

2019

Automatic Detection of Insecure Codes in Stack Overflow

Shifu Hou

West Virginia University, shhou@mix.wvu.edu

Follow this and additional works at: <https://researchrepository.wvu.edu/etd>



Part of the [Information Security Commons](#)

Recommended Citation

Hou, Shifu, "Automatic Detection of Insecure Codes in Stack Overflow" (2019). *Graduate Theses, Dissertations, and Problem Reports*. 4069.

<https://researchrepository.wvu.edu/etd/4069>

This Thesis is protected by copyright and/or related rights. It has been brought to you by the The Research Repository @ WVU with permission from the rights-holder(s). You are free to use this Thesis in any way that is permitted by the copyright and related rights legislation that applies to your use. For other uses you must obtain permission from the rights-holder(s) directly, unless additional rights are indicated by a Creative Commons license in the record and/ or on the work itself. This Thesis has been accepted for inclusion in WVU Graduate Theses, Dissertations, and Problem Reports collection by an authorized administrator of The Research Repository @ WVU. For more information, please contact researchrepository@mail.wvu.edu.

Automatic Detection of Insecure Codes in Stack Overflow

Shifu Hou

Thesis submitted to the
Benjamin M. Statler College of Engineering and Mineral Resources
at West Virginia University
in partial fulfillment of the requirements
for the degree of
Master of Science in
Computer Science

Yanfang Ye, Ph.D., Chair of Committee

Xin Li, Ph.D.

Elaine M. Eschen, Ph.D.

Lane Department of Computer Science and Electrical Engineering
Morgantown, West Virginia
2019

Keywords: Social Code platform; LSTM; HIN

Copyright 2019 Shifu Hou

ABSTRACT

Automatic Detection of Insecure Codes in Stack Overflow

Shifu Hou

As the popularity of modern social coding paradigm such as Stack Overflow grows, its potential security risks increase as well (e.g., insecure codes could be easily embedded and distributed). To address this largely overlooked issue, we bring a new insight to exploit social coding properties in addition to code content for automatic detection of insecure code snippets in Stack Overflow. To determine if the given code snippets are insecure, we not only analyze the code content, but also utilize various kinds of relations among users, badges, questions, answers, code snippets and keywords in Stack Overflow. To model the rich semantic relationships, we first introduce a structured heterogeneous information network (HIN) for representation and then use meta-path based approach to incorporate higher-level semantics to build up relatedness over code snippets. Later, we propose two different novel network embedding models named Snippet2vec and CodeHin2Vec for representation learning in HIN to automate the insecure code snippet detection in Stack Overflow. More specifically, Snippet2vec learns the low dimensional representations for the nodes (i.e., code snippets) in the HIN where both the HIN structures and semantics are maximally preserved, while CodeHin2Vec utilizes HIN to depict relatedness over code snippets to generate code-to-code sequences, based on which sequence-to-sequence (seq2seq) concept in machine translation is further leveraged to learn representations of code snippets. Accordingly, we developed systems ICSD and iTrustSO which integrate our proposed methods respectively in insecure code snippet detection in Stack Overflow. Comprehensive experiments on the data collections from Stack Overflow are conducted to validate the effectiveness of our developed systems by comparisons with the state-of-the-art baseline methods.

Acknowledgments

I would first like to express my greatest gratitude to my committee chair and advisor, Dr. Yanfang Ye, for her guidance and support not only for this thesis but throughout the time of my whole master study. Her passion, vision, attitude, and love for research is always an inspiration source and influence to me; her expertise, understanding, generous guidance, suggestions, valuable comments and revisions make it possible for me to work on such an exciting topic; her devotion of significant time and efforts on mentoring my research has resulted in seven publications by the date of this thesis.

I would also like to thank my committee members, Dr. Li and Dr. Eschen, for their time and help for my research work; I am very fortunate to work with a cheerful group members, including Lingwei Chen, Yujie Fan, Yiming Zhang, and Aaron Saas, who exchanged ideas about machine learning related research work and provided useful suggestions on my thesis.

I am highly thankful to be blessed by amazing and talented family members and friends, who have made such a positive impact on my daily life, study, and research.

This work is partially supported by the NSF under grants OAC-1839909, CNS-1618629, CNS-1814825, and CNS-1845138, the DoJ/NIJ under grant NIJ 2018-75-CX-0032, the WV HEPC Grant (HEPC.dsr.18.5), and the WVU RSA grant (R-844).

Table of Contents

Abstract	ii
Acknowledgments	iii
Table of Contents	iv
List of Figures	v
List of Tables	vi
Chapter 1 Introduction	1
1.1 Background and Motivation	1
1.2 Research Objective	2
1.3 Major Contributions	4
1.4 Thesis Organization	5
Chapter 2 Related Work	6
2.1 Research on Stack Overflow	6
2.2 HIN and its Representation Learning	7
Chapter 3 ICSD:HIN Model for Insecure Code Snippet Detection	8
3.1 System Architecture	8
3.2 Proposed Method	9
3.3 Experimental Results and Analysis	19
Chapter 4 iTrustSO: Code-to-code Sequential Model over HIN for Insecure Code Snippet Detection	26
4.1 System Architecture	26
4.2 Proposed Method	27
4.3 Experimental Results and Analysis	34
Chapter 5 Conclusion and Future Work	39
Publications	40
Bibliography	40

List of Figures

1.1	Example of code security attacks in Stack Overflow.	2
1.2	An example of relatedness over code snippets.	3
3.1	System architecture of <i>ICSD</i>	9
3.2	Network schema for HIN in our application.	12
3.3	Meta-paths built for insecure code snippet detection (The symbols are the abbreviations shown in Figure 3.2).	13
3.4	Random walk guided by single meta-path vs. random walk guided by multiple meta-paths.	15
3.5	Parameter sensitivity evaluation.	24
3.6	Scalability evaluation.	24
3.7	Stability evaluation.	25
4.1	System architecture of <i>iTrustSO</i>	26
4.2	Network schema for HIN in our application.	27
4.3	Different contexts among code snippets.	28
4.4	Architecture of LSTM using hierarchical attention.	30
4.5	Attention evaluation.	35
4.6	Parameter sensitivity evaluation.	37

List of Tables

3.1	Performance indices of code snippet detection	20
3.2	Detection Results of different meta-paths	20
3.3	Comparisons with other network representation learning methods in insecure code snippet detection	22
3.4	Comparisons of other machine learning methods	23
4.1	Detection Results of different meta-paths	35
4.2	Comparisons of <i>CodeHin2Vec</i> with other network representation learning methods in insecure code snippet detection	36
4.3	Comparisons of other machine learning methods	38

Chapter 1

Introduction

1.1 Background and Motivation

Nowadays, as computing devices and Internet become increasingly ubiquitous, software has played a vital role in modern society covering many corners of our daily lives, such as Instant Message (IM) tools of WhatsApp and WeChat. In recent years, there has been an exponential growth in the number of software; it's estimated that the global software market reached approximately \$406.6 billions in 2017 [1]. Unlike conventional approaches (e.g., code handbook based), modern software developers heavily engage in a social coding environment, i.e., they tend to reuse code snippets and libraries or adapt existing ready-to-use projects during the process of software development [2]. In particular, Stack Overflow [3], as the largest online programming discussion platform, has attracted 8.9 million registered developers [4]. The vibrant discussions and ready-to-use code snippets make it one of the most important information sources to software developers [5]. Despite the apparent benefits of such social coding environment, its profound implications into the security of software remain poorly understood [6, 7]. For example, *can one trust code snippets posted in Stack Overflow?*

As the popularity of Stack Overflow grows, the incentive of launching a large-scale security attack by exploiting the vulnerability of posted code snippets increases as well. According to a recent study [8], collected question-answer samples from Stack Overflow contain various security-related issues such as encryption with insecure mode, insecure Application Programming Interface (API) usage and so on. Those innocent-looking yet insecure code snippets - if not properly handled and directly transplanted to production

software - could cause severe damage or even a disaster (e.g., disrupting system operations, leaking sensitive information) [8, 9]. For example, as shown in Figure 1.1, since cryptocurrency has grown popular, attackers have injected malicious mining code such as *Coinhive* - a cryptocurrency mining service - into Stack Overflow; once innocent developers reuse or copy-paste such code snippets to generate the production software, the software users' devices could be compromised (e.g., processing power would be stolen to mine bits of cryptocurrency).

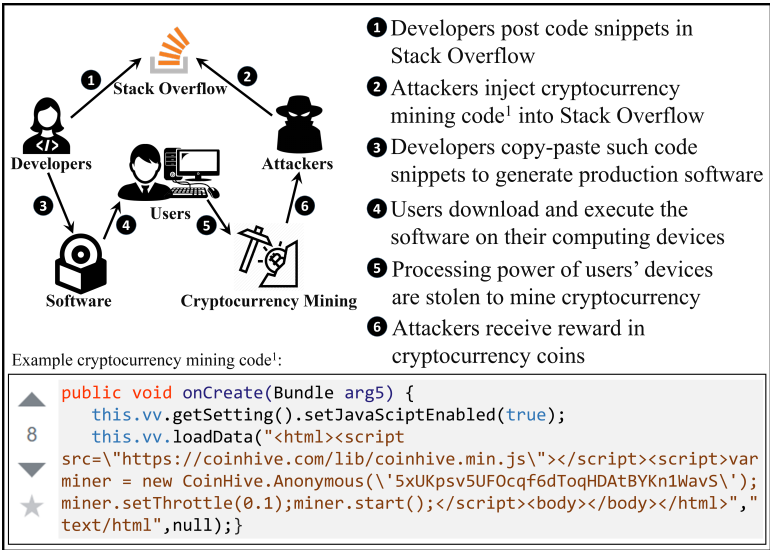


Figure 1.1: Example of code security attacks in Stack Overflow.

Stack Overflow has been aware of the negative impacts of insecure code infiltrations; unfortunately there has been no principled way of dealing with insecure code snippets included in the posted questions/answers other than labeling the moderator flag, down-voting those threads or warning in the comments [8]. Given the rich structure and information of Stack Overflow with ever-evolving programming languages, there is apparent and imminent need to develop novel and sound solutions to address the issue of code snippet security in Stack Overflow.

1.2 Research Objective

To address the above challenge of code security in Stack Overflow, an important new insight brought by this work is to exploit social coding properties in addition to code

content for automatic detection of insecure code snippets. As a social coding environment, Stack Overflow is characterized by user communication through questions and answers [10], that is, a rich source of heterogeneous information are available in Stack Overflow including users, badges, questions, answers, code snippets, and the rich semantic relationships among them. For example, as shown in Figure 1.2, to detect if a code snippet (*Code-2*) is insecure, using the code content (e.g., methods, functions, APIs, etc.) alone may not be sufficient; however, other rich information provided in Stack Overflow could be valuable for the prediction, such as (1) the same user (*User-1*) may be prone to post different insecure code snippets (*Code-1* and *Code-2*) due to his/her coherent code writing style, or (2) similar insecure code snippets (*Code-2* and *Code-3*) may be posted by a group of inexperienced users (*User-1* and *User-2* both only had the bronze badge of “commentator” that could be gained by leaving 10 comments in Stack Overflow). To utilize the social coding properties of Stack Overflow data (i.e., in-

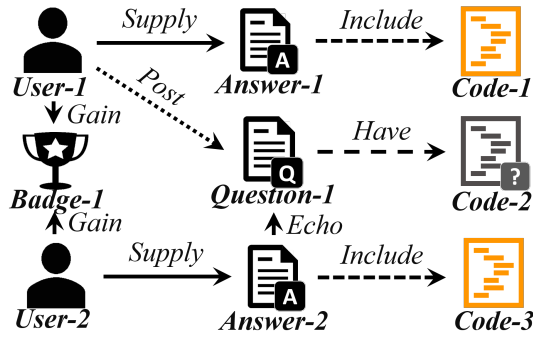


Figure 1.2: An example of relatedness over code snippets.

cluding different entities of users, badges, questions, answers and code snippets, as well as the rich semantic relationships among them) in addition to code content (i.e., keywords extracted from code snippets such as function names, methods and APIs), in our research, we propose to introduce a heterogeneous information network (HIN) [11, 12] as an abstract representation. Then we use meta-path [12] to incorporate higher-level semantic relationships to build up relatedness over the code snippets. Afterwards, to reduce the high computation and space cost, we further propose two different novel network embedding models named *snippet2vec* and *CodeHin2Vec* for node (i.e., code snippet) representation learning in the HIN, where both HIN structure and semantics are maximally preserved. After that, a classifier is constructed for automatic detection of insecure code snippets in Stack Overflow.

1.3 Major Contributions

The major contributions of our work can be summarized as follows:

- ***Novel feature representation of Stack Overflow data.*** Security risks arising from the new paradigm of social coding are more sophisticated than those from conventional wisdom, which requires a deeper understanding and a greater modeling effort. In addition to code content, a rich source of heterogeneous information in Stack Overflow including users, badges, questions, answers, code snippets, and the semantic relations among them is also available. To utilize such social coding properties (e.g., *question-code*, *answer-code*, *code-keywords*, *user-question*, *user-answer*, *question-answer*, and *user-badge* relations), we propose to introduce HIN as an abstract representation of Stack Overflow data. Then a meta-path based approach is exploited to characterize the relatedness over code snippets. The proposed solution provides a natural way of expressing complex relationships in social coding platforms such as Stack Overflow, which has not been studied in the open literature to our best knowledge.
- ***snippet2vec: an effective representation learning integrating node content and HIN-based relations*** Based on a set of built meta-path schemes, to reduce the high computation and space cost, a new network embedding model named *snippet2vec* is proposed to learn the low-dimensional representations for the nodes (i.e., code snippets) in the HIN, which are capable to preserve both the semantics and structural correlations between different types of nodes. Then, given different sets of meta-path schemes, different kinds of node (i.e., code snippet) representations will be learned by using *snippet2vec*. To aggregate these different learned node representations, we propose a multi-view fusion classifier to learn importance of them and thus to make predictions (i.e., a given code snippet will be labeled as either insecure or not).
- ***CodeHin2Vec: a code-to-code sequence modeling with LSTM for node embedding.*** More specifically, a new model *CodeHin2Vec* is proposed to seamlessly combine code content and HIN-based relations to learn representations of code snippets, in which code sequences are first generated based on the walk paths

guided by different meta-paths; in each code sequence, its elements are represented by the code content feature vectors; then, LSTM using hierarchical attention mechanism is leveraged for code sequence modeling. *CodeHin2Vec* is a generic framework which can also be applicable for other representation learning task.

