

Time for Addressing Software Security Issues: Prediction Models and Impacting Factors

Lotfi Ben Othmane¹ · Golriz Chehrazi¹ · Eric Bodden¹ · Petar Tsalovski² · Achim D. Brucker²

Received: 24 March 2016/Revised: 3 August 2016/Accepted: 29 August 2016/Published online: 27 September 2016
© The Author(s) 2016. This article is published with open access at Springerlink.com

Abstract Finding and fixing software vulnerabilities have become a major struggle for most software development companies. While generally without alternative, such fixing efforts are a major cost factor, which is why companies have a vital interest in focusing their secure software development activities such that they obtain an optimal return on this investment. We investigate, in this paper, quantitatively the major factors that impact the time it takes to fix a given security issue based on data collected automatically within SAP's secure development process, and we show how the issue fix time could be used to monitor the fixing process. We use three machine learning methods and evaluate their predictive power in predicting the time to fix issues. Interestingly, the models indicate that vulnerability type has less dominant impact on issue fix time than previously believed. The time it takes to fix an issue instead seems much more related to the component in which the potential vulnerability resides, the project related to the issue, the development groups that address the issue, and the closeness of the software release date. This indicates that the software structure, the fixing processes, and the development groups are the dominant factors that impact the time spent to address security issues. SAP can use the models to implement a continuous improvement of its secure software development process and to measure the impact of individual improvements. The development teams at SAP develop different types of software, adopt different internal development processes, use different programming languages and

platforms, and are located in different cities and countries. Other organizations, may use the results—with precaution—and be learning organizations.

Keywords Human factors · Secure software · Issue fix time

1 Introduction

Fixing vulnerabilities, before and after a release, is one of the most costly and unproductive software engineering activities. Yet, it comes with few alternatives, as code-level vulnerabilities in the application code are the basis of increasingly many exploits [1]. Large software development enterprises, such as SAP, embed in their development process activities for identifying vulnerabilities early, such as dynamic and static security testing [2]. Next to that, SAP's security development lifecycle (see, e.g., [3] for Microsoft's security development lifecycle) includes a process for fixing vulnerabilities after a software release.

Analyzing and fixing security issues is a costly undertaking that impacts a software's time to market and increases its overall development and maintenance cost. In result, software development companies have an interest to determine the factors that impact the effort, and thus, the time it takes to fix security issues, in particular to:

- identify time-consuming factors in the secure development process,
- better understand affecting factors,
- focus on important factors to enhance software's security level,
- accelerate secure software development processes, and
- enhance security cost planning for software development projects.

✉ Lotfi Ben Othmane
lotfi.ben.othmane@sit.fraunhofer.de

¹ Fraunhofer Institute for Secure Information Technology, Darmstadt, Germany

² SAP SE, Walldorf, Germany

In a previous study, Othmane et al. [4] conducted expert interviews at SAP to identify factors that impact the effort of fixing vulnerabilities. SAP collects data about fixing security issues (potential vulnerabilities that need to be analyzed further manually to ensure whether they are vulnerabilities or false positive issues) both during a software's development and after its release. With this study we supplement the previous qualitative, interview-based results with objectively gathered system data. In this study, we used these data to identify and quantify, using machine learning, to what extent automatically measured factors impact a given issue's fix time. By *issue fix time*, we mean the duration between the time at which a security issue is reported to SAP and the time at which the issue is marked as closed in number of days. For simplicity, we use the term issue to refer to a security issue in the remaining of the paper.

Vulnerabilities are subset of software defects; they allow violation of constraints that can lead to malicious use of the software. The information, tools, and expertise that help to analyze faults (functionalities errors) apply with limited efficacy to analyze vulnerabilities. Zimmermann et al. [5] reported, for example, that the number of vulnerabilities is highly correlated with code dependencies, while the metrics that are correlated with faults such as size of changed code have only small effect. In addition, detecting faults requires exercising the specified functionalities of the software, while vulnerabilities analysis requires the software developers to have the knowledge and expertise to think like attackers [6]. Moreover, vulnerabilities detection tools do not provide often sufficient information to locate the issue easily, besides that they report high number of false positive. Thus, it is believed that "finding vulnerabilities is akin to finding a needle in a haystack" [6]. Also vulnerabilities issues occur much less frequently than faults. Thus, models derived from data related to faults and vulnerabilities issues have to deal with unbalanced datasets, c.f., [7]. Therefore, better prediction models of issue fix time should use only security issues data and consider the characteristics of issues.

For the analysis, we use five data sources based on distinct system tools available at SAP. The first three main data sources relate to security issues; issues found by code scanners for the programming language ABAP [8] (Data source 1) and for Java, JavaScript, and C (Data source 2), as well as issues found in already released code, which are communicated through so-called security messages, for instance reported by customers, security experts or SAP's own security team (Data source 3). The other two data sources comprise support data. They describe the components, i.e., a group of applications that performs a common business goal such as sales order or payroll (Support data 1), and the projects (Support data 2).

After cleaning the data, we used three methods to develop prediction models, based on (1) linear regression, (2) Recursive PARTitioning (RPART), and (3) Neural Network Regression (NNR). Next, we measured the models' accuracy using three different metrics. Interestingly, the models indicate that the impact of a vulnerability's type (buffer overflow, cross-site-scripting, etc.) is less dominant than previously believed. Instead, the time it takes to fix an issue is more related to the component in which the vulnerability resides, the project related to the issue, to the development groups that address the issue, and to the closeness of the software release date.

SAP can use the results of this study to identify costly pain points and important areas in the secure development process and to prioritize improvements to this process. Such models can be used to establish a learning organization, which learns and improves its processes based on the company-specific actual facts reflected in the collected data [9]. Since SAP collects the models' input data continuously, the models can be used to analyze the company's processes and measure the impact of enhancements over time.

This paper is organized as follows. First, we give an overview of related work (Sect. 2), discuss SAP's approach to secure software development (Sect. 3), and provide an overview of the regression methods and model accuracy metrics that we use in the study (Sect. 4). Next, we describe the research methodology that we applied (Sect. 5), report about our findings (Sect. 6), and analyze the factors that impact the issue fix time (Sect. 7). Subsequently, we discuss the impacts and the limitations of the study (Sect. 8) and the main lessons (and surprises) that we learned (Sect. 9), and conclude the paper.

2 Related Work

There is related work on prediction models for development efforts and time to fix bugs but work in the area of effort estimation for fixing security issues is scarce. Thus, we discuss in this section related work that investigates influencing factors on issue fix time or vulnerability fix time and also the development of prediction models for effort estimations, and differentiates them from our work.

Cornell measured the time that the developers spent fixing vulnerabilities in 14 applications [10]. Table 1 shows the average time the developers in the study take to fix vulnerabilities for several vulnerability types. Cornell found that there are vulnerability types that are easy to fix, such as dead code, vulnerability types that require applying prepared solutions, such as a lack of authorization, and vulnerability types that, although simple conceptually, may require a long time to fix for complex cases, such as SQL

Table 1 Examples of time required for fixing vulnerabilities [10]

Vulnerability type	Average fix time (min)
Dead code (unused methods)	2.6
Lack of authorization check	6.9
Unsafe threading	8.5
XSS (stored)	9.6
SQL injection	97.5

injection. The vulnerability type is thus one of the factors that indicate the vulnerability fix time but is certainly not the only one [4].

In previous work, Othmane et al. [4] reported on a qualitative study conducted at SAP to identify the factors that impact the effort of fixing vulnerabilities and thus the vulnerability fix time. The study involved interviews with 12 security experts. Through these interviews, the authors identified 65 factors that include, beside the vulnerabilities characteristics, the structure of the software involved, the diversity of the used technologies, the smoothness of the communication and collaboration, the availability and quality of information and documentation, the expertise and knowledge of developers and security coordinators, and the quality of the code analysis tools.

Several studies aim at predicting the time to fix bugs [11–17]. Zhang et al. [18] conducted an empirical study on three open-source software to examine what factors affect the time between bug assignment to a developer and the time bug fixing starts, that is the developer's delay (when fixing bugs), along three dimensions: bug reports, source code, and code changes. The most influencing factor found was the issue's level of severity. Other factors are of technical nature, such as sum of code churn, code complexity or number of methods in changed files as well as the maximum length of all comments in a bug report. Similar to our study, Zhang et al. were interested in revealing the factors that impact time, but as opposed to them we focus on security issues, not on bugs, and include in our analysis not only automatically collected information about security issues before and after release, but additionally component- and project-related factors from which human-based and organizational factors can be derived. In contrast to Zhang et al., we consider the overall fix time that starts at the time when a security issue is reported and ends when the issue is marked as closed.

