

Detecting Java Software Similarities by using Different Clustering Techniques

Andrea Capiluppi^{a,*}, Davide Di Ruscio, Juri Di Rocco, Phuong T. Nguyen^b,
Nemitari Ajenka^c

^a*Dept. of Computer Science, Brunel University London, UK*

^b*Dept. of Information Engineering, Computer Science and Mathematics, University of
L'Aquila, Italy*

^c*Dept. of Computer Science, Edge Hill University, UK*

Abstract

Background – Research on empirical software engineering has increasingly been conducted by analysing and measuring vast amounts of software systems. Hundreds, thousands and even millions of systems have been (and are) considered by researchers, and often within the same study, in order to test theories, demonstrate approaches or run prediction models. A much less investigated aspect is whether the collected metrics might be context-specific, or whether systems should be better analysed in clusters.

Objective – The objectives of this study are (i) to define a set of clustering techniques that might be used to group similar software systems, and (ii) to evaluate whether a suite of well-known object-oriented metrics is context-specific, and its values differ along the defined clusters.

Method – We group software systems based on three different clustering techniques, and we collect the values of the metrics suite in each cluster. We then test whether clusters are statistically different between each other, using the Kolmogorov-Smirnov (KS) hypothesis testing.

Results – Our results show that, for two of the used techniques, the KS null hypothesis (e.g., the clusters come from the same population) is rejected for

*Corresponding author

Email addresses: `andrea.capiluppi@brunel.ac.uk` (Andrea Capiluppi),
`{davide.diruscio,juri.dirocco,phuong.nguyen}@univaq.it` (Davide Di Ruscio, Juri Di Rocco, Phuong T. Nguyen), `nemitari.ajienka@edgehill.ac.uk` (Nemitari Ajenka)

most of the metrics chosen: the clusters that we extracted, based on application domains, show statistically different structural properties.

Conclusions – The implications for researchers can be profound: metrics and their interpretation might be more sensitive to context than acknowledged so far, and application domains represent a promising filter to cluster similar systems.

Keywords: FOSS (Free and open-source software), Application Domains, Latent Dirichlet Allocation, Machine Learning, Expert Opinions, OO (object-oriented)

1. Introduction

Research on empirical software engineering has increasingly used data made available in online repositories or collective efforts. The latest trends for researchers is to gather “as much data as possible” to (i) prevent bias in the representation of a small sample, (ii) work with a sample as close as the population itself, and (iii) showcase the performance of existing or new tools in treating vast amount of data.

Considering the MSR¹ series of events as an example, its researchers have constantly grown the number of systems analysed in their papers. During its 10 2017 edition, for instance, the joint set of papers of the main track (i.e., 64 papers overall) collected and analysed altogether over 3 million software systems. A 10-year trend with the number of software systems jointly analysed by the MSR papers is shown in Figure 1. One of the papers alone amassed some 900K systems as its case studies [1]. Figure 2 shows the average number of analysed FOSS projects per year while Figure 3 shows the median number of analysed FOSS projects per year from 2008 to 2017 for the MSR series of events. The highest average is observed in the 2016 edition where on average, 93,553 FOSS projects were analysed. However, the highest median value is the 2017 edition,

¹<http://www.msrconf.org/>

with 532 projects, following a rise to 33 in 2016. The largest study containing
 20 FOSS projects was registered in 2016, with 3,182,590 analysed projects. This
 study accounts for around 60% of the total and the massive growth as shown in
 Figure 1.

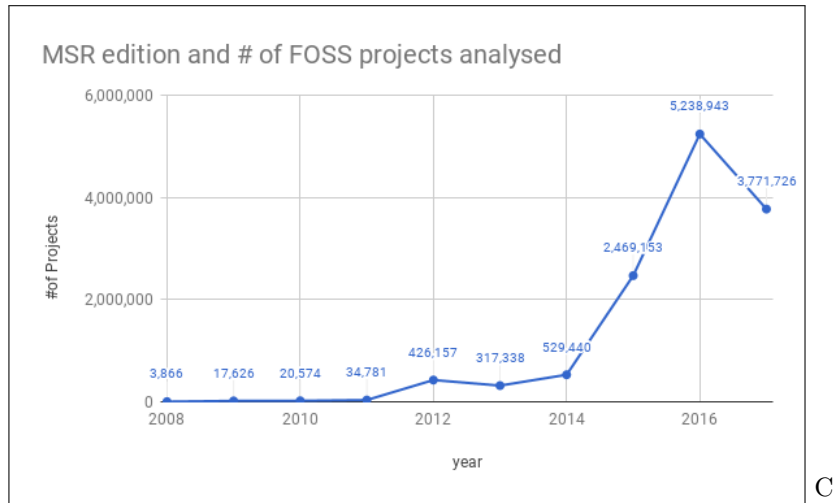


Figure 1: Cumulative number of FOSS projects per year.

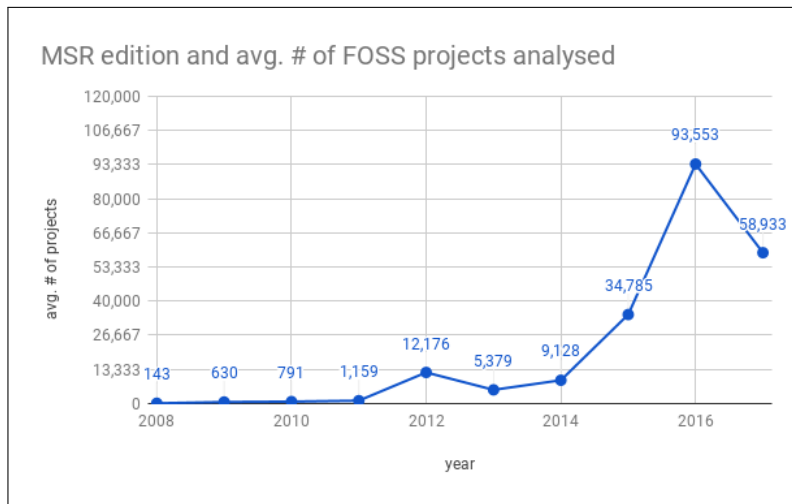


Figure 2: Average number of FOSS projects per year.

These always larger samples of systems have mostly overlooked their primary

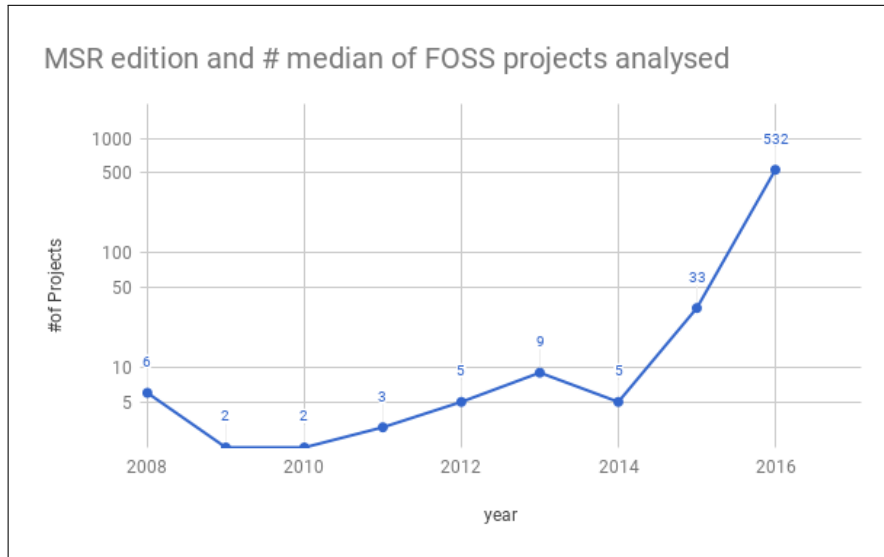


Figure 3: Median number of FOSS projects per year.

distinctive characteristics, their diversity, context, uniqueness and application domain. Very few works have clearly stated the similarity (or differences) between systems in the interpretation of the results, either by explicitly proposing explanations based on application domains [2, 3, 4], or by sampling the projects to be analysed from a specific, restricted topic [5].

This present paper is based on the assumptions that a specific software system might be *similar* to others to some degree, and that there are different approaches to defining their similarity. By applying one of those approaches, a sample of software systems might get divided into subsets (or *clusters*), each containing similar systems, and showing differences with other clusters.

Understanding the similarities among software projects allows for reusing of source code and prototyping, or choosing alternative implementations [6, 7], thereby improving software quality [8]. Having access to similar software projects helps developers speed up their development process. By looking at similar Open Source Software (OSS) projects, for example, developers are able to learn how relevant classes are implemented, and in some certain extent, to reuse useful source code [6, 7, 9].

In the past, two software projects have been considered to be similar if they implement some features being described by the same abstraction, even though they may contain various functionalities for different domains [10]. In this paper, we group similar systems into clusters based on different approaches: first, we group them based on the similarities detected by the CrossSim algorithm [11, 12], which has been developed as part of the EU H2020 CROSSMINER project². Systems are similar, or connected, if they have a limited *distance* [11]. Second, we use the clusters as manually extracted by [13, 14], that have grouped 5,000 software systems into 6 clusters. Third, we use a Python implementation of
50 the Latent Dirichlet Allocation (LDA) approach to automatically extract the descriptions of a project, and we group similar systems based on that extraction.

For all the clusters identified in this paper, we evaluated the metrics of an object-oriented suite, based on the work by [Chidamber and Kemerer \(CK\)](#) [15] and additional OO software metrics used in prior studies to complement the CK metrics [16, 17, 18, 19].

The aim of this paper is to answer the following research question (RQ):

RQ: are OO metrics sensitive to the context of their clusters?

In other words, do different software clusters exhibit different OO metrics? To answer that research question, we articulate this paper in the following parts:
60 Section 2 proposes a meta-model of the reasoning behind the need for clustering software systems. Section 3 deals with the related work. Section 5 describes the three approaches used to define an ecosystem that we use throughout this paper. Section 6 illustrates the datasets used. Section 7 provides the results of the statistical tests, that aim to reject, for every OO metric m , the null hypothesis $H_{0,m}$: *the samples are drawn from the same population*.

Section 8 discusses the findings, while Section 9 describes the threats to validity. Finally, Section 10 summarises and concludes the paper.

²<https://www.crossminer.org>

2. Reasons for clustering

Clustering is deemed to be among the fundamental techniques in knowledge
70 mining and information retrieval [20, 21]. In the context of software engineer-
ing, clustering has been used in the past, for reverse engineering and software
maintenance tasks, with the aim to categorise software artifacts [22, 23]. The
concept of *similarity* is a fundamental building block for any clustering tech-
nique, as well as a key issue in various contexts, such as detecting cloned code
[24, 25, 26, 27], software plagiarism [28], or reducing test suite size in model-
based testing [29, 30]. According to *Walenstein et al.* [31], a workable common
understanding for software similarity is as follows: “*the degree to which two*
distinct programs are similar is related to how precisely they are alike.” Never-
theless, a globally exact and shared definition of similarity has not been agreed
80 upon yet: depending on the method used to compare items, various types of
similarity may be identified.

Clustering techniques have been exploited in other fields including biology to
classify plants or animals according to their properties [32], and geology to clas-
sify observed earthquake epicenters and thus to identify dangerous zones [33]. A
clustering algorithm attempts to distribute objects into groups of similar objects
so as the similarity between one pair of objects in a cluster is higher than that
between one of the objects to any objects in a different cluster [34, 35]. In recent
years, several clustering methods have been developed to solve a wide range of
issues [36]. Among others, there are hierarchical and partitional clustering algo-
90 rithms [35]. The former use a criterion function to identify partitions while the
latter try to group similar partitions. Among partitioning-based algorithms,
there are K-Means, K-Medoids, CLARA, CLARANS [37, 38, 39]. Among
hierarchical-based algorithms, there are BIRCH [40], CURE [41], ROCK [42],
Chameleon [43], to name a few.

Several existing clustering techniques share the property that they can be
applied when it is possible to specify a *proximity (or distance) measure* that al-
lows one to assess if elements to be clustered are mutually similar or dissimilar.

The basic idea is that the *similarity level of two elements is inversely proportional to their distance*. The definition of the proximity measure is a key issue
100 in almost all clustering techniques and it depends on many factors including the considered application domain, available data, and goals. Once the proximity measure has been defined, it is possible to produce a proximity matrix for the related objects. Given that there are n objects to be clustered, an $n \times n$ proximity matrix needs to be generated containing all the pairwise similarities or dissimilarities between the considered objects.

Recommender systems rely heavily on similarity metrics to suggest suitable and meaningful items for a given item [6, 44]. For example, for third-party library recommendation, it is important to find similar projects to a given project, and mine libraries from the most similar projects. Similarities are used as a
110 base by both content-based and collaborative-filtering recommender systems to choose the most suitable and meaningful items for a given user [6]. In this sense, failing to compute similarities means concurrently adding a decline in the overall performance of these systems.

