEEE 8th International Workshop on Managing Technical Debt. MTD, pp. 28–31. <http://dx.doi.org/10.1109/MTD.2016.11>.
- Fowler, M., 1999. *Refactoring: Improving the Design of Existing Code*. Addison-Wesley Professional.
- Fowler, M., 2003. Technical debt. <http://www.martinfowler.com/bliki/TechnicalDebt.html>, (Accessed 30 July 2018).
- Fowler, M., 2009. Technical debt quadrant. <http://www.martinfowler.com/bliki/TechnicalDebtQuadrant.html>, (Accessed 30 July 2018).
- Friedman, J., Hastie, T., Tibshirani, R., 2001. *The Elements of Statistical Learning*, Vol. 1. In: Springer series in statistics New York.
- Gall, H.C., Lanza, M., 2006. Software evolution: analysis and visualization. In: Proceedings of the 28th International Conference on Software Engineering. ICSE, <http://dx.doi.org/10.1145/1134285.1134502>.
- Giger, E., Pinzger, M., Gall, H.C., 2012. Can we predict types of code changes? An empirical analysis. In: 2012 9th IEEE Working Conference on Mining Software Repositories. MSR, pp. 217–226. <http://dx.doi.org/10.1109/MSR.2012.6224284>.
- Godfrey, M.W., German, D.M., 2008. The past, present, and future of software evolution. In: Frontiers of Software Maintenance. FoSM, pp. 129–138. <http://dx.doi.org/10.1109/FOSM.2008.4659256>.
- Gondra, I., 2008. Applying machine learning to software fault-proneness prediction. *J. Syst. Softw.* 81 (2), 186–195. <http://dx.doi.org/10.1016/j.jss.2007.05.035>.
- Goulão, M., Fonte, N., Wermelinger, M., e Abreu, F.B., 2012. Software evolution prediction using seasonal time analysis: a comparative study. In: 16th European Conference on Software Maintenance and Reengineering. CSMR, pp. 213–222. <http://dx.doi.org/10.1109/CSMR.2012.30>.
- Goyal, P.K., Joshi, G., 2014. QMOOD metric sets to assess quality of Java program. In: 2014 International Conference on Issues and Challenges in Intelligent Computing Techniques. ICICT, pp. 520–533. <http://dx.doi.org/10.1109/ICICT.2014.6781337>.
- Griffith, I., Reimanis, D., Izurieta, C., Codabux, Z., Deo, A., Williams, B., 2014. The correspondence between software quality models and technical debt estimation approaches. In: Sixth International Workshop on Managing Technical Debt. MTD, pp. 19–26. <http://dx.doi.org/10.1109/MTD.2014.13>.

- Ho, T.K., 1998. The random subspace method for constructing decision forests. *IEEE Trans. Pattern Anal. Mach. Intell.* 20 (8), 832–844. <http://dx.doi.org/10.1109/34.709601>.
- Hoerl, A.E., Kennard, R.W., 2000. Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics* 42 (1), 55–67. <http://dx.doi.org/10.1080/00401706.2000.10485983>.
- ISO/IEC, 2011. ISO/IEC 25010 – Systems and Software Engineering – Systems and Software Quality Requirements and Evaluation (SQuaRE) – System and Software Quality Models. ISO/IEC.
- Izurieta, C., Vetrò, A., Zazworska, N., Cai, Y., Seaman, C., Shull, F., 2012. Organizing the technical debt landscape. In: Proceedings of the Third International Workshop on Managing Technical Debt. MTD, Zurich, Switzerland, pp. 23–26. <http://dx.doi.org/10.5555/2666036.2666040>.
- Jin, C., Liu, J., 2010. Applications of support vector machine and unsupervised learning for predicting maintainability using object-oriented metrics. In: International Conference on Multimedia and Information Technology, Vol. 1. MITA, pp. 24–27. <http://dx.doi.org/10.1109/MMIT.2010.10>.
- Jureczko, M., Spinellis, D., 2010. Using object-oriented design metrics to predict software defects. In: *Models and Methodology of System Dependability*. Oficyna Wydawnicza Politechniki Wrocławskiej, Wrocław, Poland, pp. 69–81.
- Karus, S., Dumas, M., 2012. Code churn estimation using organisational and code metrics: An experimental comparison. *Inf. Softw. Technol.* 54 (2), 203–211. <http://dx.doi.org/10.1016/j.infsof.2011.09.004>.
- Kaur, K., Anand, S., 2013. A maintainability estimation model and metrics for object-oriented design (MOOD). *Int. J. Adv. Res. Comput. Eng. Technol.* 2 (5).