Menzies et al. [19] estimated projects development effort, using project-related data, such as the type of teams involved, the development time of the projects, and the number of high-level operations within the software. They found that it is better to use local data based on related projects instead of global data, which allows to account for project-related particularities that impact the development

effort. Their data sample is a “global dataset” that includes data from several research software projects conducted by different entities. We believe that the issue is related to contextual information that are not captured by the data and are related to the entities. Instead, our data are related to thousands of projects developed by quasi-independent development teams that use, for example, different programming languages and platforms (e.g., mobile, cloud, Web), adopt different internal processes, and are located in several countries and cities.

In another study, Menzies et al. [20] reassured the usefulness of static code attributes to learn defect predictors. They showed that naive Bayes machine learning methods outperform rule-based or decision tree learning methods and they showed, on the other hand, that the choice of learning methods used for defect predictions can be much more important than used attributes. Unlike this previous work, we use static code attributes to predict issue fix time and we use neural networks as additional method for prediction.

Following the objective to reduce effort for security inspection and testing, Shin et al. [21] used in their empirical study code complexity, code churn, and developer activity metrics obtained to predict vulnerable code locations with logistic regressions. They also used J48 decision trees, random forest, and Bayesian network classification techniques based on data obtained from two large-scale open-source projects using code characteristics and version control data. They found out that the combination of these metrics is effective in predicting vulnerable files. Nevertheless, they state that further effort is necessary to characterize differences between faults and vulnerabilities and to enhance prediction models. Unlike Shin et al., our empirical research focuses on predictions using system-based data to predict vulnerability fix time.

Hewett and Kijisanayothin [17] developed models for defect repair time prediction using seven different machine learning algorithms, e.g., decision trees and support vector machines. Their predictive models are based on a case study with data from a large medical software system. Similar to our approach they consider the whole repair time including all phases of a defect lifecycle. They use twelve defect attributes selected by domain experts for their estimations such as component, severity, start and end date, and phase. Unlike them we are interested in estimating vulnerability fix time not defect fix time.

In contrast to prior work, which often is based on open-source software, we estimate the vulnerability fix time based on an industrial case study of a major software development company, based on distinct datasets that include security issues before and after release and combine them with project and component-related data. Our objective is to identify the impacting strength of the factors

on vulnerability fix time as well as to predict issue fix time in general.

3 Secure Software Development at SAP

To ensure a secure software development, SAP follows the SAP Security Development Lifecycle (S²DL). Figure 1 illustrates the main steps in this process, which is split into four phases: preparation, development, transition, and utilization.

To allow the necessary flexibility to adapt this process to the various application types (ranging from small mobile apps to large-scale enterprise resource-planning solutions) developed by SAP as well as the different software development styles and cultural differences in a worldwide distributed organization, SAP follows a two-staged security expert model:

1. a central security team defines the global security processes, such as the S²DL, provides security trainings, risk identification methods, offers security testing tools, or defines and implements the security response process;
2. local security experts in each development area/team are supporting the developers, architects, and product owners in implementing the S²DL and its supporting processes.

For this study, the *development* and *utilization* phases of the S²DL are the most important ones, as the activities carried out during these phases detect most of the vulnerabilities that need to be fixed:

- during the actual software development (in the steps *secure development* and *security testing*) vulnerabilities are detected, e.g., by using manual and automated as well as static and dynamic methods for testing application security [2, 22]. Most detected vulnerabilities are found during this step, i.e., most vulnerabilities are fixed in unreleased code (e.g., in newly developed code that is not yet used by customers);
- *security validation* is an independent quality control that acts as “first customer” during the transition from software development to release, i.e., security validation finds vulnerabilities after the code freeze (called correction close) and the actual release;

- *security response* handles issues reported after the release of the product, e.g., by external security researchers or customers.

If an issue is confirmed (e.g., by an analysis of a security expert), from a high-level perspective developers and their local security experts implement the following four steps: 1. analyze the issue, 2. design or select a recommended solution, 3. implement and test a fix, and 4. validate (e.g., by retesting the fixed solution) and release this fix. Of course, the details differ depending on the development model of the product team and, more importantly, depending on whether the issue is detected in code that is used by customers or not.

While the technical steps for fixing an issue are the same regardless of whether the issue is in released code or currently developed code, the organizational aspects differ significantly: For vulnerabilities in unreleased development code, detecting, confirming, and fixing vulnerabilities are lightweight process defined locally by the development teams. Vulnerabilities detected by security validation, e.g., after the code freeze, even if in unreleased code, involve much larger communication efforts across different organizations for explaining the actual vulnerabilities to development as well as ensuring that the vulnerability is fixed before the product is released to customers.

Fixing vulnerabilities in released code requires the involvement of yet more teams within SAP, as well as additional steps, e.g., for back-porting fixes to older releases and providing patches (called *SAP Security Notes*) to customers.

Let us have a closer look on how an externally reported vulnerability in a shipped software version is fixed: First, an external reporter (e.g., customer or independent security researcher) contacts the *security response team*, which assigns a case manager. The case manager is responsible for driving the decision if a reported problem is a security vulnerability that needs to be fixed, and for ensuring that the confirmed vulnerability is fixed and that a patch is released. After vulnerability is confirmed, the case manager contacts the development team and often also a dedicated maintenance team (called IMS) to ensure that a fix is developed and back-ported to all necessary older releases (according to SAP’s support and maintenance contracts). The developed fixes are subject to a special security test by the security validation team and, moreover, the response



Fig. 1 Overview of the SAP Security Development Lifecycle (S²DL)

teams reviews the SAP Security Note. If the technical fix as well as the resulting Security Note passes the quality checks, the Security Note is made available to customers individually and/or in the form of a support package (usually on the first Tuesday of a month). Support packages are functional updates that also contain the latest security notes.

4 Background

Assume a response variable y and a set of independent variables x_i such that $y = f(x_1, x_2, \dots, x_n)$ where f represents the systematic information that the variables x_i provide about y [23]. Regression models relate the *quantity* of a response factor, i.e., dependent variable to the independent variables.

Different regression models have different capabilities, e.g., in terms of their resistance to outliers, their fit for small datasets, and their fit for a large number of predicting factors [24]. However, in general, a regression model is assumed to be good, if it predicts responses close to the actual values observed in reality. In this study, the performance of a given model is judged by its prediction errors; the low are the errors, the better is the performance of the model.

This section provides background about the regression methods, model's performance metrics, and a metric for measuring the relative importance of the prediction factors used in the models.

4.1 Overview of Used Regression Methods

We give next an overview of the three methods used in this study.

Linear Regression This method assumes that the regression function is linear in the input [25], i.e., in the prediction factors. The linear method has the advantage of being simple and allows for an easy interpretation of the correlations between the input and output variables.

Tree-Based Regression This method recursively partitions the observations, i.e., the data records of the object being analyzed, for each of the prediction factors (aka features) such that it reduces the value of a metric that measures the information quantity of the splits [26]. In this study, we use the method *recursive partitioning and regression trees (RPART)* [27].

Neural Network Regression This method represents functions that are nonlinear in the prediction variables. It uses a multi-layer network that relates the input to the output through intermediate nodes. The output of each intermediate node is the sum of weighted input of the nodes of the previous layer. The data input is the first layer [28].

These three regression methods are the basic ones that are commonly used in data analytics. In this study, we use their implementations in packages for the statistics language R:¹ *rpart*² for RPART, and *nnet*³ for NNR. The implementation “lm” of the linear regression (LR) is already contained within the core of R.

4.2 Model Performance Metrics

Regression methods infer prediction models from a given set of training data. Several metrics have been developed to compare the performance of the models in terms of accuracy of the generated predictions [29]. The metrics indicate how well the models predict accurate responses for future inputs. Next, we describe the three metrics that we used in this work, the coefficient of determination (R^2) [30], the Akaike information criterion (AIC) [29] and the prediction at a given level (PRED) [30].⁴

Coefficient of Determination (R^2) This metric “summarizes” how well the generated regression model fits the data. It computes the proportion of the variation of the response variable as estimated using the generated regression compared to the variation of the response variable computed using the null model, i.e., the mean of the values [29]. The following equation formulates the metric.

$$R^2 = 1 - \frac{\sum_{i=0}^n (x_i - \hat{x}_i)^2}{\sum_{i=0}^n (x_i - \bar{x})^2} \quad (1)$$

Here n is the number of observations, x_i is the actual value for observation i , \hat{x}_i is the estimated value for observation i , and \bar{x} is the mean of x_i values.

The LR method focuses on minimizing R^2 . Thus, Spiess and Neumeier, for example, consider that the metric is not appropriate for evaluating nonlinear regression models [32]. Nevertheless, the metric is often used to compare models, e.g., [29]. In this study, we use the metric to evaluate the performance of the prediction models in predicting the test dataset and not the training dataset. The metric provides a “summary” of the errors of the predictions.