Nevertheless, measuring similarities between software systems has been considered as a daunting task [10, 45]. Furthermore, considering the heterogeneous nature of artifacts in open source software repositories, similarity computation becomes more complicated as many artifacts and several cross relationships prevail. As a result, similarity computation among software and projects has attracted considerable interest from many research groups. In recent years, several approaches have been proposed to solve the problem of software similarity
120 computation. Many of them deal with similarity for software systems, others are designed for computing similarities among open source software projects. Such approaches are domain specific and they can be classified according to the set of mined features. In particular, there are two main types of software similarity computation techniques as follows [45]. The first is called *low-level similarity* and it is calculated by considering low-level data, e.g., source code, byte code, function calls, API reference, etc. Meanwhile, *high-level similarity* is based on the metadata of the analysed projects e.g., similarities in readme files, textual

descriptions, star events, to name a few.

130 **3. Related Work**

While the primary goal of empirical papers is to achieve the generality of the results, the domain, context and uniqueness of a software system have not been considered very often by empirical software engineering research. As in the example reported in [46], the extensive study of all JSON parsers available would find similarities between them or common patterns. That type of study would focus on one particular language (JSON), one specific domain (parsers) and inevitably draw limited conclusions. On the other hand, considering the “parsers” domain (but without focusing on one single language) would show the common characteristics of developing that type of systems irrespective of their
140 language.

So far, several tools that capture the topics of software systems have been proposed. Among others, CLAN [47], CrossSim [11], MUDAbblue [48] and RepoPal [49] are some of the most notable tools, with various levels of precision and accuracy. CLAN outperforms MUDAbblue, furthermore it has also been reported that the similarity scores of CLAN reflect the perception of humans of similarity better than those of MUDAbblue [47]. Meanwhile, RepoPal has been proven to obtain a better performance compared to CLAN. Eventually, CrossSim surpasses RepoPal with respect to various quality metrics, including computation time [11]. Even if such tools are available for researchers and practitioners,
150 their usage to practically inform development and project clustering (finding similar projects) has been so far quite limited. For example, out of the 26 articles that have cited RepoPal [49], only two of them [11, 50] have actually adopted RepoPal for project clustering. The same applies to CrossSim [11] where two among the 7 studies that cite the paper have adopted the tool for project similarity detection [9, 51]. In addition, these few studies have adopted the tools when proposing or comparing a new software project similarity detection technique.

Wermelinger and Yu [52] posit that presenting two datasets from the same domain allows for future comparative studies and facilitates the reuse of data ex-
traction and processing scripts. On increasing the external validity of empirical
160 result findings, German *et al.* [53] have also highlighted the need to investigate
in particular systems belonging to different domains.

Prior research has shown that the number and size of open-source projects
are growing exponentially and open-source projects are becoming more diverse
by expanding into different domains [54, 55]. In view of this and to reduce the
effort required in manual categorisation of software projects, Tian *et al.* [56]
proposed a technique based on text mining to categorise software projects irre-
spective of the programming language used in their development.

The following prior studies highlight the need for investigating empirical
170 software engineering research results by application domains as well as the adop-
tion of software categorisation or clustering techniques when evaluating software
quality.

Callau *et al.* [57] studied the use of dynamic programming features such
as method and class creation and removal at run-time e.g., during testing and
how much these features are actually used in practice, whether some are used
more than others, and in which kinds of projects. Their results revealed three
application domains that rely heavily on the usage of dynamic programming
features: (i) *user interface applications*, which make heavy usage of dynamic
method invocation as a lightweight form of an event notification system, (ii)
180 *frameworks* that communicate with databases or implement object databases,
which make heavy usage of serialisation and de-serialisation of objects and (iii)
low-level system *support code* that uses object field reads and writes to imple-
ment copy operations, saving the state of the system to disk, and converting
numbers and strings from objects to compact bit representation.

Linares-Vasquez *et al.* [58] analysed the energy usage of API method calls in
55 different Android apps from different categories such as Tools, Music, Media
and others. Results revealed GUI and image manipulation apps made use of
the highest number of energy-greedy API method calls followed by Database

apps. Both categories represented 60% of the energy-greedy APIs in the studied
190 sample. In [57], authors made a study on the usage of dynamic programming
features in terms of the significant energy and memory usage of software in
the Databases category wherein database-based software made heavy usage of
serialisation and de-serialisation of objects.

In a related study, the focus is on energy management in Android applica-
tions [59] with an analysis of different power management commits (including
Power Adaptation, Power Consumption Improvement, Power Usage Monitoring,
Optimizing Wake Lock, Adding Wake Lock and Bug Fix and Code Refinement).
The studied projects were clustered into 15 categories. The top three categories
in terms of the number of power management commits were found to be Con-
200 nectivity, Development and Games.

Previous studies [60, 61, 62] revealed that projects from different domains
use exception handling differently, and that poor practices in writing excep-
tion handling code are widespread. In a study on Java projects by Osman
et al. [2] they aimed to answer the following research question: “Is there any
difference in the evolution of exception handling between projects belonging
to different domains?” The researchers manually categorised 30 projects into
6 domains, namely compilers, content management systems, editors/viewers,
web frameworks, testing frameworks, and parser libraries. Their observations
showed significant distinctions in the evolution of exception handling between
210 these domains, like the usage of `java.lang.Exception` and custom exceptions
in catch blocks. Concretely, content management systems consistently have
more exception handling code and throw more custom exceptions, as opposed
to editors/viewers, which have less error handling code and mainly use standard
exceptions instead.

In general, different results, trends or patterns have been observed in prior
empirical studies when more attention is paid to the categorisation of analysed
software projects [based on predefined software categories or domains](#).

[Regarding the effect of application domains on software metric values or
thresholds, Alves *et al.* \[63\] emphasise that the effective use of software metrics](#)

220 is hindered by the lack of meaningful thresholds. The authors designed a method that determines metric thresholds empirically from measurement data. Their results did not focus on application domains, however in their sample selection they did select 100 OO software systems from various application domains and a combination of open and closed source.

On the other hand, Ferreira *et al.* [64] mentioned that software metrics are not widely adopted in industry and one possibility could be the lack of availability of reference values for most metrics. Thresholds can aid specialists to apply metrics in their tasks and Medicine is an example of field in which the work of the specialist is strongly supported by metrics and their thresholds. The authors presented a study on the structure of 40 open-source programs developed
230 in Java, of varying sizes and ensured that the projects were from 11 application domains. They aimed to define thresholds for a set of OO metrics (e.g., LCOM, DIT, number of public methods, etc.) and evaluated the practicality of the proposed thresholds. Four types of analyses were carried out: with the entire dataset, which leads to general thresholds; by software system sizes; by application domains; and by types of software systems (tool, library and framework). Derived thresholds were broken down, for example for the LCOM metric the derived thresholds for the entire dataset are as follows: Good: 0; regular: 1-20; bad: greater than 20. These thresholds applied to only 2 (Clustering and
240 Database domains) out of the 11 studied domains. However, for the Hardware domain, the LCOM thresholds are as follows: Good: 0; regular: 1-80; bad: >80 while the LCOM thresholds for the Games domain are as follows: Good: 0; regular: 1-35; bad: >35. Their results show that the thresholds vary when considering the application domains and that the proposed thresholds could be used to support both the identification of classes which violate design principles, as well as well-designed classes.

Oliveira *et al.* [65] rightly emphasised that “*establishing credible thresholds is a central challenge for promoting source code metrics as an effective instrument to control the internal quality of software systems.*” Thus, in an attempt to
250 resolve this challenge, they investigated 106 software projects and proposed the

concept of relative thresholds for evaluating software metrics data given that source code metrics usually follow heavy-tailed distributions and it is natural to have some outlier artefacts that will not strictly follow a specified threshold, e.g., p% of the entities should have $M_{metric} \leq k$, where M is a source code metric calculated for a given software entity (method, class, etc.), k is the upper limit, and p is the minimal percentage of entities that should follow this upper limit. The application of their approach was evaluated on a sub-corpus of the original corpus including systems sharing a common functional domain. They identified slight variances in the relative thresholds when looking at the Tools domain. For example, the original threshold regarding all systems in the Qualitas Corpus³ for the FAN-OUT metric is as follows: 80% of the classes should have FAN-OUT ≤ 15 . However, the relative threshold for the same metric considering only the *Tools* software category which consists of 26 projects making up 24.5% of the original corpus is as follows: 85% of the classes should have FAN-OUT ≤ 20 . This was observed because in the *Tools* sub-corpus, the tail classes tend to be bigger than the typical tail classes in the whole corpus leading to a small difference only in the last quantiles.

In a different study by De Souza and Maia [66], software coupling metrics were studied based on software categories to identify any impact of software categories on coupling metrics (CBO⁴, DAC⁵, ATFD⁶ and AC⁷). The authors emphasised the need to pay a special attention to software categories when comparing systems in distinct categories with predefined thresholds already available in the literature. For example, empirical results from the study revealed that out of ten distinct categories selected (including *Audio and Video*, *Graphics*, *Security and Games*) there is a different level of coupling among the different categories. Games had a higher coupling level while the Development category

³<http://qualitascorpus.com/>

⁴Coupling between Objects

⁵Data Abstraction Coupling

⁶Access to Foreign Data

⁷Afferent Coupling

showed less coupling than others. Statistical tests conducted at a 0.01 significance level supported these results which indicate the importance of analysing software engineering research results by domain or category.

280 Recently, Mori *et al.* [67] argued that while deriving reliable thresholds for software metrics has been an ongoing research concern, there is still a lack of evidence about threshold variation across different software domains. In an attempt to address this limitation, the authors whether and how thresholds vary across domains in a study on 3,107 software systems from 15 domains (including communication, development, education, games and others) with a focus on 8 well-known software metrics in terms of size (LOC, NOA, NOM), complexity (WMC, LCOM, CBO) and inheritance (DIT, NOC). Presented results, showed that software domain and size affect thresholds and the following metrics: LCOM, WMC, and LOC are highly sensitive to the software domain. Furthermore, the authors investigated benchmark OO software datasets 290 namely Qualitas Corpus, AllSystems, MediumSystems and SmallSystems. Results showed that Qualitas Corpus and AllSystems have similar thresholds, although they do not have any system in common. Observed results suggested that if the benchmark is composed of a high enough number of heterogeneous systems (i.e., from different domains and sizes), the metrics thresholds tend to be comparable.

300 In general, earlier studies have shown that application domains indeed have an effect on software metric values and thresholds. Therefore, developers should pay attention to this variation when using thresholds for these metrics in software quality evaluation.

4. Experimental Design

In this section we detail the experimental design of our study, articulating it in techniques of clustering, research objectives, data sources, statistical tests performed and null hypothesis.

- *Clustering techniques*: in this paper we compare three different techniques to cluster software systems: the *CrossSim* tool, a manual clustering approach, and a LDA-based approach.
- *Research objective*: considering the clusters generated by the three techniques above, this study evaluates whether they differ with a statistical significance, and based on Object-Oriented (OO) metrics.
- *Data sources*: as the data sources for this paper we consider:
 - [1] the project clusters derived with the *CrossSim* tool by using a sample of 12 projects (6 pairs), from a larger population of 5,000 projects extracted as part of the CROSSMINER project [68];
 - [2] the project clusters obtained by a *manual* classification, using a sample of 500 Java projects, as presented in [13, 14];
 - [3] the project clusters obtained by the LDA-based approach, using a sample of 100 GitHub Java projects, but also considering the Java samples from [2] and [3] above.
- *Measurements*: for every software project analysed, we collected 9 well-known OO attributes (NOC, DIT, CBO, RFC, WMC, LCOM, NIM, IFANIN, NIV). We describe each attribute in more detail in section 3. To evaluate the statistical hypotheses, we assigned the OO measurements to the same pot, if the projects they belong are in to the same cluster.
- *Statistical tests*: for each dataset and clustering technique, we adopted the Kolgomorov-Smirnov (KS) test [69] to detect whether the distribution of software metrics in the compared domains are from the same population. The appropriate Bonferroni correction [70] was applied, due to the multiple tests being carried out at the same time. As a consequence, the base α value (i.e., 0.05) for each test varies, based on the number of clusters considered for that test.

- *Null hypothesis*: given the clustering technique c , two project clusters C_1 and C_2 , an Object-Oriented metric m , the null hypothesis states: “using c as clustering technique, the values of the metric m applied on the projects in C_1 and C_2 come from the same population”

4.1. Metrics extracted

The metrics extracted for the projects are well-known structural OO attributes (NOC, DIT, CBO, RFC, WMC, LCOM, NIM, IFANIN, NIV) [71] described as follows:

- NOC: Number of Children (CK): number of direct sub-classes of a class;
- DIT: Depth of Inheritance Tree (CK), is the length of the longest path from a given class to the root class in the inheritance hierarchy;
- CBO: Coupling between Objects (CK), two classes are coupled if one acts on the other⁸. Hence, CBO is the count of other classes coupled to a class.
- RFC: Response For a Class (CK), the sum of the number of methods defined in a class and the cardinality of the set of methods called by them and belonging to external classes;
- WMC: Weighted methods per class (CK), a weighted sum of all the methods defined in a class;
- LCOM: Lack of Cohesion of Methods in a class (CK), the LCOM metric is based on the concept of the similarity of methods in a class. The degree of similarity of two methods M_1 and M_2 is the intersection set of instance variables⁹ used by both methods for functionality. Based on this notion, the LCOM of a class is the count of method pairs where the intersection

⁸If methods in a class use methods or instance variables defined by another class.

