

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 Information 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.....	7
3.1.3 Static code analysis Phase.....	9
3.2 Evaluation of an Automated Information Privacy Risk Assessment Tool.....	11
4. Feasibility of Automated Information Privacy Risk Assessment	12
4.1 The Automated Information Privacy Risk Assessment of Free Android mHealth Apps	12
4.1.1 Download Phase	12
4.1.2 Decompilation Phase.....	12
4.1.3 Static code analysis Phase.....	12
4.2 Evaluation of the Automated Information Privacy Risk Assessment Tool.....	12
5. Discussion.....	13
5.1 Principle Findings.....	13
5.2 Contributions.....	13
5.3 Limitations	13
5.4 Future Research	13
5.5 Conclusion.....	13
References.....	16
Declaration of Good Scientific Conduct	17
Curriculum Vitae	18

Index of Abbreviations

API	Application Programming Interface
APK	Android Application Package
JAR	Java Archive
mHealth	Mobile Health
URL	Uniform Resource Locator

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 Enck et al. (2011), p. 1.

² See Chen et al. (2012), p. 2.

³ See He et al. (2014), p. 652.

⁴ This paragraph follows Dehling, Gao, Sunyaev (2014), p. 11.

⁵ See Pollach (2007), 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 Enck et al. (2011), p. 1.

⁷ See Baca, Carlsson, Lundberg (2008), p. 79.

mation privacy practices that Dehling, Sunyaev (2016) 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 Dehling, Sunyaev (2016) identified¹³ and try to implement static code analysis strategies to identify those risks au-

⁸ See Khalid et al. (2015), p. 5.

⁹ For the previous two sentences, see He et al. (2014), p. 645-646.

¹⁰ See McClurg (2012), p. 1, 5., Enck et al. (2011), p. 1. and Mitchell et al. (2013), p.6-7.

¹¹ Bardas, Others (2010), p. 1.

¹² See Bardas, Others (2010), p. 5.

¹³ See Dehling, Sunyaev (2016), p. 8-17.

tomatically. This 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

For this thesis, we used the set of information privacy practices extracted from literature and ... by Dehling, Sunyaev (2016) as a source to derive information privacy risk factors from. Since not all information privacy practices that an app provider can include in his privacy policy may express an information privacy risk, we extract the information privacy practices that are relevant in terms of posing and expressing a potential information privacy risk. The full list of information privacy practices including a comment, whether

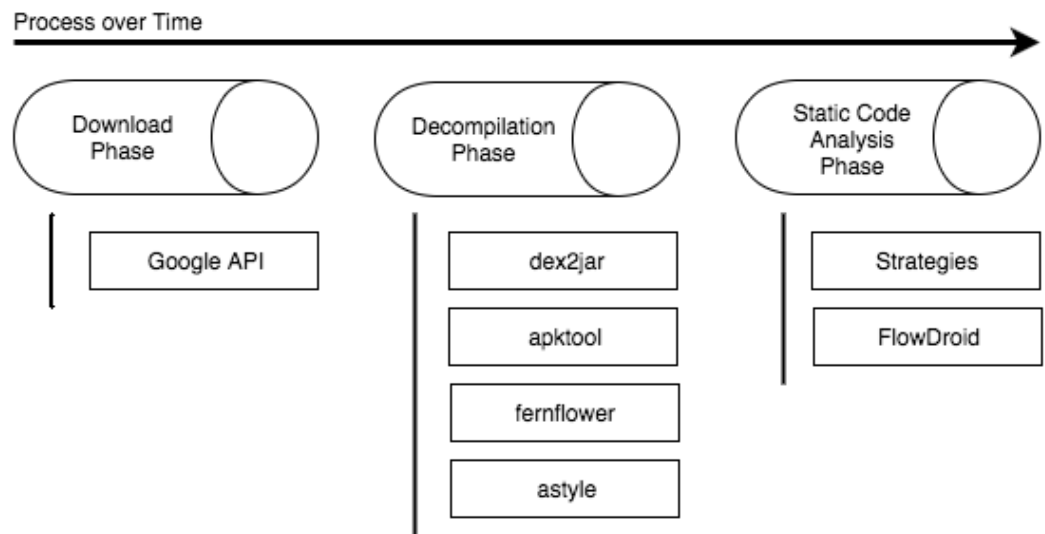
3. Implementation and Evaluation of an Automated Information Privacy Risk Assessment Tool

3.1 Implementation of an Automated Information Privacy Risk Assessment Tool

The implementation of an automated information privacy risk assessment tool is structured in three phases. In the first phase, Android APK files need to be downloaded to acquire the foundation of a static code analysis: the source code. While APK files are binary representations of source code, it is necessary, in a second phase, to decompile to binary code back into actual source files. The third phase is the analysis phase, where the information privacy risk assessment takes place.