Akaike Information Criterion This metric estimates the information loss when approximating reality. The following equation formulates the metric [29].

$$AIC = n \times \log \left(\sum_{k=0}^n (x_i - \hat{x}_i)^2 / n \right) + 2(k + 2) \quad (2)$$

¹ <https://www.r-project.org/about.html>.

² <https://cran.r-project.org/web/packages/rpart/rpart>.

³ <https://cran.r-project.org/web/packages/nnet/nnet>.

⁴ We avoided the metric Mean of the Magnitude of the Relative Error (MMRE) as it was shown to be misleading [31].

Here n is the number of observations, x_i is the actual value for observation i , \hat{x}_i is the estimated value for observation i , and k is the number of variables.

A smaller AIC value indicates a better model.

Prediction at a Given Level This metric computes the percentage of prediction falling within a threshold h [33]. The following equation formulates the metric

$$PRED(h) = \frac{100}{n} \times \sum_{i=1}^n \begin{cases} 1 & \text{if } \frac{x_i - \hat{x}_i}{x_i} \leq h \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

Here n is the number of observations, x_i is the actual value for observation i , \hat{x}_i is the estimated value for observation i , and h is the threshold, e.g., 25 %.

The perfect value for the PRED metric is 100 %.

4.3 Variable Importance

This metric measures the relative contributions of the different predicting factors used by the regression method to the response variable. For statistical use, such metric could be, for example, the statistical significance while for business use, the metric could be the “impact on the prediction factor” on the (dependent) response variable. In this work, we use the variable importance metric employed in the RPART regression method.⁵ The metric measures the sum of the weighted reduction in the impurity method (e.g., the Shannon entropy and the variance of response variable) attributed to each variable [34, 35].⁶ It associates scores with each variable, which can be used to rank the variables based on their contribution to the model.

5 Methodology

Figure 2 depicts the process that we used in this study; a process quite similar to the one used by Bener et al. [36]. First, we define the goal of the data analytics activity, which is to develop a function for predicting the issue fix time using the data that SAP collects on its processes for fixing vulnerabilities in pre-release and post-release software. The following steps are: collect data that could help achieve the goal; prepare the data to be used to derive insights using statistical methods; explore the collected datasets to understand the used coding scheme, its content, and the relationships between the data attributes; develop prediction models for each of the collected datasets; compute metrics on the model; and try to find explanations

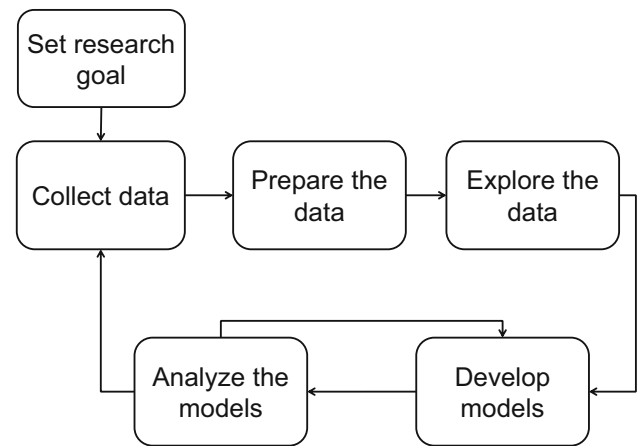


Fig. 2 Analysis method

and arguments for the results. The results of the models analysis were used to identify ways to improve the models. The improvements included the collection of new datasets for dependent information, e.g., about projects. We discuss next the individual steps in more details.

5.1 Data Collection

SAP maintains three datasets on fixing security vulnerabilities, which we refer to as our main *data sources*. In addition, it maintains a dataset about components, and a dataset about projects, which we call *support data*. Table 2 lists the different datasets we use. The datasets used in our study span over distinct time periods for each dataset (e.g., about 5 years).

The security testing process records data about fixing issues in two datasets. First, ABAP developers use SLINT for security code analysis. In Data source 1, the tool records data related to a set of attributes about each of the issues it discovers and the tasks performed on these issues. Table 3 lists these attributes.

Second, Java and JavaScript developers use Fortify⁷ and C++ developers use Coverity to analyze software for security issues. In Data source 2, these tools record data related to a set of attributes about each of the vulnerabilities they discover and the tasks performed on these vulnerabilities. Table 4 lists these attributes.

In Data source 3, the security response process maintains data about fixing issues discovered in released software. The data are collected and maintained through a Web form; it is not collected automatically as in the case of Data sources 1 and 2. The attributes of this Data source are listed in Table 5.

⁵ We use function `varImp()`.

⁶ There are other methods for ranking variables. We choose this method, which comes with RPART, because most the variables are categorical.

⁷ Since 2013, SAP uses Checkmarx for analyzing JavaScript. Thus, the use of Fortify by JavaScript developers declines since then.

Table 2 Datasets collected from SAP's tools

Dataset	Description
Data source 1	Vulnerabilities found in ABAP code
Data source 2	Vulnerabilities found in Java and C++ code
Data source 3	Security messages
Support data 1	Components
Support data 2	Projects
Extended data source 2	Extend data source 2 with information about the projects (support data 2)
Extended data source 3	Extend data source 3 with information about the components (support data 1)

Table 3 List of the attributes of ABAP issue fixing (Data source 1)

Attribute	Description
Date_found	Date on which the issue was found
Date_solved	Date on which the issue was closed
Vulnerability_name	Vulnerability types such as memory corruption and buffer overflow
Project_ID	Project identifier
Priority	The priority of fixing the vulnerability. Range: 1 to 4, with 1 highest, 4 lowest priority
Vulnerability_count	Number of issues of the same vulnerability found at once. This indicates that the issues might be related to the same problem

Table 4 List of the attributes of Java and C++ issue fixing (Data source 2)

Attribute	Description
Date_found	Date on which the issue was found
Date_solved	Date on which the issue was closed
Vulnerability_name	Vulnerability types such as memory corruption and buffer overflow
Scan_source	Tool that performed the scan, i.e., Coverity (for C++ code) or Fortify (for Java code)
Project_name	Project identifier
Folder_name	Indicates the required behavior of the developer toward the issue, e.g., must fix, fix one of the sets, optional, etc
Scan_status	Status of the issues, i.e., new, updated, removed, and reintroduced (i.e., removed but reopened). It allows to identify whether the issue is addressed or not, and is a false positive or not
Vulnerability_count	Number of issues of the same vulnerability found at once
Priority	The priority of fixing the vulnerability. Range: 1 to 4, with 1 highest and 4 lowest priority

Each issue can relate to a concrete component. Components are groups of applications that perform a common business goal. A system consists of a set of components. Table 6 lists the components attributes.

A software is developed in the context of a project. Table 7 lists the attribute of the projects dataset (Support data 2). We Extended data source 2 with project descriptions data; we joined Data source 2 and Support data 2. We also added three computed fields to the dataset:

1. FixtoRelease_period: The time elapsed from fixing the given issue to releasing the software.
2. Dev_period: The time elapsed from starting the development to closing the development of the software that contains the issue.

3. FoundtoRelease_period: The time elapsed from discovering the issue to the releasing of the software that contains the issue.

The number of records for each of the basic datasets ranges from thousands of records to hundred of thousands of records. We did not provide the exact numbers to avoid their misuse (in combination with potentially other public data) to derive statistics about vulnerabilities in SAP products, which would be outside the scope of this work.

5.2 Data Preparation

Using the collected data required us to prepare them for the model generation routines. The preparation activities

Table 5 List of the attributes for security messages (Data source 3)

Attribute	Description
CVSS_score	Common Vulnerability Scoring System (CVSS). The score indicates also the urgency of fixing the vulnerability
Processor	Identifier of development team/area and, thus, implicitly for the local instantiation of the S ² DL
Reporter	Identifier of the external researcher/company who reported the issue
Source	The source of the reported issue such as internal, security testing tool, customers
Vulnerability_type	Vulnerability type
Priority	Priority of the issue to be fixed: low, medium, or high
Component	Group of applications that perform a common business goal such as sales order or payroll

Table 6 List of the attributes for the components (support data 1)

Attribute	Description
PTU_area	The area of the component, e.g., CRM, IMS, ERP
Gr_component	Component group, i.e., semantic aggregation of components based on superordinate level
Language	The language(s) used to develop the component: ABAP, Java, ABAP and Java, or unknown
PPMS_product	The name of product that the component is part of, as stated in PPMS (Projects Management System)
Comp_owner	The component's development group
Product_owner	The product's development group

Table 7 List of the attributes for projects (support data 2)

Attribute	Description
Project_name	The name of the project (internal program name)
Prg_typ_id	Release related vs release unrelated (RR / UR)
Rel_type_id	Project type (standard, etc.)
Rel_typ_id	Release type ID (standard, pilot, etc.)
Delivery_mode_id	Mode of delivery to the customer. Values are on premise, on demand, on mobile, etc
Maintstrategy_id	Maintenance strategy. There is a codification for the strategies
Deploy_type	Deployment type. There is a codification for the deployment
D2t_date	Planned end of the test period. The period starts after the development closes
Devclose_date	Closing date of the development
P2d_date	Planned development starting date
P2r_date	Planned release date
Prg_lead_resp	Development team responsible for the project
Risk_expert	Identifier of risk expert (anonymized data)

required cleaning the data and transforming them as needed for processing.

