

Thomas Brüggemann

**Master Thesis
im Fach Information Systems**

Automated Information Privacy Risk Assessment of Android Health Applications

Themensteller: Prof. Dr. Ali Sunyaev

Vorgelegt in der Masterprüfung
im Studiengang Information Systems
der Wirtschafts- und Sozialwissenschaftlichen Fakultät
der Universität zu Köln

Köln, September 2016

Contents

Index of Abbreviations	III
Table of tables.....	IV
1. Introduction	1
1.1 Problem Statement.....	1
1.2 Objectives	2
1.3 Structure	3
2. Combining Source Code Analysis with Information Privacy Risk Assessment	4
2.1 Information Privacy Risk Assessment	5
2.2 Static Code Analysis	5
2.3 Relevant Information Privacy Risk Factors	5
3. Implementation and Evaluation of an Automated In-formation Privacy Risk Assessment Tool	6
3.1 Implementation of an Automated Information Privacy Risk Assessment Tool	6
3.1.1 Download Phase	6
3.1.2 Decompilation Phase.....	6
3.1.3 Static code analysis Phase.....	6
3.2 Evaluation of an Automated Information Privacy Risk Assessment Tool.....	6
4. Feasibility of Automated Information Privacy Risk Assessment	7
4.1 The Automated Information Privacy Risk Assessment of Free Android mHealth Apps	7
4.1.1 Download Phase	7
4.1.2 Decompilation Phase.....	7
4.1.3 Static code analysis Phase.....	7
4.2 Evaluation of the Auto-mated Information Privacy Risk Assessment Tool	7
5. Discussion.....	8
5.1 Principle Findings.....	8
5.2 Contributions.....	8
5.3 Limitations	8
5.4 Future Research	8
5.5 Conclusion.....	8
Declaration of Good Scientific Conduct	9
References.....	9
Curriculum Vitae	10

Index of Abbreviations

mHealth

Mobile Health

List of Tables

1. Introduction

1.1 Problem Statement

The market for mobile phone and tablet applications (apps) has grown extensively since recent years.¹ It is increasingly easier for companies or even single developers to create unique apps that reach millions of users around the planet via digital app stores. This market growth affected mobile health (mHealth) apps as well. More and more mHealth apps are available that support the users in resolving their health-related issues and that try to remedy health-related information deficiencies.

But receiving personal health-related information yields information privacy risks to users. Users are asked to expose personal health-related information, e.g. information on disease symptoms or medical appointments in order to receive a tailored app that fits their needs.² It remains however unclear how and where the vulnerable user information is sent, processed and stored.³

The information about these privacy related practices of app providers and their offered apps should be stated in the privacy policy document provided by the app provider.⁴ Processing these privacy policies requires a higher level of education and time to read through large bodies of text, in order to find the relevant information. Additionally, the important information is hidden in legal language or is insufficiently addressed, if at all.⁵ Aside from data usage beyond the control of the users, it is also challenging to assess what kind of private information an app asks for, prior to the app usage. Users have to download the apps of interest and try them out, before it becomes clear what health-related information is processed by the app and in which way. This leads to low comparability between apps. When users are looking for specific functionality in an mHealth app, it is challenging to find the app that offers the desired functionality at an acceptable information privacy risk. Even if users would pursue the task of finding and comparing mHealth

¹ See for this and the following sentence **Enck2011** p. 1.

² See **Chen2012** p. 2.

³ See **He2014** p. 652.

⁴ This paragraph follows **Dehling2014** p. 11.

⁵ See **Pollach2007** p. 104.

apps of similar functionality, the high volume of apps available in the app stores⁶ makes it laborious to review all of them by hand. One way to assess information privacy risks of the large amount of mHealth apps is to automate the review process of each individual app. The assessment automation can be done by downloading and analyzing the source code of each app and by tracing data leaks. Static code analysis is used in the field of informatics to analyze application source code and detect faults or vulnerabilities.⁷ It is yet unclear how and to what degree the concepts of static code analysis and information privacy risk assessment can be combined in order to automate app assessment. A static code analysis could, in theory, be used to automatically assess some of the information privacy risks that mHealth apps pose. Previous research has not shown how and to what degree the combination of static code analysis and information privacy risks assessment is feasible in the field of mHealth app information privacy risk assessment and therefore the aim of this study is to explore the possibilities of static code analysis for information privacy risk assessment. This leads to the research question: How and to what degree can the information privacy risks of mHealth apps be automatically assessed? The 'degree' refers to the amount and the level of detail that information privacy risk factors can be automatically assessed.

The automated process furthermore can help to drastically reduce the effort of reviewing each individual app and can enhance the information experience users receive while looking for mHealth apps. Additionally, it exposes new possibilities for research in the information privacy risks area. The research could be conducted on providing solutions and best practices for further enhancing the information privacy risks communication of apps.

1.2 Objectives

The main objective of this study is to ascertain how and to what degree the assessment of information privacy risk factors for mHealth apps can be automated. In order to reach this objective, the following sub-objectives have to be met.

The first sub-objective is to extract information privacy risk factors from the infor-

⁶ See **Enck2011** p. 1.

⁷ See **Baca2008** p. 79.

mation privacy practices that **Dehling2016** identified and that are relevant for automated information privacy risk assessment. As a second sub-objective we will develop strategies to identify the information privacy risk factors within the source code of mHealth apps via static code analysis. This is necessary since it is yet unclear how and to what degree the static code analysis can help to identify information privacy risk factors of mHealth apps. Finally we will evaluate how well the automated information privacy risk assessment tool can identify information privacy risk factors in comparison to two human reviewers. In order to fully ascertain the degree static code analysis can identify information privacy risk factors, a manual review of the results of the static code analysis is necessary.