- ***Two practical systems for automatic detection of insecure code snippets.*** Based on the collected and annotated data from Stack Overflow, we develop Two practical systems named *ICSD* and *iTrustSO* integrating our proposed methods for automatic detection of insecure code snippets. We provide comprehensive experimental studies to validate the performance of our developed systems in comparisons with alternative approaches. This work is the first attempt utilizing both code content and social coding properties for automatic analysis of code security in Stack Overflow. The proposed method and developed system can also be easily expanded to code security analysis in other social coding platforms, such as GitHub and Stack Exchange.

1.4 Thesis Organization

The remainder of this paper is organized as follows. Chapter 2 discusses the related work. Chapter 3 presents our developed systems *ICSD* in detail. Chapter 4 presents our developed systems *iTrustSO* in detail. Finally, Chapter 5 presents the concludes and our future work.

Chapter 2

Related Work

2.1 Research on Stack Overflow

There have been many works on knowledge discovery from Stack Overflow data [13, 14, 15, 16, 17, 18, 19, 20, 2, 21] - from gamification motivation for voluntary contributions [19], discussion interest trend [14, 15], patterns of questions/answers [17] and project-specific language differences [18], to developer interaction [20], dynamics of the participation [21], repair patterns from extracted code samples [16] and interplay between platform activities and development process [2]. However, most of these works have focused in Stack Overflow semantics and users behavior but rarely addressed the issue of code security analysis. The only exceptions appear to be [6] and [7] which both exploited Android app codes as a case study to evaluate the security of information source in Stack Overflow. Though those research results are promising, [6] only performed empirical studies while [7] merely analyzed the code snippet itself without considering any relationship to other Stack Overflow data (i.e., without utilizing the social coding properties in this platform). Different from the existing works, in this paper, to detect the insecure code snippets in Stack Overflow, we propose to utilize not only the code content, but also various kinds of relationships among users, badges, questions, answers, and code snippets. Based on the extracted relation features, the code snippets are depicted by a structured HIN.

2.2 HIN and its Representation Learning

HIN is used to model different types of entities and relations [22], which has been intensively deployed to various applications, such as scientific publication network analysis [11, 12], document analysis based on knowledge graph [23], social network analysis [24, 25], and malware detection [26, 27]. To reduce the high computation and space cost in network mining, many efficient network embedding methods have been proposed, including homogeneous network representation learning (e.g., DeepWalk [28], node2vec [29], PTE [30], and LINE [31]) and HIN representation learning (e.g., ESim [32], metapath2vec [33] and HIN2vec [34]). Unfortunately, these methods cannot be directly employed in our application, which is to exploit social coding properties in addition to code content for automatic detection of insecure code snippets. To tackle this challenge, in this paper, we propose a novel learning model named *snippet2vec* for node (i.e., code snippet) representation learning in HIN where both the HIN structures and semantics are maximally preserved; after that, a multi-view fusion classifier is constructed for insecure code snippet detection.

Chapter 3

ICSD:HIN Model for Insecure Code Snippet Detection

3.1 System Architecture

The system architecture of *ICSD* is shown in Figure 3.1, which is developed for insecure code snippet detection in Stack Overflow. It consists of the following major components:

- **Data collector.** A set of crawling tools are developed to collect the data from Stack Overflow. The collected data includes users' profiles, their posted questions and answers, and the code snippets embedded in the questions/answers.
- **Feature extractor.** Resting on the data collected from the previous module, to depict the code snippets, it first extracts the content-based features from the collected code snippets (i.e., keywords such as function names, methods and APIs), and then analyzes various relationships among different types of entities (i.e., user, badge, question, answer, code snippet, keyword), including i) *question-have-code*, ii) *answer-include-code*, iii) *code-contain-keyword*, iv) *user-post-question*, v) *user-supply-answer*, vi) *answer-echo-question*, and vii) *user-gain-badge* relations. (See Section 3.2 for details.)
- **HIN constructor.** In this module, based on the features extracted from the previous component, a structured HIN is first presented to model the relationships among different types of entities; and then different meta-paths are built from the

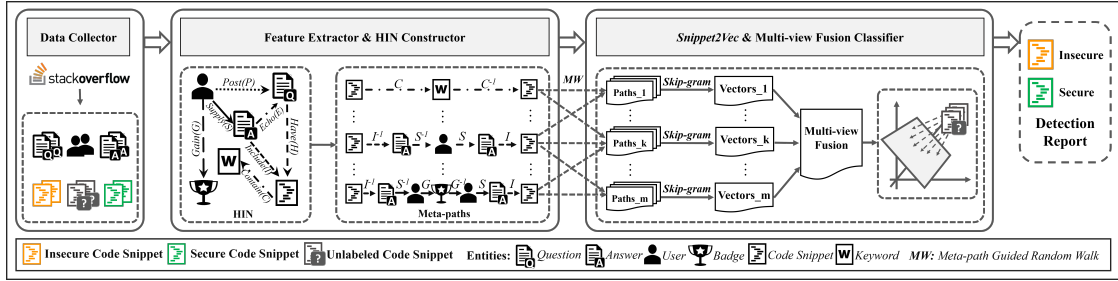


Figure 3.1: System architecture of *ICSD*.

HIN to capture the relatedness over code snippets from different views (i.e., with different semantic meanings). (See Section 3.2 for details.)

- ***snippet2vec***. Based on the built meta-path schemes, to reduce the high computation and space cost, a new network embedding model *snippet2vec* is proposed to learn the low-dimensional representations for the nodes in HIN, which are capable to preserve both the semantics and structural correlations between different types of nodes. In *snippet2vec*, given a set of different meta-path schemes, a meta-path guided random walk strategy is first proposed to map the word-context concept in a text corpus into a HIN; then skip-gram is leveraged to learn effective node representation for a HIN. (See Section 4.2 for details.)
- ***Multi-view fusion classifier***. Given different sets of meta-path schemes, different kinds of node (i.e., code snippet) representations will be learned by using *snippet2vec*. To aggregate these different representations, a multi-view fusion classifier is constructed to learn importance of them and thus to make predictions (i.e., the unlabeled code snippets will be predicted if they are insecure or not). (See Section 3.2 for details.)

3.2 Proposed Method

In this section, we present the detailed approaches of how we represent the code snippets in Stack Overflow utilizing both code content and social coding properties simultaneously, and how we solve the insecure code snippet detection problem based on the representation.

Feature Extraction

Code snippets. Stack Overflow provides the discussion platform for software developers to post their questions and answers about ever-evolving programming languages including Java, JavaScript, C/C++/C#, Python, PHP, perl, etc. In this paper, we will focus on Java programming language for Android application (app) development as a showcase for the following reasons: (1) Java is one of the most popular programming languages in Stack Overflow [17]. (2) Due to the mobility and ever expanding capabilities, mobile devices have recently surpassed desktop and other media - it is estimated that 77.7% of all devices connected to the Internet will be smart phones in 2019 [35, 36] (leaving PCs falling behind at 4.8%). Android, as an open source and customizable operating system for mobile devices, is currently dominating the smart phone market by 82.8% [37]. (3) Billions of mobile device users with millions of Android apps installed have attracted more and more developers; however, most of these Android mobile apps have poorly implemented security mechanisms partially because developers are inexperienced, distracted or overwhelmed [38, 6]. Indeed developers tend to request more permissions than what are actually needed, often use insecure options for Inter Component Communication (ICC), and fail to store sensitive information in private areas [17]. Code snippets in Stack Overflow are surrounded by `<code>` `</code>` tags, and they can thus easily be separated from accompanying texts before being extracted. Then, content-based features will be further extracted from the collected code snippets: we will first remove all the punctuations and stopwords; and then we will extract the keywords including function names, methods, APIs and parameters to represent the content of code snippets.

Social coding properties. To depict a code snippet in Stack Overflow, we not only utilize its above extracted content-based features, but also consider its social coding properties including followings.

- **R1:** To describe the relation that a question thread has a code snippet included, we generate the **question-have-code** matrix \mathbf{H} where each element $h_{i,j} \in \{0, 1\}$ indicates whether question i has code snippet j .
- **R2:** To denote the relation that an answer thread includes a code snippet, we generate the **answer-include-code** matrix \mathbf{I} where each element $i_{i,j} \in \{0, 1\}$ means if answer i includes code snippet j .

- **R3**: To represent the relation that a code snippet contains a specific keyword (e.g., function name of “Coinhive”), we build the *code-contain-keyword* matrix **C** whose element $c_{i,j} \in \{0, 1\}$ denotes whether code snippet i contains keyword j .
- **R4**: To describe the relation between a user and a question he/she posts, we generate the *user-post-question* matrix **P** where each element $p_{i,j} \in \{0, 1\}$ denotes if user i posts question j .
- **R5**: To represent the relation of a user and an answer he/she supplies, we generate the *user-supply-answer* matrix **S** where each element $s_{i,j} \in \{0, 1\}$ denotes whether the user i supplies answer j .
- **R6**: To denote the Q&A relationship, we build the *answer-echo-question* matrix **E** whose element $e_{i,j} \in \{0, 1\}$ denotes whether answer i echoes/responds to question j .
- **R7**: In order to encourage engagement, Stack Overflow has adopted a strategy of *gamification* [10] - users will be rewarded for their valued contributions to the forum. For example, “illuminator” badge (gold level in answer badges) will be awarded to the users who edit and answer 500 questions (both actions within 12 hours, answer score > 0). This can be seen as a measure of a user’s expertise by potential recruiters [39]. In Stack Overflow, there are different kinds of badges (e.g., question badges, answer badges, etc.) with different levels (i.e., gold, silver, and bronze). To describe the relationship between a user and a specific badge he/she gains, we build the *user-gain-badge* matrix **G** whose element $g_{i,j} \in \{0, 1\}$ denotes if user i gain badge j .

HIN Constructor

In order to depict users, badges, questions, answers, code snippets, keywords as well as the rich relationships among them (i.e., **R1-R7**), it is important to model them in a proper way so that different kinds of relations can be better and easier handled. We introduce how to use HIN, which is capable to be composed of different types of entities and relations, to represent the code snippets in Stack Overflow by using the features extracted above. We first present the concepts related to HIN as follows.

Definition 1. Heterogeneous information network (HIN) [40]. A HIN is defined as a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with an entity type mapping $\phi: \mathcal{V} \rightarrow \mathcal{A}$ and a relation type mapping $\psi: \mathcal{E} \rightarrow \mathcal{R}$, where \mathcal{V} denotes the entity set and \mathcal{E} is the relation set, \mathcal{A} denotes the entity type set and \mathcal{R} is the relation type set, and the number of entity types $|\mathcal{A}| > 1$ or the number of relation types $|\mathcal{R}| > 1$. The **network schema** [40] for a HIN \mathcal{G} , denoted as $\mathcal{T}_{\mathcal{G}} = (\mathcal{A}, \mathcal{R})$, is a graph with nodes as entity types from \mathcal{A} and edges as relation types from \mathcal{R} .

HIN not only provides the network structure of the data associations, but also provides a high-level abstraction of the categorical association. For our case, i.e., the detection of insecure code snippets in Stack Overflow, we have six entity types (i.e., user, badge, question, answer, code snippet, keyword) and seven types of relations among them (i.e., **R1-R7**). Based on the definitions above, the network schema for HIN in our application is shown in Figure 3.2, which enables the code snippets in Stack Overflow to be represented in a comprehensive way that utilizes both their content-based information and social coding properties.

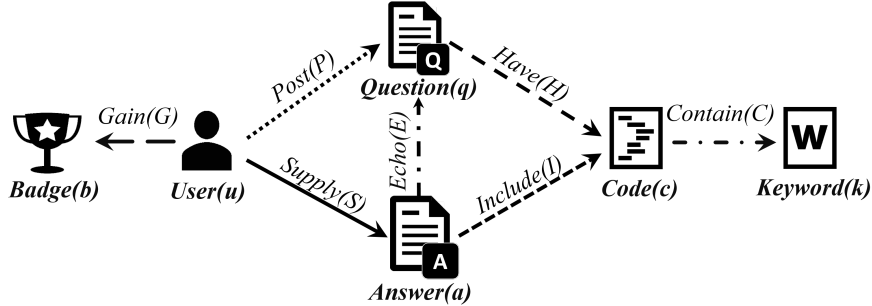


Figure 3.2: Network schema for HIN in our application.

The different types of entities and relations motivate us to use a machine-readable representation to enrich the semantics of relatedness among code snippets in Stack Overflow. To handle this, the concept of meta-path has been proposed [12] to formulate the higher-order relationships among entities in HIN. Here, we follow this concept and extend it to our application of insecure code snippet detection in Stack Overflow.