⁹Member variables declared in a class for which instances of the class own a separate copy.

set is equal to zero (i.e., a null set) minus the count of method pairs whose similarity is not zero¹⁰.

- NIM: Number of instance methods, methods defined in a class that are only accessible through an object of that class;
- IFANIN: Number of immediate base classes of a class;
- NIV: Number of instance variables, variables defined in a class that are only accessible through an object of that class.

We extracted the metrics using the Understand©tool available from `scitools.com`. We were careful to only extract individual methods data, but not class-wide or package-wide metrics. The latter ones contain aggregated metrics that are less relevant for our analysis. Also, we considered both test and non-test classes in our analysis.

5. Techniques of clustering

In this section we present the three approaches that we used to cluster the systems in our samples: the one illustrated in Section 5.1 clusters projects based on how they are linked to external libraries and components.

The technique described in Section 5.2 is based on 6 pre-determined categories, and the subjective attribution of each project to a cluster. The approach is based on the work proposed by the authors of [13, 14].

The last technique, illustrated in Section 5.3, uses the LDA algorithm to first extract the topics of the software systems, and then it assigns each project to one cluster.

We provide the replication packages, and the datasets with our results, in three collections available online¹¹.

¹⁰If the number of similar methods is more than the non-similar methods, then the class is more cohesive.

¹¹CrossSim collection: <https://figshare.com/s/0cd6925bf1601e2f0b74>,

5.1. Pairwise similarity via CrossSim

Linked Data is a representation method that allows for the interlinking and semantic querying of data [72]. The core of Linked Data is an RDF¹² graph that is made up of several nodes and oriented links to represent the semantic relationships among various artifacts. Thanks to this feature, the representation paves the way for various computations. One of the main applications of RDF is similarity computation for supporting recommender systems [73]. We designed and implemented CrossSim [11] (Cross Project Relationships for Computing Open Source Software Similarity), a tool for computing similarity among OSS projects
390 by considering the analogy of typical applications of RDF graphs and the problem of detecting the similarity of open source projects. CrossSim exploits graphs for representing different types of relationships in the OSS ecosystem. In particular, with the adoption of the graph representation, CrossSim transforms the relationships among non-human artifacts, e.g., API utilisations, source code, interactions, and humans (e.g., developers) into a mathematically computable format, i.e., one that facilitates various types of computation techniques.

The architecture of CrossSim is depicted in Figure 4: the rectangles represent artifacts, whereas the ovals represent activities that are automatically performed by the developed CrossSim tooling. In particular, the approach imports project
400 data from existing OSS repositories and represents them in a graph-based representation by means of the *OSS Ecosystem Representation* module. Depending on the considered repository (and thus to the information that is available for each project) the graph structure to be generated has to be properly configured. For instance, in case of GitHub, specific configurations have to be specified in order to enable the representation in the target graphs of the stars assigned to each project. Such a configuration is repository specific, e.g., SourceForge does not provide the star based system available in GitHub.

Categorised collection: <https://figshare.com/s/bd506aff8bff2b6ce29f>,

LDA-based collection: <https://figshare.com/s/93f5d997484705a937de>

¹²<https://www.w3.org/TR/2014/REC-rdf11-concepts-20140225/>

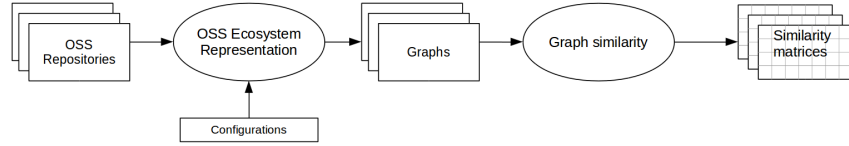


Figure 4: Overview of the CrossSim approach [11].

The *Graph similarity* module implements the similarity algorithm that applied on the source graph-based representation of the input ecosystems generates
 410 matrices representing the similarity value for each pair of the input projects.

To demonstrate the utilisation of graphs in an OSS ecosystem, we consider an excerpt of the dependencies for a pair of OSS projects, namely **project#1** and **project#2** in Figure 5. Using dependency information extracted from source code and the corresponding metadata, this graph can be properly built to represent the two projects as a whole. In this figure, **project#1** contains code file `HttpSocket.java` and **project#2** contains `FtpSocket.java` with the corresponding edges being marked with the semantic predicate `hasSourceCode`. Both source code files implement `interface#1` being marked by the semantic

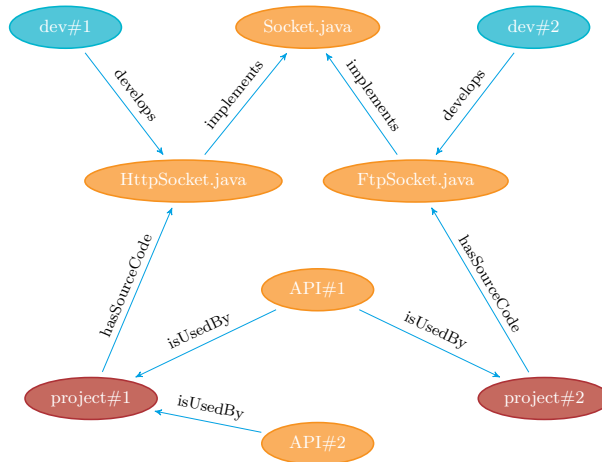


Figure 5: Sample graph-based representation of OSS ecosystems [11].

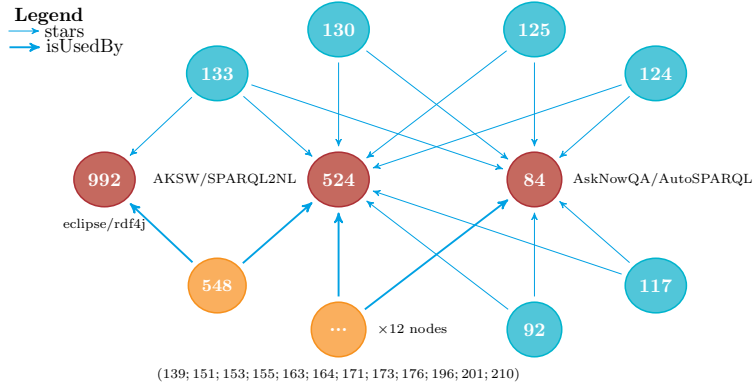


Figure 6: Sub-graph showing a fragment of the representation for three projects, i.e., AskNowQA/AutoSPARQL, AKSW/SPARQL2NL, and eclipse/rdf4j [11].

420 predicate `implements`. `Project#1` and `project#2` are also connected via other semantic paths, such as `API isUsedBy`.

Based on the graph structure, one can exploit nodes, links and the mutual relationships to compute similarity using existing graph similarity algorithms. To the best of our knowledge, there exist several metrics for computing similarity in graphs [74]. CrossSim adopts SimRank [75] as the mechanism for computing similarities among OSS graph nodes. SimRank has been developed to calculate similarities based on mutual relationships between graph nodes. Considering two nodes, the more similar nodes point to them, the more similar the two nodes are. We can compute the similarity between `project#1` and `project#2` with regards to the semantic paths between them, e.g., the two-hop path using `hasSourceCode` and `implements`, or the one-hop path using `API isUsedBy`. For example, concerning `isUsedBy`, the two projects are considered to be similar since with the predicate both projects originate from `API#1`. The hypothesis is based on the fact that the projects are aiming at creating common functionalities by using common libraries [10, 76].

Figure 6 is a concrete example to explain how the graph representation is exploited in CrossSim. The sub-graph represents the relationships between

two projects AskNowQA/AutoSPARQL and AKSW/SPARQ-L2NL. The orange nodes are dependencies whose real names are listed in Table 1. The turquoise nodes are developers who already starred the repositories. Every node is encoded using a unique number across the whole graph. The final result that lies between 0 and 1 is the similarity between two projects AskNowQA/AutoSPARQL and AKSW/SPARQL2NL.

ID	Name
139	org.apache.jena:jena-arq
151	org.dllearner:components-core
153	net.didion:jwnl:jwnl
155	net.sourceforge.owlapi:owlapi-distribution
163	net.sf.jopt-simple:jopt-simple
164	jaws:core
171	com.aliasi:lingpipe
173	org.dllearner:components-ext
176	org.apache.opennlp:opennlp-tools
196	org.apache.solr:solr-solrj
201	org.apache.commons:commons-lang3
210	javax.servlet:servlet-api
548	org.slf4j:log4j-over-slf4j

Table 1: Shared dependencies in Figure 6 [11].

5.2. Clustering based on projects descriptions

As the second approach to clustering, we consider the work proposed in [13, 14]. The authors state that “*we manually classified the domain of each system in our dataset. Initially, the first author of this paper inspected the description of the top-200 repositories to provide a first list of application domains, distributed over six domain types.*”

As the clustering unit, this approach uses the collection of software systems sharing the same application domain. The relevance of application domains as a driving factor for specific development approaches has been long recognised: in Glass and Vessey’s seminal paper it is mentioned that “*it has become clear*

that application-independent techniques and tools must be supplemented with an application-specific approach” [77].

The projects here are the Java subset of the data dump available at https://zenodo.org/record/804474#.XD11S9_njCK. The dataset contains 5,000 project names hosted on GitHub, and the authors pre-determined the following six categories:

- Application Software (AS)
- 460 • Documentation (D)
- Non Web Libraries And Frameworks (NW)
- Software Tools (ST)
- System Software (SS)
- Web Libraries And Frameworks (WL)

In the original paper, each project was assigned, by one of the authors, to one of those six categories, based on its characteristics. The assignment was later validated by a second author. Albeit the process was triaged, it is difficult to exactly reproduce the original classification.

470 The work done in [13, 14] is a unique case of third-party assignment of software projects to categories. Therefore, we make use of the sample, and the classification, and treat it as a ‘manual’ clustering technique. On the other hand, we cannot fully replicate their clustering approach in the other two samples: the knowledge to assign the projects to the six categories is mostly informal, and not fully reproducible.

5.3. LDA-informed clustering

This third clustering technique is similar to what presented in Section 5.2, but it adds an automated step to extract the topics contained in a software system. This should be helpful for the reproducibility of the approach.

For this purpose, we extracted the lexical content (e.g., its *corpus*) of each
480 Java class in two ways: (i) by considering their class names; and (ii) by parsing
their code and considering the variable names, comments and keywords. It
is worth noting that the parser that we built to extract the lexical content is
based on Java software, and its keywords. It could be further expanded into
other programming languages, by taking into account their specific constructs
and keywords. To facilitate future research, we made the parser available online
at <https://figshare.com/account/projects/56009/articles/9785861>.

The code of each Java class is converted into a *text corpus* where each line
contains elements of the implementation of a class. The corpus in this case
("dictionary" of terms derived from comments, keywords in source code) is
490 built at the *class* level of granularity [78]. The corpus includes the class name,
variable and method names and comments for each class.

Pre-processing of the system corpus is needed to eliminate common key-
words, stop words, split and to stem class names [79]. Also, we do not consider
as a term any of the Java-specific keywords (e.g., *if*, *then*, *switch*, etc.)¹³. Ad-
ditionally, the *camelCase* or *PascalCase* notations are first decoupled in their
components (e.g., the class constructor *InvalidRequestTest* produces the terms
invalid, *request* and *test*). Second, each term goes through a phase of *stemming*
and *lemmatisation*, to extract its root.

As an example of this lexical extraction, for the lines of code shown in Figure
500 7 (the *UrSQLEntry.java* class from the *UrSQL* project), we derive the following
corpus: {*ur sql entri kei valu kei valu ur sql entiti entiti ur sql entri ur sql entri*
queri split queri split ur sql control kei valu separ kei split valu split kei kei valu
valu}.

For each system, all the Java classes were reduced to a corpus of terms. All
these terms were then considered to create a model implementing the Latent

¹³The complete list of Java reserved words that we considered is available at https://en.wikipedia.org/wiki/List_of_Java_keywords. The `String` keyword was also considered as a reserved word, and excluded from the text parsing.

```

1 package tmacsoftware.ursql;
2
3 public class UrSQLEntry
4 {
5
6     private String key;
7     private String value;
8     private String firstKey;
9     private String firstValue;
10    private UrSQLEntity entity;
11
12    public UrSQLEntry()
13    {
14    }
15
16    public UrSQLEntry(String query)
17    {
18        String[] split = query.split
19        (UrSQLController.KEY_VALUE_SEPARATOR);
20        this.key = split[0];
21        this.value = split[1];
22        this.firstKey = this.key;
23        this.firstValue = this.value;
24    }

```

Figure 7: UrSQLEntry.java Source Code Snippet.

Dirichlet Allocation (LDA) algorithm. Python is the programming language used to program the models, and the *gensim* NLP package¹⁴ helped by the machine learning side.

