

Entwicklung einer Anwendung zur automatisierten Beschaffung von personenbezogenen Daten im Internet und deren Integration in Phishing-Mails

Bachelorarbeit

Wintersemester 2018/2019

im Studiengang Angewandte Informatik

an der Hochschule Ravensburg - Weingarten

von

Marco Lang Matr.-Nr.: 27416

Abgabedatum : 14. April 2019

Erklärung

Hiermit erkläre ich, dass ich die vorliegende Arbeit mit dem Titel

**Entwicklung einer Anwendung zur automatisierten Beschaffung von
personenbezogenen Daten im Internet und deren Integration in
Phishing-Mails**

selbstständig angefertigt, nicht anderweitig zu Prüfungszwecken vorgelegt, keine anderen als die angegebenen Hilfsmittel benutzt und wörtliche sowie sinngemäße Zitate als solche gekennzeichnet habe.

Weingarten, 14. April 2019

Autor Name

Inhaltsverzeichnis

Kurzfassung	IV
Danksagung	V
1 Einleitung	1
1.1 Motivation	1
1.2 Zielsetzung und Forschungsfragen	2
1.3 Eigene Leistung	3
1.4 Aufbau der Arbeit	3
2 Grundlagen	4
2.1 Personenbezogene Daten	4
2.2 Social Engineering	4
2.2.1 Phishing	5
2.2.2 Spear-Phishing	6
2.3 Open Source Intelligence	7
2.3.1 Definition OSINT	7
2.3.2 Web Crawler	7
2.3.3 Web Scraper	8
3 Problembeschreibung	10
4 Ethische und rechtliche Betrachtung	11
5 Anforderungsanalyse	12
5.1 Anforderung an das OSINT einer ausgewählten Person	12
5.2 Anforderung an die Generierung einer Phishing-Mail	13
5.2.1 Anforderung an die Generierung der E-Mail-Adressen	13
5.2.2 Anforderung an die Erstellung der E-Mail-Texte	13
6 Lösungsideen	14
6.1 Methoden für das OSINT einer ausgewählten Person	14
6.1.1 Verwendung von OSINT-Tools	14

6.1.2	Algorithmus für das OSINT entwickeln	14
6.2	Konzept für die Erstellung einer Phishing-Mail	15
6.2.1	Methoden zur Generierung der E-Mail-Adresse	15
6.2.2	Methoden zur Generierung des E-Mail-Textes	16
7	Bewertung der Lösungsideen anhand der Anforderung	17
7.1	Bewertung der OSINT-Methoden für eine ausgewählte Person	17
7.2	Bewertung der Methoden zur Erstellung einer Phishing-Mail	18
8	OSINT einer ausgewählten Person	20
8.1	Auswahl der Programmiersprache	20
8.2	Methoden zur Suche nach einer Person im Internet	21
8.2.1	Personensuche mit Hilfe einer Suchmaschine	21
8.2.2	Personensuche auf festgelegten Webseiten	22
8.3	Bewertung der Methoden zur Personensuche	22
8.3.1	Auswahl der Suchmaschine	22
8.4	Implementierung der Personensuche mit Hilfe der Google-Suchmaschine im Internet	23
8.4.1	Eingabe der bekannten Daten	23
8.4.2	Generierung der Google-Such-URLs	24
8.4.3	Mit welcher Bibliothek werden Serveranfragen umgesetzt?	28
8.4.4	Web Crawler erstellen	29
8.5	Methoden zum Erkennen von wichtigen Informationen auf einer Webseite	34
8.5.1	RAKE	35
8.5.2	Automatic Keyword Extraction mit NLP	38
8.6	Bewertung der Methoden zum Herausfiltern von wichtigen Informationen auf einer Webseite	38
8.7	Implementierung der Methoden zum Herausfiltern von wichtigen Informationen auf einer Webseite	39
8.7.1	Text formatieren	39
8.7.2	Erstellung der Wortsammlungen	39
8.7.3	Automatic Keyword Extraction mit NLP	40
8.7.4	Suche nach dem Geburtsjahr der Zielperson	41
8.7.5	E-Mail-Adressen erkennen und herauslesen	42
8.7.6	Auswahl der gewonnenen Information	43
8.8	Methoden zum Erkennen einer Person	44
8.8.1	Identifikationsschlüssel verwenden	45
8.8.2	Kontaktanalyse	46
8.9	Bewertung der Methoden zur Personenidentifizierung	46
8.10	Implementierung der Personenidentifizierung	47
8.10.1	Identifikationsschlüssel verwenden	47

8.10.2	Kontaktanalyse	47
8.11	Speicherung der gewonnenen Daten	51
8.11.1	Implementierung der Personenklasse	51
9	Generierung der Phishing-E-Mail	53
9.1	Implementierung der Methode zur Generierung der E-Mail-Adressen . . .	53
9.1.1	Funktion des eigenen Algorithmus	53
9.2	Implementierung der E-Mail-Muster	56
9.2.1	Kategorien erstellen	56
9.3	Versenden einer Phishing-E-Mail	61
10	Evaluation der Implementation	63
11	Schlussbemerkungen und Ausblick	64
11.1	Wie kann eine Person weiter identifiziert werden?	64
11.1.1	Zeitraumen wird mit Beachtet	64
11.1.2	Zeitraum beachten	65
11.2	Adressgenerierung	65
11.2.1	Wenn Firma bekannt	65
11.3	Keyword Extraction mit Hilfe von Machine Learning	65
11.4	Wie werden Wortsammlungen am effektivsten verglichen?	65
11.5	Email-adressen	66
11.6	Absender-Adresse	66
11.7	Validität der generierten Mail-Adressen prüfen	66
11.7.1	Methoden zum Prüfen der Validität	66
11.7.2	Bewertung: Validität Prüfen	67
A	Ein Kapitel des Anhangs	68
	Literatur	69
	Stichwortverzeichnis	72

Kurzfassung

Es wird gezeigt, wie eine automatisierte Suche nach personenbezogenen Daten im Internet aussehen kann und wie diese Daten für einen Phishing-Mail-Angriff verwendet werden können.

Danksagung

1 Einleitung

1.1 Motivation

70% der Internetnutzer sehen sich laut einer Umfrage durch das Risiko einer missbräuchlichen Verwendung ihrer Daten nach einem Hack nicht gefährdet [Ang18]

Das Ergebnis dieser Umfrage spricht für die Behauptung, dass viele Personen Informationen über die eigenen Person im Internet preis geben, da keine Ängste vorhanden sind. Doch diese Informationspreisgabe kann in den falschen Händen schwerwiegende Folgen haben. So kann beispielsweise bei einem Phishing-Mail-Angriff diese Art von Information genutzt werden, um ein potentielles Opfer zu täuschen oder zu manipulieren. Ein Beispiel dafür, sind die gefälschten DSGVO-E-Mails, bei denen der Angreifer das Opfer durch scheinbar echte Mails der Sparkasse täuscht. Dabei wird die Zielpersonen persönlich mit ihrem Namen angesprochen, wodurch die Mail an Glaubwürdigkeit gewinnt. [NW18]

Solch ein Angriff benötigt allerdings im Voraus eine ausführliche Recherche über das Opfer. Als Informationsquelle für die Recherche dienen beliebig viele Medien. Doch in der heutigen Zeit ist das Internet die meistgenutzte Informationsquelle für Menschen und birgt dadurch Gefahren für jeden einzelnen Internetnutzer, der personenbezogene Daten im Internet teilt. [All18] Diese Gefahr wird unter anderem durch die Entwicklung von kostenlosen OSINT-Tools erhöht. Diese Tools sammeln Informationen über Opfer von öffentlichen und frei zugänglichen Medien. Dadurch wird die Recherche im Internet nach persönlichen Informationen deutlich einfacher. Das hat zu Folge, dass jeder Internetnutzer ohne großen Aufwand OSINT im Internet betreiben kann.

1.2 Zielsetzung und Forschungsfragen

Ziel ist es eine OSINT-Anwendung zu entwickeln, welche automatisiert nach personenbezogenen Daten im Internet sucht. Die gewonnenen Daten werden anschließend in eine Phishing-Mail integriert. Dabei soll der Fokus auf der automatisierten Informationsbeschaffung liegen.

Unter anderem sollen Antworten auf die folgenden Fragen gefunden werden. Mit welchem Aufwand ist ein automatisierte Spear-Phishing-Mail-Angriff verbunden? Ist es möglich ein Personenprofil zu erstellen, bei dem ausschließlich korrekte Informationen vorhanden sind? Wie glaubwürdig sind automatisierte Phishing-E-Mails mit integrierten personenbezogenen Daten?

Ziel 1 *Informationen zu einer ausgewählten Person im Internet suchen.*

Diese Suchfunktion beinhaltet die Suche nach Informationen einer bestimmten Person. Dadurch können bereits bekannte Daten über die Person angegeben und somit die Suche verfeinert beziehungsweise verbessert werden. Hierbei besteht die Herausforderung zu erkennen, wann es sich um eine Information der gesuchten Person handelt.

Ziel 2 *E-Mail-Adressen finden oder aus den gewonnenen Daten generieren.*

Wenn eine E-Mail-Adresse zu einer gesuchten Person nicht gefunden werden kann, soll diese mit Hilfe der gewonnenen Daten generiert werden. Durch die Zusammensetzung von Vorname, Name und Geburtsjahr können die möglichen E-Mail-Adressen einer Zielperson erzeugt werden. Des Weiteren kann die Institution der gesuchten Person, falls diese bekannt ist, mit in den Generierungsprozess einfließen.

Ziel 3 *Phishing-Mail erzeugen.*

Mit den gewonnenen Informationen soll eine Phishing-E-Mail erzeugt werden. Dabei wird der Inhalt dieser Mail, abhängig von den gewonnenen Informationen erstellt. Dabei ist das Ziel, dass eine glaubhafte und sinnvolle Spear-Phishing-Mail generiert und versendet werden kann.

1.3 Eigene Leistung

In dieser Arbeit wird eine Anwendung erstellt, welche personenbezogene Daten zu einer gesuchten Person automatisiert aus dem Internet heraus sucht. Die gewonnenen Daten werden in einem potentiellen Opferprofil gespeichert und anschließend in eine personalisierte Phishing-E-Mail integriert. Für einen höheren Erfolg der Phishing-Mails werden Methoden für die Generierung des Mailtextes herausgearbeitet und realisiert.

Damit ein kompletter Ablauf eines Phishing-Mail-Angriffs simuliert werden kann, wird zu jeder Personensuche eine passende E-Mail-Adresse benötigt. Allerdings kann nicht bei jeder Suche eine korrekte E-Mail gefunden werden. Aus diesem Grund wird zusätzlich ein Algorithmus entwickelt, der im Fall, dass keine E-Mail-Adresse zu der Zielperson gefunden wurde, ein Pool aus möglichen Mail-Adressen mit Hilfe der gefundenen Informationen generiert.

1.4 Aufbau der Arbeit

Die Arbeit gliedert sich in einen theoretischen und praktischen Teil auf. Die Theorie beginnt im zweiten Kapitel und beschreibt die Grundlagen im Bereich von personenbezogenen Daten, Social Engineering und der Informationsbeschaffung im Internet. In Kapitel 3 wird das Problem aufgezeigt, auf welches in dieser Arbeit eingegangen wird. Darauf folgt die ethische und rechtliche Betrachtung in Kapitel 4. Die Anforderungsanalyse 5 beschreibt das nächste Kapitel, in welchem die Anforderungen und Prioritäten der Arbeit festgelegt werden. Darauf folgen die Lösungsvorschläge im Kapitel 6 und die Auswahl der Lösung anhand der Anforderungen im Kapitel 7. Anschließend wird bei der Umsetzung auf den Praktischen Teil eingegangen. Dieser unterteilt sich in die Themen OSINT einer ausgewählten Person 8 und die Generierung einer Phishing-Mail 9. Am Ende dieser Arbeit befindet sich die Evaluation der Implementation in Kapitel 10 sowie die Schlussbemerkung und der Ausblick in Kapitel 11.

2 Grundlagen

2.1 Personenbezogene Daten

Laut der DSGVO sind **personenbezogene Daten**, alle Informationen, die sich auf eine identifizierbare Person beziehen. Als identifizierbar wird eine natürliche Person angesehen, die mittels einem oder mehreren Merkmalen direkt oder indirekt identifiziert werden kann. Mögliche Kennungen für die Unterscheidung der Merkmale sind der Name, eine Kennnummer, Standortdaten, eine Online-Kennung, et cetera von der Person. Dabei dienen diese Kennungen als Ausdruck der physischen, physiologischen, genetischen, psychischen, wirtschaftlichen, kulturellen oder sozialen Identitäten dieser natürlichen Person. [DSG]

2.2 Social Engineering

Die Definition von Social Engineering ist nicht eindeutig, da es sehr verschiedene Ansichten davon gibt. Jedoch ist der Grundgedanke von Social Engineering, eine Zielperson so zu manipulieren, damit sie für den Angreifer bessere Entscheidung trifft. [Had11]

Kevin D. Mitnick definiert Social Engineering wie folgt:

“Social Engineering uses influence and persuasion to deceive people by convincing them that the social engineer is someone he is not, or by manipulation. As a result, the social engineer is able to take advantage of people to obtain information with or without the use of technology“ [Mit01]

Social Engineering wird Menschen von Geburt an beigebracht und begegnet einem beinahe jeden Tag. Schon ein Baby muss wissen wie es die Eltern manipulieren kann, damit es Dinge

wie Essen, Zuneigung, oder ähnliches bekommt. Darüber hinaus ist Social Engineering in vielen Berufen ein täglicher Bestandteil.

Im Bereich der Informationssicherheit, wird von Social Engineering gesprochen, wenn Angreifer durch die Manipulierung und Täuschung von Menschen vertrauliche Informationen oder Zugänge zu Systemen bekommen. Die bekanntesten Angriffsmethoden sind Phishing, Pretexting, Baiting und Quid Pro Quo. Bei dieser Arbeit wird hauptsächlich auf das Thema E-Mail-Phishing eingegangen.

Der Aufbau eines Social Engineering-Angriffes ist definiert in mehrere Phasen. Das wohl bekannteste Modell für einen Social Engineering-Angriffszyklus ist in dem Buch von Kevin D. Mitnicks [Mit01] definiert. Dieser Zyklus besteht aus den 4 Phasen **Research**, **Developing rapport and trust**, **Exploiting trust** und **Utilize information**.

In der **Research-Phase** geht es um die Informationsbeschaffung. Bei dieser Phase will der Angreifer möglichst viele Informationen über das Ziel herausfinden. Die **Developing Rapport and Trust-Phase** beschreibt den Kontaktaufbau zum Ziel, da wenn das Opfer dem Angreifer vertraut, hat dieser ein leichteres Spiel in den kommenden Phasen. Das nun erzeugte Vertrauen wird in der **Exploitation Trust-Phase** ausgenutzt. Hier will der Angreifer die eigentlich Information vom Opfer herausfinden. Dies geschieht einerseits durch bestimmtes Nachfragen oder durch Manipulation. **Utilize Information** ist die letzte Phase. Dort wird die gewonnene Information genutzt um das eigentliche Ziel des Angreifers zu erreichen.

Grundsätzlich werden bei einem Social Engineering Angriff menschliche Wünsche, Ängste und verbreitete Verhaltensmuster verwendet um ein Opfer zu manipulieren. [uDsine15]

2.2.1 Phishing

Das Wort Phishing wird von dem Wort “fishing“ abgeleitet, da die Angreifer nach Informationen fischen. Das “Ph“ kommt von “sophisticated“ und meint damit, dass die Angreifer ausgeklügelte Techniken verwenden um an Informationen heranzukommen. [Jam05]

Die wohl bekannteste Angriffsmethode von Phishing ist das E-Mail-Phishing. Bei diesem Verfahren, versendet ein Angreifer meist eine gefälschte E-Mail, um ein Opfer zu täuschen

und dadurch sein Ziel zu erreichen. Die sogenannten Phishing-Mails enthalten meist eine Aufforderung einen Link zu öffnen und sehen täuschend echt aus.

Ein reales Beispiel könnte sein, dass der Angreifer eine gefälschte E-Mail von Amazon an das Opfer versendet und es dabei auffordert, einen Link in der Mail zu öffnen. Nachdem die Zielperson auf den Link geklickt hat, muss Sie sich anmelden. Hier könnte der Angreifer ein täuschend echtes Anmeldeformular erstellt haben, um die Anmeldedaten der Zielperson zu bekommen. Sobald die Anmeldedaten eingegeben wurden, könnte eine Fehlermeldung erscheinen, die einen Authentifizierungsfehler beinhaltet und das Opfer auffordert sich erneut anzumelden. Jedoch wird während diesem Prozess das originale Anmeldeformular geladen und das Opfer kann sich korrekt bei der entsprechenden Webseite anmelden.

