

OSSEAN: Mining Crowd Wisdom in Open Source Communities

Gang Yin, Tao Wang, Huaimin Wang, Qiang Fan, Yang Zhang, Yue Yu, Cheng Yang

National Laboratory of Parallel and Distributed Computing
School of Computer, National University of Defense Technology
ChangSha, Hunan, China

yingang@nudt.edu.cn, taowang.2005@outlook.com, whm_w@163.com

Abstract—Nowadays open source software represents a successful crowd-based software production model and is becoming an ecosystem combining huge amounts of software producers (such as software developers) and consumers (such as software users and customers). Lots of research work has been conducted on analyzing software artifacts created by producers, but few of them reveal the power of feedback from consumers which we believe is very important for the evaluation and evolution of open source software. This paper introduces OSSEAN, a platform for Open Source Software Evaluating, Analyzing and Networking. OSSEAN divides the open source communities into two groups: software production communities and software consumption communities. The former contain structured software artifacts such as projects, source code and issues, while the latter are full of textual documents with rich semantics of user feedback. We show the power of OSSEAN with some interesting demos by analyzing more than 200 thousands of open source projects and 10 million documents.

Keywords—Open Source; Crowd Wisdom; Software Production Communities, Software Consumption Communities; OSSEAN

I. INTRODUCTION

Open source communities successfully leverage the power of the crowd in software production, and have great impacts on many stages of software development and applications in a global open source ecosystem. With the development and division of open source communities, software programmers and testers are attracted by software *production communities* such as SourceForge and Github, while the newcomers, users and customers are more likely going to software *consumption communities* such as StackOverflow, Slashdot and CSDN (<http://www.csdn.net>, the biggest IT community site for Chinese speaking users in the world). These two kinds of communities complement each other and greatly expand the scope of traditional software development activities to global software evaluation and evolution.

The production communities mainly help software developers manage their development processes and artifacts. For example, Github provides development tools such as version control and issue tracking, and social communication tools such as @mention [1]. SourceForge provides more complete toolkits for distributed collaboration and management, including mailing lists, feature request, etc. These tools store huge amounts of software engineering data for structured software artifacts. On the other hand, the consumption communities are usually the crowd-oriented knowledge sharing platforms attracting tens of millions of users. The data generated in in these communities are usually textual posts which are reacting more quickly than the development requests

submitted in production communities. For example, StackOverflow has an answer rate above 90% and a median answer time of only 11 minutes [2], while the average responding time for an issue in Android issue tracking system (a typical production community) is about 31 days. The consumption communities are becoming the source of the crowd wisdom for the evaluation and evolution of open source software in production communities.

We propose a new approach, OSSEAN (Open Source Software Evaluating, Analyzing and Networking), to leverage the crowd wisdom in consumption communities to support the software development in production communities. OSSEAN is composed of two steps: firstly collects the set of documents in consumption communities for each software in production communities, then use the documents to make evaluation, comparison and ranking for the software. According to our experiments, OSSEAN successfully discovers more than 8 million documents for more than 238 thousands projects. In this paper, we use three promising demos to show the potential applications of OSSEAN.

The structure of this paper is as follows. In Section II some closely related systems are introduced and discussed. In Section III key mechanisms of our approach are described. In Section IV, we show some preliminary but promising demos of OSSEAN. We conclude the paper in Section V.

II. RELATED WORK

A lot of research work has been conducted on collecting and mining data in open source communities.

The large amounts of high-quality source code publicly available over the internet have attracted great attention from researchers. Sushil et. al constructed an Internet-scale software repository Sourceer [7]. It employs the structural information like reference in source code, the dependences between libraries and so on to achieve large scale source code indexing and searching. OCEAN [8] present a federated search engine that simultaneously retrieves source code from existing open source code search engine sites including Koders, Krugle, Merobase and Google Code.

In industry, many analysis tools and services have been provided by different companies. Coverity Scan [9] mainly focus on analyzing the quality and security of open source software by providing scan and test services. Coverity Scan helps developers to identify critical quality and security defects that are hard to find by other methods, and can provide valuable information for users to locate and fix the identified defects. Until March 2014, the lines of source code they scan have reached 300 million. Many famous open source projects

such as FreeBSD, HBase and CloudStack benefit much from this service.

As a famous open source software management solution provider, Black Duck [10] has built a massive open source knowledge base named BlackDuck KnowledgeBase. It collects more than 1 million open source software from 7.5 thousands open source sites. In this knowledge base, they cover more than 100 million lines of source code and more than 2 thousands of different licenses. They build code search engine Koders and open source community OpenHub.

Most of these works focus on the static source code of open source software, and little attention has been paid on the user feedback in consumption communities.