1.3 Structure

2. Combining Source Code Analysis with Information Privacy Risk Assessment

mHealth apps have been examined in various research studies that aim at providing insights for developers as well as users into how private information is processed. Privacy issues are the most impactful user complaint while using mobile apps.⁸ This encourages research to address information privacy risks.

Research focus has been put on the technical side of information privacy breach. It has been analyzed, to what degree the data storage in internal Android log files or on the memory card within a phone or tablet poses a threat to users information privacy.⁹ Technical evaluation of mobile apps even goes further into the topics of decompilation to analyze device identification or geolocation data leaks.¹⁰ Decompilation reveals to be a feasible assessment technique for information privacy risks and data leaks.

In informatics and software development contexts, static code analysis has been used to analyze source code and provide feedback on coding styles to the users while programming or "to find defects in programs"¹¹. Static code analysis provides a fast way to analyze source code¹², which makes it suitable to automate the assessment of large datasets. A further benefit of using static code analysis to retrieve information from software is that the software does not need to be executed during the analyzation process. This additionally supports the development of fast performing assessment tools that are suitable for application on large datasets of source code since there is no need to wait for the application runtime to execute the software.

Our study will use the benefits of static code analysis and apply them to the assessment of mHealth information privacy risks. It is unclear if static code analysis is a viable tool to analyze and identify information privacy risk factors. We will use the comprehensive privacy-risk-relevant information privacy practices that **Dehling2016** identified¹³ and try to implement static code analysis strategies to identify those risks automatically. This

⁸ See **Khalid2015** p. 5.

⁹ For the previous two sentences, see **He2014** p. 645-646.

¹⁰ See **Mcclurg2012** p. 1, 5., **Enck2011** p. 1. and **Mitchell2013** p.6-7.

¹¹ **Bardas2010** p. 1.

¹² See **Bardas2010** p. 5.

¹³ See **Dehling2016** p. 8-17.

will be a vital addition to current research, since there is yet no holistic approach to apply static code analysis to information privacy risks detection that takes an ample amount of information privacy risk factors into account.

2.1 Information Privacy Risk Assessment

2.2 Static Code Analysis

2.3 Relevant Information Privacy Risk Factors

3. Implementation and Evaluation of an Automated Information Privacy Risk Assessment Tool

3.1 Implementation of an Automated Information Privacy Risk Assessment Tool

3.1.1 Download Phase

3.1.2 Decompilation Phase

3.1.3 Static code analysis Phase

3.2 Evaluation of an Automated Information Privacy Risk Assessment Tool

4. Feasibility of Automated Information Privacy Risk Assessment

4.1 The Automated Information Privacy Risk Assessment of Free Android mHealth Apps

4.1.1 Download Phase

4.1.2 Decompilation Phase

4.1.3 Static code analysis Phase

4.2 Evaluation of the Auto- mated Information Privacy Risk Assessment Tool

5. Discussion

5.1 Principle Findings

5.2 Contributions

5.3 Limitations

5.4 Future Research

5.5 Conclusion

Declaration of Good Scientific Conduct

Hiermit versichere ich an Eides Statt, dass ich die vorliegende Arbeit selbstständig und ohne die Benutzung anderer als der angegebenen Hilfsmittel angefertigt habe. Alle Stellen, die wörtlich oder sinngemäß aus veröffentlichten und nicht veröffentlichten Schriften entnommen wurden, sind als solche kenntlich gemacht. Die Arbeit ist in gleicher oder ähnlicher Form oder auszugsweise im Rahmen einer anderen Prüfung noch nicht vorgelegt worden.

Köln, den 01. September 2016

I hereby attest that I completed this work on my own and that I did not employ any tools other than those specified. All texts literally or semantically copied from other works are attributed with proper citations. This work has not been submitted in identical or similar form for any other exam, assessment, or assignment.

Cologne, September 1st, 2016

Curriculum Vitae



Persönliche Angaben

Name: Thomas Brüggemann
 Anschrift: Hoferkamp 9, 41751 Viersen
 Geburtsdatum und -ort: 31.08.1989 in Viersen
 Familienstand: verheiratet

Schulische Ausbildung

1997 - 2001 Katholische Grundschule Boisheim
 2001 - 2009 Bischöfliches Albertus-Magnus-Gymnasium in Viersen,
 Abschluss: Abitur

Grundwehrdienst

07/2009 - 04/2010 Wehrdienstleistender, Luftwaffe -
 Jagdbombergeschwader 31 "Boelke", KvD für das
 Wachpersonal, Fliegerhorst Nörvenich

Studium

10/2010 - 03/2014 Universität zu Köln, Wirtschaftsinformatik, Bachelor of
 Science
 10/2014 - 09/2016 Universität zu Köln, Information Systems, Master of
 Science

Beruflicher Werdegang

05/2010 - 09/2012 Thomas Trefz Consulting, Köln, Softwareentwicklung
 im Bereich Microsoft .NET
 10/2012 - 10/2014 Beister Software GmbH, Aschaffenburg, Softwareen-
 twicklung im Bereich Microsoft .NET
 10/2014 - heute Selbstständiger Softwareentwickler und IT-Berater