- Kenmei, B., Antoniol, G., Di Penta, M., 2008. Trend analysis and issue prediction in large-scale open source systems. In: 12th European Conference on Software Maintenance and Reengineering. CSMR, pp. 73–82. <http://dx.doi.org/10.1109/CSMR.2008.4493302>.
- Khoshgoftaar, T.M., Allen, E.B., Deng, J., 2002. Using regression trees to classify fault-prone software modules. *IEEE Trans. Reliab.* 51 (4), 455–462. <http://dx.doi.org/10.1109/TR.2002.804488>.
- Khoshgoftaar, T.M., Munson, J.C., 1990. Predicting software development errors using software complexity metrics. *IEEE J. Sel. Areas Commun.* 8 (2), 253–261. <http://dx.doi.org/10.1109/49.46879>.
- Kitchenham, B.A., Pfleeger, S.L., 2002. Principles of survey research: Part 3: Constructing a survey instrument. *SIGSOFT Softw. Eng. Notes* 27 (2), 20–24. <http://dx.doi.org/10.1145/511152.511155>.
- Kohavi, R., et al., 1995. A study of cross-validation and bootstrap for accuracy estimation and model selection. In: Proceedings of the 14th International Joint Conference on Artificial Intelligence, Vol. 2. IJCAI, pp. 1137–1145. <http://dx.doi.org/10.5555/1643031.1643047>.
- Kruchten, P., Nord, R.L., Ozkaya, I., 2012. Technical debt: From metaphor to theory and practice. *IEEE Softw.* 29 (6), 18–21. <http://dx.doi.org/10.1109/MS.2012.167>.
- Lapedes, A., Farber, R., 1987. Nonlinear signal processing using neural networks: Prediction and system modelling. United States.
- Lehman, M.M., 1980. Programs, life cycles, and laws of software evolution. *Proc. IEEE* 68 (9), 1060–1076. <http://dx.doi.org/10.1109/PROC.1980.11805>.
- Lenarduzzi, V., Lomio, F., Taibi, D., Huttunen, H., 2019. On the fault proneness of sonarqube technical debt violations: A comparison of eight machine learning techniques. *Comput. Res. Repos. CoRR abs/1907.00376*, [Online]. Available: <http://arxiv.org/abs/1907.00376>.
- Letouzey, J.-L., 2012. The SQALE method for evaluating technical debt. In: Third International Workshop on Managing Technical Debt. MTD, pp. 31–36. <http://dx.doi.org/10.1109/MTD.2012.6225997>.
- Letouzey, J.-L., Ilkiewicz, M., 2012. Managing technical debt with the sqale method. *IEEE Softw.* 29 (6), 44–51. <http://dx.doi.org/10.1109/MS.2012.129>.
- Li, Z., Avgeriou, P., Liang, P., 2015. A systematic mapping study on technical debt and its management. *J. Syst. Softw.* 193–220. <http://dx.doi.org/10.1016/j.jss.2014.12.027>.
- Lucca, G.A.D., Fasolino, A.R., Tramontana, P., Visaggio, C.A., 2004. Towards the definition of a maintainability model for web applications. In: Proceedings of the Eighth Euromicro Working Conference on Software Maintenance and Reengineering. CSMR'04, USA, p. 279. <http://dx.doi.org/10.1109/CSMR.2004.1281430>.
- Makridakis, S., Spiliotis, E., Assimakopoulos, V., 2018. Statistical and machine learning forecasting methods: Concerns and ways forward. *PLoS One* 13 (3), e0194889. <http://dx.doi.org/10.1371/journal.pone.0194889>.
- Malhotra, R., Lata, K., 2018. On the application of cross-project validation for predicting maintainability of open source software using machine learning techniques. In: 2018 7th International Conference on Reliability, Infocom Technologies and Optimization. ICRITO, pp. 175–181. <http://dx.doi.org/10.1109/ICRITO.2018.8748749>.
- Mamun, M.A.A., Martini, A., Staron, M., Berger, C., Hansson, J., 2019. Evolution of technical debt: An exploratory study. In: 2019 Joint Conference of the International Workshop on Software Measurement and the International Conference on Software Process and Product Measurement. IWSM, pp. 87–102.
- Marinescu, R., 2012. Assessing technical debt by identifying design flaws in software systems. *IBM J. Res. Dev.* 56 (5), 9:1–9:13. <http://dx.doi.org/10.1147/JRD.2012.2204512>.
- Marquardt, D.W., 1970. Generalized inverses, ridge regression, biased linear estimation, and nonlinear estimation. *Technometrics* 12 (3), 591–612. <http://dx.doi.org/10.1080/00401706.1970.10488699>.