Dieser Verfahren ermöglicht Angreifern die Anmeldedaten von einer Zielperson ohne großen Aufwand zu beschaffen. Allerdings benötigt der Angreifer für diese Methode nicht nur Social Engineering sondern auch technische Fähigkeiten. [CH15]

2.2.2 Spear-Phishing

Das Spear-Phishing ist eine erweiterte Methode des herkömmlichen E-Mail-Phishings. Hierbei wird anstatt das Versenden etlicher Phishing-Mails an unbekannte Opfer, eine gezielte Mail an eine ausgewählte Person versendet. [Fir]

Bei dieser Form von E-Mail-Phishing spielt die Opferauswahl und die Informationsbeschaffung eine sehr große Rolle, da diese Information später für personalisierte E-Mails oder vorgetäuschte Identitäten verwendet werden können. Durch diese Art von Täuschung kann ein Opfer dazu bewegt werden auf einen Link zu klicken und dadurch eine Schadsoftware herunterzuladen. [Fir]

Der Aufwand für die Informationsbeschaffung wird oft in Kauf genommen, da der Erfolg bei dieser Methode vielversprechender ist als beim herkömmlichen E-Mail-Phishing.

91% der Advanced Persistent Threat (APT) Angriffe auf Firmen beginnen mit einer Spear-Phishing-E-Mail. Die Schadsoftware wird meistens als Remote Access Trojans (RATs) in einer Zip-Datei überliefert. [Cal13]

2.3 Open Source Intelligence

2.3.1 Definition OSINT

Open Source Intelligence (OSINT) ist definiert in eine Intelligenz, welche aus öffentlich zugänglichen Informationen gewonnen wird. Allerdings kann sich die Bedeutung fallspezifisch ändern. So bedeutet OSINT für die CIA die Informationsgewinnung aus ausländischen Nachrichtensendungen. Doch für die meisten Menschen bedeutet OSINT die Gewinnung eines öffentlichen Inhalts aus dem Internet. [Baz18]

Unter Open Source wird die öffentlich zugängliche Information, die in gedruckter oder elektronischer Form vorliegt, bezeichnet. [Ste96] Eine Verbindung mit dem Begriff Open-Source-Software besteht nicht.

2.3.2 Web Crawler

Web Crawler, auch Robot oder Spider genannt, sind Computerprogramme, die mit Hilfe der Hypertextstruktur das Internet durchlaufen. [The01] Dabei können sie in einen **internen** und **externen Web Crawler** unterschieden werden. Der interne Web Crawler durchsucht ausschließliche interne Seiten einer Webseite und der externe Web Crawler durchsucht unbekannte Webseiten im ganzen Netz. [SG12]

In anderen Worten besteht die Funktionsweise darin, dass in den meisten Fällen ein automatisiertes Programm, Web Crawler, erstellt wird. Dieser lädt Webinhalte herunter und durchsucht den Inhalt nach Hyperlinks. Den gefundenen Links wird gefolgt, um neue Webseiten mit weiteren Links zu laden. So handelt sich ein Web Crawler von Link zu Link durch das Internet. [Mit15] Dieser Ablauf ist in dem Bild 2.1 noch einmal verdeutlicht.

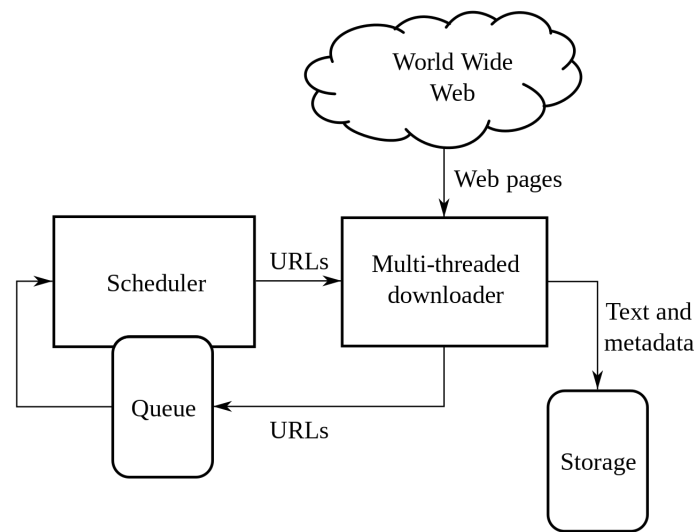


Bild 2.1: Architektur eines Web Crawlers

2.3.3 Web Scraper

In der Theorie bedeutet *web scraping* die Informationsbeschaffung im Internet mit unterschiedlichsten Mitteln. [Mit15]

Meist wird dies mit einem automatisierten Programm realisiert, welches Daten von einem Webserver anfragt, entgegen nimmt, analysiert und auswertet. In der Praxis gibt es ein großes Feld von Programmiertechniken und Einsatzmöglichkeiten. Mit Hilfe eines Web Scrapers ist es möglich, große Datenmengen zu erfassen und zu verarbeiten. [Mit15]

Natural Language Processing

Natural Language Processing (NLP) beschreibt eine Technologie, für die Kommunikation zwischen Mensch und Computer. Mit dem Ziel, dass ein Computer die natürliche Sprache verstehen und verarbeiten kann. Dafür werden verschiedenste Methoden aus der Sprach- und Computerwissenschaft sowie aus der künstliche Intelligenz verwendet. Unter anderem hat eine NLP-Anwendung die Aufgabe von **Stemming**. [Lit16]

Stemming ist eine Methode der Wortstandardisierung, bei der verwandte Wörter auf ihrer Stammform reduziert werden. Dabei wird bei dem Rechengang auf den Stamm und die

Semantik eines Wortes geachtet. Aus diesem Grund fällt der Name Stammformreduktion öfters in Verbindung mit Stemming. [EAD09] Ein Beispiel hierfür wären die Worte “Wetter“ und “Wetten“, welche auf den Stamm “Wett“ reduziert werden könnten. [PH]

Die Verwendung von Stemming, kann bei der Schlüsselwortgenerierung von Texten sehr hilfreich sein, da die Anzahl der möglichen Schlüsselwörter reduziert werden können.

3 Problembeschreibung

Persönliche Daten sind im Internet oft frei zugänglich. Das heißt, dass unterschiedlichste Webseiten persönliche Information von Menschen öffentlich bereitstellen. Die bekanntesten Webseiten sind die Social Media Seiten wie Twitter, Facebook und Instagram. Allerdings wird auch auf anderen Webseiten personenbezogene Daten in großen Mengen bereitgestellt. Ein Beispiel dafür ist Berufsportal LinkedIn oder XING. Diese Art von Webseiten sind perfekte Informationsquellen für Phisher, da im Bereich von Social Engineering, diese Informationen oft genutzt werden um ein Opfer zu täuschen oder zu manipulieren.

Dass hier beschriebene Problem zeigt, dass der Zugang für persönliche Information durch das Internet für die Öffentlichkeit einfacher gemacht wird. Es soll mit einem kritisch Blick darauf gezeigt werden, wie diese Daten für einen böswilligen Social Engineering-Angriff missbraucht werden können.

4 Ethische und rechtliche Betrachtung

Das Sammeln von personenbezogenen Daten auf sozialen Netzwerken ist ethisch und rechtlich gesehen ein sehr sensibles Thema. Jedoch werden in dieser Arbeit ausschließlich die Daten verwendet, die öffentlich frei zugänglich sind. Das heißt, unter den Informationen befinden sich keine Passwörter oder Informationen die nicht an die Öffentlichkeit gehören. Des Weiteren ist der hier verwendete Crawler nicht stark genug, um die Leistung eines Servers von einem sozialen Netzwerk zu beeinflussen.

Mit diesem realen Experiment, soll die Privatsphäre der Benutzer geschützt werden, indem aufgezeigt wird, wozu veröffentlichte Daten über eine Person im negativen Sinn verwendet werden können. Genau aus diesem Grund ist es wichtig, dass das Experiment in der realen Welt durchgeführt wird.

5 Anforderungsanalyse

Die im Kapitel 1.2 definierten Ziele sollen mit den folgenden Anforderungen gewährleistet werden.

5.1 Anforderung an das OSINT einer ausgewählten Person

Bei dieser Informationsbeschaffung soll eine Suchfunktion entwickelt werden, welche Daten zu einer angegebenen Person im Internet sucht. Hierbei sollen so viele Daten wie möglich gefunden und gespeichert werden.

Die zu entwickelnde Anwendung soll für die Suche bekannte Daten wie Vorname, Nachname, Geburtsjahr, Wohnort, E-Mail-Adresse und Benutzernamen von Social Media Plattformen einlesen können. Die Eingabe kann mit Hilfe einer Konsole oder einer grafische Oberfläche realisiert werden.

Die Herausforderung besteht darin, zu erkennen wann es sich um die gesuchte Person handelt. Aus diesem Grund werden Methoden zur Identifizierung einer Person entwickelt und umgesetzt. Des Weiteren werden die herausgelesenen Daten analysiert und interpretiert. Dadurch sollen wichtige Informationen über die Person erkannt.

5.2 Anforderung an die Generierung einer Phishing-Mail

Die Phishing-Mails sollen automatisiert erstellt werden. Dafür wird vorausgesetzt, dass E-Mail-Adressen und E-Mail-Texte passend zu der gesuchten Person ebenfalls automatisiert erzeugt werden.

5.2.1 Anforderung an die Generierung der E-Mail-Adressen

Da nicht zu jeder Suche eine E-Mail-Adresse im Internet gefunden werden kann, muss die E-Mail-Adresse aus den vorhandenen Informationen generiert werden. Es ist möglich eine größere Anzahl von möglichen E-Mail-Adressen zu erzeugen. Durch den Pool an generierten E-Mail-Adressen soll die Wahrscheinlichkeit erhöht werden, dass die richtige E-Mail-Adresse dabei ist. Darüber hinaus können die Adresse validiert werden.

5.2.2 Anforderung an die Erstellung der E-Mail-Texte

Hierbei handelt es sich ausschließlich um das Erstellen potentieller Inhalte einer E-Mail. Zur Generierung der Texte sollen die gewonnen Informationen verwendet werden. Für jede Art von Opfer wird ein übereinstimmender E-Mail-Text generiert. Die Texte sollen mit den gefunden Daten Sinn ergeben und eine korrekte Grammatik beinhalten. Weiterführend können Social Engineering-Fähigkeiten genutzt werden um die Zielperson tatsächlich zu manipulieren und zu täuschen. Hierfür können beispielsweise Gefühle wie Freude und Angst ausgenützt oder gefälschte E-Mails von bekannten Firmen in Betracht gezogen werden.

6 Lösungsideen

In diesem Kapitel werden die Lösungsideen für die Umsetzung der im Kapitel 1.2 definierten Ziele beschrieben.

6.1 Methoden für das OSINT einer ausgewählten Person

6.1.1 Verwendung von OSINT-Tools

Die Personensuche wird durch die Verwendung kostenloser OSINT-Tools durchgeführt. Eine entsprechende Webseite die mehrere OSINT-Methoden bereit stellt, ist unter dem URL "<https://inteltechniques.com/index.html>" erreichbar. Sie stellt Methoden zur Suche nach E-Mail-Adressen, Benutzernamen, Social-Media-Profilen, und noch viele mehr zu Verfügung. Allerdings werden nicht nur selbstentwickelt OSINT-Methoden von Michael Bazzell bereitgestellt, sondern auch andere Webseiten mit weiteren OSINT-Tools vorgeschlagen.

6.1.2 Algorithmus für das OSINT entwickeln

Es wird ein Algorithmus für das OSINT entwickelt, der aus einem Web Crawler und Web Scraper besteht. Mit diesem ist es möglich eigenständig nach Information zu suchen. Hierfür wird eine Suchmaschine, wie die von Google, verwendet.

Die Suchergebnisse können mit Hilfe des Web Crawlers verfolgt werden. Anschließend wird

der Webseitentext, durch den Web Scraper, ausgelesen. Im letzten Schritt, wird der Text analysiert und interpretiert.

All diese Prozesse laufen unabhängig von den vorgeschlagenen Webseiten voll automatisiert ab.

6.2 Konzept für die Erstellung einer Phishing-Mail

Die Generierung einer realen Phishing-Mail benötigt eine korrekte E-Mail-Adresse der Zielperson. Darüber hinaus sollten die gewonnen Informationen in einem sinnvollen E-Mail-Text eingebunden werden. Die Generierung einer Phishing-Mail läuft voll automatisch ab. Das bedeutet, dass das Programm eigenständig die E-Mail-Adressen generiert und passende E-Mail-Muster auswählt.

6.2.1 Methoden zur Generierung der E-Mail-Adresse

Beim OSINT einer ausgewählten Person wird bereits nach E-Mail-Adressen der Zielperson gesucht. Dadurch kann eine bis jetzt unbekannte Anzahl von Adressen gefunden werden. Die Methoden zur Generierung einer E-Mail-Adresse muss dadurch nicht für jede Zielperson durchgeführt werden. Für den Fall, dass keine E-Mail-Adressen gefunden wurde, werden die folgenden Methoden vorgeschlagen.

Algorithmus entwickeln zum generieren

Es kann ein Algorithmus entwickelt werden, der mögliche E-Mail-Adressen aus den gewonnen Daten generiert. Dies ist durch die Kombination aus Vorname, Nachname, Geburtsjahr und den bekanntesten E-Mail-Providern realisierbar. Für den Fall, dass der Arbeitgeber der Zielperson bekannt ist, kann auf der Firmenwebseite nach E-Mail-Adressen gesucht werden. Dadurch ist es möglich die Domain einer Firmen-Mailadresse zu bestimmen und eine Anzahl möglicher Firmenadressen für die Zielperson zu generieren. Durch diese Methode wird eine Pool mit möglichen Mailadressen erstellt. Dabei muss jede einzelne E-Mail-Adresse auf Validität geprüft werden.

Automatisierbare OSINT-Tools verwenden

Für die Generierung der E-Mail-Adressen kann ein kostenloses OSINT-Tools von Michael Bazzel verwendet werden. Diese Tool ermöglicht es, die gewonnenen Informationen über eine Formular einzugeben. Anschließend werden draus mögliche E-Mail-Adressen generiert. Auch hier entsteht ein Adresspool, bei dem die E-Mail-Adressen auf Validität geprüft werden. Zu dem bringt das Tool eine weitere Funktion mit sich. Es wird automatisch nach Einträgen, der generierten E-Mail-Adressen, im Internet gesucht und angezeigt. [Baz]

6.2.2 Methoden zur Generierung des E-Mail-Textes

Muster für den E-Mail-Text erstellen

Die zu erstellenden E-Mail-Muster entsprechen hier kategorisierten Lückentexten. Abhängig von den gefundenen Daten, wird ein Lückentext ausgewählt und anschließend mit den Daten an den passenden Stellen ergänzt.

Die Lückentexte werden so kategorisiert, dass für jede gefundene Information ein passender Lückentext vorhanden ist. Eine denkbare Unterteilung wären die Kategorien Privat und Geschäftlich.

E-Mail-Text aus Fragmenten erzeugen

Bei dieser Methode besteht der E-Mail-Text aus zusammengesetzten Fragmenten. Dafür wird zu jeder gefundenen Information ein Fragment erstellt. Anschließend werden alle Fragmente zu einem Text zusammengefügt. Der Unterschied zur Methode 6.2.2 besteht darin, dass der E-Mail-Text dynamisch erzeugt wird. Das bedeutet, der endgültige Text ist nicht vorgeben. Er kann aus einer variierenden Anzahl von Fragmenten besteht. Diese Anzahl kann variieren, da sie abhängig von der gefundenen Information über die Zielperson ist.

7 Bewertung der Lösungsideen anhand der Anforderung

7.1 Bewertung der OSINT-Methoden für eine ausgewählte Person

Hierfür gibt es zwei verschiedene Methoden um OSINT zu betreiben. Die erste Lösungsidee beschreibt die Verwendung von einem öffentlich frei zugänglichen OSINT-Tool. Diese Tool bietet sehr viele Möglichkeiten um eine Person beziehungsweise Daten über eine Person zu finden. Allerdings ist es auf dieser Webseite nicht möglich ein zu suchendes Profil anzugeben, um eine Person zu finden. Die Suchen sind aufgeteilt in verschiedenste Daten wie Name, E-Mail, et cetera. Aus diesem Grund wird bei einer Suche ausschließlich nach dem Namen oder einer E-Mail gesucht. Dadurch ist das Suchergebnis am Ende kein vollständiges Personen-Profil, sondern lediglich Verweise auf weiterer Webseiten mit möglichen Einträgen. Dazu ist die Eingabemöglichkeiten der im Voraus bekannten, Daten begrenzt, da die Formulare nicht individuell erweiterbar sind.

Im Gegensatz zu diesem Tool, nutzt der eigenen Algorithmus alle im Vorfeld bekannten Daten für eine Suche. Des Weiteren kann die Laufzeit verbessert werden und bekannten Suchtechniken dieses Tools, mit Hilfe des Buches [Baz18] verwendet werden. Durch die eigene Anwendung wird die Suche beliebig erweiterbar programmiert. Dadurch kann jede Information zur Personensuche verwendet werden.

