

Stochastic Models for Evaluating Software Modularization Recovery Algorithms

Rodrigo Souza

*Departamento de Sistemas e Computação
Universidade Federal de Campina Grande
Campina Grande, Brazil
rodrigors@gmail.com*

Abstract—Software modularization recovery algorithms help to understand how software systems decompose into modules, but they

Software modularization recovery algorithms... bla bla

Keywords—reverse engineering; software modularization; empirical evaluation; complex networks;

I. ABSTRACT

Software modularization recovery algorithms automatically recognize a system’s modular structure by analyzing its implementation. Due to the lack of well document software systems, though, the issue of testing these algorithms is still underexplored, limiting both their adoption in the industry and the development of better algorithms. We propose to rely on software models to produce arbitrarily large test sets. In this paper we consider three such models and analyze how similar the artifacts they produce are from artifacts from real software systems.

II. INTRODUCTION

Development of large-scale software systems is a challenge.

A key to success is the ability to decompose a system into weakly-coupled modules, so each module can be developed by a distinct team. Failing to do so results in duplicated code, non-parallelism, one’s work impacting another’s work etc.

The ability do modularize depends decisively on a vast knowledge about the system, how its different parts interact to accomplish the system’s goal.

Unfortunately, in the case of legacy systems, such knowledge isn’t available. Depending on its size, it might take months to understand the system so well as to find a good modularization.

POR ISSO SURGIRAM software modularization recovery algorithms, also known as software clustering algorithms or software architecture recovery algorithms. In its most common flavor, these algorithms analyze the dependencies between implementation components, such as classes, and then group them into modules such as there are few dependencies between classes in distinct modules.

Software modularization recovery algorithms can, therefore, do in minutes what a person would spend weeks or months. The question is: are the found modularizations good? Are they similar to what a person would find? To answer this question it’s essential to perform empirical evaluations involving systems with known reference modularizations.

The empirical evaluations consist of selecting a collection of systems with known reference modularizations and then applying the algorithms to the systems. The modularizations found by the algorithms are then compared to the reference decompositions by a metric such as MoJo CITE or PrecisionRecall CITE.

Unfortunately there are few systems with known reference modularizations and, because to obtain reference modularizations is costly, there are few empirical studies, and most of them consider a couple of small and medium systems.

We therefore propose to use synthetic, i.e., computer-generated, software dependency networks, to evaluate software modularization recovery algorithms. These networks are generated by parametrizable models and have an embedded reference modularizations. The goal of an algorithm is, thus, to find modularizations that are similar to the reference modularization embedded in the network. With this approach we can CONTAR COM a large volume of test data that is composed of networks of different sizes and controllable characteristics.

Of course the success of this approach depends on the realism of the synthetic networks, ie, how well they resemble networks extracted from real software systems. In this paper we study three models and show that all of them are, by means of a careful parameter choosing, capable of producing realistic software networks.

The remaining sections are organized as Section 2, ...

III. SOFTWARE NETWORKS

directed graph, (un)weighted

IV. COMPLEX NETWORKS

Complex network theory found many scale-free networks

Software dependency networks are scale free. CITE Valverde, Myers

Scale free means ... $N(k) \sim k^{-\gamma}$

V. MODELS

Many scale-free models have been proposed. Only a few, though, produce modular networks.

A. LF

Directed weighted networks with overlapping modules.

B. BCR plus

We propose an extension to BCR model...

The BCR+ model builds networks that respect a module-dependency network, that is, a network that contains modules and allowed dependencies between modules. In the component dependency network, an edge between two vertices is only allowed if there exist an edge between their respective modules.

First of all, the model generates a network that is a copy of the module dependency network, each vertex belonging to its respective module. Then the algorithm runs iteratively, and on each iteration it picks a event that evolves the network. Each event has an associated probability.

TODO: present in pseudo-language?