- McConnell, S., 2012. How to categorize and communicate technical debt. <https://www.castsoftware.com/blog/steve-mcconnell-on-categorizing-managing-technical-debt>, (Accessed 30 July 2018).
- Mens, T., 2008. Introduction and roadmap: History and challenges of software evolution. In: *Introduction and Roadmap: History and Challenges of Software Evolution*. Software Evolution. Springer (Chapter 1).
- Munson, J., Khoshgoftaar, T., 1990. Regression modelling of software quality: empirical investigation. *Inf. Softw. Technol.* 32 (2), 106–114. [http://dx.doi.org/10.1016/0950-5849\(90\)90109-5](http://dx.doi.org/10.1016/0950-5849(90)90109-5).
- Nagappan, N., Ball, T., 2005. Static analysis tools as early indicators of pre-release defect density. In: Proceedings of the 27th International Conference on Software Engineering. ICSE, pp. 580–586. <http://dx.doi.org/10.1109/ICSE.2005.1553604>.
- Nagappan, N., Ball, T., Zeller, A., 2006. Mining metrics to predict component failures. In: Proceedings of the 28th International Conference on Software Engineering. ICSE, pp. 452–461. <http://dx.doi.org/10.1145/1134285.1134349>.
- Nugroho, A., Visser, J., Kuipers, T., 2011. An empirical model of technical debt and interest. In: Proceedings of the 2nd Workshop on Managing Technical Debt. pp. 1–8. <http://dx.doi.org/10.1145/1985362.1985364>.
- Olbrich, S.M., Cruzes, D.S., Sjöberg, D.I., 2010. Are all code smells harmful? A study of God Classes and Brain Classes in the evolution of three open source systems. In: 2010 IEEE International Conference on Software Maintenance. ICSM, pp. 1–10. <http://dx.doi.org/10.1109/ICSM.2010.5609564>.
- Palit, A.K., Popovic, D., 2006. *Computational Intelligence in Time Series Forecasting: Theory and Engineering Applications*. Springer Science & Business Media.
- Palomba, F., Bavota, G., Di Penta, M., Fasano, F., Oliveto, R., De Lucia, A., 2018. On the diffuseness and the impact on maintainability of code smells: a large scale empirical investigation. *Empir. Softw. Eng.* 23 (3), 1188–1221. <http://dx.doi.org/10.1007/s10664-017-9535-z>.
- Raja, U., Hale, D.P., Hale, J.E., 2009. Modeling software evolution defects: a time series approach. *J. Softw. Maint. Evol.* 21 (1), 49–71. <http://dx.doi.org/10.1002/sm.398>.
- Riaz, M., Mendes, E., Tempero, E., 2009. A systematic review of software maintainability prediction and metrics. In: Proceedings of the 2009 3rd International Symposium on Empirical Software Engineering and Measurement. ESEM, pp. 367–377. <http://dx.doi.org/10.1109/ESEM.2009.5314233>.
- Roumani, Y., Nwankpa, J.K., Roumani, Y.F., 2015. Time series modeling of vulnerabilities. *Comput. Secur.* 51, 32–40. <http://dx.doi.org/10.1016/j.cose.2015.03.003>.
- Schmid, K., 2013. A formal approach to technical debt decision making. In: Proceedings of the 9th International ACM Sigsoft Conference on Quality of Software Architectures. QoSA, New York, NY, USA, pp. 153–162. <http://dx.doi.org/10.1145/2465478.2465492>.
- Seaman, C., Guo, Y., 2011. Measuring and monitoring technical debt. In: *Advances in Computers*, Vol. 82. Elsevier, pp. 25–46.
- Shatnawi, R., Li, W., 2008. The effectiveness of software metrics in identifying error-prone classes in post-release software evolution process. *J. Syst. Softw.* 81 (11), 1868–1882. <http://dx.doi.org/10.1016/j.jss.2007.12.794>.
- Shin, Y., Meneely, A., Williams, L., Osborne, J.A., 2011. Evaluating complexity code churn and developer activity metrics as indicators of software vulnerabilities. *IEEE Trans. Softw. Eng.* 37 (6), 772–787. <http://dx.doi.org/10.1109/TSE.2010.81>.
- Siavvas, M.G., Chatzidimitriou, K.C., Symeonidis, A.L., 2017a. QATCH—an adaptive framework for software product quality assessment. *Expert Syst. Appl.* 86, 350–366. <http://dx.doi.org/10.1016/j.eswa.2017.05.060>.
