



Universität Hamburg

DER FORSCHUNG | DER LEHRE | DER BILDUNG

Bachelor thesis

Privacy implications of exposing Git meta data

presented by

Arne Beer

born on the 21th of December 1992 in Hadamar

Matriculation number: 6489196

Department of Computer science

submitted on May 6, 2018

Supervisor: Dipl.-Inform. Christian Burkert

Primary Referee: Prof. Dr.-Ing. Hannes Federrath

Secondary Referee: Prof. Dr. Dominik Herrmann

Abstract

‘As long as we use modern technologies we always expose data about ourselves’. This is a statement I truly believe in.

Recent events, such as the Facebook scandal in which the data of several million people has been exposed to a consulting company ¹ show how large amounts of data can be abused to extract valuable knowledge and used for malicious purposes.

This thesis aims to give an example of how much information can be exposed by simply using the popular version control system Git. Simple meta data such as UNIX timestamps and email addresses might be enough to extract sensitive information about Git users or organizations using Git. This paper covers the whole process of gathering the data from a vast amount of git repositories through to preprocessing and interpreting the results of the analyses. With this thesis I hope to raise the awareness how dangerous it can be to expose even simple meta data and that it might be used maliciously.

¹‘Facebook scandal hits 87 million users’ BBC.com, <http://www.bbc.com/news/technology-43649018> (accessed, 24.04.2018)

Even if you're not doing anything wrong, you are being watched and recorded.

Edward Snowden

Contents

1	Introduction	1
1.1	Contents	1
1.2	Motivation	2
1.3	Leading Questions and Goals	2
1.4	Methods and approach	2
1.5	Attack models	2
1.5.1	The Employer	3
1.5.2	The Individual	4
1.5.3	The Industrial Spy	5
2	Data	7
2.1	Git	7
2.1.1	Introduction to Git	7
2.1.2	Git User Roles	8
2.1.3	Internal Representation	9
2.1.4	Features	10
2.1.5	More features	11
2.2	Data source	11
2.2.1	Requirements	12
2.2.2	Github	13
2.2.3	Github's features	13
2.3	Data Aggregation	14
2.3.1	Existing solutions	15
2.3.2	Database	16
2.3.3	Gitalizer	18
2.3.4	Database optimization	19
2.3.5	Incremental aggregation	20
2.3.6	Problems	21
3	Implementation	23
3.1	Holiday and Sick Leave Detection	23
3.1.1	Implementation	23
3.2	Working hours	26
3.2.1	Implementation	26
3.2.2	Punchcard Clustering	27
3.3	Geolocation	30
3.3.1	Implementation	30
4	Evaluation and Interpretation	36
4.1	Holiday and Sick Leave Detection	36
4.2	Working hours	37
4.2.1	Sleep rythm	37

4.2.2	Employee or Open-Source Contributor	38
4.2.3	Bot detection	40
4.2.4	Fingerprinting	41
4.3	Geolocation	41
5	Conclusion and Outlook	45

Acronyms

API application programming interface

CPU central processing unit

DBSCAN Density-based spatial clustering of applications with noise

DST daylight saving time

EU European Union

FS file system

GDPR General Data Protection Regulation

GB Gigabyte

HTTP hypertext transfer protocol

IANA Internet Assigned Numbers Authority

JSON JavaScript Object Notification

MIT Massachusetts Institute of Technology

ORM Object-Relational Mapping

OS operating system

SHA-1 secure hash algorithm 1

SSH secure shell

SQL Structured Query Language

URL uniform resource locator

UTC Coordinated Universal Time

VCS version control system

CHAPTER 1

Introduction

Git is a **VCS!** (**VCS!**) used by most programmers on a daily basis these days. It is supposed help developers to manage their According to the Eclipse Community Survey about 42.9% of professional software developers used git in 2014 with an upward tendency [12]. It is deployed in many if not most commercial and private projects and generally valued by its users. It allows to collaborate with thousands of contributors on the same project whilst maintaining order.

Several million users send new commits to their Git repositories every day. On Github alone, the currently biggest open source platform, there exist about 25 million active repositories, a total of 67 million repositories and about 24 million users [5].

Some well known projects and organizations use Git, for example Linux, Microsoft, Ansible and Facebook [5]. But every of those repository contains the complete contribution history of every contributing user. Each of those contributions contains all changes, a timestamp, a message from the author and their email address.

This raises the question how much information is hidden in the metadata of a Git repository and which attack vectors could be introduced by analysing this information. Could this information be used to harm or manipulate a contributor or maybe even a company.

The newly gained knowledge could be utilized by employers to spy on their employees. It could be used by an unknown attacker who aims to obtain sensitive information about a company and its employees through their open-source projects. It is even possible that a privat individual wants to monitor another person, which regularly contributes to open-source repositories.

As there have not been any papers published about this specific topic yet or at least no public paper and as Git plays such a crucial role in todays information technology, I want to investigate and evaluate this potential threat. Furthermore I want to create a foundation for future research and provide a first example on how such attacks might look like.

1.1 Contents

This thesis examines the process and capabilities of building a data mining software based on git metadata. After the introduction explaining the motivation, the approach and the goals of this thesis, Chapter 2 will explain the requirements to the data, existing solutions as well as process of building an own aggregator, which collects data from

Github. Chapter 3 shows the actual implementation of three different algorithms and how their gained information could be used in a real world scenario. Chapter 4 evaluates the results of the previously presented algorithms by comparing it with real world ground truth and conducting small surveys. Chapter 5 then concludes and discusses the overall results of this research.

1.2 Motivation

Each year more and more data is collected by employers. A study conducted by Massachusetts Institute of Technology (MIT) students shows, that data-driven companies, that collect data about everything in their company, are about 5% more productive [10], but this does not necessarily justify an invasion into employees' privacy. Privacy violation already goes as far as the tracking of medical records of employees, which is a service provided by the company *Castlight* [13]. Some parties warn against surveillance of employees by the management and demand stricter handling of employee data [8] Generally it can be said, that we need to take better care about information of ourselves and that people need to know in which ways they can leak information. This thesis aims to show this on the example of the version control system (VCS) Git, which is technically a very practical tool, which was developed with no malicious intend, but it might be used for such.

1.3 Leading Questions and Goals

The primary research objective of this thesis is finding out how feasible a data mining approach on git metadata would be. This question is explored by looking at various data mining techniques and applying them to aggregated data from Github.

1.4 Methods and approach

To gather data and perform analysis on the data for the purpose of this thesis, a program called *Gitalizer* is developed. It is capable of continuously collecting data from Github in several ways. Furthermore three different analysis methodologies and several evaluation methods for these are implemented.

1.5 Attack models

To show possible scenarios, this section introduces three attacker models and their respective goals. These attack models serve as a guideline for the data aggregation process,

which will be covered in the next chapter. Only a small subset of these listed attacks will be actually performed by me, but this listing also serves as a exhibition of some possible attacks, for anyone that wants to further investigate this topic.

1.5.1 The Employer

This attack model deals with the scenario of an employer, which wants to monitor their employees. The attacker's motivation is to spot irregularities in working behavior as well as unmotivated or unproductive employees.

Productivity of Employees

An employer wants to ensure that their employees produce sufficient amount of code. For this purpose the changes in lines of code over a specific time span will be evaluated.

Required data:

- Commit timestamps
- Additions of each commit
- Deletions of each commit

Compliance of Working Hours

Check if an employee is productive in the defined working hours. This is especially useful to supervise employees, which work remotely.

Required data:

- Commit timestamps

External Projects during Working Hours

Inspect if an employee is working on an external project during working hours. This only works if the employer has access to the external project, for example open source projects.

Required data:

- Commit timestamps

Code Quality Between Employees

Compare the quality of contributed code between different employees. With this metric the quality of an employee could be measured. To compare the quality we would need an external tool for code analysis.

Required data:

- Complete commit patch
- Commit timestamps

1.5.2 The Individual

This scenario describes a single person, which wants to harm, monitor or gain information about an open-source developer.

A possible goal of the attacker could be to either stalk the victim, harm him in any way or to manipulate him or one of his acquaintances. The motivation of this attacker is mostly personal and on an emotional level. Another non emotional attacker motive could be a robber trying to find the perfect time window to rob a house or the tracking of a high profile target. A third attacker motive could be a headhunter which tries to get information about the skills and reliability of a developer.

Sleeping Rhythm and Daily Routine

Learn about the persons sleep rhythm and obvious patterns in his daily routine. This attack aims to understand and predict the victim's behaviour.

Required data:

- Commit timestamps.

Personal Relationships to Various Programmers

Detect which contributor work together and try to discover the relationships between them.

Required data:

- Git history.
- Commit timestamps.

Sick Leave and Holiday

Detect breaks in his typical work behaviour. Those could represent holiday breaks or sick leave. This attack could give information about whether a developer is at home right now or if they tends to be sick a lot.

Required data:

- Commit timestamps.

Geolocation

Pinpoint the location of an contributor or at least narrow it down to a timezone or country could be interesting. Doing so could provide the attacker with a history of the target's travels.

Required data:

- Commit timestamps.

1.5.3 The Industrial Spy

This attack model covers the scenario of an external person, which wants to gain as much private or malicious information about a company as possible. The attacker's motivation is either to harm the company, gain an advantage as a competitor or in the stock market or to sell secret information to a third party. This attack vector only works if the targeted company is providing their full product or at least parts of their product as open-source software.

Company Employees

The most important target is to detect the company's employees as three other goals for this attacker model depend on this information. Another motivation could be to detect company members for further social engineering attacks or to headhunt the company's employees.

Required data:

- Commit history graph.

Employee History

Detect the timespan for which an employee worked at a given company. This could be interesting, as it shows the average employment duration and the employee amount over the history of the company, which could be an indicator of its current financial growth. Social engineering or headhunting could be a motivation here as well.

Required data:

- Company Employees
- Commit timestamps

Global Workforce Distribution

Detect the timezone of all employees and create a global distribution graphic by timezones. This graphic allows you to guess the location of a company's workforce. It is also possible to create this statistic for all contributors, which could show a trend which countries or at least continents are interested the most for the company's product.

Required data:

- Company Employees
- All Commits
- Commit timestamps

Internal Team Structures

Try to predict different teams, the role of each team and the respective team members.

Required data:

- Company Employees.
- Commit history graph.

CHAPTER 2

Data

This chapter will attend to the collection of required data as stated in Section 1.5. At first the VCS *Git* will be introduced and its functionalities explained. The actual source of the data *Github* will then be evaluated in terms of amount of ground truth and availability. At last the methodology used for aggregation and exploration of Github will be explained.

2.1 Git

This section introduces the VCS *Git*, as it plays a fundamental role in this thesis. In the following the most relevant parts of Git will be explained such as user roles, technologies and internal data representations. I will also talk about the current cases of application and some interesting scenarios which might be interesting for this thesis.

2.1.1 Introduction to Git

At its core, Git is a tool, which is used to manage different versions of files in a specific directory. A directory managed by Git is called a *repository*. Each version of the project is saved as a so called *commit*, which represents a specific state of all files and directories in the project. Users are able to meticulously specify files or changes in files that should be added to a commit. For instance a developer can only commit a subset of the changes, which were applied to a repository. By doing so, one can split a large set of possibly completely unrelated changes into several commits, where each commit in itself forms a set of logically related changes. After creating at least two commits, Git is capable of showing the exact changes between any two commits, which is called a *diff* and it allows to jump between different commits of the project, which is called a *checkout*.

Git is the currently most popular tool to control a project's code with a still upward tendency [12]. It enables to work with multiple developers on a single code base, as it provides several techniques to prevent, detect and resolve conflicts of changes at the same code, namely the *history tree*, the *branch* and the *merge*. The commit history of Git is internally represented as an directed, non-cyclic, connected graph of commits. The commits act as *nodes* and the connection to their parent commits as *edges*. Every time two edges leave a single node, a new branch is created. Git provides the feature to name branches, whereas the main branch is per default named *master*.