7.2 Bewertung der Methoden zur Erstellung einer Phishing-Mail

Generierung der E-Mail-Adressen

Bei der Verwendung eines bereits fertigen OSINT-Tools wird keine große Arbeit mehr benötigt, es ist ein komplettes System was funktioniert. Lediglich die Automatisierung muss entwickelt werden. Allerdings kann nicht jede individuelle Information für die Generierung genutzt werden. Dies ist für die zu entwickelnde Anwendung ein großer Nachteil. Die Wahrscheinlichkeit, dass sich die richtige E-Mail-Adresse darunter befindet, wird dadurch kleiner.

Im Gegensatz dazu, kann ein eigener Algorithmus jegliche Information mit in die Generierung einer E-Mail-Adresse einfließen lassen. Ein Beispiel hierfür wäre das Geburtsjahr einer Zielperson. Das OSINT-Tool [Baz] verwendet das nicht. Allerdings können die möglichen Adressen, welche von dem OSINT-Tool generiert wurden, als Anregung und Ideengeber für den eigenen Algorithmus dienen.

Für die erfolgreiche Simulation eines Phishing-Mail-Angriffes, wird die richtige E-Mail-Adresse benötigt. Aus diesem Grund wird der eigene Algorithmus verwendet, damit die Wahrscheinlichkeit erhöht wird, dass sich die korrekte Mailadresse in dem Pool befindet.

E-Mail-Text

Der Inhalt einer E-Mail ist sehr wichtig für die Glaubwürdigkeit einer Phishing-Mail. Aus diesem Grund ist es von Bedeutung, dass der E-Mail-Text Sinn ergibt und mit einer guten Grammatik geschrieben wurde. Bei der dynamischen Texterzeugung mit der Fragment-basierten Methode 6.2.2, kann die Grammatik und der Zusammenhang des Textes zu einer Problematik führen. Durch die Verkettung von verschiedensten Fragmenten kann ein Text erzeugt werden, welcher kein sinnvoller Zusammenhang hat. Allerdings muss hierbei nicht für jede Kombination aus gewonnen Daten ein vollständiges Fragment-Muster erstellt werden. Wogegen bei der Verwendung von fertigen Lückentexten, ein Muster für jede Kombination aus gewonnen Daten vorhanden sein muss. Dennoch ist die Glaubwürdigkeit

durch einen sinnvollen E-Mail-Text höher. Aus diesem Grund werden die vollständigen E-Mail-Muster umgesetzt.

8 OSINT einer ausgewählten Person

8.1 Auswahl der Programmiersprache

Damit das Programm anhand den Lösungsideen umgesetzt werden kann, ist der erste Schritt die Auswahl der Programmiersprache.

Hierbei wird keine Anforderung an die Geschwindigkeit der Sprache gestellt, da beim web scraping das Internet den zeitlichen Engpass darstellt. Allerdings wäre es von Vorteil wenn bereits entwickelte Bibliotheken für das OSINT vorhanden sind. Die Eingabe der Information für die Suche kann über eine Konsole oder über eine graphische Benutzeroberfläche möglich sein.

Als mögliche Programmiersprachen zählen Python, Ruby, C++.

Für web-basierende Anwendung eignet sich eine dynamische Programmsprache. Im Gegensatz zu Python und Ruby zählt C++ nicht zur Familie der dynamischen Programmiersprachen und fällt aus diesem Grund als mögliche Lösung heraus.

Python und Ruby können beide Webseiten, die JavaScript zum rendern benötigen, laden. Dies ist mit Hilfe eines automatisierten Webbrowsers möglich. Des Weiteren lässt sich die Anwendung durch beide Sprachen, entsprechend den Anforderungen entwickeln. Es kann sowohl eine Oberflächenanwendung als auch eine Konsolenanwendung programmiert werden. Zusätzlich bringen beide Sprachen Module mit sich, um das Projekt mit den vorgegebenen Zielen umzusetzen. Somit haben beide Programmiersprachen die Voraussetzungen für die Entwicklung der Anwendung. Allerdings bietet Python in diesem Bereich eine große Community und eignet sich sehr gut für die Bearbeitung von linguistischen Daten. [BKL09] Aus diesen Gründen wird die zu erstellende Anwendung mit der Programmiersprache Python entwickelt.

8.2 Methoden zur Suche nach einer Person im Internet

Die Art der Personensuche wird abhängig von den eingegeben Daten variiert. Das heißt, dass die eingegebenen Daten über die Zielperson vor der Suche analysiert werden und dementsprechend angepasst wird. Nachstehen werden zwei grundsätzliche Methoden für die Art der Personensuche beschrieben.

8.2.1 Personensuche mit Hilfe einer Suchmaschine

Hier wird mit Hilfe einer Suchmaschine nach Informationen gesucht. Mögliche Suchmaschinen sind die von Google und Bing. Allerdings muss nicht für jede Suche eine Suchmaschine verwendet werden. Die nachfolgenden Fälle sollen diesen Ansatz verdeutlichen.

Im Fall, dass der Vorname, Nachname und Wohnort der gesuchten Person eingegeben wird, kann mit Hilfe der festgelegten Suchmaschine nach Information gesucht werden. Die von der Suchmaschine vorgeschlagenen Seiten werden anschließend analysiert, ausgelesen und gespeichert. Dadurch können weitere Informationen gewonnen werden. Falls Benutzernamen von anderen Webseiten wie Instagram, Facebook oder ähnliches vorgeschlagen werden, kann somit die Suche mit diesen Daten speziell auf den entsprechenden Seiten erweitert werden.

Ein weiterer Fall beschreibt das Szenario, wenn ein Benutzername der gesuchten Person in das Programm eingegeben wird. Hierbei handelt es sich um einen Benutzernamen von Social-Media-Webseiten wie Facebook, Instagram, LinkedIn, et cetera.

Zuallererst, wird hier nach Einträgen auf der entsprechende Webseite zu dem angegebenen Benutzernamen gesucht. Dadurch können zusätzliche Daten herausgefunden werden. Diese sind bei der weiteren Suche von Vorteil.

Sobald die Webseite mit Hilfe des Nutzernamens durchsucht und ausgewertet wurde, kann die Suche mit einer Suchmaschine und den gewonnen Daten erweitert werden.

8.2.2 Personensuche auf festgelegten Webseiten

Unabhängig von den eingegebenen Daten, wird eine festgesetzte Anzahl von Webseiten durchsucht. Als potentielle Kandidaten-Webseiten eignen sich die Social-Media-Seiten wie Facebook, Instagram, Twitter, LinkedIn, et cetera. Diese Art der Personensuche arbeitet allerdings ohne die Verwendung einer Suchmaschine.

8.3 Bewertung der Methoden zur Personensuche

Um möglichst viele Informationen über eine Person im Internet zu finden, bietet die Personensuche mit der Verwendung einer Suchmaschine die beste Lösung. Es wird anstatt ausschließlich festgelegten Seiten das ganze Internet durchsucht. Dadurch können wesentlich mehr individuelle Einträge gefunden werden. Des Weiteren wird keine Logik zur Suche nach Einträgen im Internet benötigt, da lediglich den vorgeschlagenen Suchergebnissen gefolgt werden kann.

Allerdings muss beachtet werden, dass Benutzer bei verschiedensten Social-Media-Seiten auswählen können, ob das Benutzerprofil von einer Suchmaschine gefunden werden kann oder nicht. Bekannte Webseiten die diese Einstellungsmöglichkeiten unterstützen sind XING und LinkedIn. Aus diesem Grund, werden zu Beginn der Suche die Social-Media-Seiten durchsucht. Dadurch können vor der Google-Suche zusätzliche Informationen herausgefunden werden, die für das später OSINT von Vorteil sind. Falls sich eine Social-Media-Seite unter den Google-Suchergebnissen befindet, kann diese nachträglich ebenfalls durchsucht werden.

8.3.1 Auswahl der Suchmaschine

Laut Expertenaussage sucht Bing tiefgreifender nach Information auf Social Media Seiten wie Facebook, Twitter und LinkedIn. Allerdings finden nur 3,5% aller Suchanfragen in Deutschland über Bing statt. Im Gegensatz dazu hat Google einen Marktanteil von 91,2% in Deutschland. Diese Zahlen sprechen eindeutig für Google. Durch die höhere Anzahl von Suchanfragen, können mehr Daten erfasst und die Ergebnislisten besser gerankt werden.

Dies hat zu Folge, dass Bing bei einer konkreten Suche schlechter abschneidet. [Boh14] Grundsätzlich stellt die Verwendung von zwei Suchmaschinen die beste Lösung dar, da die Wahrscheinlichkeit für einen Suchtreffer erhöht wird. Dennoch wird in dieser Arbeit ausschließlich die Suchmaschine von Google verwendet, da sie gegenüber dem Konkurrenten keine Nachteile hat. Selbst die detailliertere Suche auf Sozialen Netzwerken, bringt bei der hier verwendeten Personensuche keinen großen Vorteil für Bing. Das heißt, durch die Analyse der Suchergebnisse, wird erkannt ob sich die bekannten Social Media Webseiten darunter befinden. Falls diese es nicht tun, wird die Suche auf den entsprechenden Sozialen Netzwerken erweitert.

8.4 Implementierung der Personensuche mit Hilfe der Google-Suchmaschine im Internet

Die Suchmaschine von Google wird für die Personensuche im Internet verwendet. Gesucht wird nach den eingegebenen Daten, welche über die Konsole eingelesen werden.

8.4.1 Eingabe der bekannten Daten

Es besteht die Möglichkeit den **Vorname**, **Nachname**, **Wohnort**, **Arbeitgeber**, **Instagram Benutzername**, **Facebook Benutzername**, **Twitter Benutzername**, **E-Mail-Adresse**, und das genaue beziehungsweise geschätzte **Geburtsjahr** der gesuchten Person über eine Konsole einzugeben. Falls der genaue Jahrgang der Zielperson nicht bekannt ist, kann ein geschätztes Geburtsjahr eingetragen werden. Dies kann später bei der Identifizierung der gesuchten Person hilfreich sein.

Zu Beginn werden alle Personen-Variablen mit einem leeren String initialisiert. Das bedeutet, alle Variablen, zu denen keine Information eingegeben wurde, enthalten einen leeren String.

Verarbeitung der Daten

Im ersten Schritt wird kontrolliert, welche Informationen vom Programm-Nutzer eingegeben wurden. Der Vorname und Nachname sind nicht ausreichend für die Suche. Es wird mindestens ein weiteres Attribut benötigt. Dagegen ist der Benutzernamen von Instagram und Twitter sowie die E-Mail-Adresse einzigartig. Dadurch kann mit einem dieser Attribute gesucht werden.

Bei der Eingabe des Wohnortes, kann dieser vor der Suche mit der entsprechenden Wortsammlung verglichen werden. Falls sich der Wohnort nicht in der Datenbank befindet, wird er nachträglich ergänzt. Für die Personenerkennung ist es wichtig, dass sich der korrekt Wohnort in der Datenbank befindet.

Daraufhin werden mit diesen Eingaben Kombinationen für die Suche und die URL-Generierung erstellt. Mögliche Such-Kombinationen für erfolgreiche Ergebnisse sind:

Vorname, Nachname, Wohnort;

Vorname, Nachname, Geburtsjahr;

Vorname, Nachname, Institution;

Vorname, Nachname, Wohnort, Geburtsjahr;

Vorname, Nachname, Wohnort, Institution;

Benutzername einer Social-Media-Seite;

Die Kombination aus vielen oder allen Daten ist ebenfalls eine mögliche Option. Allerdings wird dadurch oft kein Ergebnis gefunden, da nicht zur jeder Information ein Eintrag im Internet besteht.

Sobald die Kombinationen aus den Daten bekannt sind, werden die Such-URLs für die Google-Suchmaschine generiert.

8.4.2 Generierung der Google-Such-URLs

Aufbau eines URLs

Ein Uniform Resource Locator (URL) lokalisiert eine Ressource, indem eine abstrakte Identifikation der Lokalisierung verwendet wird. Dabei wird ein URL grundsätzlich im

folgenden Format angegeben. [RFC94]

< scheme > : < scheme – specific – part > [RFC94]

Das Schema gleicht hierbei meist dem verwendeten Protokoll wie HTTP oder FTP. Der Doppelpunkt stellt die Trennung zum Schema-spezifischen Teil dar. Ein Beispiel für ein HTTP-URL-Aufbau ist im Folgenden definiert. [RFC94]

http : // < host > : < port > / < path > ? < searchpart > [RFC94]

Hier wird das Protokoll HTTP als Schema verwendet, wobei sich der Aufbau bei der Verwendung des HTTPS-Protokolls kaum unterscheidet. Lediglich das Schema und der Port verändert sich.

Für den <host> kann der FQDN oder die IP-Adresse des Hostrechners eingetragen werden. Wenn der Port nicht angegeben wird, ist der Standardport voreingestellt. Bei HTTP wäre dies Port 80 und bei HTTPS Port 443. Der <path> stellt ein HTTP-Selektor dar und ist mit einem Fragezeichen von der Suchzeichenkette getrennt. [RFC94]

Im Bereich des <searchpart> lassen sich URL-Parameter einfügen um Informationen an die entsprechende Webseite mitzugeben. Die Parameter bestehen aus einem Schlüssel und aus einem Wert, welche durch ein Gleichheitszeichen getrennt werden. Um mehrere Parameter hinzuzufügen und zu kombinieren wird das kaufmännische Und-Zeichen verwendet. [AH19] Ein URL für die Google-Suche von *Max Mustermann* ist in dem folgenden Beispiel gegeben.

https : //www.google.com/search?q = Max + Mustermann

Allerdings können URLs nur mit ASCII-Zeichen erzeugt und versendet werden. Aus diesem Grund müssen Zeichen die nicht im ASCII vorkommen, in ein gültiges Format umgewandelt werden. Dies wird realisiert, indem die URL-Kodierung das nicht enthaltende ASCII-Zeichen durch ein “%”, gefolgt von zwei Hexadezimalen Ziffern, ersetzt. Beispielsweise repräsentiert “%20” ein Leerzeichen und “%22” ein Anführungszeichen. [W3S]

Erstellen der Such-URLs

Dieser Absatz beschreibt die Erstellung der Such-URLs für Google, mit dem Wissen aus Kapitel 8.4.2.

Für jede genannte Kombination aus den eingegebenen Daten werden Link-Muster erzeugt. Diese entsprechen einem Lückentext. Sobald die entsprechenden Muster ausgewählt wurden, werden die Lücken mit den Daten befüllt. Dadurch wird eine Liste mit einer variierende Menge von Suchlinks erstellt. Diese Liste wird anschließend von dem Web Crawler verwendet um die Suche zu starten. Ein URL für die Suche nach Information auf beliebigen Webseiten wird wie folgt dargestellt.

<https://www.google.com/search?q=%22Max+Mustermann%22+%22Weingarten%22>

Wenn allerdings der Benutzername einer Social-Media-Seite bekannt ist, werden zwei unterschiedliche URLs verwendet. Mit Hilfe des ersten URLs, wird speziell nach Einträgen auf der entsprechenden Webseite gesucht. Dazu kann der Operator “site“ verwendet werden. Dieser beschränkt die Suchergebnisse soweit, dass die vorgeschlagenen Einträge ausschließlich auf einer festgelegten Webseite vorkommen. Das folgende Beispiel beschreibt die Suche nach dem Benutzer “Mustermann“ auf der Webseite “Instagram.com“. Dabei ersetzt die ASCII-Zeichenkette “%3A“ den Doppelpunkt. [W3S]

<https://www.google.com/search?q=site%3Ainstagram.com+%22Mustermann%22>

Der zweite URL wird für eine Social-Media-Suche verwendet. Bei dieser Suche werden Social-Media-Seiten nach Einträgen durchsucht. Dafür wird kein zusätzlicher Operator benötigt. Es wird lediglich ein @-Zeichen, welches mit der Zeichenkette “%40“ dargestellt wird, vor dem zu suchenden Wort eingefügt. Die Social-Media-Suche nach dem Benutzernamen “Mustermann“ sieht folgendermaßen aus. [Goo19]

<https://www.google.de/search?q=%40Mustermann>

Optimierung der Such-URLs

Um die Suchergebnisse von Google zu verbessern, können die Suchbegriffe in Anführungszeichen gesetzt werden. Dadurch wird eine Phrasensuche gestartet, die nach einer

Zeichenfolge sucht. Das bedeutet, es wird ausschließlich nach diesen Zeichenfolgen gesucht und nicht nach einer Abwandlung. Ein Beispiel hierfür ist die Suche nach “Mike Bazzell“. Wenn diese Suche ohne Anführungszeichen durchgeführt wird, werden zusätzlich Webseiten vorgeschlagen die den Namen Mike Bazzell anstatt Micheal Bazzell beinhalten. Diese erweiterte Suche kann dazu führen, dass unzählige Webseiten vorgeschlagen werden, die nicht unbedingt was mit dem Thema der Suchbegriffe zu tun hat. Um dem vorzubeugen können Anführungszeichen verwendet werden, welche die Anzahl der Suchergebnisse um einen sehr großen Teil verringern. [Baz18]

Für die Suche nach **Marco Lang** werden ungefähr **96.400.000** Ergebnisse mit Hilfe der Google-Suchmaschine gefunden. Wird die Suche mit den Anführungszeichen verfeinert indem nach “**Marco**“ “**Lang**“ gesucht wird, werden etwa **55.600.000** Ergebnisse gefunden. Allerdings werden hier Webseiten vorgeschlagen, welche die Wörter “Marco“ und “Lang“ beinhalten, jedoch müssen diese nicht direkt nebeneinander und auch nicht in der Reihenfolge vorkommen. Es wäre Möglich, dass bei dieser Suche, Webseite mit Verweisen auf die Namen “Marco Mustermann“ und “Max Lang“ beinhaltet. Aus diesem Grund kann nach “**Marco Lang**“ gegoogelt werden. Dadurch wird die Anzahl der Suchergebnisse auf **45.500** Ergebnisse reduziert. Der Grund für die starke Reduzierung ist, dass ausschließlich die Webseiten vorgeschlagen werden, die den kompletten String “Marco Lang“ beinhalten. Für eine weitere Optimierung der Ergebnisse, wird der Wohnort hinzugefügt, wie in dem Beispiel “**Marco Lang**“ “**Tett nang**“. Dadurch werden die Suchvorschläge auf lediglich **95** Ergebnisse reduziert. Die URL zu dieser optimierten Suche lautet:

<https://www.google.com/search?q=%22Marco+Lang%22+%22Tett nang%22>

Nicht nur die Reduzierung der Suchergebnisse, sondern auch das herausfiltern von unerwünschten Webseiten hat einen positiven Effekt auf die zu erstellende Anwendung, da die vorgeschlagenen Seiten in den folgenden Schritten analysiert werden müssen. Das bedeutet, dass jede unerwünschte Seite die allein durch die Suche herausgefiltert werden kann, einen großen Laufzeitvorteil mit sich bringt.