Definition 2. Meta-path [12]. A meta-path \mathcal{P} is a path defined on the graph of network schema $\mathcal{T}_{\mathcal{G}} = (\mathcal{A}, \mathcal{R})$, and is denoted in the form of $A_1 \xrightarrow{R_1} A_2 \xrightarrow{R_2} \dots \xrightarrow{R_L} A_{L+1}$, which defines a composite relation $R = R_1 \cdot R_2 \cdot \dots \cdot R_L$ between types A_1 and A_{L+1} , where \cdot denotes relation composition operator, and L is the length of \mathcal{P} .

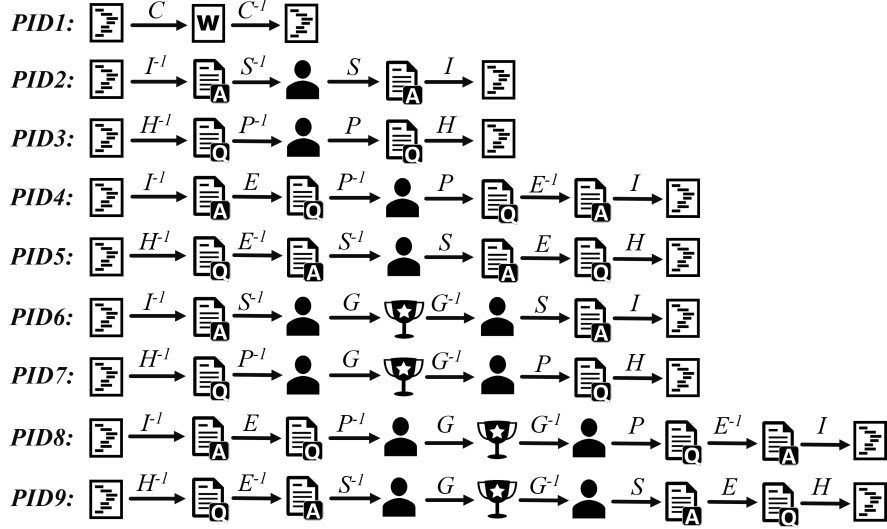


Figure 3.3: Meta-paths built for insecure code snippet detection (The symbols are the abbreviations shown in Figure 3.2).

Given a network schema with different types of entities and relations, we can enumerate a lot of meta-paths. In our application, based on the collected data, resting on the seven different kinds of relationships, we design nine meaningful meta-paths for characterizing relatedness over code snippets in Stack Overflow, i.e., **PID1-PID9** shown in Figure 3.3. Different meta-paths depict the relatedness between two code snippets at different views. For example, the meta-path **PID2** formulates the relatedness over code snippets in Stack Overflow: $code \xrightarrow{Include^{-1}} answer \xrightarrow{Supply^{-1}} user \xrightarrow{Supply} answer \xrightarrow{Include} code$ which means that two code snippets can be connected as they are included in the answers supplied by the same user; while another meta-path **PID6**: $code \xrightarrow{Include^{-1}} answer \xrightarrow{Supply^{-1}} user \xrightarrow{Gain} reputation \xrightarrow{Gain^{-1}} user \xrightarrow{Supply} answer \xrightarrow{Include} code$ denotes that two code snippets are related as they are included in the answers supplied by the users with the same kind of badge (e.g., “illuminator” badge) indicating similar expertise or contribution. In our application, meta-path is a straightforward method to connect code snippets via different relationships among different entities in HIN, and enables us to depict the relatedness over code snippets in Stack Overflow utilizing both their content-based information and social coding properties in a comprehensive way.

***snippet2vec*: HIN Representation Learning**

To measure the relatedness over HIN entities (e.g., code snippets), traditional representation learning for HIN [41, 42, 12, 43] mainly focuses on factorizing the matrix (e.g., adjacency matrix) of a HIN to generate latent-dimension features for the nodes (e.g., code snippets) in the HIN. However, the computational cost of decomposing a large-scale matrix is usually very expensive, and also suffers from its statistical performance drawback [29]. To reduce the high computation and space cost, it calls for scalable representation learning method for HIN. Given a HIN $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, the ***HIN representation learning*** task [34, 33] is to learn a function $f : \mathcal{V} \rightarrow \mathbb{R}^d$ that maps each node $v \in \mathcal{V}$ to a vector in a d -dimensional space \mathbb{R}^d , $d \ll |\mathcal{V}|$ that are capable to preserve the structural and semantic relations among them.

To solve the problem of HIN representation learning, due to the heterogeneous property of HIN (i.e., network consisting of multi-typed entities and relations), it is difficult to directly apply the conventional homogeneous network embedding techniques (e.g., DeepWalk [28], LINE [31], node2vec [29]) to learn the latent representations for HIN. To address this issue, HIN embedding methods such as metapath2vec [33] was proposed. In metapath2vec, given a meta-path scheme, it employs meta-path based random walk and heterogeneous skip-gram to learn the latent representations for HIN such that the semantic and structural correlations between different types of nodes could be persevered. The metapath2vec was proposed to support one meta-path scheme to guide the walker traversing HIN; however, in our application, the code snippets in Stack Overflow can be connected through nine different meta-path schemes. It may not be feasible to directly employ metapath2vec in our case for insecure code snippet detection. To put this into perspective, as shown in Figure 3.4, we gain further insight into Stack Overflow data and have following interesting findings:

- ***Finding 1***: Both insecure *Code-1* and *Code-2* (i.e., they can both cause potential confidential information leakage) are posted by *User-1* “Ke****a” (we here anonymize his user name) answering the questions about string access for Android app. Actually, *Code-1* and *Code-2* can be connected by the *Path-A* guided by the designed meta-path *PID2*.
- ***Finding 2***: The insecure codes of *Code-3* (i.e., it may allow users to remotely execute the malicious code) and *Code-4* (i.e., it can cause potential data breach)

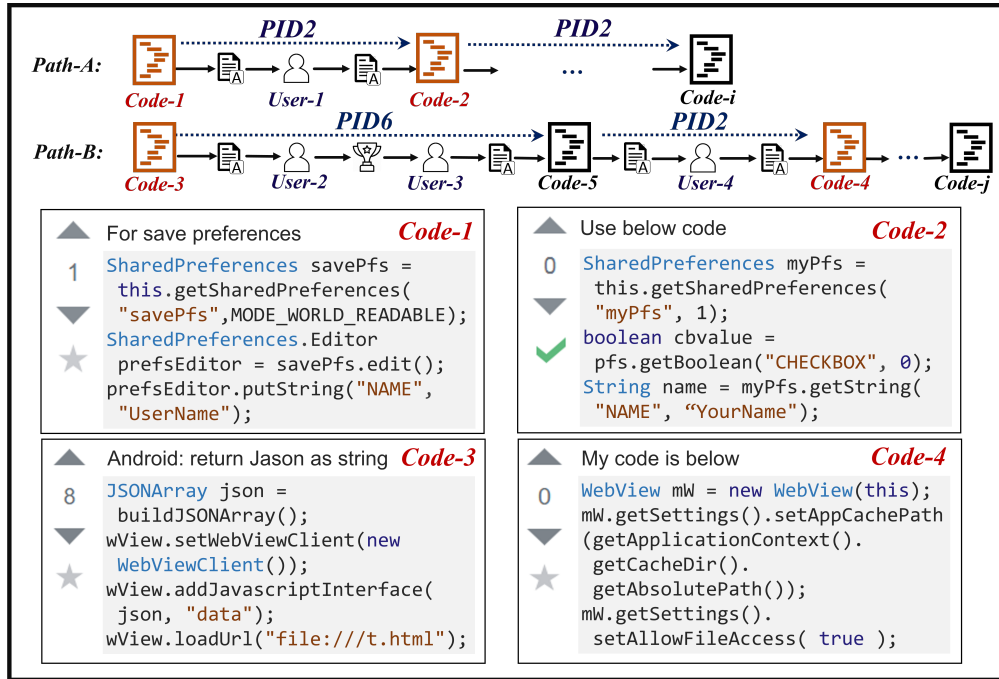


Figure 3.4: Random walk guided by single meta-path vs. random walk guided by multiple meta-paths.

are connected in the way that (1) *Code-3* and *Code-5* are related as they were posted by *User-2* and *User-3* who only had the bronze badge of “student” (i.e., first question with score of 1 or more); and then (2) *User-4* copied and pasted *Code-5* while also provided *Code-4* to answer another user’s posted question. In this way, *Code-3* and *Code-4* can be connected by the *Path-B* guided by meta-paths of *PID6* and *PID2*.

Based on the above observations, *metapath2vec* [33] fails to generate the path such as *Path-B* to represent the relatedness between code snippets like *Code-3* and *Code-4*. To address this issue, we design a new network embedding model *snippet2vec* to learn desirable node representations in HIN: first, a new random walk method guided by different meta-paths is proposed to map the word-context concept in a text corpus into a HIN; then skip-gram is leveraged to learn effective node representation for a HIN.

Random walk guided by different meta-paths. Given a source node v_j in a homogeneous network, the traditional random walk is a stochastic process with random variables $v_j^1, v_j^2, \dots, v_j^k$ such that v_j^{k+1} is a node chosen at random from the neighbors of node v_k . The transition probability $p(v_j^{i+1}|v_j^i)$ at step i is the normalized probabilit-

ity distributed over the neighbors of v_j^i by ignoring their node types. However, this mechanism is unable to capture the semantic and structural correlations among different types of nodes in a HIN. Here, we show how we use different built meta-paths to guide the random walker in a HIN to generate the paths of multiple types of nodes. Given a HIN $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with schema $\mathcal{T}_{\mathcal{G}} = (\mathcal{A}, \mathcal{R})$, and a set of different meta-paths $\mathcal{S} = \{\mathcal{P}_j\}_{j=1}^n$ (e.g., in *Finding2*, $\mathcal{S} = \{PID6, PID2\}$), each of which is in the form of $A_1 \rightarrow \dots A_t \rightarrow A_{t+1} \dots \rightarrow A_l$, we put a random walker to traverse the HIN. The random walker will first randomly choose a meta-path \mathcal{P}_k from \mathcal{S} and the transition probabilities at step i are defined as follows:

$$p(v^{i+1}|v_{A_t}^i, \mathcal{S}) = \begin{cases} \frac{\lambda}{|\mathcal{S}|} \frac{1}{|N_{A_{t+1}}(v_{A_t}^i)|} & \text{if } (v^{i+1}, v_{A_t}^i) \in \mathcal{E}, \phi(v_{A_t}^i) = A_c, \phi(v^{i+1}) = A_{t+1} \\ \frac{1}{|N_{A_{t+1}}(v_{A_t}^i)|} & \text{if } (v^{i+1}, v_{A_t}^i) \in \mathcal{E}, \phi(v_{A_t}^i) \neq A_c, \\ & \phi(v^{i+1}) = A_{t+1}, (A_t, A_{t+1}) \in \mathcal{P}_k \\ 0 & \text{otherwise,} \end{cases} \quad (3.1)$$

where ϕ is the node type mapping function, $N_{A_{t+1}}(v_{A_t}^i)$ denotes the A_{t+1} type of neighborhood of node $v_{A_t}^i$, A_c is entity type of *Code*, and λ is the number of meta-paths starting with $A_c \rightarrow A_{t+1}$ in the given meta-path set \mathcal{S} . The walk paths generated by the above strategy are able to preserve both the semantic and structural relations between different types of nodes in the HIN, and thus will facilitate the transformation of HIN structures into skip-gram.

Skip-gram. After mapping the word-context concept in a text corpus into a HIN via meta-path guided random walk strategy (i.e., a sentence in the corpus corresponds to a sampled path and a word corresponds to a node), skip-gram [44, 28] is then applied on the paths to minimize the loss of observing a node’s neighbourhood (within a window w) conditioned on its current representation. The objective function of skip-gram is:

$$\arg \min_Y \sum_{-w \leq k \leq w, j \neq k} -\log p(v_{j+k}|Y(v_j)), \quad (3.2)$$

where $Y(v_j)$ is the current representation vector of v_j , $p(v_{j+k}|Y(v_j))$ is defined using

the softmax function:

$$p(v_{j+k}|Y(v_j)) = \frac{\exp(Y(v_{j+k}) \cdot Y(v_j))}{\sum_{q=1}^{|\mathcal{V}|} \exp(Y(v_q) \cdot Y(v_j))}. \quad (3.3)$$

Due to its efficiency, we first apply hierarchical softmax technique [45] to solve Eq. 3.3; then the stochastic gradient descent [46] is employed to train the skip-gram.

Multi-view Fusion Classifier

Given a set of different meta-path schemes, by using the above proposed *snippet2vec*, the node (i.e., code snippet) representations will be learned in the HIN. In our application, as described in Section 3.2, we have nine meta-paths (i.e., *PID1–MID9*) which characterize the relatednesses over code snippets at different views (i.e., with different semantic meanings). Based on our observations on the Stack Overflow data and leveraging the domain expertise, we generate m sets of meta-path schemes $\mathbf{S} = \{\mathcal{S}_i\}_{i=1}^m$ for *snippet2vec* to learn the node representations in the HIN, where $m = 4$ and $\mathbf{S} = \{(PID1, PID2, PID6), (PID1, PID3, PID7), (PID1, PID4, PID8), (PID1, PID5, PID9)\}$. Given these different sets of meta-paths, using *snippet2vec*, different node representations will be learned in the HIN. Here, we propose to use multi-view fusion to aggregate these different learned node representations for code snippet classification.