- Siavvas, M., Kehagias, D., Tzovaras, D., 2017b. A preliminary study on the relationship among software metrics and specific vulnerability types. In: 2017 International Conference on Computational Science and Computational Intelligence. CSCI, pp. 916–921. <http://dx.doi.org/10.1109/CSCI.2017.159>.
- Siavvas, M., et al., 2019. An empirical evaluation of the relationship between technical debt and software security. In: 9th International Conference on Information Society and Technology, ICIST 2019.
- Singh, Y., Saha, A., 2012. Prediction of testability using the design metrics for object-oriented software. *Int. J. Comput. Appl. Technol.* 44 (1), 12–22. <http://dx.doi.org/10.1504/IJCAT.2012.048204>.
- Sjöberg, D.I., Yamashita, A., Anda, B.C., Mockus, A., Dybå, T., 2012. Quantifying the effect of code smells on maintenance effort. *IEEE Trans. Softw. Eng.* 39 (8), 1144–1156. <http://dx.doi.org/10.1109/TSE.2012.89>.
- Skourletopoulos, G., Mavromoustakis, C.X., Bahsoon, R., Mastorakis, G., Pallis, E., 2014. Predicting and quantifying the technical debt in cloud software engineering. In: 19th International Workshop on Computer Aided Modeling and Design of Communication Links and Networks. CAMAD, pp. 36–40. <http://dx.doi.org/10.1109/CAMAD.2014.7033201>.

- Spearman, C., 1987. The proof and measurement of association between two things. *Am. J. Psychol.* 100 (3/4), 441–471. <http://dx.doi.org/10.2307/1422689>.
- Stone, M., 1974. Cross-validatory choice and assessment of statistical predictions. *J. R. Stat. Soc. Ser. B Stat. Methodol.* 36 (2), 111–133. <http://dx.doi.org/10.1111/j.2517-6161.1974.tb00994.x>.
- Suryanarayana, G., Samarthyan, G., Sharma, T., 2014. Refactoring for Software Design Smells: Managing Technical Debt. Morgan Kaufmann.
- Tan, J., Lungu, M., Avgeriou, P., 2018. Towards studying the evolution of technical debt in the python projects from the apache software ecosystem. In: 17th Belgium-Netherlands Software Evolution Workshop, BENEVOL, pp. 43–45.
- Tibshirani, R., 1996. Regression shrinkage and selection via the lasso. *J. R. Stat. Soc. Ser. B Stat. Methodol.* 58 (1), 267–288. <http://dx.doi.org/10.1111/j.2517-6161.1996.tb02080.x>.
- Tsoukalas, D., Jankovic, M., Siavvas, M., Kehagias, D., Chatzigeorgiou, A., Tzovaras, D., 2019. On the applicability of time series models for technical debt forecasting. In: 15th China-Europe International Symposium on Software Engineering Education (CEISSE 2019). (In press).
- Tsoukalas, D., Siavvas, M., Jankovic, M., Kehagias, D., Chatzigeorgiou, A., Tzovaras, D., 2018. Methods and tools for TD estimation and forecasting: A state-of-the-art survey. In: IEEE International Conference on Intelligent Systems (IS 2018). pp. 698–705. <http://dx.doi.org/10.1109/IS.2018.8710521>.
- Van Koten, C., Gray, A., 2006. An application of Bayesian network for predicting object-oriented software maintainability. *Inf. Softw. Technol.* 48 (1), 59–67. <http://dx.doi.org/10.1016/j.infsof.2005.03.002>.
- Vetro', A., 2012. Using automatic static analysis to identify technical debt. In: Proceedings of the 34th International Conference on Software Engineering, ICSE, Zurich, Switzerland, pp. 1613–1615. <http://dx.doi.org/10.5555/2337223.2337499>.
- Wagner, S., 2009. A Bayesian network approach to assess and predict software quality using activity-based quality models. In: Proceedings of the 5th International Conference on Predictor Models in Software Engineering, Vol. 1. PROMISE, <http://dx.doi.org/10.1145/1540438.1540447>.
- Wagner, S., et al., 2015. Operationalised product quality models and assessment: The Quamoco approach. *Inf. Softw. Technol.* 62, 101–123. <http://dx.doi.org/10.1016/j.infsof.2015.02.009>.
- Walden, J., Stuckman, J., Scandariato, R., 2014. Predicting vulnerable components: Software metrics vs text mining. In: International Symposium on Software Reliability Engineering, ISSRE, pp. 23–33. <http://dx.doi.org/10.1109/ISSRE.2014.32>.
- Werbos, P., 1974. Beyond Regression: New Tools for Prediction and Analysis in the Behavioral Sciences. Harvard University, Cambridge.