III. THE OSSEAN PLATFORM

To make better use of the data generated in open source ecosystems, we design and implement OSSEAN—an Open Source Software Evaluating, Analyzing and Networking platform. OSSEAN aims to bridge the production communities and consumption communities to retrieve the knowledge hidden among different communities, which are very important for software evaluation and evolution in Internet. In this section we introduce some important design considerations and the data source of the platform.

A. Architecture

The OSSEAN architecture mainly consists of three layers: acquisition layer, analysis layer and presentation layer. Through the cooperation of the three parts, OSSEAN achieves the automatic data processing and mining for multiple production communities and consumption communities. Figure 1 presents an overview of the architecture.

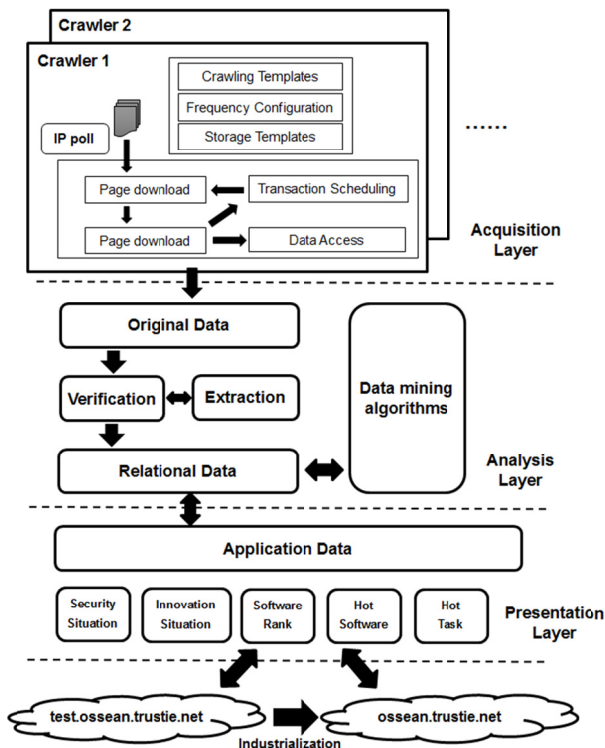


Figure 1: The architecture of OSSEAN

OSSEAN cuts the correlation between different modules in the three layers to a fullest extent. The communication between different modules is achieved through event service and database operations, so that each module can be developed and deployed by independent teams. Each module retrieves its input data from database and then notifies the next module by sending a specified event, with a relatively stable interface between each other. All modules in OSSEAN cooperate as in an assembly line.

The crawler module in acquisition layer is the engine of OSSEAN by crawling raw web pages from the two kinds of communities. A crawler will be created and deployed when a community is added as a new data source. Each crawler monitors a specified community site and retrieves data (raw web pages or API data) according to the predefined data model. Finally, the data crawled by crawler will wait for being analyzed by data extraction module. Here, the wrong data will be moved and the missing data will be re-crawled. Then, the data analysis module will construct the mappings between software projects and documents, and use the result to evaluate each project, such as analyzing the trend of projects and ranking them by popularity.

The modules in presentation layer are mainly designed for demonstration of the above results, including operations such as searching and ranking. Each module obtains data from the database schema of previous module, and saves result in its own schema. The results of data analysis will be placed in the data pool, prepared for presenting in platform. In order to provide a high quality data service, we have two online systems for OSSEAN: one is a testing system, and the other is a production system.

B. Data Source

OSSEAN collects data from both production communities and consumption communities. The first problem we met is what communities should be selected as data sources. Take development communities as example, Wikipedia lists more than 20 open source hosting websites such as SourceForge and Github, and there are thousands of community websites for interdependent open source organizations such as Linux, Apache, Eclipse, OW2, etc. As to consumption communities, there are far more candidates for consideration, such as StackOverflow, Slashdot, CSDN and OsChina (the last two are famous in China). We rank the communities according to their popularity.

(1) Data source from production communities:

We select two of the most popular development communities (SourceForge and GitHub) and one directory community (OpenHub). The number of projects and Alexa rank are shown in table 1. GitHub is a repository hosting community. It hosts 10,392,468 repositories up to August 2014. However, usually a software project can be forked many times in Github, and therefore a lot of repositories can be created for one project. So we use the number of root repositories as the number of projects in Github.

Among the total projects in these communities, there are lots of projects that are created just for a try of the community tools, these projects have few value and are the main cause of disturbance in later data processing. Another data disturbance is the duplicates of projects among different communities. To make sure that only active and useful projects are selected as

data source, we implement the following filtering policies: (1) remove the projects that either of its number of stars or downloads is zero; (2) rank the root repositories by their scores in Github and chose the first 116,850, where the score of each repository is the sum of its forks, watchers and commits in Github; (3) rank the projects by their scores in OpenHub and chose the first 98,833, where the score of each repository is the sum of its number of users, reviews and contributors in OpenHub. The statistical results are shown in table 1.