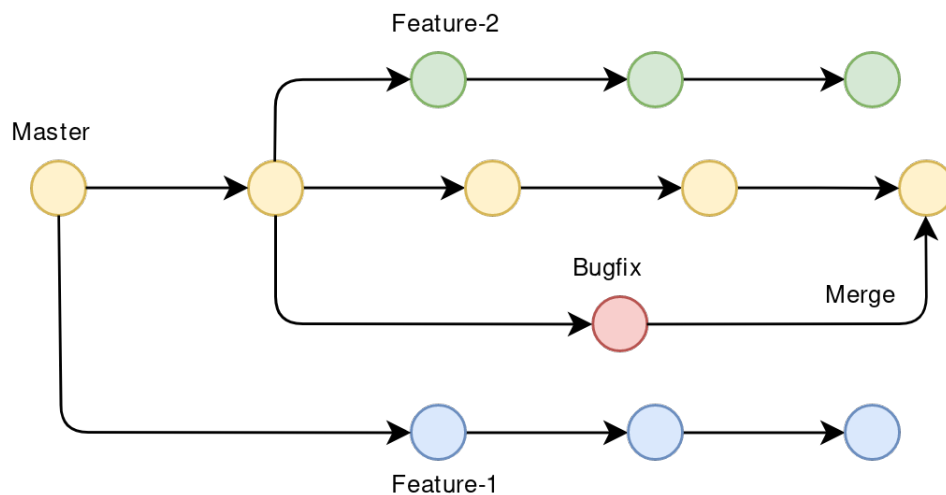


Figure 2.1: A Git commit history tree.

As shown in figure 2.1 two developer can for example create their own branch on which they can work unimpeded. If they finished their tasks and want to add their work to the master branch, they can now merge their changes. Git then tries to automatically resolve any conflicts which might have emerged from editing the same lines in a file. If that is not possible, it marks the conflicts and allows the user to manually correct them. After this resolution a new *merge commit* is created. A merge commit represents the merge of changes from two different branches.

With this methodology it is possible to work with many people or teams on the same project, without accidentally overwriting changes of another developer, whilst maintaining a clear history of all changes in the project.

Another important feature of Git is the possibility to set up a *remote*. A remote acts as a centralized repository developers can *push* their changes to or *pull* changes from other developers. It can for example be a distinct server, which is attached to some kind of network that is accessible by the developers. This feature allows developers to contribute to a project, as long as they have access to the remote's network. Git also supports several protocols such as hypertext transfer protocol (HTTP) or secure shell (SSH) to connect to a remote and thereby provide a simple user management layer.

2.1.2 Git User Roles

There exist two roles in Git, namely the *committer* and the *author*. Every commit in Git contains the email addresses and the names of those two people. The author of a commit is the person which actually contributed the changes in the files. The committer is the

person, which created the git commit. This is important to keep track of the original author of the changes.

Lets look at the case of an author contributing code to a project in an email with an attached patch file. If a maintainer of the project now applies the patch file and commits without setting the author, the information about the original author would be lost. Collected data indicates that in about 89% the author and the committer are the same person.

2.1.3 Internal Representation

Git's underlying storage and management solution for files is commonly described as a mini file system (FS) [1, p. 9] In the following I will thereby refer to its file management layer as *git-fs* and explain the most important aspects of it.

The representation of a single file in Git is called a *blob* object [1, p. 56]. A blob object is a file, which has been added to a *git-fs*. It is compressed and saved in the `.git/objects` directory under the respective secure hash algorithm 1 (SHA-1) hash of the uncompressed file. As a result, there exists a blob object for every version of every file of the project.

The SHA-1 hashing for unique file identifier might seem unsafe at first, but the probability of a SHA-1 collision is really low, roughly 10^{-45} . In 2017 Google managed to force a collision in a controlled environment, but it is really unlikely to encounter such a collision under normal circumstances [6]. This characteristic of SHA-1 hashing will become quite important in the design of the database later on.

As mentioned in the introduction 2.1.1 Git is used to store the state of a specific directory of an actual FS running under an operating system (OS). To represent a FS or to simply bundle multiple Git blob objects together, Git uses the tree object.

A tree object is a file, which has a SHA-1 hash reference to all underlying blob and tree objects as well as their names and FS permissions. To represent a subdirectory, a tree simply holds a reference to another tree object. With these tools git manages to build its own basic representation of a FS.

```
1      100644 blob 11d1ee77f9a23ffcb4afa860dd4b59187a9104e9 .gitignore
2      040000 tree ac0f5960d9c5f662f18697029eca67fcea09a58c expose
3      100644 blob 61b5b2808cc2c8ab21bb9caa7d469e08f875277a install.sh
4      040000 tree 8aaf336db307bdcab2f082bd710b31ddb5f9ebd4 thesis
```

Listing 1: A tree file example.

As stated before the commit is utilized to provide an exact representation of a state of the repository's files and directories. Just as blob object, the tree and commit files are also stored in the `.git/objects` directory under their respective hash.

```

1      tree        cd7d001b696db430b898b75c633686067e6f0b76
2      parent      c19b969705e5eae0ccca2cde1d8a98be1a1eab4d
3      author      Arne Beer <arne@twobeer.de> 1513434723 +0100
4      committer   Arne Beer <arne@twobeer.de> 1513434723 +0100
5
6      Chapter 2, acronyms

```

Listing 2: A commit file example.

As you can see in listing 2, the commit is just another kind of file utilized by Git, which contains some metadata about a repository version:

- The reference to a tree object, which represents the root directory of the commit's version of the project.
- A reference to one or multiple parent commits, to maintain a version history.
- The name and email address of the author.
- The name and email address of the committer.
- The Coordinated Universal Time (UTC) timestamps with UTC offset for the commit and author date.
- The commit message. A message with arbitrary text from the committer.

The commit is the most important object for this thesis. It contains crucial information such as the date of the commit as well as the email, which is needed to identify a contributor across several commits.

2.1.4 Features

Git provides a collection of high level abstraction tools to work with its underlying FS. In the following the most important features for data aggregation will be listed and briefly explained.

HEAD

To mark the current checkout of an repository, Git utilizes a special marker called *HEAD*. For instance, it is possible, due to this feature, to jump to the previous commit in history with `git checkout HEAD 1`.

diff

The diff is a tool to compare the state of two different commits in a repository. It allows to list all files which changed between those commits as well as showing the exact changes in the files

force push

Git allows to push to a remote with the `-force` flag, which is called a *force push*. This enables the users to rewrite every commit in the whole history tree. If another user has a git repository version

2.1.5 More features

Git provides many more features, which are not necessarily important for data analysis, but which might need be taken into account when aggregating the data. In the following some functionalities will be shown, which can lead to problems or irregularities in the gathered data.

rebase

It is possible to *rebase* branches. For instance a rebase can rewrite the commit history and change the branch point of a branch to another commit. This is for example a very powerful but also dangerous tool, as it implicitly changes the timestamps of the commits of the rebased branch.

filter-branch

In case somebody pushes a huge file, which significantly increases the size of the repository, or adds a file with secret information, for example a password file, git provides the *filter-branch* command. The filter-branch command removes a specified file from every commit in the history and thereby rewrites the history to the point, where this file was introduced to the repository. This command leads to similar problems as the rebase command, since it can severely change the history of a repository.

force push

Git allows to push to a remote with the `-force` flag, which is called a *force push*. This enables the users to rewrite every commit in the whole history tree. If another user has the old git repository version before the force push, they would now be incapable of simply pulling the newest version. Instead they would need to manually checkout the newest commit of the rewritten remote branch.

2.2 Data source

The biggest initial task for this thesis was the acquisition of data. Selecting a data source was a crucial step, as good data for analysis and evaluation is the backbone of this thesis. This section will list all requirements in detail and evaluate why I chose to use Github as a data source. Furthermore some functionalities of Github will be explained and a brief overview of the data provided by Github's application programming interface (API) will be given.

2.2.1 Requirements

The data source had to satisfy as many requirements, which were specified in Chapter 1.5, as possible.

To accomplish a meaningful analysis of commit conduct one needs a sufficient amount of commits. For instance it is necessary to have a few commits per weekday over a timespan of a month for a simple sleep rhythm analysis. If there are only 20 commits for a user over the past month there is probably not enough data for a representative analysis. To gather as many commits as possible we have to get access to as many repositories to which the targeted user contributed to as possible. Thereby the data source has to provide a way to dynamically explore repositories around a single user.

For analysis of companioned persons as described in Section 1.5.3 it is crucial to find users, which are likely to know each other. Optimally the data source provides a functionality for users to actively mark other users as their friends or colleagues.

To attack a company, as described in Section 1.5.3, or to spy on company members, as described in Section 1.5.1, the best case scenario would be to have full access to all repositories owned by the company. The data source thereby needs to provide some kind of representation for a company. Ideally there should also be a list of all company members for evaluation purposes of data mining findings.

A summary of the requirements to the data source:

- Realistic noise
- Real world data
- Large amount of repositories
- Access to all commits of each repository
- Access complete metadata for each commit
- Email to user association
- Methods to discover repositories a user contributed to
- Methods to discover possibly companioned contributor
- A representation of a company
- Access to members of a company

2.2.2 Github

I decided to use Github as a data source, for it is not only convenient to find uniform resource locators (URLs) for cloning repositories, but also provides solutions for most of the other requirements. It hosts one of the biggest collections of open source projects [11] with 64 million repositories, 24 million users and 1,5 million organizations [5]. Github also provides a well documented API for querying its metadata and there are libraries for most major languages, which provide an abstraction layer for this API. This API is publicly available and can be used by anyone registered on Github.

For instance Gitlab, one of Github's competitors, has much less data to offer. Gitlab doesn't provide detailed usage statistics, but they state that they only host about 100000 organizations, which is remarkably less than Github [4]. As Gitlab is an open-source project, there also are an unknown number of privately hosted Gitlab instances, which are completely inaccessible for unauthorized users.

On the other hand one of the downsides of using Github is, that we don't have access to all metadata. For example the full list of members for organizations is often inaccessible, as users actively need to opt-in to be publicly displayed as a member of the team. The internal team structures of organizations are not visible, as one needs to be a member of the organization to access those. Another problem are dangling email addresses, which are not related to any account any more. All commits made with this email address cannot be used any more for any analysis which requires a user to commit relationship. Even though some ground truth is missing, I decided to use this approach as it still was the most promising way to gather as much data and real world noise as possible, compared to other open source hosting services.

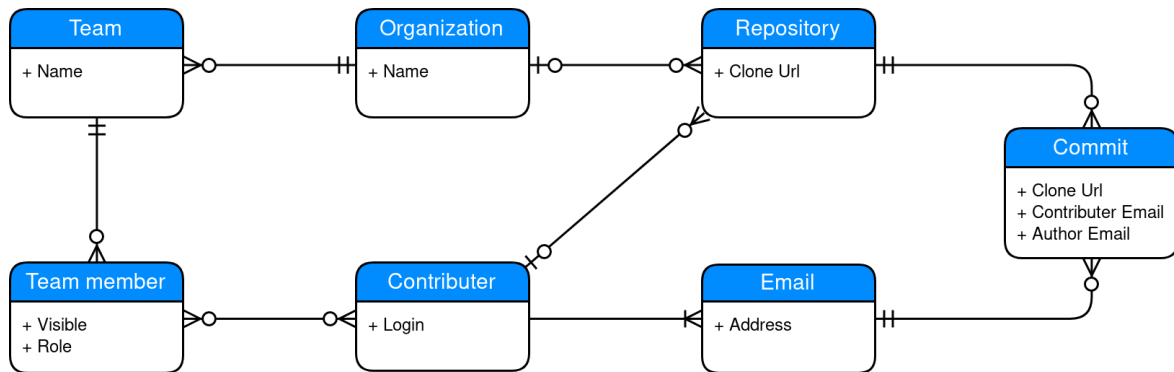


Figure 2.2: Simplified Github relationships in Crow's foot notation.