* With probability p , a new vertex is added to the network, together with an edge from the new vertex to an existing vertex, chosen preferentially according to Considering $d_{in} = 0$, this means that a vertex with indegree = 8 is twice more likely to receive the edge than a vertex with indegree = 4. The parameter d_{in} can be used to alleviate the handicap. (If $d_{in} = 4$, the first vertex will only be 3/2 more likely to receive the edge). The new vertex is put on the same cluster as the old vertex. * With probability q , a new vertex is added to the network, together with an edge from an existing vertex, chosen preferentially accordingly to $d_{out} + \text{outdegree}$, to the new vertex. The new vertex is put on the same cluster as the old vertex * With probability r , a new edge is added between two vertices, v and w . v is chosen according to After that, with prob X , w is chosen from one of the modules connected to v 's module, else it's chosen from v 's own module. In any case, the exact choice is made according to ...

It's a growth model, that is, it can take as input an existing network and evolve it.

C. CGW

The CGW model is similar to BCR+ in which it is a growth model. The initial network is composed of two vertices with a directed edge between, belonging to the same module. The other $M - 1$ modules are initially empty.

Then, at each iteration, the algorithm executes one of the following events:

* With probability p_1 , one vertex is added to a randomly-chosen module, together with e_1 edges from p_1 to vertices chosen according to module-based preferential attachment.

* With probability

Accounts for the removal of edges. Growth model.

VI. EXPERIMENTAL SETUP

We want to show investigate if the models can produce networks that resemble software networks. We know that they share with software networks the scale-free property. This is not enough, since many real networks share this property. So we looked for a method to differentiate between software networks and other networks.

What we are looking for, after all, is an oracle that accepts software networks and rejects non-software networks. If the oracle has these two properties, we can be confident that it'll accept only synthetic networks that resemble software networks.

In a recent work, Milo et al. proposed the study of triads in order to characterize different classes of networks. We thus follow their work here.

Triads are...

Figure 1a show triads for a software system... Figure 2a show triads for network from domain X.

We've collected 65 systems written in Java and X networks from many domains, such as sociology, biology, technology and linguistics. Table 1 ...

We then used Pearson's correlation coefficient as a similarity measure between two networks. For each software network, We then computed the average correlation coefficient (ACC) to the other software networks. We've observed that among software networks the ACC is $X \pm Y$.

We then computed, for each non-software network, the ACC to the 65 software networks. By the 3-sigma rule, we use X as the threshold for realistic software networks: networks whose ACC is below this threshold are rejected.

The oracle has X precision and X recall...

We generated networks with many combinations of parameters...

BCR: 5 different architectures, p, q, r in $(0.0, 0.2, 0.4, 0.6, 0.8, 1.0)$, with $p+q+r = 1$ and $p + q \neq 0$ Total: X networks

LR: ...

CGW: ...

VII. EXPERIMENTAL RESULTS

All models produce networks that resemble software networks.

For some parameters, though, the networks are not realistic.

We cannot blame one single parameter for the non-realism.

VIII. CONCLUSION AND FUTURE WORK

The use of models lead to a more controlled way of experimented, and allows us to play with a large data set and a wide variety of characteristics. Although we have shown that

the models studied can produce realistic networks according to one criterium, we can't just discard the experimentation with real software systems, as they are the subjects in which algorithms will be used in the real world. In particular, it's an open question whether the models can reproduce any software system or if they just generate a subset of systems. We believe that there are software systems that cannot be simulated by none of these models.

Although common in distributed computing research, simulation is underexplored in software engineering. This work opens a door to the usage of simulation in the field of reverse engineering of software in order to evaluate RE algorithms.

Future work: to apply algorithms to the networks and compare the results with results found in the literature.

IX. CONCLUSION

The conclusion goes here. this is more of the conclusion

ACKNOWLEDGMENT

Dalton, Jorge, Christina, Garcia, Charles, Fabiola, Italo, Roberto, Jemerson, Sandra.