Given m kinds of node representations $Y_i (i = 1, \dots, m)$ learned based on m sets of meta-path schemes, the incorporated node representations can be denoted as: $Y' = \sum_{i=1}^m (\alpha_i \times Y_i)$, where $\alpha_i (i = 1, \dots, m)$ is the weight of Y_i . To determine the weight of α_i for each mapped low-dimensional vector space Y_i , we measure the geometric distances among them. The distance measure based on the principal angles between two vector spaces is significant if and only if the vector spaces have the same dimensions [47]. In our case, the m mapped vector spaces are all with the same dimensions of d . Therefore, we apply the geodesic distance based on principal angles [48] to measure the geometric distances between the mapped vector spaces. The principal angle between space Y_i and Y_j is defined as the number $0 \leq \theta \leq \frac{\pi}{2}$ that satisfies:

$$\cos \theta = \max_{\mathbf{y} \in Y_i, \mathbf{y}' \in Y_j} \mathbf{y}^T \mathbf{y}'. \quad (3.4)$$

The angle θ is 0 if and only if $Y_i \cap Y_j \neq 0$, while $\theta = \frac{\pi}{2}$ if and only if $Y_i \perp Y_j$. Let $\theta_1, \theta_2, \dots, \theta_d$ be the d principal angles between space Y_i and Y_j , the geodesic distance between them is formulated as:

$$d(Y_i, Y_j) = \sqrt{\theta_1^2 + \theta_2^2 + \dots + \theta_d^2}. \quad (3.5)$$

Thus, we compute α_i for each mapped vector space Y_i as:

$$\alpha_i = \frac{\sum_{j=1, i \neq j}^m d(Y_i, Y_j)}{\sum_{i=1}^m \sum_{j=1, i \neq j}^m d(Y_i, Y_j)}. \quad (3.6)$$

To this end, the incorporated node representations Y' will be fed to the Support Vector Machine (SVM) to train the classification model, based on which the unlabeled code snippets can be predicted if they are insecure or not. Algorithm 2 shows the implementation of the our developed insecure code snippet detection system *ICSD*.

Algorithm 1 *ICSD* – Insecure code snippet detection in Stack Overflow based on structured HIN

Input: The HIN $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, network schema $\mathcal{T}_{\mathcal{G}} = (\mathcal{A}, \mathcal{R})$, m sets of meta-path schemes $\mathbf{S} = \{\mathcal{S}_i\}_{i=1}^m$, number of walk paths per node r , walk length l , and vector dimension d , training data set D_t , testing data set D_e **Output:** \mathbf{f} : The labels for the testing code snippets.

- 1: **for** $i = 1 \rightarrow m$ **do**
 - 2: **for** $j = 1 \rightarrow r$ **do**
 - 3: get l -length random walks using Eq. 3.1 guided by \mathcal{S}_i
 - 4: **end for**
 - 5: Generate $Y_i \in \mathbb{R}^{|\mathcal{V}| \times d}$ using skip-gram in Eq. 3.2
 - 6: **end for**
 - 7: **for** $i = 1 \rightarrow m$ **do**
 - 8: Calculate α_i for Y_i using Eq. 3.4–Eq. 3.6
 - 9: **end for**
 - 10: Get incorporated node representations $\mathbf{Y}' = \sum_{i=1}^m (\alpha_i \times Y_i)$
 - 11: Train SVM using \mathbf{Y}'_{D_t}
 - 12: **for** $k = 1 \rightarrow |D_e|$ **do**
 - 13: Generate the label f_k using trained SVM
 - 14: **end for**
 - 15: return \mathbf{f}
-

3.3 Experimental Results and Analysis

Experimental Setup

We develop a set of crawling tools to collect the data from Stack Overflow. As stated in Section 3.2, we consider Java programming language for Android app as a case study to evaluate our developed system. Note that it’s also applicable to other kinds of programming languages in Stack Overflow. We use our developed crawling tools to collect users’ profiles, question threads, answer threads, and code snippets in Stack Overflow in a period of time. By the date, we have collected 429,523 question threads and 623,746 answer threads posted by 213,560 users including 737,215 code snippets, through March 2010 to May 2018. To obtain the ground truth for the evaluation of different detection methods, we need to prelabel a fraction of code snippets (i.e., either secure or insecure). We first categorize code security risks and vulnerabilities for Android apps into six categories: (1) Android Manifest configuration, (2) WebView component, (3) data security, (4) file directory traversal, (5) implicit intents, and (6) security checking; and then we leverage our domain expertise and follow the principles such as least permission request, correct usage of HTTPS and TLS for networking, secure inter-component communication, secure storage to manually label a filtered set of 20,137 code snippets (i.e., 9,054 code snippets are labeled as *insecure* while 11,083 are *secure*). After feature extraction and based on the designed network schema, the constructed HIN has 80,405 nodes (i.e., 20,137 nodes with type of code snippet, 24,286 nodes with type of answer, 13,924 nodes with type of question, 21,471 with type of user, 94 with type of badges, and 493 with type of selected keywords) and 592,082 edges including relations of *R1-R7*. We use the performance indices shown in Table 3.1 to quantitatively validate the effectiveness of different methods in insecure code snippet detection.

snippet2vec based on Different Sets of Meta-path Schemes

In this set of experiments, based on the dataset described in Section 3.3, we first evaluate the performance of different kinds of relatedness over code snippets depicted by different sets of meta-path schemes. In the experiments, given a specific set of meta-path schemes, we use *snippet2vec* to learn the latent representations of the nodes (i.e., code snippets) in the HIN, which are then fed to SVM to build the classification model for in-

Table 3.1: Performance indices of code snippet detection

Indices	Description
TP	# of code snippets correctly classified as insecure
TN	# of code snippets correctly classified as secure
FP	# of code snippets mistakenly classified as insecure
FN	# of insecure mistakenly classified as secure
$Precision$	$TP/(TP + FP)$
$Recall/TPR$	$TP/(TP + FN)$
ACC	$(TP + TN)/(TP + TN + FP + FN)$
$F1$	$2 \times Precision \times Recall/(Precision + Recall)$

secure code snippet detection. For SVM, we use LibSVM and the penalty is empirically set to be 10 while other parameters are set by default. As described in Section 3.2, we generate four sets of meta-path schemes (denoted as $\mathcal{S}_1, \mathcal{S}_2, \mathcal{S}_3,$ and \mathcal{S}_4) for *snippet2vec* to learn the node representations in the HIN. We conduct 10-fold cross validations for evaluation. The performances of four different sets of meta-path schemes (i.e., \mathcal{S}_1 - \mathcal{S}_4) in comparison with nine individual meta-paths (i.e., $PID1$ - $PID9$) in insecure code snippet detection are shown in Table 4.1.

Table 3.2: Detection Results of different meta-paths

ID	Meta-paths included	Precision	Recall	ACC	F1
\mathcal{S}_1	(PID1,PID2,PID6)	0.9065	0.8887	0.8883	0.8975
\mathcal{S}_2	(PID1,PID3,PID7)	0.8899	0.8678	0.8682	0.8787
\mathcal{S}_3	(PID1,PID4,PID8)	0.9028	0.8834	0.8834	0.8930
\mathcal{S}_4	(PID1,PID5,PID9)	0.8922	0.8709	0.8710	0.8814
\mathcal{S}'_5	(PID1)	0.8795	0.8561	0.8562	0.8676
\mathcal{S}'_6	(PID2)	0.8340	0.7988	0.8018	0.8160
\mathcal{S}'_7	(PID3)	0.8017	0.7657	0.7668	0.7833
\mathcal{S}'_8	(PID4)	0.8463	0.8179	0.8180	0.8318
\mathcal{S}'_9	(PID5)	0.8312	0.8001	0.8006	0.8153
\mathcal{S}'_{10}	(PID6)	0.8449	0.8119	0.8145	0.8281
\mathcal{S}'_{11}	(PID7)	0.8108	0.7708	0.7748	0.7903
\mathcal{S}'_{12}	(PID8)	0.8020	0.7642	0.7664	0.7826
\mathcal{S}'_{13}	(PID9)	0.7897	0.7518	0.7532	0.7703

From Table 4.1, we can see that different sets of meta-path schemes indeed show

different performances in insecure code snippet detection, since each of them represents specific semantics in insecure code snippet detection. From Table 4.1, we can also observe that: (1) *PID1* outperforms the other individual meta-paths (i.e., *PID2–PID9*), which indicates that the semantics of this meta-path reflect the problem of insecure code snippet detection better than the others. (2) The meta-paths of *PID2*, *PID4*, *PID6*, and *PID8* perform better than *PID3*, *PID5*, *PID7*, and *PID9* respectively; the reason behind this is that the code snippets posted in the answer threads are more likely to be reused by the developers than the ones posted in question threads, and thus they have closer connections. (3) Obviously, \mathcal{S}_1 , \mathcal{S}_2 , \mathcal{S}_3 , and \mathcal{S}_4 utilizing different meta-paths built from HIN are more expressive than each individual meta-path (i.e., *PID1–PID9*) in depicting the code snippets in Stack Overflow and thus achieve better detection performance. It will be interested to see the detection performance if different sets of meta-paths are further aggregated. This will be evaluated in the next set of experiments.

Comparisons with Different Network Embedding Models

In this set of experiments, we evaluate our developed system *ICSD* integrating our proposed method described in Section 3.2 by comparisons with several network representation learning methods: (1) DeepWalk [28] and LINE [31] which are homogeneous network embedding methods; and (2) metapath2vec [33] which is a HIN embedding model. For DeepWalk and LINE, we ignore the heterogeneous property of HIN and directly feed the HIN for representation learning; in metapath2vec, a walk path will be generated only based on a single meta-path scheme; while in our proposed *snippet2vec*, a walk path will be guided by a set of different meta-path schemes. The parameter settings used for *snippet2vec* are in line with typical values used for the baselines: vector dimension $d = 200$ (LINE: 200 for each order (1st- and 2nd-order)), walks per node $r = 10$, walk length $l = 80$, and window size $w = 10$. To facilitate the comparisons, we use the experimental procedure as in [28, 31, 33]: we randomly select a portion of labeled code snippets described in Section 3.3 (ranging from 10% to 90%) for training and the remaining ones for testing. For all the baselines, the SVM is used as the classification model; for *ICSD*, based on the four given sets of meta-path schemes, it will generate four different kinds of node representations using *snippet2vec* and then use multi-view fusion classifier proposed in Section 3.2 to train the classification model. Table 3.3 il-

Table 3.3: Comparisons with other network representation learning methods in insecure code snippet detection

Metric	Method	10%	30%	50%	70%	90%
<i>ACC</i>	DeepWalk	0.6085	0.6550	0.6810	0.7148	0.7279
	LINE	0.6347	0.6847	0.7268	0.7475	0.7732
	metapath2vec	0.7772	0.8197	0.8490	0.86632	0.8826
	<i>ICSD</i>	0.7973	0.8384	0.8771	0.8953	0.9123
<i>F1</i>	DeepWalk	0.6308	0.6764	0.7006	0.7329	0.7461
	LINE	0.6569	0.7047	0.7451	0.7644	0.7892
	metapath2vec	0.7932	0.8332	0.8609	0.8765	0.8921
	<i>ICSD</i>	0.8121	0.8508	0.8871	0.9036	0.9197

illustrates the detection results of different network representation learning models. From Table 3.3, we can see that *ICSD* integrating the proposed *snippet2vec* model consistently and significantly outperforms all baselines for insecure code snippet detection in terms of *ACC* and *F1*. That is to say, *snippet2vec* learns significantly better code snippet representation than current state-of-the-art methods. The success of *snippet2vec* lies in the proper consideration and accommodation of the heterogeneous property of HIN (i.e., the multiple types of nodes and relations), and the advantage of random walk guided by different meta-paths for sampling the node paths. Furthermore, from Table 4.1 and Table 3.3, we can also observe that using the multi-view fusion classifier proposed in Section 3.2 to aggregate different node representations learned based on different sets of meta-graph schemes can significantly improve the detection performance.

Comparisons with Traditional Machine Learning Methods

In this set of experiments, based on the dataset described in Section 3.3, we compare *ICSD* which integrates our proposed method with other traditional machine learning methods by 10-fold cross validations. For these methods, we construct three types of features: *f-1*: content-based features (i.e., keywords extracted from code snippets described in Section 3.2); *f-2*: two relation-based features associated with code snippets (i.e., *R1* and *R2* introduced in Section 3.2); *f-3*: augmented features of content-based features and *R1-R2*. Based on these features, we consider two typical classification models, i.e., Naive Bayes (NB) and SVM. The experimental results are illustrated in

Table 3.4. From the results we can observe that feature engineering (*f-3*: concatenation of different features altogether) helps the performance of machine learning, but *ICSD* added the knowledge represented as HIN significantly outperforms other baselines. This again demonstrates that, to detect the insecure code snippets, *ICSD* utilizing both code content and social coding properties represented by the HIN is able to build the higher-level semantic and structural connection between code snippets with a more expressive and comprehensive view and thus achieves better detection performance.

Table 3.4: Comparisons of other machine learning methods

Metric	NB			SVM			<i>ICSD</i>
	<i>f-1</i>	<i>f-2</i>	<i>f-3</i>	<i>f-1</i>	<i>f-2</i>	<i>f-3</i>	
<i>ACC</i>	0.7757	0.6597	0.8161	0.8064	0.6904	0.8494	0.9118
<i>F1</i>	0.8002	0.6914	0.8372	0.8278	0.7208	0.8675	0.9190

Evaluation of Parameter Sensitivity, Scalability, and Stability

In this set of experiments, based on the dataset described in Section 3.3, we first conduct the **sensitivity** analysis of how different choices of parameters (i.e., walks per node r , walk length l , vector dimension d , and neighborhood size w) will affect the performance of *ICSD* in insecure code snippet detection. From the results shown in Figure 3.5(a) and 3.5(b), we can observe that the balance between computational cost (number of walks per node r and walk length l in x -axis) and efficacy (F1 in y -axis) can be achieved when $r = 10$ and $l = 60$ for insecure code snippet detection. We also examine how vector dimension (d) and neighborhood size (w) affect the performance. As shown in Figure 3.5(c), we can see that the performance tends to be stable once d reaches around 300; similarly, from Figure 3.5(d) we can find that the performance inclines to be stable when w increases to around 8. Overall, *ICSD* is not strictly sensitive to these parameters, and is able to reach high performance under a cost-effective parameter choice.