2.2.3 Github's features

In the following I will explain some of the features provided by Github, which cover the requirements listed in Section 2.2.1. Github offers some features, which are convenient

to for example find repositories a specific user contributed to or to find contributors which likely personally know each other.

Stars

A very crucial feature is *starring*. Every user has the possibility to star a repository to show appreciation or interest in this specific project. Hence popular repositories usually have a comparative big amount of stars. For instance the Github Linux kernel mirror has a star count of over 58000 ¹. Even though Github allows to query all repositories, which are owned or forked by a user, their API doesn't provide a method to get all repositories a user ever contributed to. However Github provides an endpoint to query all starred repositories of an user. In case a user stars a repository he contributed to, whilst not owning it, it is now possible to get this repository with this feature. Of course it is still not a reliable way to get all repositories a user ever contributed to, but it is a viable approach to get at least a few of them.

Follower

Another important feature is *following*. Every user can follow any other user to get informed, if they do specific things, like creating new repositories or starring repositories, or to simply show interest in or respect for their work. By getting all followed or following users, one might catch some friends of the user. It is also possible that a user follows the owner of a repository he contributed to. By using this feature it is thereby possible to get some additional contributed repositories as well some friends of the user.

Organizations

The third feature are *organizations*. An organization is used to host projects under an account, which is not necessarily led by a single natural person, but rather supports roles with different permissions and team structures. Github allows to query all repositories of an organization via their API. This enables us to link an organization to its owned repositories and as a result to perform analyses for users on a specific organization repository subset.

Generally organizations provide us with some important ground truth, even though the information might not be complete. Despite not knowing all members of an organization, we still get some useful information to estimate the tendency of precision of our knowledge extraction algorithms.

2.3 Data Aggregation

As mentioned in Section 2.2.2, I decided to use Github as my primary data source and to utilize their *Github APIv3* for this purpose. The aggregator and analysis program

¹'Linux kernel source tree' Github.com, <https://github.com/torvalds/linux> (accessed, 24.04.2018)

written for this thesis is named *Gitalizer*. In this section I will explain the technologies and methods used in the data aggregation process, the database structure and the interaction with Github's API. In the end some problems which occurred during the data collection will be shown as well.

2.3.1 Existing solutions

There are a number of existing solutions for accessing and collecting git metadata. In the following the practicability of these solutions will be evaluated based on our requirements.

GHTorrent

The *GHTorrent* project aims to provide Github's metadata to elude the limitations Github's rate limiting [7]. It provides representations for followers and commits, as well as organizations and organization members, but there are some crucial informations missing. GHTorrent only stores the main email address of a user and does thereby not support the handling of multiple emails, as commits are directly assigned to their respective Github user id. Commits miss information about additions and deletions in lines of code, which implicates that each commit would need to be scanned by a separate program again. Furthermore GHTorrent does not have the concept of *starring*, which makes it hard to reduce the amount of repositories to scan to a manageable size. Their database provides about 4 **TB!** (TB!) of data according to their website, which is too much information without very precise limitation ² It provides a vast amount of data, but at the same time it cannot be ensure, that the data is as complete as possible for a specific user. Modifying the GHTorrent code base and extending their database schema has thereby been judged as impractical.

ghcrawler

Microsoft provides an open-source crawler called *ghcrawler*, which is supposed to continuously fetch data from Github ³. Sadly their documentation is very bad and after diving into their source code, it appears that their crawler is for Github entities only and not for the underlying data of git repos.

²'The GHTorrent project' ghtorrent.com, <http://ghtorrent.org/> (accessed, 05.05.2018)

³'Github crawler' github.com, <http://github.com/> (accessed, 05.05.2018)

Alitheia-Core

Alitheia-Core is a Java data collector for git repositories. It is not actively maintained since more than three years and their documentation website is offline. Using this library seemed unpromising and unviable.

<https://github.com/mauricioaniche/repodriller>

RepoDriller

The *RepoDriller* project is project that aims to support researchers by providing easy access to repository data from Github ⁴. Despite providing a good solution for getting all necessary information from a repository, it provides no way to explore Github using *stars* or *following*. These features, as well as assignment of emails to a contributor via the Github API, would have to be added. As the program is also written in Java and I am no longer familiar with the language and its ecosystem, I decided to stick to writing my own solution.

2.3.2 Database

To store the gathered Information I chose a Structured Query Language (SQL) based solution. PostgreSQL provides excellent tools to ensure a high consistency in your database, namely check constraints, as well as great support for working with times, time zones and locations. The SQL database is used in combination with the Object-Relational Mapping (ORM) library *SQLAlchemy*.

The basic schema created for the purpose of this thesis consists of five ORM models. A model for commits, emails, repositories, contributors and organizations has been created. The latter is only important to validate results and is not actually used for knowledge discovery, as this is Github specific data.

⁴‘A tool to support researchers on mining software repositories studies’ github.com, <http://github.com/> (accessed, 05.05.2018)

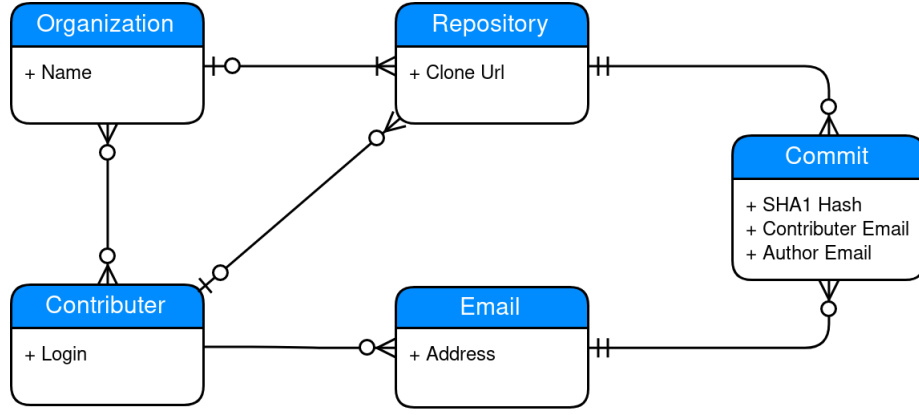


Figure 2.3: Gitalizer database relationships.

Every commit of each repository is saved in the database along with its SHA-1 hash and the two email addresses as in Listing 2. It is important to note that there is a many-to-many relationship in figure 2.5 between commits and repositories. This feature prevents duplication of the same commits from forked repositories. It is, for instance, a common practice to create a fork of a repository to develop without intervening with the main git repository of the project. As described in Section ?? the probability of a SHA-1 collision is highly improbable. By exploiting this feature, it is possible to enforce a unique constraint on the commit hash column, assuming that any duplicated commit hash actually results from a forked or copied repository. The formula for calculating the probability of such a collision is:

$$p \leq \frac{n(n-1)}{2} * \frac{1}{2^b} \quad (2.1)$$

Without this mechanism it could be possible that the same commit of a contributor could be used multiple times as a result of repository forking. After collecting 43 million commits from Github and actively ignoring obvious project forks, there are still 49 million references between commits and repositories. This means that about 13% of gathered commits result from forked repositories which can not easily be identified as such. Considering Formula 2.1, the probability of a collision for 49 million hashes on the 160 bit SHA-1 hash would be about $8.21 * 10^{-34}$.

As stated above each commit is also saved with its respective email addresses. There exists a one-to-many relationship between contributors and emails, as every contributor can obtain an unlimited amount of email addresses. It is necessary to connect all email addresses commit to a specific contributor, to unambiguously assign all commit to their respective contributor. This relationship does not have a `NOT NULL` constraint as it happens quite often that an email address can not be assigned to any person. Looking at the collected data it appears that roughly 20% of all collected email addresses from Github are no longer connected to an active user.

As stated in Section 2.2.3 Github provides a way to organize several people in organizations and teams. As one of the goals of this thesis is to see if it is possible to detect member of an organization in open-source projects, a model for organization has been created as well. This data can then be used to check against the results of this research's results.

2.3.3 Gitalizer

The Program written for this thesis features data aggregation, preprocessing, knowledge extraction and visualization. Gitalizer uses a PostgreSQL database for data storage and data consistency checks as described in 2.3.2. For interaction with the Github API the *PyGithub* library is used, which provides a convenient abstraction layer for requests and automatically maps JavaScript Object Notification (JSON) responses to *Python* objects.

The data aggregation module of Gitalizer is capable of several approaches for gathering data. In the following we will look at those approaches in detail.

Git repository

Gitalizer can scan any git repository from a SSH or HTTP URL as long as the current user has access to it. The repository is cloned into a local directory. After the cloning finished the actual scanning process begins. During the scan, we git checkout the HEAD of the current default branch for this repository and walk down every reachable commit of the Git history. The program saves all available metadata for each commit in its database, namely the emails, timestamps as well as additions and deletions to the project in lines of code.

After this scan we are still missing a lot of information. There exists no unique identifier of an author or committer, as names may change or can be ambiguous and a single person can have multiple email addresses. The problem with the simplicity of Git is that there exists no concept of an user. Thereby we cannot easily link several email addresses to a specific contributor without additional metadata.

Github Repository

To tackle the problem of missing metadata in 2.3.3, I used the Github API to get some of the missing metadata. The general approach is the same as in the previous scan method. The repository is cloned and locally scanned. However, a request to Github is issued every time an email is found, which we do not already have linked to a contributor. Github allows to link multiple email addresses with a single user account and automatically references the respective user in their own API commit representation. With this additional metadata we gain ground truth about the identity of an author or committer.

Anyway this approach does not work, if the user of a commit removes the email, which has been used for the commit, from his account or if the user completely

deletes their account. If this happens and the email contributor relationship has not already been created, there is nothing that can be done and these commits need to be handled later on in the preprocessing of the data.

Github User

To try getting all repositories of a specific user, a new functionality has been added, which highly utilizes the Github API. At first several requests are issued to get all repositories of the specified user, as well as all starred repositories of this user. During the repository exploration, every relevant repository is added to a shared queue, let's call it "repo-queue", which is then processed by a multiprocessing pool of workers. Each worker process pops another entry from the "repo-queue" and scans a single repository as described in 2.3.3.

Connected users and organizations

For detection and analysis of connections between contributors over multiple repositories additional user repository discovery as described in 2.2.1, another feature has been added to Gitalizer. Gitalizer is able to achieve this by not just scanning a single user, but also scanning the repositories of all following and followed users as described in ???. For this task two different worker pools are utilized. The user discovery pool is initialized with a shared queue, let's call it "user-queue", of all users we need to look at. This worker pool simply searches for relevant repositories of each user and passes the repository URL to a second shared queue. The second pool then processes this "repo-queue" as described in 2.3.3.

For organizations it is nearly the same approach. Initially all repositories, which are owned by the organization, are added to the "repo-queue". All publicly visible organization members are then added to the "user-queue" and processed as described above.

2.3.4 Database optimization

As the database kept growing, it became the bottleneck in the aggregation process several times. As a result, adjustments in the database schema, PostgreSQL parameter tweaking and migration to better hardware were necessary. The first considerable slowdown occurred after reaching about 12 gigabytes (GBs) of data. At this stage the database write and read operations took longer than the actual aggregation process and the whole machine started to become unresponsive because of high I/O load.

The performance of the database then needed to be continuously tweaked in several steps. The first countermeasure was the reduction of commits using deduplication by using the features of the SHA-1 hash as stated in Section ?? The ignoring of forked repositories reduced the amount of entries in the relation table between commits and repositories by another 26%.

The next performance boosts were achieved by disabling or loosening several fail-safe mechanisms of PostgreSQL, namely ‘synchronous commit’ and ‘write ahead’ parameter, which are designed to save data on a system crash. As there is no important or critical data handled it was acceptable to pass on these mechanisms, and trade safety for performance.