To extract the main topics emerging from the corpus of a software project, we
510 utilize LDA which is a topic-modeling technique [80]. Each document is featured using a Natural Language Processing approach termed the Term-frequency-inverse document frequency (TF-IDF). In NLP, TF-IDF [81] is used to measure the weight of a term within documents (in our case, the source code of a class). With TF-IDF, words are assigned weights, as the product of term frequency and inverse document frequency. We use TF-IDF as a pre-processing step to LDA: the result is a representation of the source code contained in the Java classes, where the same terms can appear multiple times (see Figure 7).

As an example, for the *okhttp* project¹⁵, the LDA model produces the following topics from the corpus of its Java classes:

¹⁴<https://radimrehurek.com/gensim/>

¹⁵<https://github.com/square/okhttp>

Topics extracted with the LDA approach

Topic 0: 0.003***stream**" + 0.003***bodi**" + 0.003***header**" + 0.003***content**" + 0.003***id**" + 0.002***benchmark**" + 0.002***type**" + 0.002***ssl**" + 0.002***socket**" + 0.002***stori**"

Topic 1: 0.002***entiti**" + 0.002***url**" + 0.002***proxi**" + 0.002***slack**" + 0.002***event**" + 0.001***frame**" + 0.001***filter**" + 0.001***client**" + 0.001***equal**" + 0.001***session**"

Topic 2: 0.005***cooki**" + 0.004***header**" + 0.004***interceptor**" + 0.003***chain**" + 0.003***url**" + 0.002***bodi**" + 0.002***certif**" + 0.002***content**" + 0.002***client**" + 0.002***timeout**"

Topic 3: 0.005***cach**" + 0.004***socket**" + 0.004***connect**" + 0.004***bodi**" + 0.003***rout**" + 0.003***server**" + 0.003***web**" + 0.003***header**" + 0.003***client**" + 0.003***url**"

Topic 4: 0.006***event**" + 0.006***socket**" + 0.005***certif**" + 0.005***address**" + 0.005***cach**" + 0.004***file**" + 0.003***deleg**" + 0.003***connect**" + 0.003***server**" + 0.003***inet**"

520

The topics extracted from the projects were finally assigned to a category by two of the authors of this paper: in case of disagreement, a discussion was held to reconcile the views. As the list of categories, we adopted in fact what has been historically used by the `SourceForge.net` repository to classify the hosted projects:

List of categories used in `SourceForge.net`

1. Communications
2. Database
3. Desktop Environment
4. Education
5. Formats and Protocols
6. Games/Entertainment
7. Internet
8. Mobile
9. Multimedia
10. Office/Business
11. Other/Nonlisted Topic
12. Printing
13. Religion and Philosophy
14. Scientific/Engineering
15. Security
16. Social sciences
17. Software Development
18. System
19. Terminals
20. Text Editors

6. Datasets used for the Empirical Analysis

This section presents the datasets that have been used for performing the empirical analysis. In particular, the dataset used by CrossSim is presented in Section 6.1. The categorised dataset as presented in [13] is summarized in Section 6.2. The dataset used to apply the LDA-based technique is presented in Section 6.3.

6.1. CrossSim dataset

The overall dataset gathered by the CROSSMINER project to evaluate the CrossSim approach consists of 580 GitHub Java projects. Such a dataset was collected from GitHub by considering the following requirements: (i) being GitHub Java projects; (ii) providing the specification of their dependencies by means of pom.xml or .gradle files (iii) for CrossSim, it is necessary to consider only projects that are of decent quality. Thus we selected only projects having at least 9 external dependencies; (iv) having the README.md file available; (v) possessing at least 20 stars. The rationale behind the selection of 9 libraries and 20 stars is as follows. By performing various empirical evaluations on the CrossSim tool [11, 51] we realized that a project including at least 9 libraries is suitable to be used as input for the similarity computation. In the data collection phase, we tried to cover a wide range of possibilities by taking into consideration that many repositories in GitHub are of low quality, which is especially true when they do not have many stars. Thus, we consider only projects that have been starred by at least 20 developers. Such a number of stars has been used in some studies [7, 82] as a sign of a decent project. The collected dataset and the CrossSim tool are available online for public usage [68].

For the purpose of our analysis, 6 pairs of software systems were extracted from the CROSSMINER dataset. The criteria of selection were:

- maximum similarity between projects within the pairs;
- minimum similarity between projects of different pairs.

Cluster	Projects	Similarity (between pairs)
AB	A = JamsMusicPlayer ¹⁶ B = ACEMusicPlayer ¹⁷	1.95e-04
CD	C = sparql-plugin ¹⁸ D = neo4j-sparql-extension ¹⁹	0.0027
EF	E = jpmml-model ²⁰ F = visitante ²¹	7.32e-04
GH	G = c2d-engine ²² H = LeanEngine-Server ²³	5.47e-04
IJ	I = RestOpenGov ²⁴ J = elasticsearch-analysis-lemmagen ²⁵	0.00112
KL	K = jcabi-email ²⁶ L = jcabi-jdbc ²⁷	0.0048

Table 2: Clusters of projects identified by CrossSim.

The outcome of the extraction satisfying such criteria is shown in Table 2.

6.2. Manually categorised dataset (from [14])

From the original, overall sample of 5,000 projects, we extracted all the Java projects (520 projects in total), while maintaining the information about their assigned category. Considering the Java subset, and using the categories provided, we have the distribution of projects within the ecosystems as shown in Table 3. We extracted the metrics for each project, then clustered those metrics based on the six domains identified above. A statistical test was run between

¹⁶ Available at [git://github.com/psaravan/JamsMusicPlayer.git](https://github.com/psaravan/JamsMusicPlayer.git)

¹⁷ Available at [git://github.com/C-Aniruddh/ACEMusicPlayer.git](https://github.com/C-Aniruddh/ACEMusicPlayer.git)

¹⁸ Available at [git://github.com/neo4j-contrib/sparql-plugin.git](https://github.com/neo4j-contrib/sparql-plugin.git)

¹⁹ Available at [git://github.com/niclashoyer/neo4j-sparql-extension.git](https://github.com/niclashoyer/neo4j-sparql-extension.git)

²⁰ Available at [git://github.com/jpmml/jpmml-model.git](https://github.com/jpmml/jpmml-model.git)

²¹ Available at [git://github.com/pranab/visitante.git](https://github.com/pranab/visitante.git)

²² Available at [git://github.com/lycyng/c2d-engine.git](https://github.com/lycyng/c2d-engine.git)

²³ Available at [git://github.com/PeterKnego/LeanEngine-Server.git](https://github.com/PeterKnego/LeanEngine-Server.git)

²⁴ Available at [git://github.com/RestOpenGov/RestOpenGov.git](https://github.com/RestOpenGov/RestOpenGov.git)

²⁵ Available at [git://github.com/vhyza/elasticsearch-analysis-lemmagen.git](https://github.com/vhyza/elasticsearch-analysis-lemmagen.git)

²⁶ Available at [git://github.com/jcabi/jcabi-email.git](https://github.com/jcabi/jcabi-email.git)

²⁷ Available at [git://github.com/jcabi/jcabi-jdbc.git](https://github.com/jcabi/jcabi-jdbc.git)

Category	Projects
Application Software	30
Documentation	48
Non Web Libraries And Frameworks	342
Software Tools	49
System Software	26
Web Libraries And Frameworks	25

Table 3: Number of Java projects in the categories extracted in [14].

each pair of domains, and per OO metric, to determine if the metrics come from the same population.

6.3. LDA-based dataset

Leveraging GitHub, we collected the project IDs of the 100 most successful Java projects hosted on GitHub as case studies.²⁸ The “success” of the projects is determined by the number of *stars* received by the community of GitHub users and developers, as a sign of appreciation. [Forking is a means to contribute to the original repositories \[83\], and there is a strong correlation between forks and stars \[84\]. In this sense, we suppose that a project with a high number of stars means that it gets attention from the OSS community, and thus being suitable to identify popular repositories \[85\]. We used this approach to stratified sampling because the projects obtained by this filter are likely to be used by a large pool of users \[86\], and potentially have a good intake of new developers \[87\].](#)

[Smaller projects are less likely to be sampled by this stratified approach. To prove this claim, let us consider the curated population of 14,118 Java projects, contained in the research published in \[88\]. When we extracted a random sample of 100 Java projects from that population, we observed that 29 of these projects contain less than or equal to 10 Java files. The violin boxplot in Figure 8 shows how the smaller projects skew the distribution, as compared to the original](#)

²⁸The list of projects is available at <https://figshare.com/s/c627af8e9e496a9025c4>

GitHub sample.

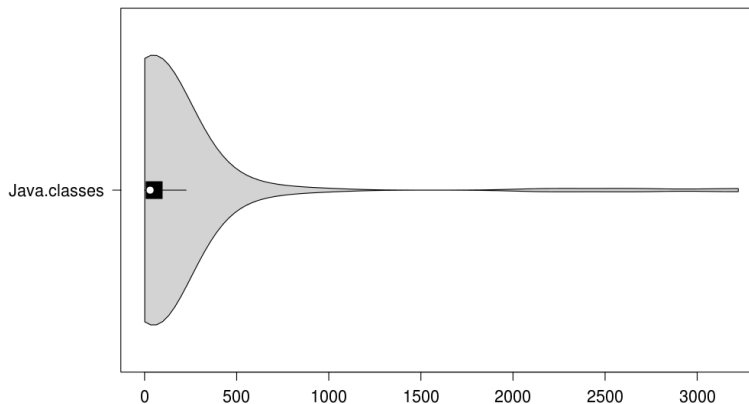


Figure 8: Number of Java files from a random sample of GitHub projects.

The source code of the selected systems was downloaded for the analysis: only the ‘master’ branch of the systems was considered. The boxplots in Figure 9 show the basic characteristics of the sample analysed: the distributions of duration of the projects (in days), the number of distinct developers (as authors) and the total number of commits are plotted. In terms of most likely value of each distribution, we observed that the median in the project’s duration is 2,002
590 days; the median number of developers is 87; while the median in the number of commits is 1,204.

Figure 10 shows the distribution of application domains in the sample of 100 Java projects, as extracted by the LDA algorithm, then assigned by the authors of the paper, finally agreed between authors to ensure consistency.

A few of the basic domains, as used by SourceForge, are completely absent from our sample: *Desktop Environment*, *Education* or *Text Editor* (and a few others) are not represented when sampling projects based on their success (e.g., usage or further development).

On the other hand, there are 4 domains that are more prominent than others: *Internet* (with 27 projects), *Mobile* (11), *Multimedia* (11) and *Software Development* (33). For the statistical analysis, we use only these 4 domains to
600 find differences between the distributions of metrics: using smaller sized clus-

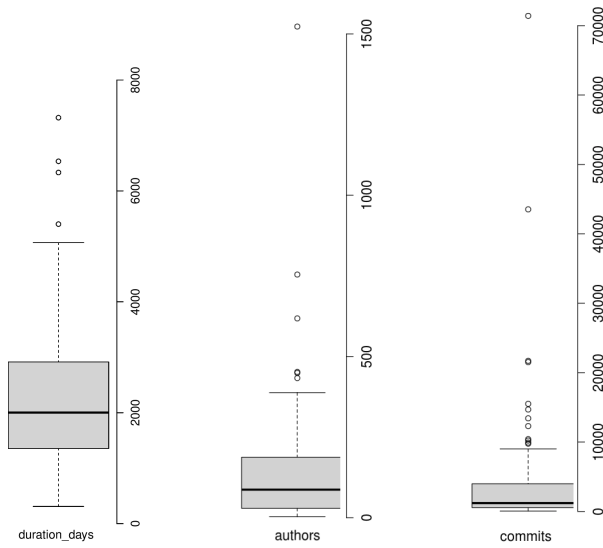


Figure 9: Distribution of duration (in days), number of developers (as authors) and number of commits in the sample analysed with the LDA approach.

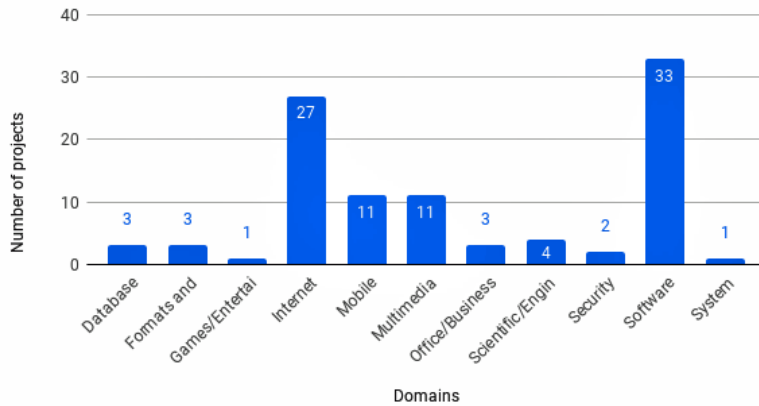


Figure 10: Domains extracted with the LDA approach.

ters would suffer from a small effect size, and the relative results would be less relevant.