Data Cleaning First, we identified the data columns where data are frequently missing. Missing values impact the results of the regression algorithms because these algorithms may incorrectly assume default values for the missing ones. We used plots such as the one in Fig. 3 to identify data columns that require attention.

Second, we developed a set of plots to check outliers—values that are far from the common range of the values of the attributes. We excluded data rows that include semantically wrong values, e.g., we removed records from

Data source 1 where the value of “Date_found” is 1 Dec. 0003.

Third, we excluded records related to issues that are not addressed yet; we cannot deduce issue fix time of such records.

Fourth, we excluded records that include invalid data. For example, the vulnerability type attribute of Data source 2 includes values such as “not assigned,” “?,” and “&novuln.” The records that have these values are excluded. There is no interpretation of prediction results that include these values.

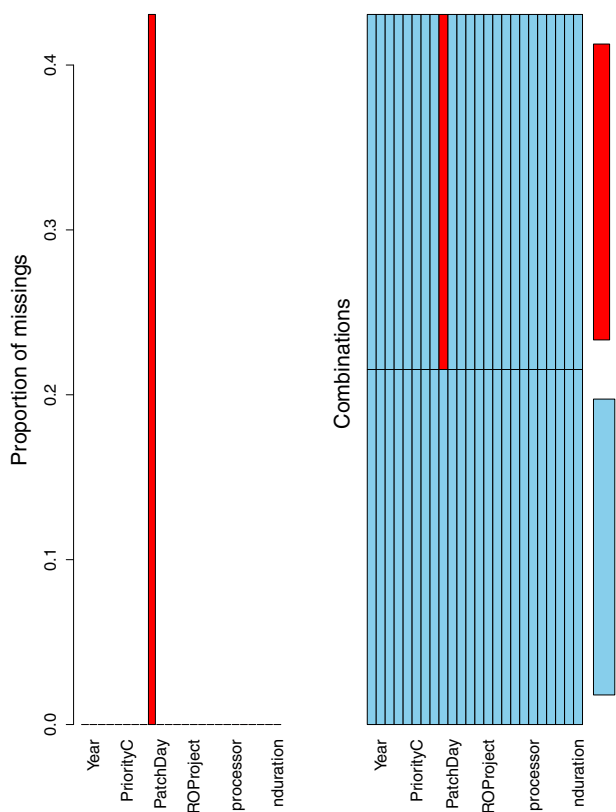


Fig. 3 Plot that visualizes missing data for Data source 3

Fifth, we excluded non-useful data attributes. These include, for example, the case where the attribute is derived from other attributes considered in the models.

Data Transformation First, we transformed the data of some columns from type text to appropriate types. For instance, we transformed the data of the CVSS column to numeric. Next, we computed new data columns from the source (original) data. For example, we computed the issue fix time from the issue closing date and issue discovery date or we performed some attributes' value transformations to obtain machine readable data for model generation. Some attributes contain detailed information that reduces the performance of the regression algorithms. We addressed this issue by developing a good level of data aggregation for the prediction algorithm. For example, the original dataset included 511 vulnerability types. We grouped the vulnerabilities types in vulnerability categories, which helps to derive better prediction models. Also, we aggregated the "component" variable to obtain "Gr_component" to include in our regression.

5.3 Data Exploration

We developed a set of plots and statistics about the frequencies of values for the factors and the relationship between the issue fix time and some of the prediction

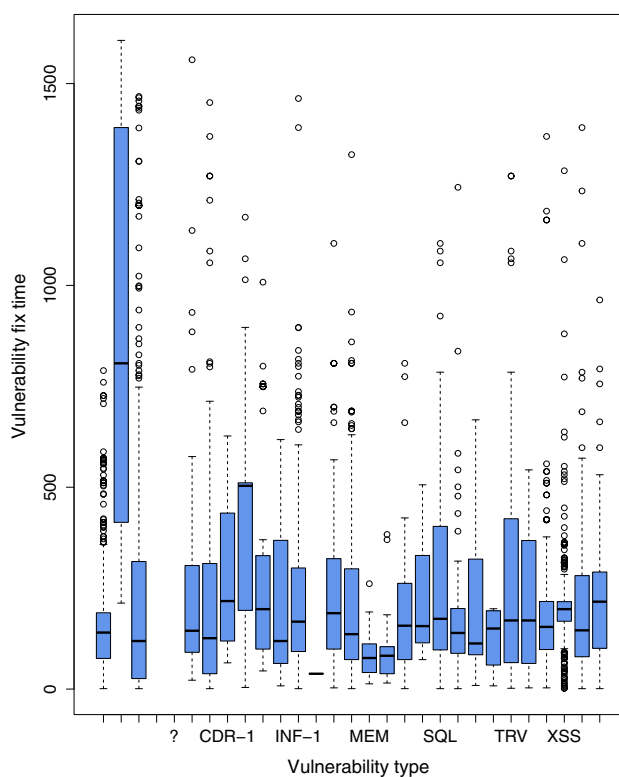


Fig. 4 Relationship between issue fix time (in days) and vulnerability types in the context of Data source 3. CDR-1, INF-1, MEM, SQL, TRV, XSS are internal codes for vulnerabilities types and code "?" indicates unknown or uncategorized type of reported vulnerabilities. (Some vulnerability types do not appear on the X axis to ensure clarity)

factors. For example, Fig. 4 shows the relationship between the issue fix time in days and vulnerability type. This gives us a first impression of the relations among the attributes of a given dataset. Also, Table 8 shows the coefficients of the linear regression (LR) of the issue fix time using the factor message source, that identifies the source of the reported issue. The table shows that the coefficients in this categorical factor indicate the different contributions of the factor on the issue fix time. The results indicate different impacting strengths of the different sources of security messages (e.g., external parties, customers, or the security department) on the issue fix time.

5.4 Models Development

We partitioned each prepared Data source into a training set that includes 80 % of the data and a test set (used to validate the developed model) that includes the remaining 20 %.⁸ We used the training set to develop the prediction

⁸ We used 80 % of the data for developing the prediction models and not 60 % of the data because the size of Extended data source 3 is limited: 380 records. (We wanted to use the same ratio for all datasets.)

Table 8 Coefficients of the linear regression of issue fix time to security message dataset (Data source 3)

Message source	Coefficient	<i>p</i> value
(Intercept)	249.17	<0.001
Code scan tool	-50.04	<0.001
Central security department	-38.05	<0.001
Customers	-60.68	<0.001
External research organizations	-102.78	<0.001
Internal development departments	-12.21	0.304
Test services	-124.74	<0.001
Validation services	-21.88	0.136

models, or fits, and the test set to assess the goodness of the generated models. The selection of the records for both sets is random.

Next, we performed three operations for each of the main data sources. First, we generated three prediction models using the training set, one using the linear regression method, one using the RPART method, and one using the NNR method. The three data sources have different data attributes and cannot be combined. Thus, we cannot use them together to develop a generic prediction model.

5.5 Models Analysis

We used the variable importance metric described in Sect. 4.3 to assess the impact of the different prediction factors on the issue fix time for each of the three data sources. The metric indicates that the factor “project name” is very important for Data source 2 and the factor “component” is very important for Data source 3. The results and their appropriateness were discussed with the security experts at SAP. We Extended data source 2 with Support data 2 (i.e., projects dataset) and we Extended data source 3 with Support data 1 (i.e., components dataset). Next, we performed the model development phase (Sect. 5.4) using the extended datasets. Then, we used each of the prediction models to predict the issue fix time for the test dataset and computed the performance metrics (see Sect. 4.2) for each model. We discuss the results in the next section.

6 Study Results

This section discusses the developed prediction models addressing issue fix time and their performance, the relative importance of the prediction factors, and the evolution of mean vulnerability fix time over time.

6.1 Issue Fix Time Prediction Models

This section aims to address the question: How well do the chosen models (LR, RPART, and NNR) predict the issue fix time from a set of given factors?