8.4.3 Mit welcher Bibliothek werden Serveranfragen umgesetzt?

Damit eine Person im Internet gesucht werden kann, muss das Programm in der Lage sein, Anfragen an einen Server zu versenden und die dazugehörigen Antwort zu empfangen. Im Folgenden werden drei Möglichkeiten beschrieben, um Anfragen an einen Server zu versenden. Zum einen ist das die Python Request-Bibliothek, welche sich optimal für HTTP-Anfragen eignet. [Mit15] Zum anderen bietet sich die Verwendung eines automatisierten Webbrowsers an, was mit Hilfe der Selenium Python API realisierbar ist. [Law15] Über diese API ist es möglich auf alle Funktionen des Selenium WebDrivers zuzugreifen. [Mut18] Eine Alternative dazu, ist das Python Framework Scrapy, welches zum Crawlen von Webseiten und Extrahieren von Daten verwendet werden kann. [dev18] Die letzte Möglichkeit stellt die Scrapy Middleware Scrapy-Selenium dar. [Fou19] Dadurch wird die Kommunikation von Scrapy und Selenium ermöglicht.

Für komplizierte Anfragen an einen Server eignet sich die Request-Bibliothek von Python sehr gut. Der Umgang mit Cookies, Header und vielem mehr ist sehr einfach gestaltet. Auch die Generierung des Such-URLs wird von dieser Bibliothek übernommen. Des Weiteren hat Requests einen großen Laufzeit-Vorteil gegenüber dem automatisierten Webbrowser und kann HTTP-Fehlermeldungen empfangen. Allerdings lässt sich mit der Request-Bibliothek keine Javascript-Seite auslesen.

Wenn das Framework Scrapy standardmäßig verwendet wird, können ebenfalls keine Javascript-Seiten ausgelesen werden. Doch in Scrapy lässt sich ein automatisierter Webbrowser einfügen, mit welchem das Auslesen von Javascript-Webseiten möglich ist. Zusätzlich lässt sich mit Scrapy ein effektiver Web Crawler und Web Scraper entwickeln, was für die nächsten Schritte ein erheblicher Vorteil ist.

Aus den erläuternden Gründen, wird das Framework Scrapy mit der Verbindung eines automatisierten Webbrowsers für die Personensuche verwendet. Der automatisierte Webbrowser muss in dem Framework implementiert werden, da auf bestimmte Webseiten mit Javascript direkt zugegriffen wird. Infolgedessen wird die Middleware Scrapy-Selenium verwendet, da sie die eine kompakte Möglichkeit bietet, den automatisierten Webbrowser in Scrapy zu implementieren. Durch diese Kombination aus Scrapy und dem Selenium WebDriver, lassen sich Javascript-Seiten problemlos auslesen.

Zusätzlich zu diesem Framework wird ein unabhängiger Selenium-Wedriver implementiert.

Dieser wird für den Umgang mit den Social-Media-Seiten benötigt, da auf diesen Seiten ein Login vollzogen werden muss. Das hat den Vorteil, dass die Anmeldung in der Session gespeichert wird. Somit muss bei einem erneuten Zugriff auf die selbe Seite keine neue Anmeldung vollzogen werden.

8.4.4 Web Crawler erstellen

Nachdem der Selenium WebDriver in das Scrapy Framework implementiert wurde, kann mit dem crawling begonnen werden. Der Web Crawler hat die Aufgabe ausgewählte Social-Media-Seiten zu durchstöbern und den von Google vorgeschlagenen Webseiten zu folgen. Wie in Bild 8.1 gezeigt, werden zuerst die Informationen über die Zielperson eingelesen. Anschließend werden die Social-Media-Seiten behandelt. Dadurch können weitere Informationen über die Person gefunden werden. Die angegebenen Daten über die Zielperson, werden zur Generierung der Google-Such-Links verwendet. Mit diesen Links und der Google-Suchmaschine werden Webseiten gesucht, die mögliche Inhalte betreffend der Zielperson enthalten. Im nächsten Schritt wird die Google-Webseite mit den Suchergebnissen analysiert und ausgelesen. Dadurch können die URLs für die entsprechenden Webseiten gewonnen werden. Diesen URLs wird anschließend gefolgt, um Informationen über die Zielperson zu gewinnen.

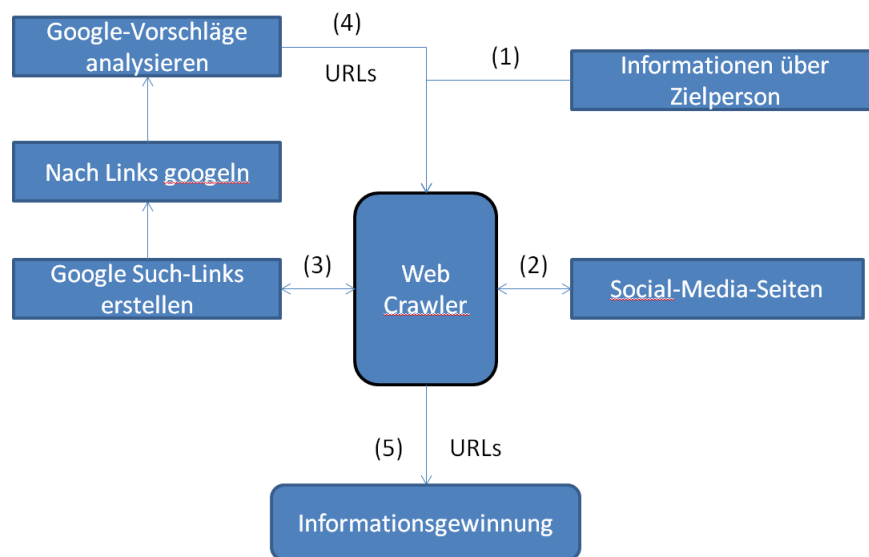


Bild 8.1: Aufbau des Web Crawlers

Social-Media-Seiten

Zu den verwendeten Social-Media-Webseiten gehören Instagram, Facebook, Twitter, Xing und LinkedIn. Für diese Webseiten werden gefälschte Accounts und ein eigener Selenium WebDriver erstellt. Dieser automatisierte Webbrowser ist ausschließlich für die Social-Media-Seiten zuständig. Damit vollständige Profile angezeigt werden können, loggt sich dieser automatisch ein. Die gewonnenen Informationen werden dem Profil der gesuchten Person hinzugefügt.

Es kann passieren, dass ausschließlich der Benutzername von einer Person bekannt ist. Weiter Attribute wie Vorname, Nachname und Wohnort sind nicht bekannt. In dem Fall, kann mit einem einzigartigen Benutzername nach diesen Attributen auf der entsprechenden Webseite gesucht werden. Jedoch verwendet nur Instagram und Twitter ein einmaligen Benutzernamen. Aus diesem Grund kann die erweiterte Suche nur auf diesen beiden Plattformen umgesetzt werden.

Login-Formulare

Der Selenium WebDriver muss sich nicht bei jeder Social-Media-Webseite einloggen. Zu Beginn wird kontrolliert, zu welcher Plattform ein Benutzername eingegeben wurde. Dazu gibt es bei dieser Anwendung die Möglichkeit einen Facebook-, Instagram- oder

Twitter-Benutzername einzugeben. Je nach Eingabe meldet sich der Browser auf den entsprechenden Seiten an. Auf den Seiten Xing und LinkedIn wird sich standardmäßig angemeldet. Das hat den Grund, dass diese beiden Seiten immer durchsucht werden sollen. Bei der Umsetzung wird im ersten Schritt die Login-Seite der entsprechenden Plattform angefordert. Die Antwort wird mit Hilfe der BeautifulSoup-Bibliothek nach dem HTML-Tag `<input>` durchsucht. Allerdings hat nicht jede Webseite den komplett identischen Aufbau. Das bedeutet, dass sich bei diesen Tags die Attribute unterscheiden können. Aus diesem Grund muss nach der Angeforderten-Webseite die Suche der `<input>`-Tags unterschieden werden.

Die Anmeldung von Instagram, Facebook und LinkedIn sind nahezu identisch. Hier können die zwei gesuchten `<input>`-Felder mit Hilfe des Attributs *type* identifiziert und gefunden werden. Bei der Twitter-Login-Seite muss anstatt dem *type* nach dem Attribute *class* gesucht werden. Andernfalls findet der Browser keine interaktiven Elemente.

Zum übertragen der Benutzernamen und Passwörter, benötigt der Selenium WebDriver ein Element mit eindeutigen Attribut zur Referenzierung. Dafür dient bei Instagram, LinkedIn und Facebook das vorher gefundene Element mit dem Attribute `id`. Bei Twitter ist das das Attribut `class[0]` und bei Xing `name`.

Instagram

Instagram und Twitter verwenden beide einen einzigartigen Benutzername. Dies ist ein riesigen Vorteil für die Identifikation einer Person, wenn der Benutzername bekannt ist. Allerdings ist es ebenfalls möglich eine Person mit ihrem offiziellen Namen zu suchen und zu finden.

Im Fall, dass der Instagram-Benutzername eingegeben wurde, wird die dazugehörige Profilseite angezeigt und nach Informationen durchsucht. Im nächsten Schritt dienen die vorgeschlagenen Freunde, welche in Beziehung zu diesem Profil stehen, als weitere Informationsquelle. Das bedeutet, dass diese Kontakte ebenfalls durchsucht werden. Bei Übereinstimmungen wie zum Beispiel der selben Universität, kann später eine glaubwürdige Phishing-Mail generiert werden.

Falls sich jedoch eine Instagram-Seite unter den Google-Vorschlägen befindet, und die Übereinstimmung des Profils mit der Zielperson nicht klar ist, können die Vorgeschlagenen Kontakte für die Identifizierung der Zielperson genutzt werden. So wird beispielsweise erkannt, wenn ein Freunden aus der selben Stadt vorgeschlagen wird, dass es sich um die

gesuchte Person handeln kann. Diese Methode erzielt kein sicheres Ergebnis. Jedoch kann die Wahrscheinlichkeit erhöht werden, dass es sich um die richtige Person handelt.

Die Profilseite wird mit dem folgenden Link <https://www.instagram.com/username/> angefordert. Weitere Seiten und Einträge der gesuchten Person, unabhängig von der Profilseite, können mit dem Suchbefehl **site:instagram.com "username"**

-site:instagram.com/username angezeigt werden. [Baz18] In der Anwendung wird dies mit dem folgenden URL umgesetzt.

<https://www.google.com/search?q=site%3Ainstagram.com+%22username%22+-site%3Ainstagram.com%2Fusername&oq=site%3Ainstagram.com+%22username%22+-site%3Ainstagram.com%2Fusername>

Die dazugehörigen Suchergebnisse werden anschließend gleich den normalen Google-Suchergebnissen behandelt.

Twitter

Auf der Webseite Twitter wird ausschließlich die Profilseite nach Informationen durchsucht. Dies ist mit dem Link **<https://twitter.com/username>** möglich.

Facebook

Facebook bietet ein großes Potential um OSINT zu betreiben. Allerdings hat Facebook optimal Algorithmen zur Erkennung von automatisierten Crawlern entwickelt. Aus diesem Grund wurde das Fake-Konto nach wenigen Versuchen gesperrt. Das hat zu Folge, dass auf dieser Plattform nur begrenzt gesucht werden kann, da keine Anmeldung vorgenommen wird. Um das Konto zu entsperren müssten eine Kopie des Ausweises, ein Bild mit erkennbarem Gesicht und eine Handynummer an Facebook übermittelt werden.

Aus diesen Gründen wird Facebook nur dann und ohne Anmeldung verwendet, wenn sich ein Vorschlag unter den Google-Suchergebnissen befindet.

LinkedIn und XING

LinkedIn und XING bieten eine optimale Informationsquelle bezüglich der schulischen und beruflichen Laufbahn der Zielperson. Allerdings gibt es hier keinen einzigartigen Benutzernamen. Demzufolge, wird eine Person mit dem vollen Namen und dem aktuellen Wohnort gesucht. Dabei wird auf LinkedIn ein Filter angewendet, bei dem ausschließlich Personen aus Deutschland angezeigt werden. Wenn genau eine Person vorgeschlagen wird, wird

dieses gefundene Personenprofil durchsucht. Bei einer Mehrzahl von gefundenen Personen werden diese nicht auf Informationen durchsucht, da keine Identifikation möglich ist.

Die Personensuche bei LinkedIn, mit angewandtem Filter, wird mit dem URL

https://www.linkedin.com/search/results/people/?facetGeoRegion=%5B%22de%3A0%22%5D&keywords=vorname%20nachname%20wohnort&origin=FACETED_SEARCH dargestellt.

Bei Xing sieht der Such-URL wie folgt aus.

<https://www.xing.com/search/old/members?hdr=1&keywords=vorname+nachname+wohnort>

Webseite mit den Suchergebnissen von Google analysieren

Zur Analyse der Webseite mit den Suchergebnissen von Google, wird der Seiten Quelltext benötigt. Mit Hilfe des Quelltextes, können die entsprechenden Links erkannt werden. Der Seiten Quelltext wird mit Hilfe der BeautifulSoup-Bibliothek angezeigt.

Das Bild 8.2 stellt ein Suchergebnis von Google dar. Der dazugehörige Quelltext befindet sich in der Darstellung 8.1.



Marco Lang - Spieler - FuPa - FuPa
<https://www.fupa.net/spieler/marco-lang-1261543.html> ▼ Translate this page
Marco Lang ist ein Fußballspieler der diese Saison noch keine Einsätze zu verzeichnen hat. ...
Geburtsdatum: 11.08.1995 (23). Nationalität ... TSV Tettnang.

Bild 8.2: Google-Suchergebnis [LLC19]

Im Ausschnitt des Seiten Quelltextes 8.1 ist zu sehen, dass der <div>-Container mit der Klasse "g" einen Hyperlink enthält. URLs oder Links werden mit dem HTML-Tag <a> dargestellt. Dieser Link wird für die Suche benötigt. Deswegen wird genau nach diesem Link gesucht.

Da der <div>-Container bei jedem Suchergebnis identisch ist, kann bei jedem Ergebnis nach dem entsprechenden <div>-Container gesucht werden. Anschließend kann der erste Link in diesem <div>-Tag ausgelesen werden. Dies wird mit Hilfe der BeautifulSoup-Bibliothek umgesetzt.

Um zu erkennen, ob mehrere Seiten mit Suchergebnissen existieren, wird nach bestimmten

Hyperlinks gesucht. Diese Links werden über das Attribut “class” identifiziert. Mit Hilfe der BeautifulSoup-Bibliothek wird nach dem Klassennamen “fl” gesucht. Falls weitere Seiten mit Suchergebnissen vorhanden sind, werden die dazugehörigen Links in einer Liste gespeichert. Anschließend wird ihnen gefolgt und die neue Seite wird nach weiteren URLs durchsucht.

```
<div class="g">
  <h3 class="r">
    <a href="/url?q=https://www.fupa.net/spieler/marco-lang-1261543.html&sa=U&ved=0ahUKEwiZ3PDGqMvhAhWtURUIHU7VAcwQFggUMAA&usg=A0vVaw2QiSMFzScB0JcvoPCisBGw"><b>Marco Lang</b>- Spieler - FuPa - FuPa
  </a>
</h3>
```

Listing 8.1: Ausschnitt des Quelltextes von einem Google-Suchergebnis [LLC19]

8.5 Methoden zum Erkennen von wichtigen Informationen auf einer Webseite

Bei der Suche nach einer ausgewählten Person können verschiedenste Arten von Webseiten gefunden werden. Aus diesem Grund muss das Programm eine gewisse Intelligenz mit sich bringen, um die wichtigsten Daten aus einer Seite herauszufiltern. Dabei ist es nicht möglich festgelegte Bereiche einer Webseite durch eine Hartkodierung auszulesen, da jede Webseite eine individuelle Struktur hat.