7. Results

In this section we report the results of the analysis using the three clustering techniques: the one implemented by the CrossSim algorithm (Section 7.1); the mostly manual one (Section 7.2); and the semi-automated one (Section 7.3).

7.1. CrossSim dataset – Results

610 In order to evaluate the rejection of (or the inability to reject) the null hypothesis, we performed the KS statistical tests for each of the 6 project pairs (AB, CD, EF, GH, IJ and KL) described in Section 6.1. We assigned a standard threshold value $\alpha = 0.05$ for the sensitivity of the test; the Bonferroni correction resulted in $\alpha_B = 0.0034137$.

In Table 4 we summarise the tests by reporting the p-value of each KS test. We reject the null hypothesis (i.e., “the two samples come from the same distribution”) if the p-value of the KS test is lower than the corrected α_B . This is done for each metric, and for each pair of ecosystems. As an example, the KS test for the IFANIN metric, and the pair of ecosystems AB and CD produces a
620 p-value = 1: we reject the null hypothesis that the values of the IFANIN metric come from the same population, when considering the AB and CD ecosystems.

The results of each KS test (in the form of p-values) are summarised in Table 4: as an example, we rejected all the null hypotheses for the EF and IJ clusters: none of the OO attributes can be considered to be sampled from different populations. The most dissimilar pair of clusters appears to be the AB and GH pairs: for most of the structural OO metrics, we could reject the relative null hypothesis, therefore the IFANIN, NIM, NIV, WMC, RFC and DIT values can be considered as drawn from different populations.

In general however, for most of the metrics, and for most of the pairs of
630 ecosystems, the null hypothesis cannot be rejected. When we grouped projects

	IFANIN	CBO	NOC	NIM	NIV	WMC	RFC	DIT	LCOM
AB v CD	1	0.000	0.882	0.000	0.083	2.95e-06	2.51e-11	0.043	0.025
AB v EF	0.874	3.02e-07	0.000	0.006	3.51e-05	1.72e-14	0	3.02e-09	0.894
AB v GH	1.11e-10	0.009	0.051	7.07e-06	3.16e-05	4.41e-07	0	0	0.009
AB v IJ	0.077	0.146	0.339	0.192	0.640	0.533	0.371	1.19e-08	0.497
AB v KL	1.33e-15	1.17e-10	0.497	0.156	0.001	0.116	0.122	0	0.067
CD v EF	1.000	0.030	0.999	0.001	0.014	0.001	0.000	0.285	0.030
CD v GH	0.197	1.95e-06	1	0.001	0.031	1.46e-05	0.013	0.000	0.002
CD v IJ	0.558	0.004	0.967	0.005	0.672	0.000	6.43e-06	0.206	0.789
CD v KL	2.13e-06	0.128	1	8.69e-05	0.126	1.14e-05	9.31e-10	9.31e-10	0.001
EF v GH	0.001	6.13e-05	0.272	0.151	0.995	0.012	0.001	0.000	0.267
EF v IJ	0.337	0.484	0.999	0.363	0.155	0.190	0.014	0.040	0.333
EF v KL	2.17e-11	0.0046	0.979	0.567	0.002	0.185	0.000	0	0.154
GH v IJ	0.926	0.753	0.984	0.475	0.294	0.539	0.000	0.000	0.095
GH v KL	3.17e-06	5.60e-08	0.999	0.084	0.006	0.093	2.14e-10	0	0.516
IJ v KL	0.001	0.003	0.885	0.128	0.303	0.647	0.367	5.67e-05	0.033

Table 4: Results of the pair-wise statistical tests of the OO metrics analysed: AB, CD, EF, GH, IJ and KL refer to the pairs of projects. Highlighted the p-values larger than α_B .

based on the type of dependencies that they have, it is in general difficult to conclude that the clusters are structurally different from each other.

The result is not surprising: the projects paired by CrossSim are associated by what type of external libraries they include, as well as by the developers who

have starred or contributed to them. The similarity computed by CrossSim takes also into account the fact that a project has been starred or committed by different developers. In this respect, the similarity among them could be low even though they share some libraries in common.

7.2. Manually categorised dataset (from [14]) – Results

640 The same approach to statistical testing was also applied for the second dataset. Table 5 shows the results of the statistical analysis: as above, each cell contains the p-value of the Kolgomorov-Smirnov (KS) test between two subsets of the dataset. For example, the ASvD row contains the p-values of the KS test between the projects in the *Application Software* domain, and the projects in the *Documentation* domain. As above, the null hypothesis is rejected when the p-value is lower than a threshold α : the Bonferroni correction produces a corrected threshold of $\alpha_B = 0.0033$ (i.e., $0.05/15$ where 15 is the number of multiple tests performed).

We have highlighted in Table 5 the specific OO attribute where we reject the null hypothesis. As can be seen, the NOC, DIT and LCOM attributes
650 can be considered from the *same* distribution (e.g., we can not reject the null

	IFANIN	CBO	NOC	NIM	NIV	WMC	RFC	DIT	LCOM
ASvD	0	0	0	0	0	0	0	0	4.76e-06
ASvNW	0	0	0	3.23e-12	0	2.22e-16	0	0	0
ASvSS	1.37e-06	0	0.0247	0	1.14e-10	0	0	0.476	0.204
ASvST	0	0	0	0	0	0	0	0	0
ASvWL	0	0	0	0	0	0	0	0	0
DvNW	0	0	1.79e-05	0	2.80e-14	0	0	0	0
DvSS	0	0	0	0	0	0	0	0	5.83e-10
DvST	0	0	4.66e-06	0	0	0	0	0	0
DvWL	0	0	4.90e-06	0	0	0	0	0	0
NWvSS	0	0	0	0	0	0	0	0	0
NWvST	0	0	0	0	6.14e-07	0	0	0	0
NWvWL	0	0	0.344	0	0	0	0	0	0
SSvWL	0	0	0	0	0	0	0	0	0
STvSS	0	0	0	0	0	0	0	0	0
STvWL	0	0	0	0	0	0	0	0	0

Table 5: Results of the pair-wise statistical tests of the OO metrics analysed: AS refers to *Application Software*, D to *Documentation*, NW to *Non Web Libraries And Frameworks*, ST to *Software Tools*, SS to *System Software*, and WL to *Web Libraries And Frameworks*. Highlighted the p-values larger than α_B .

hypothesis), but only considering the subsets of projects in the *Application Software domain* and those in the *System Software domain*. Similarly, for the NOC attribute, and considering the *Non Web Libraries And Frameworks* and *Web Libraries And Frameworks* clusters.

Considering all the other tests, the resulting p-values allow for a general rejection of the null hypothesis: we can assume that the OO attributes come from different populations.

This is an interesting result: the clusters generated by a manual inspection
660 of the projects' characteristics result in pools of attributes that are structurally different from each other. As an example, the *Software Tools* (ST) category rejects all the null hypotheses, for all the selected OO attributes, when compared to any other category. This places the Software Tools cluster as a standalone category, with specific (and unique) characteristics. In order to visualise this result, we plotted in Figure 11 the distributions of three of the OO metrics analysed when comparing the projects from the *Software Tools* category and the category with most of the projects in that sample (i.e., the *Non Web Libraries And Frameworks* category).

670 7.2.1. LDA-informed categories for the manually categorised dataset

The manually categorised dataset was originally curated by the authors reading the documentation and assigning each project to one of six categories.

In order to test the same project sample with the LDA-informed clustering technique, the following steps were performed:

1. we extracted the *corpora* of the 500 Java projects from that sample;
2. we derived their *topics* using the LDA-based technique;
3. we assigned each project's topics to one of the 20 SourceForge *categories*;
4. we gathered the projects' OO data into their relative categories;
5. we re-run the KS pair-wise tests to check whether each OO attribute can
680 be considered from the same population.

Software Tools (ST) vs Non Web Libraries And Frameworks (NW)

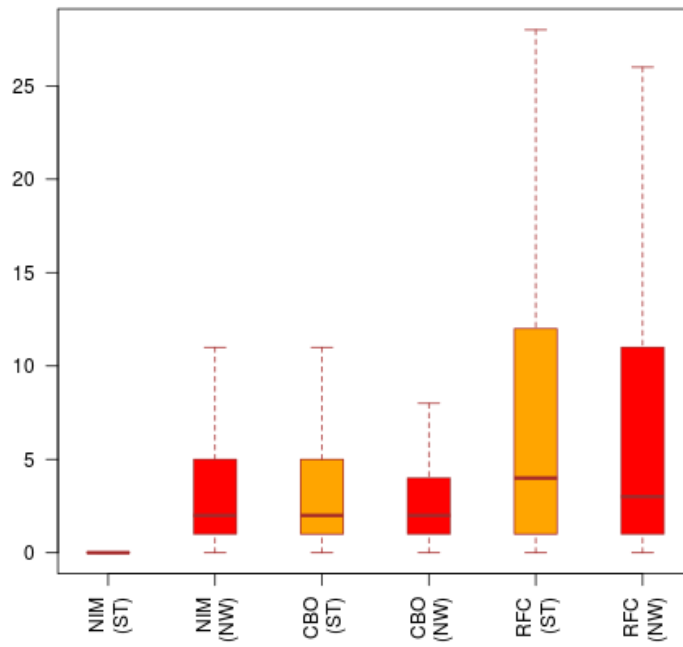


Figure 11: Comparison of OO metrics distributions for the *Software Tools* (ST) and *Non Web Libraries And Frameworks* (NW) categories.

Figure 12 shows how the original categories are composed by SourceForge categories found via the LDA-based technique. As visible, each ‘manual’ category contains more than one of the SourceForge categories, when analysing it via the LDA technique.

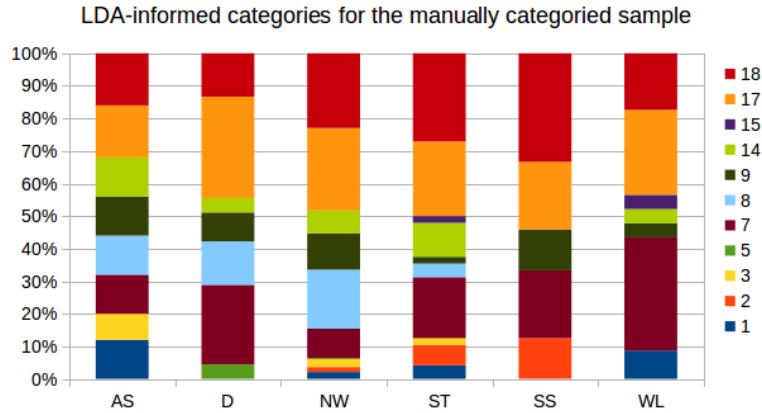


Figure 12: Original manually extracted categories broken down by the LDA-based technique. “AS” stands for ApplicationSoftware, “D” stands for Documentation, “NW” stands for Non-WebLibrariesAndFrameworks, “ST” stands for SoftwareTools, “SS” stands for SystemSoftware, “WL” stands for WebLibrariesAndFrameworks. 1 to 20 are the SourceForge categories, available in the list in Section 5.3.

From the gathered results it is possible to notice that some categories (e.g., *Religion*, *Social Sciences*, *Formats and Protocols*) are not detectable via the LDA-based techniques, whereas other categories (e.g., *Software Development*, *Mobile*) are most easily found. This means that also the SourceForge taxonomy is too coarse in some parts (e.g., *System*, *Software Development*, etc), and too fine in others (e.g., *Religion*).

690

7.3. LDA-based dataset – Results

Similarly to what has been done in Section 7.2, we clustered the projects in the categories proposed in SourceForge. The OO data was then extracted, per cluster, and one KS test executed per pair of clusters, and for all the metrics.

	IFANIN	CBO	NOC	NIM	NIV	WMC	RFC	DIT	LCOM
I v SD	0	0	0	0	0	0	0	0	0
I v Mob	0	0	0	0	0	0	0	0	0
I v MM	0	0	0	0	0	0	0	0	0
SD v Mob	3.29e-09	0	0.868	9.24e-4	1.9e-4	3.27e-05	0	5.55e-16	7.12e-4
MM v Mob	6.68e-06	0	0.999	1.36e-11	0	1.59e-09	2.62e-13	0.57	4.08e-10
MM v SD	0.015	8.42e-09	4.56e-05	0	0	0	0	0	0

Table 6: Results of the pair-wise statistical tests of the OO metrics analysed: I refers to Internet, SD to Software Development, Mob to Mobile and MM to Multimedia. Highlighted the p-values larger than α_B .

Table 6 summarises the tests by reporting the p-value of each test. Highlighted are the metrics for which the p-value does not allow us to reject the null hypothesis (*‘the two samples come from the same distribution’*) based on $\alpha = 0.05$ and a Bonferroni correction at $\alpha_B = 0.0083333$.

As shown in Table 6, and similarly to the results in Table 5, the NOC
700 attribute does not always allow for a clear differentiation between clusters. In two cases, we could not reject the null hypothesis: the values of the NOC attribute can be considered as drawn from the same population, at least in the cases of the SD and Mob comparison, and the MM and Mob comparison. For the other comparisons, we can reject all the other null hypotheses: we can conclude that OO attributes are extracted from different populations, when considering the clusters formed by an LDA-informed approach.