We then further evaluate the **scalability** of *ICSD* which can be parallelized for optimization. We run the experiments using the default parameters with different number of threads (i.e., 1, 4, 8, 12, 16), each of which utilizes one CPU core. Figure 3.6(a) shows the speed-up of *ICSD* deploying multiple threads over the single-threaded case, which reveals that the model achieves acceptable sub-linear speed-ups as the line is close to

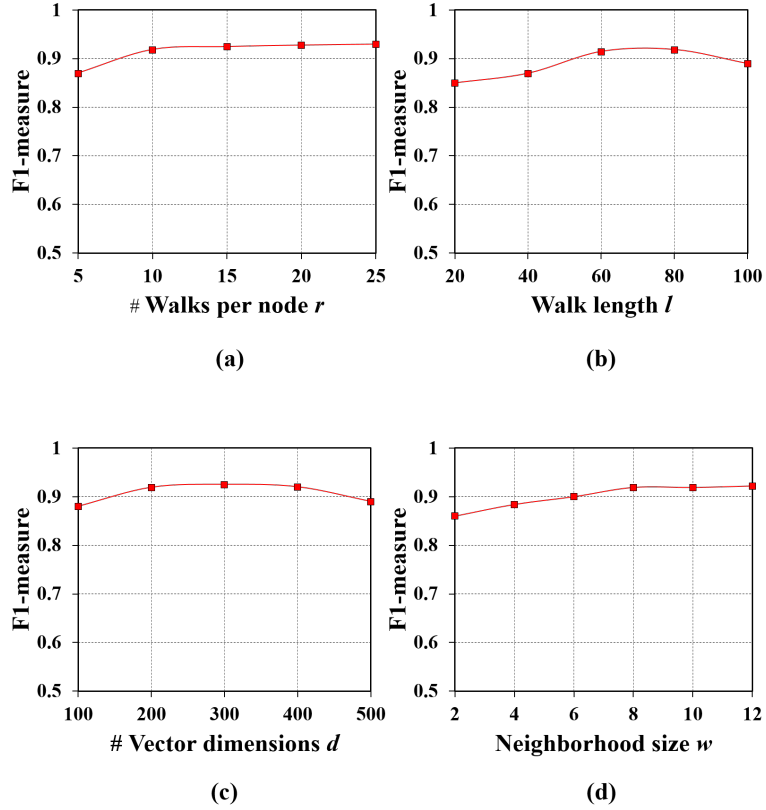


Figure 3.5: Parameter sensitivity evaluation.

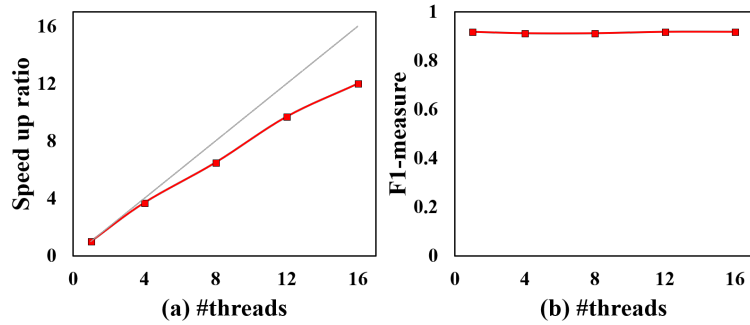


Figure 3.6: Scalability evaluation.

the optimal line; while Figure 3.6(b) shows that the performance remains stable when using multiple threads for model updating. Overall, the proposed system are efficient and scalable for large-scale HIN with large numbers of nodes. For **stability** evaluation, Figure 3.7 shows the receiver operating characteristic (ROC) curves of *ICSD* based on the 10-fold cross validations; it achieves an average 0.9094 TP rate (*TPR*) at the 0.0851

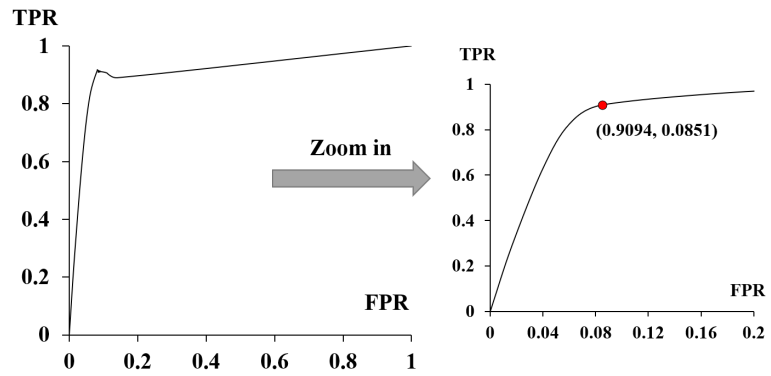


Figure 3.7: Stability evaluation.

FP rate (*FPR*) for insecure code snippet detection.

Chapter 4

iTrustSO: Code-to-code Sequential Model over HIN for Insecure Code Snippet Detection

4.1 System Architecture

The system architecture of *iTrustSO* is shown in Figure 4.1, which is developed for insecure code snippet detection in Stack Overflow. It consists of the following major components:

- **Data collector, Feature extractor and HIN constructor.** *iTrustSO* uses the same modules with *ICSD* which are displayed in Figure 3.1, except that *iTrustSO* do not consider the keyword entity.

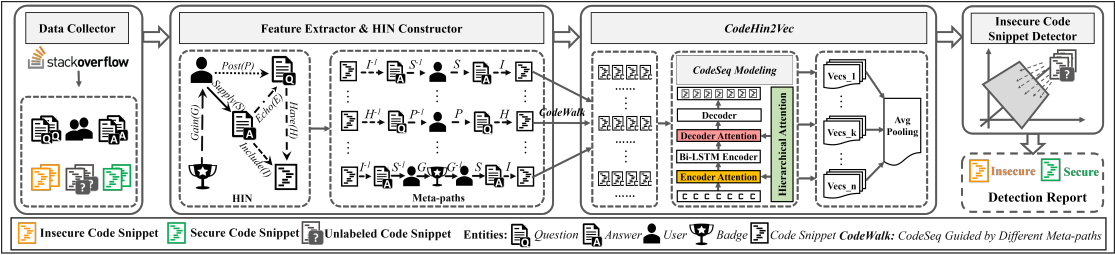


Figure 4.1: System architecture of *iTrustSO*.

- **CodeHin2Vec.** Based on the built meta-path schemes, to devise a comprehensive solution to seamlessly combine both node content, a new network embedding model *CodeHin2Vec* is proposed to learn the low-dimensional representations for the nodes in HIN, which are capable to use the content feature vector to represent each code snippet in the Insecure code snippet detector. Using *CodeHin2vec*, the mapped feature vectors of code snippets, encoding the information of code content and HIN-based relations, will be fed to a Support Vector Machine (SVM) to train the classification model, based on which the unlabeled code snippets can be predicted if they are insecure or not.

4.2 Proposed Method

In this section, we present the detailed approaches of how we represent the code snippets in Stack Overflow, and how we solve the insecure code snippet detection problem based on the representation.

Feature Extraction

iTrustSO uses the similar feature extraction module with *ICSD*(See Section 3.2 or details). The only different is that in *iTrustSO*, it do not extract keywords as content feature. We use a new approach to represent content feature.

HIN Constructor

4.2 HIN Constructor in *iTrustSO* is very similar with *ICSD*(See Section 3.2 or details).

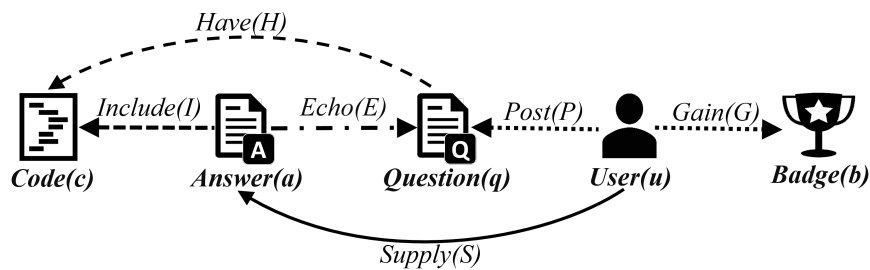


Figure 4.2: Network schema for HIN in our application.

We do not use keywords in heterogeneous information network. The network schema for HIN in *iTrustSO* is shown in Figure

CodeHin2Vec: HIN Representation Learning

To devise a comprehensive solution to seamlessly combine both node (i.e., code snippet) content and HIN-based relations for insecure code snippet detection in Stack Overflow, we gain further insight into Stack Overflow data; as observed in our previous work [49], the HIN-based neighborhood relationships among code snippets can be represented by the code sequences (denoted as CodeSeq) based on different meta-paths. In this way, the generated CodeSeqs can preserve both semantic and structure information of HIN. To further couple CodeSeqs with code content, a straightforward yet novel way is to use the content feature vector x_c to represent each code snippet in the CodeSeq. To this end, the representation learning of code snippets can be viewed as a sequence modeling task. As LSTM has shown significant improvement in language modeling [50], we leverage its power to seamlessly integrate code content and HIN structure into hidden layer vectors that can be used as the representations of code snippets [51].

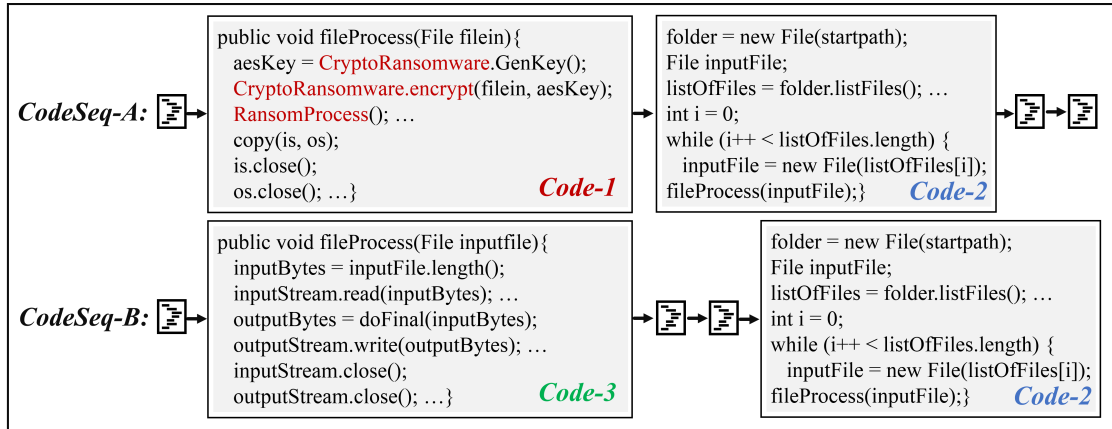


Figure 4.3: Different contexts among code snippets.

Although it is promising to comprehensively utilize LSTM to learn the mapping from the code content sequence to code identity sequence, it still faces the following two challenges: (1) *word2vec* assigns each code snippet a static embedding vector based on code content which is not context-aware to different sequences it interacts with. For example, as illustrated in Figure 4.3, guided by the designed meta-paths, we may generate *CodeSeq-A* and *CodeSeq-B*. With function *fileProcess* defined, *Code-1* in *CodeSeq-A*

performs as file encryption for Ransomware while *Code-3* in *CodeSeq-B* implements the regular file reading and writing; in this respect, even though *Code-2* listed in both sequences calls the same function *fileProcess*, its embedding vector should be significantly different which may demonstrate insecure potential when interacting with *Code-1* and normal aspect when related to *Code-3*. LSTM is known to learn the sequential dependencies [52], but strict to align the positions of the input sequence; therefore, contextualized code content embeddings may help to refine the hidden-layer information in the early stage. (2) Since LSTM needs to read the whole input sequence to further generate the output sequence, its performance using a basic encoder-decoder architecture may degrade as the length of an input sequence increases [53, 50] which may in turn degenerate the representations learned from hidden layers, especially in our case that code sequences are much longer than the sentences.

Attention mechanism has shown remarkable effectiveness in various sequence modeling tasks, allowing models to learn alignments between different modalities [54, 50, 55, 56]. In this work, to address the challenges above, we propose *CodeHin2Vec* to elaborate a hierarchical attention mechanism into LSTM to fully exploit code content and HIN structure to learn effective representations of code snippets, which first generates CodeSeqs based on the walk paths guided by different meta-paths; and then leverages LSTM with hierarchical attention mechanism for CodeSeq modeling.