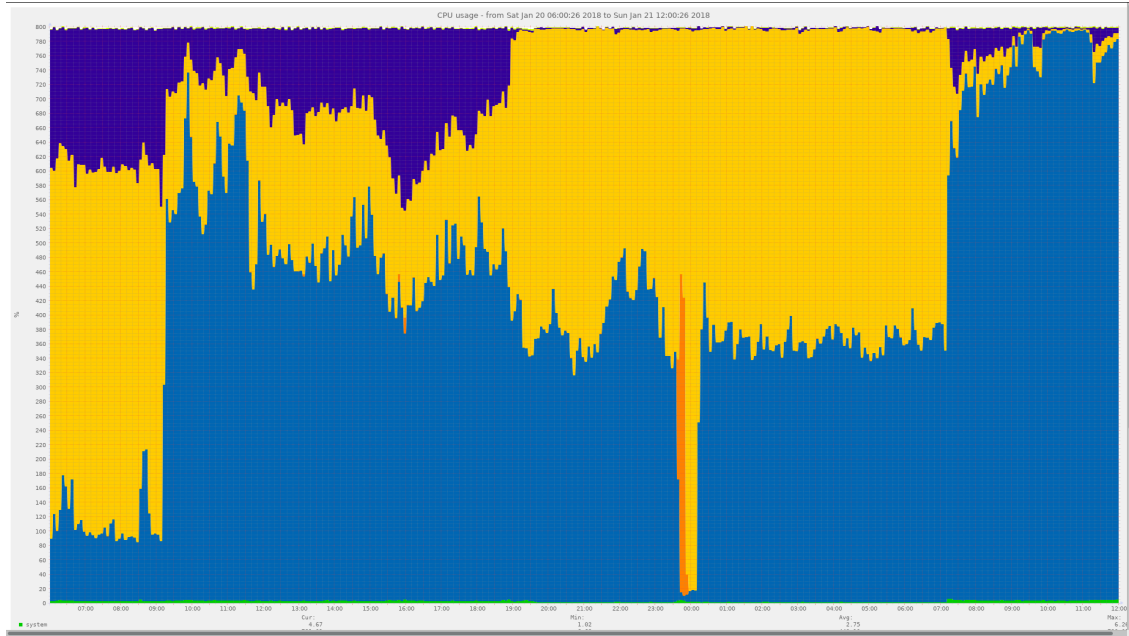


Figure 2.4: The CPU load of the aggregation server during optimization.

After renting a root server and deploying Gitalizer onto it, the aggregation process worked really well, until the database size hit about 25 GB. In Figure you see the overall central processing unit (CPU) load right before optimizing several SQL queries by caching intermediate results and increasing the cache size of PostgreSQL. The dark blue represents the I/O wait time while the light blue represents the actually used processor capacity by the aggregator. Due to complex and numerous SQL queries the server became partly unresponsive and the aggregation process began to stall.

After many improvements the server can now run with 16 worker threads and roughly 38 GB of data without any signs of slowdown whilst using the rate limit to capacity.

2.3.5 Incremental aggregation

To ensure a constantly up to date database, Gitalizer needed to be capable of fast rescans of repositories. The initial scan of a repository always includes cloning, scanning the whole repository and writing it into the database. After a repository is scanned completely at least once it is marked as as such and will never be completely scanned again. All following runs then only clone the repository and scan the newest unknown

commits. These are discovered by performing a breadth first search until no new nodes are found.

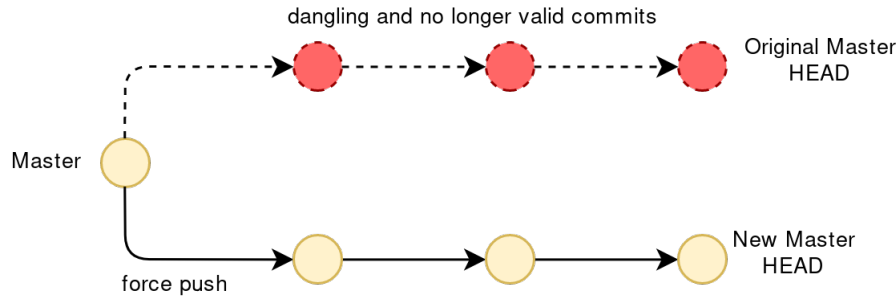


Figure 2.5: Gitalizer database relationships.

As explained in Section 2.1.5 it is possible to rewrite commits and force push them. This scenario needs to be explicitly handled since force pushes can completely alter the history of a git repository, which can subsequently lead to a split in the Git history and leaves dangling commits. As the complete history of a repository is stored inside the database, Gitalizer needs to detect a force push by walking down the git history tree until it finds known commits. If any of these commits has children, which are not in the newly scanned commits, a force push took place and the old commit history has to be truncated, since it is now outdated and irrelevant. An example scenario can be seen in Figure 2.5, where all red commits represent the old commit history, which needs to be truncated.

2.3.6 Problems

During the development of the data aggregator I experienced a few problems and edge cases which needed to be handled. The earliest and most delaying problem was the rate limit of the Github API, which limits to 5000 requests per hour. Beside this rate limiting there also is an abuse detection mechanism, which triggers, if too many requests are fired in a short amount of time. The solution for this problem resulted in various hacks, which include random wait times to detain those mechanisms from triggering.

The first version of the aggregator didn't clone and scan the repository locally, but rather gathered all information from the Github API endpoints. This approach worked well until the aggregator hit the official repository of *Nmap*, which has about 11.000 commits and took over three hours to scan. Soon I realized that this would severely slow down my research and I then started to continuously minimize the amount API calls issued by Gitalizer. A connected user scan of my own Github account led to about 600.000 commits and took about one and a half day on the final working version of Gitalizer, to provide you with a reference of scale.

After implementing multiprocessing, I managed to hit the rate limit again, as I was now issuing requests to the API with multiple threads. To fix this issue I implemented a wait and retry wrapper around every single function call or object access, which triggered a call to the Github API. Afterwards the aggregator was capable of running multiple days without worker processes silently dying or finishing with incomplete data.

Fine tuning the edge cases and the handling of the API took about 3 months, since there were many problems such as unpredictable error responses from Github, missing data in queries or simply unknown or broken encodings in Github's metadata.

A big throwback became apparent as I discovered that PostgreSQL automatically normalizes UTC timestamps with any offset to the `(UTC +0)` timezone. As a result of this normalization, the exact time of the commit admittedly stays the same, but the crucial metadata about the offset is lost. As a consequence the commit model needed to be adjusted, as the UTC offset had to be stored explicitly, and the whole commit aggregation process was started from scratch.

Another problem occurred during the local scanning of the repositories. The library used for interaction with Git *libgit2* issued *stat* Linux syscalls during a diff operation for each file, which changed between those commits, to check if there were any local uncommitted changes. Anyhow the repositories, which were locally scanned, were cloned with *bare* mode. This means that there exists no project root directory, but rather only the git internal representations of those files, which makes the behaviour stated above unnecessary and unwanted. As a result all processes slowed severely down due to high I/O wait times, because of *stat* syscalls on non existent files. Luckily after reporting the issue ⁵ it was resolved in a week and I was able to continue developing with my own compiled version of the *libgit2* library.

⁵'Unnecessary syscalls on bare repository' github.com, <https://github.com/libgit2/libgit2/issues/4480> (accessed, 25.04.2018)

CHAPTER 3

Implementation

After collecting all necessary data as shown in Chapter 2, I will now begin to analyse this data. In this chapter the approach for several attacks, as listed in Section 1.5 will be introduced and the attack’s goals recapitulated. The possible applications for the gained knowledge will be stated and the implementation of and requirements to the respective algorithm will be explained for each attack.

3.1 Holiday and Sick Leave Detection

The information about anomalies in the regular work pattern can be a valuable information for several parties. Usually only few parties know about the holiday or sick leave times of a person. To know if a persons tends to become sick more often or for long times is a dangerous intrusion into a persons privacy. For instance this could be abused by head hunters or personnel managers to cull possible employees with too high sick leave rates and thereby reduce the job prospects of the target.

For employers this might be convenient to detect anomalies in the productivity of an employee. In case an employee doesn’t commit on a regular basis for several days, this behaviour would be instantly visible with this method.

Another attack vector could be to look at the correlation of miss-out between several persons. This attack could even be performed by an outsider on a commercial open-source project, if the employees of the targeted company are known. The information gained by this attack could be quite delicate, as it could reveal relationships between persons. This attack is heavily inspired by an article about data mining articles from the popular German weekly magazine *Der Spiegel* written by the David Kriesel [9].

3.1.1 Implementation

One requirement for this algorithm is the detection of a regular work pattern for a given interval. It must have the ability to adjust to a changing work pattern, but at the same time needs to be capable of detecting anomalies in this pattern. If the attacker wants to look at multiple people, some kind of measure for similarity in the missing time patterns has to exist.

The input for this analysis is the intersection between all commits from the considered repositories and all commits from the considered contributors. The commits’ meta data

used for this analysis are time stamps as well as additions and deletions in lines of code.

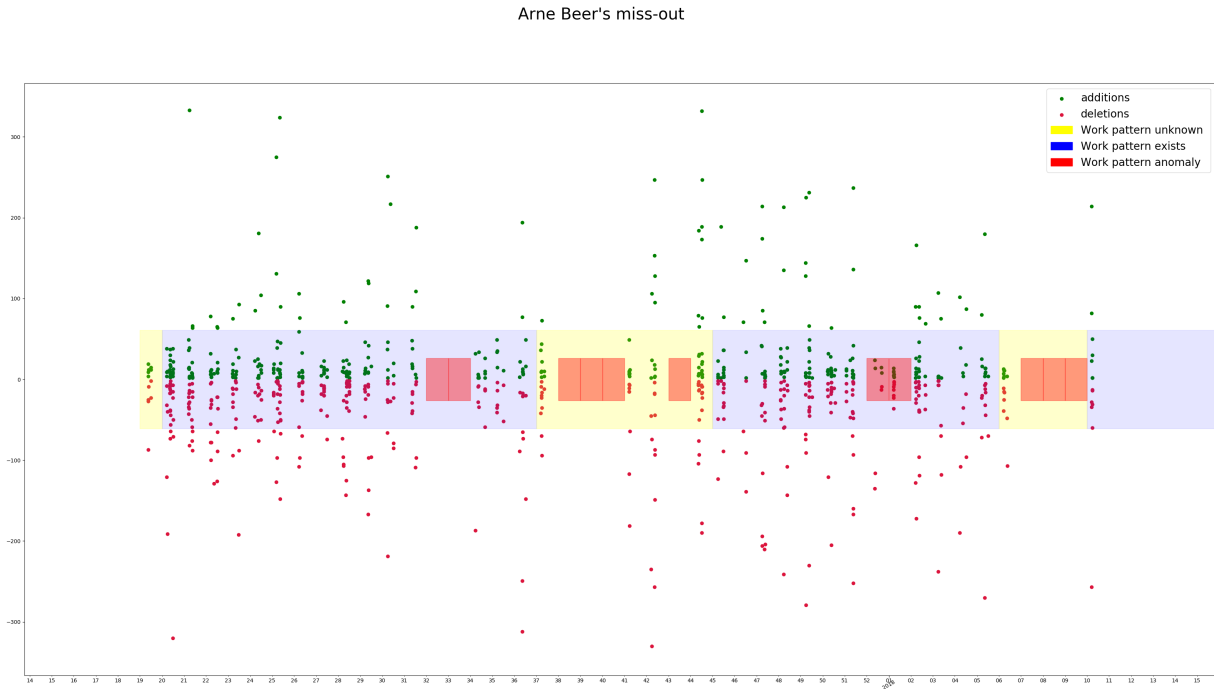


Figure 3.1: The work time analysis of the author.