It is worth noting that for both Categorized dataset, and the LDA-informed dataset, the DIT and NOC attributes are less sensitive to identify differences between samples. This is partly due to the distribution of values for the two
710 OO attributes in the sample: we observed an average value of DIT at 1.77 (median = 1), and an average value of NOC at 1.17 (median = 0). These low values reflect a similar characteristic of the structure of Java software (i.e.,

inheritance): shared design principles suggest that the DIT attribute should be kept low. This is true for two main reasons:

- The deeper a class is in the hierarchy, the greater the total number of methods it is likely to inherit [89], making its behaviour less predictable [90].
- Khalid *et al.* [91] state that “DIT is directly proportional to complexity” (i.e., an increased DIT will lead to higher maintenance efforts).

On the other hand, the NOC attribute should also be kept low: a large CBO increases the complexity of the system, and it adversely affects other quality factors, such as maintainability, testability and reusability [92].

7.3.1. LDA-based results for the dataset in [14]

In order to cross-examine the three datasets with multiple techniques, we ran the pair-wise statistical tests for the categories found via the LDA-based technique, this time applied to the dataset from [14]. As before, we assigned a standard threshold value $\alpha = 0.05$ for the sensitivity of the test; we performed 36 parallel tests, so the Bonferroni correction resulted in $\alpha_B = 0.00139$. We gathered the results of the tests in Table 7.

Similarly to the manual clusters, presented in Table 5, the LDA-based technique applied to this 500-project sample does not always clearly reject the null hypothesis regarding the NOC attribute. We could always reject the null hypothesis for the DIT attribute, that was highlighted in Table 5. On the other hand, we could not reject the null hypothesis, for a couple of pair-wise comparisons, regarding the IFANIN attribute. Considering the SourceForge categories, the ID=3 (Desktop Environment) contains projects that overall are less clearly differentiated from other categories (i.e., there is a larger likelihood not to reject the null hypothesis).

8. Discussion

In this section we add further insights as part of our discussion: in light of our evidence, we studied past research works, specifically focused on the OO

	IFANIN	CBO	NOC	NIM	NIV	WMC	RFC	DIT	LCOM
1 V 2	0	0	1.23e-10	0	0	0	0	0	0
1 V 3	0	0	0	9.31e-12	1.61e-07	1.47e-13	0	0	0.05
1 V 7	0	0	0	2.46e-10	0	9.23e-14	0	0	1.45e-07
1 V 8	0	0	2.22e-16	0	0	0	0	0	0
1 V 9	0	0	0	0	0	0	0	0	0
1 V 14	0	0	0	0	3.99e-11	0	0	0	0
1 V 17	0	0	0	0	0	0	0	0	0
1 V 18	0	0	0	0	5.40e-12	0	0	0	0
2 V 14	0	0	0	0	6.48e-09	0	0	0	0.0001
2 V 17	0	0	0	0	0	0	0	0	0
2 V 18	0.06	0	8.82e-13	0	3.78e-05	0	0	0	1.57e-11
2 V 3	0	0	5.04e-09	0	0	0	0	5.42e-06	0
2 V 7	5.45e-12	0	1.11e-16	0	0	0	0	0	0
2 V 8	0	0	1.68e-06	0	0	6.61e-05	0	0	0
2 V 9	0	0	0	0	0	0	0	0	0
3 V 7	0	0	0.003	1.27e-06	8.40e-10	1.79e-09	0	0	1.60e-06
3 V 8	0	0	0.005	0	0	0	0	0	0
3 V 9	0	0	0.005	0	0	0	0	0	0
3 V 14	0.26	0	6.28e-06	0	7.77e-16	0	9.66e-15	0	0
3 V 17	4.84e-10	0	0.99	2.33e-10	1.18e-06	0	0	0	0
3 V 18	0	0	5.81e-06	0	0	0	0	6.46e-14	0
7 V 8	0	0	0.79	0	0	0	0	0	0
7 V 9	0	0	0	0	0	0	0	0	0
7 V 14	0	0	0	0	0	0	0	1.11e-16	0
7 V 17	0	0	0	0	9.74e-08	0	0	0	0
7 V 18	1.38e-12	0	6.69e-10	0	0	0	0	0	0
8 V 9	0	0	0	0	0	0	0	0	0
8 V 14	0	0	0	0.042	0	0.0015	0	0	0
8 V 17	0	0	1.84e-06	0	0	0	0	0	0
8 V 18	0	0	0.24	1.73e-10	0	0	9.75e-14	0	0
9 V 14	0	0	0.006	0	0	0	0	0	0
9 V 17	0	0	0	0	0	0	0	0	0
9 V 18	0	0	0	0	0	0	0	0	0
14 V 17	0	5.15e-13	0	0	0	0	0	0	0
14 V 18	0	0	0	5.15e-05	2.69e-09	2.60e-10	0	0	0.0003
17 V 18	0	0	0	0	0	0	0	0	0

Table 7: Results of the pair-wise statistical tests using the LDA-informed categories on the manually clustered projects. The first column categories are taken from the list in Section 5.3. Highlighted the p-values larger than α_B .

metrics that we used, and we reflected on the importance of clustering.

In a prior study [93], we collected empirical results showing that projects from the same domain exhibit common structural properties in terms of the C&K metrics. In this work, we have expanded on those results: the results that we have gathered above indicate that the projects clustered around the domains (i.e., 2 of the 3 clustering techniques presented above) show indeed a difference in the structural metrics. For interested stakeholders, this can imply that the structure of a software system (and its [development and future](#)

750 *maintenance*) depends on domain-based factors common to projects in the same domain. For example, projects from different domains making use of exception handling differently [60, 62].

The discussion that we present here takes into consideration the correlations between the OO metrics that we utilised above. The work reported by [94] has already shown some correlation between pairs of metrics from the C&K suite, for example, CBO and RFC, and RFC and LCOM. We want to know whether these correlations change sensibly, when considering specific clusters. If indeed there were a correlation in all the clusters analysed above, it would suggest an increased probability of falsely rejecting a null hypothesis within the clusters
760 shown in Figure 10.

In the analysis below, we report on the correlation study between pairs of OO metrics, and when clustered by application domain. The value of the correlation coefficient lies in the range $[-1; 1]$, where -1 indicates a strong negative correlation and 1 indicates a strong positive correlation. We adapt the categorisation for correlation coefficients used in [95] ($[0 - 0.1]$ to be *insignificant*, $[0.1 - 0.3]$ *low*, $[0.3 - 0.5]$ *moderate*, $[0.5 - 0.7]$ *large*, $[0.7 - 0.9]$ *very large*, and $[0.9 - 1]$ *almost perfect*) if the rank correlation coefficient proves to be statistically significant at the $\alpha = 0.01$ level.

Table 8 shows the correlations among pairs of OO metrics, when considering
770 the projects in the domain-driven clusters identified by the LDA technique. It becomes clear that the metrics show collinearity, but depending on the cluster considered, this collinearity could be stronger or weaker. An example of this is between the RFC and WMC attributes: for the projects in the *Internet* cluster, this association has a medium (M) strength; the association becomes large (L) when considering the projects in the *Mobile* cluster; for the projects in the *MultiMedia* category, the association becomes almost perfect (AP), thus being larger than 0.9; when considering the projects in the *Software Development* cluster, on the other hand, such collinearity becomes small (s), hence isolating these projects, and their characteristics, from the rest of the sample.

780 Considering the definition of the RFC and WMC attributes, a stronger cor-

Internet								
	IFANIN	CBO	NOC	NIM	NIV	WMC	RFC	DIT
CBO	-i							
NOC	-i	i						
NIM	s	M	i					
NIV	i	M	i	M				
WMC	i	M	i	AP	M			
RFC	-i	M	i	M	s	M		
DIT	-s	i	-i	i	-i	i	L	
LCOM	-i	s	i	M	M	M	M	s

Mobile								
	IFANIN	CBO	NOC	NIM	NIV	WMC	RFC	DIT
CBO	i							
NOC	-i	i						
NIM	s	M	s					
NIV	i	M	i	XL				
WMC	i	M	s	AP	XL			
RFC	-i	s	s	L	M	L		
DIT	-s	s	s	s	s	s	L	
LCOM	-i	M	s	L	L	L	M	s

Multimedia								
	IFANIN	CBO	NOC	NIM	NIV	WMC	RFC	DIT
CBO	i							
NOC	i	i						
NIM	i	s	i					
NIV	i	M	i	M				
WMC	-i	L	i	s	i			
RFC	-i	L	i	s	i	AP		
DIT	-M	s	i	s	i	i	s	
LCOM	i	s	i	M	L	i	i	s

Software Development								
	IFANIN	CBO	NOC	NIM	NIV	WMC	RFC	DIT
CBO	i							
NOC	-i	i						
NIM	i	M	i					
NIV	i	i	i	i				
WMC	i	M	i	AP	i			
RFC	-i	s	i	s	i	s		
DIT	-s	i	-i	i	i	i	M	
LCOM	-i	M	i	M	i	M	s	s

Table 8: Correlation between OO pairs, after grouping projects within domain-driven clusters. *i* stands for insignificant correlation; *M* for medium; *L* for large; *XL* for very large; and *AP* for almost perfect.

relation implies a larger complexity of the code: when the number of methods (i.e., WMC) grows in a Java class, the response for that class (i.e., the RFC) also grows. This is also an indication that more testing will be needed for that class: the projects in the MultiMedia category show a different behaviour to those belonging to the Software Development category.

Based on these results, it is possible to summarise our correlation findings as follows:

The correlation among OO metrics can be extremely sensitive to application domains

Thus, according to such a finding, evaluating the quality of a system becomes also dependent on what type of domain a project belongs to. For instance, the metrics one should consider to analyse gaming software should be different from those used to assess the quality of security software. For instance, the latter is mainly characterized by incremental contributions aiming at fixing already existing functionalities and making them more stable and secure. Such a particular attention to stability and security aspects can be less peculiar for gaming software and consequently OO metrics, e.g., LOCs are less appropriate for assessing the quality of security software or in general of mission critical software systems.

8.1. Clusters and their maturity

The maturity of the considered application domain (as cluster) is another factor that can interfere with the analysis performed by means of OO metrics. Emerging domains are characterized by a plethora of new applications even though while the domain becomes more “stable” and mature, only those applications that managed to create a community (and thus that are actually used and maintained) eventually survive. Thus, the population of the resulting application domain is characterized by OO metrics that might be different from those of the initial population:

Emerging new application domains might show a larger variability in their structural metrics, in comparison to established domains.

Finally, some domains consist of reusable libraries and frameworks instead of ready-to-use applications. For instance, in the domain of automation testing, we can have JUnit, Selenium, Jasmin, etc. Thus, the domain consists of a population that overall shares only the same goal, but it is likely to be characterized by an insignificant correlation in terms of OO metrics.

This research is part of a wider context, and it requires further multidisciplinary investigation: the similarity of software systems should follow the same approach of compiling a biological taxonomy, where systems (or parts of) are given a place (or rank) in a hierarchy. Lower levels share (OO-related) attributes with higher levels in the same branch, whereas different branches of the taxonomy show the highest degree of dissimilarity. Such a taxonomy could have a massive impact in how practitioners and researchers work and develop software systems: ad-hoc techniques and tools would be needed and tailored to domain-specific constraints. This has already started to emerge for the software engineering techniques around gaming development [96]: we envisage that a similar approach would be needed when new domains become established, and their boundaries are clearer.

9. Threats to validity

This section distinguishes between threats to external, internal, and construct validity as follows.

9.1. External validity

Threats to external validity refer to the extent to which the results of our study can be generalised. The projects we have analysed come from publicly available repositories and the used clustering techniques also come from already validated works. However, we cannot claim that the resulting conclusion

concerning the considered null hypotheses is generalisable, even though the performed experiments provide us with an acceptable confidence about the general validity of the reached conclusions.

9.2. *Internal validity*

Threats to internal validity concern any confounding factor that could influence our findings. We attempted to avoid any bias in the building of the ecosystems. To this end, we applied CrossSim and the LDA-based approach 840 by means of the corresponding supporting tools without any manual intervention. However, the tools we have used and implemented could be defective. To counter this threat, we have run several manual assessments and counter-checks. Concerning the approach in [14] we used the categories as they are presented in the original work without any change.

9.3. *Construct validity*

It pertains to any factor that can compromise the validity of inferences that observations or measurement tools actually represent or measure the construct being investigated. A potential threat to construct validity is related to the 850 size of the analysed data. However, the datasets that have been considered for the experiments come from the original works proposing CrossSim, the LDA-based approach, and the categorization in [14]. Thus, we built on those datasets that were considered to be big enough for experimenting and validating such approaches.

The second threat to construct validity is based on the fact that the three approaches proposed use three different samples of data. We could not use the same sample because the approach described in Section 5.2 above is not fully reproducible. On the other hand, the LDA-based approach was replicated for the CrossSim sample, and we concluded that the clusters grouped by CrossSim 860 represent the type of external libraries they use, while they could be implementing very different features. They could be based on very different domains, but still use a similar set of external libraries.