CodeSeq generation guided by different meta-paths. Given a source node v_j in a homogeneous network, the traditional random walk is a stochastic process with random variables $v_j^1, v_j^2, \dots, v_j^k$ such that v_j^{k+1} is a node chosen at random from the neighbors of node v_k . The transition probability $p(v_j^{i+1}|v_j^i)$ at step i is the normalized probability distributed over the neighbors of v_j^i by ignoring their node types. However, this mechanism is unable to capture the semantic and structural correlations among different types of nodes in a HIN. In our application, given a HIN $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with schema $\mathcal{T}_{\mathcal{G}} = (\mathcal{A}, \mathcal{R})$, and a set of different meta-paths $\mathbf{P} = \{\mathcal{P}_j\}_{j=1}^n$, each of which is in the form of $A_1 \rightarrow \dots A_t \rightarrow A_{t+1} \dots \rightarrow A_l$, we put a random walker to traverse the HIN. The random walker first randomly chooses a meta-path \mathcal{P}_k from \mathbf{P} and the transition

probabilities at step i are defined as follows:

$$p(v^{i+1}|v_{A_t}^i, \mathbf{P}) = \begin{cases} \frac{\lambda}{|\mathbf{P}| |N_{A_{t+1}}(v_{A_t}^i)|} & \text{if } (v^{i+1}, v_{A_t}^i) \in \mathcal{E}, \phi(v_{A_t}^i) = A_c, \phi(v^{i+1}) = A_{t+1} \\ \frac{1}{|N_{A_{t+1}}(v_{A_t}^i)|} & \text{if } (v^{i+1}, v_{A_t}^i) \in \mathcal{E}, \phi(v_{A_t}^i) \neq A_c, \\ & \phi(v^{i+1}) = A_{t+1}, (A_t, A_{t+1}) \in \mathcal{P}_k \\ 0 & \text{otherwise} \end{cases} \quad (4.1)$$

where ϕ is the node type mapping function, $N_{A_{t+1}}(v_{A_t}^i)$ denotes the A_{t+1} -type neighborhood of node $v_{A_t}^i$, A_c is entity type of *Code*, and λ is the number of meta-paths starting with $A_c \rightarrow A_{t+1}$. For each walk path, the nodes whose entity types are not *Code* will be removed; then the remaining ones form a CodeSeq, whose element is represented by the content feature vector \mathbf{x}_c . In such way, given walk path length l , a CodeSeq is presented as $(\mathbf{x}_{c_1}, \mathbf{x}_{c_2}, \dots, \mathbf{x}_{c_l})$.

CodeSeq modeling with LSTM. LSTM learns a mapping from an input sequence to an output sequence. As intermediate states, a hidden vector is generated for each timestep; we can extract it as the embedding vector for the input at that timestep. In our application, we employ an encoder-decoder LSTM architecture [57] for CodeSeq modeling in which two attention layers are elaborately added to improve the quality of representation learning (as illustrated in Figure 4.4).

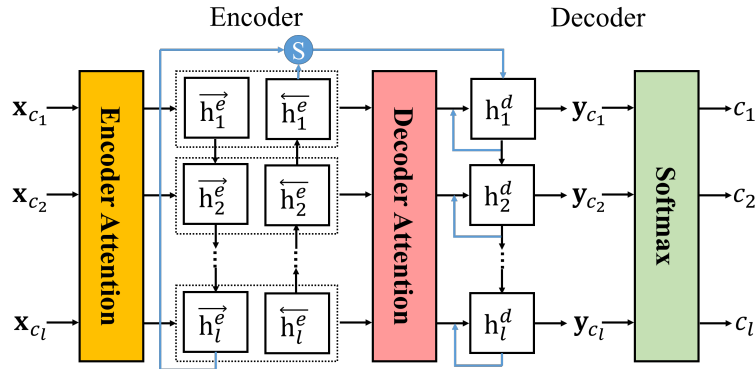


Figure 4.4: Architecture of LSTM using hierarchical attention.

Encoder attention: Resting on all the content vectors in the input sequence, the encoder attention layer computes the contextualized embedding for each code snippet

as a weighted sum where the weight, also called context score, assigned to each content vector is computed by a dot product of the corresponding pair of content vectors [54]. Specifically, given an input CodeSeq $(\mathbf{x}_{c_1}, \mathbf{x}_{c_2}, \dots, \mathbf{x}_{c_l})$, for any two code snippets c_t and c_i , the context score can be calculated as

$$\mathcal{S}(\mathbf{x}_{c_t}, \mathbf{x}_{c_i}) = \mathbf{x}_{c_t} \top \mathbf{x}_{c_i}, \quad (4.2)$$

where \top denotes the dot product, and thus the contextualized embedding for code snippet c_t can be computed as

$$\tilde{\mathbf{x}}_{c_t} = \sum_{i=1}^l \frac{\exp(\mathcal{S}(\mathbf{x}_{c_t}, \mathbf{x}_{c_i}))}{\sum_{j=1}^l \exp(\mathcal{S}(\mathbf{x}_{c_t}, \mathbf{x}_{c_j}))} \mathbf{x}_{c_i}. \quad (4.3)$$

In this sense, a CodeSeq can be refined as $(\tilde{\mathbf{x}}_{c_1}, \tilde{\mathbf{x}}_{c_2}, \dots, \tilde{\mathbf{x}}_{c_l})$, which will be used as the actual input sequence.

Encoder: The encoder reads $(\tilde{\mathbf{x}}_{c_1}, \tilde{\mathbf{x}}_{c_2}, \dots, \tilde{\mathbf{x}}_{c_l})$ through the hidden layer function \mathcal{H} so that each hidden layer vector \mathbf{h}_t^e at timestep t can be denoted as

$$\mathbf{h}_t^e = \mathcal{H}(\tilde{\mathbf{x}}_{c_t}, \mathbf{h}_{t-1}^e), \quad (4.4)$$

where \mathcal{H} is implemented using memory cells to store information, which can be formulated as the following composite functions [58]:

$$\mathbf{i}_t = \sigma(\mathbf{W}_{xi}\tilde{\mathbf{x}}_{c_t} + \mathbf{W}_{hi}\mathbf{h}_{t-1}^e + \mathbf{W}_{ci}\mathbf{c}_{t-1} + \mathbf{b}_i) \quad (4.5)$$

$$\mathbf{f}_t = \sigma(\mathbf{W}_{xf}\tilde{\mathbf{x}}_{c_t} + \mathbf{W}_{hf}\mathbf{h}_{t-1}^e + \mathbf{W}_{cf}\mathbf{c}_{t-1} + \mathbf{b}_f) \quad (4.6)$$

$$\mathbf{c}_t = \mathbf{f}_t \circ \mathbf{c}_{t-1} + \mathbf{i}_t \circ \tanh(\mathbf{W}_{xc}\tilde{\mathbf{x}}_{c_t} + \mathbf{W}_{hc}\mathbf{h}_{t-1}^e + \mathbf{b}_c) \quad (4.7)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}_{xo}\tilde{\mathbf{x}}_{c_t} + \mathbf{W}_{ho}\mathbf{h}_{t-1}^e + \mathbf{W}_{co}\mathbf{c}_{t-1} + \mathbf{b}_o) \quad (4.8)$$

$$\mathbf{h}_t^e = \mathbf{o}_t \circ \tanh(\mathbf{c}_t) \quad (4.9)$$

where σ is the logistic sigmoid function, \mathbf{i}_t , \mathbf{f}_t , \mathbf{o}_t , \mathbf{c}_t are the input gate, forget gate, output gate, and cell activation vectors respectively, \mathbf{W} s are the weight matrices, \mathbf{b} s are the bias vectors, and \circ is the point-wise product between two vectors. Since the input sequence has no direction, in order to learn both the forward and backward sequential dependency information, we utilize bidirectional encoder so that hidden layer vector \mathbf{h}_t^e

at timestep t can be concatenated as $\mathbf{h}_t^e = [\overrightarrow{\mathbf{h}}_t^e; \overleftarrow{\mathbf{h}}_t^e]$. After forward and backward reading CodeSeq, the concatenation of the last two hidden states $[\overrightarrow{\mathbf{h}}_1^e; \overleftarrow{\mathbf{h}}_1^e]$ is used as the summary vector \mathbf{s} of the whole input sequence.

Decoder attention: The decoder attention layer exploits all the hidden states of the encoder to compute the aligned and joint information as the context vector [50, 55], which is integrated with the summary vector \mathbf{s} to extract the target code identity. Similar to the encoder attention, the alignment scores need to be first defined to formulate such context vector as a weighted sum. Note that, unlike the dot product attention, decoder attention should allow the gradient of the cost function to be backpropagated through [50]. We accordingly use a simple feed-forward neural network to compute the alignment score

$$\alpha_t = \mathbf{W}_\alpha^2 \text{ReLU}(\mathbf{W}_\alpha^1 \mathbf{h}_t^e + \mathbf{b}_\alpha^1) + \mathbf{b}_\alpha^2, \quad (4.10)$$

where \mathbf{W}_α s and \mathbf{b}_α s denote the weight matrices and the bias vectors, and the alignment score vector α_t trained by all the other hidden states of the encoder reflects the importance of \mathbf{h}_t^e in generating \mathbf{y}_t . The context vector for \mathbf{h}_t^e can thus be

$$\tilde{\mathbf{h}}_t^e = \sum_{i=1}^l \frac{\exp(\alpha_{t,i})}{\sum_{j=1}^l \exp(\alpha_{t,j})} \mathbf{h}_i^e. \quad (4.11)$$

Decoder: The decoder takes the summary vector \mathbf{s} as input (i.e., $\mathbf{h}_0^d = \mathbf{s}$) and generates a sequence of target hidden states; each hidden state \mathbf{h}_t^d at timestep t can be calculated as

$$\mathbf{h}_t^d = \mathcal{H}(\mathbf{0}, \mathbf{h}_{t-1}^d), \quad (4.12)$$

where $\mathbf{0}$ is an all-zero vector. Given the target hidden state \mathbf{h}_t^d and the context vector $\tilde{\mathbf{h}}_t^e$, we concatenate them to formulate an attentional hidden state $\tilde{\mathbf{h}}_t^d = [\tilde{\mathbf{h}}_t^e; \mathbf{h}_t^d]$ [55]. Accordingly, the output vector $\mathbf{y}_t \in \mathbb{R}^{|\mathcal{V}|}$ can be generated as follows [58]

$$\mathbf{y}_t = \sigma(\mathbf{W}_{hy} \tilde{\mathbf{h}}_t^d + \mathbf{b}_y). \quad (4.13)$$

\mathbf{y}_t is capable to predict the real code snippet c_t through a softmax layer. The sequence

loss \mathcal{L} is adopted to measure the correctness of decoding, which is computed as

$$\mathcal{L} = - \sum_{t=1}^l \log p(c_t | \mathbf{y}_t) = - \sum_{t=1}^l \log \frac{\exp(y_t^{c_t})}{\sum_{i=1}^{|\mathcal{V}|} \exp(y_t^{c_i})}. \quad (4.14)$$

The weights can be efficiently calculated with backpropagation through time [59, 58], and the LSTM model can then be trained using Adam optimization algorithm.

For the generated CodeSeqs guided by different meta-paths, each code snippet may appear in multiple CodeSeqs. Suppose that code snippet c_t exists in $|c_t|$ CodeSeqs, by doing avg pooling over all \mathbf{h}_i^e 's for code snippet c_t , $\forall i = 1, \dots, |c_t|$, we obtain an embedding \mathbf{h} for each code snippet

$$\mathbf{h} = \text{avgPooling}(\{\mathbf{h}_i^e : i = 1, \dots, |c_t|\}). \quad (4.15)$$

Using *CodeHin2Vec*, the mapped feature vectors of code snippets, encoding the information of code content and HIN-based relations, can be fed to a classifier to train the classification model, based on which the unlabeled code snippets can be predicted if they are insecure or not.

Algorithm 2 *iTrustSO* – Code-to-code Sequential Model over HIN for Insecure Code Snippet Detection

Input: The HIN $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, network schema $\mathcal{T}_{\mathcal{G}} = (\mathcal{A}, \mathcal{R})$, number of meta-paths schemes m , number of walk paths per node r , walk length l , and vector dimension d , training data set D_t , testing data set D_e **Output:** The labels for the testing code snippets.

- 1: **for** $i = 1 \rightarrow |D_t|$ **do**
 - 2: Generate content feature vector x_{ci}
 - 3: **for** $j = 1 \rightarrow r$ **do**
 - 4: get l -length random walks using Eq. 4.1
 - 5: **end for**
 - 6: **end for**
 - 7: Generate $y_t \in \mathbb{R}^{|\mathcal{V}| \times d}$ using LSTM in Eq. 4.13
 - 8: Train SVM using \mathbf{y}_{D_t}
 - 9: **for** $k = 1 \rightarrow |D_e|$ **do**
 - 10: Generate the label f_k using trained SVM
 - 11: **end for**
 - 12: return \mathbf{f}
-

4.3 Experimental Results and Analysis

In this section, we fully evaluate the performance of *iTrustSO* in insecure code snippet detection. We consider Java programming language for Android app as a case study. Based on our prior work *ICSD* [49], in this paper, we further expand our data collection and annotation from Stack Overflow: (1) using our developed crawlers, we collect 505,548 question threats and 719,430 answer threats posted by 229,394 users including 821,792 code snippets, through March 2010 to October 2018; (2) we also expand our annotated data in [49] to finally obtain 21,989 labeled code snippets (10,013 are insecure while 11,976 are secure) as the ground truth to evaluate different detection methods. To quantitatively validate the effectiveness of different methods, we use accuracy (ACC) and F1 measure (F1) as the performance measures.

Evaluation of Different Meta-paths

In this set of experiments, given a specific meta-path scheme, we use a basic LSTM to learn the latent representations of code snippets in HIN, which is then fed to SVM for detection. Here we perform 10-fold cross validations for evaluation. The experimental results are shown in Table 4.1, from which we can see that different meta-paths indeed show different performances: (1) *PID1*, *PID3*, *PID5*, and *PID7* perform better than *PID2*, *PID4*, *PID6*, and *PID8*; the reason behind this is that the code snippets posted in the answer threads are more likely to be reused by the developers than the ones posted in question threads, and thus they have closer connections. (2) *PID3* outperforms the others, which indicates that its semantics reflecting the insecure code snippet detection problem is better than the others. (3) *PID9* using different meta-paths is more expressive than individuals in depicting the code snippets and thus achieve better performance.