Figure 3-1 shows the implementation phases over time including the tools used within each phase. The tools will be described in greater detail within the following chapters.

Figure 3-1: Diagram of implementation phases over time.



3.1.1 Download Phase

The download phase is the first of the three implementation phases and comprises the acquisition of Android APK files. The APK files hold the necessary Java source code that we will perform the static code analysis on. Since this thesis emphasizes on Android mHealth apps, we used the repository database of Xu, Liu (2015) as our main app datasource.¹⁴

¹⁴ This paragraph follows Xu, Liu (2015).

Xu, Liu (2015) list mHealth apps from the Apple AppStore and Android PlayStore and update their repository quarterly by scraping the app stores. The list contains information for example on the app's id, category in the app stores, description, email address of the developer, price and the user rating of the app.

We used the repository database to loop over the available mHealth app listings and filtered out the apps that were available for free, indicated by a price of \$0.00. As soon as the package name of an app is gathered, the download of the APK file can begin. While there is no official source to download APK files for Android apps, a multitude of websites exist that host copies of APK files to download for free. Unfortunately, all of these websites implement mechanisms that make it impossible to browse and download the APK files programatically within a download script. Instead, we used a open source Python implementation of an undocumented Google PlayStore API.¹⁵ The undocumented part of the Google API allows users to download APK files Even though the project has not been maintained for four years, the software is still in working order. The Python script authenticates to the Google API via the hardware ID of an Android smartphone or tablet and pretends to request data from this smartphone or tablet, even though the requests are sent from a desktop computer. We used a real Android tablet to detect its hardware ID and authenticate the Google API requests with this hardware ID. The main issue that has to be taken care of during the download phase is not to run into Google API limitations. Google allows an API user to only request a certain amount requests per time unit. After this limit is exceeded, the requests will just return a HTTP error code and no APK file will be downloaded. In order to mitigate this circumstance, we ran our download script multiple times, always until the Google API returns error codes. We then waited a couple of hours and tried the download script again, which would pick up the download process where it had stopped on the last run.

3.1.2 Decompilation Phase

In order to decompile the amount of APK files available, it is necessary to automate the process. The automation script¹⁶ uses a chain of tools to gather access to the source code

¹⁵ <https://github.com/egirault/googleplay-api>, visited 05/12/2016

¹⁶ <https://github.com/thomasbrueggemann/AIPRAT/blob/master/decompile/decompile.sh>

files from an APK file. The tools used to decompile the APK files follow closely the tools described and used by Enck et al. (2011).¹⁷

In a first step, we use the tool *dex2jar* to extract the JAR file from the APK file. The JAR file contains the java bytecode representations of the app which is just one part of the contents of an APK file. The next step is to extract resource files, such as the *Android Manifest* file from the APK file. The *Android Manifest* contains meta information about the app in a structured XML format.¹⁸ The meta information include the package name of the app, the permissions the app requests, e.g. camera usage, internet access or geolocation usage. The *Android Manifest* file is therefore an important indicator for high level activities within the given app. In order to extract the *Android Manifest* file from the APK file along with other resources such as images, icons, xml files or other files used within the app, we use the *apktool*¹⁹. *apktool* is a frequently updated Android reverse engineering tool that is used to extract resources from APK files. At the core of the decompilation process is the usage of *fernflower*²⁰. *fernflower* is the recommended java decompiler by Enck et al. (2011). They used the tool to decompile a test sample of apps and gained a significantly higher code recovery rate than by using other decompiling tools.²¹ An obstacle in decompiling java source code is obfuscation. Java developers can make use of a security feature called obfuscation that aims at hiding away the logic of java classes by renaming classes, variables and method names and disassembling the code into pieces that are difficult to read for an human interpreter. The idea is to make it more difficult to retrieve and make sense of the original source code by decompiling the byte code. *fernflower* uses a renaming approach by assigning every obfuscated class with a new naming pattern. Member variables and methods will be automatically renamed and therefore provide an easier and more unique way of reading the source code. Optionally the decompilation process can use an automatic code formatting tool called *astyle*²² to format the source code. This helps humans to read the source code files more easily,

¹⁷ See Enck et al. (2011), p. 5.