10. Conclusion

This paper presented three techniques to cluster software projects, and evaluated how the clusters differ from each other, when comparing their internal characteristics. The null hypothesis that we attempted to reject is that all the clusters come from the same population. As metrics of assessment, we adopted a suite of well known OO attributes.

As the first technique we adopted the CrossSim algorithm [11], that draws
870 a similarity between projects based on their usage of external libraries. Interestingly enough, with this technique we could not reject the null hypothesis: projects might still be *internally* similar even if the CrossSim algorithms detect a large distance between them.

The second technique was based on a subjective classification, based on the description of a software project. With this technique we found a strong foundation for rejecting the null hypothesis, although the steps of the clustering process are not easily reproducible.

The third technique was also based on categories, but it automatically extracted topics from a software systems, before assigning them to categories. Also
880 in this case, we could reject the null hypothesis for most of the OO metrics.

The implications of these findings can be profound: software projects can manifest different characteristics, based on the domain or cluster that they belong to. Empirical findings might need readjustment, and tailoring to one or another cluster or domain, to fit with other findings in the same cluster.

Acknowledgements

The research described in this paper has been carried out as part of the CROSSMINER Project, which has received funding from the European Union's Horizon 2020 Research and Innovation Programme under Grant 732223.

References

- 890 [1] D. Yang, P. Martins, V. Saini, C. Lopes, Stack overflow in github: any snippets there?, in: Proceedings of the 14th International Conference on Mining Software Repositories, IEEE Press, 2017, pp. 280–290.
- [2] H. Osman, A. Chiş, C. Corrodi, M. Ghafari, O. Nierstrasz, Exception evolution in long-lived java systems, in: Proceedings of the 14th International Conference on Mining Software Repositories, IEEE Press, 2017, pp. 302–311.
- [3] R. Kikas, M. Dumas, D. Pfahl, Using dynamic and contextual features to predict issue lifetime in github projects, in: Proceedings of the 13th International Conference on Mining Software Repositories, ACM, 2016, pp. 291–302.
- 900 [4] M. Linares-Vásquez, A. Holtzhauer, C. Bernal-Cárdenas, D. Poshyvanyk, Revisiting android reuse studies in the context of code obfuscation and library usages, in: Proceedings of the 11th Working Conference on Mining Software Repositories, ACM, 2014, pp. 242–251.
- [5] G. Bulet, A. Hindle, An empirical study of end-user programmers in the computer music community, in: Proceedings of the 12th Working Conference on Mining Software Repositories, IEEE Press, 2015, pp. 292–302.
- [6] J. B. Schafer, D. Frankowski, J. Herlocker, S. Sen, The adaptive web, Springer-Verlag, Berlin, Heidelberg, 2007, Ch. Collaborative Filtering Recommender Systems, pp. 291–324.
URL <http://dl.acm.org/citation.cfm?id=1768197.1768208>
- 910 [7] Y. Zhang, D. Lo, P. S. Kochhar, X. Xia, Q. Li, J. Sun, Detecting similar repositories on github, 2017 IEEE 24th International Conference on Software Analysis, Evolution and Reengineering (SANER) 00 (2017) 13–23. doi:[doi.1109/SANER.2017.7884605](https://doi.org/10.1109/SANER.2017.7884605).
- [8] D. Spinellis, C. Szyperski, How is open source affecting software development?, IEEE Software 21 (1) (2004) 28–33. doi:[10.1109/MS.2004.1259204](https://doi.org/10.1109/MS.2004.1259204).
- 920 [9] P. T. Nguyen, J. Di Rocco, D. Di Ruscio, L. Ochoa, T. Dagueule, M. Di Penta, FOCUS: A Recommender System for Mining API Function Calls and Usage Patterns, in: Proceedings of the 41st International Conference on Software Engineering, ICSE '19, IEEE Press, Piscataway, NJ, USA, 2019, pp. 1050–1060. doi:[10.1109/ICSE.2019.00109](https://doi.org/10.1109/ICSE.2019.00109).
URL <https://doi.org/10.1109/ICSE.2019.00109>

- [10] C. McMillan, M. Grechanik, D. Poshyvanyk, Detecting similar software applications, in: Proceedings of the 34th International Conference on Software Engineering, ICSE '12, IEEE Press, Piscataway, NJ, USA, 2012, pp. 364–374.
URL <http://dl.acm.org/citation.cfm?id=2337223.2337267>
- [11] P. T. Nguyen, J. Di Rocco, R. Rubei, D. Di Ruscio, CrossSim: Exploiting Mutual Relationships to Detect Similar OSS Projects, in: 2018 44th Euromicro Conference on Software Engineering and Advanced Applications (SEAA), 2018, pp. 388–395. doi:10.1109/SEAA.2018.00069.
- [12] P. T. Nguyen, J. Di Rocco, R. Rubei, D. Di Ruscio, An Automated Approach to Assess the Similarity of GitHub Repositories, Software Quality Journal To appear.
- 930 [13] H. Borges, A. Hora, M. T. Valente, Understanding the factors that impact the popularity of github repositories, in: 2016 IEEE International Conference on Software Maintenance and Evolution (ICSME), IEEE, 2016, pp. 334–344.
- [14] H. Borges, M. T. Valente, Whats in a github star? understanding repository starring practices in a social coding platform, Journal of Systems and Software 146 (2018) 112–129.
- [15] S. R. Chidamber, C. F. Kemerer, Towards a metrics suite for object oriented design, SIGPLAN Not. 26 (11) (1991) 197–211. doi:10.1145/118014.117970.
URL <http://doi.acm.org/10.1145/118014.117970>
- 940 [16] G. Destefanis, S. Counsell, G. Concas, R. Tonelli, Software metrics in agile software: An empirical study, in: International Conference on Agile Software Development, Springer, 2014, pp. 157–170.
- [17] R. Amelia, H. Washizaki, Y. Fukazawa, K. Oshima, R. Mibe, R. Tsuchiya, An empirical study on relationship between requirement traceability links and bugs., JSW 12 (5) (2017) 315–325.
- [18] E. A. AlOmar, M. W. Mkaouer, A. Ouni, M. Kessentini, On the impact of refactoring on the relationship between quality attributes and design metrics, in: 2019 ACM/IEEE International Symposium on Empirical Software Engineering and Measurement (ESEM), IEEE, 2019, pp. 1–11.
- 950 [19] J. A. Reshi, S. Singh, Investigating the role of code smells in preventive maintenance, Journal of Information Technology Management 10 (4) (2019) 41–63.
- [20] C. D. Manning, P. Raghavan, H. Schütze, Introduction to Information Retrieval, Cambridge University Press, New York, NY, USA, 2008.

- [21] M. Steinbach, G. Karypis, V. Kumar, A comparison of document clustering techniques, in: 6th ACM SIGKDD, World Text Mining Conference, 2000.
- [22] A. Kuhn, S. Ducasse, T. Girba, Enriching reverse engineering with semantic clustering, in: 12th Working Conference on Reverse Engineering (WCRE'05), IEEE, 2005, pp. 10–pp.
- [23] O. Maqbool, H. Babri, Hierarchical clustering for software architecture recovery, IEEE Transactions on Software Engineering 33 (11) (2007) 759–780.
- [24] W. S. Evans, C. W. Fraser, F. Ma, Clone detection via structural abstraction, Software Quality Journal 17 (4) (2009) 309–330. doi:10.1007/s11219-009-9074-y.
URL <https://doi.org/10.1007/s11219-009-9074-y>
- [25] C. Ragkhitwetsagul, J. Krinke, D. Clark, A comparison of code similarity analysers, Empirical Software Engineering 23 (4) (2018) 2464–2519. doi:10.1007/s10664-017-9564-7.
URL <https://doi.org/10.1007/s10664-017-9564-7>
- [26] C. Ragkhitwetsagul, J. Krinke, B. Marnette, A picture is worth a thousand words: Code clone detection based on image similarity, in: 2018 IEEE 12th International Workshop on Software Clones (IWSC), 2018, pp. 44–50. doi:10.1109/IWSC.2018.8327318.
- [27] R. Tiarks, R. Koschke, R. Falke, An extended assessment of type-3 clones as detected by state-of-the-art tools, Software Quality Journal 19 (2) (2011) 295–331. doi:10.1007/s11219-010-9115-6.
URL <https://doi.org/10.1007/s11219-010-9115-6>
- [28] A. M. Leitão, Detection of redundant code using r2d2, Software Quality Journal 12 (4) (2004) 361–382. doi:10.1023/B:SQJ0.0000039793.31052.72.
URL <https://doi.org/10.1023/B:SQJ0.0000039793.31052.72>
- [29] A. E. V. B. Coutinho, E. G. Cartaxo, P. D. de Lima Machado, Analysis of distance functions for similarity-based test suite reduction in the context of model-based testing, Software Quality Journal 24 (2014) 407–445.
- [30] S. U. R. Khan, S. P. Lee, R. W. Ahmad, A. Akhunzada, V. Chang, A survey on test suite reduction frameworks and tools, Int. J. Inf. Manag. 36 (6) (2016) 963–975. doi:10.1016/j.ijinfomgt.2016.05.025.
URL <https://doi.org/10.1016/j.ijinfomgt.2016.05.025>
- [31] A. Walenstein, M. El-Ramly, J. R. Cordy, W. S. Evans, K. Mahdavi, M. Pizka, G. Ramalingam, J. W. von Gudenberg, Similarity in programs, in: Duplication, Redundancy, and Similarity in Software, 23.07. - 26.07.2006, 2006.
URL <http://drops.dagstuhl.de/opus/volltexte/2007/968>

- [32] U. von Luxburg, R. C. Williamson, I. Guyon, Clustering: Science or art?, in: I. Guyon, G. Dror, V. Lemaire, G. Taylor, D. Silver (Eds.), Proceedings of ICML Workshop on Un-supervised and Transfer Learning, Vol. 27 of Proceedings of Machine Learning Research, 990 PMLR, Bellevue, Washington, USA, 2012, pp. 65–79.
URL <http://proceedings.mlr.press/v27/luxburg12a.html>
- [33] C. Rodolfo, M. Murru, F. Catalli, G. Falcone, Real time forecasts through an earthquake clustering model constrained by the rate-and-state constitutive law: Comparison with a purely stochastic etas model, *Seismological Research Letters* 78 (2007) 49–56. doi: 10.1785/gssrl.78.1.49.
- [34] P. Berkhin, A survey of clustering data mining techniques, in: J. Kogan, C. Nicholas, M. Teboulle (Eds.), *Grouping Multidimensional Data*, Springer, 2006, pp. 25–71. doi: 10.1007/3-540-28349-8_2.
1000 URL http://dx.doi.org/10.1007/3-540-28349-8_2
- [35] A. K. Jain, M. N. Murty, P. J. Flynn, Data clustering: a review, *ACM computing surveys (CSUR)* 31 (3) (1999) 264–323.
- [36] D. Xu, Y. Tian, A comprehensive survey of clustering algorithms, *Annals of Data Science* 2 (2) (2015) 165–193. doi:10.1007/s40745-015-0040-1.
URL <https://doi.org/10.1007/s40745-015-0040-1>
- [37] L. Kaufman, P. J. Rousseeuw, *Finding groups in data : an introduction to cluster analysis*.
URL <http://opac.inria.fr/record=b1087461>
- [38] R. T. Ng, J. Han, Clarans: A method for clustering objects for spatial data mining, 1010 *IEEE Trans. on Knowl. and Data Eng.* 14 (5) (2002) 1003–1016. doi:10.1109/TKDE.2002.1033770.
URL <http://dx.doi.org/10.1109/TKDE.2002.1033770>
- [39] O. Maimon, L. Rokach (Eds.), *Clustering Methods*, Springer US, Boston, MA, 2005, pp. 321–352. doi:10.1007/0-387-25465-X_15.
URL http://dx.doi.org/10.1007/0-387-25465-X_15
- [40] T. Zhang, R. Ramakrishnan, M. Livny, Birch: A new data clustering algorithm and its applications, *Data Min. Knowl. Discov.* 1 (2) (1997) 141–182. doi:10.1023/A:1009783824328.
URL <https://doi.org/10.1023/A:1009783824328>
- 1020 [41] S. Guha, R. Rastogi, K. Shim, Cure: An efficient clustering algorithm for large databases, *SIGMOD Rec.* 27 (2) (1998) 73–84. doi:10.1145/276305.276312.
URL <http://doi.acm.org/10.1145/276305.276312>