Evaluation of Attentions

In this set of experiments, we'd like to assess whether the hierarchical attention mechanism devised in our model is meaningful for representation learning. To this end, we explore the performances of basic LSTM without attention (LSTM-b), LSTM with encoder attention (LSTM-e), LSTM with decoder attention (LSTM-d), and *CodeHin2Vec*. The better detection result implies that the learn representations take better advantage

Table 4.1: Detection Results of different meta-paths

ID	Meta-paths included	Recall	Precision	ACC	F1
PID1	–	0.8481	0.7956	0.8316	0.8210
PID2	–	0.8098	0.7491	0.7899	0.7783
PID3	–	0.8596	0.8119	0.8454	0.8351
PID4	–	0.8344	0.7769	0.8155	0.8046
PID5	–	0.8605	0.8086	0.8437	0.8337
PID6	–	0.8140	0.7588	0.7975	0.7854
PID7	–	0.8042	0.7444	0.7851	0.7731
PID8	–	0.7843	0.7203	0.7631	0.7509
PID9	$\mathbf{P} = (\text{PID1}, \dots, \text{PID8})$	0.8785	0.8415	0.8693	0.8596

of the corresponding sequence learning architecture. From the results illustrated in Figure 4.5, we have the following observations: (1) LSTM-e and LSTM-d with single attention layer both outperform LSTM-b without attention; (2) *CodeHin2Vec* achieves the most promising performance for fully utilizing the contextualized input embeddings and the aligned information from the hidden states of the encoder. In other words, *CodeHin2Vec* has potential to let LSTM learn better sequential dependencies and code better with the sequence extraction from the proper context information, which in turn generates better representations for code snippets.

Evaluation of *CodeHin2Vec*

Here, *CodeHin2Vec* is evaluated by comparisons with several representation learning methods: (1) word2vec [44] is a baseline using code content information; (2) Deep-

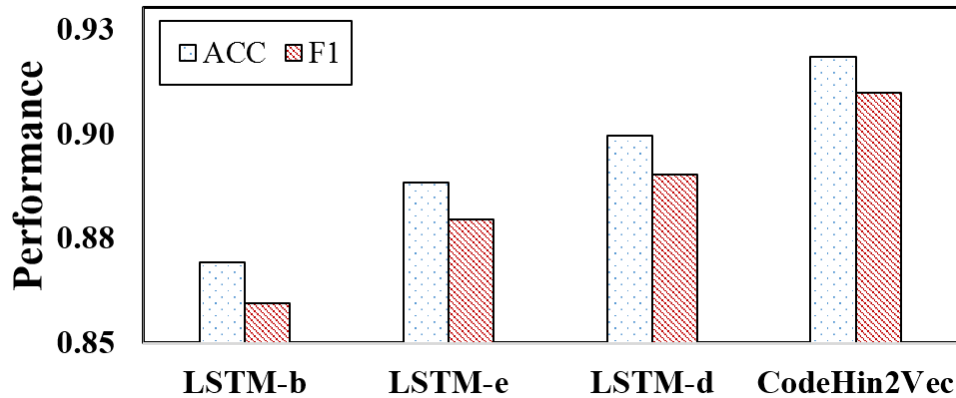


Figure 4.5: Attention evaluation.

Table 4.2: Comparisons of *CodeHin2Vec* with other network representation learning methods in insecure code snippet detection

Metric	Method	Feature	10%	30%	50%	70%	90%
<i>ACC</i>	word2vec	Content	0.6554	0.6989	0.7379	0.7725	0.7753
	DeepWalk	Relation	0.6263	0.6678	0.7087	0.7349	0.7430
	metapath2vec	Relation	0.7241	0.7562	0.7898	0.8035	0.8312
	TADW	Content&Relation	0.7659	0.7902	0.8144	0.8394	0.8537
	ICSD	Content&Relation	0.8026	0.8487	0.8783	0.8968	0.9107
	<i>CodeHin2Vec</i>	Content&Relation	0.7983	0.8630	0.8752	0.8975	0.9223
<i>F1</i>	word2vec	Content	0.6292	0.6756	0.7166	0.7519	0.7560
	DeepWalk	Relation	0.6023	0.6439	0.6871	0.7139	0.7233
	metapath2vec	Relation	0.7005	0.7356	0.7711	0.7853	0.8147
	TADW	Content&Relation	0.7446	0.7717	0.7977	0.8239	0.8390
	ICSD	Content&Relation	0.7855	0.8338	0.8662	0.8866	0.9015
	<i>CodeHin2Vec</i>	Content&Relation	0.7831	0.8502	0.8624	0.8873	0.9160

Walk [60] is a homogeneous network embedding method leveraging relation information; (3) metapath2vec [61] is a HIN embedding model utilizing HIN-based relations; (4) TADW [62] considers both content and relation information for homogeneous network representation learning; (5) ICSD [49] takes content and relation into account in HIN. For DeepWalk and TADW, we ignore the heterogeneous property of HIN and directly feed the HIN for embedding; in metapath2vec, a walk path is generated based on a single meta-path scheme; in ICSD, code content is extracted as keywords to be devised to HIN. The parameter settings used for *CodeHin2Vec* are in line with typical values used for the baselines: content dimension $c = 300$, vector dimension $d = 200$, walks per node $r = 10$, walk length $l = 80$ (TADW: walk steps are set to 2), and window size $w = 10$. To facilitate the comparisons, we randomly select a portion of labeled code snippets (ranging from 10% to 90%) for training and the remaining ones for testing. SVM is used as the classification model for all the methods. Table 4.2 illustrates the detection results: *CodeHin2Vec* outperforms all baselines in terms of *ACC* and *F1* in most cases. That is to say, *CodeHin2Vec* learns significantly better code snippet representation than current state-of-the-art methods. The success of *CodeHin2Vec* lies in the seamless integration of code content with HIN-based relations for representation learning, which leverages the advantage of (1) CodeSeq generation based on the different meta-paths and (2) the CodeSeq modeling power of LSTM using hierarchical attentions.

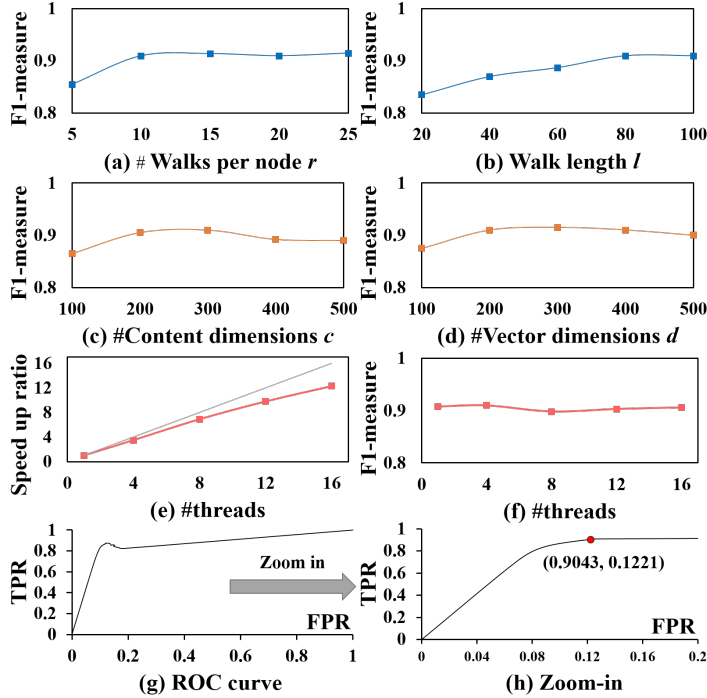


Figure 4.6: Parameter sensitivity evaluation.

Evaluation of Parameters

In this set of experiments, we first conduct the **sensitivity** analysis of how different choices of parameters will affect the performance of *CodeHin2Vec*. From the results shown in Figure 4.6(a) and 4.6(b), we can observe that the balance between computational cost (number of walks per node r and walk length l in x -axis) and efficacy (F1 in y -axis) can be achieved when $r \geq 10$ and $l \geq 80$. As shown in Figure 4.6(c), we can see that the performance tends to be stable once content vector dimension c reaches around 200 to 300; similarly, from Figure 4.6(d) we can find that the performance inclines to be stable when vector dimensions d increases to around 200 to 400. Overall, *CodeHin2Vec* is not strictly sensitive to these parameters, and is able to reach high performance under a cost-effective parameter choice. We then further evaluate the **scalability** of *CodeHin2Vec* which can be parallelized for optimization. We run the experiments using the default parameters with different number of threads (i.e., 1, 4, 8, 12, 16), each of which utilizes one CPU core. Figure 4.6(e) shows the speed-up of *CodeHin2Vec* deploying multiple threads over the single-threaded case, which reveals that the model achieves acceptable sub-linear speed-ups as the line is close to the optimal line; while Figure 4.6(f) shows that the performance remains stable when using multiple threads for model up-

dating. Overall, the proposed system are efficient and scalable for large-scale HIN with large numbers of nodes. For **stability** evaluation, Figure 4.6(g) shows the ROC curves of *CodeHin2Vec* based on the 10-fold cross validations; it achieves an average 0.9043 TPR at the 0.1221 FPR for detection.

Comparisons with Traditional Machine Learning Methods

In this set of experiments, *iTrustSO* is compared with other traditional machine learning methods. For these methods, we construct three types of features: *f-1*: content-based features (i.e., \mathbf{x}_c); *f-2*: two original relation-based features (i.e., *R1* and *R2*); *f-3*: augmented features of content-based features and *R1-R2*. Based on these features, we consider two typical classification models, i.e., Naive Bayes (NB) and SVM. The experimental results shown in Table 4.3 illustrates that feature engineering (*f-3*) helps the performance of machine learning, but *iTrustSO* leveraging the knowledge represented as HIN and the long-range influence among code snippets learned from LSTM with attentions significantly outperforms other baselines. This again demonstrates that, to detect the insecure code snippets, *iTrustSO* using *CodeHin2Vec* to seamlessly integrate node content with HIN relations is able to build the higher-level semantic and structural connection between code snippets with a more expressive and comprehensive view and thus achieves better detection performance.

Table 4.3: Comparisons of other machine learning methods

Metric	NB			SVM			<i>iTrustSO</i>
	<i>f-1</i>	<i>f-2</i>	<i>f-3</i>	<i>f-1</i>	<i>f-2</i>	<i>f-3</i>	
<i>ACC</i>	0.7493	0.6854	0.7952	0.7753	0.7034	0.8415	0.9184
<i>F1</i>	0.7284	0.6613	0.7834	0.7560	0.6793	0.8317	0.9098

Chapter 5

Conclusion and Future Work

To address the code security issue in modern social coding platforms, in this paper, we bring an important new insight to exploit social coding properties in addition to code content for automatic detection of insecure code snippets in Stack Overflow. To depict the code snippets, we not only analyze the code content, but also utilize various kinds of relations among users, badges, questions, answers and code snippets in Stack Overflow. To model the rich semantic relationships, we first introduce a structured HIN for representation and then use meta-path based approach to incorporate higher-level semantics to build up relatedness over code snippets. Later, we propose two different novel network embedding models named *Snippet2vec* and *CodeHin2Vec* for representation learning in HIN to automate the insecure code snippet detection in Stack Overflow. After that, a classifier is built for insecure code snippet detection. Though it's proposed for code security analysis, the embedding methods are general framework which are able to learn desirable node representation in HIN and thus can be further applied to various network mining tasks, such as node classification, clustering and similarity search. The experimental results based on the data collections from Stack Overflow demonstrate that the developed systems *ICSD* and *iTrustSO* integrating our proposed methods outperform alternative approaches in insecure code snippet detection.

In our future work, we will continue to improve our system to automate analysis in other social coding platform (e.g., Github, Reddit, etc) for insecure code snippet detection. On the other hand, the study such as computational cost and incremental learning over heterogeneous information networks is still worth exploring.

Publications

1. Yanfang Ye, Lingwei Chen, **Shifu Hou**, William Hardy, Xin Li. “DeepAM: a heterogeneous deep learning framework for intelligent malware detection” *In Knowledge and Information Systems, 2018*.
2. Yujie Fan, **Shifu Hou**, Yiming Zhang, Yanfang Ye, Melih Abdulhayoglu. “Gotcha - Sly Malware! Scorpion: A Metagraph2vec Based Malware Detection System”, *In ACM SIGKDD, 2018*.
3. Yanfang Ye, **Shifu Hou**, Lingwei Chen, Xin Li, Liang Zhao, Shouhuai Xu, Jiabin Wang, Qi Xiong. “ICSD: An Automatic System for Insecure Code Snippet Detection in Stack Overflow over Heterogeneous Information Network” *In ACSAC, 2018*.
4. **Shifu Hou**, Yanfang Ye, Yangqiu Song, Melih Abdulhayoglu. “Make Evasion Harder: An Intelligent Android Malware Detection System.” *In IJCAI, 2018*.
5. **Shifu Hou**, Yanfang Ye, Yangqiu Song, Melih Abdulhayoglu. “Hindroid: An intelligent android malware detection system based on structured heterogeneous information network” *In SIGKDD, 2017*.
6. Lingwei Chen, **Shifu Hou**, Yanfang Ye. “Securedroid: Enhancing security of machine learning-based detection against adversarial android malware attacks” *In Proceedings of the 33rd Annual Computer Security Applications Conference, 2017*.
7. William Hardy, Lingwei Chen, **Shifu Hou**, Yanfang Ye, Xin Li. “DL4MD: A deep learning framework for intelligent malware detection” *In DMIN, 2016*.
8. **Shifu Hou**, Aaron Saas, Lifei Chen, Yanfang Ye. “Deep4maldroid: A deep learning framework for android malware detection based on linux kernel system call graphs” *In WIW, 2016*.
9. **Shifu Hou**, Aaron Saas, Yanfang Ye, Lifei Chen. “Droiddelver: An android malware detection system using deep belief network based on api call blocks” *In International Conference on Web-Age Information Management , 2016*.
10. **Shifu Hou**, Lifei Chen, Egemen Tas, Igor Demihovskiy, Yanfang Ye. “Cluster-oriented ensemble classifiers for intelligent malware detection” *In IEEE ICSC , 2015*.