The analysis of the data is a chronological scan of all commits for specific user. Before performing the actual analysis, the data is converted into an easier to handle format. It is really difficult to measure productivity in lines of code committed or in the amount of commits made by a person, as they do not necessarily display the amount of work that have been put into those commits. As a result I decided, that a day counts as a work day as long as at least single commit has been made during the day. The preprocessed data is thereby equivalent to the days a contributor worked on, ordered by the week of the year.

```

1 def analyse(weeks):
2     prototype = None
3     for index, week in weeks.items():
4         next_six_weeks = weeks[index:index+future_lookup]
5         if not prototype:
6             # See if there is a prototype in the next few weeks.
7             prototype = find_prototype(next_six_weeks)
8
9             # Check if this specific week is a anomaly

```

```

10         check_anomaly(prototype, week)
11
12         continue
13
14     prototype_exists = prototype_exists_in_next_weeks(next_six_weeks)
15     if not prototype_exists:
16         # We couldn't find the prototype in the next few rows
17         # Try to find a new prototype
18         prototype = find_prototype(next_six_weeks)
19
20     check_anomaly(prototype, week)
21
22
23 def check_anomaly(prototype, week):
24     if week.working_days == 0:
25         save_anomaly(week)
26
27     if prototype is not None:
28         different_days = week.working_days - prototype.working_days
29         // A single day variance is acceptable
30         if different_days >= 1:
31             save_anomaly(week)

```

Listing 3: Miss-out analysis algorithm.

The algorithm inspects every week work pattern of a given interval. At the beginning a new *prototype* is tried to be found, which is performed in the function `find_prototype` in Listing3. A prototype is a representative week work pattern, which resembles the average work day pattern of the next weeks.

The function performs a simple iteration over a given interval to find a work day pattern, which occurs more often than a given threshold. If a prototype is found, we are capable of identifying anomalies that deviate from this pattern.

For each following week it is firstly checked if this week is a anomaly for this prototype. Anomalies are simply detected by comparing the amount of working days of the prototype and the currently looked at week. The real difference in the working pattern is not suitable for this analysis, as it produces too many false positives for employees with flexible work time.

Secondly it is checked if there exists a week in the near future, which is identical to the prototype. If there is no week identical to the prototype in the near future, the current prototype is reset and a new prototype needs to be found.

In case no prototype can be found, anomalies cannot be easily identified, as there exists no pattern to check against. Only obvious anomalies, namely weeks without a single

work day, will then be marked as such.

3.2 Working hours

This attack aims to extract information about the working hour behaviour of a target. This should be achieved by displaying the pattern of the target in form of a weighted scatter plot sorted by hour per weekday and by comparing those patterns between several targets. There are several possible vectors for this attack:

- Information about the sleep rhythm of the target.
- Detect whether the target is a person working regular shifts from Monday to Friday or rather an open-source contributor working in their leisure time.
- Detect anomalies such as automated programs, which contribute to a project on a regular basis.
- Compare the working hour patterns of several people in the same project or organization. This could, for instance, be used to infer relationships between colleagues, based on an equal working shifts.

Additionally a clustering will be performed to find anomalies, common patterns and to evaluate the results of this analysis. As we are only interested in contributors with a representative amount of commits, all contributors with less than one hundred commits in the last year have been excluded. This reduced the amount of considered contributors from 175.000 to about 10.300.

3.2.1 Implementation

The initial data for this analysis are the commit timestamps of the target, as well as the Github employee information for verification. These commit timestamps are then converted into a different format, which represents the occurrences of commit per hour per weekday over the last year. The result is a simple vector with length 168, which corresponds seven days with 24 hours each. I will refer to this representation hereafter as a *punchcard*.

```
1 def preprocess(commits):
2     punchcard = [0] * 168
3     for commit in commits:
4         hour = commit.commit_time.hour # returns 0-23
5         weekday = commit.commit_time.weekday() # returns 0-6
6
7         index = weekday*24 + hour
8         punchcard[index] += 1
```


Listing 4: Data preprocessing

The data transformation is achieved by incrementing the field of the respective weekday and hour by one for each commit, as can be seen in 4. The resulting punchcard vector is then stored in the database for faster and easier analysis in the next steps.

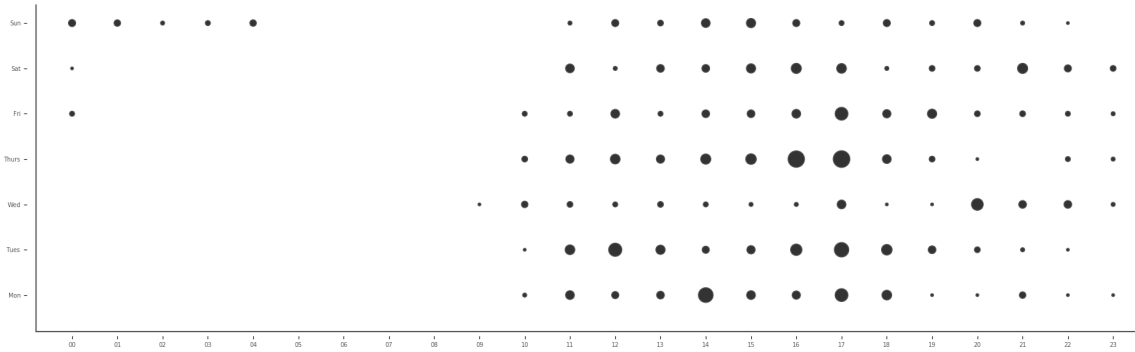


Figure 3.2: Punchcard of the author.

3.2.2 Punchcard Clustering

To find common work patterns, several cluster cluster algorithms have been performed on the aggregated data. The Python *scikit* framework has been used for this purpose, as it features nine different clustering methods and provides good documentation and abstraction from the underlying clustering logic ¹. For the task of finding similar punchcard patterns in the data, a clustering algorithm is required, which can operate on a high-dimensional dataset with an unknown amount of clusters Scikit provides three different clustering algorithms, which can handle an unknown amount of clusters.

Mean Shift

Mean shift is a clustering methods, which performs an operation similar to a gradient descent, which shifts all adjacent data points to their center [2]. The goal of this algo-

¹‘Clustering’ [scikit-learn.org, http://scikit-learn.org/stable/modules/clustering.html](http://scikit-learn.org/stable/modules/clustering.html) (accessed, 24.04.2018)

rithm is to find a representative centroid for each cluster and assign each data point to a cluster. This methodology proved to be too aggressive for the current data.

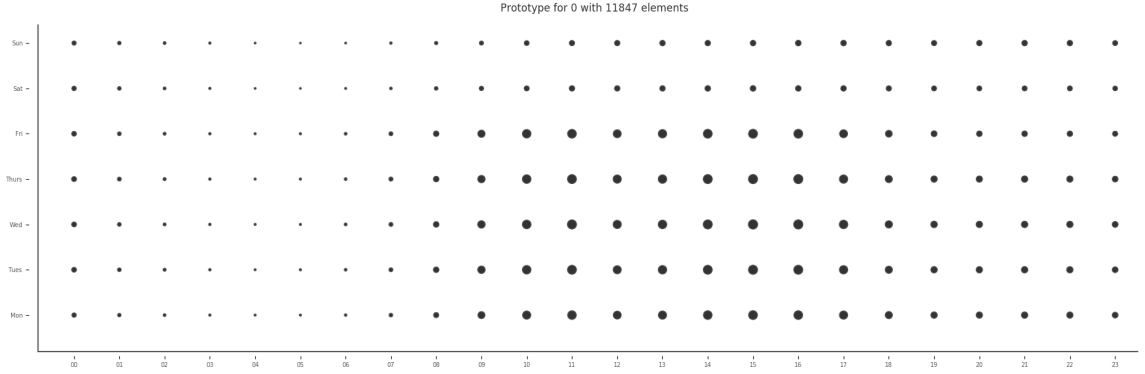


Figure 3.3: Super cluster centroid found by mean-shift clustering.

Despite trying a wide range of values for the bandwidth, which is the measure of distance used for detecting adjacent points, this algorithm always created a super cluster, which contained more than 89% of all data points. Such a super cluster can be seen in Figure 3.3. The other 11% were invariably small clusters representing extreme outlier or strange patterns, which do not resemble any of the expected patterns for normal work shift or leisure time developers. An example for an extreme outlier can be seen in Figure 3.4.

My assumption is that the density for the major part of the provided data is too high and equally distributed around the centroid of the super cluster. Thereby all those data points are slowly shifted to this single centroid. As it is difficult to debug 168-dimensional space, I decided that an profound analysis would be too time consuming and to try the next solution.

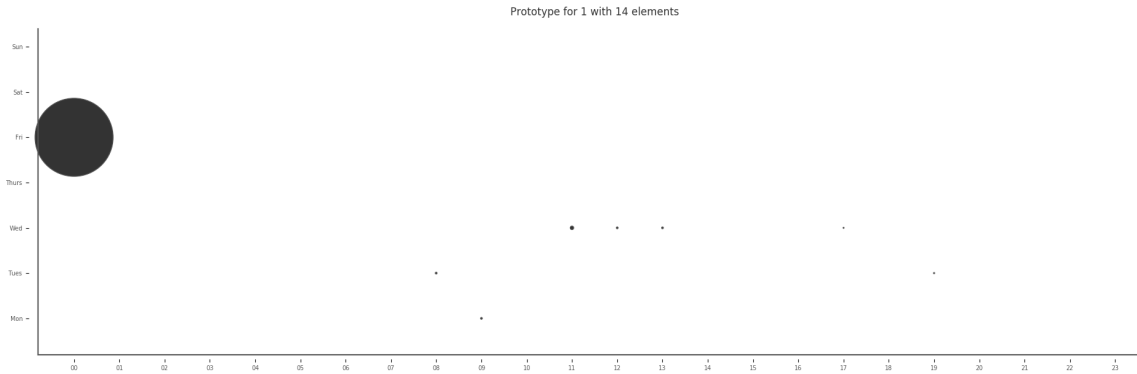


Figure 3.4: Extreme outlier centroid found by mean-shift clustering.

DBSCAN

The Density-based spatial clustering of applications with noise (DBSCAN) clustering algorithm operates by creating clusters of transitively connectable data points with a maximal distance between each point. ?? It is highly scalable and performant, even for large data sets, because of this scalability it was my first choice. Unfortunately it produces very similar results to the mean shift clustering algorithm 3.2.2 as it finds a super cluster very similar to Figure 3.3.

I assume that this algorithm suffers from similar problems as the mean shift approach, which are high density of data points without clear borders. Thereby this methodology manages to most part of all data points transitively from a single starting point, but creates a huge amount of mini clusters, when supplied with smaller values for the maximum distance between data points.

This algorithm also manages to find extreme outlier clusters, but it is not suitable for the purpose of this thesis, due to the extremely low granularity on the data inside the super cluster.

Affinity Propagation

The Affinity Propagation algorithm considers similarities between all data points to find clusters [3]. This clustering algorithm features a promising approach, as it utilizes a

method similar to message passing, to find an *exemplar*, which resembles the representative of a cluster, and its surrounding cluster member. Affinity Propagation was the only available clustering method that was detailed enough to find interesting patterns without creating a super cluster. About 200 different patterns have been discovered using this methodology. However it has to be noted, that this clustering method is sometimes a little too detailed, as it, for instance, split very similar patterns into two or more different clusters. Additionally the memory requirements for this method scale quadratically, for non-sparse sets, with the number of the data points [3, p. ii]. About 13 GB memory have already been used with a sample of roughly 10.000 data points. This algorithm becomes thereby impractical for analyses on the whole dataset, but it works for smaller analyses and is thereby suitable for the validation of this thesis.

3.3 Geolocation

This attack aims to extract information about the actual location on the globe of a target. The algorithm tries to firstly determine the exact timezone of the target and then to exclude any countries or states, which do not match the observed timestamps. There are several possible vectors for this attack:

- Detect the main country or timezone of a target.
- Analyse where and how often a target travels to a specific location.
- Determine the location of a workforce of a company.

3.3.1 Implementation

The data used for this analysis are all commit timestamps of the target, as well as the target's Github location for verification. In the following I will explain the algorithm used to detect the different timezones and determine the main location of the contributor. Several external data sources have been used to accomplish this:

IANA Database

The Internet Assigned Numbers Authority (IANA) provides a free to use database, with all timezones and the respective daylight saving times (DSTs) switches for each year. It also provides the exact UTC offset and offset switches for all timezones and thereby a relation to all countries.