Most of the data that we use are not numeric; they are categorical variables, e.g., vulnerability types and component IDs. The number of categories in each of these variables can be very high. For instance, there are about 2300 components. This makes the prediction models large, e.g., in the order of a couple of hundred of nodes for the tree-based model and few thousands for the neural network model. This problem of large sets for the categorical factors limits the ability to generate accurate prediction models.

The regression algorithms cluster the elements of the categorical variables automatically; this clustering does not follow a given semantics, such as aggregation on superordinate component level, i.e., “Gr_component” in Support data 1. Because of this, it is impractical to plot the prediction models. Figure 5, for instance, shows a prediction model that we generated from Data source 3 by using the RPART method. The big number of (categorical) values for many of the prediction variables makes visualizing the generated models clearly difficult.

Interestingly, we observe that the component factor is built upon a set of distinct component classes (i.e., the first three digits indicate the superordinate component level, e.g., CRM for customer relationship management). An investigation of underlying reasons for such clustering might reveal, e.g., coherence with process-related factors.

6.2 Performance of Selected Regression Methods on the Prediction of Issue Fix Time

This subsection addresses the question: Which of the developed regression models gives the most accurate predictions? It reports and discusses the measurements of the performance metrics (introduced in Sect. 4.2) that we performed on the models that we generated for predicting the issue fix time. Table 9 summarizes the measurements that we obtained.

Coefficient of Determination Metric We observe that the LRs method outperforms the RPART and NNR methods for the five datasets. The metric values indicate that the prediction models generated using LR explain about half of the variation of the real values for Data source 1 and for Data source 2 and explains most of the variations for the remaining data sources. Indeed, the estimates of the model for the Extended data source 2 perfectly match the observed values. We note also that the residues metric values indicate that the prediction models generated using

84) vulnerabilitytype=,&OTHER,ACI-1,CDR-1,INF-1,MAC-1,MEM,XSS,XS
S-2 270 5063771.00 286.53700

168) Component=AP-RC-ANA-UI-XLS,BC-BSP,BC-CST-DP,BC-C
ST-IC,BC-CTS-SDM,BC-CTS-TMS,BC-DOC-HLP,BC-DOC-TTL,BC-I18,BC
-JAS-ADM-MON,BC-JAS-DPL,BC-SEC,BC-SEC-DIR,BC-SRV-ARL,BC-SR
V-FSI,BC-UPG-SLM,BC-UPG-TLS-TLJ,BC-WD-CMP-FPM,BC-XI-CON-AX
S,BC-XI-IBD,BC-XI-IBF,BI-BIP-AUT,BI-OD-STW,BI-RA-WBI,BW-BEX-OT-
MDX,CA-GTF-IC-BRO,CA-GTF-IC-SCR,CA-GTF-RCM,CRM-BF,CRM-BF-
SVY,CRM-CIC,CRM-IC-EMS,CRM-IC-FRW,CRM-IPS-BTX-APL,CRM-ISA,
CRM-ISA-AUC,CRM-ISE,CRM-LAM-BF,CRM-MD-PRO,CRM-MKT-DAM,C
RM-MKT-MPL,CRM-MSA,FS-CM,FS-SR,IS-A-DP,IS-U-CS-ISS,LO-AB-BS
P,LO-GT,MFG-ME,MOB-APP-EMR-AND,PA-GE,PLM-PPM-PDN,PLM-WUI
-RCP,PSM-GPR-SN,SBO-INT-B1ISN,SCM-EWM-RF,XAP-IC-IDM,XX-PRO
J-CDP-TEST-296 119 1015233.00 205.82350 *

169) Component=AP-CFG,AP-LM-MON-HC,AP-LM-SUP,AP-RC-
ANA-RT-MDA,AP-RC-RSP,AP-RC-UIF-RT,AP-SDM-EXC,BC-CCM-MON-
OS,BC-CCM-SLD-JAV,BC-CST,BC-CUS-TOL-CST,BC-DB-ORA-INS,BC-D
OC-TER,BC-ESI-WS-ABA,BC-ESI-WS-JAV-RT,BC-FES-BUS-RUN,BC-JA
S-ADM-ADM,BC-JAS-COR,BC-JAS-SEC-UME,BC-MID-RFC,BC-SEC-SA
L,BC-SRV-COM,BC-SRV-COM-FTP,BC-SRV-KPR-CS,BC-SRV-MCM,BC-
SRV-SSF,BC-WD-ABA,BC-WD-

Fig. 5 Part of the prediction model generated from Data source 3 using RPART method. (It shows that there is a long list of component IDs (among a set of about 2300 components) for the selection of nodes 168 and 169 while also the parent node uses a set of vulnerability types)

the NNR perform worse than the null model, that is, taking the average of the values.

AIC Metric We observe that the LR method outperforms the RPART and NNR methods for two datasets and that the RPART method outperforms the LR and NNR methods for the remaining three datasets. Thus, this metric gives mixed results with respect to performance of the three regression methods.

PRED Metric We observe that the LR method outperforms the RPART and NNR methods for two datasets, the RPART method outperforms the LR and NNR methods for

one dataset, and the NNR method outperforms the RPART and LR methods for two datasets. This gives mixed results with respect to performance of the three regression methods. However, the NNR performance improves when the dataset is extended with related data. For instance, the PRED value increased from 0.73 % in the case of Data source 3–65.05 % for the case of the Extended data source 3. We acknowledge that the PRED value improved also for the LR method for the case of Data source 2 and Extended data source 2. However, the number of records (*N*) for the Extended data source 2 is low (*N* = 380); the result should be taken with caution.

Different regression methods have shown conflicting performance measurements toward the problem of effort estimation. For example, Gray and MacDonell [24] compared a set of regression approaches using MMRE and PRED metrics. The methods have shown conflicting results; their rank change based on the used performance metrics. For example, they found that based on the MMRE metric, LR outperforms NNR and based on the PRED metric, NNR outperforms LR. This finding was confirmed by Wen et al. [37] who analyzed the performance of several other regression methods. The regression methods have different strengths and weaknesses. Most importantly, they perform differently in the presence of small datasets, outliers, categorical factors, and missing values. We found in this study that it is not possible to claim that a regression method is better than the other in the context of predicting the issue fix time. This result supports the findings of Gray and MacDonnell [24] and of Wen et al. [37].

Table 9 Measurement of the performance metrics of the prediction models

Dataset	LR	RPART	NN	Best method
<i>Residuals metric</i>				
Data source 1	0.526	0.498	-1.252	LR (0.526)
Data source 2	0.461	0.44	-0.294	LR (0.461)
Extended data source 2	1	0.956	-0.587	LR (1)
Data source 3	0.944	0.6585	1.944	LR (0.944)
Extended data source 3	0.909	0.701	1.97	LR (0.909)
<i>AIC metric</i>				
Data source 1	122,465	123,157	141,462	LR (122,465)
Data source 2	334,565	335,936	365,665	LR (334,565)
Extended data source 2	-4877	463	793	RPART (463)
Data source 3	6632	6507	6958	RPART (6507)
Extended data source 3	6581	6421	7057	RPART (6421)
<i>PRED metric</i>				
Data source 1	31.81 %	31.74 %	0.156 %	LR (31.81 %)
Data source 2	14.93 %	13.96 %	33.81 %	NN (33.81 %)
Extended data source 2	100 %	30.32 %	39.40 %	LR (100 %)
Data source 3	33.98 %	34.71 %	0.73 %	RPART (34.71 %)
Extended data source 3	35.41 %	34.52 %	65.05 %	NN (65.05 %)

Table 10 Importance factors for the main datasets

Factor	Data source 1		Data source 2		Data source 3	
	Rank	Factor	Metric	Factor	Metric	Factor
1	Project_ID	0.408	Scan_status	1.231	Component	3.653
2	Vulnerability_name	0.179	Project_name	0.956	Processor	3.352
3	Vulnerability_count	0.097	Vulnerability_name	0.841	Reporter	1.813
4	Priority	0.042	Folder_name	0.522	Vulnerability_type	0.879
5			Scan_source	0.238	Source	0.443
6			Priority	0.358	CVSS_score	0.230
7			Vulnerability_count	0.078		

6.3 Relative Importance of the Factors Contributing to Issue Fix Time

This section aims to address the question: What are the main factors that impact the issue fix time? To answer this question, we computed the variable importance metric (see Sect. 4.3) for each factors [38]. Given that the factors used in the datasets are different, we present, and shortly discuss, each dataset separately, and we used the variable importance metric discussed in Subsection 4.3 to compute the importance of each factor.

Data source 1 Table 10 reports the importance of the factors used in Data source 1 on issue fix time. The most important factor in this dataset is “Project_ID,” followed by “Vulnerability_name.” This implies that there is a major contribution of the project characteristics to issue fix time. Unfortunately, there were no additional metadata available on the projects that could have been joined with Data source 1 to allow us to further investigate aspects of projects that impact the fixing time.