Die Grundidee zur Lösung dieser Probleme ist die Analyse des vorliegenden Webseiten-Textes. Eine Methode zur Textanalyse ist die automatisierte Schlüsselwort-Gewinnung. Hierbei wird die HTML-Seite zu einem verwendbaren Text formatiert, wobei alle Sonderzeichen herausgefiltert werden. Im nächsten Schritt werden Schlüsselwörter aus dem formatierten Webseitentext generiert. Möglichkeiten zur automatisierten Schlüsselwortgenerierung sind die Verfahren RAKE 8.5.1 und die Automatic Keyword Extraction mit NLP 8.5.2, welche im Laufe dieser Arbeit detailliert beschrieben werden.

Eine weitere Methode zur Textanalyse wäre der Vergleich von Zeichenketten. Dabei wird der vorliegende Webseitentext in einen String umgewandelt. Nachdem die Schlüsselwörter generiert und in Listen gespeichert wurden, werden Wortsammlungen erstellt. Diese Wortsammlungen sind Listen, welche aussagekräftige Schlüsselwörter enthalten und nach Themen kategorisiert sind. Beispiele für den Inhalt der Listen sind alle Hochschulen und Universitäten in Deutschland, Berufsbezeichnungen und Tätigkeiten, Studiengänge, Hobbybezeichnungen und alle Städte und Gemeinden in Deutschland.

Mit diesen Wortsammlungen kann nun die Liste mit den bereits generierten Schlüsselwörtern aus dem Webseitentext verglichen werden. Bei einer Übereinstimmung eines Schlüsselwortes wird das Wort mit der entsprechenden Kategorie vorgemerkt und später in die verwendete Speicherstruktur eingetragen.

Die Wortsammlungen werden mit Hilfe von bekannten Listen im Internet eigenständig befüllt. Als Informationsquelle dienen alle öffentlich frei zugänglichen Listen, die hilfreiche Informationen enthalten.

8.5.1 RAKE

RAKE steht für *Rapid Automatic Keyword Extraction* und stellt eine sehr effiziente Methode zur Schlüsselwortgenerierung dar. Die Funktion von RAKE basiert darin, dass Schlüsselwörter mehrere Wörter mit inhaltlicher Relevanz enthalten, allerdings selten Stoppwörter und Sonderzeichen. [RECC10]

Als Stoppwörter werden Wörter bezeichnet, die sehr oft auftreten und keinen großen Informationsgewinn mit sich bringen. Beispiele dafür sind *und*, *weil*, *der* oder *als*. [Sla]

In einer jungen Wissenschaft wie der Informatik mit ihrer Vielschichtigkeit und ihrer unüberschaubaren Anwendungsvielfalt ist man oftmals noch bestrebt, eine Charakterisierung des Wesens dieser Wissenschaft und Gemeinsamkeiten und Abgrenzungen zu anderen Wissenschaften zu finden. Etablierte Wissenschaften haben es da leichter, sei es, dass sie es aufgegeben haben, sich zu definieren, oder sei es, dass ihre Struktur und ihre Inhalte allgemein bekannt sind.

Bild 8.3: Beispieltext [SS11]

Zu Beginn wird der zu analysierende Text, hier der Beispieltext in Bild 8.3, durch einen Worttrenner in ein Array, bestehen aus möglichen Schlüsselwörtern, aufgeteilt. Das erzeugte Array wird anschließend in Sequenzen von zusammenhängenden Wörtern unterteilt. Dabei erhalten die Wörter in einer Sequenz die gleiche Position und Reihenfolge wie im Ursprungstext und dienen gemeinsam als Kandidatenschlüsselwort. [RECC10]

Nachdem die möglichen Schlüsselwörter identifiziert sind, wird für jeden einzelnen Kandidaten ein Score ausgerechnet. Dieser besteht aus dem Quotient des Grades $deg(w)$ und der Häufigkeit des Vorkommens eines Wortes innerhalb der Kandidaten $freq(w)$. Daraus ergibt sich die Formel:

$$deg(w)/freq(w)$$

Dabei beschreibt der Grad eines Wortes, dass gemeinsame Auftreten mit sich selbst und anderen Schlüsselwörtern. In der Tabelle 8.5.1 ist der Grad für jedes Wort ablesbar, indem die Einträge in der entsprechenden Reihe summiert werden. Beispielsweise beträgt der Grad des Wortes “Wissenschaft” den Wert 3. Dies ergibt sich aus der Rechnung:

$$2 + 1 = 3$$

Das Wort “Wissenschaft” kommt hier selbst zweimal in dem Kandidaten-Array vor und davon einmal in Verbindung mit dem Worten “jungen”.

Die Häufigkeit des Vorkommens eines Wortes lässt sich ebenfalls in der Tabelle 8.5.1 ablesen. Allerdings muss hier in der Reihe und Spalte des jeweiligen Wortes nachgeschaut werden. Für das Wort “Wissenschaft” beträgt die Häufigkeit des Vorkommens den Wert 3. Zusammenfassend kann gesagt werden, dass $deg(w)$ die Kandidaten bevorzugt, welche oft und in langen Schlüsselwörtern, die mehrere Wörter enthalten, vorkommen. Dies bedeutet, dass beispielsweise $deg(etabliert)$ eine höhere Bewertung als $deg(informatik)$ bekommt, obwohl beide Wörter gleich oft im Text vorkommen. Dagegen wird bei $freq(w)$, ausschließlich die Häufigkeit des Vorkommens bewertet. Bei der Formel $deg(w)/freq(w)$ werden die Wörter bevorzugt, welche überwiegend in langen Kandidatenwörtern vorkommen. Diese Formel bietet dadurch einen guten Mittelweg zur Schlüsselwortgewinnung. Ein Beispiel dafür sind die Wörter “Wissenschaften und “allgemein“. Hier ist der Quotient von $deg(allgemein)/freq(allgemein)$ höher als von $deg(Wissenschaften)/freq(Wissenschaften)$,

obwohl die Häufigkeit des Wortes “*Wissenschaften*“ höher und der Grad gleich hoch ist. [RECC10]

Durch das genannte Verfahren und der Formel $deg(w)/freq(w)$ für die Bewertung, ergeben sich die im Bild 8.4 befindenden Kandidaten mit den dazugehörigen Endbewertungen. [RECC10]

	wissenschaften	wissenschaft	sei	etablierte	informatik	aufgegeben	gemeinsamkeiten	oftmals	charakterisierung	jungen	inhalte	allgemein	bekannt	struktur	wesens	bestrebt	unüberschaubaren	anwendungsvielfalt	definieren	abgrenzungen	leichter	finden	vielschichtigkeit
wissenschaften	2			1																			
wissenschaft		2								1													
sei			1																				
etablierte	1			1																			
informatik					1																		
aufgegeben						1																	
gemeinsamkeiten							1																
oftmals								1															
charakterisierung									1														
jungen		1								1													
inhalte											1	1	1										
allgemein											1	1	1										
bekannt											1	1	1										
struktur														1									
wesens															1								
bestrebt																1							
unüberschaubaren																	1	1					
anwendungsvielfalt																	1	1					
definieren																			1				
abgrenzungen																				1			
leichter																					1		
finden																						1	
vielschichtigkeit																							1

Tabelle 8.1: Co-occurrence

inhalte allgemein bekannt (9.0), unüberschaubaren anwendungsvielfalt (4.0), jungen wissenschaft(3.5), etablierte wissenschaften (3.5), wissenschaften (1.5), wissenschaft (1.5), wesens (1.0), vielschichtigkeit (1.0), struktur (1.0), sei (1.0), oftmals (1.0), leichter (1.0), informatik (1.0), gemeinsamkeiten (1.0), finden (1.0), definieren (1.0), dass (1.0), charakterisierung (1.0), bestrebt (1.0), aufgegeben (1.0), abgrenzungen (1.0)
--

Bild 8.4: Schlüsselwörter mit zugehörigem Score

8.5.2 Automatic Keyword Extraction mit NLP

Bei dieser Methode wird der vorliegende Text in die einzelnen Wörter unterteilt. Dabei wird eine Liste mit potentiellen Schlüsselwörtern erstellt, in der *Stoppwörter* und Sonderzeichen herausgefiltert werden. Bei den Schlüsselwörtern handelt es sich nicht ausschließlich um ein Wort sondern auch um Wortsequenzen. Sogenannte N-Gramme bestehen aus einer festgelegten Anzahl von Wörtern. Dies hat den Vorteil, dass nicht nur Schlüsselwörter bestehend aus einem Wort erstellt werden können, sondern auch Schlüsselwörter mit Fragmenten eines Textes. Diese Art von Schlüsselwort wird benötigt um Informationen wie *Hochschule Ravensburg-Weingarten* herauszulesen.

Erweiternd kann die Anzahl der Schlüsselwörter mit dem Verfahren von Stemming reduziert werden. Durch die Verwendung von ergänzende Regeln wie, eine Mindestanzahl von Buchstaben in einem Wort, können die Schlüsselwörter weiter begrenzen.

8.6 Bewertung der Methoden zum Herausfiltern von wichtigen Informationen auf einer Webseite

RAKE stellt eine fertige Methode dar, um Schlüsselwörter, die den Inhalt eines Textes in kurz wiedergeben, zu erstellen. Dabei hat ein Anwender kaum Möglichkeiten eigene Implementierungen vorzunehmen, da vieles vorgegeben ist. In der zu erstellenden Anwendung soll jedoch nicht der Inhalt eines Textes in Schlüsselwörter zusammengefasst werden, sondern es wird nach informationsreichen Wörtern gesucht. Aus diesem Grund ist jedes einzelne Wort aus dem Webseiten-Text von Bedeutung. Dies spricht gegen RAKE, da es nur die selbst errechnenden Favoriten-Schlüsselwörter zur Verfügung stellt. Dadurch werden viele Wörter nicht in Betracht gezogen oder für weiterführende Bearbeitungen

nicht bereitgestellt. Darüber hinaus ist die Berechnung eines Scores für diese Anwendung nicht notwendig.

Die Methode zur automatisierten Schlüsselwortgenerierung mit NLP bringt dagegen ein eigene Implementationsmöglichkeit mit sich. Das bedeutet, es kann selbst festgelegt werden, aus wie vielen Wörtern die Schlüsselwörter bestehen sollen. Des Weiteren wird jedes einzelne Wort in Betracht gezogen und verwendet.

Die Suche nach einer E-Mail-Adresse im Text lässt sich bei beiden Methoden hinzufügen. Jedoch wird aus den eben genannten Vorteilen, die Information mit Hilfe der Methode zur automatisierten Schlüsselwortgewinnung mit NLP herausgefiltert.

8.7 Implementierung der Methoden zum Herausfiltern von wichtigen Informationen auf einer Webseite

8.7.1 Text formatieren

Bevor die Schlüsselwörter generiert werden können, muss der Text in ein verwertbares Format umgewandelt werden. Aus diesem Grund wird der Seitenquelltext zuallererst mit Hilfe des Python-Skripts `html2text` zu einem ASCII Plaintext umgewandelt. [Fou18] Anschließend werden Zeilenumbrüche und Sonderzeichen aus diesem Text herausgefiltert. Einzelne Wörter und Zahlen die weniger als 2 Zeichen beinhalten, werden ebenfalls aussortiert. Nachdem der Text in ein verwertbares Format umgewandelt wurde, kann mit der Umsetzung für die automatisierte Schlüsselwortgenerierung mit NLP begonnen werden.

8.7.2 Erstellung der Wortsammlungen

Die Informationsgewinnung ist abhängig von diesen Wortsammlungen. Es können nur die Informationen herausgefunden werden, welche auch in einer dieser Listen vorkommen. Infolgedessen ist die Umsetzung der Wortsammlungen sehr wichtig.

In welchem Format werden die Wortsammlungen gespeichert?

Die Sammlung mit Schlüsselworten soll manuell erstellbar und beliebig erweiterbar sein. Die Anwendung muss ohne aufwendige Zugriffe von diesen Listen lesen können. Als Möglichkeit dafür zählt eine CSV-Datei oder eine SQL-Datenbank.

Die SQL-Datenbank ist ein komplexeres System. Infolgedessen werden aufwendigere Zugriff benötigt. Die CSV-Datei bringt alle Anforderungen mit sich. Es ist unkompliziert, diese manuell zu befüllen und beliebig zu erweitern. Darüber hinaus kann eine CSV-Datei ohne großen Aufwand mit Hilfe eines Python-Skriptes ausgelesen werden. Die erwähnten Gründe sprechen für die Verwendung einer CSV-Datei.

Generierung der Wortsammlungen

Für eine sinnvolle Informationsgewinnung werden die Wortsammlungen kategorisiert. Dabei entsprechen die Kategorien einem Teil der zu suchenden Personenattributen, wie Darstellung 8.5 gezeigt wurde. Demzufolge gibt es die Kategorie “Tätigkeiten“, “Hobbys“, “Institutionen“ und “Städte und Gemeinden“. Dabei enthalten die Wortsammlungen möglichst alle Bezeichnungen und Namen dieser Kategorien. In Wortsammlung für Institution sind sowohl Firmennamen als auch Universitäts- und Hochschulnamen aufgelistet. Zur Veranschaulichung ist die Liste mit allen Städten und Gemeinden im Anhang beigefügt.

8.7.3 Automatic Keyword Extraction mit NLP

Durch das *Natural Language Toolkit* (NLTK) von Python ist es möglich, den vorhandenen Webseitentext zu analysieren.

Zu Beginn wird der vorhandene Text in einzelne Wörter zerlegt und in eine Liste gespeichert. Aus diesen Wörtern werden die *stopwords* der deutschen als auch der englischen Sprache herausgefiltert. Dadurch verringert sich die Anzahl der gesamten Wörter im Text um einen sehr großen Teil.

Im nächsten Schritt kann die Liste mit den entsprechenden Wortsammlungen verglichen werden. Die Wortsammlung, welche die möglichen Institutionen enthält, wird nicht mit

dieser Liste vergleichen. Die Problematik besteht darin, dass sich die Anzahl der Wörter für die Institutionen variieren kann. Für diesen bestimmten Fall, wird das Wort in dem Webseitentext ohne Fragmentierung und Formatierung gesucht. Dadurch kann bei dieser Suche ein Laufzeitnachteil entstehen, welcher aber nicht von Bedeutung ist.

8.7.4 Suche nach dem Geburtsjahr der Zielperson

Das Geburtsjahr ist für die Generierung der E-Mail wichtig. Viele Personen verwenden eine Kombination aus dem bürgerlichen Namen und dem Geburtsjahr als lokalen Teil der E-Mail-Adresse. Aus diesem Grund wird speziell nach dem Geburtsjahr in den generierten Schlüsselwörtern aus Kapitel 8.7.3 gesucht.