Natural Earth

The Natural Earth organization provides extremely detailed free to use and up-to-date geological data with timezones, countries and even states on a resolution down to 1:10m.

Pycountry

As the codes and names used by the IANA database and Natural Earth don't always are identical, another layer with more information about countries and country codes was necessary. To match non-assignable timezones to their respective country on the map, the library *pycountry* was added.

OpenStreetMap and geocoders

As the codes and names used for states and provinces by the IANA database are not necessarily identical to the Natural Earth database, another solution for getting the relation between timezone and state was necessary. To match non-assignable timezone strings to their respective state, Gitalizer occasionally issues requests to the *OpenStreetMap* API with help of the python *geocoder* library.

To assign UTC offsets to their possible timezones a special reverse mapping of the existing IANA database was necessary. The python library *tzinfo* provides interaction with the IANA database, but this adapter is only capable of resolving timezones to their respective UTC offset.

As a result I wrote my own adapter, which extractes the data from IANA, with help of *tzinfo*, and saves it into the Gitalizer database. The database model is named *TimezoneInterval* and contains the timezone identifier, the UTC offset and the exact start and end of this specific timezone interval. This table only needs to be populated once, but it needs to be updated as soon as a new version of the IANA database is released.

```
1 def get_travel_path(commits):
2     travel_path = []
3     current_location = None
4     last_valid_location = None
5     change_at_day = None
6     location_candidate = None
7
8     for commit in commits:
9         commit_time = commit.commit_time
10        zones = find_timezones(commit_time, commit.commit_time_offset)
11
12        # Create the initial timezone
13        if current_location is None:
14            current_location = {
15                "set": set(zones),
16                "start": commit_time.date(),
17                "end": commit_time.date(),
18            }
19            last_valid_location = commit_time.date()
20
21        continue
```

22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

```
# Get possible timezone candidates for this commit and intersect them with
location = set(zones)
intersection = location & current_location["set"]

# Check if the possible timezones of this commit matches any timezone of th
if len(intersection) > 0:
    # By reassigning the intersected set we gain additional precision by co
    current_location["set"] = intersection
    current_location["end"] = commit_time.date()
    last_valid_location = commit_time.date()

# There is no match between the possible timezones and the current set.
# In this case we need to check if this is a single occurrence (anomaly) or
# if this is an actual change.
else:
    # No change_at_day exists, but we detected a change
    # Remember the change. If this change lasts for at least a day it will
    if change_at_day is None:
        change_at_day = commit.commit_time.date()
        location_candidate = {
            "set": set(zones),
            "start": commit_time.date(),
            "end": commit_time.date(),
        }

    # No change detected
    if change_at_day is None:
        continue

    # There was an anomaly, but not for a whole day.
    # This could for instance be a developer committing from a remote server.
    if change_at_day <= last_valid_location:
        change_at_day = None
        location_candidate = None

        continue

    # The change is not older than a day
    # ignore it until the change lasts for longer than a day
    if change_at_day <= last_valid_location:
        continue

    # There exists a change from the last day.
```

```

66     duration = current_location["end"] - current_location["start"]
67
68     # The current_location set only existed for a single day.
69     # This is most likely an outlier. Thereby drop it and restore the previous
70     if duration < timedelta(days=1) and len(travel_path) > 0:
71         last_location = travel_path.pop()
72         last_location["end"] = current_location["end"]
73         current_location = last_location
74
75         # Check if the old location and the current candidate actually match
76         # If that is the case drop the candidate and completely replace the cur
77         intersection = location_candidate["set"] & current_location["set"]
78         if len(intersection) > 0:
79             # Update current_timezone
80             current_location["set"] = intersection
81             current_location["end"] = commit_time.date()
82
83             # Reset candidate and last_valid_location occurrence
84             last_valid_location = commit_time.date()
85             change_at_day = None
86             location_candidate = None
87
88             continue
89
90         # We detected a change and it seems to be valid.
91         # Save the current timezone and set the candidate as the current timezone.
92         travel_path.append(current_location)
93         current_location = location_candidate
94         change_at_day = None
95         location_candidate = None
96         last_valid_location = commit_time.date()
97
98     current_location["end"] = datetime.now().date()
99     travel_path.append(current_location)
100     return travel_path

```

Listing 5: Algorithm used to detect changes in the target's location by analysing the UTC offsets of git commit timestamps

The algorithm in Listing 5 iterates through every commit and determines in which timezone the contributor could have been at commit time. For each following commit it is checked, if there is a intersection between the possible timezones of the last commits and the current commit. This is usually the case, if the contributor did not travel to another timezone. But it is possible that a change in the timezone happens, even though

the contributor did not travel. This is due to DST, which is something that can be used to improve the precision of the location.

Germany for instance enforces DST and switches between the UTC offsets +1 and +2. Angola on the other hand does not have DST and thereby has a continuous offset of +1. For instance, if commits during a small time interval in the winter from a German contributor are considered, it cannot be determined whether he lives in South-Africa or in Western Europe. But if the commits of a whole year are considered, it can be concluded, due to the intersection of the possible timezones of previous commits, that the contributor has to be in a country, which enforce DST and switches between the offsets +1 and +2.

In case no intersection between the timezones can be found, it needs to be determined, whether the contributor actually committed or the change happened through some other event, such as the commit from a remote server in a different timezone. For this purpose, all timezone switches, which do not stay longer than a day or happen at the same day as a commit from the previous location, are marked as insignificant and ignored.

The algorithm returns a chronological list of all detected and as significantly ranked locations with their respective time interval.

```
1 def find_home_location(travel_path):
2     home_location_candidates = []
3     home_location = None
4     found = False
5
6     # Try to find the current home location and narrow it down as good as possible:
7     for location in travel_path:
8         duration = location["end"] - location["start"]
9
10        # Try to find a set which intersects with the current set
11        for candidate in home_location_candidates:
12            intersection = location["set"] & candidate["set"]
13
14            # Found an intersection, set the new intersection and increment days
15            if len(intersection) > 0:
16                candidate["set"] = intersection
17                candidate["days"] += duration.days
18                found = True
19                if candidate["days"] > home_location["days"]:
20                    home_location = candidate
21
22            break
23
24        # Found no matching location, create a new candidate
```



```

25     if not found:
26         location["days"] = duration.days
27         home_location_candidates.append(location)
28     else:
29         found = False
30
31     if not home_location:
32         home_location = location
33
34     return home_location

```

Listing 6: Methodology used to determine the main location of a target, based on the information gained from the function in Listing 5

To detect the home location of a contributor, the algorithm in Figure 6 is used. The parameter provided to this function are the results of function `get_travel_path` in Figure 5. The algorithm simply tries to determine the best possible set of timezones, which exists for the longest duration.

During this process further intersections of matching timezone sets are made to further increase the precision of the home location. The result of this function is a non empty intersection of timezone sets, which persisted for the longest time period, compared to all other possible non empty intersection sets.

CHAPTER 4

Evaluation and Interpretation

In the last chapter I showed the implementation of several possible attacks, which could be performed on the gathered data. This Chapter will now attend to the evaluation of all results gained from these attacks. I will present the extracted information from each algorithm and compare it to real world ground truth. This information will be then be explained and audited in terms of precision and reliability.

4.1 Holiday and Sick Leave Detection

Figure 3.1 shows the analysis of the author for his work repositories. The y-axis shows the additions or deletions per commit, the x-axis shows the week of a year. For better verification and evaluation of the results, a scatter plot with the additions and deletions per commit has been added on top of the miss-out graph.

The evaluation of this algorithm turned out to be quite difficult, as there is no publicly available information about sick leave or holiday. For the purpose of this thesis I had to use anonymous statistics of several friends and colleagues to evaluate the algorithm. The algorithm successfully manages to find all anomalies, which occurred in the last year, for all seven regarded users.

An unexpected side effect of detecting prototypes is that the algorithm also finds inconsistencies in the work routine. For instance between week 37 to 45 in Figure 3.1 I was forced to reduce my working hours due to legal questions and continuously shift hours and working days for several weeks. It is hard to interpret those inconsistencies without more contextual information, but nevertheless it provides the fact that something happened during this time.

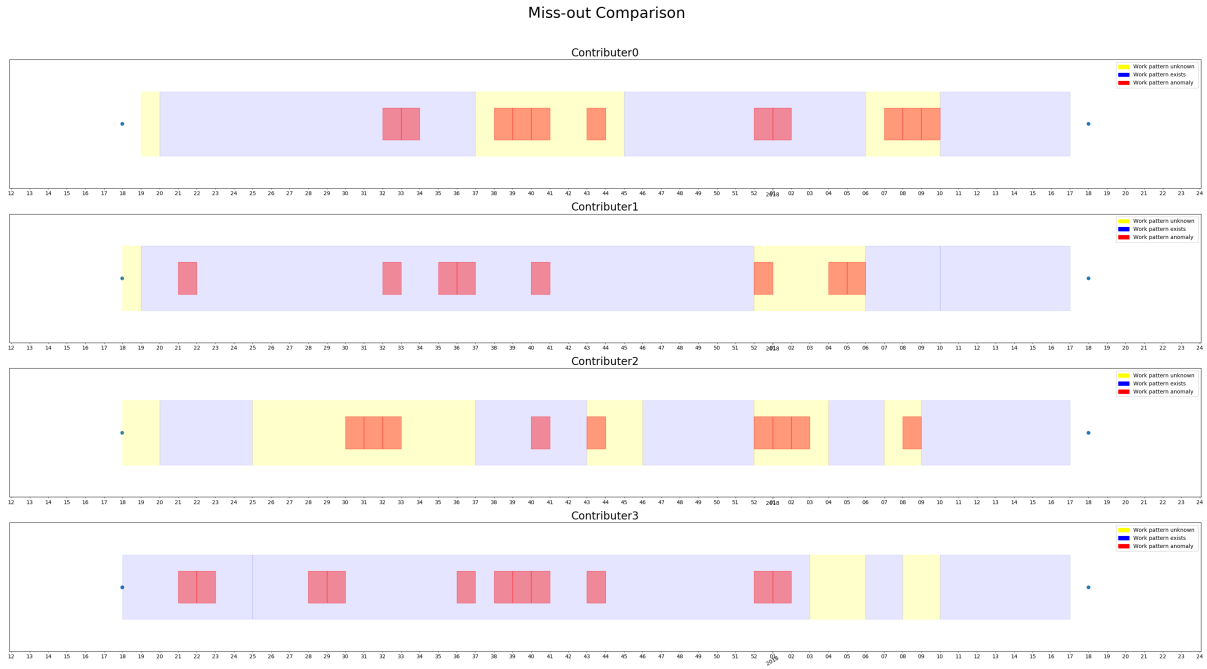


Figure 4.1: The miss-out analysis of several employees.

In Figure 4.1 the comparison between multiple employees can be seen. Contributor0 and Contributor2 are working on flexible work time, which reflects in the inconsistencies of those contributors, while the other two contributors have regular working hours.

4.2 Working hours

4.2.1 Sleep rythm

To evaluate the significance of the punchcard in terms of sleep rhythm analysis, a small survey in a closed community has been made. A selected group of ten people, who know each other well has been selected, for this purpose. Furthermore a subset of four people has been chosen, which were going to be evaluated. All ten people then needed to assign the punchcard of the chosen subset to a specific person.

To this end, three quite similar patterns, where one contributor has a very regular sleep rhythm, Graph 2 in Figure 4.5, and one pattern of a practically inexistent sleep rhythm, Graph 1 in Figure 4.5, have been selected.

The results of this survey where just as expected. For the three similar graphs, no significant results could be assessed, as the assignments where more or less random. The contributor with the irregular sleep rhythm on the other hand got correctly assigned in all in every survey.