Data source 2 Table 10 reports the importance of the factors used in Data source 2 on issue fix time. The most important factor in this dataset is “Scan_status.” This shows that depending on whether the issue is false positive or not impacts the issue fix time.⁹ The second ranking factor is “Project_name,” followed by “Vulnerability_name.” This results support the observation we had with Data source 1. We observe also that the factor “Scan_source,” which indicates the static code analysis tool used to discover the vulnerabilities (i.e., Fortify or Coverity) is ranked at place 5.

Extended data source 2 We Extended data source 2 with data that describe the projects and computed three additional variables: the time elapsed between fixing the vulnerability and releasing the software, called FixtoRelease_period; the development period, called Dev_period; and the time elapsed between discovering the vulnerability and releasing the software, called FoundtoRelease_period.

⁹ As indicated before, issues marked as, e.g., new and updated are not considered in the models; they are for issues that are not addressed yet.

Table 11 reports the importance of the factors used in the Extended data source on issue fix time. We observe that the most important factor is FixRelease_period while a related factor, FoundtoRelease_period has less importance (rank 6). The other main important factors include the development period (Dev_period), the project (Project_name), the risk expert (Risk_expert), the project development team (Prg_lead_resp), and the vulnerability name (Vulnerability_name). We observe that vulnerability name is ranked at the seventh position.

Data source 3 Table 10 reports the importance of the factors used in Data source 3 on the issue fix time. The most important factor in this dataset is the software component that needs to be changed (Component) followed by the development team who addresses the issue (Processor). We observe that the vulnerability name (Vulnerability_type) has a moderate importance, ranked 4th, while the CVSS (CVSS_score) is ranked on the sixth position.

Extended data source 3 We Extended data source 3 with data that describe the components. Table 11 reports the importance of the factors used in the Extended data source 3 on the issue fix time. The most important factors in this extended dataset is the component (Component), followed by the development team (Processor), the development team responsible for the component (Dev_comp_owner), the reporter of the vulnerability (Reporter), and a set of other factors. We observe that the vulnerability name (Vulnerability_type) has a moderate importance, ranked eighth, and the CVSS score (CVSS_score) decreased considerably to be 0.

6.4 Evolution of the Issue Fix Time

This section aims to address the question: Is the company improving in fixing security issues? The tendency of the issue fix time could be used as “indicator” of such improvement. For instance, increasing time indicates deteriorating capabilities and decreasing time indicates improving capabilities. The information should not be used as an evidence but as indicator of a fact that requires further investigation.

Table 11 Importance factors for the extended datasets

Rank	Extended data source 2		Extended data source 3	
	Factor	Metric	Factor	Metric
1	FixtoRelease_period	0.946	Component	2.434
2	Dev_period	0.634	Processor	2.363
3	Project_name	0.634	Dev_comp_owner	1.380
4	Risk_expert	0.634	Reporter	1.207
5	Prg_lead_resp	0.672	Dev_product_owner	0.633
6	FoundtoRelease_period	0.433	PPMS_product	0.407
7	Vulnerability_name	0.429	PTU_area	0.407
8	Priority	0.045	Vulnerability_type	0.058
9	Vulnerability_count	0.045		
10	Folder_name	0.045		

We modeled the evolution of the mean issue fix time for the resolving (closing) issue month¹⁰ for the three data sources using the LR, which shows the trend of the response variable over time. Figure 6 depicts, respectively, the mean issue fix time for (a) Data source 1 (pre-release ABAP-based code), (b) Data source 2 (pre-release Java, C++, and JavaScript-based code), and (c) Data source 3 (post-release security issues).

The figure indicates a fluctuation of the mean issue fix time but with an increasing trend. This trend indicates a deteriorating performance with respect to fixing security issues. A close look at the figure shows that there is a recent reverse in the tendency, which indicates a response to specific events such as dedicated quality releases or the development of new flag ship products. So-called quality releases are releases that focus on improving the product quality instead of focusing on new features. To ensure a high level of product quality and security of SAP products, top-level managements plans, once in a while, for such quality releases. Also the development of new flagship products that change the development focus of a large fraction of all developers at SAP can have an influence. Such a shift might result in significant code simplifications of the underlying frameworks.

Figure 6 shows that the increasing global trend applies to pre-release as well as post-release issues. We believe that this indicates that the causes of the increase in the mean issue fix time apply to both cases. We again see that the management actions impacted both cases.

Berner et al. [36] advise that models are sensitive to time. This work supports the claim because it shows that the issue fix time is sensitive to the month of closing the issue.

¹⁰ Compute the mean issue fix time for the vulnerabilities resolved (addressed) in the specified month.

7 Analysis of the Impacting Factors

We observe from Data source 1 and Data source 2 that projects (represented by, e.g., Project_ID, and Project_name data attributes) have major contributions to issue fix time for the case of pre-release issue fixing. The extension of Data source 2 with project-related data confirmed our observation: the most impacting factors of pre-release issues on the time to fix are project characteristics. Among these characteristics, we find the time between issue fixing and software release, project development duration, and the development team (data attribute Project_name). We believe that the factor time between issue fixing and software release indicates that developers tend to fix vulnerabilities as the release date becomes close. This is not surprising, since they must address all open issues before the software can pass the quality gate to be prepared for release. We expect that the factor project development duration is related to, e.g., the used development models, and the component-related characteristics. Further data analysis shall provide insights about the impact of the factors related to the project development duration (Dev_period), such as component complexity. For instance, updating smaller component could be easy and be performed in short development cycles while updating complex components requires long development cycles. In addition, we believe that the factor development team (Processor) indicates the level of expertise of the developers and the smoothness of communication and collaboration among the team. Nevertheless, it is interesting to observe that the influence of vulnerability type decreases when project-related factors are included.

There are two additional dominant factors for the issue fix time, based on the analysis of Data source 2: scan status and folder name. We believe that the factor scan status indicates that the developers address issues based on their perception of whether the given issue is a false positive or

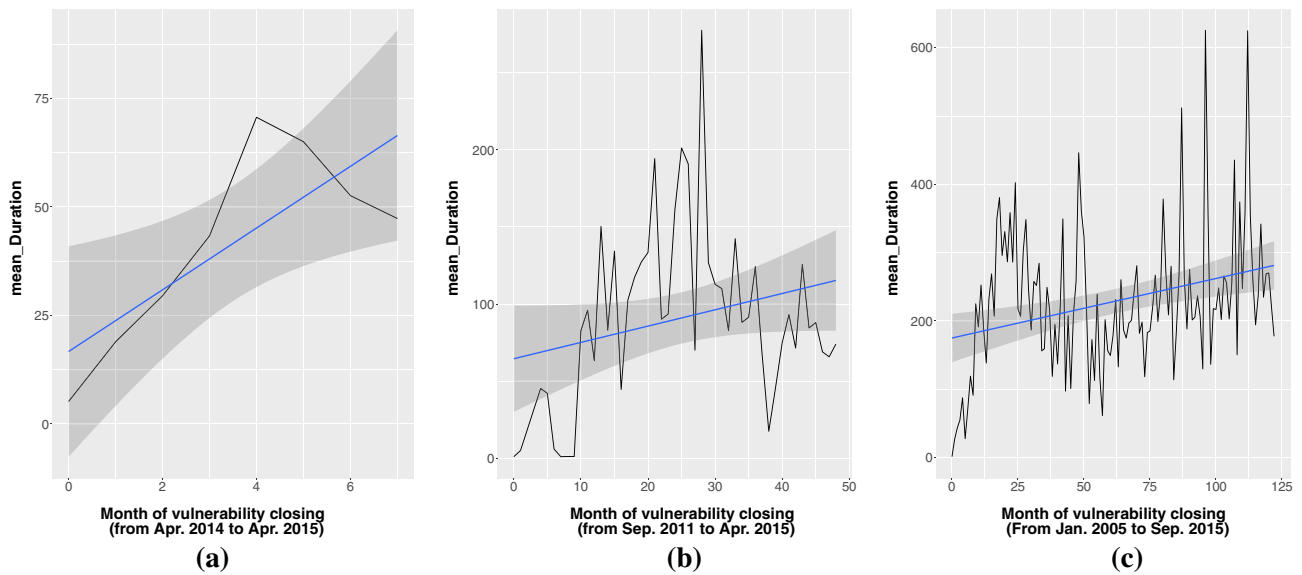


Fig. 6 Trend of the issue fix effort by closing month. The x axis indicates the number of months elapsed since the start date of the data. The y axis indicates the mean issue fix time in number of days. **a** Data source 1. **b** Data source 2. **c** Data source 3

not and whether it is easy to address or not. For example, they may close false positive issues that are easy to analyze and postpone addressing issues that are difficult to analyze and/or fix to, e.g., when the time for the quality gates becomes close. We also expect that the factor “folder name” indicates that the developers behave differently toward issues flagged must fix, fix one of the sets, or optional to fix.