Dazu wird eine Suche nach einer vierstelligen Zahl, welche größer als 1900 und kleiner-gleich 2019 ist, durchgeführt. Beim Fund einer Zahl, werden fünfzehn Schlüsselwörter vor und hinter der vermutlichen Jahreszahl kontrolliert. Falls dabei das Wort "Geburtsdatum", "Alter", "geboren", "Geburtsort", "Geburtsstag", "born", oder "birth" vorkommt, wird das entsprechende Jahr als Geburtsjahr der Zielperson festgelegt. Die Implementierung zur Erkennung eines Geburtsjahrs ist in Listing 8.2 dargestellt.

```
regex_string = "(geburtsdatum)|(alter)|(geboren)|(geburtsort)|  
                (geburtsstag)|(born)|(birth)"  
for year in all_years_in_text:  
    vistited_elements = 0  
    max_number_of_visited_elements = 15  
    # to get all occurrences of this year  
    occurrences = [i for i, x in enumerate(keywords) if x == year]  
    for position_of_year in occurrences:  
        while (position_of_year+vistited_elements) < len(keywords)-1 and  
            vistited_elements <= max_number_of_visited_elements and  
            position_of_year is not -1:  
            index_behind = position_of_year+ vistited_elements  
            index_front = position_of_year - vistited_elements  
            if re.match(r"+regex_string", keywords[index_behind]):  
                print("Behind: Geburtsjahr wurde gefunden", year)  
                return year  
            elif re.match(r"+regex_string", keywords[index_front]):
```



```
        print("Front: Geburtsjahr wurde gefunden", year)
        return year
    visited_elements += 1
return -1
```

Listing 8.2: Algorithmus zur Suche nach dem Geburtsjahr

8.7.5 E-Mail-Adressen erkennen und herauslesen

Zu Beginn wird der unformatierte Webseitentext in Textfragmente zerlegt. Getrennt wird der Text bei einem Leerzeichen. Anschließend werden die erzeugten Fragmente mit einem regulären Ausdruck nach einer gültigen E-Mail-Adresse durchsucht. Bei einer Übereinstimmung des regulären Ausdrucks, wird der korrekte Teilstring, somit die E-Mail-Adresse, ausgelesen. Der Algorithmus zu diesem Vorgang ist in Listing 8.3 aufgezeigt.

```
for fragment in email_words:
    mail_regex = re.search('(. *((@)|(\(at\))) .*\. (de|com|net)) .*',
        fragment)
    if mail_regex:
        print("Email found:", mail_regex.group(1))
```

Listing 8.3: Teil des Algorithmuses zum Auslesen einer E-Mail-Adresse

Es werden nur die E-Mail-Adressen herausgesucht, welche einen Bezug zur Zielperson haben. Aus diesem Grund wird der lokale Teil aller gefundenen Adressen mit dem Vor- und Nachnamen der Zielperson verglichen. Mit Hilfe der “difflib” und dem implementierten “SequenceMatcher” von Python, lassen sich diese beide Sequenzen vergleichen und es wird ein prozentuale Übereinstimmung berechnet. Zur Differenzierung, ob eine E-Mail-Adresse eine Verbindung zum Opfer hat oder nicht, wird eine Prozent-Grenze bestimmt. Die Grenze wurde aus den Ergebnissen von zahlreichen Tests auf die Zahl 0,5 % festgelegt. Das folgende Beispiel soll die Methode zur Erkennung von korrekten E-Mail-Adressen verdeutlichen.

In diesem Beispiel heißt die Person “Max Mustermann“ und es werden zwei E-Mail-Adresse gefunden. Die erste Adresse lautet *MusterMax@gmail.com* und die zweite *MartaFrau@gmx.de*. Im ersten Schritt wird der Name “Max Mustermann“ zu einem String “maxmustermann“ umgewandelt. Im nächsten Schritt werden die lokalen Namen aus den E-Mail-Adressen herausgelesen und gleichzeitig in Kleinbuchstaben umgewandelt. In diesem Fall wäre das “mustermx“ und “marta frau“. Anschließend werden die lokalen Namen der E-Mail-Adressen mit dem erzeugten Namensstring der Zielperson verglichen. Dabei erreicht die lokale Namen *mustermx* eine prozentuale Übereinstimmung von 0,73 % mit dem Namensstring und *marta frau* 0,27 %. Da die Prozent-Grenze bei 0,5 % beträgt, wird die zweite E-Mail-Adresse verworfen.

8.7.6 Auswahl der gewonnenen Information

Die gefundenen Schlüsselwörter einer Webseite, werden in einer Liste gespeichert. Nachdem eine Seite vollständig durchsucht wurde, wird mit der Formel 8.1 eine prozentuale Wertung für das Vorkommen eines Wortes in der Liste berechnet.

$$\frac{\text{Vorkommen eines Wortes}}{\text{Anzahl aller gefundenen Wörter in der Liste}} \quad (8.1)$$

Die Schlüsselwörter werden anschließend mit dem dazugehörigen Score in einer neuen Liste gespeichert. Jedes Wort kommt dabei nur einmal vor. Eine beispielhafte Liste ist nachstehend dargestellt.

```
[['fussball', 0.7], ['basketball', 0.2], ['fechten', 0.1]]
```

Hierbei ist zu sehen, dass das Wort “Fußball“ sieben Mal öfter als das Wort “Fechten“ auf der Webseite vorgekommen ist. Für jede durchsuchte Seite wird solch eine Liste erstellt und anschließend zu einer großen Liste zusammengefügt. Dabei bleibt die Struktur bestehen, damit erkannt wird, welche Wörter von unterschiedlichen Webseiten kommen. Ein Beispiel hierfür ist die folgende Liste.

```
[[['fussball', 0.7], ['basketball', 0.2], ['fechten', 0.1]], [['fussball', 0.5], ['volleyball', 0.5]]
```

In dieser Liste befinden sich die gewonnenen Informationen aller Webseiten für eine Kategorie. Hier wäre es die Kategorie “Hobby“. Für jede dieser kategorisierten Listen, muss nun ein Element bestimmt werden, welches am wahrscheinlichsten eine Verbindung zu der Zielperson hat. Dazu wird die Formel 8.2 verwendet. Wobei beachtet werden muss, dass nur die Elemente summiert werden, bei denen das Schlüsselwort identisch ist.

$$\sum_{i=0}^{n-1} \sum_{j=0}^{m-1} \frac{liste[i][j][1]}{n}$$

mit

n = Anzahl der Elemente von Liste (8.2)

m = Anzahl der Elemente von Liste[i]

i = Element in Liste

j = Element von Liste[i]

In dem aufgezeigten Beispiel würde das zu der folgenden Liste führen.

```
[['fussball', 0.6], ['basketball', 0.1], ['fechten', 0.05], ['volleyball', 0.25]]
```

Aus dieser Liste kann nun das Schlüsselwort mit der höchsten Wertung gewählt und dem Personenobjekt hinzugefügt werden. In diesem Fall wäre dass das Wort “Fußball“. Falls zwei Wertungen gleich hoch sind, wird das erste Wort ausgewählt.

8.8 Methoden zum Erkennen einer Person

Bei jeder einzelnen Suche, besteht die Herausforderung darin, zu erkennen, wann es sich um die gesuchte Person handelt. Durch die große Anzahl an verfügbaren Informationen im Internet, besteht eine hohe Wahrscheinlichkeit, dass Personen mit sehr ähnlichen Profilen gefunden werden.

Aus diesem Grund werden Maßnahmen getroffen um die gesuchte Person zu erkennen. Dafür ist der erste Schritt die Anzahl der Suchergebnisse zu reduzieren. Dies ist durch den Ansatz der Personensuche im Kapitel 8.2 möglich. Dabei wird abhängig von der

einggegebenen Information die Suche variiert. Des Weiteren kann durch eine Optimierung des Such-URLs 8.4.2, die Personensuche verfeinert und somit die Ergebnisse verbessert werden. Durch diese Maßnahmen steigt die Wahrscheinlichkeit, dass es sich um die richtige Person handelt.

Im zweiten Schritt können die folgenden Methoden angewendet werden.

8.8.1 Identifikationsschlüssel verwenden

Bei der Personensuche wird mit Hilfe der eingegebenen Daten nach einer Person gesucht. Dabei können fehlerhafte Webseiten von Google vorgeschlagen werden. Fehlerhaft bedeutet hier, dass die Webseiten einen Inhalt repräsentieren, welcher nicht mit der gesuchten Person übereinstimmt.

Um dem entgegenzuwirken können bekannte Informationen als Identifikationsschlüssel verwendet werden. Allerdings müssen diese einzigartige Daten sein. Dazu zählt beispielsweise die E-Mail-Adresse oder Benutzernamen von den Plattformen Instagram und Twitter. Der vollständige Name ist nicht einzigartig und dient deswegen nicht als Identifikationsschlüssel. Dass bedeutet, dass es mehrere Personen mit dem selben vollständigen Namen geben kann.

Um eine Person zu identifizieren, zur welcher keine einzigartigen Informationen bekannt sind, können Kombinationen aus den angegebenen Daten erstellt werden. Diese Kombinationen dienen in dem Fall als Identifikationsschlüssel. Im folgenden sind alle möglichen Kombinationen aufgelistet.

Vorname, Nachname, Wohnort;

Vorname, Nachname, Geburtsjahr;

Vorname, Nachname, Institution;

Der Webseitentext kann anschließend auf das Vorkommen des Identifikationsschlüssels kontrolliert werden. Wenn der Text nur eine dieser Kombination beinhaltet, wird diese Seite für die Informationsgewinnung verwendet. Andernfalls wird die Webseite verworfen.

8.8.2 Kontaktanalyse

Hier kann die Suche erweitert werden, indem auf soziale und berufliche Verbindungen der Zielperson eingegangen wird. Das heißt, dass bekannte Kontakte der gesuchten Person ebenfalls durchsucht und ausgewertet werden. Als Kontaktquellen können Facebook-Freunden, FuPa-Teammitglieder, Instagram-Follower oder Xing-Kontakte dienen.

Durch die erwähnte Methode können weitere Informationen gewonnen werden. Diese sind zur Unterscheidung von Profilen nützlich.

8.9 Bewertung der Methoden zur Personenidentifizierung

Beide Methoden zur Identifizierung einer Person bringen eine Verbesserungen der Ergebnisse mit sich. Die Wahrscheinlichkeit wird erhöht, dass es sich um die korrekte Person handelt.

Die Methoden unterscheiden sich in der Wirksamkeit und in der Laufzeit. Durch die Verwendung von Identifikationsschlüsseln wird die Anzahl von Fehlinformationen in dem Profil der gesuchten Person reduziert. Allerdings können gleichzeitig wichtige Informationsquellen ignoriert werden, wenn diese den Kriterien nicht entsprechen. Bei der Kontaktanalyse werden jedoch keine Informationsquellen ignoriert. Es werden weitere Informationen gesammelt. Diese sind zusätzlich zur E-Mail-Generierung von Vorteil. Das Ergebnis bei der Verwendung der Kontaktanalyse ist nicht optimal. Es kann nicht davon ausgegangen werden dass das Ergebnis für unmittelbar für die Person spricht. Beim betrachten der Laufzeit, kann davon ausgegangen werden, dass die Kontaktanalyse deutlich mehr Zeit und Ressourcen benötigt.

Es werden beide Methoden umgesetzt, da sie einen positiven Effekt auf die Anwendung haben.

8.10 Implementierung der Personenidentifizierung

8.10.1 Identifikationsschlüssel verwenden

Zu Beginn der vorläufigen Inhaltskontrolle werden die Eingaben abgefragt. Dadurch wird erkannt, zu welchen Daten Informationen vom Benutzer eingegeben wurden. Anschließend werden mit diesen Daten alle möglichen Kombinationen aus Kapitel 8.8.1 erstellt. Es sind allerdings nur die Kombinationen möglich, für die die Daten bekannt sind.

Für die Suche des Vornamen und Nachnamen wird ein String erzeugt, der beide Attribute kleingeschrieben beinhaltet. Ein korrekter String ist "max mustermann". Infolgedessen wird der Webseitentext zu einem String umgewandelt. Anschließend wird kontrolliert, ob sich der String bestehend aus Vornamen und Nachnamen und das entsprechende Attribut, beispielsweise der Wohnort, in dem Webseitentext befindet. Wenn diese Abfrage korrekt ist, wird die Webseite weiter behandelt und es kann nach Information gesucht werden.

8.10.2 Kontaktanalyse

Welche Seiten eignen sich zur Kontaktanalyse?

Diese Methode funktioniert auf der Webseite LinkedIn nicht. Es gibt dort keine Möglichkeit, die Kontakte der gesuchten Person anzuzeigen. Bei Xing kann ein Nutzer einstellen, ob diese Kontaktanzeige freigegeben wird oder nicht. Dadurch sind die Kontakte bei vielen Usern nicht erkennbar. Facebook, Twitter und Instagram bieten die Möglichkeit, die Kontakte der gesuchten Person anzuzeigen. Allerdings wird dafür ein Account benötigt. Für diese Methode eignen sich somit die Seiten Twitter, Xing, Facebook und Instagram. Wie in Kapitel 8.4.4 beschrieben, wird für Facebook kein Account angelegt. Dadurch ist es nicht möglich, Kontakte auf dieser Webseite anzuzeigen. Um die Funktion der Methode aufzuzeigen, wird ausschließlich die Webseite Instagram verwendet.

Instagram Kontakte durchsuchen

Zuallererst wird unterschieden, ob das Profil der gesuchten Person privat oder öffentlich ist. Bei einem öffentlichen Profil, können alle Abonnenten und abonnierte Profile angezeigt werden. Die Abonnenten und abonnierte Profile können sich unterscheiden. Im Gegensatz dazu, werden bei einem privaten Profil, nur eine begrenzte Anzahl von Profilen vorgeschlagen. Des Weiteren kann bei einem privaten Profil nicht unterschieden werden, ob die Abonnenten oder die abonnierten Profile angezeigt werden sollen.

Von den gefunden Followern wird jedes einzelne Profil durchsucht, bis eine Übereinstimmung mit der Zielperson gefunden wurde. Eine Übereinstimmung bedeutet, dass auf diesem Profil ein Teil mit dem Opferprofil identisch ist. Beispielsweise kann das die selbe Universität oder der selbe Wohnort sein. Sobald dies gefunden wurde, kann die Suche beendet werden. Wenn keine Profilinformation übereinstimmt, wird nicht ausgeschlossen, dass es sich trotzdem um die gesuchte Person handelt.

Ein Fund einer identischen Information ist keine vollständiger Beweis, dass es sich um die richtige Person handelt. Allerdings erhöht sich die Wahrscheinlichkeit für die Aussage, dass es sich bei diesem Profil um die gesuchte Person handelt.

Wie werden Kontakte ausgelesen

Im ersten Schritt entscheidet der Algorithmus, ob es sich um ein privates oder öffentliches Profil handelt. Dies wird realisiert, indem nach einem String auf der Webseite gesucht wird. Der String lautet "Diese Konto ist privat". Wenn diese Zeichenfolge gefunden wird, handelt es sich um ein privates Konto. Andernfalls um ein öffentliches.

Damit die Links zu den Kontakt-Profilseiten auf einer privaten Seite herausgelesen werden können, wird ein scrollbarer Container ausgelesen. Dieser Container beinhaltet die vorgeschlagenen Kontakte und zwei Buttons. Wie im Bild 8.5 zu sehen, kann mit den beiden Buttons nach rechts und links gewischt werden. Sobald die Links zu den aktuell angezeigten Profilen ausgelesen wurden, wird auf den rechten Button geklickt. Dies wird mit einem vorgetäuschten Mausklick des Selenium WebDrivers realisiert. Durch diese Schritt-für-Schritt-Methode können alle vorgeschlagenen Kontakte ausgelesen werden. Andernfalls werden nur die aktuell angezeigten Profile geladen und gefunden.

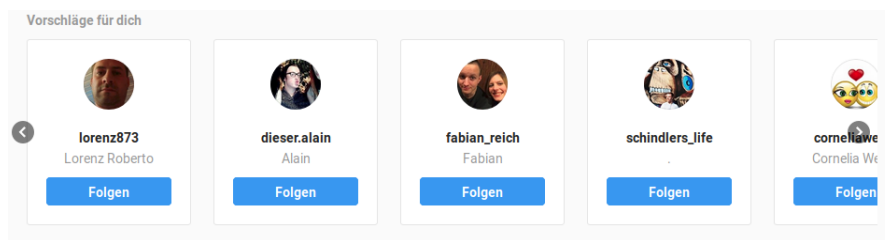


Bild 8.5: Container mit Profil-Vorschlägen

Falls es sich um ein öffentlich frei zugängliches Profil handelt, kann eine Liste der abonnierten Kontakte angezeigt werden. Hierbei handelt es sich um ein scrollbares Pop-Up-Fenster 8.6. Vergleichbar zur Methode bei einer privaten Profilseite, wird hier ebenfalls Schritt-für-Schritt durchgescrollt. Dadurch wird jedes einzelne Profil geladen und der dazugehörig Link, zur dieser Profilseite, kann dadurch ausgelesen werden.

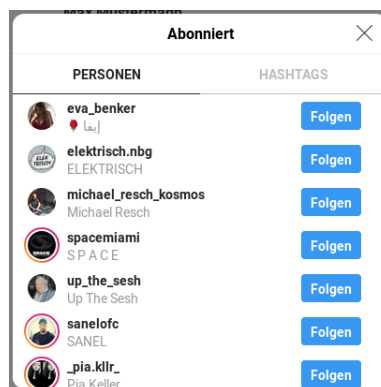


Bild 8.6: Pop-up-Fenster mit abonnierten Profilen

Die Herausforderung besteht darin, dass nicht zu schnell gescrollt werden darf. Aus diesem Grund wird ein Algorithmus 8.4 verwendet, welcher einem menschlichen Verhalten ähneln soll. Hierbei wird zuallererst das Pop-up-Fenster gesucht und festgelegt. Anschließend wird die Anzahl der abonnierten Profile gezählt. Die Anzahl der Profile wird dazu verwendet, dass der Algorithmus weiß, wie weit nach unten geblättert werden muss, um alle Profile zu laden.

Im ersten Schritt wird das Fenster nur ein sechstel des möglichen Bereichs nach unten gescrollt. Dadurch werden weitere Profile geladen. Wenn direkt nach ganz unten geblättert

wird, wäre dies beim ersten Scrollvorgang zu schnell. In diesem Fall werden keine Kontakte geladen. Es werden lediglich Profil von sehr bekannten Instagram-Usern angezeigt, die von dieser Person abonniert wurden.