¹⁸ This and the next sentence follow Xu (2013), p. 7.

¹⁹ <http://ibotpeaches.github.io/Apktool/>

²⁰ <https://github.com/fesh0r/fernflower>

²¹ See Enck et al. (2011), p. 6.

²² <http://astyle.sourceforge.net/>

since the formatting and indentation of all source code files is identical and therefore very structured. Formatting the source code will help in the evaluation phase of this thesis to support the manual inspection the source code for information privacy risks by human researchers.

The expected result of the decompilation phase is a directory named after the package name of a given app that contains the resource files, including the *Android Manifest* and the decompiled source code of the app.

3.1.3 Static code analysis Phase

The static code analysis phase is the main analysis phase of the thesis and uses the output of the previous decompilation phase to perform the static code analysis. The static code analysis tool is implemented as a Java software project, since the used analysis libraries are implemented in Java and Android source code is written in Java too. The output of the static code analysis Java project is an executable Java archive file called *AIPRAT.jar* that can be executed in the command line terminal. In order for *AIPRAT.jar* to perform the static code analysis on APK files, two preparation steps are required.

The first preparation step is to run an Android data flow analysis tool over the APK files that extract potential data flows. The data flow analysis is achieved with an open source tool called *FlowDroid*, introduced by Arzt et al. (2014).²³ *FlowDroid* extends the Java optimization framework *Soot*, which was already used by Enck et al. (2011) for post-decompilation optimization tasks.²⁴ The data flow is analysed by scanning an intermediate byte code format provided by *Soot* for so called 'sources' and 'sinks'.²⁵ A source is the origin of a data flow, e.g. the user input of data via a textfield and a sink is the destination that data flows. An example for a sink is a HTTP internet connection or a local log file. *FlowDroid* is also able to emulate Android lifecycle entry points. While a regular Java program has a single entry point to start the application from, the *main()* function, Android apps provide multiple entry points. The entry points of an Android app are determined by the states an app can be in. It can e.g. return from being in the background, do a fresh start

²³ See Arzt et al. (2014), p. 259-269.

²⁴ See Enck et al. (2011), p. 5.

²⁵ See Arzt et al. (2014), p. 264.

and return from being offline. All these entry points are being emulated by *FlowDroid* into a single *main()* function call. The output of the data flow analysis is one XML file per analysed APK file that contains a list of sinks and the coherent sources of data flows to that sink. The XML file will be parsed by the main static code analysis tool later on and the sink and source methods will be interpreted in the context of information privacy risks.

The machine learning text classifiers will be trained within the second preparation step. During the static code analysis phase of this study, we will be making great use of the naive Bayes classifier. A machine learning text-classifier classifies text segments into distinct categories. The categories are predefined in the training phase of the classifier, since every trained text segment is assigned with a training category. These training categories are the categories the classifier can assign to new, previously unseen, text segments. The incisive feature of a Bayes classifier is the fact, that it chooses to classify a category to a new segment of text by picking the most probable category.²⁶ A naive Bayes classifier furthermore assumes that all categories are distinct and independent of each other. Even though this might not always be the case in a real life usage scenario, the naive Bayes classifier still performs well enough for a wide range of use cases.

In the case of the static code analysis in this study, we will be using the naive Bayes classifier to classify URLs into categories. The categories that URLs can belong to, in the context of this study, are: advertisement, delivery services, government, instant-messaging, (data-) aggregation services, search engines and social networks. While the set of categories might not complete in terms of all possible and available categories, it is sufficient for the classification of URLs within this static code analysis to classify into the mentioned category-set. In order for the naive Bayes classifier to classify text into categories we trained a naive Bayes classifier implementation with meta-information about URLs from the previously mentioned categories. First, it was necessary to collect URLs for the categories to train the naive Bayes classifier and we used a collection of URLs from *URLBlacklist.com*²⁷. *URLBlacklist.com* provides URL lists for the categories advertisement, government, instant-messaging, search engines and social networks. *pro-*

²⁶ For this and the following two sentences see Rish (2001), p. 41.

²⁷ <http://www.urlblacklist.com/?sec=download>, visited 05/30/2016

grammableweb.com catalogues API descriptions including the service providers' URL. We developed a program to automatically download and store the API directory for the two remaining categories, from *programmableweb.com*.