The analysis of Data source 2 reveals that the security scan tools (represented by the data attribute *Scan_source*) are not a leading factor of issue fix time. It is possible that the developers rely on their expertise in analyzing security issues and not on the tool features as they get experts in addressing security issues. Further analysis may explain the finding better.

The results obtained from the analysis of Data source 3 (and its extension) suggest that the software structure and development team characteristics are the dominant factors that impact the issue fix time for the case of post-release issue fixing. (Note that issues for post-release are not related to projects but to released components.). The analysis results show that the component factor is among the most impacting factors on the issue fix time, which indicates the impact of software structure characteristics. Unfortunately, we do not have, at this moment, data that describe the components, such as the component’s complexity, which could be used to get detailed insights about these characteristics.

The results obtained from the analysis of Data source 3 shows the dominance of the impact of processor and reporter on the issue fix time and thus the importance of the human-related factors. The importance of the reporter

factor is aligned with the results of Hooimeijer and Weimer [39], who found a correlation between a bug reporter’s reputation and triaging time: We confirm the importance of the reporting source on vulnerability fix time. Unlike Zhang et al. [18] finding, we found that severity, represented by the “*CVSS_score*” in our study, is not a leading factor.

The qualitative study [4], which was based on expert interviews at SAP, revealed several factors that impact the issue/vulnerability fix time, such as communication and collaboration issues, experience and knowledge of the involved developers and security coordinators, and technology diversification. The results of this study confirm the impact of some of these factors and show their importance. For example, the category technology diversification included factors related to technologies and libraries supported by the components associated with the given vulnerability. The impact of the component, found in our current study, reflects these underlying factors. Unfortunately, it was only possible to relate components’ attributes to security messages, i.e., post-release issues, not to pre-release issues to further investigate the components’ impact on these. The impact of the development groups reflects the importance of the experience and vulnerability- and software-related knowledge of the teams as well as the importance of the smoothness of communication and collaboration between the involved stakeholders.

At SAP, the project development teams work independently; e.g., they choose their own development model and tools, as long as they confirm to the corporate requirements, such as the global security policies. Further investigation of component-, project-, human-, and process-

related characteristics of the development teams might reveal more insights on the underlying factors that impact the issue fix time. Such investigation may reveal why certain products/teams are more efficient than others. Reasons may, for example, include the local setup of the communication structures, the used development model (i.e., SCRUM, DevOps.) and the security awareness of teams. Another potential factor to check its impact is the number of people involved in fixing the given issue. This factor was found to impact the fix time of bugs [15, 40]. Controlling these factors allows to control the issue fix time and thus the cost of addressing security issues.

A question worth also investigating is: Are the factors that impact the time for addressing pre-release and post-release issues similar? We argue in Sect. 3 that the processes for fixing pre- and post-release issues are different, which shall impact the issue fix time for both cases. Nevertheless, acquiring evidence to answer this question requires using the same data attributes for both cases, which may not be possible, at the moment, with data collected at SAP.

8 Study Validity and Impacts

This section discusses the impacts of the finding and the limitations of the study.

8.1 Impacts of the Findings

This study showed that the models generated using the LR, RPART, and NNR methods have conflicting accuracy measurements in predicting the issue fix time. This implies that the conflict in the performance measurements in estimating software development effort, e.g., in [37], applies to security issues. We infer from this result that there is no better regression method, from the analyzed ones, when it comes to predicting security issue time. We believe that more work needs to be done to develop regression methods appropriate for predicting issue fix time.

The second main finding of this study is that vulnerability type is not the dominant impacting factors for issue fix time. Instead, the dominant factors are the component in which the vulnerability resides, the project related to the issue, the development groups that address the issue, and the closeness of the software release date, which are process-related information. This result implies that we should focus on the impact of software structure, developers' expertise and knowledge, and secure software development process when investigating ways to reduce the cost of fixing issues.

The third main finding is that the monthly mean issue fix time changes with the issue closing month. We can infer

from this result that the prediction models are time sensitive; that is, they depend on the data collection period. This supports Berner et al. advice to consider recently collected data when developing prediction models [36]. We infer, though, that prediction models are not sufficient for modeling issue fix time since they provide a static view. We believe that prediction methods should be extended to consider time evolution, that is, combine prediction and forecasting.

Finally, SAP can use the models to implement a continuous improvement of its secure software development process and to measure the impact of individual improvements. Other companies can use similar models and mechanism to realize a learning organization.

8.2 Limitations

There is a consensus among the community that there are many "random" factors involved in software development that may impact the results of data analytics experiments [36]. This aligned with Menzies et al.'s [19] findings about the necessity to be careful about generalization of results related to effort estimations in a global context.

The data analysis described in this report suffers from the two common threats to validity that apply to effort estimation [20]. First, the conclusions are based on the data that SAP collects about fixing vulnerabilities in its software. Changes to the data collection processes, such as changes to the attributes of the collected data, could impact the predictions and the viability of producing predictions in the first place. Second, the conclusions of this study are based on the regression methods we used, i.e., LR, RPART, and NNR. There are many other single and ensemble regression methods that we did not experiment with. We note that performance issues due to the size of the datasets inhibit us from using random forest [23] and boosting [23], two ensemble regression methods.

In addition, the data are collected over 5 years. During that time SAP refined and enhanced its secure software development processes. This could bias our results. The identification of major process changes along with the times of the changes and a partitioning of the data accordingly might reduce such bias and reveal measurable insights about impacts of process changes on issue fix time.

On the positive side, the conclusions are not biased by the limited data size and the subjectivity in the responses. First the number of records of each of the dataset was high enough to derive meaningful statistics. Second, the data are generated automatically and do not include subjective opinions, except the CVSS score of Data source 3. This score is generated based on issue-related information that is assessed by the security coordinator responsible for the issue.

Our findings are based on particular datasets of SAP and might mirror only the particularities of time to fix issues for this organization. However, SAP has a diversified software portfolio, the development teams are highly independent in using development processes and tools (as long as they follow generic guidelines such as complying with corporate security requirements), teams are located in different countries, and software are developed using several programming languages (e.g., ABAP, C++, and Java). These characteristics encourage us to believe that the findings apply to industrial companies in general and therefore contribute to the discussion about predicting the issue fix time.

Vulnerabilities, such as SQL injection, cross-site scripting (XSS), buffer overflow, and directory path traversal are commonly identified using the same techniques, such as taint analysis [41], but by applying different patterns. This may explain why vulnerability type is not a dominant factor of issue fix time—because the time to analyze many of the vulnerability types using the techniques is the same. These techniques may be used to detect other defect types besides vulnerabilities, but not all (or most) defects. Thus, the fact that vulnerability type is not a dominant factor of issue fix time does not imply that the study results apply to defects in general.

9 Lessons Learned

Data analytics methods are helpful tools to make generalizations about past experience [36]. These generalizations require considering the context of the data being used. In our study, we learned few lessons in this regard.

Anonymization Companies prefer to provide anonymized data for data analytics experiments and keep the anonymization map to trace the results to the appropriate semantics. There is a believe that the analyst would develop models and the data expert (from the company) would interpret them using the anonymization map. We initially applied the technique and we found that it prevents the analyst from even cleaning the data correctly. We worked closely with the owner of the data to understand them, interpret the results, and correct or improve the models. The better the data analyst understands the data, the more they are able to model them.

Prediction Using Time Series Data We initially sliced the data sequentially into folds (sliced them based on their order in the dataset) and used the cross-validation method in the regression.¹¹ We found that the performance metrics of the generated prediction models deviate considerably.

¹¹ This slicing method allows for easily splitting all the data among the folds.

Table 12 Predicted values for automatically clustered component factor and Gr_component

Error metric	LR		RPART		NNR	
	AC	MC	AC	MC	AC	MC
RSQ	0.98	0.76	0.80	0.7045	2.02	1.92
AIC	6586	6187	6461	6139	7033	6733
PRED	33.88	34.12	33.48	32.6915	0.48	0.67

AC is for automatic clustering of components and MC is for manual clustering of components

To explore this further, we developed the tendency of the mean issue fix time shown in Fig. 6. The figure shows a fluctuation of the issue fix time over time. This leads to believe that the prediction models are of temporal relevance as claimed by Berner et al. [36]. The lesson warns to check whether the data are time series or not when using cross-validation with sequential slicing of the data in the regression.