TABLE I. SOFTWARE PRODUCTION COMMUNITIES

	Site Name	Total Projects	Selected Projects	Alexa Popularity Ranking
1	GitHub	3,099,089	116,850	98
2	SourceForge	411,716	48,712	217
3	OpenHub	664,620	98,833	28,419
	Total	4,175,425	264,395	—
	Without Duplicates	—	238,467	—

(2) Data source from consumption communities:

There are much more consumption communities than production communities in Internet. We select most popular and hottest communities as our data source, mainly referring to their Alexa Rank (a websites ranking system provided by alexa.com according to the frequency of visits). Table 2 shows the number of posts collected from popular communities up to December 2014, including some famous Chinese communities. The posts are various kinds of documents generated in these communities for Q&A, discussion and sharing among users.

TABLE II. SOFTWARE CONSUMPTION COMMUNITIES

	Sites	Posts	Alexa Popularity Ranking
1	http://stackoverflow.com/	7,205,198	59
2	http://www.csdn.net/	848,998	442
3	http://www.cnblogs.com/	4140	1,241
4	http://www.codeproject.com/	187,692	1,346
5	http://slashdot.org/	45,203	1,484
6	http://www.oschina.net/	382,200	1,632
7	http://www.iteye.com/	92,545	2,441
8	http://www.dewen.io/	16,229	159,508
9	http://www.zdnet.com/	43,644	18,400
10	http://www.lupaworld.com/	29,541	210,871
	Total	8,855,390	—

C. Data Acquisition

The two kinds of communities are publicly available, but distributed dispersedly over the whole Internet. It is a great challenge to retrieve such kind of data continuously and automatically. The task of acquisition layer is obtaining data from these community sites. OSSEAN uses APIs provided by community web site and web crawlers to construct a data flow storage and management platform. For data crawling, it mainly consists of two interdependent processes: firstly crawls the raw web pages, and then extracts the attributes from the web pages. We know that most data collecting platform extracts the page information simultaneously while downloads web pages, but this is not the good solution for the changing and growing sites.

This design is necessary for the insurance of data quality, which is very important for the retrieval of open-source data. We cannot identify whether an error is created in the process of crawling or extraction when extracted information is wrong. Because of this, the conventional methods may introduce a lot

of low-quality data and data checking task is very difficult, and thus easily bring a lot of problems to the next stage of data analyzing. In order to improve the quality of acquired data, we add a new verification step into extracting process. The following are the details about these three processes.

Crawling process: A crawler in OSSEAN uses target-oriented retrieval policies. This process contains two steps. A crawler firstly crawls all the lists from a community and extracts the links of every item (here an item is an open source project). Then, the target pages referred by these links are downloaded and stored in database. Here we use multi-threaded techniques to improve performance. However, in order to prevent disturbing the performance of target communities, we set our crawling speed to a very low level.

Extracting process: This process extracts key attributes from raw web pages. Usually, different community sites have totally different page templates. Thus, we design corresponding extraction templates separately for each site. OSSEAN will automatically select right template to extract the input page according to its “site” name. The extracted attributes will be stored in database (here we use MySQL to stall these structured data).

Verifying process: This process verifies the potential errors exist in the crawling and extracting process. Verification of the extracted attributes will be invoked just after all pages are extracted. Different attributes have different verification methods. For example, a “time” related attribute (such as the “publish time” of a post) will be checked whether it has a right time format; while for the title of a post, it will be checked whether it is a valid textual string. In this process, the error handing module of each site will also locate the position of the error, either in download module or extraction module. The verification results will be logged for providing further bug-fixing instructions.

OSSEAN provides a plugin-based structure for adding new communities. The code framework for the three processes will be created and developers need only to change the extraction templates and verification templates.

D. Data Bridging

OSSEAN builds connections between production communities and consumption communities. Bridging two kinds of communities can be seen as a classification of posts to software project names, but we cannot use the conventional method to complete this classification. First, there is no training set for us to train a classification model, and clustering algorithm has the problem of discerning the body of each class. Second, there are many projects germinating every day, which makes it hard to build a stable model. Therefore, we use text matching method to solve this problem. For each post, we consider three attributes – tags, title and content. OSSEAN performs the matching in three steps – data preprocess, index building and entity matching. The following is the detail about the three steps.

Data preprocess: This step integrates the data from different communities. The data from different communities have different structure, so we define a standard structure for two kinds of communities. For example, the basic data structure of a post is a 13-tuple (subject, abstract, content, author, number of replies, number of views, type, category, created time, updated time, crawled time, URL, site), while that

of a project is a 13-tuple (name, description, language, number of followers, number of downloads, number of views, number of committers, category, created time, updated time, crawled time, URL, site).