Bibliography

- [1] LUCINTEL (2017) “Growth Opportunities in the Global Software Market,” in <http://www.lucintel.com/software-market-2017.aspx>.
- [2] VASILESCU, B., V. FILKOV, and A. SEREBRENIK (2013) “StackOverflow and GitHub: Associations Between Software Development and Crowdsourced Knowledge,” in *International Conference on Social Computing (SocialCom)*, pp. 188–195.
- [3] OVERFLOW, S. (2018) “Stack Overflow,” in <https://stackoverflow.com/>.
- [4] STACKEXCHANGE (2018) “StackExchange Statistics,” in <https://stackexchange.com/sites#traffic>.
- [5] COOGLE, J., J. GAJJAR, and C. GRECO (2017) “StackInTheFlow: StackOverflow Search Engine,” in *VCU Capstone Design Expo Posters*.
- [6] ACAR, Y., M. BACKES, S. FAHL, D. KIM, M. L. MAZUREK, and C. STRANSKY (2016) “You Get Where You’re Looking For The Impact of Information Sources on Code Security,” in *IEEE Symposium on Security and Privacy (SP)*, pp. 289–305.
- [7] FISCHER, F., K. BOTTINGER, H. XIAO, C. STRANSKY, Y. ACAR, M. BACKES, and S. FAHL (2017) “Stack Overflow Considered Harmful? The Impact of Copy and Paste on Android Application Security,” in *IEEE Symposium on Security and Privacy (SP)*, pp. 121–136.
- [8] ATTACKFLOW (2017) “Watch Out For Insecure StackOverflow Answers,” in <https://www.attackflow.com/Blog/StackOverflow>.
- [9] YE, Y., T. LI, D. ADJEROH, and S. S. IYENGAR (2017) “A survey on malware detection using data mining techniques,” *ACM Computing Surveys (CSUR)*, **50**(3), p. 41.
- [10] DETERDING, S. (2012) “Gamification: designing for motivation,” *Interactions*, **19**(4), pp. 14–17.

- [11] SUN, Y., R. BARBER, M. GUPTA, C. C. AGGARWAL, and J. HAN (2011) “Co-author relationship prediction in heterogeneous bibliographic networks,” in *ASONAM*, IEEE, pp. 121–128.
- [12] SUN, Y., J. HAN, X. YAN, P. S. YU, and T. WU (2011) “Pathsim: Meta path-based top-k similarity search in heterogeneous information networks,” *VLDB*, **4**(11), pp. 992–1003.
- [13] CZYCZYN-EGIRD, D. and R. WOJSZCZYK (2016) “Determining the Popularity of Design Patterns Used by Programmers Based on the Analysis of Questions and Answers on Stackoverflow.com Social Network,” in *Communications in Computer and Information Science (CCIS)*, pp. 421–433.
- [14] LINARES-VASQUEZ, M., G. BAVOTA, M. D. PENTA, and R. OLIVETO (2014) “How Do API Changes Trigger Stack Overflow Discussions? A Study on the Android SDK,” in *ICPC*, pp. 83–94.
- [15] LEE, R. K.-W. and D. LO (2017) “GitHub and Stack Overflow: Analyzing developer interests across multiple social collaborative platforms,” in *International Conference on Social Informatics*, Springer, pp. 245–256.
- [16] LIU, X. and H. ZHONG (2018) “Mining StackOverflow for Program Repair,” in *IEEE 25th International Conference on Software Analysis, Evolution and Reengineering (SANER)*, pp. 118–129.
- [17] TREUDE, C., O. BARZILAY, and M.-A. STOREY (2011) “How do programmers ask and answer questions on the web?: Nier track,” in *33rd International Conference on Software Engineering (ICSE)*, pp. 804–807.
- [18] KAVALER, D., S. SIROVICA, V. HELLEDOORN, R. ARANOVICH, and V. FILKOV (2017) “Perceived Language Complexity in GitHub Issue Discussions and Their Effect on Issue Resolution,” in *ASE*, pp. 72–83.
- [19] CAVUSOGLU, H., Z. LI, and K.-W. HUANG (2015) “Can Gamification Motivate Voluntary Contributions? The Case of StackOverflow Q&A Community,” in *Proceedings of the 18th ACM conference companion on computer supported cooperative work & social computing*, pp. 171–174.
- [20] AHMED, T. and A. SRIVASTAVA (2017) “Understanding and evaluating the behavior of technical users. A study of developer interaction at StackOverflow,” *Hum. Cent. Comput. Inf. Sci.*, **7**(8).
- [21] CALEFATO, F., F. LANUBILE, and N. NOVIELLI (2018) “How to ask for technical help? Evidence-based guidelines for writing questions on Stack Overflow,” *Information and Software Technology*, **94**, pp. 186–207.

- [22] SHI, C., Y. LI, J. ZHANG, Y. SUN, and S. Y. PHILIP (2017) “A survey of heterogeneous information network analysis,” *TKDE*, **29**(1), pp. 17–37.
- [23] WANG, C., Y. SONG, H. LI, M. ZHANG, and J. HAN (2015) “Knowsim: A document similarity measure on structured heterogeneous information networks,” in *ICDM*, IEEE, pp. 1015–1020.
- [24] FAN, Y., Y. ZHANG, Y. YE, X. LI, and W. ZHENG (2017) “Social Media for Opioid Addiction Epidemiology: Automatic Detection of Opioid Addicts from Twitter and Case Studies,” in *CIKM*, ACM, pp. 1259–1267.
- [25] FAN, Y., Y. ZHANG, Y. YE, and X. LI (2018) “Automatic Opioid User Detection from Twitter: Transductive Ensemble Built on Different Meta-graph Based Similarities over Heterogeneous Information Network.” in *IJCAI*, pp. 3357–3363.
- [26] HOU, S., Y. YE, Y. SONG, and M. ABDULHAYOGLU (2017) “HinDroid: An Intelligent Android Malware Detection System Based on Structured Heterogeneous Information Network,” in *KDD*, ACM, pp. 1507–1515.
- [27] FAN, Y., S. HOU, Y. ZHANG, Y. YE, and M. ABDULHAYOGLU (2018) “Gotcha-Sly Malware! Scorpion: A Metagraph2vec Based Malware Detection System,” in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD)*, ACM, pp. 253–262.
- [28] PEROZZI, B., R. AL-RFOU, and S. SKIENA (2014) “DeepWalk: Online Learning of Social Representations,” in *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining (KDD)*, pp. 701–710.
- [29] GROVER, A. and J. LESKOVEC (2016) “node2vec: Scalable feature learning for networks,” in *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, ACM, pp. 855–864.
- [30] TANG, J., M. QU, and Q. MEI (2015) “PTE: Predictive Text Embedding through Large-scale Heterogeneous Text Networks,” in *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, pp. 1165–1174.
- [31] TANG, J., M. QU, M. WANG, M. ZHANG, J. YAN, and Q. MEI (2015) “Line: Large-scale information network embedding,” in *WWW '15 Proceedings of the 24th International Conference on World Wide Web*, pp. 1067–1077.
- [32] SHANG, J., M. QU, J. LIU, L. M. KAPLAN, J. HAN, and J. PENG (2016) “Meta-Path Guided Embedding for Similarity Search in Large-Scale Heterogeneous Information Networks,” in *arXiv:1610.09769*.

- [33] DONG, Y., N. V. CHAWLA, and A. SWAMI (2017) “metapath2vec: Scalable representation learning for heterogeneous networks,” in *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD’2017)*, pp. 135–144.
- [34] FU, T.-Y., W.-C. LEE, and Z. LEI (2017) “HIN2Vec: Explore Meta-paths in Heterogeneous Information Networks for Representation Learning,” in *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management (CIKM)*, pp. 1797–1806.
- [35] HOU, S., A. SAAS, Y. YE, and L. CHEN (2016) “DroidDelver: An Android Malware Detection System Using Deep Belief Network Based on API Call Blocks,” in *International Conference on Web-Age Information Management (WAIM)*, pp. 54–66.
- [36] HOU, S., A. SAAS, L. CHEN, and Y. YE (2016) “Deep4MalDroid: A Deep Learning Framework for Android Malware Detection Based on Linux Kernel System Call Graphs,” in *WIW ’16*.
- [37] IDC (2018) “International Data Corporation (IDC),” in <http://www.idc.com>.
- [38] POEPLAU, S., Y. FRATANONIO, A. BIANCHI, C. KRUEGEL, and G. VIGNA (2014) “Execute This! Analyzing Unsafe and Malicious Dynamic Code Loading in Android Applications,” in *NDSS*, pp. 23–26.
- [39] CAPILUPPI, A., A. SEREBRENIK, and L. SINGER (2013) “Assessing technical candidates on the social web,” in *IEEE Software*, pp. 45–51.
- [40] SUN, Y. and J. HAN (2012) “Mining heterogeneous information networks: principles and methodologies,” *Synthesis Lectures on DMKD*, **3**(2), pp. 1–159.
- [41] HOFF, P. D., A. E. RAFTERY, and M. S. HANDCOCK (2002) “Latent space approaches to social network analysis,” *Journal of the American Statistical Association*, **97**(460), pp. 1090–1098.
- [42] YAN, S., D. XU, B. ZHANG, H.-J. ZHANG, Q. YANG, and S. LIN (2007) “Graph embedding and extensions: A general framework for dimensionality reduction,” *IEEE transactions on pattern analysis and machine intelligence (TPAMI)*, **29**(1), pp. 40–51.
- [43] ZHAO, H., Q. YAO, J. LI, Y. SONG, and D. L. LEE (2017) “Meta-graph based recommendation fusion over heterogeneous information networks,” in *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD’2017)*, pp. 635–644.

- [44] MIKOLOV, T., K. CHEN, G. CORRADO, and J. DEAN (2013) “Efficient estimation of word representations in vector space,” in *arXiv preprint arXiv:1301.3781*.
- [45] MIKOLOV, T., I. SUTSKEVER, K. CHEN, G. S. CORRADO, and J. DEAN (2013) “Distributed representations of words and phrases and their compositionality,” in *NIPS*, pp. 3111–3119.
- [46] BOTTOU, L. (1991) “Stochastic gradient learning in neural networks,” *Proceedings of Neuro-Nimes*, **91**(EC2).
- [47] ZUCCON, G., L. A. AZZOPARDI, and C. VAN RIJSBERGEN (2009) “Semantic spaces: Measuring the distance between different subspaces,” in *International Symposium on Quantum Interaction*, Springer, pp. 225–236.
- [48] IPSEN, I. C. and C. D. MEYER (1995) “The angle between complementary subspaces,” *American Mathematical Monthly*, pp. 904–911.
- [49] YE, Y., S. HOU, L. CHEN, X. LI, L. ZHAO, S. XU, J. WANG, and Q. XIONG (2018) “ICSD: An Automatic System for Insecure Code Snippet Detection in Stack Overflow over Heterogeneous Information Network,” in *ACSAC*, pp. 542–552.
- [50] BAHDANAU, D., K. CHO, and Y. BENGIO (2015) “Neural Machine Translation by Jointly Learning to Align and Translate,” in *ICLR*.
- [51] LIU, J., Z. HE, L. WEI, and Y. HUANG (2018) “Content to node: Self-translation network embedding,” in *KDD*, pp. 1794–1802.
- [52] PETERS, M. E., W. AMMAR, C. BHAGAVATULA, and R. POWER (2017) “Semi-supervised sequence tagging with bidirectional language models,” *arXiv preprint arXiv:1705.00108*.
- [53] CHO, K., B. VAN MERRIËNBOER, D. BAHDANAU, and Y. BENGIO (2014) “On the properties of neural machine translation: Encoder-decoder approaches,” *arXiv preprint arXiv:1409.1259*.
- [54] VASWANI, A., N. SHAZEER, N. PARMAR, J. USZKOREIT, L. JONES, A. N. GOMEZ, Ł. KAISER, and I. POLOSUKHIN (2017) “Attention is all you need,” in *Advances in Neural Information Processing Systems*, pp. 5998–6008.
- [55] LUONG, M.-T., H. PHAM, and C. D. MANNING (2015) “Effective approaches to attention-based neural machine translation,” *arXiv preprint arXiv:1508.04025*.
- [56] KIM, Y., C. DENTON, L. HOANG, and A. M. RUSH (2017) “Structured attention networks,” *arXiv preprint arXiv:1702.00887*.

- [57] CHO, K., B. VAN MERRIENBOER, C. GULCEHRE, D. BAHDANAU, F. BOUGARES, H. SCHWENK, and Y. BENGIO (2014) “Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation,” in *arXiv:1406.1078*.
- [58] GRAVES, A. (2013) “Generating sequences with recurrent neural networks,” in *Arxiv preprint arXiv:1308.0850*.
- [59] WILLIAMS, R. and D. ZIPSER (1995) “Gradient-based learning algorithms for recurrent networks and their computational complexity,” in *Back-propagation: Theory, Architectures and Applications*, pp. 433–486.
- [60] PEROZZI, B., R. AL-RFOU, and S. SKIENA (2014) “DeepWalk: Online Learning of Social Representations,” in *KDD '14*, pp. 701–710.
- [61] DONG, Y., N. V. CHAWLA, and A. SWAMI (2017) “metapath2vec: Scalable representation learning for heterogeneous networks,” in *KDD*.
- [62] YANG, C., Z. LIU, D. ZHAO, M. SUN, and E. Y. CHANG (2015) “Network Representation Learning with Rich Text Information,” in *IJCAI'15*, pp. 2111–2117.