In den nächsten Schritte wird das Fenster jeweils ganz nach unten verschoben. Dadurch werden alle Profile geladen. Sobald alle Links bekannt sind, wird eine URL zu den entsprechenden Profilseiten erstellt. Diese Seiten werden anschließend wie jede andere Seite ausgelesen und nach Information durchsucht. Infolgedessen wird die gewonnene Information jedes Profils mit der Information der Zielperson verglichen. Die Suche wird bei einem beliebigen Treffer beendet. Anschließend wird die gefunden Information mit dem Namen des Benutzers der Profilseite gespeichert. Diese abgespeicherten Daten können später zur E-Mail-Generierung verwendet werden.

```
# Find the pop-up window
pop_up = self.browser.find_element_by_xpath('/html/body/div[2]
      /div/div[2]')
# find number of followers
all_following = int(self.browser.find_element_by_xpath("//li[2]
      /a/span").text)
# scroll down the page
for i in range(int(all_following / 6)):
    if i == 0:
        self.browser.execute_script("arguments[0].scrollTop =
            arguments[0].scrollHeight/5", pop_up)
        time.sleep(2)
    else:
        self.browser.execute_script("arguments[0].scrollTop =
            arguments[0].scrollHeight", pop_up)
        time.sleep(random.randint(500, 1000) / 1000)
```

Listing 8.4: Herunterscrollen des Pop-up Fensters

8.11 Speicherung der gewonnenen Daten

Die gespeicherten Daten werden von verschiedenen Klassen benötigt. Aus diesem Grund muss es möglich sein, dass andere Klassen auf die Speicherstruktur zugreifen können. Des Weiteren wird eine gute Struktur vorausgesetzt, damit auf einzelne Attribute der Person zugegriffen werden kann. Es ist nicht notwendig, dass die Daten nach Programmende abrufbar sind. Infolgedessen wird keine externe Speicherung in einer Datenbank oder in einer Datei vorausgesetzt.

Eine mögliche Speicherung der Daten wäre in einer SQL-Datenbank. Alternativ könnten die Personendaten in einer externen Datei oder mit Hilfe einem Personenobjekt gespeichert werden.

Eine SQL-Datenbank bringt eine gute Speicherstruktur mit sich. Allerdings sind mit einer SQL-Datenbank aufwendigere Speicher- und Abrufvorgänge verbunden. Die externe Speicherung in einer Datei wie CSV oder TXT ist keine Anforderung. Auf eine SQL-Datenbank sowie auf eine externe Datei, lässt sich mit jeder Klasse darauf zugegriffen. Dennoch wird ein Personenobjekt verwendet. Dadurch werden keine unnötigen Speicher- und Lesezugriffe benötigt. Darüber hinaus lässt sich das Personenobjekt an die entsprechenden Klassen übergeben.

8.11.1 Implementierung der Personenklasse

Die gewonnenen Daten werden in einem Personenobjekt, wie in Listing 8.5 dargestellt, gespeichert. Dabei werden die vom Anwender eingegebenen Daten direkt in das Personenobjekt übertragen. Die gewonnenen Informationen werden in Form einer Liste hinzugefügt. Falls bei der Kontaktanalyse ein Treffer gemacht wurde, kann diese Information in dem Attribut "contacts_information" gespeichert werden. Hierbei wird zuerst der vollständige Kontaktname und anschließend die übereinstimmende Information gespeichert. Ein Beispiel hierfür wäre ["Max Mustermann", "Fußball"]

```
class Person(object):  
    def __init__(self):  
        self.first_name = input("Vorname: ")
```

```
self.second_name = input("Nachname: ")
self.place_of_residence = input("Wohnort: ")
self.year_of_birth = input("Geburtsjahr: ")
self.institution = input("Institution: ")
self.instagram_name = input("Instagram Benutzername: ")
self.facebook_name = input("Facebook Benutzername: ")
self.twitter_name = input("Twitter Benutzername: ")
self.input_email = input("E-Mail-Adresse: ")
self.occupation = []
self.hobbies = []
self.universities = []
self.founded_mails = []
self.locations = []
self.contacts_information = []
```

Listing 8.5: Personklasse

9 Generierung der Phishing-E-Mail

In diesem Kapitel wird die Umsetzung zur Erstellung einer Phishing-Mail beschrieben.

9.1 Implementierung der Methode zur Generierung der E-Mail-Adressen

Für den zu entwickelnden Algorithmus wird eine eigene Klasse erstellt. Diese Klasse ist ausschließlich für die Generierung der E-Mail-Adressen zuständig. Diese Methode zeigt, wie aus personenbezogenen Daten eine mögliche E-Mail-Adresse zu dieser Person generiert werden kann. Jedoch wird in der Anwendung keine Phishing-Mail an eine dieser Adresse, welche mit der Methode erzeugt wurden, versendet. Die Ausgaben dienen lediglich zum aufzeigen der Funktion dieser Methode.

9.1.1 Funktion des eigenen Algorithmus

Der lokale Teil einer E-Mail-Adresse befindet sich vor dem At-Zeichen. Dieser kann aus verschiedensten Daten bestehen. Allerdings wird in den meisten Fällen der bürgerlichen Namen verwendet. [Med17] Aus diesem Grund verwendet der Algorithmus die Personenattribute Vorname, Nachname und das Geburtsjahr.

Im ersten Schritt wird kontrolliert, welche Daten bekannt sind. Im Idealfall sind das alle drei Attribute. Im zweiten Schritt wird festgelegt aus welchen Daten der lokale Teil bestehen kann. Im Folgenden sind möglichen Kombinationen aufgezeigt.

Vorname;

Nachname;

Vorname, Nachname;

Vorname, Nachname, vollständiges Geburtsjahr;

Vorname, Nachname, Kurzform von Geburtsjahr;

Ein lokaler Teil kann somit aus mehreren Daten bestehen. Es kann vorkommen, dass anstatt “Max Mustermann” “Mustermann Max” als lokaler Namen verwendet wird. Aus diesem Grund wird für jeden lokalen Teil, der aus mehreren Daten besteht, eine Permutation ohne Wiederholung angewendet. Dadurch werden alle möglichen Kombinationen aus den Daten gewonnen, da bei der Zusammensetzung der Daten zusätzlich auf die Reihenfolge geachtet wird. Außerdem werden bei der Zusammensetzung der Daten die bekannten Trennzeichen “.”, “_” und “-” hinzugefügt. Jedoch gibt es ebenfalls jede Kombination ohne Trennzeichen. Die lokale Namen werden anschließend in einer Liste gespeichert.

Für den Domainteil werden die bekannte Mailprovider in Deutschland verwendet. Dazu gehören die Provider GMX, WEB.DE, Gmail, T-Online, Freenet und 1&1. [Anb19]. Das bedeutet, es wird für jeden lokalen Namen eine E-Mail-Adresse mit den jeweiligen Mail Providern und der Landeskenntung “de” erzeugt. Die folgende Tabelle zeigt die erzeugten E-Mail-Adressen des Algorithmus für die Daten “Marco”, “Lang” und “1995”. Es sind allerdings nur die Mailadressen für die Provider WEB.DE, Gmail und Freenet aufgelistet.

marco@web.de	marco@gmail.com	marco@freenet.de
lang@web.de	lang@gmail.com	lang@freenet.de
marcolang@web.de	marcolang@gmail.com	marcolang@freenet.de
marco.lang@web.de	marco.lang@gmail.com	marco.lang@freenet.de
marco_lang@web.de	marco_lang@gmail.com	marco_lang@freenet.de
marco-lang@web.de	marco-lang@gmail.com	marco-lang@freenet.de
langmarco@web.de	langmarco@gmail.com	langmarco@freenet.de
lang.marco@web.de	lang.marco@gmail.com	lang.marco@freenet.de
lang_marco@web.de	lang_marco@gmail.com	lang_marco@freenet.de
lang-marco@web.de	lang-marco@gmail.com	lang-marco@freenet.de
marcolang1995@web.de	marcolang1995@gmail.com	marcolang1995@freenet.de
marco.lang.1995@web.de	marco.lang.1995@gmail.com	marco.lang.1995@freenet.de
marco_lang_1995@web.de	marco_lang_1995@gmail.com	marco_lang_1995@freenet.de
marco-lang-1995@web.de	marco-lang-1995@gmail.com	marco-lang-1995@freenet.de
marco1995lang@web.de	marco1995lang@gmail.com	marco1995lang@freenet.de
marco.1995.lang@web.de	marco.1995.lang@gmail.com	marco.1995.lang@freenet.de
marco_1995_lang@web.de	marco_1995_lang@gmail.com	marco_1995_lang@freenet.de
marco-1995-lang@web.de	marco-1995-lang@gmail.com	marco-1995-lang@freenet.de
langmarco1995@web.de	langmarco1995@gmail.com	langmarco1995@freenet.de

lang.marco.1995@web.de	lang.marco.1995@gmail.com	lang.marco.1995@freenet.de
lang_marco_1995@web.de	lang_marco_1995@gmail.com	lang_marco_1995@freenet.de
lang-marco-1995@web.de	lang-marco-1995@gmail.com	lang-marco-1995@freenet.de
lang1995marco@web.de	lang1995marco@gmail.com	lang1995marco@freenet.de
lang.1995.marco@web.de	lang.1995.marco@gmail.com	lang.1995.marco@freenet.de
lang_1995_marco@web.de	lang_1995_marco@gmail.com	lang_1995_marco@freenet.de
lang-1995-marco@web.de	lang-1995-marco@gmail.com	lang-1995-marco@freenet.de
1995marcolang@web.de	1995marcolang@gmail.com	1995marcolang@freenet.de
1995.marco.lang@web.de	1995.marco.lang@gmail.com	1995.marco.lang@freenet.de
1995_marco_lang@web.de	1995_marco_lang@gmail.com	1995_marco_lang@freenet.de
1995-marco-lang@web.de	1995-marco-lang@gmail.com	1995-marco-lang@freenet.de
1995langmarco@web.de	1995langmarco@gmail.com	1995langmarco@freenet.de
1995.lang.marco@web.de	1995.lang.marco@gmail.com	1995.lang.marco@freenet.de
1995_lang_marco@web.de	1995_lang_marco@gmail.com	1995_lang_marco@freenet.de
1995-lang-marco@web.de	1995-lang-marco@gmail.com	1995-lang-marco@freenet.de
marcolang95@web.de	marcolang95@gmail.com	marcolang95@freenet.de
marco.lang.95@web.de	marco.lang.95@gmail.com	marco.lang.95@freenet.de
marco_lang_95@web.de	marco_lang_95@gmail.com	marco_lang_95@freenet.de
marco-lang-95@web.de	marco-lang-95@gmail.com	marco-lang-95@freenet.de
marco95lang@web.de	marco95lang@gmail.com	marco95lang@freenet.de
marco.95.lang@web.de	marco.95.lang@gmail.com	marco.95.lang@freenet.de
marco_95_lang@web.de	marco_95_lang@gmail.com	marco_95_lang@freenet.de
marco-95-lang@web.de	marco-95-lang@gmail.com	marco-95-lang@freenet.de
langmarco95@web.de	langmarco95@gmail.com	langmarco95@freenet.de
lang.marco.95@web.de	lang.marco.95@gmail.com	lang.marco.95@freenet.de
lang_marco_95@web.de	lang_marco_95@gmail.com	lang_marco_95@freenet.de
lang-marco-95@web.de	lang-marco-95@gmail.com	lang-marco-95@freenet.de
lang95marco@web.de	lang95marco@gmail.com	lang95marco@freenet.de
lang.95.marco@web.de	lang.95.marco@gmail.com	lang.95.marco@freenet.de
lang_95_marco@web.de	lang_95_marco@gmail.com	lang_95_marco@freenet.de
lang-95-marco@web.de	lang-95-marco@gmail.com	lang-95-marco@freenet.de
95marcolang@web.de	95marcolang@gmail.com	95marcolang@freenet.de
95.marco.lang@web.de	95.marco.lang@gmail.com	95.marco.lang@freenet.de
95_marco_lang@web.de	95_marco_lang@gmail.com	95_marco_lang@freenet.de
95-marco-lang@web.de	95-marco-lang@gmail.com	95-marco-lang@freenet.de
95langmarco@web.de	95langmarco@gmail.com	95langmarco@freenet.de
95.lang.marco@web.de	95.lang.marco@gmail.com	95.lang.marco@freenet.de
95_lang_marco@web.de	95_lang_marco@gmail.com	95_lang_marco@freenet.de
95-lang-marco@web.de	95-lang-marco@gmail.com	95-lang-marco@freenet.de

9.2 Implementierung der E-Mail-Muster

Ein E-Mail-Muster entspricht einem Lückentext, bei dem die entsprechenden Lücken mit den gewonnenen Daten ergänzt werden. Die Texte müssen so erstellt werden, dass sie die Zielperson ansprechen. Aus diesem Grund, muss für jede Kombination der gewonnenen Daten ein Muster zur Verfügung stehen. Infolgedessen, stellt sich die Frage, wie die E-Mail-Texte möglichst passend kategorisiert werden können.

9.2.1 Kategorien erstellen

Die Muster können in zwei große Kategorien unterteilt werden. Es gibt eine private und eine berufliche Kategorie. Der Unterschied zwischen privat und beruflich, besteht in der Art und Weise wie ein Text geschrieben wird. Genaugenommen bedeutet das, dass ein privates Muster in einer Alltagssprache und ein berufliches in einer formelleren Sprache erstellt wird. Diese beiden Kategorien haben weitere Unterkategorien, welche verschiedene Kombinationen aus den personenbezogenen Daten verwenden.

Um die Kategorie zu erkennen, werden zu Beginn Abfragen gestartet. Dadurch wird kontrolliert, welche Daten bekannt sind. Im Fall, dass die Institution oder die Tätigkeit der Zielperson bekannt ist, wird ein berufliches Muster gewählt. Wenn keines dieser Attribute bekannt ist, wird ein privates Muster verwendet.

Berufliche E-Mail-Muster

Die beruflichen E-Mail-Muster sind in den folgenden Bildern aufgezeigt. Die Reihenfolge zur Auswahl der Muster im laufenden Programm entspricht der Reihenfolge der Bilder. Dabei werden im normalen Anwendungsablauf die kursiv geschriebenen Wörter mit den gewonnenen Daten über die Zielperson ersetzt.

Das Bild 9.2.1 beschreibt ein berufliches Muster, welches die Personenattribute Nachname und Institution verwendet. Im Bild 9.2.1 ist ein Muster, mit den Attributen Nachname und Tätigkeit, aufgezeigt. Das Bild 9.2.1 zeigt das Muster mit den selben Attributen, wobei die Tätigkeit zu Beginn abgefragt wurde. Infolgedessen kann dieses Datenelement als Entscheidungskriterium für die Auswahl eines Musters verwendet werden. In diesem

bestimmten Fall wurde das Muster mit dem festgesetzten Wort “Professor“ gewählt. Im letzten Bild wird das berufliche Muster für die Daten Nachname, Tätigkeit und Institution beschrieben.

SUBJECT: *Musterinstitution* - Netzwerkänderungen

Hallo Herr *Mustermann*,
wir bauen unsere Netzwerkstruktur um. Bitte registrieren Sie sich unter der folgenden Webseite, damit wir Sie in das neue System aufnehmen können.

<https://badlink.com>

Mit freundlichen Grüßen

Ihr IT-Team der *Musterinstitution*

Bild 9.1: Ein berufliches E-Mail-Muster mit den Personenattributen Nachname und Institution

SUBJECT: *Mustertätigkeit* bei der ZF Friedrichshafen AG gesucht

Hallo Herr *Mustermann*,
wir, die ZF Friedrichshafen AG suchen einen kompetenten *Mustertätigkeit*. Im Anhang befindet sich die Stellenausschreibung mit allen Anforderungen und den vorstellbaren Gehaltsstufen.

Ihr Karriere Team der ZF Friedrichshafen AG

Bild 9.2: Ein berufliches E-Mail-Muster mit den Personenattributen Nachname und Tätigkeit

Private E-Mail-Muster

Im nachfolgenden werden die privaten E-Mail-Muster aufgezeigt. Die Ersetzung der kursiven Wörter, sowie die Reihenfolge der Muster-Auswahl ist identisch zum vorherigen Kapitel 9.2.1.

Im ersten Bild wird ein privates E-Mail-Muster beschrieben, welche die gewonnenen Kontaktinformationen verwendet. Genaugenommen ist das der Kontaktnamen und die

SUBJECT: Feedback zur Ausarbeitung

Hallo Herr Professor *Mustermann*,
wie besprochen befindet sich im Anhang meine vorläufige Ausarbeitung.
Könnten Sie diese erneut überprüfen und mir ein Feedback geben?