- Werbos, P.J., 1988. Generalization of backpropagation with application to a recurrent gas market model. *Neural Netw.* 1 (4), 339–356. [http://dx.doi.org/10.1016/0893-6080\(88\)90007-X](http://dx.doi.org/10.1016/0893-6080(88)90007-X).
- Xuan, J., Hu, Y., Jiang, H., 2017. Debt-prone bugs: Technical debt in software maintenance. *Comput. Res. Repos. CoRR*, abs/1704.04766, [Online]. Available: <http://arxiv.org/abs/1704.04766>.
- Yazdi, H.S., Mirbolouki, M., Pietsch, P., Kehrer, T., Kelter, U., 2014. Analysis and prediction of design model evolution using time series. In: International Conference on Advanced Information Systems Engineering, CAiSE, pp. 1–15. http://dx.doi.org/10.1007/978-3-319-07869-4_1.
- Zazworka, N., Seaman, C., Shull, F., 2011. Prioritizing design debt investment opportunities. In: Proceedings of the 2nd Workshop on Managing Technical Debt, MTD, pp. 39–42. <http://dx.doi.org/10.1145/1985362.1985372>.
- Zazworka, N., Spínola, R.O., Vetro, A., Shull, F., Seaman, C., 2013. A case study on effectively identifying technical debt. In: Proceedings of the 17th International Conference on Evaluation and Assessment in Software Engineering, EASE, <http://dx.doi.org/10.1145/2460999.2461005>.
- Zazworka, N., et al., 2014. Comparing four approaches for technical debt identification. *Softw. Qual. J.* 22 (3), 403–426. <http://dx.doi.org/10.1007/s11219-013-9200-8>.
- Zhou, Y., Leung, H., 2007. Predicting object-oriented software maintainability using multivariate adaptive regression splines. *J. Syst. Softw.* 80 (8), 1349–1361. <http://dx.doi.org/10.1016/j.jss.2006.10.049>.
- Zhou, Y., Xu, B., 2008. Predicting the maintainability of open source software using design metrics, Wuhan university. *J. Nat. Sci.* 13 (1), 14–20. <http://dx.doi.org/10.1007/s11859-008-0104-6>.
- Zhou, Y., et al., 2012. An in-depth investigation into the relationships between structural metrics and unit testability in object-oriented systems. *Sci. China Inf. Sci.* 55, 2800–2815. <http://dx.doi.org/10.1007/s11432-012-4745-x>.

Dimitrios Tsoukalas holds a B.Sc. in “Applied Informatics” and a M.Sc. in “Computer Systems” from the Department of Applied Informatics - University of Macedonia (UoM). He also holds a M.Sc. in “Intelligent Systems – Computational Intelligence Methodologies and Applications” from the Department of Electrical and Computer Engineering – Aristotle University of Thessaloniki (A.U.Th.). Since May 2018, he is a Ph.D. candidate at the Department of Applied Informatics of the University of Macedonia (UoM). His main research interests lie in the areas of Software Engineering and Intelligent Systems. He is currently a Research Associate at the Centre for Research and Technology Hellas (CERTH/ITI).

Dionyios Kehagias received the Diploma and Ph.D. degrees in electrical and computer engineering from the Aristotle University of Thessaloniki, Thessaloniki, Greece, in 1999 and 2006, respectively. He is currently a Researcher Grade B with the Information Technologies Institute of the Centre for Research and Technology Hellas (CERTH). His research interests include software technologies, algorithms, data mining, machine learning, time-series analysis, big data analytics, service-oriented architectures and ontology-based knowledge engineering.

Miltiadis Siavvas received the Diploma degree in electrical and computer engineering from the Aristotle University of Thessaloniki in 2016, and the Ph.D. degree in Software Security and Reliability from the Intelligent Systems and Networks Group, Imperial College London, in 2019. He is currently a Research Associate at the Information Technologies Institute of the Centre for Research and Technology Hellas (CERTH/ITI). His main research interests lie in the areas of software engineering (with emphasis on software quality, reliability, and security), machine learning and data science.

Dr. Alexander Chatzigeorgiou is a Professor of Software Engineering in the Department of Applied Informatics at the University of Macedonia, Thessaloniki, Greece and since 2017 he serves as the Dean of the School of Information Sciences. He received the Diploma in Electrical Engineering and the Ph.D. degree in Computer Science from the Aristotle University of Thessaloniki, Greece, in 1996 and 2000, respectively. From 1997 to 1999 he was with Intracom, as a software designer. His research interests include object-oriented design, software maintenance and evolution. He has published more than 140 articles in international journals and conferences.