A more generic lesson that we learned concerns *attribute values clustering*. We found in this study an insignificant small difference in the performance of the prediction models that automatically cluster components and the ones that use semantically clustered components instead. The latter aggregates components based on a semantic based on superordinate level, i.e., Gr_component. Manual investigation is necessary to infer the component characteristics that the algorithms silently used in the clustering. Table 12, for example, shows that the performance of the prediction models using the automated clustering and using the manual clustering are similar. This implies that manual clustering does not provide additional information.

10 Conclusions

We developed in this study prediction models for issue fix time using data that SAP, one of the largest software vendors worldwide, and the largest in Germany, collects about fixing security issues in the software development phase and also after release. The study concludes that none of the regression methods that we used (linear regression (LR), Recursive PARTitioning (RPART), Neural Network Regression (NNR)) outperforms the others in the context of predicting issue fix time. Second, it shows that vulnerability type does not have a strong impact on the issue fix time. In contrast, the development groups involved in processing the issue, the component, the project, and the closeness of the release date have strong impact on the issue fix time.

We also investigated in this study the evolution of the mean issue fix time as time progresses. We found that the issue fix time fluctuates over time. We suggest that better models for predicting issue fix time should consider the temporal aspect of the prediction models; they shall combine both prediction technique and forecasting techniques.

Acknowledgments This work was supported by SAP SE, the BMBF within EC SPRIDE, the Scientific and Economic Excellence in Hesse (LOEWE), and a Fraunhofer Attract grant.

Open Access This article is distributed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made.

References

- McGraw, G.: Software security: building security. In: Addison-Wesley software security series. Pearson Education Inc, Boston (2006)
- Bachmann R, Brucker AD (2014) Developing secure software: a holistic approach to security testing. *Datenschutz und Datensicherheit (DuD)* 38:257–261
- Howard M, Lipner S (2006) The security development lifecycle: SDL—a process for developing demonstrably more secure software. Microsoft Press
- ben Othmane L, Chehraz G, Bodden E, Tsalovski P, Brucker A, Miseldine P (2015) Factors impacting the effort required to fix security vulnerabilities. In: Proceedings of information security conference (ISC 2015), Trondheim, Norway, pp 102–119
- Zimmermann T, Nagappan N, Williams L (2010) Searching for a needle in a haystack: predicting security vulnerabilities for windows vista. In: Proceedings of the 2010 third international conference on software testing, verification and validation, Washington, DC, pp 421–428
- Shin Y, Williams L (2013) Can traditional fault prediction models be used for vulnerability prediction? *Empir Softw Eng* 18:25–59
- Morrison P, Herzig K, Murphy B, Williams L (2015) Challenges with applying vulnerability prediction models. In: Proceedings of the 2015 symposium and bootcamp on the science of security, pp 4:1–4:9
- Keller H, Krüger S (2007) ABAP objects. SAP Press
- Chehraz G, Schmitz C, Hinz O (2015) QUANTSEC—ein modell zur nutzenquantifizierung von it-sicherheitsmaßnahmen. In: Smart enterprise engineering: 12. Internationale Tagung Wirtschaftsinformatik, WI 2015, Osnabrück, Germany, March 4–6, 2015, pp 1131–1145
- Cornell D (2012) Remediation statistics: what does fixing application vulnerabilities cost? In: RSAConference, San Francisco, CA
- Zeng H, Rine D (2004) Estimation of software defects fix effort using neural networks. In: Proceedings of the 28th annual international computer software and applications conference (COMPSAC 2004), vol 2, Hong Kong, China, pp 20–21
- Weiss C, Premraj R, Zimmermann T, Zeller A (2007) How long will it take to fix this bug? In: Proceedings of the fourth international workshop on mining software repositories. MSR '07, Washington, DC, p 1
- Panjer LD (2007) Predicting eclipse bug lifetimes. In: Proceedings of the fourth international workshop on mining software repositories. MSR '07, Washington, DC, IEEE Computer Society, p 29
- Bhattacharya P, Neamtii I (2011) Bug-fix time prediction models: can we do better? In: Proceedings of the 8th working conference on mining software repositories. MSR '11, ACM, New York, NY, pp 207–210
- Giger E, Pinzger M, Gall H (2010) Predicting the fix time of bugs. In: Proceedings of the 2nd international workshop on recommendation systems for software engineering. RSSE '10, ACM, New York, NY, pp 52–56
- Hamill M, Goseva-Popstojanova K (2014) Software faults fixing effort: analysis and prediction. Technical Report 20150001332, NASA Goddard Space Flight Center, Greenbelt, MD USA
- Hewett R, Kijisanayothin P (2009) On modeling software defect repair time. *Empir Softw Eng* 14:165–186
- Zhang F, Khomh F, Zou Y, Hassan A (2012) An empirical study on factors impacting bug fixing time. In: 19th Working conference on reverse engineering (WCRE), Kingston, Canada, pp 225–234
- Menzies T, Butcher A, Marcus A, Zimmermann T, Cok D (2011) Local versus global models for effort estimation and defect prediction. In: Proceedings of the 2011 26th IEEE/ACM international conference on automated software engineering. ASE '11, Washington, DC, pp 343–351
- Menzies T, Greenwald J, Frank A (2007) Data mining static code attributes to learn defect predictors. *IEEE Trans Softw Eng* 33:2–13
- Shin Y, Meneely A, Williams L, Osborne J (2011) Evaluating complexity, code churn, and developer activity metrics as indicators of software vulnerabilities. *IEEE Trans Softw Eng* 37:772–787
- Brucker AD, Sodan U (2014) Deploying static application security testing on a large scale. In: GI Sicherheit 2014, vol 228 of lecture notes in informatics, pp 91–101
- James G, Witten D, Hastie T, Tibshirani R (2013) An introduction to statistical learning with applications in R. Springer, New York
- Gray AR, MacDonell SG (1997) A comparison of techniques for developing predictive models of software metrics. *Inf Softw Technol* 39:425–437
- Hastie T, Tibshirani R, Friedman J (2013) The elements of statistical learning, 2nd edn. Springer, Berlin
- Menzies T (2013) Data mining: a tutorial. In: Robillard MP, Maalej W, Walker RJ, Zimmermann T (eds) Recommendation systems in software engineering. Springer, Berlin, pp 39–75
- Breiman L, Friedman J, Stone CJ, Olshen R (1984) Classification and regression trees. Chapman and Hall/CRC, Belmont
- Specht DF (1991) A general regression neural network. *IEEE Trans Neural Netw* 2:568–576
- Hyndman R, Athanasopoulos G (2014) Forecasting: principles and practice. Otexts
- Menzies EKT, Mendes E (2015) Transfer learning in effort estimation, empirical software engineering. *Empir Softw Eng* 20:813–843
- Foss T, Stensrud E, Kitchenham B, Myrteit I (2003) A simulation study of the model evaluation criterion mmre. *IEEE Trans Softw Eng* 29:985–995
- Spieß ANN, Neumeyer N (2010) An evaluation of R2 as an inadequate measure for nonlinear models in pharmacological and biochemical research: a Monte Carlo approach. *BMC Pharmacol* 10:6

33. Kocaguneli E, Menzies T, Keung J (2012) On the value of ensemble effort estimation. *IEEE Trans Softw Eng* 38:1403–1416
34. Louppe G, Wehenkel L, Sutera A, Geurts P (2013) Understanding variable importances in forests of randomized trees. In: Burges C, Bottou L, Welling M, Ghahramani Z, Weinberger K (eds) *Advances in neural information processing systems*, vol 26, pp 431–439
35. Eisenhardt KM (1989) Building theories from case study research. *Acad Manag Rev* 14:532–550
36. Bener A, Misirli A, Caglayan B, Kocaguneli E, Calikli G (2015) Lessons Learned from software analytics in practice. In: *The art and science of analyzing software data*, 1st edn. Elsevier, Waltham, pp 453–489
37. Wen J, Li S, Lin Z, Hu Y, Huang C (2012) Systematic literature review of machine learning based software development effort estimation models. *Inf Softw Technol* 54:41–59
38. Therneau TM, Atkinson EJ (2011) An introduction to recursive partitioning using the rpart routines. Technical Report 61, Mayo Foundation for Medical Education and Research; Mayo Clinic; and Regents of the University of Minnesota, Minneapolis, USA
39. Hooimeijer P, Weimer W (2007) Modeling bug report quality. In: *Proceedings of the twenty-second IEEE/ACM international conference on automated software engineering. ASE '07*, ACM, New York, NY, pp 34–43
40. Guo PJ, Zimmermann T, Nagappan N, Murphy B (2011) “not my bug!” and other reasons for software bug report reassignments. In: *Proceedings of the ACM 2011 conference on computer supported cooperative work. CSCW '11*, ACM, New York, NY, pp 395–404
41. Chess B, West J (2007) *Secure programming with static analysis*, 1st edn. Addison-Wesley, Reading