Mit freundlichen Grüßen

Max Mustermann

Bild 9.3: Ein berufliches E-Mail-Muster mit dem Personenattribut Nachname und einer festgesetzten Tätigkeit

SUBJECT: *Mustertätigkeit* bei der *Musterinstitution*

Hallo Herr *Mustermann*,
als *Mustertätigkeit* bei der *Musterinstitution*, stehen Ihnen nun alle
Möglichkeiten offen. Sehen Sie nun Ihre neuen Möglichkeiten unter folgen-
dem Link an.
<https://badlink.com>

Mit freundlichen Grüßen

Ihr Team der *Musterinstitution*

Bild 9.4: Ein berufliches E-Mail-Muster mit den Personenattributen Nachname, Tätigkeit und Institution

Kontaktinformation. Diese Kontaktinformation ist identisch zu einer Information über die Zielperson. Das nächste Bild beschreibt ein Muster, bei dem die Personendaten Vorname und Hobby verwendet werden. Das Bild 9.2.1 zeigt die Verwendung des Vornamens und Geburtsjahrs in einem E-Mail-Muster. Das letzte muss verwendet die Personenattribute Vorname und Ort. Dabei wird hierfür der gefundene Ort verwendet.

SUBJECT: Fragen bzgl. *Kontaktinformation*

Hi *Max*,
hier ist *Kontaktname*. Bezüglich *Kontaktinformation* hätte ich noch ein paar fragen an dich...
Könntest du zufällig in den Anhang schauen und bewerten was ich dazu so rausgesucht habe?
Vielen Danke im Voraus!

Kontaktname

Bild 9.5: Ein privates E-Mail-Muster mit den Personenattributen Vorname, Kontaktname und Kontaktinformation

SUBJECT: Verbessere deine Technik im *Musterhobby*

Hi *Max*,
damit du deine Leistung im *Musterhobby* verbessern kannst, musst du unbedingt die Techniken deiner Vorbilder anschauen!
Im Anhang befindet sich darüber eine kleine Übersicht.

Dein Team der deutschen Förderung

Bild 9.6: Ein privates E-Mail-Muster mit den Personenattributen Vorname und Hobby

SUBJECT: Jahrgang *Geburtsjahr*

Hi *Max*,
dieses Jahr findet ein Treffen für alle Personen, die Geburtsjahr geboren sind, statt.
Im Anhang befindet sich eine Liste mit den Leuten die bereits zugesagt haben.

Dein Orga-Team

Bild 9.7: Ein privates E-Mail-Muster mit den Personenattributen Vorname und Jahrgang

SUBJECT: Streetfood-Festival in *Musterort*

Hi *Max*,
dieses Jahr findet das erste STREETFOOD-FESTIVAL in *Musterort* statt. Im
Anhang befindet sich der Plan, auf dem alles weitere erklärt wird.
Wir freuen uns auf dich!

Dein Streefood-Team aus *Musterort*

Bild 9.8: Ein privates E-Mail-Muster mit den Personenattributen Vorname und Ort

9.3 Versenden einer Phishing-E-Mail

Damit eine Phishing-Mail beispielhaft versendet werden kann, wird eine Sender-E-Mail-Adresse benötigt. Sehr große Provider wie Gmail oder Yahoo sind dafür nicht optimal. Das hat den Grund, dass viele Spammer diese Provider verwendet haben. Dadurch werden diese Adressen gerne öfter überprüft. Kleine Provider wie GMX sind dagegen nicht so bekannt. Ein weiterer Vorteil von GMX ist, dass ein Account ohne eine weitere gültige E-Mail-Adresse erzeugt werden kann. Das spricht für GMX, da bei einem gefälschten Account keine Verbindung zu einem benutzten Account oder zu einer Person besteht. Des Weiteren ist dieser kleine Provider großen Plattformen wie Facebook fremd und wird dadurch weniger streng überprüft. [Baz18] Aus den erwähnten Gründen wird eine E-Mail-Adresse bei GMX erstellt.

Mit Hilfe der erzeugten Adresse und einem Python Skript kann eine Phishing-Mail beispielhaft versendet werden. Dazu wird die Python Bibliothek “smtplib“ verwendet. Im ersten Schritt werden die erstellten E-Mail-Daten wie Sender-Adresse, Ziel-Adresse, Betreff und Inhalt einer mehrteiligen Nachricht hinzugefügt. Im nächsten Schritt wird die Verbindung mit dem GMX-Mailserver aufgebaut. Nach einer erfolgreichen Anmeldung kann eine E-Mail versendet werden. In der Darstellung 9.1 ist der dazugehörige Programmcode aufgezeigt.

```
msg = multipart.MIMEMultipart()
msg['From'] = source_email
msg['To'] = destination_email
msg['Subject'] = subject

emailText = email_text
msg.attach(text.MIMEText(emailText, 'html'))

server = smtplib.SMTP('mail.gmx.net', 587)
server.ehlo()
server.starttls()
server.login(source_email, password)
text = msg.as_string()
server.sendmail(source_email, destination_email, text)
```

```
server.quit()
```

Listing 9.1: Implementierung des versenden einer Phishing-Mail

10 Evaluation der Implementation

11 Schlussbemerkungen und Ausblick

11.1 Wie kann eine Person weiter identifiziert werden?

Durch die Google Bildersuche ist es möglich, anstatt einem Suchbegriff ein Bild zu verwenden und nach diesem zu suchen. Dabei kann ein zu suchendes Bild selbst hochgeladen oder ein URL angegeben werden. Bei dem Ergebnis kann es sich um ein ähnliches Bild oder eine Webseite, die das Bild enthält, handeln.

Als Alternative zur Google-Bildersuche kann eine Bilderkennungssoftware verwendet werden um Personen zu identifizieren bzw. zu unterscheiden.

11.1.1 Zeitrahmen wird mit Beachtet

Wie kann Alter der Webseite herausgefunden werden

Der Webseitentext kann nach Datums suchen und diese mit dem angegebenen Geburtsjahr verglichen werden. Dabei kann erkannt werden, ob das theoretische Alter des Artikels mit dem Alter der Person übereinstimmen kann. Möglicherweise können Metadaten von der Webseite ausgelesen werden. möglicherweise über domain

Bereits umgesetzt

Jahr nach copyright wird ausgelesen, wenn das nicht vorhanden werden alle Jahreszahlen genommen und ein durchschnitt ausgerechnet

11.1.2 Zeitraum beachten

Eine Methode für das Erkennen von Personen kann das Beachten von Zeiträumen sein. Dabei fließt das Alter der Zielperson mit in die Suche ein. Das bedeutet, dass nach dem Alter der Webseite gesucht wird, indem Jahreszahlen aus dem Webseitentext ausgelesen werden. Dadurch wird erkannt, ob der Zeitrahmen des Artikels oder das Erstellungsdatum einer Webseite mit dem Alter der Person grundsätzlich übereinstimmt.

11.2 Adressgenerierung

11.2.1 Wenn Firma bekannt

11.3 Keyword Extraction mit Hilfe von Machine Learning

In der Theorie ist es möglich, ein Neuronales Netz mit den Begriffen zu trainieren und eine Kategorisierung durchzuführen. Dabei entsteht ein Netz, welches selbst entscheiden würde, in welche Kategorie ein Wort fällt. Das Wort "Fußball" müsste dadurch in die Kategorie Hobby eingeordnet werden.

11.4 Wie werden Wortsammlungen am effektivsten verglichen?

Es gibt keine Methode die den Vergleich der Schlüsselwörter mit den Wörtern der Datenbanken verbessern kann, da jedes einzelne Schlüsselwort mit jedem einzelnen Wort aus der Datenbank verglichen werden muss. Suchalgorithmen

<https://softwareengineering.stackexchange.com/questions/280361/list-comparing-techniques-for-faster-performance>

11.5 Email-adressen

Adressen von Micheal Bazzel mit verwenden wie ml@web.de

11.6 Absender-Adresse

Spoofing, Kontakte von Fupa oder Instagram nutzen.

11.7 Validität der generierten Mail-Adressen prüfen

11.7.1 Methoden zum Prüfen der Validität

Die erzeugten Adressen werden anschließend auf Validität geprüft. Hierfür gab es früher eine *VERFY* Anfrage von SMTP. Mit dieser Anfrage konnte eine angegebene E-Mail-Adresse überprüft werden. Allerdings wurde der Dienst von Spammern ausgenutzt und wird dadurch von den meisten SMTP-Servern nicht mehr zu Verfügung gestellt. [BPH⁺10] Demnach muss die Validität auf einem anderen Weg geprüft werden. Eine Möglichkeit zur Prüfung ist die Verwendung bereitgestellter Webseiten, bei der die zu prüfenden E-Mail-Adresse angegeben werden kann. Eine anschließende Rückmeldung verrät dann, ob die Adresse verwendet wird oder nicht. Eine Webseite dafür wäre "<https://centralops.net/co/>". Als Alternative dazu, ist die Entwicklung eines Skriptes, welches die Validität der Adresse prüft.

Im Fall, dass mehrere Adressen von diesem Adresspool gültig sind, kann nach mit Hilfe dieser Mail-Adressen nach Einträgen im Internet gesucht werden. Wenn es eine Übereinstimmung mit der Zielperson gibt, wird diese E-Mail ausgewählt. Andernfalls wird an jede gültige Adresse eine Phishing-Mail gesendet.

11.7.2 Bewertung: Validität Prüfen

Für eine bessere Laufzeit des Programms, wird ein Skript zur Überprüfung der Adressen auf Verfügbarkeit und Gültigkeit, verwendet.

A Ein Kapitel des Anhangs

Literatur

- [AH19] ADS-HILFE, GOOGLE: *URL-Parameter*. <https://support.google.com/google-ads/answer/6277564?hl=de>, 2019. Abrufdatum: 26.02.2019.
- [All18] ALLENSBACH, IFD: *Meistgenutzte Informationsquellen der Bevoelkerung in Deutschland im Jahr 2018*. <https://de.statista.com/statistik/daten/studie/171257/umfrage/normalerweise-genutzte-quelle-fuer-informationen/>, 2018. Abrufdatum: 18.01.2019.
- [Anb19] *Bei welchem Anbieter haben Sie Ihr Haupt-E-Mail-Postfach?* <https://de.statista.com/statistik/daten/studie/170371/umfrage/nutzung-von-e-mail-domains/>, 2019. Abrufdatum: 04.02.2019.
- [Ang18] *Haben Sie gro Angst davor, dass Sie Opfer von Datendiebstahl im Internet, also der missbrhlichen Verwendung Ihrer persnlichen Daten durch Dritte, werden?* <https://de.statista.com/statistik/daten/studie/886892/umfrage/angst-vor-einem-datendiebstahl-im-internet-in-deutschland/>, 2018. Abrufdatum: 22.02.2019.
- [Baz] BAZZELL, MICHAEL: *Email Assumptions*. <https://inteltechniques.com/osint/email.html>. Abrufdatum: 01.02.2019.
- [Baz18] BAZZELL, MICHAEL: *Open Source Intelligence Techniques: Resources for Searching and Analyzing Online Information*. CreateSpace Independent Publishing Platform, USA, 6th , 2018.
- [BKL09] BIRD, STEVEN, EWAN KLEIN EDWARD LOPER: *Natural language processing with Python: analyzing text with the natural language toolkit*. Ö'Reilly Media, Inc., 2009.
- [Boh14] BOHNENSTEFFEN, MARCEL: *Die alternativlose Suchmaschine*. <https://www.handelsblatt.com/unternehmen/it-medien/google-die-alternativlose-suchmaschine/11061626-all.html>, 2014. Abrufdatum: 24.02.2019.

- [BPH⁺10] BALDUZZI, MARCO, CHRISTIAN PLATZER, THORSTEN HOLZ, ENGIN KIRDA, DAVIDE BALZAROTTI CHRISTOPHER KRUEGEL: *Abusing social networks for automated user profiling. International Workshop on Recent Advances in Intrusion Detection*, 422–441. Springer, 2010.
- [Cal13] CALDWELL, TRACEY: *Spear-phishing: how to spot and mitigate the menace. Computer Fraud & Security*, 2013(1):11–16, 2013.
- [CH15] CHRISTOPHER HADNAGY, MICHELE FINCHER: *Phishing Dark Waters: The Offensive and Defensive Sides of Malicious E-mails*. 2015.
- [dev18] DEVELOPERS, SCRAPY: *Scrapy at a glance*. <http://doc.scrapy.org/en/latest/intro/overview.html>, 2018. Abrufdatum: 28.02.2019.
- [DSG] DSGVO: *Art. 4 DSGVO Begriffsbestimmungen*. <https://dsgvo-gesetz.de/art-4-dsgvo/>. Abrufdatum: 09.01.2019.
- [EAD09] ELDESOUKI, MOHAMED I, W ARAFA K DARWISH: *Stemming techniques of Arabic language: Comparative study from the information retrieval perspective. The Egyptian Computer Journal*, 36(1):30–49, 2009.
- [Fir] FIREEYE, INC: *Spear-Phishing-Angriffe ? Warum sie erfolgreich sind und wie sie gestoppt werden knnen*.
- [Fou18] FOUNDATION, PYTHON SOFTWARE: *html2text 2018.1.9*. <https://pypi.org/project/html2text/>, 2018. Abrufdatum: 15.03.2019.
- [Fou19] FOUNDATION, PYTHON SOFTWARE: *scrapy-selenium 0.0.7*. <https://pypi.org/project/scrapy-selenium/>, 2019. Abrufdatum: 16.03.2019.
- [Goo19] GOOGLE: *Refine web searches*. <https://support.google.com/websearch/answer/2466433?hl=en>, 2019. Abrufdatum: 27.02.2019.
- [Had11] HADNAGY, CHRISTOPHER: *Social Engineering: The Art of Human Hacking*. 2011.
- [Jam05] JAMES, LANCE: *Phshing Exposed: Uncover Secrets from the Dark Side*. 2005.
- [Law15] LAWSON, RICHARD: *Web scraping with Python*. Packt Publishing Ltd, 2015.
- [Lit16] LITZEL, NICO: *Was ist Natural Language Processing?* <https://www.bigdata-insider.de/was-ist-natural-language-processing-a-590102/>, 2016. Abrufdatum: 10.02.2019.

- [LLC19] LLC, GOOGLE: *marco lang tettnang - Google-Suche*. <https://www.google.com/search?q=marco+lang+tettnang>, 2019. Abrufdatum: 12.04.2019.
- [Med17] MEDIA, UNITED INTERNET: *Bürgerlicher Name als E-Mail-Adresse in Österreich und der Schweiz 2017*. <https://de.statista.com/statistik/daten/studie/745611/umfrage/buergerlicher-name-als-e-mail-adresse-in-oesterreich-und-der-schweiz/>, 2017. Abrufdatum: 31.10.2018.
- [Mit01] MITNICK, KEVIN D.: *The art of deception:controlling the human element of security*. 2001.
- [Mit15] MITCHELL, RYAN: *Web Scraping with Python: Collecting Data from the Modern Web*. 2015.
- [Mut18] MUTHUKADAN, BAIJU: *Selenium with Python*. <https://selenium-python.readthedocs.io/installation.html#introduction>, 2018. Abrufdatum: 27.02.2019.
- [NW18] NORDRHEIN-WESTFALEN, VERBRAUCHERZENTRALE: *Phishing-Radar: Aktuelle Warnungen*. <https://www.verbraucherzentrale.nrw/wissen/digitale-welt/phishingradar/phishingradar-aktuelle-warnungen-6059>, 2018. Abrufdatum: 29.10.2018.
- [PH] PHILIPP, JONAS NATHANAL GERHARD HEYER: *Multi-Label Klassifikation am Beispiel sozialwissenschaftlicher Texte*.
- [RECC10] ROSE, STUART, DAVE ENGEL, NICK CRAMER WENDY COWLEY: *Automatic keyword extraction from individual documents*. Text Mining: Applications and Theory, 1–20, 2010.
- [RFC94] *Uniform Resource Locators (URL)*. <https://tools.ietf.org/html/rfc1738#section-3.1>, 1994. Abrufdatum: 27.02.2019.
- [SG12] SHARMA, ARVIND KUMAR PC GUPTA: *Study & Analysis of Web Content Mining Tools to Improve Techniques of Web Data Mining*. International Journal of Advanced Research in Computer Engineering & Technology (IJARCET), 1(8):pp-287, 2012.
- [Sla] SLAVIN, TIM: *Stop Words*. <https://www.kidscodexs.com/stop-words/>. Abrufdatum: 29.01.2019.
- [SS11] SCHUBERT, SIGRID ANDREAS SCHWILL: *Didaktik der Informatik. Didaktik der Informatik*, 1–30. Springer, 2011.

-
- [Ste96] STEELE, ROBERT DAVID: *Open Source Intelligence: What Is It? Why Is It Important to the Military?* American Intelligence Journal, 35–41, 1996.
- [The01] THELWALL, MIKE: *A web crawler design for data mining.* Journal of Information Science, 27(5):319–325, 2001.
- [uDsiNe15] NETZ E.V., DATEV UND DEUTSCHLAND SICHER IM: *Verhaltensregeln zum Thema “Social Engineering”.* 2015.
- [W3S] W3SCHOOLS: *HTML URL Encoding Reference.* https://www.w3schools.com/tags/ref_urlencode.asp. Abrufdatum: 27.02.2019.