Sadly, as such a survey needs very specific targeting and a long lead time for data collection, it could only be executed on a small set on subjects. Anyhow the result of this survey proofs, that there exists a correlation between the structure of the punchcard and the actual commit and sleep behaviour of an contributor, even if it can only be accurately assigned in extreme cases.

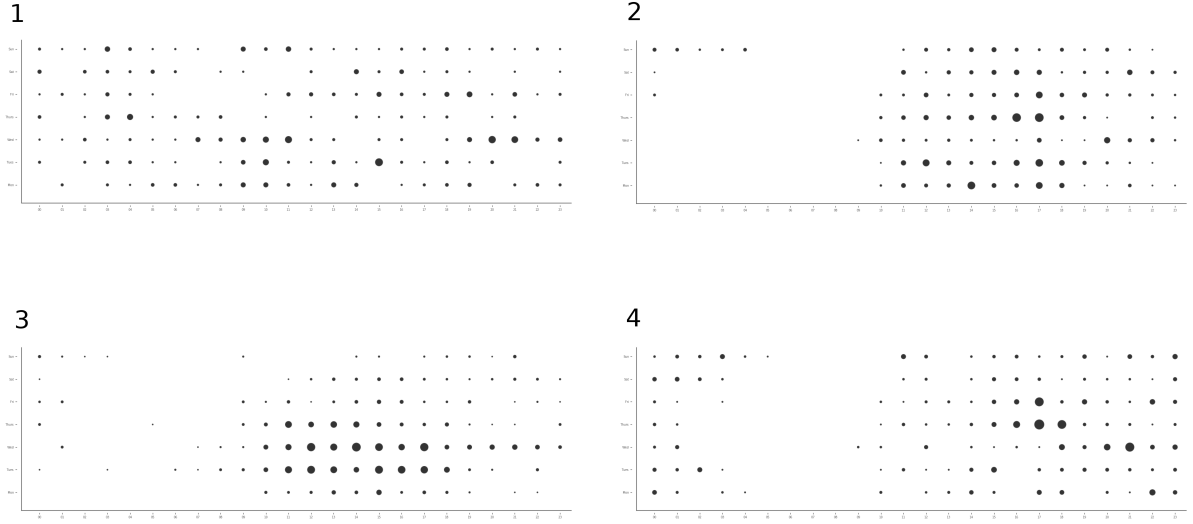


Figure 4.2: Punchcard of an user without an irregular sleep rhythm.

4.2.2 Employee or Open-Source Contributor

To determine whether a punchcard could be used to distinguish between an employee or an open-source contributor, the results of the clustering described in Section 3.2.1 has been utilized. For this approach two assumptions have been made. An usual employee works between Monday and Friday during the day and only as an exception at the weekend. An open-source developer works outside of the usual work shifts, which means early and late during weekdays and at the weekend.

For each assumption two representative clusters have been chosen and ten random persons have been selected for each cluster. The manual verification is conducted by checking if the contributor mainly contributes to repositories which belong to the registered employee. In case no employee exists, it is examined whether the contributor pushes to their own or open-source projects or rather to the repositories of a specific company.

4,156 contributions in the last year

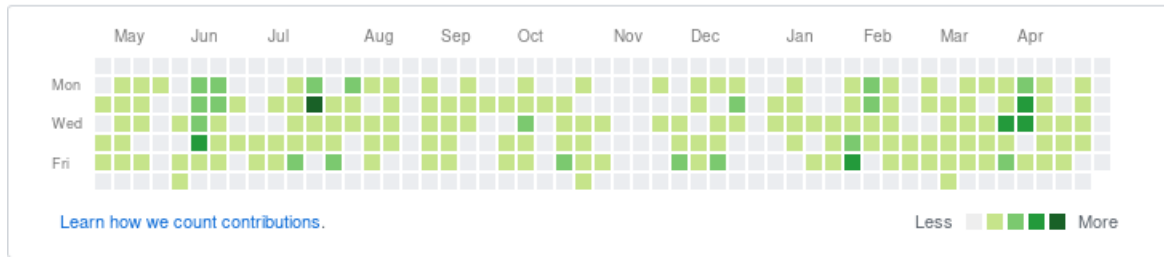


Figure 4.3: Github contribution overview a Google developer with the Github nickname alxhub.

It provides a good overview of the usual weekday work pattern over the last year and allows to quickly inspect the repositories a contributor committed to at a specific date.

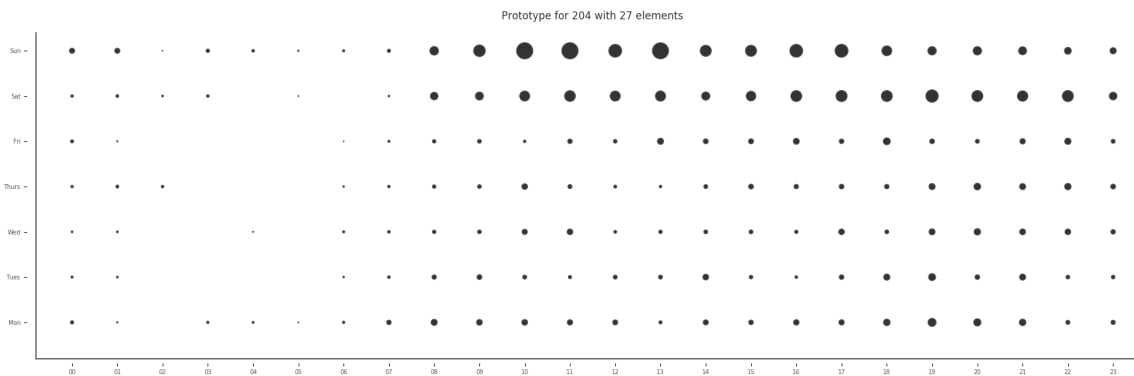


Figure 4.4: Punchcard of an example from an affinity propagation cluster with a weekend tendency.

The representatives for the usual five-day week commit behaviour were surprisingly accurate. About 93% of considered contributors were mainly working on projects of their companies, with occasional commits to their open source projects. The remaining 7% were either open-source developer with a very consistent commit behaviour, or it could not be determined if they work for a company.

The representatives for the leisure time commit behaviour are mostly correct as well. About 88% of considered contributors were irregularly contributing to either work unrelated open-source projects or their own projects. Approximately the half of those

contributors specified their employee and didn't work on their employees' projects. The other half worked either on their own projects or on open-source projects on an irregular basis. The remaining 12% were either contributors working and committing to their employee's projects, but also to their own and open source projects, and employees with an untypical commit behaviour.

This analysis shows quite well, that there is a correlation between the assumed patterns and the github commit behaviour or the employment status of contributor. Unfortunately, the evaluation process for these results is very time consuming and thereby only a relatively small sample group has been chosen. As it is not trivial to link the employee of an contributor to all their funded projects, all verification needed to be conducted manually.

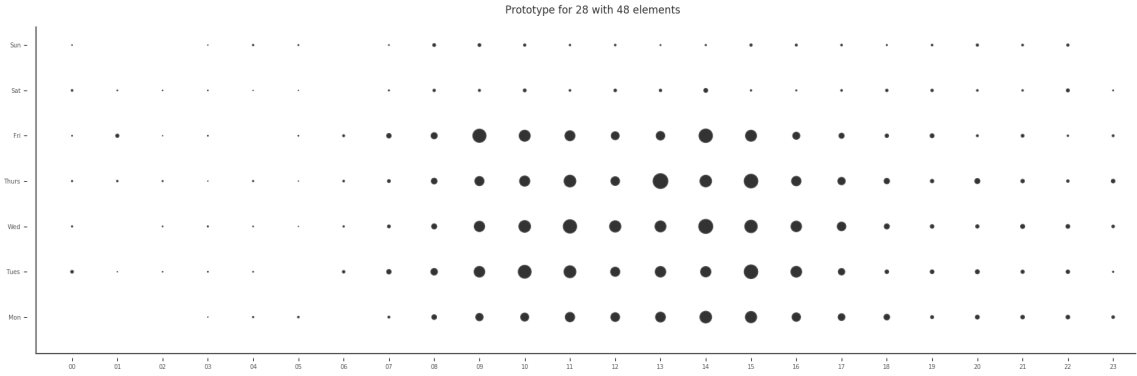


Figure 4.5: Punchcard of an example from an affinity propagation cluster with normal work shifts.

4.2.3 Bot detection

Another possible attack was the detection of automatically committing programs, so called *bots*. This algorithm simply detected centroids with an extremely equally distributed pattern or patterns with a spike at a specific hour. After a manual revision of the by the algorithm detected clusters, it became apparent, that only a small subset of those clusters actually contained bots. And even if the cluster contained bots, there were usually only one or two of them from a much larger pool of cluster members.

Detection of bots in the outliers, which were not assigned to any cluster, did not seem to be promising as well. Manual revision of over 50 possible candidates lead to not a single

bot. After reviewing these results, I decided, that there is currently no viable approach for my data to this problem.

4.2.4 Fingerprinting

Another possible attack was to fingerprint a contributor, by analysing their commit behaviour and thereby creating a unique identifier.

This attack soon proved to be unfeasible, as the pattern of a contributor significantly differs from month to month. There even are significant changes, if the interval of a year is considered and the shift is only by a month. The commit behaviour of people seem to be too inconsistent for a fingerprint.

4.3 Geolocation

First of all, it needs to be clarified, that parts of this attack only works under specific circumstances. Git commit timestamps identical to the current local time of the underlying OS. If one wants to show the travel path of a target, the target's OS needs to automatically adjust the UTC accordingly to the current geolocation of the device.

This feature is available for the newer versions of big OSs, such as Windows ¹ and Mac, but it is not necessarily enabled by default. It is also available for Linux, for instance with the *tzupdate* package ², but it needs to be installed and activated manually.

¹Ivan Jenic, 'Your Time Zone Can Now Switch Automatically in Windows 10', windowsreport.com, <https://windowsreport.com/time-zone-automatic-switch-windows-10> (accessed, 24.04.2018)

²'Set the system timezone based on IP geolocation', github.com, <https://github.com/cdown/tzupdate> (accessed, 24.04.2018)

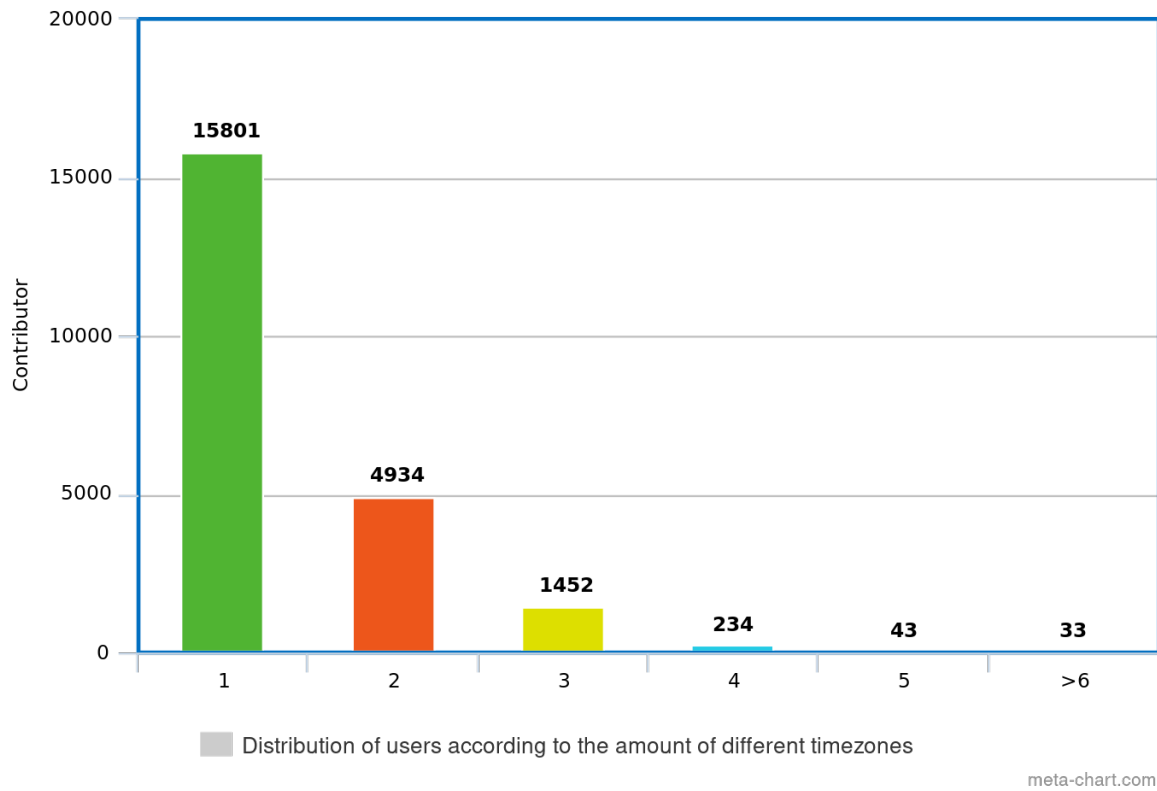


Figure 4.6: Distribution of users according to the amount of different timezones.