Next, to acquire meta-information for all the URLs, we implement a downloader for the HTML source-code of all URLs and store the 'description' HTML-meta tag content in a file. The 'description' meta-tag contains a small amount of text, provided by the website owner, that describes the content or function of the website topic. We use this 'description' meta-information to train the classifier with the associated categories.

As soon as the preparation steps are finished, the main static code analysis tool is ready to run.

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 Automated Information Privacy Risk Assessment Tool

5. Discussion

5.1 Principle Findings

5.2 Contributions

5.3 Limitations

5.4 Future Research

5.5 Conclusion

References

Arzt et al. (2014)

Steven Arzt, Siegfried Rasthofer, Christian Fritz, Eric Bodden, Alexandre Bartel, Jacques Klein, Yves Le Traon, Damien Octeau, Patrick McDaniel. “Flowdroid: Precise context, flow, field, object-sensitive and lifecycle-aware taint analysis for android apps”. In: *ACM SIGPLAN Notices*. Vol. 49. 6. ACM. 2014, pp. 259–269.

Baca, Carlsson, Lundberg (2008)

Dejan Baca, Bengt Carlsson, Lars Lundberg: Evaluating the cost reduction of static code analysis for software security. In: Proceedings of the third ACM SIGPLAN workshop on Programming languages and analysis for security - PLAS '08. 2008, p. 79

Bardas, Others (2010)

Alexandru G Bardas, Others: Static code analysis. In: Journal of Information Systems & Operations Management. No. 2, Vol. 4, 2010, pp. 99–107

Chen et al. (2012)

Connie Chen, David Haddad, Joshua Selsky, Julia E Hoffman, Richard L Kravitz, Deborah E Estrin, Ida Sim: Making sense of mobile health data: an open architecture to improve individual- and population-level health. In: Journal of medical Internet research. No. 4, Vol. 14, 2012, e112

Dehling, Gao, Sunyaev (2014)

Tobias Dehling, Fangjian Gao, Ali Sunyaev: Assessment Instrument for Privacy Policy Content: Design and Evaluation of PPC. In: WISP 2014 Proceedings. 2014,

Dehling, Sunyaev (2016)

Tobias Dehling, Ali Sunyaev: “Designing for Privacy: A Design Theory for Transparency of Information Privacy Practices”. 2016

Enck et al. (2011)

William Enck, Damien Oteau, Patrick McDaniel, Swarat Chaudhuri: A Study of Android Application Security. In: Proceedings of the 20th USENIX Conference on Security. No. August, Vol. SEC'11, 2011, pp. 1–21

He et al. (2014)

Dongjing He, Muhammad Naveed, Carl A Gunter, Klara Nahrstedt: Security Concerns in Android mHealth Apps. In: AMIA ... Annual Symposium proceedings / AMIA Symposium. AMIA Symposium. Vol. 2014, 2014, pp. 645–54

Khalid et al. (2015)

Hammad Khalid, Emad Shihab, Meiyappan Nagappan, Ahmed E. Hassan: What Do Mobile App Users Complain About? In: IEEE Software. No. 3, Vol. 32, 2015, pp. 70–77

Mcclurg (2012)

Jedidiah Mcclurg: Android Privacy Leak Detection via Dynamic Taint Analysis. In: . 2012,

Mitchell et al. (2013)

Stacy Mitchell, Scott Ridley, Christy Tharenos, Upkar Varshney, Ron Vetter, Ulku Yaylacicegi: Investigating privacy and security challenges of mhealth applications. In: 19th Americas Conference on Information Systems, AMCIS 2013 - Hyperconnected World: Anything, Anywhere, Anytime. Vol. 3, 2013, pp. 2166–2174

Pollach (2007)

Irene Pollach: What's Wrong With Online Privacy Policies? In: Communications of the ACM. No. 9, Vol. 50, 2007, pp. 103–108

Rish (2001)

Irina Rish. “An empirical study of the naive Bayes classifier”. In: *IJCAI 2001 workshop on empirical methods in artificial intelligence*. Vol. 3. 22. IBM New York. 2001, pp. 41–46.

Xu (2013)

Liang Xu: “Techniques and Tools for Analyzing and Understanding Android Applications”. PhD thesis. 2013

Xu, Liu (2015)

Wenlong Xu, Yin Liu: mHealthApps: A Repository and Database of Mobile Health Apps. In: *JMIR mHealth and uHealth*. No. 1, Vol. 3, 2015, e28

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