For production communities, we also aggregate the multiple occurrence of a same project across different communities into one project.

Index Building: All the posts in OSSEAN are indexed according to their three attributes (tags, title and content). Each post is segmented and then its inverted index is created. Therefore, it will be very fast to answer queries such as whether a post contains a project name.

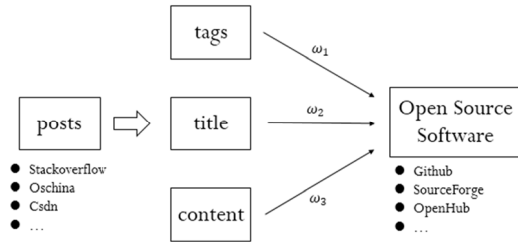


Figure 2. The flow of the matching process

Entity Matching: Given a post, each kind of its attributes is attached a weight ω_i , as shown in figure 2. OSSEAN uses string comparison to match tags matching, i.e., if the tag list of a post contains the software name then we get a positive match. OSSEAN uses index-based search for title and content matching, i.e., the name of the software will be used as a query term to search in the index. The score of the match is the sum of ω_i , where the match i should be successful. Currently, OSSEAN set ω_1 , ω_2 and ω_3 to 1, 0.8 and 0.5 respectively.

Text matching method has many advantages in data bridging, but it also performs poorly in some cases. For example, it cannot confirm whether a term in a post refers to a software project even when the term is identical to the name of that project. Another case is that different projects may share similar or even same names. We are using some techniques to resolve these problems, such as Apriori and SVM.

IV. PRELIMINARY RESULTS

With the continuous growth of community data, OSSEAN can provide lots of interesting services by mining the crowd wisdom in open source communities. This section shows some preliminary results and demos.

A. Ecosystem Measurement

OSSEAN can answer some basic questions related to the scope and scale of the global open source ecosystem, such as: (1) How many open source organizations and software projects in production communities? (2) How many posts (or documents) in consumption communities are related to open source projects?

These questions are related to the global development of open source software ecosystem. As shown in figure 3, from 2008 to 2014, the number of organizations in Github goes up to the highest point 118845. Furthermore, we investigate the quantity of answer posts and question posts in StackOverflow. Both the number of answers and questions are increasing

steadily to the highest points 12517744 and 7148960 as shown in figure 4. This indicates that the open source communities attract more and more organization as well as users to contribute and participate.

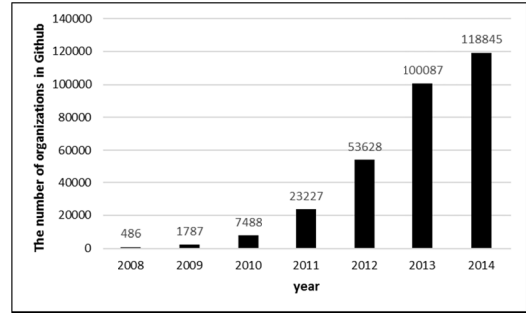


Figure 3: The number of organizations in Github in different years

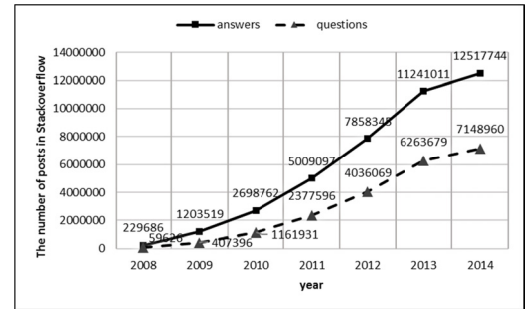


Figure 4: The number of posts in StackOverflow in different years

B. Software Ranking

OSSEAN can help users find the favorable software based on the popularity of its related posts. Some key issues should be properly addressed, such as: (1) How to compute the popularity of a post? (2) What is the rank of a software project in the global open source ecosystem?

Compared with the traditional software ranking mechanism, we sort the software projects by the set of posts related to them. The basic idea is that, the more posts a software project successfully matches, the more popular this project is and should be sorted forward. In our observation, we find that different posts have different impacts. In order to distinguish these differences, we use the number of comments for judging the popularity of a post. We denote the popularity of a post by the score of this post. The following is the formula for computing the score of post p :

$$\begin{cases} Score(p) = comments(p) / timediff(p) \\ timediff(p) = time_{now}(p) - time_{open}(p) \end{cases}$$

Where $Score(p)$ is the score of post p . $comments(p)$ is the number of comments of post p . $timediff(p)$ is the time interval between post p is opened and now.

The rank of a software project is determined by the sum of scores of its posts. However, we find that in many cases one post may discuss more than one project. To avoid bias, we use the average score to calculate the software rank.

$$Rank(s) = \sum_{i=1}^n \frac{Score(p_i)}{software(p_i)}$$