Figure 4.6 shows the amount of contributors in relation to the amount of detected timezones. On about 70% of considered contributors only a single timezone has been detected, looking at the last year. These 70% do either not commit, when they travel, their OS does not synchronize the timezone accordingly to their location or they simply do not traveled in the last two years.

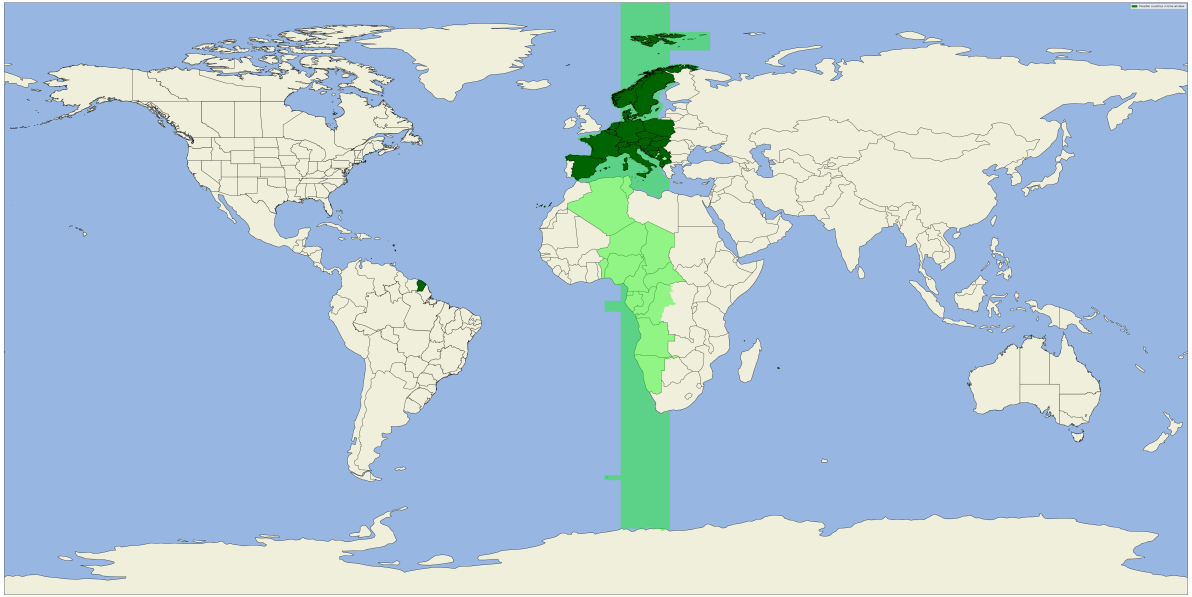


Figure 4.7: Home location analysis of the author.

In Figure 4.7 the visualized home location analysis of the author can be seen. Regions marked in dark green are regions, in which the contributor is likely to live. The light green region represents the timezone of the home location. As you can see in Figure 4.7 the country French Guiana is also marked as a possible home location. This problem occurs due to the several conversions between country names and codes, which were necessary as stated ???. This misassignment only happens during the visualization of the results and thereby doesn't affect the results of the analysis.

To evaluate the overall precision of the geolocation results, the correctness of the determined home location is checked. Github allows users to specify a string for their current home, which is collected during the aggregation process. Unfortunately there are no conventions on how this string has to look like. Initially I tried to pass these strings to OpenStreetMap, but this resulted in too many wrongly assigned locations. The data provided by the users was obviously too arbitrary and full of mistakes for the OpenStreetMap API to handle.

As a result I decided to manually choose a subset of locations by looking for distinct identifiers in the location strings. For instance, every home location of contributor, which contained *Germany* or *Deutschland* in their location string, should be in the UTC +1 timezone and contain the specific timezone `Europe/Berlin`. I created 17 such rules and was thus able to validate the home location of about 2000 contributors. The assignment of the contributors home location was correct in about 76% of the

considered contributors. It needs to be noted, that the accuracy of this result is quite certainly deteriorated by location strings, which contain ambiguous information, and maybe by contributors, which moved to another country in the last year, as we cannot detect those for sure.

Nevertheless an accuracy of 76% clearly shows, that it is possible to narrow down the location of a contributor to a timezone and even to a subset of countries, by simply looking at the their git commit timestamps.

CHAPTER 5

Conclusion and Outlook

The study set out to determine how feasible and precise data mining attacks on simple git metadata. All three performed attacks, lead to promising results and showed potential for actual usage in a professional environment.

However, it must be noted that in all attacks only a small amount of the available data was used. Simply using the git commit timestamps allowed us to perform analyses like narrowing down the location of a contributor. The possible applications for the remaining data, such as actual changes in code, references of contributors between repositories or commit messages, are extensive.

If one would add additional data from Github, such as followers, stars or information from their issues system, the results could further become substantially more accurate. Developers around the world provide metadata about themselves on a daily basis, probably without knowing how much they are actually exposing. To prevent the unauthorized usage and abuse of this data, we need to create counter-measures or prevent exposing this data in the first place.

Luckily the European Union (EU) set an example by enforcing the General Data Protection Regulation (GDPR), which is a regulation that strictly rules the handling of any user data. But there are still many countries left in the world, that do not have such strict rules and that might need ways to protect their privacy from being invaded.

While *Gitalizer* is a great foundation for data aggregation and performing rather simple analyses, there is a necessity for more detailed research with better sources for ground truth and more manpower. *Gitalizer* is a quite complex program, but it is well documented and should allow other people to easily jump into using it. Many of the attacks mentioned in Section 1.5 are not implemented, as they did not fit in the scope of this thesis. Implementing those could be the topic of another bachelor thesis or for a subsequent master thesis.

Furthermore it would be interesting to explore the possibilities of counter measures such as obfuscating git commit timestamps.

List of Figures

2.1	A Git commit history tree.	8
2.2	Simplified Github relationships in Crow's foot notation.	13
2.3	Gitalizer database relationships.	17
2.4	The CPU load of the aggregation server during optimization.	20
2.5	Gitalizer database relationships.	21
3.1	The work time analysis of the author.	24
3.2	Punchcard of the author.	27
3.3	Super cluster centroid found by mean-shift clustering.	28
3.4	Extreme outlier centroid found by mean-shift clustering.	29
4.1	The miss-out analysis of several employees.	37
4.2	Punchcard of an user without an irregular sleep rhythm.	38
4.3	Github contribution overview a Google developer with the Github nickname alxhub. . .	39
4.4	Punchcard of an example from an affinity propagation cluster with a weekend tendency. .	39
4.5	Punchcard of an example from an affinity propagation cluster with normal work shifts. .	40
4.6	Distribution of users according to the amount of different timezones.	42
4.7	Home location analysis of the author.	43

List of Listings

1	A tree file example.	9
2	A commit file example.	10
3	Miss-out analysis algorithm.	25
4	Data preprocessing	27
5	Algorithm used to detect changes in the target's location by analysing the UTC offsets of git commit timestamps	33
6	Methology used to determine the main location of a target, based on the information gained from the function in Listing 5	35

List of Tables

Bibliography

- [1] Scott Chacon and Ben Straub. *Pro Git, 2nd Edition*. Apress, 2014. ISBN: 978-1484200773.
- [2] Yizong Cheng. “Mean shift, mode seeking, and clustering”. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 17.8 (Aug. 1995), pp. 790–799. ISSN: 0162-8828. DOI: 10.1109/34.400568.
- [3] Brendan J. Frey and Delbert Dueck. “Clustering by Passing Messages Between Data Points”. In: *Science* 315.5814 (2007), pp. 972–976. ISSN: 0036-8075. DOI: 10.1126/science.1136800. eprint: <http://science.sciencemag.org/content/315/5814/972.full.pdf>. URL: <http://science.sciencemag.org/content/315/5814/972>.
- [4] GitLab. *GitLab Help*. 2017. URL: <https://gitlab.com/help> (visited on Jan. 22, 2018).
- [5] Google. *Version Control Systems Popularity in 2016*. 2017. URL: <https://octoverse.github.com> (visited on Jan. 22, 2018).
- [6] Google. *Version Control Systems Popularity in 2016*. 2017. URL: <https://security.googleblog.com/2017/02/announcing-first-sha1-collision.html> (visited on Dec. 16, 2016).
- [7] Georgios Gousios et al. “Lean GHTorrent: GitHub Data on Demand”. In: *Proceedings of the 11th Working Conference on Mining Software Repositories*. MSR 2014. Hyderabad, India: ACM, 2014, pp. 384–387. ISBN: 978-1-4503-2863-0. DOI: 10.1145/2597073.2597126. URL: <http://doi.acm.org/10.1145/2597073.2597126>.
- [8] Heinz-Peter Höller and Peter Wedde. *Die Vermessung Der Belegschaft*. 2018. URL: https://www.boeckler.de/112779_112796.htm (visited on Apr. 28, 2018).
- [9] David Kriesel. *SpiegelMining: Wer, wann, was, mit wem? Das soziale Netz der SpiegelOnline-Redakteure*. 2016. URL: http://www.dkriesel.com/blog/2016/0814_spiegelmining_soziales_netz_redakteure (visited on Apr. 28, 2018).
- [10] Andrew McAfee and Erik Brynjolfsson. “Big Data: The Management Revolution”. In: (2012). URL: <https://hbr.org/2012/10/big-data-the-management-revolution> (visited on May 5, 2018).
- [11] Cade Metz. *Version Control Systems Popularity in 2016*. 2015. URL: <https://www.wired.com/2015/03/github-conquered-google-microsoft-everyone-else/> (visited on Apr. 20, 2018).
- [12] rhodecgcode.com. *Version Control Systems Popularity in 2016*. 2016. URL: <https://rhodecode.com/insights/version-control-systems-2016> (visited on Jan. 22, 2018).
- [13] Rachel Emma Silverman. “Bosses Tap Outside Firms to Predict Which Workers Might Get Sick”. In: (2016). URL: <https://www.wsj.com/articles/bosses-harness-big-data-to-predict-which-workers-might-get-sick-1455664940?mod=e2tw> (visited on May 5, 2018).

Eidesstattliche Erklärung

„Hiermit versichere ich an Eides statt, dass ich die vorliegende Arbeit im Studiengang Informatik selbstständig verfasst und keine anderen als die angegebenen Hilfsmittel – insbesondere keine im Quellenverzeichnis nicht benannten Internet-Quellen – benutzt habe. Alle Stellen, die wörtlich oder sinngemäß aus Veröffentlichungen entnommen wurden, sind als solche kenntlich gemacht. Ich versichere weiterhin, dass ich die Arbeit vorher nicht in einem anderen Prüfungsverfahren eingereicht habe und die eingereichte schriftliche Fassung der auf dem elektronischen Speichermedium entspricht.“

Ort, Datum

Unterschrift