- [42] Rock: A robust clustering algorithm for categorical attributes, in: Proceedings of the 15th International Conference on Data Engineering, ICDE '99, IEEE Computer Society, Washington, DC, USA, 1999, pp. 512–. URL <http://dl.acm.org/citation.cfm?id=846218.847264>
- [43] G. Karypis, E.-H. S. Han, V. Kumar, Chameleon: Hierarchical clustering using dynamic modeling, *Computer* 32 (8) (1999) 68–75. doi:10.1109/2.781637. URL <http://dx.doi.org/10.1109/2.781637>
- 1030 [44] B. Sarwar, G. Karypis, J. Konstan, J. Riedl, Item-based collaborative filtering recommendation algorithms, in: Proceedings of the 10th International Conference on World Wide Web, WWW '01, ACM, New York, NY, USA, 2001, pp. 285–295. doi:10.1145/371920.372071. URL <http://doi.acm.org/10.1145/371920.372071>
- [45] N. Chen, S. C. Hoi, S. Li, X. Xiao, Simapp: A framework for detecting similar mobile applications by online kernel learning, in: Proceedings of the Eighth ACM International Conference on Web Search and Data Mining, WSDM '15, ACM, New York, NY, USA, 2015, pp. 305–314. doi:10.1145/2684822.2685305. URL <http://doi.acm.org/10.1145/2684822.2685305>
- 1040 [46] M. Nagappan, T. Zimmermann, C. Bird, Diversity in software engineering research, in: Proceedings of the 2013 9th Joint Meeting on Foundations of Software Engineering, ACM, 2013, pp. 466–476.
- [47] C. McMillan, M. Grechanik, D. Poshyvanyk, Detecting similar software applications, in: Proceedings of the 34th International Conference on Software Engineering, IEEE Press, 2012, pp. 364–374.
- [48] S. Kawaguchi, P. K. Garg, M. Matsushita, K. Inoue, Mudablue: An automatic categorization system for open source repositories, *Journal of Systems and Software* 79 (7) (2006) 939–953.
- 1050 [49] Y. Zhang, D. Lo, P. S. Kochhar, X. Xia, Q. Li, J. Sun, Detecting similar repositories on github, in: Software Analysis, Evolution and Reengineering (SANER), 2017 IEEE 24th International Conference on, IEEE, 2017, pp. 13–23.
- [50] K. W. Nafi, B. Roy, C. K. Roy, K. A. Schneider, Crolsim: Cross language software similarity detector using api documentation, in: 2018 IEEE 18th International Working Conference on Source Code Analysis and Manipulation (SCAM), IEEE, 2018, pp. 139–148.

- [51] P. T. Nguyen, J. D. Rocco, D. D. Ruscio, M. D. Penta, CrossRec: Supporting software developers by recommending third-party libraries, *Journal of Systems and Software*[doi: https://doi.org/10.1016/j.jss.2019.110460](https://doi.org/10.1016/j.jss.2019.110460).
URL <http://www.sciencedirect.com/science/article/pii/S0164121219302341>
- 1060 [52] M. Wermelinger, Y. Yu, An architectural evolution dataset, in: *Mining Software Repositories (MSR)*, 2015 IEEE/ACM 12th Working Conference on, IEEE, 2015, pp. 502–505.
- [53] D. M. German, M. Di Penta, Y.-G. Gueheneuc, G. Antoniol, Code siblings: Technical and legal implications of copying code between applications, in: *Mining Software Repositories*, 2009. MSR'09. 6th IEEE International Working Conference on, IEEE, 2009, pp. 81–90.
- [54] A. Deshpande, D. Riehle, The total growth of open source, in: *IFIP International Conference on Open Source Systems*, Springer, 2008, pp. 197–209.
- [55] S. Karus, H. Gall, A study of language usage evolution in open source software, in: *Proceedings of the 8th Working Conference on Mining Software Repositories*, ACM, 2011, pp. 13–22.
- 1070 [56] K. Tian, M. Reville, D. Poshyvanyk, Using Latent Dirichlet Allocation for automatic categorization of software, in: *Mining Software Repositories*, 2009. MSR'09. 6th IEEE International Working Conference on, IEEE, 2009, pp. 163–166.
- [57] O. Callaú, R. Robbes, É. Tanter, D. Röthlisberger, How (and why) developers use the dynamic features of programming languages: the case of smalltalk, *Empirical Software Engineering* 18 (6) (2013) 1156–1194.
- [58] M. Linares-Vásquez, G. Bavota, C. Bernal-Cárdenas, R. Oliveto, M. Di Penta, D. Poshyvanyk, Mining energy-greedy api usage patterns in android apps: an empirical study, in: *Proceedings of the 11th Working Conference on Mining Software Repositories*, ACM, 2014, pp. 2–11.
- 1080 [59] L. Bao, D. Lo, X. Xia, X. Wang, C. Tian, How android app developers manage power consumption?: An empirical study by mining power management commits, in: *Proceedings of the 13th International Conference on Mining Software Repositories*, ACM, 2016, pp. 37–48.
- [60] S. Nakshatri, M. Hegde, S. Thandra, Analysis of exception handling patterns in java projects: An empirical study, in: *Proceedings of the 13th International Conference on Mining Software Repositories*, ACM, 2016, pp. 500–503.
- [61] M. Asaduzzaman, M. Ahasanuzzaman, C. K. Roy, K. A. Schneider, How developers use exception handling in java?, in: *Proceedings of the 13th International Conference on Mining Software Repositories*, ACM, 2016, pp. 516–519.

- 1090 [62] B. Cabral, P. Marques, Exception handling: A field study in java and. net, in: European Conference on Object-Oriented Programming, Springer, 2007, pp. 151–175.
- [63] T. L. Alves, C. Ypma, J. Visser, Deriving metric thresholds from benchmark data, in: 2010 IEEE International Conference on Software Maintenance, IEEE, 2010, pp. 1–10.
- [64] K. A. Ferreira, M. A. Bigonha, R. S. Bigonha, L. F. Mendes, H. C. Almeida, Identifying thresholds for object-oriented software metrics, *Journal of Systems and Software* 85 (2) (2012) 244–257.
- [65] P. Oliveira, M. T. Valente, F. P. Lima, Extracting relative thresholds for source code metrics, in: 2014 Software Evolution Week-IEEE Conference on Software Maintenance, Reengineering, and Reverse Engineering (CSMR-WCRE), IEEE, 2014, pp. 254–263.
- 1100 [66] L. B. L. De Souza, M. D. A. Maia, Do software categories impact coupling metrics?, in: Proceedings of the 10th working conference on mining software repositories, IEEE Press, 2013, pp. 217–220.
- [67] A. Mori, G. Vale, M. Vigiato, J. Oliveira, E. Figueiredo, E. Cirilo, P. Jamshidi, C. Kastner, Evaluating domain-specific metric thresholds: an empirical study, in: 2018 IEEE/ACM International Conference on Technical Debt (TechDebt), IEEE, 2018, pp. 41–50.
- [68] P. T. Nguyen, J. Di Rocco, R. Rubel, D. Di Ruscio, CrossSim tool and evaluation data, <https://doi.org/10.5281/zenodo.1252866> (2018).
- [69] R. H. Lopes, Kolmogorov-smirnov test, *International encyclopedia of statistical science* (2011) 718–720.
- 1110 [70] E. W. Weisstein, Bonferroni correction.
- [71] M. Lorenz, J. Kidd, *Object-oriented software metrics*, Vol. 131, Prentice Hall Englewood Cliffs, 1994.
- [72] C. Bizer, T. Heath, T. Berners-Lee, Linked data - the story so far, *Int. J. Semantic Web Inf. Syst.* 5 (3) (2009) 1–22.
- [73] P. T. Nguyen, P. Tomeo, T. Di Noia, E. Di Sciascio, An evaluation of simrank and personalized pagerank to build a recommender system for the web of data, in: Proceedings of the 24th International Conference on World Wide Web, WWW '15 Companion, ACM, New York, NY, USA, 2015, pp. 1477–1482. doi:10.1145/2740908.2742141.
- 1120 URL <http://doi.acm.org/10.1145/2740908.2742141>

- [74] V. D. Blondel, A. Gajardo, M. Heymans, P. Senellart, P. V. Dooren, A measure of similarity between graph vertices: Applications to synonym extraction and web searching, *SIAM Rev.* 46 (4) (2004) 647–666. doi:10.1137/S0036144502415960.
URL <http://dx.doi.org/10.1137/S0036144502415960>
- [75] G. Jeh, J. Widom, Simrank: A measure of structural-context similarity, in: *Proceedings of the Eighth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '02*, ACM, New York, NY, USA, 2002, pp. 538–543. doi:10.1145/775047.775126.
URL <http://doi.acm.org/10.1145/775047.775126>
- 1130 [76] F. Thung, D. Lo, J. Lawall, Automated library recommendation, in: *2013 20th Working Conference on Reverse Engineering (WCRE)*, 2013, pp. 182–191. doi:10.1109/WCRE.2013.6671293.
- [77] R. L. Glass, I. Vessey, Contemporary application-domain taxonomies, *Software*, IEEE 12 (1995) 63 – 76. doi:10.1109/52.391837.
- [78] H. Kagdi, M. Gethers, D. Poshvanyk, Integrating conceptual and logical couplings for change impact analysis in software, *Empirical Software Engineering* 18 (5) (2013) 933–969.
- [79] A. Marcus, A. Sergeyev, V. Rajlich, J. Maletic, et al., An information retrieval approach to concept location in source code, in: *Reverse Engineering, 2004. Proceedings. 11th Working Conference on*, IEEE, 2004, pp. 214–223.
1140
- [80] D. M. Blei, A. Y. Ng, M. I. Jordan, Latent Dirichlet Allocation, *Journal of Machine Learning Research* 3 (Jan) (2003) 993–1022.
- [81] J. Ramos, Using tf-idf to determine word relevance in document queries (1999).
- [82] A. Sharma, F. Thung, P. S. Kochhar, A. Sulistya, D. Lo, Cataloging github repositories, in: *Proceedings of the 21st International Conference on Evaluation and Assessment in Software Engineering, EASE'17*, ACM, New York, NY, USA, 2017, pp. 314–319. doi:10.1145/3084226.3084287.
URL <http://doi.acm.org.univaq.clas.cineca.it/10.1145/3084226.3084287>
- [83] J. Jiang, D. Lo, J. He, X. Xia, P. S. Kochhar, L. Zhang, Why and How Developers Fork What from Whom in GitHub, *Empirical Software Engineering* 22 (1) (2017) 547–578.
1150 doi:10.1007/s10664-016-9436-6.
URL <https://doi.org/10.1007/s10664-016-9436-6>

- [84] H. Borges, A. C. Hora, M. T. Valente, Understanding the Factors That Impact the Popularity of GitHub Repositories, in: 2016 IEEE International Conference on Software Maintenance and Evolution, ICSME 2016, Raleigh, NC, USA, October 2-7, 2016, 2016, pp. 334–344. doi:10.1109/ICSME.2016.31.
URL <https://doi.org/10.1109/ICSME.2016.31>
- [85] H. Borges, A. Hora, M. T. Valente, Predicting the popularity of github repositories, in: Proceedings of the The 12th International Conference on Predictive Models and Data Analytics in Software Engineering, PROMISE 2016, ACM, New York, NY, USA, 2016, pp. 9:1–9:10. doi:10.1145/2972958.2972966.
URL <http://doi.acm.org/10.1145/2972958.2972966>
- [86] K. Crowston, H. Annabi, J. Howison, Defining open source software project success, ICIS 2003 Proceedings (2003) 28.
- [87] G. Von Krogh, S. Spaeth, K. R. Lakhani, Community, joining, and specialization in open source software innovation: a case study, Research policy 32 (7) (2003) 1217–1241.
- [88] M. Allamanis, C. Sutton, Mining Source Code Repositories at Massive Scale using Language Modeling, in: The 10th Working Conference on Mining Software Repositories, IEEE, 2013, pp. 207–216.
- [89] S. R. Chidamber, C. F. Kemerer, A metrics suite for object oriented design, IEEE Transactions on software engineering 20 (6) (1994) 476–493.
- [90] R. Subramanyam, M. S. Krishnan, Empirical analysis of ck metrics for object-oriented design complexity: Implications for software defects, IEEE Transactions on software engineering 29 (4) (2003) 297–310.
- [91] S. Khalid, S. Zehra, F. Arif, Analysis of object oriented complexity and testability using object oriented design metrics, in: Proceedings of the 2010 National Software Engineering Conference, ACM, 2010, p. 4.
- [92] U. L. Kulkarni, Y. Kalshetty, V. G. Arde, Validation of ck metrics for object oriented design measurement, in: Emerging Trends in Engineering and Technology (ICETET), 2010 3rd International Conference on, IEEE, 2010, pp. 646–651.
- [93] A. Capiluppi, N. Ajiienka, The relevance of application domains in empirical findings, in: Proceedings of the 2nd International Workshop on Software Health, SoHeal@ICSE 2019, Montreal, Canada, May 30, 2019, 2019.
- [94] G. Succi, W. Pedrycz, S. Djokic, P. Zuliani, B. Russo, An empirical exploration of the distributions of the chidamber and kemerer object-oriented metrics suite, Empirical Software Engineering 10 (1) (2005) 81–104.

[95] A. Marcus, D. Poshyvanyk, The conceptual cohesion of classes, in: 21st IEEE International Conference on Software Maintenance (ICSM'05), IEEE, 2005, pp. 133–142.

[96] J. P. Flynt, O. Salem, Software engineering for game developers, Course Technology PTR, 2005.

1190